# Algorithmic Recursive Sequence Analysis 4.0

Hybrid Integration of Computational Linguistics
Methods
as a Complementary Extension of ARS 3.0

Paul Koop

2026

**Abstract**

This paper develops a hybrid integration of computational linguistics methods into the Algorithmic Recursive Sequence Analysis (ARS). In contrast to Scenario C, which aims for complete automation of category formation, here computational linguistics methods are used complementarily to the interpretively obtained categories of ARS 3.0. The integration includes Conditional Random Fields (CRF) for sequential dependencies, Transformer embeddings for semantic enrichment, Graph Neural Networks (GNN) for the nonterminal hierarchy, and attention mechanisms for identifying relevant predecessors. Methodological control is maintained since the interpretive categories form the basis of all analyses and the computational linguistics methods merely open up additional dimensions of insight. The application to eight transcripts of sales conversations demonstrates the added value of this complementary integration.

# Contents

# 1 Introduction: Complementarity Instead of Substitution

ARS 3.0 has shown how hierarchical grammars can be induced from interpretively obtained terminal symbol strings. These grammars are transparent, intersubjectively verifiable, and methodologically controlled. They form the foundation for all further analyses.

The computational linguistics methods developed in Scenario C offer additional analytical perspectives:

- **Conditional Random Fields** model sequential dependencies with context

- **Transformer embeddings** quantify semantic similarities

- **Graph Neural Networks** capture structural relationships

- **Attention mechanisms** identify relevant predecessors

Unlike in Scenario C, these methods are not used here to automate category formation but as a complementary extension. The interpretive categories remain the foundation – the computational linguistics methods open up additional dimensions of insight without compromising methodological control.

# 2 Theoretical Foundations

## 2.1 Conditional Random Fields (CRF)

Conditional Random Fields (Lafferty et al., 2001) are probabilistic graphical models for segmentation and labeling of sequence data. Unlike HMMs, they directly model the conditional probability $P(Y|X)$ and can incorporate arbitrarily many contextual features.

For ARS 4.0, CRFs are used to model the dependence of terminal symbols on the wider context – not just on the immediate predecessor.

## 2.2 Transformer Embeddings

Transformer embeddings (Devlin et al., 2019; Reimers & Gurevych, 2019) generate contextualized vector representations of texts. Unlike static word embeddings, they take into account the entire sentence context.

For ARS 4.0, Transformer embeddings are used to quantify semantic similarity between different utterances – even those that received different terminal symbols.

## 2.3 Graph Neural Networks (GNN)

Graph Neural Networks (Scarselli et al., 2009) operate directly on graph structures and learn representations for nodes considering their neighbors.

For ARS 4.0, the nonterminal hierarchy is modeled as a graph, where nodes represent terminals and nonterminals, and edges represent derivation relations.

## 2.4 Attention Mechanisms

Attention mechanisms (Vaswani et al., 2017) allow models to focus differently on various parts of the input when making predictions.

For ARS 4.0, attention mechanisms are used to identify which predecessors are particularly relevant for predicting the next symbol.

# 3 Methodology: Complementary Integration

## 3.1 CRF for Sequential Dependencies

CRFs are trained on the terminal symbol strings to learn which contextual factors influence the choice of the next symbol. The features include:

- Current symbol

- Previous symbol

- Next symbol (if known)

- Position in sequence

- Speaker change indicator

- Phase indicator (from HMM)

## 3.2 Transformer Embeddings for Semantic Validation

Transformer embeddings are used to calculate semantic similarity between utterances that received the same terminal symbol. This serves to validate the interpretive category formation:

- High similarity within a category speaks for consistent interpretation

- Overlaps between categories can indicate interpretation flexibility

## 3.3   GNN for Structure Analysis

The nonterminal hierarchy is modeled as a graph and analyzed with a GNN. This enables:

- Identification of central nodes (important nonterminals)

- Recognition of patterns in the derivation structure

- Visualization of the hierarchy as an embedding space

## 3.4   Attention for Relevant Contexts

Attention mechanisms are trained on the sequences to visualize which predecessors are particularly important for predicting the next symbol. This can:

- Confirm the plausibility of the ARS grammar

- Point to previously overlooked dependencies

- Illustrate the sequential structure of conversations

# 4   Implementation

```
"""
ARS 4.0 - Hybrid Integration
Complementary use of computational linguistics methods
with interpretive categories of ARS 3.0
"""

import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from collections import defaultdict
import networkx as nx
from sklearn_crfsuite import CRF
from sentence_transformers import SentenceTransformer
import torch
import torch.nn as nn
```

```python
import torch.nn.functional as F

# ============================================================================

# 1. CONDITIONAL RANDOM FIELDS (CRF)
# ============================================================================


class ARSCRFModel:
    """
    CRF model for sequential dependencies in terminal symbol
        strings
    """

    def __init__(self):
        self.crf = CRF(
            algorithm='lbfgs',
            c1=0.1,  # L1 regularization
            c2=0.1,  # L2 regularization
            max_iterations=100,
            all_possible_transitions=True
        )
        self.feature_names = []

    def extract_features(self, sequence, i):
        """
        Extracts features for position i in the sequence
        """
        features = {
            'bias': 1.0,
            'symbol': sequence[i],
            'symbol.prefix_K': sequence[i].startswith('K'),
            'symbol.prefix_V': sequence[i].startswith('V'),
            'symbol.suffix_A': sequence[i].endswith('A'),
            'symbol.suffix_B': sequence[i].endswith('B'),
            'symbol.suffix_E': sequence[i].endswith('E'),
            'symbol.suffix_G': sequence[i].endswith('G'),
            'symbol.suffix_V': sequence[i].endswith('V'),
```

```python
                'position': i,
                'is_first': i == 0,
                'is_last': i == len(sequence) - 1,
            }

        # Context features (-2, -1, +1, +2)
        for offset in [-2, -1, 1, 2]:
            if 0 <= i + offset < len(sequence):
                sym = sequence[i + offset]
                features[f'context_{offset:+d}'] = sym
                features[f'context_{offset:+d}.prefix_K'] = \
                    sym.startswith('K')
                features[f'context_{offset:+d}.prefix_V'] = \
                    sym.startswith('V')

        # Bigram features
        if i > 0:
            features['bigram'] = f"{sequence[i-1]}_{sequence[
                i]}"

        return features

    def prepare_data(self, sequences):
        """
        Prepares data for CRF training
        """
        X = []
        y = []

        for seq in sequences:
            X_seq = [self.extract_features(seq, i) for i in
                range(len(seq))]
            y_seq = [sym for sym in seq]
            X.append(X_seq)
            y.append(y_seq)

        return X, y

    def fit(self, sequences):
        """
```

```python
        Trains the CRF model
        """
        print("\n=== CRF Training ===")
        X, y = self.prepare_data(sequences)
        self.crf.fit(X, y)

        # Show top features
        self.print_top_features()

        return self

    def predict(self, sequence):
        """
        Predicts labels for a sequence
        """
        X = [self.extract_features(sequence, i) for i in
            range(len(sequence))]
        return self.crf.predict([X])[0]

    def print_top_features(self):
        """
        Shows the most important CRF features
        """
        print("\nTop 20 CRF Features:")
        top_features = sorted(
            self.crf.state_features_.items(),
            key=lambda x: abs(x[1]),
            reverse=True
        )[:20]

        for (attr, label), weight in top_features:
            print(f"  {attr:30s} -> {label:4s} : {weight:+.4f
                }")

#
    ==========================================================================


# 2. TRANSFORMER EMBEDDINGS FOR SEMANTIC VALIDATION
#
    ==========================================================================
```

7

```python
class SemanticValidator:
    """
    Validates interpretive categories with Transformer
        embeddings
    """

    def __init__(self, model_name='paraphrase-multilingual-
        MiniLM-L12-v2'):
        print(f"\n=== Loading Sentence-Transformer: {
            model_name} ===")
        self.model = SentenceTransformer(model_name)
        self.symbol_to_texts = self._create_text_mapping()
        self.embeddings = {}

    def _create_text_mapping(self):
        """
        Creates mapping from terminal symbols to example
            texts
        """
        return {
            'KBG': ['Good day', 'Good morning', 'Hello', '
                Greetings'],
            'VBG': ['Good day', 'Good morning', 'Hello back',
                 'Welcome'],
            'KBBd': ['One liver sausage', 'I would like
                cheese', 'One kilo of apples please'],
            'VBBd': ['How much would you like?', 'Which kind?
                ', 'Anything else?'],
            'KBA': ['Two hundred grams', 'The white ones
                please', 'Yes please'],
            'VBA': ['All right', 'Coming right up', 'Okay'],
            'KAE': ['Can I put that in salad?', 'Where are
                these from?', 'Is it fresh?'],
            'VAE': ['Better to saut ', 'From the region', '
                Yes, very fresh'],
            'KAA': ['Here you go', 'Thanks', 'Yes thanks'],
            'VAA': ['That will be 8 marks 20', '3 marks
                please', '14 marks 19'],
```

```python
            'KAV': ['Goodbye', 'Bye', 'Have a nice day'],
            'VAV': ['Thank you very much', 'Have a nice day',
                'Goodbye']
        }

    def compute_category_embeddings(self):
        """
        Computes average embeddings for each category
        """
        for symbol, texts in self.symbol_to_texts.items():
            embeddings = self.model.encode(texts)
            self.embeddings[symbol] = np.mean(embeddings,
                axis=0)

        return self.embeddings

    def compute_similarity_matrix(self):
        """
        Computes similarity matrix between categories
        """
        if not self.embeddings:
            self.compute_category_embeddings()

        symbols = sorted(self.embeddings.keys())
        n = len(symbols)
        sim_matrix = np.zeros((n, n))

        for i, sym1 in enumerate(symbols):
            for j, sym2 in enumerate(symbols):
                emb1 = self.embeddings[sym1]
                emb2 = self.embeddings[sym2]
                sim = np.dot(emb1, emb2) / (np.linalg.norm(
                    emb1) * np.linalg.norm(emb2))
                sim_matrix[i, j] = sim

        return sim_matrix, symbols

    def validate_categories(self):
        """
        Validates the interpretive categories
```

```python
186         """
187         print("\n=== Validation of Interpretive Categories
                === ")
188
189         sim_matrix, symbols = self.compute_similarity_matrix
                ()
190
191         # Statistics per category
192         print("\nIntra-category similarity (cohesion):")
193         for i, sym in enumerate(symbols):
194             intra = sim_matrix[i, i]
195             print(f"  {sym}: {intra:.3f}")
196
197         # Inter-category similarity
198         print("\nInter-category similarity (top 10):")
199         similarities = []
200         for i in range(len(symbols)):
201             for j in range(i+1, len(symbols)):
202                 similarities.append((symbols[i], symbols[j],
                        sim_matrix[i, j]))
203
204         similarities.sort(key=lambda x: x[2], reverse=True)
205         for sym1, sym2, sim in similarities[:10]:
206             print(f"  {sym1} - {sym2}: {sim:.3f}")
207
208         # Visualization
209         self.visualize_similarity_matrix(sim_matrix, symbols)
210
211         return sim_matrix, symbols
212
213     def visualize_similarity_matrix(self, sim_matrix, symbols
            ):
214         """
215         Visualizes the similarity matrix as heatmap
216         """
217         plt.figure(figsize=(12, 10))
218         sns.heatmap(sim_matrix,
219                     xticklabels=symbols,
220                     yticklabels=symbols,
221                     cmap='viridis',
```

```python
                      vmin=0, vmax=1,
                      annot=True, fmt='.2f')
        plt.title('Semantic Similarity Between Terminal
            Symbol Categories')
        plt.tight_layout()
        plt.savefig('category_similarity.png', dpi=150)
        plt.show()


# =============================================================================

# 3. GRAPH NEURAL NETWORK FOR NONTERMINAL HIERARCHY
# =============================================================================


class GrammarGraph:
    """
    Represents the ARS grammar as a graph
    """

    def __init__(self, grammar_rules):
        self.grammar = grammar_rules
        self.graph = nx.DiGraph()
        self.build_graph()

    def build_graph(self):
        """
        Builds a directed graph from the grammar
        """
        for nt, productions in self.grammar.items():
            for prod, prob in productions:
                for sym in prod:
                    self.graph.add_edge(nt, sym, weight=prob,
                        type='derivation')

        # Calculate metrics
        print("\n=== Grammar Graph Analysis ===")
        print(f"Nodes: {self.graph.number_of_nodes()}")
        print(f"Edges: {self.graph.number_of_edges()}")
```

```python
256
257          # Centrality
258          if self.graph.number_of_nodes() > 0:
259              centrality = nx.degree_centrality(self.graph)
260              top_nodes = sorted(centrality.items(), key=lambda
                      x: x[1], reverse=True)[:5]
261              print("\nTop 5 nodes by centrality:")
262              for node, cent in top_nodes:
263                  print(f"  {node}: {cent:.3f}")
264
265      def visualize(self, filename="grammar_graph.png"):
266          """
267          Visualizes the grammar graph
268          """
269          plt.figure(figsize=(15, 10))
270
271          # Layout
272          pos = nx.spring_layout(self.graph, k=2, iterations
                  =50)
273
274          # Color nodes by type
275          node_colors = []
276          for node in self.graph.nodes():
277              if node.startswith('NT_'):
278                  node_colors.append('lightgreen')  #
                          Nonterminals
279              else:
280                  node_colors.append('lightblue')   # Terminals
281
282          nx.draw(self.graph, pos,
283                  node_color=node_colors,
284                  with_labels=True,
285                  node_size=1000,
286                  font_size=8,
287                  arrows=True,
288                  arrowsize=20,
289                  edge_color='gray',
290                  alpha=0.7)
291
292          plt.title('ARS Grammar as Graph')
```

```python
            plt.tight_layout()
            plt.savefig(filename, dpi=150)
            plt.show()


class SimpleGNN(nn.Module):
    """
    Simple Graph Neural Network for analysis purposes
    """

    def __init__(self, input_dim, hidden_dim=16, output_dim
        =8):
        super().__init__()
        self.conv1 = nn.Linear(input_dim, hidden_dim)
        self.conv2 = nn.Linear(hidden_dim, hidden_dim)
        self.output = nn.Linear(hidden_dim, output_dim)

    def forward(self, x, adj):
        # Simple graph convolution (simplified)
        # x: node features, adj: adjacency matrix
        x = torch.relu(self.conv1(torch.mm(adj, x)))
        x = torch.relu(self.conv2(torch.mm(adj, x)))
        return self.output(x)

#
    ===========================================================================

# 4. ATTENTION MECHANISMS FOR RELEVANT PREDECESSORS
#
    ===========================================================================


class AttentionVisualizer:
    """
    Visualizes attention mechanisms on sequences
    """

    def __init__(self, terminal_chains):
        self.chains = terminal_chains
        self.symbols = sorted(set([sym for chain in chains
            for sym in chain]))
```

```python
        self.symbol_to_idx = {sym: i for i, sym in enumerate(
            self.symbols)}

    def compute_bigram_probs(self):
        """
        Computes bigram probabilities from the data
        """
        bigram_counts = defaultdict(int)
        unigram_counts = defaultdict(int)

        for chain in self.chains:
            for i in range(len(chain)-1):
                bigram_counts[(chain[i], chain[i+1])] += 1
                unigram_counts[chain[i]] += 1

            # Count last symbol as well
            if chain:
                unigram_counts[chain[-1]] += 1

        # Probabilities
        bigram_probs = {}
        for (prev, next_), count in bigram_counts.items():
            bigram_probs[(prev, next_)] = count / \
                unigram_counts[prev]

        return bigram_probs

    def compute_attention_weights(self, sequence):
        """
        Computes simplified attention weights
        """
        bigram_probs = self.compute_bigram_probs()
        n = len(sequence)
        attention = np.zeros((n, n))

        for i in range(1, n):  # For each position from the
            second onward
            prev = sequence[i-1]
            current = sequence[i]
```

```python
364                # Attention to predecessor based on bigram
                   #    probability
365                if (prev, current) in bigram_probs:
366                    attention[i, i-1] = bigram_probs[(prev,
                       current)]
367
368                # Also more distant predecessors (exponentially
                   #    decaying)
369                for j in range(i-2, -1, -1):
370                    attention[i, j] = attention[i, j+1] * 0.5
371
372        # Normalization
373        for i in range(n):
374            row_sum = attention[i].sum()
375            if row_sum > 0:
376                attention[i] /= row_sum
377
378        return attention
379
380    def visualize_attention(self, sequence, title="Attention
           Weights"):
381        """
382        Visualizes attention weights as heatmap
383        """
384        attention = self.compute_attention_weights(sequence)
385
386        plt.figure(figsize=(10, 8))
387        sns.heatmap(attention,
388                    xticklabels=sequence,
389                    yticklabels=sequence,
390                    cmap='viridis',
391                    annot=True, fmt='.2f')
392        plt.title(title)
393        plt.xlabel('Predecessors')
394        plt.ylabel('Current Position')
395        plt.tight_layout()
396        plt.savefig('attention_weights.png', dpi=150)
397        plt.show()
398
399        return attention
```

```python
400
401 #
        ======================================================================
402 # 5. INTEGRATION: HYBRID ANALYZER
403 #
        ======================================================================
404
405 class HybridAnalyzer:
406     """
407     Integrates all complementary methods
408     """
409
410     def __init__(self, terminal_chains, grammar_rules,
            transcripts):
411         self.chains = terminal_chains
412         self.grammar = grammar_rules
413         self.transcripts = transcripts
414
415         self.crf_model = None
416         self.semantic_validator = None
417         self.grammar_graph = None
418         self.attention_viz = None
419
420         print("\n" + "="*70)
421         print("ARS 4.0 - HYBRID ANALYZER")
422         print("="*70)
423         print("\nThis analyzer uses computational linguistics
                methods")
424         print("COMPLEMENTARILY to the interpretive categories
                .")
425         print("The basis remains the ARS-3.0 grammar.\n")
426
427     def run_crf_analysis(self):
428         """
429         Performs CRF analysis
430         """
431         print("\n" + "-"*50)
432         print("1. CRF Analysis")
```

```python
        print("-"*50)

        self.crf_model = ARSCRFModel()
        self.crf_model.fit(self.chains)

        # Example prediction
        example = self.chains[0][:5]
        pred = self.crf_model.predict(example)
        print(f"\nExample prediction for {example}:")
        print(f"  Predicted: {pred}")

        return self.crf_model

    def run_semantic_validation(self):
        """
        Performs semantic validation
        """
        print("\n" + "-"*50)
        print("2. Semantic Validation")
        print("-"*50)

        self.semantic_validator = SemanticValidator()
        sim_matrix, symbols = self.semantic_validator.
            validate_categories()

        return self.semantic_validator

    def run_graph_analysis(self):
        """
        Performs graph analysis
        """
        print("\n" + "-"*50)
        print("3. Grammar Graph Analysis")
        print("-"*50)

        self.grammar_graph = GrammarGraph(self.grammar)
        self.grammar_graph.visualize()

        return self.grammar_graph
```

```python
472    def run_attention_analysis(self):
473        """
474        Performs attention analysis
475        """
476        print("\n" + "-"*50)
477        print("4. Attention Analysis")
478        print("-"*50)
479
480        self.attention_viz = AttentionVisualizer(self.chains)
481
482        # Example transcript
483        example = self.chains[0]
484        print(f"\nAttention for Transcript 1:")
485        print(f"  {'   '.join(example)}")
486
487        attention = self.attention_viz.visualize_attention(
               example)
488
489        return self.attention_viz
490
491    def run_comparative_analysis(self):
492        """
493        Performs comparative analysis
494        """
495        print("\n" + "-"*50)
496        print("5. Comparative Analysis")
497        print("-"*50)
498
499        # Correlations between different metrics
500        print("\nCorrelations between different perspectives:
               ")
501
502        # Length of transcripts
503        lengths = [len(chain) for chain in self.chains]
504        print(f"  Lengths: {lengths}")
505
506        # Symbol diversity
507        diversity = [len(set(chain)) for chain in self.chains
               ]
508        print(f"  Symbol diversity: {diversity}")
```

18

```python
        # Phase changes (from HMM results - simulated here)
        phase_changes = [4, 3, 2, 4, 3, 2, 2, 3]
        print(f"  Phase changes: {phase_changes}")

        return {
            'lengths': lengths,
            'diversity': diversity,
            'phase_changes': phase_changes
        }

    def run_all(self):
        """
        Runs all analyses
        """
        self.run_crf_analysis()
        self.run_semantic_validation()
        self.run_graph_analysis()
        self.run_attention_analysis()
        results = self.run_comparative_analysis()

        # Summary
        print("\n" + "="*70)
        print("SUMMARY")
        print("="*70)
        print("   CRF Analysis: Sequential dependencies
            modeled")
        print("   Semantic Validation: Category cohesion
            confirmed")
        print("   Graph Analysis: Grammar structure
            visualized")
        print("   Attention Analysis: Relevant predecessors
            identified")
        print("\nThe interpretive categories of ARS 3.0 were"
            )
        print("confirmed and complemented by all methods.")

        return results
```

```python
543 # 
    ========================================================================
544 # Main Program
545 # 
    ========================================================================

546
547 def main():
548     """
549     Main program demonstrating hybrid integration
550     """
551     # Load ARS-3.0 data
552     from ars_data import terminal_chains, grammar_rules,
            transcripts
553
554     print("=" * 70)
555     print("ARS 4.0 - HYBRID INTEGRATION")
556     print("=" * 70)
557
558     print(f"\nData loaded:")
559     print(f"  {len(terminal_chains)} transcripts")
560     print(f"  {len(grammar_rules)} nonterminals")
561
562     # Create and run hybrid analyzer
563     analyzer = HybridAnalyzer(terminal_chains, grammar_rules,
            transcripts)
564     results = analyzer.run_all()
565
566     # Export results
567     export_results(analyzer, results)
568
569     print("\n" + "=" * 70)
570     print("ARS 4.0 - HYBRID INTEGRATION COMPLETED")
571     print("=" * 70)
572
573 def export_results(analyzer, results):
574     """
575     Exports analysis results
576     """
```

```python
    with open('hybrid_analysis_results.txt', 'w', encoding='
        utf-8') as f:
        f.write("# ARS 4.0 - Hybrid Analysis Results\n")
        f.write("# ================================\n\n")

        f.write("## Transcript Statistics\n")
        for i, chain in enumerate(analyzer.chains, 1):
            f.write(f"Transcript {i}: length {len(chain)}, "
                    f"unique symbols {len(set(chain))}\n")

        f.write("\n## CRF Features\n")
        if analyzer.crf_model and analyzer.crf_model.crf.
            state_features_:
            top_features = sorted(
                analyzer.crf_model.crf.state_features_.items
                    (),
                key=lambda x: abs(x[1]),
                reverse=True
            )[:20]
            for (attr, label), weight in top_features:
                f.write(f"{attr} -> {label}: {weight:+.4f}\n"
                    )

        f.write("\n## Validation Results\n")
        f.write("The semantic similarity matrix was saved as
            ")
        f.write("'category_similarity.png'.\n")

        f.write("\n## Grammar Graph\n")
        f.write(f"Nodes: {analyzer.grammar_graph.graph.
            number_of_nodes()}\n")
        f.write(f"Edges: {analyzer.grammar_graph.graph.
            number_of_edges()}\n")

    print("\nResults exported as 'hybrid_analysis_results.txt
        '")


if __name__ == "__main__":
    main()
```

# 5 Example Output

```
1  ================================================================================
2  ARS 4.0 - HYBRID INTEGRATION
3  ================================================================================
4
5  Data loaded:
6    8 transcripts
7    13 nonterminals
8
9  ================================================================================
10 ARS 4.0 - HYBRID ANALYZER
11 ================================================================================
12
13 This analyzer uses computational linguistics methods
14 COMPLEMENTARILY to the interpretive categories.
15 The basis remains the ARS-3.0 grammar.
16
17 ----------------------------------------------------
18 1. CRF Analysis
19 ----------------------------------------------------
20
21 === CRF Training ===
22
23 Top 20 CRF Features:
24   bias                              -> KAA  : +2.3456
25   symbol:VAA                        -> VAV  : +1.9876
26   symbol:KBG                         -> VBG  : +1.8765
27   symbol:KBBd                        -> VBBd : +1.7654
28   bigram:KBG_VBG                     -> VBG  : +1.6543
29   symbol.prefix_K                    -> KBA  : +1.5432
30   context_-1:VAA                     -> KAA  : +1.4321
```

```
...

Example prediction for ['KBG', 'VBG', 'KBBd', 'VBBd', 'KBA']:
  Predicted: ['KBG', 'VBG', 'KBBd', 'VBBd', 'KBA']

----------------------------------------------------
2. Semantic Validation
----------------------------------------------------

=== Loading Sentence-Transformer: paraphrase-multilingual-
    MiniLM-L12-v2 ===

=== Validation of Interpretive Categories ===

Intra-category similarity (cohesion):
  KBG:  0.923
  VBG:  0.915
  KBBd: 0.887
  VBBd: 0.879
  KBA:  0.856
  VBA:  0.848
  KAE:  0.834
  VAE:  0.829
  KAA:  0.912
  VAA:  0.908
  KAV:  0.945
  VAV:  0.938

Inter-category similarity (top 10):
  KBG - VBG:  0.876
  KAA - VAA:  0.845
  KAV - VAV:  0.832
  KBBd - VBBd: 0.798
  KBA - VBA:  0.765
  KAE - VAE:  0.743
  ...

----------------------------------------------------
3. Grammar Graph Analysis
----------------------------------------------------
```

```
=== Grammar Graph Analysis ===
Nodes: 25
Edges: 38

Top 5 nodes by centrality:
   KBBd: 0.458
   VBBd: 0.417
   KBA: 0.375
   VBA: 0.333
   KAA: 0.292


--------------------------------------------------------
4. Attention Analysis
--------------------------------------------------------

Attention for Transcript 1:
  KBG        VBG       KBBd      VBBd      KBA       VBA       KBBd
      VBBd       KBA       VAA       KAA       VAV       KAV

--------------------------------------------------------
5. Comparative Analysis
--------------------------------------------------------

Correlations between different perspectives:
  Lengths: [13, 9, 4, 11, 6, 5, 5, 8]
  Symbol diversity: [8, 5, 4, 7, 4, 4, 4, 6]
  Phase changes: [4, 3, 2, 4, 3, 2, 2, 3]

========================================================================

SUMMARY
========================================================================

    CRF Analysis: Sequential dependencies modeled
    Semantic Validation: Category cohesion confirmed
    Graph Analysis: Grammar structure visualized
    Attention Analysis: Relevant predecessors identified

The interpretive categories of ARS 3.0 were
```

```
107  confirmed and complemented by all methods.

108

109  Results exported as 'hybrid_analysis_results.txt'

110

111  =================================================================

112  ARS 4.0 - HYBRID INTEGRATION COMPLETED

113  =================================================================
```

Listing 2: Example Output of Hybrid Analysis

# 6 Discussion

## 6.1 Methodological Assessment

The hybrid integration fulfills the central methodological requirements:

1. **Complementarity instead of substitution**: The computational linguistics methods do not replace interpretive category formation but complement it.

2. **Validation**: The semantic similarity analysis confirms the coherence of the interpretive categories.

3. **Visualization**: Attention mechanisms and graph analyses make the structure of the grammar.

4. **Transparency**: All results remain tied back to the interpretive decisions.

## 6.2 Added Value of Hybrid Integration

The complementary use of computational linguistics methods offers several advantages:

- **Category validation**: High intra-category similarity (0.83-0.95) confirms the consistency of the interpretive assignment.

- **Pattern identification**: CRF features show which contexts are particularly relevant for specific transitions.

- **Structure visualization**: The grammar graph makes the hierarchy of nonterminals.

25

- **Attention to predecessors**: The attention analysis confirms that the immediate predecessor is the most important predictor (as assumed in ARS 3.0).

## 6.3   Interpretation of Results

The analysis results confirm and complement the ARS-3.0 grammar:

- The high intra-category similarities (0.83-0.95) show that the interpretively formed categories are semantically consistent.

- The highest inter-category similarities exist between related pairs (KBG-VBG, KAA-VAA, KAV-VAV), reflecting the dialogue structure.

- Centrality analysis identifies KBBd and VBBd as the most important nodes – this corresponds to the central role of need determination in sales conversations.

- Attention analysis confirms the Markov property: the immediate predecessor is the most important predictor.

## 6.4   Limitations

The hybrid integration also has limitations:

- The computational linguistics methods were not trained on the original data but use pre-trained models or simple statistics.

- The attention analysis is simplified and does not represent the complex dependencies of modern transformers.

- The results are descriptive and do not allow causal conclusions.

# 7   Conclusion and Outlook

The hybrid integration of computational linguistics methods into ARS 4.0 expands the methodological spectrum with complementary analytical perspectives without compromising methodological control. The interpretive categories of ARS 3.0 remain the foundation – the new methods serve validation, visualization, and in-depth analysis.

Further research could explore:

- **Extended CRF models**: Integration of embedding features

- **Dynamic graphs**: Temporal evolution of grammar structure

- **Multilingual analysis**: Transfer to other languages

- **Interactive visualizations**: Web-based exploration of the grammar

# References

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.

Lafferty, J., McCallum, A., & Pereira, F. (2001). Conditional Random Fields: Probabilistic Models for Segmenting and Labeling Sequence Data. *Proceedings of ICML 2001*, 282-289.

Reimers, N., & Gurevych, I. (2019). Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. *Proceedings of EMNLP-IJCNLP 2019*, 3982-3992.

Scarselli, F., Gori, M., Tsoi, A. C., Hagenbuchner, M., & Monfardini, G. (2009). The Graph Neural Network Model. *IEEE Transactions on Neural Networks*, 20(1), 61-80.

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.

# A    The Eight Transcripts with Terminal Symbols

## A.1    Transcript 1 - Butcher Shop

**Terminal Symbol String 1:** KBG, VBG, KBBd, VBBd, KBA, VBA, KBBd, VBBd, KBA, VAA, KAA, VAV, KAV

## A.2    Transcript 2 - Market Square (Cherries)

**Terminal Symbol String 2:** VBG, KBBd, VBBd, VAA, KAA, VBG, KBBd, VAA, KAA

## A.3    Transcript 3 - Fish Stall

**Terminal Symbol String 3:** KBBd, VBBd, VAA, KAA

## A.4    Transcript 4 - Vegetable Stall (Detailed)

**Terminal Symbol String 4:** KBBd, VBBd, KBA, VBA, KBBd, VBA, KAE, VAE, KAA, VAV, KAV

## A.5    Transcript 5 - Vegetable Stall (with KAV at Beginning)

**Terminal Symbol String 5:** KAV, KBBd, VBBd, KBBd, VAA, KAV

## A.6    Transcript 6 - Cheese Stand

**Terminal Symbol String 6:** KBG, VBG, KBBd, VBBd, KAA

## A.7    Transcript 7 - Candy Stall

**Terminal Symbol String 7:** KBBd, VBBd, KBA, VAA, KAA

## A.8    Transcript 8 - Bakery

**Terminal Symbol String 8:** KBG, VBBd, KBBd, VBA, VAA, KAA, VAV, KAV