

REAL-TIME HALLUCINATION DETECTION METHOD FOR STREAMING AI RESPONSES

PROVISIONAL PATENT APPLICATION

Inventor: Kinan Lemberg

Address: 270 Bolton Rd, Koah, 4881, Australia

Filing Date: June 3, 2025

FIELD OF THE INVENTION

This invention relates to real-time hallucination detection systems for artificial intelligence applications, and more specifically to methods for detecting AI hallucinations in streaming responses before completion through pattern recognition, statistical anomaly detection, and consistency checking algorithms.

BACKGROUND OF THE INVENTION

Current AI systems generate responses token-by-token, but hallucination detection typically occurs only after complete response generation. This post-hoc approach allows hallucinated content to reach users before detection, creating risks in time-sensitive applications. Existing systems lack mechanisms to detect hallucinations during the streaming generation process.

Current limitations include:

- No real-time detection during token-by-token generation
- Lack of pattern recognition in partial AI outputs
- Absence of statistical anomaly detection for streaming responses
- No consistency checking mechanisms during response generation
- Limited early warning systems for emerging hallucinations

There exists a need for a comprehensive real-time hallucination detection system that can identify potential hallucinations during the streaming generation process and intervene before problematic content reaches users.

SUMMARY OF THE INVENTION

The present invention provides a real-time hallucination detection method that monitors streaming AI responses token-by-token, implementing pattern recognition and anomaly detection to identify hallucinations before response completion.

The invention comprises:

1. Token Stream Analysis Engine - Real-time analysis of streaming tokens for hallucination patterns and anomalies.
2. Pattern Recognition System - Detection of known hallucination patterns in partial outputs using machine learning models.
3. Statistical Anomaly Detector - Mathematical detection of statistical deviations indicating potential hallucinations.
4. Consistency Checking Module - Real-time verification of internal consistency within streaming responses.
5. Early Warning and Intervention System - Automated alerts and intervention mechanisms when hallucinations are detected.

The system provides significant advantages by detecting hallucinations during generation rather than after completion, enabling immediate intervention and preventing delivery of hallucinated content.

DETAILED DESCRIPTION OF THE INVENTION

System Architecture

The Real-Time Hallucination Detection Method operates as a streaming analysis layer that monitors AI responses token-by-token during generation to identify emerging hallucinations.

1. Token Stream Analysis Engine

The Analysis Engine implements sophisticated real-time processing of streaming tokens to detect hallucination indicators.

Streaming Token Processing:

...

```
Token_Stream_Buffer = {  
    sliding_window: Last_N_Tokens_Generated,  
    window_sizes: [10, 50, 100, 500],  
    update_frequency: Every_Single-Token,  
    context_preservation: Full_Conversation_History,  
    processing_mode: Parallel_Multi_Window_Analysis  
}
```

```
Token_Features_Extraction = {  
    linguistic_features: {  
        perplexity_score: Token_Sequence_Perplexity,  
        entropy_measure: Information_Entropy_Calculation,  
        coherence_score: Semantic_Coherence_Metric,  
        grammatical_validity: Syntax_Tree_Consistency  
    },  
    statistical_features: {  
        token_probability: Model_Assigned_Probability,  
        sequence_likelihood: Joint_Probability_Score,  
        distribution_deviation: KL_Divergence_from_Expected,  
        repetition_metrics: N_Gram_Repetition_Scores  
    },  
    semantic_features: {  
        embedding_distance: Semantic_Vector_Distances,  
        topic_drift: Topic_Coherence_Degradation,  
        factual_grounding: Knowledge_Base_Alignment,  
        entity_consistency: Named_Entity_Tracking  
    }  
}
```

```
}  
}  
...
```

Real-Time Feature Computation:

```
...
```

```
Incremental_Perplexity = {  
  current_perplexity:  $-\log(P(\text{token}_t \mid \text{tokens}_{1\dots t-1}))$ ,  
  running_average:  $(1/t) \times \sum(\text{perplexity}_i)$ ,  
  windowed_average: Average_Over_Recent_Window,  
  trend_analysis: Rate_of_Perplexity_Change,  
  spike_detection: Sudden_Perplexity_Increases  
}
```

```
Streaming_Entropy_Calculation = {  
  token_entropy:  $-\sum(p_i \times \log(p_i))$ ,  
  sequence_entropy: Joint_Entropy_Estimation,  
  conditional_entropy:  $H(X_t \mid X_{1\dots t-1})$ ,  
  entropy_rate: Limit_of_Conditional_Entropy,  
  anomaly_threshold: 2_Standard_Deviations_Above_Mean  
}
```

```
Coherence_Tracking = {  
  local_coherence: Adjacent-Token_Relationships,  
  global_coherence: Full_Response_Consistency,  
  topic_coherence: Latent_Topic_Stability,  
  entity_coherence: Entity_Reference_Consistency,  
  temporal_coherence: Time_Reference_Validity  
}  
...
```

2. Pattern Recognition System

The Pattern Recognition System identifies known hallucination patterns using trained machine learning models on partial outputs.

Hallucination Pattern Database:

...

Known_Hallucination_Patterns = {

 fabrication_patterns: {

 fake_citations: "According to [Non-existent Source]",

 invented_statistics: "[Specific Number]% of [Unverified Claim]",

 false_quotes: "'[Fabricated Quote]' - [Real Person]",

 fictional_events: "In [Year], [Never Happened Event]"

 },

 confidence_without_basis: {

 unsupported_certainty: "definitely", "certainly", "always" without evidence,

 false_expertise: "As an expert in [Unqualified Domain]",

 overreach_patterns: "The only way", "Must always", "Never fails"

 },

 logical_inconsistencies: {

 self_contradiction: Statement_A followed by NOT(Statement_A),

 circular_reasoning: A_implies_B_implies_A patterns,

 non_sequiturs: Conclusion_not_following_from_Premises,

 category_errors: Mixing_incompatible_concept_types

 }

}

...

Neural Pattern Detection:

...

Pattern_Detection_LSTM = {

```
input_layer: Token_Embeddings_Sequence,
lstm_layers: [
    LSTM(256, return_sequences=True),
    Dropout(0.3),
    LSTM(128, return_sequences=True),
    Attention_Mechanism(),
    LSTM(64)
],
output_layer: {
    hallucination_probability: Sigmoid_Activation,
    pattern_type: Softmax_Multi_Class,
    confidence_score: Linear_Activation
},
training_data: Labeled_Hallucination_Examples,
online_learning: Continuous_Model_Updates
}
```

```
Transformer_Based_Detection = {
    architecture: BERT_Style_Encoder,
    attention_mechanism: Multi_Head_Self_Attention,
    hallucination_heads: Specialized_Attention_for_Patterns,
    input_processing: {
        token_embeddings: Pretrained_Embeddings,
        position_encodings: Absolute_Position_Information,
        segment_embeddings: Response_vs_Context_Segments
    },
    detection_strategy: {
        early_layers: Surface_Pattern_Detection,
        middle_layers: Semantic_Inconsistency_Detection,
        late_layers: Complex_Hallucination_Patterns
    }
}
```

```
}  
...
```

Pattern Matching Optimization:

```
...
```

```
Efficient_Pattern_Search = {  
    trie_structure: Prefix_Tree_for_Known_Patterns,  
    boyer_moore: Efficient_String_Matching,  
    aho_corasick: Multiple_Pattern_Simultaneous_Search,  
    regex_engine: Compiled_Regular_Expressions,  
    fuzzy_matching: Approximate_Pattern_Detection  
}
```

```
Incremental_Matching = {  
    state_preservation: Maintain_Matching_State_Between_Tokens,  
    partial_match_tracking: Score_Incomplete_Pattern_Matches,  
    branch_prediction: Anticipate_Likely_Completions,  
    early_termination: Stop_on_High_Confidence_Detection  
}
```

```
...
```

3. Statistical Anomaly Detector

The Anomaly Detector implements mathematical methods for identifying statistical deviations indicating hallucinations.

Statistical Deviation Metrics:

```
...
```

```
Probability_Distribution_Analysis = {  
    expected_distribution: Model_Training_Distribution,  
    observed_distribution: Current-Token_Distribution,
```

```

divergence_metrics: {
  kl_divergence: KL(Observed || Expected),
  js_divergence: JS(Observed, Expected),
  wasserstein_distance: Earth_Movers_Distance,
  total_variation: TV(Observed, Expected)
},
threshold_determination: {
  static_threshold: Pre_Computed_Critical_Values,
  dynamic_threshold: Adaptive_Based_on_Context,
  percentile_based: 99th_Percentile_of_Training,
  bayesian_threshold: Posterior_Probability_Based
}
}

```

```

Sequence_Likelihood_Analysis = {
  log_likelihood:  $\Sigma(\log(P(\text{token}_i | \text{context}_i)))$ ,
  normalized_likelihood: Log_Likelihood / Sequence_Length,
  likelihood_gradient: Rate_of_Likelihood_Change,
  sudden_drops: Detect_Sharp_Likelihood_Decreases,
  recovery_pattern: Monitor_Likelihood_Recovery
}
...

```

Multivariate Anomaly Detection:

...

```

Feature_Vector = [
  perplexity,
  entropy,
  likelihood,
  coherence_score,
  repetition_count,

```

```
semantic_drift,  
grammatical_score  
]
```

```
Anomaly_Detection_Methods = {
```

```
isolation_forest: {  
    n_estimators: 100,  
    contamination: 0.1,  
    random_state: 42,  
    online_update: True
```

```
},
```

```
one_class_svm: {  
    kernel: 'rbf',  
    gamma: 'auto',  
    nu: 0.05,  
    incremental: True
```

```
},
```

```
autoencoder: {  
    architecture: [128, 64, 32, 64, 128],  
    reconstruction_error_threshold: 2_sigma,  
    adaptation_rate: 0.01
```

```
},
```

```
statistical_process_control: {  
    control_charts: EWMA_and_CUSUM,  
    violation_rules: Western_Electric_Rules,  
    multivariate_extension: Hotelling_T_Squared
```

```
}
```

```
}
```

```
'''
```

Temporal Anomaly Patterns:

...

Time_Series_Analysis = {

 trend_detection: {

 moving_average: Exponential_Weighted_MA,

 trend_strength: Mann_Kendall_Test,

 change_points: PELT_Algorithm,

 seasonality: STL_Decomposition

 },

 anomaly_types: {

 point_anomalies: Single-Token-Outliers,

 contextual_anomalies: Normal-Token-Wrong-Context,

 collective_anomalies: Abnormal-Token-Sequences,

 pattern_anomalies: Unusual-Generation-Patterns

 }

}

...

4. Consistency Checking Module

The Consistency Module performs real-time verification of internal consistency within the streaming response.

Multi-Level Consistency Verification:

...

Factual_Consistency = {

 entity_tracking: {

 named_entities: Track_All_Mentioned_Entities,

 entity_properties: Maintain_Consistent_Attributes,

 entity_relationships: Verify_Relationship_Consistency,

 entity_timeline: Check_Temporal_Consistency

 },

```
numerical_consistency: {
  number_tracking: Store_All_Numerical_Values,
  calculation_verification: Check_Mathematical_Consistency,
  statistical_validation: Verify_Statistical_Claims,
  unit_consistency: Ensure_Unit_Compatibility
},
reference_consistency: {
  citation_tracking: Monitor_Source_References,
  quote_verification: Check_Quote_Consistency,
  definition_stability: Maintain_Term_Definitions,
  cross_reference_validation: Verify_Internal_References
}
}
```

```
Logical_Consistency = {
  proposition_tracking: {
    stated_facts: Extract_Factual_Propositions,
    logical_relationships: Map_Implication_Structure,
    contradiction_detection: Find_Logical_Conflicts,
    inference_validation: Check_Deductive_Validity
  },
  argument_structure: {
    premise_tracking: Identify_Argument_Premises,
    conclusion_mapping: Link_Conclusions_to_Premises,
    fallacy_detection: Identify_Logical_Fallacies,
    coherence_scoring: Measure_Argument_Coherence
  }
}
```

```
Semantic_Consistency = {
  concept_stability: {
```

```
    definition_tracking: Monitor_Concept_Definitions,
    usage_consistency: Verify_Consistent_Usage,
    semantic_drift: Detect_Meaning_Changes,
    context_appropriateness: Check_Contextual_Fit
  },
  topic_coherence: {
    topic_modeling: Dynamic_Topic_Extraction,
    topic_stability: Monitor_Topic_Persistence,
    relevance_scoring: Measure_Topic_Relevance,
    drift_detection: Identify_Off_Topic_Shifts
  }
}
...

Real-Time Consistency Scoring:
...

Consistency_Score_Computation = {
  incremental_calculation: Update_Score_Per_Token,
  weighted_components: {
    factual_weight: 0.4,
    logical_weight: 0.3,
    semantic_weight: 0.3
  },
  violation_penalties: {
    minor_inconsistency: -0.1,
    major_inconsistency: -0.3,
    direct_contradiction: -0.5,
    impossible_claim: -1.0
  },
  recovery_mechanism: Allow_Score_Recovery_with_Corrections
}
```

...

5. Early Warning and Intervention System

The Warning System provides immediate alerts and intervention when hallucinations are detected during generation.

Detection Threshold Framework:

...

```
Multi_Tier_Alert_System = {
  green_zone: {
    threshold: All_Metrics_Within_Normal_Range,
    action: Continue_Normal_Generation,
    monitoring: Standard_Tracking
  },
  yellow_zone: {
    threshold: 1-2_Metrics_Showing_Anomalies,
    action: Increase_Monitoring_Sensitivity,
    intervention: Add_Uncertainty_Markers
  },
  orange_zone: {
    threshold: Multiple_Metrics_Anomalous,
    action: Slow_Generation_Rate,
    intervention: Insert_Verification_Prompts
  },
  red_zone: {
    threshold: Critical_Hallucination_Detected,
    action: Pause_Generation,
    intervention: Trigger_Regeneration
  }
}
```

...

Intervention Strategies:

...

Automated_Interventions = {

 uncertainty_injection: {

 trigger: Moderate_Hallucination_Risk,

 action: Add_Phrases_Like_"possibly", "might", "unclear",

 implementation: Modify-Token-Probabilities

 },

 fact_checking_prompt: {

 trigger: Factual_Claim_Anomaly,

 action: Insert_"Let_me_verify_that",

 implementation: Inject_Verification_Tokens

 },

 regeneration_trigger: {

 trigger: High_Hallucination_Confidence,

 action: Stop_and_Regenerate_from_Last_Safe_Point,

 implementation: Rollback_to_Clean_State

 },

 model_switching: {

 trigger: Persistent_Hallucination_Pattern,

 action: Switch_to_More_Conservative_Model,

 implementation: Dynamic_Model_Substitution

 }

}

...

Real-Time Dashboard:

...

Monitoring_Interface = {

```

live_metrics: {
  perplexity_gauge: Real_Time_Perplexity_Display,
  anomaly_indicators: Multi_Metric_Status_Lights,
  consistency_tracker: Consistency_Score_Timeline,
  pattern_alerts: Detected_Pattern_Notifications
},
intervention_controls: {
  manual_pause: Human_Override_Button,
  threshold_adjustment: Dynamic_Sensitivity_Control,
  regeneration_trigger: Force_Regeneration_Option,
  model_selection: Choose_Alternative_Model
},
audit_trail: {
  detection_log: All_Anomaly_Detections,
  intervention_log: All_Automatic_Interventions,
  effectiveness_metrics: Intervention_Success_Rates,
  false_positive_tracking: Monitor_Over_Sensitivity
}
}
...

```

System Integration and Performance

Real-Time Performance Requirements:

- Token Processing Latency: < 5ms per token
- Pattern Detection Speed: < 10ms for pattern matching
- Anomaly Detection: < 15ms for statistical analysis
- Consistency Checking: < 20ms for comprehensive check
- Total Stream Delay: < 50ms maximum added latency

Optimization Strategies:

- GPU Acceleration: Parallel processing of detection algorithms
- Caching Mechanisms: Reuse computed features across windows
- Incremental Computation: Update metrics rather than recompute
- Pipeline Parallelism: Concurrent execution of detection modules

ADVANTAGES OVER PRIOR ART

The present invention provides significant advantages over existing hallucination detection approaches:

1. Real-Time Detection: Unlike post-generation analysis, the invention detects hallucinations during streaming generation.
2. Early Intervention: The system can intervene before hallucinated content completes rather than after delivery.
3. Pattern-Based Detection: Recognizes known hallucination patterns in partial outputs rather than requiring complete responses.
4. Statistical Rigor: Implements multiple statistical methods for robust anomaly detection rather than simple thresholds.
5. Consistency Verification: Maintains real-time consistency checking rather than post-hoc analysis.
6. Adaptive Intervention: Provides graduated responses based on hallucination severity rather than binary decisions.

CLAIMS

Claim 1: A real-time hallucination detection method for streaming AI responses comprising:

- a token stream analysis engine configured to analyze streaming tokens for hallucination indicators;
- a pattern recognition system detecting known hallucination patterns in partial outputs;
- a statistical anomaly detector implementing mathematical detection of deviations;

- a consistency checking module performing real-time internal consistency verification; and
- an early warning and intervention system providing automated alerts and interventions.

Claim 2: The method of claim 1, wherein the token stream analysis engine computes incremental perplexity, entropy, and coherence metrics for sliding windows of tokens.

Claim 3: The method of claim 1, wherein the pattern recognition system uses LSTM and transformer models trained on labeled hallucination examples.

Claim 4: The method of claim 1, wherein the statistical anomaly detector implements isolation forests, one-class SVM, and statistical process control for multivariate anomaly detection.

Claim 5: The method of claim 1, wherein the consistency checking module tracks factual, logical, and semantic consistency across entity mentions, numerical values, and propositions.

Claim 6: The method of claim 1, wherein the early warning system implements multi-tier alerts with graduated intervention strategies including uncertainty injection and regeneration triggers.

Claim 7: A method for detecting AI hallucinations during response generation comprising:

- monitoring token-by-token generation in real-time;
- extracting linguistic, statistical, and semantic features from partial outputs;
- comparing features against known hallucination patterns;
- detecting statistical anomalies indicating potential hallucinations;
- verifying internal consistency of streaming content; and
- triggering interventions when hallucination indicators exceed thresholds.

Claim 8: The method of claim 7, further comprising implementing sliding window analysis with multiple window sizes for comprehensive pattern detection.

Claim 9: The method of claim 7, wherein anomaly detection includes probability distribution analysis, sequence likelihood tracking, and temporal pattern recognition.

Claim 10: The method of claim 7, wherein interventions include slowing generation rate, adding uncertainty markers, prompting for verification, and triggering regeneration from safe points.

ABSTRACT

A real-time hallucination detection method monitors streaming AI responses token-by-token to identify hallucinations before completion. The system comprises: (1) token stream analysis computing incremental metrics, (2) pattern recognition detecting known hallucination patterns, (3) statistical anomaly detection using multivariate methods, (4) consistency checking for factual and logical coherence, and (5) early warning system with graduated interventions. The method enables detection and intervention during generation rather than post-completion, providing advantages through real-time monitoring, early intervention, pattern-based detection, statistical rigor, and adaptive response strategies.

END OF PATENT SPECIFICATION