

Between Interpretation and Computation

Didactic Exploration of Computational Linguistics

Methods

with Augmented Transcripts of Sales Conversations

Paul Koop

Teaching and Learning Materials 2026

Abstract

This teaching and learning material serves the didactic exploration of computational linguistics methods based on eight transcripts of sales conversations. In contrast to the previous ARS versions 2.0 and 3.0, which were based on interpretively formed terminal symbols, this material takes the step toward automatic language processing. The methods are trained on augmented data for demonstration purposes to make their functioning transparent. The focus is on didactic knowledge acquisition, not on empirical validity. Scenarios C (Computational Linguistics Integration) and D (Hybrid Modeling) are developed step by step and compared with each other.

Contents

1	Introduction: Didactic Goals and Methodological Reflection	2
2	The Eight Transcripts: Raw Data and Terminal Symbols	3
2.1	The Raw Data	3
2.1.1	Transcript 1 - Butcher Shop	3
2.1.2	Transcript 2 - Market Square (Cherries)	3
2.1.3	Transcript 3 - Fish Stall	4
2.1.4	Transcript 4 - Vegetable Stall (Detailed)	4
2.1.5	Transcript 5 - Vegetable Stall (with KAV at Beginning)	4
2.1.6	Transcript 6 - Cheese Stand	4
2.1.7	Transcript 7 - Candy Stall	5
2.1.8	Transcript 8 - Bakery	5
2.2	The Terminal Symbol Strings (ARS 3.0)	6
3	Scenario C: Computational Linguistics Integration	7
3.1	Didactic Augmentation	7
3.2	Speech Act Recognition with Transformer Models	8
3.3	Word Embeddings and Semantic Similarity	12
3.4	Topic Modeling with BERTopic	15
3.5	Rhetorical Structure Theory (RST)	18
3.6	Integration of Components in Scenario C	23
4	Scenario D: Hybrid Modeling	28
4.1	CRF for Sequential Dependencies	28
4.2	Transformer Embeddings as Supplement	32
4.3	Graph Neural Networks for the Nonterminal Hierarchy	35
4.4	Attention Mechanisms for Relevant Predecessors	38
4.5	Integration of Components in Scenario D	42
5	Comparison of Scenarios and Methodological Reflection	47
5.1	Comparison of Approaches	47
5.2	Didactic Insights from Scenario C	47
5.3	Didactic Insights from Scenario D	47
5.4	Conclusion for Teaching Practice	48
6	Outlook	48

1 Introduction: Didactic Goals and Methodological Reflection

The previous versions of Algorithmic Recursive Sequence Analysis (ARS 2.0 and 3.0) have shown how formal grammars can be induced from interpretively obtained terminal symbol strings. These methods remain methodologically controlled: the category formation occurs through qualitative interpretation, the formal models merely explicate the observable regularities.

The following scenarios C and D venture a step beyond this methodological boundary. They explore how computational linguistics methods – especially neural networks, word embeddings, and topic models – could be applied to the eight transcripts if they were augmented for demonstration purposes.

This document is conceived as teaching and learning material. It pursues the following didactic goals:

1. **Understanding neural architectures:** How do transformers, LSTM networks, and attention mechanisms work on sequence data?
2. **Data augmentation as a technique:** How can one handle small datasets to demonstrate the functioning of methods?
3. **Comparison of different modeling levels:** What differences exist between purely computational linguistics (C) and hybrid (D) approaches?
4. **Methodological reflection:** Where are the limits of automatic methods compared to interpretive category formation?

All implementations presented here work with augmented data – the eight original transcripts were artificially multiplied to enable the training of neural networks. The results are therefore not empirically valid but serve exclusively for didactic illustration.

2 The Eight Transcripts: Raw Data and Terminal Symbols

2.1 The Raw Data

The following eight transcripts document sales conversations at Aachen market square in June/July 1994. They form the empirical basis for all subsequent analyses.

2.1.1 Transcript 1 - Butcher Shop

Date: June 28, 1994, **Location:** Butcher Shop, Aachen, 11:00 AM

```
1 Customer: Good day!
2 Salesperson: Good day!
3 Customer: One of the coarse liver sausage, please.
4 Salesperson: How much would you like?
5 Customer: Two hundred grams.
6 Salesperson: Anything else?
7 Customer: Yes, then also a piece of the Black Forest ham.
8 Salesperson: How large should the piece be?
9 Customer: Around three hundred grams.
10 Salesperson: That will be eight marks twenty.
11 Customer: Here you go.
12 Salesperson: Thank you and have a nice day!
13 Customer: Thanks, you too!
```

Listing 1: Transcript 1 - Raw Data

2.1.2 Transcript 2 - Market Square (Cherries)

Date: June 28, 1994, **Location:** Market Square, Aachen

```
1 Seller: Everyone can try cherries here!
2 Customer 1: Half a kilo of cherries, please.
3 Seller: Half a kilo? Or one kilo?
4 Seller: Three marks, please.
5 Customer 1: Thank you very much!
6 Seller: Everyone can try cherries here!
7 Customer 2: Half a kilo, please.
8 Seller: Three marks, please.
9 Customer 2: Thank you very much!
```

Listing 2: Transcript 2 - Raw Data

2.1.3 Transcript 3 - Fish Stall

Date: June 28, 1994, **Location:** Fish Stall, Market Square, Aachen

```
1 Customer: One pound of saithe, please.
2 Seller: Saithe, all right.
3 Seller: Four marks nineteen, please.
4 Customer: Thank you very much!
```

Listing 3: Transcript 3 - Raw Data

2.1.4 Transcript 4 - Vegetable Stall (Detailed)

Date: June 28, 1994, **Location:** Vegetable Stall, Aachen, Market Square, 11:00 AM

```
1 Customer: Listen, I'll take some mushrooms with me.
2 Seller: Brown or white?
3 Customer: Let's take the white ones.
4 Seller: They're both fresh, don't worry.
5 Customer: What about chanterelles?
6 Seller: Ah, they're great!
7 Customer: Can I put them in rice salad?
8 Seller: Better to briefly saut them in a pan.
9 Customer: Okay, I'll do that.
10 Seller: Have a nice day!
11 Customer: Likewise!
```

Listing 4: Transcript 4 - Raw Data

2.1.5 Transcript 5 - Vegetable Stall (with KAV at Beginning)

Date: June 26, 1994, **Location:** Vegetable Stall, Aachen, Market Square, 11:00 AM

```
1 Customer 1: Goodbye!
2 Customer 2: I would like a kilo of the Granny Smith apples here.
3 Seller: Anything else?
4 Customer 2: Yes, another kilo of onions.
5 Seller: Six marks twenty-five, please.
6 Customer 2: Goodbye!
```

Listing 5: Transcript 5 - Raw Data

2.1.6 Transcript 6 - Cheese Stand

Date: June 28, 1994, **Location:** Cheese Stand, Aachen, Market Square

```
1 Customer 1: Good morning!
2 Seller: Good morning!
3 Customer 1: I would like five hundred grams of Dutch Gouda.
4 Seller: In one piece?
5 Customer 1: Yes, in one piece, please.
```

Listing 6: Transcript 6 - Raw Data

2.1.7 Transcript 7 - Candy Stall

Date: June 28, 1994, **Location:** Candy Stall, Aachen, Market Square, 11:30 AM

```
1 Customer: I would like one hundred grams of the assorted ones.
2 Seller: For home or to take away?
3 Customer: To take away, please.
4 Seller: Fifty pfennigs, please.
5 Customer: Thanks!
```

Listing 7: Transcript 7 - Raw Data

2.1.8 Transcript 8 - Bakery

Date: July 9, 1994, **Location:** Bakery, Aachen, 12:00 PM

```
1 (Footsteps audible, background noises, partially unintelligible)
2 Customer: Good day!
3 (Unintelligible greeting in the background)
4 Salesperson: One of our best coffee, freshly ground, please.
5 (Noises of coffee grinder, packaging sounds)
6 Salesperson: Anything else?
7 Customer: Yes, also two pieces of fruit salad and a small bowl of
   cream.
8 Salesperson: All right!
9 (Noises of coffee grinder, paper sounds)
10 Salesperson: A small bowl of cream, yes?
11 Customer: Yes, thanks.
12 (Door noise, laughter, paper sounds)
13 Salesperson: Nobody takes care of oiling the doors.
14 Customer: Yes, that's always the case.
15 (Laughter, sounds of coins and packaging)
16 Salesperson: That will be fourteen marks and nineteen pfennigs,
   please.
17 Customer: I'll pay in small change.
18 (Laughter and sounds of coins)
19 Salesperson: Thank you very much, have a nice Sunday!
```

```
20 Customer: Thanks , you too!
```

Listing 8: Transcript 8 - Raw Data

2.2 The Terminal Symbol Strings (ARS 3.0)

For ARS 3.0, these raw data were converted into terminal symbol strings, which served as the basis for hierarchical grammar induction:

Table 1: Terminal Symbol Strings of the Eight Transcripts

Transcript	Terminal Symbol String
1 (Butcher)	KBG, VBG, KBBd, VBBd, KBA, VBA, KBBd, VBBd, KBA, VAA, KAA, VAV
2 (Cherries)	VBG, KBBd, VBBd, VAA, KAA, VBG, KBBd, VAA, KAA
3 (Fish)	KBBd, VBBd, VAA, KAA
4 (Vegetable)	KBBd, VBBd, KBA, VBA, KBBd, VBA, KAE, VAE, KAA, VAV, KAV
5 (Vegetable KAV)	KAV, KBBd, VBBd, KBBd, VAA, KAV
6 (Cheese)	KBG, VBG, KBBd, VBBd, KAA
7 (Candy)	KBBd, VBBd, KBA, VAA, KAA
8 (Bakery)	KBG, VBBd, KBBd, VBA, VAA, KAA, VAV, KAV

The meaning of the terminal symbols:

- **KBG**: Customer greeting
- **VBG**: Salesperson greeting
- **KBBd**: Customer need (concrete)
- **VBBd**: Salesperson inquiry
- **KBA**: Customer response
- **VBA**: Salesperson reaction
- **KAE**: Customer inquiry
- **VAE**: Salesperson information
- **KAA**: Customer completion
- **VAA**: Salesperson completion
- **KAV**: Customer farewell
- **VAV**: Salesperson farewell

3 Scenario C: Computational Linguistics Integration

Scenario C implements a fully computational linguistics modeling of the eight transcripts. It comprises four components:

1. **Speech Act Recognition:** Automatic recognition of speech acts from raw data
2. **Word Embeddings:** Vector representations of utterances
3. **Topic Modeling:** Identification of thematic shifts
4. **Rhetorical Structure Theory (RST):** Analysis of argumentative structure

3.1 Didactic Augmentation

Since neural networks require large amounts of data for training, the eight transcripts are augmented for demonstration purposes:

```
1 def augment_transcripts_for_teaching(transcripts, factor=20):
2     """
3     Augments the eight transcripts for didactic purposes.
4
5     Didactic note: This augmentation serves exclusively for
6     illustrating
7     the methodology. The resulting data are not empirically
8     valid but
9     merely enable demonstration of how neural methods
10    function.
11    """
12    augmented = []
13
14    # 1. Basic augmentation: simple copying
15    for _ in range(factor):
16        augmented.extend(transcripts)
17
18    # 2. Syntactic variations (didactically controlled)
19    import copy
20    import random
21
22    for transcript in transcripts:
```

```

20     for _ in range(factor // 4):
21         var = copy.deepcopy(transcript)
22         # Swap two adjacent utterances (rarely)
23         if len(var) > 3 and random.random() < 0.1:
24             idx = random.randint(0, len(var)-2)
25             var[idx], var[idx+1] = var[idx+1], var[idx]
26         augmented.append(var)
27
28     # 3. Lexical variations (synonyms)
29     synonyms = {
30         'Good day': ['Good morning', 'Hello', 'Good evening'
31                    ],
32         'Thanks': ['Thank you', 'Thank you very much', 'Merci
33                  '],
34         'Please': ['Please', 'You're welcome']
35     }
36
37     # Further variations could be implemented here
38
39     return augmented

```

Listing 9: Data Augmentation for Teaching Purposes

3.2 Speech Act Recognition with Transformer Models

Automatic recognition of speech acts is performed with a fine-tuned BERT model:

```

1  """
2  Speech Act Recognition with transformer-based models
3  Didactic implementation for teaching purposes
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 for 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-based classifier for speech acts
47     Didactically simplified architecture
48     """
49     def __init__(self, num_classes=12, dropout=0.3):
50         super().__init__()
51         self.bert = BertModel.from_pretrained('bert-base-
52             german-cased')
53         self.dropout = nn.Dropout(dropout)
54         self.classifier = nn.Linear(768, num_classes)

```

```

52
53     # Freeze BERT layers for didactic purposes (faster
54     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     Prepares data for speech act training
69     """
70     utterances = []
71     labels = []
72
73     # Extract all utterances from raw data
74     # Simplified: use terminal symbols directly for didactic
75     purposes
76     for trans, chain in zip(transcripts, terminal_chains):
77         for symbol in chain:
78             utterances.append(f"Example utterance for {symbol
79                             }")
80             labels.append(symbol)
81
82     # Label encoding
83     label_encoder = LabelEncoder()
84     y_encoded = label_encoder.fit_transform(labels)
85
86     return utterances, y_encoded, label_encoder
87
88 def train_speech_act_model(utterances, labels, epochs=10):
89     """
90     Trains the speech act recognition model
91     """

```

```

88     tokenizer = BertTokenizer.from_pretrained('bert-base-
      german-cased')
89     dataset = SpeechActDataset(utterances, labels, tokenizer)
90     dataloader = DataLoader(dataset, batch_size=8, shuffle=
      True)
91
92     model = BertSpeechActClassifier(num_classes=len(set(
      labels)))
93     optimizer = torch.optim.Adam(model.parameters(), lr=2e-5)
94     criterion = nn.CrossEntropyLoss()
95
96     print("\n=== Speech Act Recognition Training (Didactic)
      ===")
97     for epoch in range(epochs):
98         total_loss = 0
99         for batch in dataloader:
100             optimizer.zero_grad()
101             outputs = model(batch['input_ids'], batch['
      attention_mask'])
102             loss = criterion(outputs, batch['label'])
103             loss.backward()
104             optimizer.step()
105             total_loss += loss.item()
106
107             print(f"Epoch {epoch+1}: Loss = {total_loss/len(
      dataloader):.4f}")
108
109     return model, tokenizer, label_encoder
110
111 # Didactic note
112 print("\n" + "="*70)
113 print("DIDACTIC NOTE ON SPEECH ACT RECOGNITION")
114 print("="*70)
115 print("The implementation shown here uses augmented")
116 print("data and serves exclusively teaching purposes.")
117 print("Automatic recognition of speech acts would in practice
      :")
118 print("        Require millions of annotated training data")
119 print("        Be fine-tuned to specific domains (sales
      conversations)")

```

```
120 print("      Be subject to considerable uncertainties")
```

Listing 10: Speech Act Recognition with BERT

3.3 Word Embeddings and Semantic Similarity

For quantifying semantic similarity, pre-trained word embeddings are used:

```
1  """
2  Word Embeddings for Semantic Similarity Analysis
3  Didactic implementation with pre-trained models
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     Analyzes semantic similarities between utterances
15     """
16     def __init__(self, model_name='paraphrase-multilingual-
17         MiniLM-L12-v2'):
18         print(f"Loading pre-trained model: {model_name}")
19         self.model = SentenceTransformer(model_name)
20         self.embeddings = {}
21
22     def encode_utterances(self, utterances):
23         """
24         Creates embeddings for a list of utterances
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     def similarity_matrix(self, utterances):
32         """
33         Calculates similarity matrix for all utterances
```

```

33     """
34     embeddings = self.encode utterances(utterances)
35     sim_matrix = cosine_similarity(embeddings)
36     return sim_matrix
37
38 def find_similar(self, query, utterances, top_k=5):
39     """
40     Finds the most similar utterances to a query
41     """
42     query_emb = self.model.encode([query])[0]
43     utt_embs = self.encode utterances(utterances)
44
45     similarities = cosine_similarity([query_emb],
46                                     utt_embs)[0]
47     top_indices = np.argsort(similarities)[-top_k:][::-1]
48
49     results = []
50     for idx in top_indices:
51         results.append({
52             'utterance': utterances[idx],
53             'similarity': similarities[idx]
54         })
55
56     return results
57
58 def visualize_similarity(self, utterances, labels=None):
59     """
60     Visualizes similarity matrix as heatmap
61     """
62     sim_matrix = self.similarity_matrix(utterances)
63
64     plt.figure(figsize=(12, 10))
65     sns.heatmap(sim_matrix,
66                 xticklabels=labels if labels else range(
67                     len(utterances)),
68                 yticklabels=labels if labels else range(
69                     len(utterances)),
70                 cmap='viridis', vmin=0, vmax=1)
71     plt.title('Semantic Similarity Between Utterances')
72     plt.tight_layout()

```

```

70     plt.savefig('semantic_similarity.png', dpi=150)
71     plt.show()
72
73 # Didactic example
74 def demonstrate_semantic_analysis():
75     """
76     Demonstrates semantic analysis with examples
77     """
78     analyzer = SemanticAnalyzer()
79
80     # Example utterances from the transcripts
81     utterances = [
82         "Good day!",
83         "Good morning!",
84         "One liver sausage, please.",
85         "I would like sausage.",
86         "Thank you!",
87         "Thanks very much!",
88         "Goodbye!",
89         "Bye!"
90     ]
91
92     print("\n=== Semantic Similarity Analysis ===")
93
94     # Calculate similarity matrix
95     sim_matrix = analyzer.similarity_matrix(utterances)
96
97     # Most similar utterances to "Good day!"
98     similar = analyzer.find_similar("Good day!", utterances,
99                                     top_k=3)
100    print("\nMost similar to 'Good day!':")
101    for r in similar:
102        print(f"    {r['utterance']}: {r['similarity']:.3f}")
103
104    # Visualization
105    analyzer.visualize_similarity(utterances, utterances)
106
107    return analyzer
108
109 # Didactic note

```

```

109 print("\n" + "="*70)
110 print("DIDACTIC NOTE ON WORD EMBEDDINGS")
111 print("="*70)
112 print("The embeddings used were pre-trained on large corpora"
      )
113 print("(Wikipedia, news, web texts). They capture general")
114 print("linguistic similarities, not the specific categories")
115 print("of sales conversations.")

```

Listing 11: Semantic Similarity with Word Embeddings

3.4 Topic Modeling with BERTopic

For identifying thematic shifts, BERTopic is used:

```

1 """
2 Topic Modeling for Identifying Thematic Shifts
3 Didactic implementation with 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     Performs topic modeling on the transcripts
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         Prepares transcripts as documents for topic modeling
23         """
24         documents = []
25         metadata = []
26

```

```

27     for i, transcript in enumerate(transcripts, 1):
28         # Each transcript as one document
29         doc = ' '.join(transcript)
30         documents.append(doc)
31         metadata.append(f'Transcript {i}')
32
33     return documents, metadata
34
35 def fit_model(self, documents):
36     """
37     Trains the topic model
38     """
39     # Custom stop words
40     stopwords = ['please', 'thanks', 'thank', 'yes', 'no'
41                 ]
42     vectorizer = CountVectorizer(stop_words=stopwords)
43
44     self.model = BERTopic(
45         embedding_model="paraphrase-multilingual-MiniLM-
46             L12-v2",
47         vectorizer_model=vectorizer,
48         verbose=True,
49         nr_topics='auto'
50     )
51
52     self.topics, self.probs = self.model.fit_transform(
53         documents)
54     return self.topics, self.probs
55
56 def visualize_topics(self):
57     """
58     Visualizes the found topics
59     """
60     if self.model is None:
61         return
62
63     fig = self.model.visualize_topics()
64     fig.write_html("topic_visualization.html")
65
66     # Statistics

```

```

64     topic_counts = pd.Series(self.topics).value_counts()
65     print("\n=== Topic Distribution ===")
66     for topic, count in topic_counts.items():
67         if topic == -1:
68             print(f"Outlier: {count} documents")
69         else:
70             words = self.model.get_topic(topic)[:5]
71             words_str = ', '.join([w for w, _ in words])
72             print(f"Topic {topic}: {count} documents - {
73                 words_str}")
74 def demonstrate_topic_modeling(transcripts):
75     """
76     Demonstrates topic modeling on the transcripts
77     """
78     modeler = TranscriptTopicModeler()
79     documents, metadata = modeler.prepare_documents(
80         transcripts)
81     print("\n=== Topic Modeling of Eight Transcripts ===")
82     topics, probs = modeler.fit_model(documents)
83
84     for i, (doc, topic, prob, meta) in enumerate(zip(
85         documents, topics, probs, metadata)):
86         if topic != -1:
87             words = modeler.model.get_topic(topic)[:3]
88             words_str = ', '.join([w for w, _ in words])
89             print(f"{meta}: Topic {topic} (Confidence: {prob
90                 :.2f}) - {words_str}")
91         else:
92             print(f"{meta}: No clear topic (Outlier)")
93
94     modeler.visualize_topics()
95     return modeler
96
97 # Didactic note
98 print("\n" + "="*70)
99 print("DIDACTIC NOTE ON TOPIC MODELING")
100 print("="*70)
101 print("Topic modeling identifies latent themes in text

```

```

        corpora.")
100 print("With only eight documents, topic finding is unstable."
        )
101 print("The results therefore serve only to illustrate the")
102 print("methodology, not for substantive analysis.")

```

Listing 12: Topic Modeling with BERTopic

3.5 Rhetorical Structure Theory (RST)

For analyzing argumentative structure, an RST parser is implemented:

```

1  """
2  Rhetorical Structure Theory (RST) Analysis
3  Didactic implementation for sequence data
4  """
5
6  import networkx as nx
7  import matplotlib.pyplot as plt
8  from collections import defaultdict
9
10 class RSTRelation:
11     """Represents an RST relation between text segments"""
12     def __init__(self, type_name, nucleus, satellite,
13                 direction='nucleus-satellite'):
14         self.type = type_name # e.g., 'Elaboration', '
15             Contrast', 'Cause'
16         self.nucleus = nucleus # Central segment
17         self.satellite = satellite # Supporting segment
18         self.direction = direction
19
20 class SimpleRSTParser:
21     """
22     Simple RST parser for didactic purposes
23     Based on cue phrases and structural patterns
24     """
25     # Cue phrases for different relations
26     cue_phrases = {
27         'Elaboration': ['for example', 'in particular', '
28             namely', 'specifically'],

```

```

27     'Contrast': ['but', 'however', 'on the other hand', '
28         conversely'],
29     'Cause': ['because', 'since', 'therefore', 'thus', '
30         hence'],
31     'Condition': ['if', 'provided that', 'as long as'],
32     'Purpose': ['in order to', 'so that'],
33     'Sequence': ['then', 'after that', 'first', 'finally'
34 ]
35 }
36
37 def __init__(self):
38     self.relations = []
39     self.graph = nx.DiGraph()
40
41 def segment_transcript(self, transcript):
42     """
43     Segments a transcript into elementary discourse units
44     (EDUs)
45     Simplified: each utterance is an EDU
46     """
47     return transcript
48
49 def identify_relations(self, segments):
50     """
51     Identifies RST relations between segments
52     Didactically simplified implementation
53     """
54     relations = []
55
56     for i in range(len(segments)-1):
57         current = segments[i]
58         next_seg = segments[i+1]
59
60         # Check for cue phrases
61         for rel_type, cues in self.cue_phrases.items():
62             for cue in cues:
63                 if cue in current.lower() or cue in
64                     next_seg.lower():
65                     relations.append(RSTRelation(
66                         type_name=rel_type,

```

```

62         nucleus=i,
63         satellite=i+1
64     ))
65     break
66
67     # Default: Sequence 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     Builds an RST tree from identified relations
80     """
81     self.graph.clear()
82
83     # Add nodes
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     # Add edges
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     Complete RST analysis of a transcript
98     """
99     segments = self.segment_transcript(transcript)
100    relations = self.identify_relations(segments)
101    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 Structure"):
109     """
110     Visualizes the RST tree
111     """
112     pos = nx.spring_layout(self.graph)
113     plt.figure(figsize=(12, 8))
114
115     # Draw nodes
116     nx.draw_networkx_nodes(self.graph, pos, node_color='
117         lightblue',
118                             node_size=500)
119
120     # Draw edges with relation type as label
121     for edge in self.graph.edges(data=True):
122         nx.draw_networkx_edges(self.graph, pos, [(edge
123             [0], edge[1])])
124         nx.draw_networkx_edge_labels(
125             self.graph, pos,
126             {(edge[0], edge[1]): edge[2]['relation']}
127         )
128
129     # Node labels
130     labels = {node: f"{node}: {self.graph.nodes[node]['
131         text ']}"}
132     for node in self.graph.nodes():
133         nx.draw_networkx_labels(self.graph, pos, labels,
134             font_size=8)
135
136     plt.title(title)
137     plt.axis('off')
138     plt.tight_layout()
139     plt.savefig('rst_structure.png', dpi=150)
140     plt.show()

```

```

137
138 def demonstrate_rst_analysis(transcripts):
139     """
140     Demonstrates RST analysis on the transcripts
141     """
142     parser = SimpleRSTParser()
143
144     print("\n=== RST Analysis of Transcripts ===")
145
146     for i, transcript in enumerate(transcripts, 1):
147         print(f"\nTranscript {i}:")
148         result = parser.parse(transcript)
149
150         # Show identified relations
151         for rel in result['relations'][:5]: # Only first 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: # Visualize only first transcript
159             parser.visualize(f"RST Structure Transcript {i}")
160
161     return parser
162
163 # Didactic note
164 print("\n" + "="*70)
165 print("DIDACTIC NOTE ON RST ANALYSIS")
166 print("="*70)
167 print("The RST analysis implemented here is greatly
168     simplified.")
169 print("A full RST parser would:")
170 print("    Require extensive manual annotation")
171 print("    Work with trained neural models")
172 print("    Consider multiple hierarchy levels of discourse
173     relations")

```

Listing 13: Rhetorical Structure Theory Parser

3.6 Integration of Components in Scenario C

The complete integration of all components in Scenario C:

```
1  """
2  Scenario C: Complete Computational Linguistics Integration
3  Didactic implementation for teaching purposes
4  """
5
6  import os
7  import json
8  from datetime import datetime
9
10 class ScenarioC:
11     """
12     Integrates all computational linguistics components:
13     - Speech Act Recognition
14     - Word Embeddings / Semantic Analysis
15     - Topic Modeling
16     - RST Analysis
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("SCENARIO C: COMPUTATIONAL LINGUISTICS
26               INTEGRATION")
27         print("="*70)
28         print("\nThis scenario demonstrates the application
29               of")
30         print("computational linguistics methods to the eight
31               ")
32         print("transcripts. All results serve didactic
33               purposes")
34         print("and are not empirically valid.\n")
35
36     def run_speech_act_recognition(self):
37         """
```

```

34     Runs speech act recognition
35     """
36     print("\n--- Speech Act Recognition ---")
37     utterances, labels, encoder = prepare_speech_act_data
38         (
39         self.transcripts, self.terminal_chains
40         )
41
42     model, tokenizer, label_encoder =
43         train_speech_act_model(
44             utterances, labels, epochs=5
45         )
46
47     self.results['speech_act'] = {
48         'model': model,
49         'tokenizer': tokenizer,
50         'label_encoder': label_encoder,
51         'num_classes': len(label_encoder.classes_)
52     }
53
54     return self.results['speech_act']
55
56 def run_semantic_analysis(self):
57     """
58     Runs semantic similarity analysis
59     """
60     print("\n--- Semantic Similarity Analysis ---")
61     analyzer = SemanticAnalyzer()
62
63     # Collect all utterances
64     all_utterances = []
65     for transcript in self.transcripts:
66         all_utterances.extend(transcript)
67
68     # Similarity matrix
69     sim_matrix = analyzer.similarity_matrix(
70         all_utterances[:20]) # Only first 20

```

```

71         'utterances': all_utterances,
72         'similarity_matrix': sim_matrix
73     }
74
75     return self.results['semantic']
76
77     def run_topic_modeling(self):
78         """
79         Runs topic modeling
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         Runs RST analysis
101         """
102         print("\n--- RST Analysis ---")
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'],

```

```

110         'relations': [(r.type, r.nucleus, r.satellite
111             ) for r in result['relations']]
112     })
113
114     if i == 1:
115         parser.visualize(f"RST Structure Transcript {
116             i}")
117
118     self.results['rst'] = rst_results
119     return rst_results
120
121 def run_all(self):
122     """
123     Runs all analyses
124     """
125     self.run_speech_act_recognition()
126     self.run_semantic_analysis()
127     self.run_topic_modeling()
128     self.run_rst_analysis()
129
130     # Summary
131     print("\n" + "="*70)
132     print("SCENARIO C SUMMARY")
133     print("="*70)
134     print(f"    Speech Act Recognition: {self.results['
135         speech_act']['num_classes']} classes")
136     print(f"    Semantic Analysis: {len(self.results['
137         semantic']['utterances'])} utterances")
138     print(f"    Topic Modeling: {len(set(self.results['
139         topic']['topics']))} topics")
140     print(f"    RST Analysis: {len(self.results['rst'])}
141         transcripts analyzed")
142
143     return self.results
144
145 # Didactic execution
146 def run_scenario_c_demonstration():
147     """
148     Runs the complete demonstration of Scenario C
149     """

```

```

144 # Load transcripts
145 from ars_data import transcripts, terminal_chains
146
147 # Augment data for didactic purposes
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("DIDACTIC AUGMENTATION")
155 print("="*70)
156 print(f"Original: {len(transcripts)} transcripts")
157 print(f"Augmented: {len(augmented_transcripts)}
158     transcripts")
159
160 # Run Scenario C
161 scenario = ScenarioC(augmented_transcripts,
162     augmented_chains)
163 results = scenario.run_all()
164
165 # Save results
166 with open('scenario_c_results.json', 'w') as f:
167     # Convert non-serializable objects
168     serializable = {
169         'speech_act': {'num_classes': results['speech_act']
170             ['num_classes']},
171         'semantic': {'num_utterances': len(results['
172             semantic']['utterances'])},
173         'topic': {'num_topics': len(set(results['topic']
174             ['topics']))},
175         'rst': results['rst']
176     }
177     json.dump(serializable, f, indent=2)
178
179 print("\nResults saved to 'scenario_c_results.json'")
180
181 return results
182
183 if __name__ == "__main__":

```

```
177 run_scenario_c_demonstration()
```

Listing 14: Scenario C - Complete Integration

4 Scenario D: Hybrid Modeling

Scenario D integrates computational linguistics methods with the interpretively formed categories of ARS 3.0. It skips the complete automation of category formation (Scenario C) and uses the new methods complementarily.

4.1 CRF for Sequential Dependencies

Conditional Random Fields model dependencies of speech acts on the wider context:

```
1  """
2  Conditional Random Fields (CRF) for Sequential Dependencies
3  Didactic implementation with sklearn-crfsuite
4  """
5
6  import sklearn_crfsuite
7  from sklearn_crfsuite import metrics
8  import numpy as np
9
10 class CRFSequenceModel:
11     """
12     CRF model for sequence modeling of terminal symbols
13     """
14
15     def __init__(self):
16         self.crf = sklearn_crfsuite.CRF(
17             algorithm='lbfgs',
18             c1=0.1, # L1 regularization
19             c2=0.1, # L2 regularization
20             max_iterations=100,
21             all_possible_transitions=True
22         )
23         self.label_encoder = None
24
25     def word2features(self, tokens, i):
26         """
```

```

27     Creates features for position i in the sequence
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     # Context 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             ,
65             '+1:word.prefix_V': word_next.startswith('V')
66             ,

```

```

63         '+1:word.suffix_A': word_next.endswith('A'),
64     })
65     else:
66         features['EOS'] = True
67
68     return features
69
70     def extract_features(self, sequences):
71         """
72         Extracts features for all sequences
73         """
74         X = []
75         for seq in sequences:
76             X.append([self.word2features(seq, i) for i in
77                       range(len(seq))])
78         return X
79
80     def fit(self, sequences, labels):
81         """
82         Trains the CRF model
83         """
84         X = self.extract_features(sequences)
85         self.crf.fit(X, labels)
86         return self
87
88     def predict(self, sequences):
89         """
90         Predicts labels for new sequences
91         """
92         X = self.extract_features(sequences)
93         return self.crf.predict(X)
94
95     def evaluate(self, test_sequences, test_labels):
96         """
97         Evaluates the model
98         """
99         pred = self.predict(test_sequences)
100
101         # Flatten for metrics
102         y_true = [label for seq in test_labels for label in

```

```

    seq]
102     y_pred = [label for seq in pred for label in seq]
103
104     return {
105         'accuracy': np.mean(np.array(y_true) == np.array(
106             y_pred)),
107         'classification_report': metrics.
108             flat_classification_report(
109                 test_labels, pred, labels=sorted(set(y_true))
110             )
111     }
112
113 def demonstrate_crf(terminal_chains):
114     """
115     Demonstrates CRF modeling on terminal symbols
116     """
117     print("\n=== CRF Modeling of Terminal Symbols ===")
118
119     # Train-test split (didactic)
120     train_size = int(len(terminal_chains) * 0.7)
121     train_chains = terminal_chains[:train_size]
122     test_chains = terminal_chains[train_size:]
123
124     # Extract features
125     model = CRFSequenceModel()
126     X_train = model.extract_features(train_chains)
127
128     # Training
129     print(f"Training CRF with {len(train_chains)} sequences
130         ...")
131     model.fit(train_chains, train_chains) # Labels are the
132         sequences themselves
133
134     # Evaluation
135     results = model.evaluate(test_chains, test_chains)
136     print(f"\nAccuracy: {results['accuracy']:.3f}")
137
138     return model

```

Listing 15: CRF for Sequential Dependencies

4.2 Transformer Embeddings as Supplement

Transformer embeddings are used in addition to categorical terminal symbols:

```
1 """
2 Transformer Embeddings as Supplement to Categorical Terminal
   Symbols
3 """
4
5 import torch
6 import numpy as np
7 from sentence_transformers import SentenceTransformer
8
9 class TerminalEmbeddingEnricher:
10     """
11     Enriches terminal symbols with semantic embeddings of
       underlying utterances
12     """
13
14     def __init__(self, model_name='paraphrase-multilingual-
       MiniLM-L12-v2'):
15         self.model = SentenceTransformer(model_name)
16         self.symbol_to_embedding = {}
17         self.symbol_to_text = self._create_symbol_mapping()
18
19     def _create_symbol_mapping(self):
20         """
21         Creates a mapping from terminal symbols to example
           texts
22         """
23         return {
24             'KBG': ['Good day', 'Good morning', 'Hello'],
25             'VBG': ['Good day', 'Good morning', 'Hello back'
26                   ],
27             'KBBd': ['One liver sausage', 'I would like
28                     cheese', 'One kilo of apples please'],
29             'VBBd': ['How much would you like?', 'Which kind?
30                     ', 'Anything else?'],
31             'KBA': ['Two hundred grams', 'The white ones
32                     please', 'Yes, please'],
33             'VBA': ['All right', 'Coming right up', 'Okay'],
```

```

30         'KAE': ['Can I put that in salad?', 'Where are
31             these from?', 'Is it fresh?'],
32         'VAE': ['Better to saut ', 'From the region', '
33             Yes, very fresh'],
34         'KAA': ['Here you go', 'Thanks', 'Yes, thanks'],
35         'VAA': ['That will be 8 marks 20', '3 marks
36             please', '14 marks 19'],
37         'KAV': ['Goodbye', 'Bye', 'Have a nice day'],
38         'VAV': ['Thank you very much', 'Have a nice day',
39             'Goodbye']
40     }
41
42     def get_embedding(self, symbol):
43         """
44         Returns the embedding for a terminal symbol
45         """
46         if symbol in self.symbol_to_embedding:
47             return self.symbol_to_embedding[symbol]
48
49         # Average of example text embeddings
50         texts = self.symbol_to_text.get(symbol, [symbol])
51         embeddings = self.model.encode(texts)
52         avg_embedding = np.mean(embeddings, axis=0)
53
54         self.symbol_to_embedding[symbol] = avg_embedding
55         return avg_embedding
56
57     def enrich_sequence(self, sequence):
58         """
59         Enriches a sequence of terminal symbols with
60         embeddings
61         """
62         symbols = sequence
63         embeddings = np.array([self.get_embedding(sym) for
64             sym in symbols])
65
66         return {
67             'symbols': symbols,
68             'embeddings': embeddings,
69             'combined': np.column_stack([

```

```

64         self._one_hot_encode(symbols),
65         embeddings
66     ]) if len(symbols) > 0 else np.array([])
67 }
68
69 def _one_hot_encode(self, symbols):
70     """
71     One-hot encoding of terminal symbols
72     """
73     unique_symbols = sorted(set(self.symbol_to_text.keys
74                                ()))
75     symbol_to_idx = {sym: i for i, sym in enumerate(
76                       unique_symbols)}
77
78     one_hot = np.zeros((len(symbols), len(unique_symbols)
79                        ))
80     for i, sym in enumerate(symbols):
81         if sym in symbol_to_idx:
82             one_hot[i, symbol_to_idx[sym]] = 1
83
84     return one_hot
85
86 def demonstrate_embedding_enrichment():
87     """
88     Demonstrates enrichment of terminal symbols with
89     embeddings
90     """
91     enricher = TerminalEmbeddingEnricher()
92
93     print("\n=== Enrichment of Terminal Symbols with
94           Embeddings ===")
95
96     # Example sequence
97     sequence = ['KBG', 'VBG', 'KBBd', 'VBBd', 'KBA']
98
99     enriched = enricher.enrich_sequence(sequence)
100
101     print(f"\nSequence: {' '.join(sequence)}")
102     print(f"Embedding dimension: {enriched['embeddings'].
103           shape[1]}")

```

```

98     print(f"One-hot dimension: {enriched['combined'].shape[1]
99           - enriched['embeddings'].shape[1]}")
100    print(f"Combined dimension: {enriched['combined'].shape
101          [1]}")
    return enricher

```

Listing 16: Transformer Embeddings for Terminal Symbols

4.3 Graph Neural Networks for the Nonterminal Hierarchy

The nonterminal hierarchy is modeled as a Graph Neural Network:

```

1  """
2  Graph Neural Network for the Nonterminal Hierarchy
3  Didactic implementation with 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 for the grammar hierarchy
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)
22         self.conv2 = GCNConv(hidden_dim, hidden_dim)
23         self.classifier = nn.Linear(hidden_dim, num_classes)
24
25     def forward(self, x, edge_index):
26         x = self.conv1(x, edge_index)
27         x = F.relu(x)
28         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     Manages the GNN for the nonterminal hierarchy
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         Builds a graph from the grammar hierarchy
46         """
47         G = nx.DiGraph()
48
49         # Nodes: terminals and nonterminals
50         all_symbols = set()
51
52         # Nonterminals as nodes
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         # Edges: derivation relations
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     Prepares data for the GNN
70     """
71     # Node indices
72     nodes = list(self.graph.nodes())
73     node_to_idx = {node: i for i, node in enumerate(nodes
74                    )}
75
76     # Feature matrix (simplified: one-hot)
77     x = torch.eye(len(nodes))
78
79     # Edge 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
85     edge_index = torch.tensor(edge_index, dtype=torch.
86                               long).t().contiguous()
87
88     return Data(x=x, edge_index=edge_index)
89
90 def train(self, epochs=100):
91     """
92     Trains the GNN
93     """
94     data = self.prepare_data()
95     self.model = GrammarGNN(input_dim=data.x.shape[1])
96
97     optimizer = torch.optim.Adam(self.model.parameters(),
98                                  lr=0.01)
99
100    print("\n=== Training Grammar GNN ===")
101    for epoch in range(epochs):
102        self.model.train()
103        optimizer.zero_grad()
104        out = self.model(data.x, data.edge_index)
105
106        # Self-supervised learning: graph reconstruction
107        # Simplified: predict neighbors

```

```

104         loss = F.nll_loss(out[data.edge_index[0]], data.
105                             edge_index[1])
106
107         loss.backward()
108         optimizer.step()
109
110         if epoch % 20 == 0:
111             print(f"Epoch {epoch}: Loss = {loss.item():.4
112                   f}")
113
114     return self.model
115
116 def demonstrate_gnn(grammar_rules):
117     """
118     Demonstrates GNN for the grammar hierarchy
119     """
120     print("\n=== Graph Neural Network for Nonterminal
121           Hierarchy ===")
122
123     gnn = GrammarHierarchyGNN(grammar_rules)
124     print(f"Graph: {gnn.graph.number_of_nodes()} nodes, "
125           f"{gnn.graph.number_of_edges()} edges")
126
127     model = gnn.train(epochs=100)
128
129     return gnn, model

```

Listing 17: Graph Neural Network for Nonterminal Hierarchy

4.4 Attention Mechanisms for Relevant Predecessors

Attention mechanisms identify particularly relevant predecessors for current decisions:

```

1 """
2 Attention Mechanisms for Identifying Relevant Predecessors
3 """
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 mechanism for sequence modeling
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 parameters
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         # Compute Query, Key, Value
34         Q = self.W_q(x) # (seq_len, batch, hidden_dim)
35         K = self.W_k(x) # (seq_len, batch, hidden_dim)
36         V = self.W_v(x) # (seq_len, batch, hidden_dim)
37
38         # Attention scores
39         scores = torch.matmul(Q.transpose(0, 1), K.transpose
40             (0, 1).transpose(1, 2))
41         scores = scores / self.scale
42
43         if mask is not None:
44             scores = scores.masked_fill(mask == 0, -1e9)
45
46         # Attention weights
47         attention_weights = F.softmax(scores, dim=-1)
48

```

```

45     # Weighted sum
46     context = torch.matmul(attention_weights, V.transpose
47         (0, 1))
48
49     return context, attention_weights
50
51 class SymbolPredictorWithAttention(nn.Module):
52     """
53     Predicts the next symbol with attention on predecessors
54     """
55     def __init__(self, num_symbols, embedding_dim=50,
56         hidden_dim=64):
57         super().__init__()
58         self.embedding = nn.Embedding(num_symbols,
59             embedding_dim)
60         self.attention = SequenceAttention(embedding_dim,
61             hidden_dim)
62         self.lstm = nn.LSTM(embedding_dim, hidden_dim,
63             batch_first=True)
64         self.classifier = nn.Linear(hidden_dim +
65             embedding_dim, num_symbols)
66
67     def forward(self, x):
68         """
69         x: (batch, seq_len) with symbol indices
70         """
71         # Embeddings
72         embedded = self.embedding(x) # (batch, seq_len,
73             embedding_dim)
74
75         # LSTM for sequential dependencies
76         lstm_out, (hidden, cell) = self.lstm(embedded)
77
78         # Attention over the sequence
79         # Transpose for attention (seq_len, batch,
80             embedding_dim)
81         context, attention_weights = self.attention(embedded.
82             transpose(0, 1))

```

```

76     # Combine last LSTM state with attention context
77     last_hidden = hidden[-1] # (batch, hidden_dim)
78     last_context = context[-1] # (batch, hidden_dim)
79
80     # Prediction
81     combined = torch.cat([last_hidden, last_context], dim
82                          =-1)
83     logits = self.classifier(combined)
84
85     return logits, attention_weights
86
87 def demonstrate_attention(terminal_chains, symbol_to_idx):
88     """
89     Demonstrates attention mechanisms on the sequences
90     """
91     print("\n=== Attention Mechanisms for Relevant
92           Predecessors ===")
93
94     # Prepare data
95     sequences = []
96     for chain in terminal_chains:
97         seq = [symbol_to_idx[sym] for sym in chain]
98         sequences.append(seq)
99
100    # Padding for batch processing
101    from torch.nn.utils.rnn import pad_sequence
102    sequences_padded = pad_sequence([torch.tensor(seq) for
103                                   seq in sequences],
104                                   batch_first=True,
105                                   padding_value=0)
106
107    # Initialize model
108    model = SymbolPredictorWithAttention(num_symbols=len(
109        symbol_to_idx))
110
111    # Forward pass
112    logits, attention_weights = model(sequences_padded[:2])
113    # Only first 2 sequences
114
115    print(f"\nInput shape: {sequences_padded[:2].shape}")

```

```

110     print(f"Attention weights shape: {attention_weights.shape
111           }")
112
113     print(f"Logits shape: {logits.shape}")
114
115     # Visualize attention weights
116     plot_attention_weights(attention_weights[0].detach().
117         numpy(),
118                             sequences[0], sequences[0])
119
120     return model
121
122 def plot_attention_weights(attention, source_tokens,
123     target_tokens):
124     """
125     Visualizes attention weights as heatmap
126     """
127     import matplotlib.pyplot as plt
128     import seaborn as sns
129
130     plt.figure(figsize=(10, 8))
131     sns.heatmap(attention[:len(target_tokens), :len(
132         source_tokens)],
133                 xticklabels=source_tokens,
134                 yticklabels=target_tokens,
135                 cmap='viridis', annot=True, fmt='.2f')
136     plt.title('Attention Weights Between Predecessors and
137         Prediction')
138     plt.xlabel('Predecessor Symbols')
139     plt.ylabel('Prediction Position')
140     plt.tight_layout()
141     plt.savefig('attention_weights.png', dpi=150)
142     plt.show()

```

Listing 18: Attention Mechanisms for Sequence Modeling

4.5 Integration of Components in Scenario D

The complete integration of all components in Scenario D:

```

1 """
2 Scenario D: Hybrid Modeling

```

```

3 Integration of computational linguistics methods with
  interpretive categories
4 """
5
6 import json
7 import numpy as np
8
9 class ScenarioD:
10     """
11     Integrates computational linguistics methods
12     complementarily to the
13     interpretively formed categories of ARS 3.0
14     """
15     def __init__(self, terminal_chains, grammar_rules,
16                 reflection_log):
17         self.terminal_chains = terminal_chains
18         self.grammar_rules = grammar_rules
19         self.reflection_log = reflection_log
20         self.results = {}
21
22         print("\n" + "="*70)
23         print("SCENARIO D: HYBRID MODELING")
24         print("="*70)
25         print("\nThis scenario integrates computational
26               linguistics")
27         print("methods COMPLEMENTARILY to the interpretive")
28         print("categories of ARS 3.0. The interpretive basis"
29               )
30         print("is preserved but enriched by new methods.\n")
31
32     def run_crf_modeling(self):
33         """
34         Runs CRF modeling on terminal symbols
35         """
36         print("\n--- CRF Modeling ---")
37         crf_model = demonstrate_crf(self.terminal_chains)
38         self.results['crf'] = {'model': crf_model}
39         return crf_model

```

```

38     def run_embedding_enrichment(self):
39         """
40         Enriches terminal symbols with transformer embeddings
41         """
42         print("\n--- Embedding Enrichment ---")
43         enricher = demonstrate_embedding_enrichment()
44
45         # Example enriched sequence
46         example_seq = self.terminal_chains[0][:5]
47         enriched = enricher.enrich_sequence(example_seq)
48
49         self.results['embeddings'] = {
50             'enricher': enricher,
51             'example': enriched
52         }
53
54         return enricher
55
56     def run_gnn_hierarchy(self):
57         """
58         Models the nonterminal hierarchy as GNN
59         """
60         print("\n--- GNN for Nonterminal Hierarchy ---")
61         gnn, model = demonstrate_gnn(self.grammar_rules)
62         self.results['gnn'] = {'gnn': gnn, 'model': model}
63         return gnn, model
64
65     def run_attention_analysis(self):
66         """
67         Analyzes attention mechanisms on the sequences
68         """
69         print("\n--- Attention Analysis ---")
70
71         # Symbol to index mapping
72         all_symbols = set()
73         for chain in self.terminal_chains:
74             all_symbols.update(chain)
75         symbol_to_idx = {sym: i for i, sym in enumerate(
76             sorted(all_symbols))}

```

```

77     model = demonstrate_attention(self.terminal_chains,
78         symbol_to_idx)
79
80     self.results['attention'] = {'model': model}
81
82     return model
83
84 def run_all(self):
85     """
86     Runs all analyses (complementary, not substitutive)
87     """
88     self.run_crf_modeling()
89     self.run_embedding_enrichment()
90     self.run_gnn_hierarchy()
91     self.run_attention_analysis()
92
93     # Summary
94     print("\n" + "="*70)
95     print("SCENARIO D SUMMARY")
96     print("="*70)
97     print("    CRF Modeling: Sequential dependencies
98         modeled")
99     print("    Embedding Enrichment: Terminal symbols
100         semantically enriched")
101     print("    GNN Hierarchy: Nonterminal structure as
102         graph")
103     print("    Attention Analysis: Relevant predecessors
104         identified")
105     print("\nThe interpretive categories of ARS 3.0
106         remain")
107     print("the foundation of all analyses. Computational")
108     print("linguistics methods serve complementary
109         insight.")
110
111     return self.results
112
113 def run_scenario_d_demonstration(terminal_chains,
114     grammar_rules, reflection_log):
115     """
116     Runs the complete demonstration of Scenario D

```

```

108     """
109     scenario = ScenarioD(terminal_chains, grammar_rules,
110                          reflection_log)
111     results = scenario.run_all()
112
113     # Save results
114     with open('scenario_d_results.json', 'w') as f:
115         # Simplified serializable version
116         serializable = {
117             'crf': {'status': 'completed'},
118             'embeddings': {'status': 'completed'},
119             'gnn': {'num_nodes': results['gnn'][0].graph.
120                   number_of_nodes()},
121             'attention': {'status': 'completed'}
122         }
123         json.dump(serializable, f, indent=2)
124
125     print("\nResults saved to 'scenario_d_results.json'")
126
127     return results
128
129 # Didactic note
130 print("\n" + "="*70)
131 print("METHODOLOGICAL NOTE ON SCENARIO D")
132 print("="*70)
133 print("Scenario D preserves the interpretive basis of ARS
134       3.0.")
135 print("The computational linguistics methods are used
136       COMPLEMENTARILY,")
137 print("not as a replacement for manual category formation.")
138 print("This corresponds to the methodological demand for")
139 print("control and transparency in the sense of XAI criteria.
140       ")

```

Listing 19: Scenario D - Complete Hybrid Integration

5 Comparison of Scenarios and Methodological Reflection

5.1 Comparison of Approaches

Table 2: Comparison of Scenarios C and D

Criterion	Scenario C	Scenario D
Category Formation	Automatic (Speech Act Recognition)	Interpretive (ARS 3.0)
Data Basis	Augmented raw data	Terminal symbol strings
Representation	Vector embeddings	Discrete categories + embeddings
Hierarchy	Automatically learned	Explicitly induced (ARS 3.0)
Transparency	Low (black box)	High (documented decisions)
Didactic Value	Functioning of neural methods	Integration of old and new methods
Empirical Validity	Not given	Limited (based on interpretation)
Methodological Control	Lost	Preserved

5.2 Didactic Insights from Scenario C

The implementation of Scenario C has shown:

- Need for large data volumes:** Neural methods require data volumes far exceeding the eight transcripts for valid results. Augmentation enables demonstration of functioning but does not replace real data.
- Opacity of decisions:** Automatically learned categories and attention weights are not easily comprehensible to third parties. The XAI criteria of meaningfulness and transparency are violated.
- Loss of interpretive basis:** Automatic speech act recognition does not capture the qualitatively meaningful distinctions of ARS (e.g., between KBA and KAA) but learns statistical correlations in vector space.

5.3 Didactic Insights from Scenario D

The implementation of Scenario D has shown:

1. **Complementarity instead of substitution:** Computational linguistics methods can provide valuable additional information (e.g., semantic similarities between different utterances) without replacing the interpretive basis.
2. **Validation possibilities:** Embedding similarities can be used to validate interpretive category formation: similar utterances should receive similar terminal symbols.
3. **Visualization of dependencies:** Attention mechanisms and CRF models visualize which predecessors are particularly relevant for current decisions – this can illustrate the sequential structure of conversations.
4. **Methodological control preserved:** Since interpretive categories form the foundation, all results remain tied back to qualitative decisions and are intersubjectively verifiable.

5.4 Conclusion for Teaching Practice

The didactic exploration of Scenarios C and D leads to the following conclusions:

1. **Scenario C is suitable for demonstrating the functioning** of neural methods but should be used with explicit note of lacking empirical validity and methodological problems.
2. **Scenario D is methodologically preferable** as it preserves the interpretive basis and uses computational linguistics methods complementarily. It conveys how old and new methods can be productively combined.
3. **Data augmentation is a valuable didactic tool** to demonstrate the functioning of methods with small datasets. The augmented nature of the data must always be made transparent.
4. **The XAI criteria** (meaningfulness, accuracy, knowledge limits) provide a suitable framework to evaluate different approaches and reflect on their strengths and weaknesses.

6 Outlook

The didactic implementations presented here can be further developed in several directions:

1. **Extension of augmentation strategies:** Beyond simple copying, more complex augmentations (paraphrasing, controlled variation) could be implemented.
2. **Integration of further methods:** e.g., Petri nets for concurrency, Bayesian networks for inference, or formal verification methods.
3. **Development of comparison metrics:** How can the results of different scenarios be compared quantitatively without losing the qualitative basis?
4. **Transfer to other datasets:** The methodology can be transferred to other interaction types (doctor-patient conversations, classroom interactions, etc.).

What remains crucial throughout is methodological control: the formal procedures must respect the interpretive character of the analysis and must not lead to its automation.

References

- Barredo Arrieta, A., Díaz-Rodríguez, N., Del Ser, J., Bennetot, A., Tabik, S., Barbado, A., Garcia, S., Gil-Lopez, 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 Social Research: An Introduction* (9th ed.). Rowohlt. [German original]
- 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 Social Research: A Workbook* (5th ed.). De Gruyter Oldenbourg. [German original]
- 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.