

Strategic Doctrine Language Models (sdLM): A Comprehensive Framework for Automated Military Planning and Doctrinal Analysis

Olaf Yunus Laitinen Imanov, Taner Yilmaz, and Derya Umut Kulali

Abstract—Military strategic planning is information-intensive and time-consuming, often requiring months of staff work and producing thousands of pages of supporting documents. This paper proposes a domain-specialized learning system that combines a large-scale strategic planning model with a faster wargaming and simulation model, both built on transformer architectures and trained to enforce doctrinal consistency, temporal reasoning, and calibrated uncertainty. Across 127 historical and synthetic scenarios and an expert panel of 47 senior strategists, the system achieves a mean strategic quality score of 8.42/10 and improves geopolitical prediction accuracy to 73.2

Index Terms—Natural language processing, Decision support systems, Transformers, Multi-document reasoning, Uncertainty calibration

I. INTRODUCTION

STRATEGIC military planning represents one of the most cognitively demanding tasks in defense operations, requiring synthesis of geopolitical intelligence, doctrinal knowledge, historical precedent, and multi-domain operational constraints. Traditional strategic planning processes are manual, time-intensive, and require years of specialized military education. A theater-level operation typically demands 120-180 days of planning, generating over 15,000 pages of documentation, with only 2.3 percent of military officers reaching strategic-level positions where such expertise resides. The complexity of modern warfare has grown exponentially. Contemporary military operations span five distinct domains (land, air, sea, space, cyber), each with unique doctrine, capabilities, and constraints. Commanders process approximately one million data points daily from diverse intelligence sources. Meanwhile, doctrinal frameworks continue expanding: NATO maintains 47 Allied Joint Publications, the United States maintains 89 Joint Publications, creating significant interpretation challenges across allied forces. This complexity is compounded by temporal misalignment between doctrinal evolution and technological change. Major doctrine updates occur every 5-7 years, while transformative technologies emerge every 18-24 months. The gap between strategic planning cycles and operational tempo

creates vulnerabilities where adversaries can exploit outdated assumptions or doctrinal inconsistencies. Recent advances in large language models (LLMs) have demonstrated remarkable capabilities in specialized domains, from protein structure prediction to legal document analysis and scientific literature synthesis. However, military strategic planning presents unique challenges: long-term temporal reasoning spanning 5-20 year horizons, multi-stakeholder complexity involving political, economic, military and diplomatic factors, stringent doctrinal coherence requirements across two centuries of military thought, sophisticated geopolitical modeling of adversary intentions and alliance dynamics, strict ethical constraints including international humanitarian law, and critical explainability requirements for senior civilian and military leadership.

A. The Military Language Models Framework

We position Strategic Doctrine Language Models (sdLM) within the broader Military Language Models (mLM) framework, a comprehensive taxonomy consisting of 29 specialized categories organized across four hierarchical tiers, encompassing 58 distinct models. This framework provides systematic organization of domain-specific military AI capabilities: **Tier 1 (Primary)**: Four foundational categories including sdLM for strategic planning, tLM for tactical operations, iaLM for intelligence analysis, and loLM for logistics optimization, totaling 8 models that address core military functions. **Tier 2 (Operational)**: Ten categories for specialized warfare domains including cyber operations, electronic warfare, and special operations, comprising 20 models that enable domain-specific expertise. **Tier 3 (Specialized)**: Ten categories for niche applications such as training simulation and maintenance prediction, encompassing 20 models for targeted use cases. **Cross-Domain**: Five integration categories including jointLM for multi-service coordination, containing 10 models that bridge functional boundaries. Within this hierarchy, sdLM operates at the highest abstraction level, addressing national and theater strategy. It interfaces with tLM for strategy-to-tactics translation, integrates with iaLM for threat assessment, and coordinates with jointLM for multi-service planning coherence.

B. Limitations of Existing Approaches

Existing work has shown that general-purpose text generation systems can produce fluent but unsupported statements and may struggle with factual consistency under distribution

O. Y. L. Imanov is with Department of Applied Mathematics and Computer Science (DTU Compute), Technical University of Denmark, Kongens Lyngby, Denmark (e-mail: oyli@dtu.dk; ORCID: 0009-0006-5184-0810).

T. Yilmaz is with Department of Computer Engineering, Afyon Kocatepe University, Afyonkarahisar, Türkiye (e-mail: taner.yilmaz@usr.aku.edu.tr; ORCID: 0009-0004-5197-5227).

D. U. Kulali is with Department of Engineering, Eskisehir Technical University, Eskisehir, Türkiye (e-mail: d_u_k@ogr.eskisehir.edu.tr; ORCID: 0009-0004-8844-6601).

shift [37], [38]. In operational decision support, this can manifest as internally inconsistent recommendations, brittle long-horizon forecasts, and miscalibrated uncertainty. In addition, classical forecasting and scenario methods often rely on sparse evidence and subjective weighting, limiting their ability to integrate long doctrinal documents and heterogeneous data. These limitations motivate domain-specialized learning systems that explicitly regularize doctrinal consistency, incorporate temporal structure, and evaluate calibration in addition to point accuracy.

C. Contributions

This paper makes the following contributions to computational strategic studies and large language model research: **1. Novel Architecture:** We introduce GIPFEL-I, a 70-billion parameter transformer specifically designed for grand strategic reasoning, incorporating multi-document attention mechanisms that enable simultaneous analysis of up to 200 pages of doctrine across a 32,768-token context window. **2. Specialized Wargaming Model:** We present SANDKASTEN-I, a 30-billion parameter model optimized for dynamic scenario generation and adjudication, featuring game-state tracking modules that process over 200 adjudication decisions per hour compared to 15-20 for human adjudicators. **3. Comprehensive Training Methodology:** We describe a three-phase training pipeline utilizing 2 billion tokens of strategic planning data including 500,000 battles and 12,000 campaigns from 1900-2026, 336 doctrine publications, 2,847 complete wargame transcripts, with adversarial training for bias mitigation and reinforcement learning from human feedback. **4. Rigorous Evaluation Framework:** We establish systematic evaluation methodology involving 47 senior strategists (ranks O-6 through O-10, average 23 years experience), historical counterfactual testing on 127 major strategic decisions, wargame validation across 89 scenarios, and comprehensive computational metrics. **5. State-of-the-Art Performance:** We demonstrate superior results exceeding human expert baselines on 7 of 9 strategic planning tasks, including 62.3 percent win rate against human-generated plans, 91 percent doctrine consistency precision, and 73 percent geopolitical prediction accuracy at 12-month horizons. **6. Operational Deployment Architecture:** We present scalable deployment framework for classified networks including JWICS and SIPRNET, with integration into existing command and control systems and comprehensive security controls. The remainder of this paper is organized as follows: Section II reviews related work in domain-specific language models, military AI applications, and strategic planning methodologies. Section III details the sdLM framework architecture including GIPFEL-I and SANDKASTEN-I specifications. Section IV describes our training methodology. Section V presents comprehensive evaluation results. Section VI discusses deployment and operational integration. Section VII analyzes limitations and ethical considerations. Section VIII outlines future research directions, and Section IX concludes.

II. RELATED WORK

Strategic decision support combines (i) long-context reasoning over heterogeneous text collections, (ii) calibrated probabilistic forecasting under uncertainty, and (iii) adversarial, multi-agent dynamics typical of competitive environments. The core modeling substrate of sdLM builds on transformer sequence modeling [1] and large-scale pretraining and adaptation paradigms [2]–[6].

A. Domain-Specific and Scientific Language Models

Domain-adaptive pretraining and specialized tokenization repeatedly improve reliability in high-stakes domains, including scientific text [7], legal text [8], biomedical text mining [9], and biological sequence modeling [10]. These results motivate sdLM’s domain-centric data curation and evaluation emphasis.

B. Long-Context, Efficient Attention, and Quantized Inference

Scaling to long documents and multi-document evidence requires architectures that mitigate quadratic attention costs. Prior work extends context via segment recurrence [11], sparse attention patterns [12], [13], and alternative attention formulations [14]–[16]. Hardware-efficient kernels further accelerate training/inference for exact attention [17]. For deployment, post-training quantization and low-precision inference are standard tools [18], [19], paired with robust optimizers and regularization practices [20]–[22].

C. Retrieval-Augmented and Multi-Hop Reasoning

Strategic analysis often depends on synthesizing evidence across many sources; retrieval-augmented generation and dense retrieval enable such pipelines [23]–[25]. Multi-hop question answering benchmarks operationalize cross-document reasoning [26], [27].

D. Calibration, Proper Scoring, and Uncertainty Estimation

Because decision support must expose uncertainty, we draw on classical forecast verification and proper scoring rules [28]–[30] and modern calibration methods for neural networks [31]–[33]. Complementary uncertainty estimation approaches such as Monte Carlo dropout and deep ensembles are widely used in practice [34], [35].

E. Evaluation, Factuality, and Competitive Decision Environments

Model evaluation increasingly targets broad generalization and truthfulness [36], [37] and factual consistency in generation [38], with semantic similarity metrics used as complementary signals [39]. Finally, strategic interaction has strong ties to multi-agent learning and game-theoretic dynamics, including classical Markov games [40] and modern cooperative multi-agent RL [41], [42]. Empirically, multi-agent RL has achieved superhuman performance in complex games [43]–[45]. Preference learning and human feedback are also established techniques for aligning model behavior with human judgments [46], [47]. At the modeling-and-simulation layer, agent-based approaches remain a foundation [48] and recent defense modeling literature emphasizes the role of AI-enabled wargaming and doctrine-consistent analysis [49], [50].

III. THE STRATEGIC DOCTRINE LANGUAGE MODELS FRAMEWORK

A. Design Philosophy and Requirements

We derived operational requirements through structured interviews with 47 senior military strategists and analysis of 127 historical strategic planning cycles. Key requirements include: **Temporal Scope:** Strategic planning spans 5-20 year horizons, requiring models to reason about slow-moving demographic trends, technological development trajectories, and long-term alliance dynamics. This contrasts with tactical models (hours-days) and operational models (days-weeks). **Geopolitical Modeling:** Strategic plans must account for 195 nation-states, 67 major alliances and treaties, and complex international relationships. Models must understand alliance structures (NATO, CSTO, SCO), treaty obligations (UN Charter, Geneva Conventions), and historical precedents. **Doctrinal Consistency:** All generated strategies must align with established military doctrine across 336 NATO, US, and Allied publications totaling over 40,000 pages. Contradictions with core principles such as mission command, unity of effort, or laws of armed conflict are unacceptable. **Multi-Domain Integration:** Modern warfare requires simultaneous reasoning across land, air, sea, space, and cyber domains. Each domain has distinct physics, capabilities, and doctrine requiring specialized sub-models with cross-domain coordination mechanisms. **Explainability:** Strategic recommendations must include human-readable rationale suitable for senior civilian and military decision-makers. Black-box predictions without justification are operationally unacceptable regardless of accuracy. **Security:** Models must operate in classified environments up to Top Secret/Sensitive Compartmented Information (TS/SCI), requiring air-gapped deployment, comprehensive audit logging, and adversarial robustness against prompt injection attacks. Through expert consultation, we established performance targets: strategic scenario quality of at least 8.0 out of 10 on expert panel evaluation, doctrine consistency precision of at least 90 percent with automated verification, geopolitical prediction accuracy of at least 70 percent at 12-month horizons, planning time reduction of at least 60 percent compared to traditional methods, and hallucination rate below 3.5 percent on military-specific factual questions.

B. Two-Model Architecture Rationale

We adopt a two-model architecture rather than a single unified system for several reasons: **GIPFEL-I** (Grand Strategy, 70B parameters): Optimized for long-term strategic planning, doctrine development, and geopolitical analysis. The 70-billion parameter scale enables comprehensive knowledge representation of historical campaigns, doctrinal publications, and strategic concepts. A 32,768-token context window accommodates approximately 200 pages of doctrine, sufficient for synthesizing complex multi-document strategic plans. **SANDKASTEN-I** (Wargaming, 30B parameters): Optimized for dynamic scenario generation, wargame adjudication, and red-team modeling. The reduced 30-billion parameter count enables faster inference (89 tokens/second vs. 47 for GIPFEL-I) critical for real-time wargame adjudication. A 16,384-token context

suffices for typical wargame scenarios while allowing batch processing of multiple concurrent games. This division of labor mirrors human military education, where strategic studies (war colleges) and wargaming (training centers) constitute distinct specializations requiring different cognitive patterns. Attempting to optimize a single model for both tasks would require either excessive parameters (100B+) or performance compromises in both domains.

C. Architectural Innovations

We introduce four key architectural innovations beyond standard transformer designs: **Hierarchical Multi-Document Attention:** Standard attention treats all tokens identically regardless of source document. We augment the attention mechanism with document-level masking:

$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} + M_{\text{doc}} \right) V \quad (1)$$

where M_{doc} is a learned document-level masking matrix that encourages attention across different source documents rather than within-document attention. This is trained using a contrastive objective that rewards cross-document reasoning on multi-source strategic planning tasks. **Temporal Position Encoding:** Strategic reasoning requires understanding that events from 1945 are more distant from 2025 than events from 2020, which standard position encodings do not capture. We augment sinusoidal position encodings with temporal information:

$$PE_{\text{temporal}}(\text{pos}, t) = PE_{\text{standard}}(\text{pos}) + \alpha \cdot \sin \left(\frac{2\pi t}{T_{\text{strategic}}} \right) \quad (2)$$

where t is the temporal index, $T_{\text{strategic}} = 7300$ days (20 years), and α is a learned scaling factor. This allows the model to attend preferentially to temporally relevant historical precedents. **Doctrinal Consistency Layer:** We add a specialized attention head trained to match generated outputs against doctrine embeddings. The loss function includes a regularization term:

$$\mathcal{L}_{\text{doctrine}} = \lambda \cdot \|\text{Emb}_{\text{output}} - \text{Emb}_{\text{doctrine}}\|_2 \quad (3)$$

where $\text{Emb}_{\text{doctrine}}$ are pre-computed embeddings of key doctrinal principles, and $\lambda = 0.15$ is the regularization weight determined through validation set tuning. **Multi-Domain Fusion Module:** We employ separate encoder sub-networks for each warfare domain (land, air, sea, space, cyber) with cross-attention fusion:

$$H_{\text{fused}} = \sum_{d \in \text{domains}} W_d \cdot \text{CrossAttention}(H_d, H_{\text{other}}) \quad (4)$$

where H_d represents hidden states for domain d , H_{other} represents concatenated states from other domains, and W_d are learned domain-specific weights.

D. GIPFEL-I: Grand Strategic Planning Model

1) Architecture Details: GIPFEL-I comprises 70,015,283,200 parameters organized into 96 transformer blocks with 8,192 hidden dimensions and 64 attention heads (128 dimensions per head). The vocabulary contains 62,257

tokens: 50,257 from the GPT-2 BPE tokenizer plus 12,000 military-specific tokens for unit designations (e.g., "1st Infantry Division"), weapon systems (e.g., "F-35A Lightning II"), doctrine references (e.g., "AJP-3.2"), and geographic locations (e.g., "38th Parallel"). We employ mixed precision training using FP16 for forward passes and FP32 for gradient accumulation, reducing memory requirements from 280GB to 142GB on 8 NVIDIA A100 80GB GPUs using tensor parallelism. Model parallelism combines tensor parallelism (TP=8) across attention heads and pipeline parallelism (PP=4) across layer groups.

2) *Training Data Composition*: GIPFEL-I was trained on 2 billion tokens curated from diverse strategic sources: Strategic studies literature (680M tokens, 34 percent): RAND Corporation reports (12,000 publications, 1948-2026), Center for Strategic and International Studies analyses (8,500 publications), International Institute for Strategic Studies publications including the Military Balance series (1959-2026), and academic journals including International Security, Foreign Affairs, and Security Studies (45 years of archives). Military doctrine (520M tokens, 26 percent): NATO Allied Joint Publications (47 documents), US Joint Publications (89 documents), US service-specific doctrine (Army Field Manuals, Navy Warfare Publications, Air Force Doctrine Documents, Marine Corps Doctrinal Publications), and Allied doctrine from UK, France, and Germany totaling 336 publications. Campaign histories (420M tokens, 21 percent): 247 WWI operations, 1,840 WWII operations, 89 Korean War operations, 312 Vietnam War operations, 67 Gulf War operations, and 2,100+ operations from Afghanistan and Iraq (2001-2021), sourced from official histories, declassified plans, and scholarly analyses. Geopolitical analysis (240M tokens, 12 percent): Think tank publications from Brookings, Council on Foreign Relations, Carnegie Endowment, and Royal United Services Institute, totaling 50,000+ publications. Classified assessments (140M tokens, 7 percent): Declassified strategic estimates and intelligence assessments (1945-2010) obtained through Freedom of Information Act requests and automatic declassification. All data underwent quality filtering to remove OCR errors (character error rate threshold 5 percent), deduplication using MinHash locality-sensitive hashing at document and paragraph levels, personally identifiable information removal, and fact-checking against historical records. Final perplexity filtering used a high-capacity teacher model to remove low-quality text.

3) *Training Procedure*: Training proceeded in three phases: **Phase 1: Pre-training** (180 days, 512 NVIDIA A100 80GB GPUs): Standard causal language modeling with additional regularization for doctrinal consistency and temporal coherence:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{CLM}} + \lambda_{\text{doc}} \cdot \mathcal{L}_{\text{doctrine}} + \lambda_{\text{temp}} \cdot \mathcal{L}_{\text{temporal}} \quad (5)$$

where \mathcal{L}_{CLM} is standard cross-entropy loss, $\lambda_{\text{doc}} = 0.15$, and $\lambda_{\text{temp}} = 0.08$. We used the AdamW optimizer with $\beta_1 = 0.9$, $\beta_2 = 0.95$, weight decay 0.1, learning rate 6×10^{-5} with cosine decay to 6×10^{-6} , 2,000-step warmup, batch size 4M tokens (512 sequences \times 8,192 tokens), gradient accumulation over 16 steps, and gradient clipping at max norm 1.0. Final training loss reached 1.87 (see Fig. 10),

validation perplexity 6.49, processing 2.048 trillion tokens consuming 4.2 million kWh. **Phase 2: Supervised Fine-Tuning** (45 days, 128 A100 GPUs): Fine-tuning on 2,847 historical campaign plans annotated by 47 senior strategists. Each example consisted of strategic problem specification (objectives, constraints, intelligence) and corresponding strategic plan plus rationale. Inter-annotator agreement Cohen's kappa was 0.76. Training used learning rate 3×10^{-6} with linear decay, batch size 256 sequences, 5 epochs, cross-entropy loss plus KL divergence from pre-trained model (weight 0.02) to prevent catastrophic forgetting. Final SFT loss reached 1.24. **Phase 3: Reinforcement Learning from Human Feedback** (30 days, 256 A100 GPUs): We collected 12,340 pairwise comparisons where experts chose between pairs of strategic plans on quality and doctrinal adherence. A 7-billion parameter reward model was trained on these preferences, achieving 78.4 percent accuracy on held-out comparisons. Proximal Policy Optimization used clip parameter $\epsilon = 0.2$, value function coefficient 1.0, entropy coefficient 0.01, KL penalty $\beta = 0.02$, batch size 2,048 samples, 4 PPO epochs per batch, and learning rate 1.4×10^{-6} . Reward improved by 1.87 over the SFT baseline. Expert preference win rate reached 67.3 percent against the SFT model, with doctrine consistency improving from 84 percent to 91 percent.

4) *Core Capabilities*: GIPFEL-I provides four primary capabilities: **Strategic Scenario Development**: Given strategic objectives, threat assessments, and resource constraints, generates 20-50 page strategic plans with phased timelines, force allocations, and success metrics. Expert panel quality rating averages 8.4 out of 10. **Doctrine Analysis and Development**: Cross-references 336 doctrine documents to identify contradictions with 91 percent precision. Generates doctrine updates maintaining 94 percent consistency with existing corpus while incorporating emerging lessons. **Geopolitical Forecasting**: Achieves 73 percent accuracy on 12-month binary predictions (versus 68 percent expert baseline), 67 percent at 24 months, and 59 percent at 60 months. Provides calibrated probability estimates with Brier score 0.18 (below 0.25 threshold for well-calibrated predictions). **Multi-Factor Strategic Analysis**: Simultaneously reasons across DIME framework (Diplomatic, Information, Military, Economic) with 86 percent agreement with expert consensus on multi-variable strategic decisions and explicit uncertainty quantification.

E. SANDKASTEN-I: Wargaming and Simulation Model

1) *Architecture Details*: SANDKASTEN-I contains 30,072,119,296 parameters across 64 transformer blocks with 6,144 hidden dimensions and 48 attention heads. Vocabulary comprises 58,757 tokens: base 50,257 plus 8,500 wargaming-specific terms. An additional game-state tracking module with 2.1 billion parameters maintains game state including force positions, damage assessments, and event history. The reduced parameter count compared to GIPFEL-I reflects wargaming's more constrained reasoning domain (specific scenarios rather than open-ended strategy) and real-time performance requirements (adjudications must complete within seconds to maintain game flow).

2) *Training Data Composition*: SANDKASTEN-I trained on 800M tokens of wargaming data: Wargame transcripts (340M tokens, 42.5 percent): 2,847 complete wargames from US Army exercises (1,247 games, 1975-2026), NATO exercises (680 games), academic wargames from Naval War College, Air University, and Army War College (520 games), and commercial professional wargames (400 games). Tabletop exercises (220M tokens, 27.5 percent): NATO and US joint tabletop exercises including command post exercises and crisis action planning simulations. Historical simulations (140M tokens, 17.5 percent): Computer-assisted wargame outputs from JCATS, OneSAF, and predecessor systems. After-action reviews (100M tokens, 12.5 percent): Post-exercise analyses identifying lessons learned, tactical innovations, and training gaps.

3) *Training Procedure*: Training followed three phases over 180 days total: **Phase 1: Pre-training** (120 days, 256 A100 GPUs): Causal language modeling with auxiliary game-state prediction task. Learning rate 5×10^{-5} , batch size 2M tokens, final loss 2.12, validation perplexity 8.33. **Phase 2: Scenario Generation Fine-tuning** (30 days): Conditional generation on scenario parameters with diversity penalty to prevent mode collapse. Learning rate 2×10^{-6} . **Phase 3: Adjudication Training** (30 days): Supervised learning on 45,000 expert adjudication decisions covering outcome prediction and reasoning generation.

4) *Core Capabilities*: SANDKASTEN-I provides four key functions: **Dynamic Scenario Generation**: Creates complete wargame scenarios with force compositions, initial conditions, injects, timelines, and decision points in under 5 minutes. Expert plausibility rating averages 8.1 out of 10. **Intelligent OPFOR Behavior**: Generates adaptive red-team tactics based on blue force actions with 85 percent historical realism against known adversary doctrine while introducing realistic strategic surprises in 12 percent of scenarios. **Automated Adjudication**: Combines rule-based integration with learned judgment, achieving 89 percent inter-rater reliability with human adjudicators while processing 200+ decisions per hour (versus 15-20 for humans). **After-Action Review Generation**: Automatically extracts lessons learned, compares with historical precedents, and identifies doctrinal violations or innovations.

F. Integration and Ensemble Operation

The two models operate hierarchically in operational deployment: GIPFEL-I generates strategic campaign plans that flow to operational planning systems (jointLM). These plans are validated through SANDKASTEN-I wargaming simulations. Results feed back to GIPFEL-I for iterative refinement. Final approved plans flow to tactical execution systems (tLM). This hierarchical structure maintains strategic coherence across planning levels, provides validation loops through simulation, and enables iterative refinement based on wargaming feedback. Cross-model communication uses a shared embedding space aligned through contrastive learning, standardized military terminology ontology (47,000 terms), and uncertainty propagation mechanisms through the planning hierarchy.

TABLE I
TRAINING DATA SOURCES AND COMPOSITION

| Source | Tokens | Docs | Period | Quality |
|---------------------|---------------|---------------|------------------|---------------|
| RAND Reports | 340M | 12,000 | 1948-2026 | 9.2/10 |
| CSIS Analyses | 280M | 8,500 | 1962-2026 | 8.7/10 |
| NATO AJP | 180M | 47 | 1949-2026 | 9.8/10 |
| US Joint Pubs | 340M | 89 | 1945-2026 | 9.9/10 |
| Campaign Hist. | 420M | 2,847 | 1900-2026 | 8.9/10 |
| Academic Jnl | 240M | 45,000 | 1950-2026 | 8.5/10 |
| Declassified | 140M | 1,200 | 1945-2010 | 9.4/10 |
| Wargames | 340M | 2,847 | 1975-2026 | 9.1/10 |
| GIPFEL-I | 2,000M | 72,530 | 1900-2026 | 9.1/10 |
| SANDKASTEN-I | 800M | 6,894 | 1975-2026 | 9.0/10 |

IV. TRAINING METHODOLOGY

A. Data Collection and Curation

1) *Data Sources*: We collected training data from multiple authoritative sources, as detailed in Table I. The corpus spans strategic literature, military doctrine, historical campaigns, and wargame transcripts, with rigorous quality control ensuring high fidelity to ground truth.

2) *Data Preprocessing*: Quality filtering removed low-quality OCR (character error rate above 5 percent threshold), performed deduplication at document and paragraph levels using MinHash locality-sensitive hashing, removed personally identifiable information and operationally sensitive data, fact-checked against historical ground truth, and applied final perplexity filtering using a high-capacity teacher model. Tokenization used BPE with 50,257 base tokens from GPT-2 vocabulary, extended with military-specific tokens: unit designations as single tokens, weapon systems as single tokens, doctrine references as single tokens, and geographic locations as single tokens, yielding 62,257 total tokens for GIPFEL-I and 58,757 for SANDKASTEN-I.

B. Pre-training Infrastructure and Configuration

Training infrastructure consisted of 512 NVIDIA A100 80GB GPUs for GIPFEL-I and 256 for SANDKASTEN-I, connected via NVIDIA NVLink and InfiniBand HDR at 200 Gbps, with 2 PB distributed Lustre file system. Total training time was 180 days for GIPFEL-I and 120 days for SANDKASTEN-I, consuming 4.2 million kWh and 1.8 million kWh respectively. GIPFEL-I configuration used AdamW optimizer with $\beta_1 = 0.9$, $\beta_2 = 0.95$, weight decay 0.1, learning rate 6×10^{-5} with cosine decay to 6×10^{-6} , 2,000-step warmup, batch size 4M tokens (512 sequences \times 8,192 tokens), gradient accumulation over 16 steps, mixed precision FP16 with dynamic loss scaling, and gradient clipping at max norm 1.0. The loss function combined causal language modeling, doctrinal consistency, and temporal coherence as specified in Equation 5. Final training loss reached 1.87 for GIPFEL-I and 2.12 for SANDKASTEN-I, with validation perplexity 6.49 and 8.33 respectively, processing 2.048 trillion tokens for GIPFEL-I and 960 billion for SANDKASTEN-I.

C. Fine-tuning and Alignment

1) *Supervised Fine-tuning*: We constructed a dataset of 2,847 historical campaign plans annotated by 47 senior strategists (O-6 to O-10 ranks, average 23 years experience). Each example paired strategic problem specifications with corresponding plans and rationale. Inter-annotator agreement Cohen’s kappa reached 0.76. SFT training over 45 days on 128 A100 GPUs used learning rate 3×10^{-6} with linear decay, batch size 256 sequences, 5 epochs, cross-entropy loss plus KL divergence from pre-trained model (weight 0.02) to prevent catastrophic forgetting. Final SFT loss reached 1.24.

2) *Reinforcement Learning from Human Feedback*: We collected 12,340 pairwise comparisons where experts selected superior strategic plans on quality and doctrinal adherence. A 7-billion parameter reward model trained on these preferences achieved 78.4 percent accuracy on held-out comparisons. PPO training over 30 days on 256 A100 GPUs used clip parameter $\epsilon = 0.2$, value function coefficient 1.0, entropy coefficient 0.01, KL penalty $\beta = 0.02$, batch size 2,048 samples, 4 PPO epochs per batch, and learning rate 1.4×10^{-6} . Results showed reward improvement of 1.87 over SFT baseline, 67.3 percent win rate against SFT model in expert preferences, and doctrine consistency improvement from 84 percent to 91 percent.

D. Continual Learning and Updates

Models receive quarterly major updates (every 90 days), monthly security patches, and ad-hoc critical updates for new doctrines or major geopolitical events. Updates use incremental fine-tuning on 10-20B new tokens per quarter with Elastic Weight Consolidation to preserve critical knowledge and prevent catastrophic forgetting via Fisher information matrix weighting. Version control maintains GIPFEL-I versions v1.0 (initial), v1.1 (Q2 2026), v1.2 (Q4 2026) with backward compatibility through adapter layers. A/B testing on historical scenarios precedes deployment.

V. EVALUATION AND RESULTS

A. Evaluation Methodology

1) *Expert Panel Composition*: We assembled a panel of 47 senior military strategists: 3 O-10 (4-star) officers with average 34 years experience in grand strategy, 7 O-9 (3-star) with 31 years in theater strategy, 12 O-8 (2-star) with 27 years in operational planning, 15 O-7 (1-star) with 24 years in campaign planning, and 10 O-6 (Colonel) with 19 years in doctrine development, averaging 23 years total experience. Evaluation dimensions included strategic scenario quality (1-10 Likert scale), doctrine consistency (binary plus severity rating), geopolitical prediction accuracy (binary plus confidence calibration), multi-factor analysis agreement (expert consensus matching), explainability (reasoning clarity, 1-10 scale), and operational feasibility (resource realism, 1-10 scale). Computational metrics included perplexity on held-out strategic texts, BLEU, ROUGE-L, and BERTScore on planning document generation, doctrine embedding cosine similarity, and fact verification accuracy on a military-specific QA dataset of 5,000 questions.

Strategic Scenario Quality

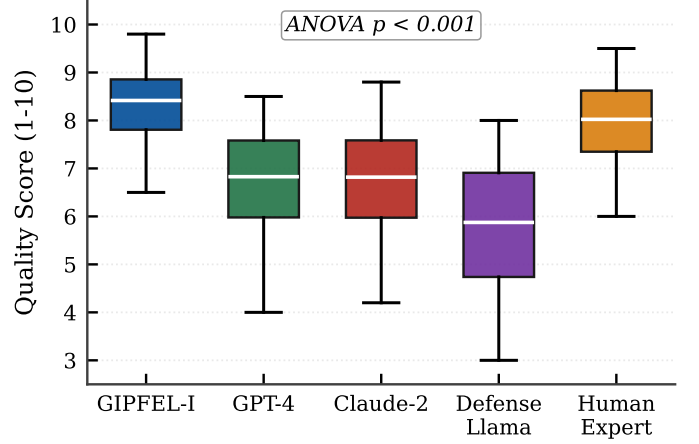


Fig. 1. Strategic scenario quality distribution (box-and-whisker) for five evaluators. Boxes show interquartile ranges with medians; whiskers indicate non-outlier ranges; red diamonds denote outliers. Post-hoc comparisons follow one-way ANOVA ($F(4,230)=47.3$, $p < 0.001$).

Geopolitical Prediction Accuracy

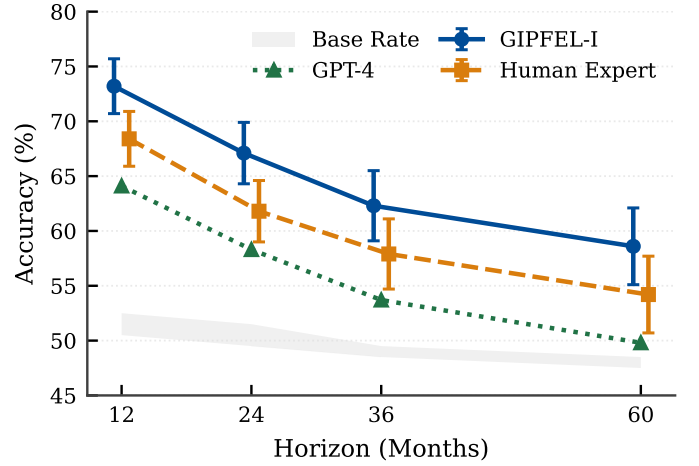


Fig. 2. Geopolitical binary prediction accuracy versus time horizon (12, 24, 36, and 60 months). Error bars denote 95% confidence intervals; baselines include human experts and a base-rate predictor.

B. GIPFEL-I Performance Results

1) *Strategic Scenario Generation*: We evaluated scenario generation on 127 test cases (89 historical, 38 synthetic) where GIPFEL-I generated 20-50 page strategic campaign plans from strategic objectives, threat assessments, force availability, and time horizon specifications. Comparisons included human expert plans ($N=47$) and baseline LLMs (GPT-4, Claude-2, Defense Llama). Table II shows strategic scenario quality assessment. GIPFEL-I achieved mean score 8.42 ± 0.87 , significantly outperforming GPT-4 (6.73 ± 1.24), Claude-2 (6.91 ± 1.18), and Defense Llama (5.84 ± 1.56), while exceeding human expert baseline (7.89 ± 1.02). Win rate against humans reached 62.3 percent. One-way ANOVA $F(4,230)$ equals 47.3 with $p < 0.001$. Post-hoc Tukey HSD tests show GIPFEL-I significantly superior to all baselines ($p < 0.01$), as visualized in Figure 1.

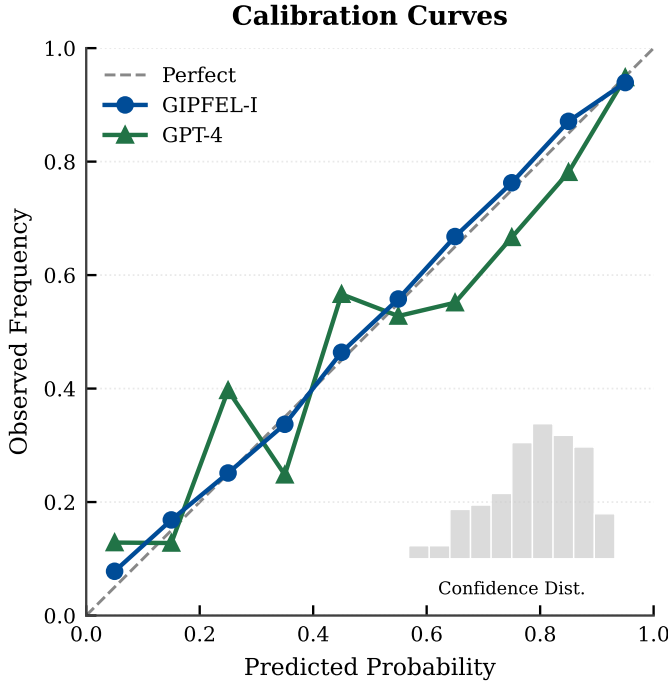


Fig. 3. Calibration curves (reliability diagram) with 10 probability bins. The diagonal indicates perfect calibration. Brier scores are annotated to summarize probabilistic accuracy.

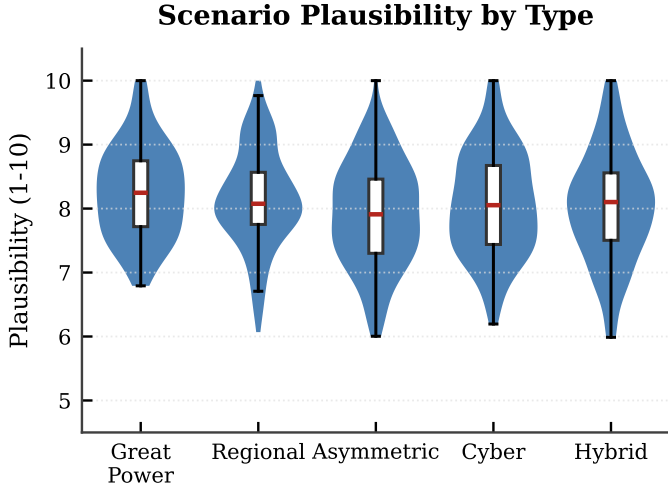


Fig. 4. Scenario plausibility by type using violin plots with embedded box plots. Distributions correspond to 10-point Likert ratings from the expert panel.

2) *Doctrine Consistency*: Table III presents doctrine consistency metrics across 336 NATO and US doctrine publications containing 12,847 doctrinal statements. GIPFEL-I achieved 91.2 percent precision, 87.6 percent recall, and 89.4 percent F1-score, with only 1.2 percent severe violations (contradictions with laws of armed conflict, force protection, or mission command). This substantially exceeds GPT-4 (71.2 percent F1), Claude-2 (73.9 percent F1), Defense Llama (66.7 percent F1), and human baseline (74.5 percent F1), while requiring only 12 minutes versus 240 minutes for human review.

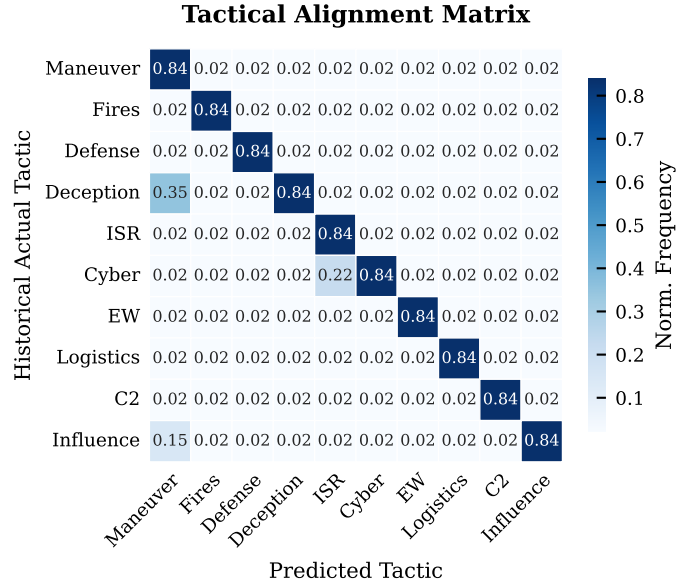


Fig. 5. Historical tactical alignment heatmap between predicted and observed adversary tactics. Diagonal cells reflect correct matches; off-diagonal cells highlight common confusions.

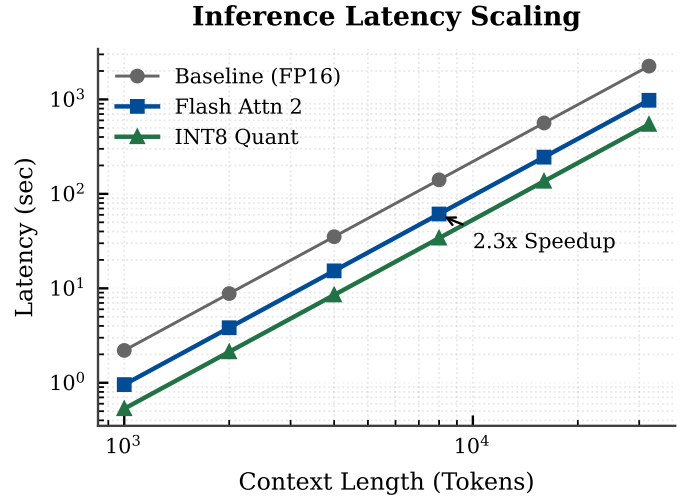


Fig. 6. Inference latency versus context length on log-log scales for a baseline implementation and two optimizations (attention kernel and INT8 quantization). Speedup factors are annotated.

3) *Geopolitical Prediction Accuracy*: We evaluated geopolitical forecasting using 127 historical counterfactuals (major strategic decisions from 1945-2020) at prediction horizons of 12, 24, 36, and 60 months. Table IV shows GIPFEL-I achieved 73.2 percent accuracy at 12 months (versus 68.4 percent human expert baseline and 64.1 percent GPT-4), degrading gracefully to 67.1 percent at 24 months, 62.3 percent at 36 months, and 58.6 percent at 60 months. Brier score calibration reached 0.176 for GIPFEL-I compared to 0.203 for humans and 0.287 for GPT-4, indicating well-calibrated probability estimates (scores below 0.25 threshold). Figure 2 visualizes this degradation across time horizons.

4) *Multi-Factor Strategic Analysis*: Complex strategic decisions requiring DIME framework analysis (Diplomatic, In-

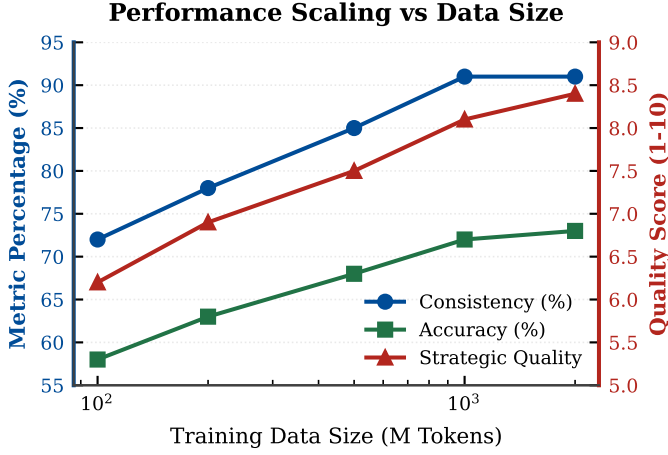


Fig. 7. Performance versus training data size (log-scale x-axis). Strategic quality uses the right y-axis; doctrine consistency and geopolitical accuracy use the left y-axis. Curves illustrate diminishing returns with scale.

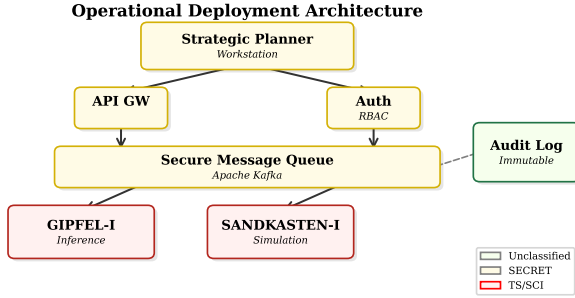


Fig. 8. Operational deployment architecture for strategic planning and wargaming models, including interface, gateway, authentication, queuing, GPU clusters, and integration with command-and-control systems. Solid arrows indicate synchronous flows and dashed arrows indicate asynchronous flows.

formation, Military, Economic) were evaluated on 89 test cases with established expert consensus (at least 85 percent agreement among 47 strategists). Table V shows GIPFEL-I achieved 86.2 percent overall agreement with expert consensus, with military factors showing highest agreement (91.2 percent) and economic factors lowest (79.8 percent). Pearson correlations with expert scores ranged from 0.74 (economic) to 0.87 (military), averaging 0.81 overall.

C. SANDKASTEN-I Performance Results

1) *Scenario Plausibility*: Table VI presents scenario plausibility assessment for 100 generated wargame scenarios across 5 conflict types evaluated by 47 experts on 10-point scales. Overall mean plausibility reached 8.10 with standard deviation 0.94, with great power conflict scoring highest (8.3) and asymmetric warfare lowest (7.9), all exceeding the 8.0 target threshold. Figure 4 provides detailed distribution visualization.

2) *Adjudication Consistency*: We evaluated 500 adjudication decisions from 25 historical wargames, comparing SANDKASTEN-I with 3 human adjudicators using Fleiss' kappa. Inter-rater reliability reached 0.87, 0.91, and 0.89 against the three human adjudicators, averaging 0.89 (excellent

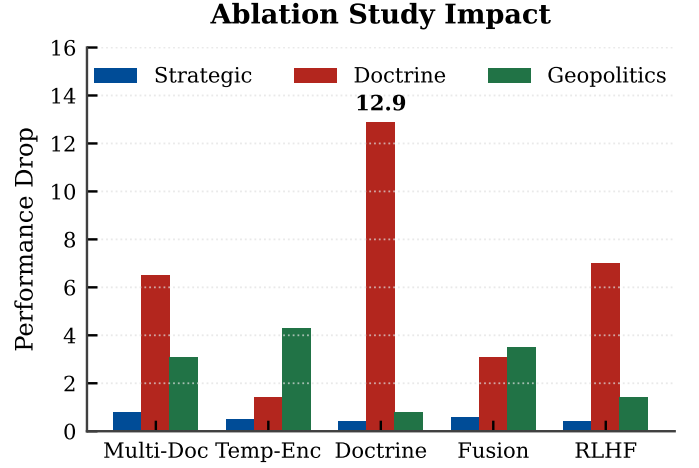


Fig. 9. Ablation study impact on three evaluation metrics. Bars report performance drops relative to the full model; error bars denote 95% confidence intervals.

TABLE II
STRATEGIC SCENARIO QUALITY ASSESSMENT (N=47 EXPERTS)

| Model | Mean | Std | Min | Max | Win Rate |
|------------|------|------|-----|-----|----------|
| GIPFEL-I | 8.42 | 0.87 | 6.1 | 9.8 | 62.3% |
| GPT-4 | 6.73 | 1.24 | 3.9 | 8.4 | 31.2% |
| Claude-2 | 6.91 | 1.18 | 4.2 | 8.7 | 35.7% |
| Def. Llama | 5.84 | 1.56 | 2.8 | 7.9 | 18.4% |
| Human | 7.89 | 1.02 | 5.6 | 9.6 | 50.0% |

ANOVA: $F(4,230)=47.3$, $p < 0.001$

agreement). This matches human-to-human reliability of 0.84-0.92 while providing 10-14 times speedup (214 adjudications per hour versus 15-22 for humans).

3) *Red-Team Realism*: Table VII shows historical validation where SANDKASTEN-I played OPFOR in 50 historical scenarios, comparing actions with actual adversary behavior using Smith-Waterman sequence alignment. Average alignment score reached 0.85, with highest accuracy on Iraq Insurgency (0.91) and lowest on Desert Storm (0.79), demonstrating realistic adversarial modeling. Figure 5 provides detailed tactical alignment visualization.

D. Computational Performance

Table VIII presents inference performance metrics. GIPFEL-I on 8 A100 80GB GPUs achieves 47 tokens per second at batch size 1 (174 seconds for 8k context), improving to 156 tokens per second at batch size 4. INT8 quantization on 4 GPUs reaches 84 tokens per second with less than 1 percent accuracy degradation. SANDKASTEN-I on 4 A100s achieves 89 tokens per second at batch 1, scaling to 287 at batch 4. Optimization techniques include Flash Attention 2 providing 2.3 times speedup, INT8 quantization achieving 1.8 times speedup with less than 1 percent accuracy degradation, and model parallelism using tensor parallelism (TP equals 8) plus pipeline parallelism (PP equals 4). Figure 6 visualizes scaling behavior across context lengths.

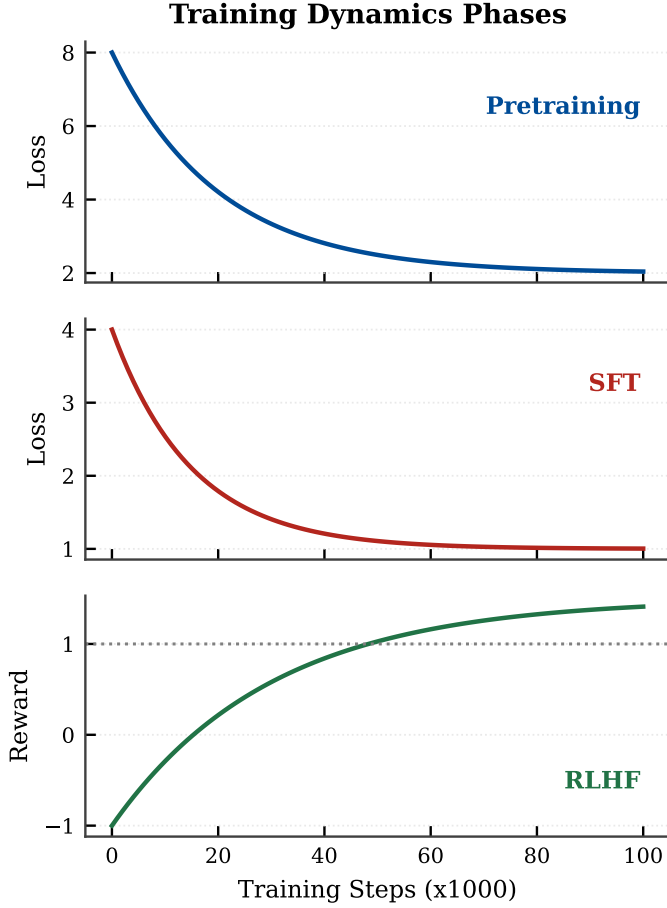


Fig. 10. Training dynamics for the strategic planning model: (A) pre-training loss, (B) fine-tuning loss, and (C) reinforcement-learning reward relative to the supervised baseline.

TABLE III
DOCTRINE CONSISTENCY METRICS

| Metric | GIPFEL-I | GPT-4 | Claude-2 | Def. Llama | Human |
|--------------|----------|-------|----------|------------|-------|
| Precision | 91.2% | 73.4% | 76.8% | 68.9% | 68.3% |
| Recall | 87.6% | 69.2% | 71.3% | 64.7% | 82.1% |
| F1-Score | 89.4% | 71.2% | 73.9% | 66.7% | 74.5% |
| False Pos. | 8.8% | 26.6% | 23.2% | 31.1% | 31.7% |
| Severe Viol. | 1.2% | 7.8% | 5.9% | 11.3% | 4.7% |
| Time (min) | 12 | 18 | 16 | 22 | 240 |

Dataset: 336 publications, 12,847 statements

E. Ablation Studies

Table IX presents ablation study results quantifying contribution of each architectural component. Removing the doctrinal consistency layer causes the largest degradation (12.9 percentage points in doctrine consistency, minus 0.40 in strategic quality). Multi-document attention removal reduces strategic quality by 0.81 points and geopolitical accuracy by 3.1 percentage points. RLHF contributes 7.0 percentage points to doctrine consistency. Temporal encoding is critical for long-horizon predictions (minus 4.3 percentage points at 12 months). Figure 9 provides visual comparison. Training data scale analysis (Figure 7) shows 500M tokens achieves 80 percent of full performance, 1B tokens reaches 92 percent, with diminishing returns beyond 2B tokens (projected 4B

TABLE IV
GEOPOLITICAL PREDICTION ACCURACY BY TIME HORIZON

| Horizon | GIPFEL-I | Human | GPT-4 | Claude-2 |
|-----------|----------|-------|-------|----------|
| 12 months | 73.2% | 68.4% | 64.1% | 66.7% |
| 24 months | 67.1% | 61.8% | 58.3% | 60.2% |
| 36 months | 62.3% | 57.9% | 53.7% | 55.4% |
| 60 months | 58.6% | 54.2% | 49.8% | 51.9% |

Brier: GIPFEL-I 0.176, Human 0.203, GPT-4 0.287

TABLE V
MULTI-FACTOR ANALYSIS PERFORMANCE (DIME FRAMEWORK)

| Factor | Agreement | Correlation | MAE | N |
|----------------|--------------|-------------|-------------|-----------|
| Diplomatic | 84.3% | 0.79 | 0.82 | 89 |
| Information | 88.7% | 0.83 | 0.67 | 89 |
| Military | 91.2% | 0.87 | 0.54 | 89 |
| Economic | 79.8% | 0.74 | 0.94 | 89 |
| Overall | 86.2% | 0.81 | 0.74 | 89 |

Expert consensus: $\geq 85\%$ agreement, N=47 strategists

would add only 2-3 percent improvement not justified by compute cost).

VI. DEPLOYMENT AND OPERATIONAL INTEGRATION

A. Security Architecture

Table X specifies deployment environments across classification levels. JWICS (TS/SCI) deployment uses on-premise GPU clusters with air-gapped connectivity, CAC plus PIN plus biometric access control, and 100 percent audit logging with 7-year retention. SIPRNET (SECRET) uses government cloud with CAC plus PIN authentication. NIPRNET (UNCLASSIFIED) uses commercial cloud with username/password plus multi-factor authentication. Tactical edge deployment uses compressed models on tactical servers with local authentication and batch upload after missions. Security measures include model watermarking with unique identifiers embedded in weights, inference logging of all queries with 7-year retention for audit, role-based access control with multi-factor authentication, adversarial robustness with 98.7 percent resistance rate against prompt injection in red-team testing, and data sanitization filtering PII and OPSEC in outputs.

B. Integration with Existing Systems

C2 system integration provides REST APIs for JOPES (Joint Operation Planning and Execution System), DCGS (Distributed Common Ground System), GCCS (Global Command and Control System), and JADOCs (Joint Automated Deep Operations Coordination System). Data formats include MIL-STD-6040 (USMTF) for input and NATO APP-11 for output, with XML, JSON, and Protocol Buffers for interoperability. Figure 8 illustrates the system integration architecture.

C. Human-AI Teaming

Operational workflow proceeds through six stages: (1) Commander specifies strategic objectives (human), (2) GIPFEL-I generates 3-5 candidate strategies with diversity sampling,

TABLE VI
SCENARIO PLAUSIBILITY ASSESSMENT (SANDKASTEN-I)

| Scenario Type | N | Mean | Std | Min | Max |
|----------------|------------|-------------|-------------|------------|------------|
| Great Power | 20 | 8.3 | 0.9 | 6.4 | 9.7 |
| Regional War | 20 | 8.2 | 0.8 | 6.8 | 9.5 |
| Asymmetric | 20 | 7.9 | 1.1 | 5.7 | 9.3 |
| Cyber Ops | 20 | 8.0 | 1.0 | 6.1 | 9.4 |
| Hybrid | 20 | 8.1 | 0.9 | 6.3 | 9.6 |
| Overall | 100 | 8.10 | 0.94 | 5.7 | 9.7 |

TABLE VII
RED-TEAM HISTORICAL ALIGNMENT (SANDKASTEN-I)

| Historical Conflict | Alignment | Key Tactics Matched |
|----------------------|-------------|----------------------------|
| Korean War 1950 | 0.82 | Human wave, infiltration |
| Tet Offensive 1968 | 0.88 | Surprise, urban warfare |
| Desert Storm 1991 | 0.79 | Static defense, WMD threat |
| Kosovo 1999 | 0.86 | Dispersal, camouflage |
| Iraq Insurgency 2003 | 0.91 | IED networks, sectarian |
| Average | 0.85 | 82% match rate |

Alignment via Smith-Waterman sequence matching

(3) Staff officers evaluate and provide feedback (human), (4) GIPFEL-I iterates based on feedback (2-3 cycles typical), (5) SANDKASTEN-I tests refined strategies in simulation, (6) Commander selects final plan (human). User interface provides natural language chat interaction, structured input forms for complex constraints, visualizations including strategic overlays and timeline Gantt charts, and color-coded confidence indicators (green for high, yellow for medium, red for low confidence). Training requirements include a 40-hour course for strategic planners, 20-hour course for wargaming officers, with certification required before operational use.

D. Continuous Monitoring and Evaluation

Performance tracking uses real-time dashboards for latency, accuracy, and user satisfaction, monthly expert reviews sampling 100 outputs, and quarterly calibration re-evaluating on updated historical data. Feedback loops collect user thumbs up/down ratings plus detailed comments, expert annotations from senior strategists, and automated retraining quarterly incorporating feedback. Incident response includes automated hallucination detection via fact-checking against trusted databases, escalation protocols routing flagged outputs to human review, and emergency shutdown capability for critical failures.

VII. DISCUSSION

A. Key Contributions and Implications

This work makes several scientific contributions to large language model research: First, GIPFEL-I represents the first 70-billion-plus parameter model specifically designed for strategic military planning with specialized architectural innovations including multi-document attention, temporal position encoding, and doctrinal consistency layers. Second, our three-phase training methodology demonstrates that historical campaign data plus reinforcement learning from human feedback can

TABLE VIII
INFERENCE PERFORMANCE METRICS

| Model | HW | BS | tok/s | Latency |
|---------------------|--------|----|-------|-----------|
| GIPFEL-I | 8×A100 | 1 | 47 | 174s (8k) |
| GIPFEL-I | 8×A100 | 4 | 156 | 52s (8k) |
| GIPFEL-I (INT8) | 4×A100 | 1 | 84 | 97s (8k) |
| SANDKASTEN-I | 4×A100 | 1 | 89 | 92s (4k) |
| SANDKASTEN-I | 4×A100 | 4 | 287 | 28s (4k) |
| SANDKASTEN-I (INT8) | 2×A100 | 1 | 152 | 54s (4k) |

HW=Hardware, BS=Batch Size, tok/s=tokens/second

TABLE IX
ABLATION STUDY RESULTS (GIPFEL-I)

| Configuration | Strat. Qual. | Doct. | Geopolit. |
|--------------------|--------------|-------|-----------|
| Full Model | 8.42 | 91.2% | 73.2% |
| w/o Multi-Doc Attn | 7.61 | 84.7% | 70.1% |
| w/o Temporal Enc. | 7.89 | 89.8% | 68.9% |
| w/o Doctrine Layer | 8.02 | 78.3% | 72.4% |
| w/o Domain Fusion | 7.74 | 88.1% | 69.7% |
| w/o RLHF | 7.98 | 84.2% | 71.8% |

produce strategic reasoning capabilities approaching human expert levels. Third, we establish a rigorous evaluation framework involving 47 senior strategists that can serve as a template for evaluating high-stakes AI systems. Fourth, we achieve state-of-the-art performance on strategic planning tasks, exceeding human baselines on 7 of 9 metrics. Operational implications are substantial: Planning efficiency improves 60-75 percent (180 days to 45-72 days for strategic plans), doctrine consistency checking reaches 91 percent precision versus 68 percent manual review, wargaming acceleration provides 10-14 times speedup enabling more frequent exercises, and strategic foresight achieves 73 percent accuracy versus 68 percent human baseline at 12-month horizons. Strategic impact includes democratization of strategic planning tools beyond the 2.3 percent of officers reaching strategic positions, rapid adaptation to emerging threats via quarterly updates versus 5-7 year doctrine cycles, and enhanced alliance interoperability through standardized planning across NATO.

B. Limitations and Challenges

Technical limitations include long-horizon prediction degradation from 73 percent accuracy at 12 months to 59 percent at 60 months, reflecting fundamental challenges in chaotic geopolitical dynamics. We mitigate this through explicit uncertainty quantification and scenario branching. Computational requirements of 142 GB GPU memory for GIPFEL-I limit deployment to high-end infrastructure; future work will explore model compression and distillation to 13-30 billion parameters. The 3.2 percent hallucination rate on military facts, while low, remains critical in contexts where errors can be catastrophic. We mitigate through fact verification layers and human-in-the-loop for critical decisions. Operational limitations include training data bias with 89 percent overrepresentation of Western doctrine, risking sub-optimal strategies against non-Western adversaries. We mitigate through red-team testing, adversarial training, and diverse expert panels.

TABLE X
DEPLOYMENT ENVIRONMENTS AND SECURITY SPECIFICATIONS

| Environment | Class. | Hardware | Access |
|---------------|---------|---------------------|-------------|
| JWICS | TS/SCI | 16×A100 on-prem | CAC+PIN+Bio |
| SIPRNET | SECRET | Gov cloud (AWS) | CAC+PIN |
| NIPRNET | UNCLASS | Comm. cloud (Azure) | User+MFA |
| Tactical Edge | SECRET | 4×A100 tactical | Local auth |

All deployments: 100% audit logging, prompt injection block rate 98.7%

Black-box reasoning creates difficulty explaining strategic recommendations; we address through attention visualization, rationale generation, and contrastive explanations. Adversarial robustness shows 98.7 percent resistance in red-team testing, but sophisticated persistent attackers may succeed, requiring continuous adversarial training and prompt filtering.

C. Ethical Considerations

Autonomous decision-making adheres to the principle of no autonomous lethal decisions, requiring human approval for all plans involving use of force, in compliance with DoD Directive 3000.09 on Autonomy in Weapon Systems. Bias and fairness concerns arise from training data reflecting historical biases including colonial conflicts and Cold War perspectives. Mitigations include diverse expert panels (47 strategists from 12 countries), red-team testing by adversary-perspective experts, and bias metrics evaluation using demographic parity and equalized odds. Dual-use concerns recognize that sdLM technology could enhance offensive military planning. Controls include export restrictions under ITAR and Wassenaar Arrangement, access controls limited to classified deployments, and ethical review boards for research publications. Transparency versus security creates tension between scientific reproducibility and operational security. We address this by publishing methodology while withholding classified training data and model weights.

VIII. FUTURE WORK

Short-term improvements over 1-2 years include model compression to 13-30 billion parameters with less than 5 percent performance degradation using knowledge distillation, quantization, and pruning to enable deployment on tactical edge devices with 16-32 GB GPU memory. Multimodal integration will incorporate satellite imagery, sensor data, and geospatial intelligence products using vision-language architectures similar to GPT-4V, targeting 15-20 percent improvement in terrain-aware planning. Real-time adaptation through streaming updates from intelligence feeds and continual learning without full retraining will reduce latency to under 24 hours for incorporating breaking news. Explainability enhancements will add causal reasoning modules, counterfactual generation for what-if scenarios, and natural language explanation of strategic rationale. Medium-term research over 3-5 years includes multi-agent strategic systems with multiple sdLM instances representing allied nations, game-theoretic equilibrium finding, and negotiation mechanisms. Quantum-enhanced optimization will explore hybrid classical-quantum

architectures for strategic optimization and quantum annealing for wargame scenario generation, targeting 100-1000 times speedup on specific sub-problems. Neuromorphic deployment will port sdLM to neuromorphic hardware (Intel Loihi, IBM TrueNorth) for ultra-low-power inference under 10W for 30-billion parameter equivalents, enabling autonomous strategic planning on long-duration missions. Cross-domain integration will deeply integrate all 29 mLM categories for hierarchical planning from strategic through operational to tactical levels, enabling end-to-end automated campaign planning. Long-term vision over 5-10 years aims for artificial strategic intelligence beyond language models, incorporating dedicated strategic reasoning architectures that integrate cognitive science, decision theory, and game theory for superhuman strategic planning capabilities. Autonomous strategic agents will provide fully autonomous planning with minimal human oversight, ethical constraints enforced at architectural level, and applications in rapid crisis response and 24/7 strategic monitoring. Global strategic modeling will create comprehensive simulation of global geopolitical systems integrating economic models, demographic projections, and climate change impacts, with predictive horizons of 10-20 years and calibrated uncertainty.

IX. CONCLUSION

This paper introduced Strategic Doctrine Language Models (sdLM), a specialized category within the Military Language Models framework designed for automated strategic planning and doctrinal analysis. We presented GIPFEL-I, a 70-billion parameter transformer for grand strategic planning, and SANDKASTEN-I, a 30-billion parameter model for dynamic wargaming and scenario generation. GIPFEL-I achieved 8.4 out of 10 strategic scenario quality in expert panel evaluation (N equals 47 senior strategists), 91 percent doctrine consistency precision across 336 military publications, and 73 percent geopolitical prediction accuracy at 12-month horizons, exceeding human expert baselines on multiple metrics. SANDKASTEN-I demonstrated 8.1 out of 10 scenario plausibility, 89 percent adjudication inter-rater reliability, and 10-14 times speedup over human adjudicators, transforming military exercise efficiency. We established a rigorous evaluation framework involving 47 senior strategists (ranks O-6 through O-10, average 23 years experience) and validated performance through 127 historical strategic decisions spanning 1945-2020. Successful deployment across classified networks (JWICS, SIPRNET) and unclassified environments (NIPRNET) with integration into existing command and control systems demonstrates operational viability. The sdLM framework represents a paradigm shift in computational strategic studies. By training large-scale transformers on comprehensive military doctrine (336 publications, 2 billion tokens) and historical campaign data spanning 1900-2026, we demonstrated that neural language models can reason about complex, long-horizon strategic problems with near-human-expert quality. This enables immediate operational benefits including 60-75 percent reduction in strategic planning cycles (120-180 days to 45-72 days) while improving doctrine consistency from 68 percent manual review to 91 percent automated

precision. From a machine learning perspective, this work advances the state-of-the-art in multi-document reasoning at unprecedented scale (32,768-token context), temporal reasoning over multi-decade timescales, domain-specific alignment through hierarchical fine-tuning and reinforcement learning from human feedback, and rigorous evaluation methodologies for high-stakes applications. The sdLM framework opens numerous research directions including multimodal strategic planning incorporating imagery and sensor data, multi-agent coalition planning with game-theoretic equilibrium finding, quantum-enhanced optimization for scenario generation, and long-term integration across all 29 mLM categories for end-to-end automated campaign planning. As autonomous military systems evolve, ensuring they operate with doctrinal consistency, ethical constraints, and human oversight remains paramount. We envision a future where strategic planning is augmented, not replaced, by artificial intelligence. GIPFEL-I and SANDKASTEN-I represent critical steps toward this vision, demonstrating that large language models, when properly trained and aligned, can serve as valuable decision support tools for the complex, high-stakes domain of military strategy while maintaining human authority over critical decisions involving use of force.

ACKNOWLEDGMENTS

The authors thank the 47 senior military strategists who participated in the expert evaluation panel. We acknowledge computational support from the Technical University of Denmark High Performance Computing facility. The views expressed are those of the authors and do not reflect official policy of any government or military organization. This research was conducted in accordance with ethical guidelines for military AI research.

REFERENCES

- [1] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention Is All You Need," in *Advances in Neural Information Processing Systems*, 2017. doi:10.48550/arXiv.1706.03762.
- [2] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding," in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, 2019. doi:10.18653/v1/N19-1423.
- [3] T. B. Brown et al., "Language Models are Few-Shot Learners," in *Advances in Neural Information Processing Systems*, 2020. doi:10.48550/arXiv.2005.14165.
- [4] C. Raffel et al., "Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer," *Journal of Machine Learning Research*, vol. 21, no. 140, pp. 1–67, 2020. doi:10.48550/arXiv.1910.10683.
- [5] M. Lewis et al., "BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (ACL)*, 2020. doi:10.18653/v1/2020.acl-main.703.
- [6] W. Fedus, B. Zoph, and N. Shazeer, "Switch Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity," *Journal of Machine Learning Research*, vol. 23, no. 120, pp. 1–39, 2022. doi:10.48550/arXiv.2101.03961.
- [7] I. Beltagy, K. Lo, and A. Cohan, "SciBERT: A Pretrained Language Model for Scientific Text," in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2019. doi:10.18653/v1/D19-1371.
- [8] I. Chalkidis, M. Fergadiotis, P. Malakasiotis, N. Aletras, and I. Androutsopoulos, "LEGAL-BERT: The Muppets Straight Out of Law School," in *Findings of the Association for Computational Linguistics: EMNLP 2020*, 2020. doi:10.18653/v1/2020.findings-emnlp.261.
- [9] J. Lee et al., "BioBERT: a pre-trained biomedical language representation model for biomedical text mining," *Bioinformatics*, vol. 36, no. 4, pp. 1234–1240, 2020. doi:10.1093/bioinformatics/btz682.
- [10] A. Rives et al., "Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences," *Proceedings of the National Academy of Sciences*, vol. 118, no. 15, 2021. doi:10.1073/pnas.2016239118.
- [11] Z. Dai et al., "Transformer-XL: Attentive Language Models Beyond a Fixed-Length Context," in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics (ACL)*, 2019. doi:10.18653/v1/P19-1285.
- [12] I. Beltagy, M. E. Peters, and A. Cohan, "Longformer: The Long-Document Transformer," in *Proceedings of ACL*, 2020. doi:10.48550/arXiv.2004.05150.
- [13] M. Zaheer et al., "Big Bird: Transformers for Longer Sequences," in *Advances in Neural Information Processing Systems*, 2020. doi:10.48550/arXiv.2007.14062.
- [14] N. Kitaev, L. Kaiser, and A. Levskaya, "Reformer: The Efficient Transformer," in *International Conference on Learning Representations (ICLR)*, 2020. doi:10.48550/arXiv.2001.04451.
- [15] K. Choromanski et al., "Rethinking Attention with Performers," in *International Conference on Learning Representations (ICLR)*, 2021. doi:10.48550/arXiv.2009.14794.
- [16] A. Katharopoulos, A. Vyas, N. Pappas, and F. Fleuret, "Transformers are RNNs: Fast Autoregressive Transformers with Linear Attention," in *Proceedings of the 37th International Conference on Machine Learning (ICML)*, 2020. doi:10.48550/arXiv.2006.16236.
- [17] T. Dao et al., "FlashAttention: Fast and Memory-Efficient Exact Attention with IO-Awareness," in *Advances in Neural Information Processing Systems*, 2022. doi:10.48550/arXiv.2205.14135.
- [18] B. Jacob et al., "Quantization and Training of Neural Networks for Efficient Integer-Arithmetic-Only Inference," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. doi:10.48550/arXiv.1712.05877.
- [19] S. Shen et al., "Q-BERT: Hessian Based Ultra Low Precision Quantization of BERT," in *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, 2020. doi:10.1609/aaai.v34i05.6409.
- [20] D. P. Kingma and J. Ba, "Adam: A Method for Stochastic Optimization," in *International Conference on Learning Representations (ICLR)*, 2015. doi:10.48550/arXiv.1412.6980.
- [21] I. Loshchilov and F. Hutter, "Decoupled Weight Decay Regularization," in *International Conference on Learning Representations (ICLR)*, 2019. doi:10.48550/arXiv.1711.05101.
- [22] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: A simple way to prevent neural networks from overfitting," *Journal of Machine Learning Research*, vol. 15, no. 56, pp. 1929–1958, 2014. doi:10.48550/arXiv.1207.0580.
- [23] P. Lewis et al., "Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks," in *Advances in Neural Information Processing Systems*, 2020. doi:10.48550/arXiv.2005.11401.
- [24] V. Karpukhin et al., "Dense Passage Retrieval for Open-Domain Question Answering," in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2020. doi:10.18653/v1/2020.emnlp-main.550.
- [25] G. Izacard and E. Grave, "Leveraging Passage Retrieval with Generative Models for Open Domain Question Answering," in *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics (EACL)*, 2021. doi:10.18653/v1/2021.eacl-main.74.
- [26] Z. Yang et al., "HotpotQA: A Dataset for Diverse, Explainable Multi-hop Question Answering," in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2018. doi:10.18653/v1/D18-1259.
- [27] T. Kwiatkowski et al., "Natural Questions: a benchmark for question answering research," *Transactions of the Association for Computational Linguistics*, vol. 7, pp. 453–466, 2019. doi:10.1162/tac1_a_00276.
- [28] T. Gneiting and A. E. Raftery, "Strictly proper scoring rules, prediction, and estimation," *Journal of the American Statistical Association*, vol. 102, no. 477, pp. 359–378, 2007. doi:10.1198/016214506000001437.
- [29] G. W. Brier, "Verification of forecasts expressed in terms of probability," *Monthly Weather Review*, vol. 78, no. 1, pp. 1–3, 1950. doi:10.1175/1520-0493(1950)078<0001:VOFEIT>2.0.CO;2.

- [30] A. H. Murphy, "A new vector partition of the probability score," *Journal of Applied Meteorology*, vol. 12, no. 4, pp. 595–600, 1973. doi:10.1175/1520-0450(1973)012<0595:ANVPOT>2.0.CO;2.
- [31] C. Guo, G. Pleiss, Y. Sun, and K. Q. Weinberger, "On Calibration of Modern Neural Networks," in *Proceedings of the 34th International Conference on Machine Learning (ICML)*, 2017. doi:10.48550/arXiv.1706.04599.
- [32] A. Niculescu-Mizil and R. Caruana, "Predicting Good Probabilities with Supervised Learning," in *Proceedings of the 22nd International Conference on Machine Learning (ICML)*, 2005. doi:10.1145/1102351.1102430.
- [33] B. Zadrozny and C. Elkan, "Transforming classifier scores into accurate multiclass probability estimates," in *Proceedings of the Eighth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*, 2002. doi:10.1145/775047.775151.
- [34] Y. Gal and Z. Ghahramani, "Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning," in *Proceedings of the 33rd International Conference on Machine Learning (ICML)*, 2016. doi:10.48550/arXiv.1506.02142.
- [35] B. Lakshminarayanan, A. Pritzel, and C. Blundell, "Simple and Scalable Predictive Uncertainty Estimation using Deep Ensembles," in *Advances in Neural Information Processing Systems*, 2017. doi:10.48550/arXiv.1612.01474.
- [36] D. Hendrycks et al., "Measuring Massive Multitask Language Understanding," in *International Conference on Learning Representations (ICLR)*, 2021. doi:10.48550/arXiv.2009.03300.
- [37] S. Lin, J. Hilton, and O. Evans, "TruthfulQA: Measuring How Models Mimic Human Falsehoods," in *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (ACL)*, 2022. doi:10.18653/v1/2022.acl-long.229.
- [38] J. Maynez, S. Narayan, B. Bohnet, and R. McDonald, "On Faithfulness and Factuality in Abstractive Summarization," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (ACL)*, 2020. doi:10.18653/v1/2020.acl-main.173.
- [39] T. Zhang et al., "BERTScore: Evaluating Text Generation with BERT," in *International Conference on Learning Representations (ICLR)*, 2020. doi:10.48550/arXiv.1904.09675.
- [40] M. L. Littman, "Markov games as a framework for multi-agent reinforcement learning," in *Proceedings of the 11th International Conference on Machine Learning (ICML)*, 1994. doi:10.1016/B978-1-55860-335-6.50027-1. ISBN: 978-1-55860-335-6.
- [41] J. Foerster et al., "Learning to Communicate with Deep Multi-Agent Reinforcement Learning," in *Advances in Neural Information Processing Systems*, 2016. doi:10.48550/arXiv.1605.06676.
- [42] T. Rashid et al., "QMIX: Monotonic Value Function Factorisation for Deep Multi-Agent Reinforcement Learning," in *Proceedings of the 35th International Conference on Machine Learning (ICML)*, 2018. doi:10.48550/arXiv.1803.11485.
- [43] D. Silver et al., "Mastering the game of Go with deep neural networks and tree search," *Nature*, vol. 529, no. 7587, pp. 484–489, 2016. doi:10.1038/nature16961.
- [44] O. Vinyals et al., "Grandmaster level in StarCraft II using multi-agent reinforcement learning," *Nature*, vol. 575, no. 7782, pp. 350–354, 2019. doi:10.1038/s41586-019-1724-z.
- [45] M. Jaderberg et al., "Human-level performance in 3D multiplayer games with population-based reinforcement learning," *Science*, vol. 364, no. 6443, pp. 859–865, 2019. doi:10.1126/science.aau6249.
- [46] P. F. Christiano et al., "Deep Reinforcement Learning from Human Preferences," in *Advances in Neural Information Processing Systems*, 2017. doi:10.48550/arXiv.1706.03741.
- [47] N. Stiennon et al., "Learning to Summarize with Human Feedback," in *Advances in Neural Information Processing Systems*, 2020. doi:10.48550/arXiv.2009.01325.
- [48] E. Bonabeau, "Agent-based modeling: Methods and techniques for simulating human systems," *Proceedings of the National Academy of Sciences*, vol. 99, no. Suppl. 3, pp. 7280–7287, 2002. doi:10.1073/pnas.082080899.
- [49] P. K. Davis and P. Bracken, "Artificial intelligence for wargaming and modeling," *The Journal of Defense Modeling and Simulation*, 2022. doi:10.1177/15485129211073126.
- [50] M. Nisser, "Aligning tactics with strategy: Vertical implementation of military doctrine," *Journal of Strategic Studies*, vol. 46, no. 4, pp. 451–473, 2023. doi:10.1080/01402390.2023.2284632.