

MULTI-MODEL CONSENSUS VALIDATION SYSTEM FOR AI-GENERATED 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 consensus-based validation systems for artificial intelligence applications, and more specifically to multi-model agreement protocols that require consensus among multiple large language models before validating AI-generated responses as accurate and appropriate for delivery.

BACKGROUND OF THE INVENTION

Current AI validation systems typically rely on single-model confidence scores or simple threshold-based validation. This creates vulnerability to model-specific biases, hallucinations, and errors that could be detected through multi-model consensus validation. Existing approaches lack systematic methods for requiring agreement among multiple AI models before accepting responses as valid.

Current limitations include:

- No mandatory multi-model consensus requirements for AI validation
- Lack of weighted voting systems based on model expertise
- Absence of systematic dispute resolution when models disagree
- No mathematical framework for consensus scoring across diverse models
- Limited integration of model-specific strengths in validation decisions

There exists a need for a comprehensive consensus validation system that requires agreement among multiple AI models with different architectures and training to ensure response accuracy before delivery.

SUMMARY OF THE INVENTION

The present invention provides a multi-model consensus validation system that requires agreement among multiple large language models before AI-generated responses are validated for delivery, implementing weighted voting and intelligent dispute resolution mechanisms.

The invention comprises:

1. Multi-Model Orchestration Engine - Automated submission of responses to multiple AI models for independent validation assessment.
2. Weighted Consensus Voting System - Mathematical voting framework that weights model opinions based on domain expertise and historical accuracy.
3. Intelligent Dispute Resolution Protocol - Automated resolution mechanisms when models disagree, including tie-breaking and escalation procedures.
4. Model Expertise Profiling System - Dynamic assessment of each model's strengths and weaknesses for optimal consensus weighting.
5. Real-Time Consensus Scoring - Mathematical consensus calculation providing instant validation decisions based on multi-model agreement.

The system provides significant advantages by requiring multi-model consensus, thereby reducing single-model vulnerabilities and improving validation accuracy through diverse AI perspectives.

DETAILED DESCRIPTION OF THE INVENTION

System Architecture

The Multi-Model Consensus Validation System operates as an orchestration layer that submits AI responses to multiple independent models and requires consensus agreement before validation approval.

1. Multi-Model Orchestration Engine

The Orchestration Engine manages parallel submission of responses to multiple AI models and coordinates their validation assessments.

Model Selection Framework:

...

Model_Portfolio = {

primary_models: [GPT-4, Claude-3, Gemini-Pro, LLaMA-3],
specialized_models: [Medical-LLM, Legal-LLM, Finance-LLM],
verification_models: [Fact-Check-LLM, Logic-LLM, Safety-LLM],
minimum_models: 3,
maximum_models: 10

}

Model_Selection_Algorithm = {

domain_matching: Select_Models_By_Domain_Expertise(Query_Domain),
diversity_requirement: Ensure_Architectural_Diversity(Selected_Models),
availability_check: Verify_Model_Availability_and_Latency,
cost_optimization: Balance_Accuracy_vs_API_Costs

}

Parallel_Submission_Protocol = {

query_transformation: Adapt_Query_for_Each_Model_API,
simultaneous_submission: Async_Parallel_API_Calls,
timeout_management: Maximum_Response_Time_Per_Model,
fallback_handling: Backup_Model_Substitution_on_Failure

}

...

Response Collection and Normalization:

...

```
Response_Normalization = {  
  format_standardization: Convert_All_Responses_to_Common_Format,  
  confidence_extraction: Extract_Model_Confidence_Scores,  
  reasoning_extraction: Parse_Model_Reasoning_and_Explanations,  
  metadata_collection: Gather_Model_Version_and_Parameters  
}
```

```
Normalized_Response = {  
  model_id: Unique_Model_Identifier,  
  response_content: Standardized_Response_Text,  
  confidence_score: Model_Reported_Confidence,  
  validation_assessment: Model_Validation_Opinion,  
  reasoning: Model_Provided_Reasoning,  
  processing_time: Response_Generation_Duration  
}  
...
```

2. Weighted Consensus Voting System

The Voting System implements sophisticated mathematical frameworks for weighted consensus calculation based on model expertise and historical performance.

Mathematical Voting Framework:

...

$$\text{Consensus_Score} = \frac{\sum(\text{Model_Vote}_i \times \text{Model_Weight}_i \times \text{Expertise_Factor}_i)}{\sum(\text{Model_Weight}_i)}$$

where:

Model_Vote_i = {

1.0: "STRONGLY_AGREE",

0.75: "AGREE",

0.5: "NEUTRAL",

```
0.25: "DISAGREE",
0.0: "STRONGLY_DISAGREE"
}
```

$Model_Weight_i = Historical_Accuracy_i \times Domain_Expertise_i \times Reliability_Factor_i$

```
Expertise_Factor_i = {
  domain_match_score: Relevance_to_Current_Query_Domain,
  historical_performance: Past_Accuracy_in_Similar_Queries,
  specialization_bonus: Additional_Weight_for_Domain_Experts
}
...
```

Dynamic Weight Calculation:

...

$Historical_Accuracy_i = Correct_Validations_i / Total_Validations_i$

```
Domain_Expertise_i = {
  general_knowledge: Base_Expertise_Score,
  medical_expertise: Specialized_Medical_Training_Score,
  legal_expertise: Legal_Domain_Knowledge_Score,
  technical_expertise: Technical_Accuracy_Score,
  financial_expertise: Financial_Domain_Accuracy_Score
}
```

$Reliability_Factor_i = Uptime_i \times Consistency_i \times Version_Stability_i$

```
Weight_Adjustment = {
  recency_bias: Recent_Performance_Weighted_Higher,
  outlier_handling: Extreme_Votes_Reduced_Weight,
  confidence_correlation: Higher_Weight_for_High_Confidence_Accurate_Votes
}
```

```
}  
...
```

Consensus Thresholds:

```
...
```

```
Validation_Decision = {  
    if Consensus_Score ≥ 0.9: "STRONG_CONSENSUS_PASS",  
    elif Consensus_Score ≥ 0.75: "CONSENSUS_PASS",  
    elif Consensus_Score ≥ 0.6: "WEAK_CONSENSUS_PASS",  
    elif Consensus_Score ≥ 0.4: "NO_CONSENSUS",  
    else: "CONSENSUS_FAIL"  
}
```

```
Minimum_Agreement_Requirement = {  
    critical_domains: 0.9, # Medical, Legal, Financial  
    high_risk_domains: 0.8, # Business Strategy, Technical  
    standard_domains: 0.75, # General Knowledge  
    low_risk_domains: 0.6 # Entertainment, Casual  
}
```

```
...
```

3. Intelligent Dispute Resolution Protocol

The Dispute Resolution Protocol implements sophisticated mechanisms for handling disagreements among models.

Disagreement Detection:

```
...
```

Disagreement_Metric = Standard_Deviation(Model_Votes) / Mean(Model_Votes)

Disagreement_Classification = {

```
if Disagreement_Metric < 0.1: "STRONG_AGREEMENT",
elif Disagreement_Metric < 0.3: "MINOR_DISAGREEMENT",
elif Disagreement_Metric < 0.5: "MODERATE_DISAGREEMENT",
else: "MAJOR_DISAGREEMENT"
}
```

```
Dispute_Analysis = {
    faction_identification: Cluster_Models_by_Vote_Similarity,
    reasoning_comparison: Analyze_Different_Reasoning_Approaches,
    fact_verification: Cross_Check_Disputed_Facts,
    assumption_identification: Extract_Underlying_Assumptions
}
```

...

Resolution Strategies:

...

```
Tie_Breaking_Protocol = {
    expert_model_preference: Prefer_Domain_Expert_Models,
    confidence_weighted: Higher_Weight_to_Confident_Models,
    historical_accuracy: Prefer_Historically_Accurate_Models,
    human_escalation: Flag_for_Human_Review_if_Unresolved
}
```

```
Dispute_Resolution_Strategies = {
    fact_checking_arbitration: {
        trigger: Factual_Disagreement_Detected,
        action: Submit_to_Specialized_Fact_Checking_Model,
        weight: Fact_Checker_Gets_2x_Vote_Weight
    },
    domain_expert_arbitration: {
        trigger: Domain_Specific_Disagreement,
```

```

    action: Consult_Domain_Specialist_Model,
    weight: Specialist_Gets_3x_Vote_Weight
  },
  consensus_building: {
    trigger: Close_Vote_Split,
    action: Request_Models_to_Review_Other_Opinions,
    iteration: Allow_Up_to_3_Consensus_Building_Rounds
  }
}
'''

```

4. Model Expertise Profiling System

The Profiling System dynamically assesses and updates each model's expertise profile based on performance.

Expertise Assessment Framework:

'''

```

Model_Expertise_Profile = {
  domain_strengths: {
    medical: Medical_Accuracy_Score,
    legal: Legal_Accuracy_Score,
    technical: Technical_Accuracy_Score,
    financial: Financial_Accuracy_Score,
    general: General_Knowledge_Score
  },
  capability_metrics: {
    factual_accuracy: Fact_Checking_Performance,
    logical_reasoning: Logic_Problem_Performance,
    creative_tasks: Creative_Response_Quality,
    safety_assessment: Risk_Detection_Accuracy
  }
}
'''

```

```
},
performance_trends: {
  improvement_rate: Performance_Change_Over_Time,
  consistency_score: Response_Consistency_Metric,
  reliability_index: Uptime_and_Error_Rate
}
}
```

```
Dynamic_Profile_Update = {
  continuous_learning: Update_After_Each_Validation,
  performance_decay: Reduce_Scores_for_Prolonged_Errors,
  breakthrough_detection: Bonus_for_Catching_Missed_Errors,
  comparative_analysis: Relative_Performance_vs_Other_Models
}
...
```

```
Expertise-Based_Model_Selection:
...
```

```
Optimal_Model_Selection = {
  primary_criterion: Match_Query_Domain_to_Model_Expertise,
  diversity_criterion: Include_Models_with_Different_Strengths,
  performance_criterion: Minimum_Historical_Accuracy_Threshold,
  cost_criterion: Balance_Expertise_with_API_Costs
}
```

```
Model_Team_Composition = {
  lead_expert: Highest_Domain_Expertise_Model,
  validators: 2-3_High_Accuracy_General_Models,
  specialists: 1-2_Domain_Specific_Models,
  safety_checker: Dedicated_Safety_Validation_Model
}
```

...

5. Real-Time Consensus Scoring

The Consensus Scoring system provides instant validation decisions through efficient parallel processing and mathematical consensus calculation.

Real-Time Processing Pipeline:

...

```
Parallel_Processing = {  
    request_distribution: Simultaneous_Model_Queries,  
    response_collection: Async_Response_Gathering,  
    incremental_scoring: Update_Consensus_as_Responses_Arrive,  
    early_termination: Stop_if_Strong_Consensus_Reached_Early  
}
```

```
Consensus_Calculation_Optimization = {  
    vector_computation: SIMD_Optimized_Vote_Calculations,  
    cache_utilization: Previous_Consensus_Pattern_Caching,  
    threshold_shortcuts: Early_Exit_on_Unanimous_Agreement,  
    parallel_aggregation: Distributed_Vote_Aggregation  
}
```

```
Real_Time_Metrics = {  
    consensus_latency: Time_to_Consensus_Decision,  
    model_response_times: Individual_Model_Latencies,  
    disagreement_resolution_time: Dispute_Resolution_Duration,  
    total_validation_time: End_to_End_Validation_Latency  
}
```

...

Streaming Consensus Updates:

...

```
Streaming_Consensus = {  
    initial_estimate: First_2_Model_Preliminary_Consensus,  
    progressive_refinement: Update_with_Each_New_Model_Response,  
    confidence_intervals: Statistical_Confidence_Bounds,  
    stability_detection: Consensus_Stabilization_Monitoring  
}
```

```
Client_Notification = {  
    preliminary_result: Early_Consensus_Indication,  
    final_result: Complete_Consensus_Decision,  
    disagreement_alert: Notification_of_Significant_Disputes,  
    resolution_status: Dispute_Resolution_Progress_Updates  
}
```

...

System Integration and Performance

Consensus System Performance Requirements:

- Model Query Latency: < 100ms parallel distribution
- Consensus Calculation: < 50ms for vote aggregation
- Dispute Resolution: < 2 seconds for standard disputes
- Total Validation Time: < 5 seconds for 5-model consensus

Scalability Architecture:

- Horizontal Scaling: Distributed model query handling
- Load Balancing: Intelligent distribution across API endpoints
- Caching Strategy: Consensus pattern caching for similar queries
- Failover Handling: Automatic model substitution on failures

ADVANTAGES OVER PRIOR ART

The present invention provides significant advantages over existing validation approaches:

1. **Multi-Model Validation:** Unlike single-model systems, the invention requires consensus among multiple independent AI models.
2. **Weighted Expertise:** The system weights model opinions based on domain expertise rather than treating all models equally.
3. **Intelligent Dispute Resolution:** Automated mechanisms resolve model disagreements rather than simple majority voting.
4. **Dynamic Expertise Profiling:** Continuous assessment of model strengths rather than static model selection.
5. **Real-Time Consensus:** Instant consensus calculation with streaming updates rather than batch processing.
6. **Reduced Single-Point Failure:** Multiple model requirement eliminates single-model vulnerabilities.

CLAIMS

Claim 1: A multi-model consensus validation system comprising:

- a multi-model orchestration engine configured to submit AI responses to multiple independent models for validation assessment;
- a weighted consensus voting system implementing mathematical voting frameworks based on model expertise and performance;
- an intelligent dispute resolution protocol configured to resolve disagreements among models through specialized arbitration;
- a model expertise profiling system configured to dynamically assess and update each model's domain strengths; and
- a real-time consensus scoring system providing instant validation decisions through parallel processing.

Claim 2: The system of claim 1, wherein the multi-model orchestration engine implements parallel submission protocols with model diversity requirements and automatic failover handling.

Claim 3: The system of claim 1, wherein the weighted consensus voting system calculates consensus scores using model weights based on historical accuracy, domain expertise, and reliability factors.

Claim 4: The system of claim 1, wherein the intelligent dispute resolution protocol implements tie-breaking mechanisms, fact-checking arbitration, and consensus-building iterations.

Claim 5: The system of claim 1, wherein the model expertise profiling system maintains dynamic profiles of domain strengths, capability metrics, and performance trends for optimal model selection.

Claim 6: The system of claim 1, wherein the real-time consensus scoring provides streaming consensus updates with progressive refinement as model responses arrive.

Claim 7: A method for multi-model consensus validation comprising:

- submitting AI-generated responses to multiple independent models in parallel;
- collecting and normalizing validation assessments from each model;
- calculating weighted consensus scores based on model expertise and performance;
- resolving disputes when models disagree through intelligent arbitration; and
- providing real-time consensus-based validation decisions.

Claim 8: The method of claim 7, further comprising dynamically selecting models based on query domain and required expertise diversity.

Claim 9: The method of claim 7, wherein dispute resolution comprises faction identification, reasoning comparison, and specialized arbitration for fact verification.

Claim 10: The method of claim 7, wherein consensus calculation implements early termination on strong agreement and streaming updates for real-time feedback.

ABSTRACT

A multi-model consensus validation system requires agreement among multiple large language models before validating AI-generated responses. The system comprises: (1) multi-model orchestration for parallel validation submission, (2) weighted consensus voting based on model expertise, (3) intelligent dispute resolution with specialized arbitration, (4) dynamic model expertise profiling, and (5) real-time consensus scoring with streaming updates. The system ensures validation accuracy through multi-model agreement, providing advantages over single-model systems through consensus requirements, weighted expertise, intelligent dispute resolution, and reduced single-point vulnerabilities.

END OF PATENT SPECIFICATION