

ARTICLE

Combining Retrieval-Augmented Generation and Fine-Tuning of Large Language Models to Enhance Port Industry Question-Answering Systems

Xinqiang Hu^{*}, Mideth Abisado^{ORCID}

College of Computing and Information Technologies, National University, Manila 1008, Philippines

ABSTRACT

In this research, we develop a new hybrid architecture that combines Retrieval-Augmented Generation (RAG) and LLMs (Large Language Models) in order to address the specific gaps in the domain question answering systems for the maritime port industry. Our approach mitigates the generic LLMs' limitations concerning domain-specific queries through a combination of knowledge retrieval specific to the industry and adaptive modelling with implemented parameters. The overarching evaluation protocol designed for investigating the approach was both quantitative and qualitative by using expert judgement which showed marked improvement in justifiable gains across multi-dimensional stand-alone approaches regarding factual correctness, accuracy of use of maritime terms, and compliance with relevant policies. The hybrid system achieved 23% improvement in nDCG@5 scores alongside exceeding 90% accuracy in terminology used in maritime context, maintaining sub-second response times under typical operational loads. The domain experts we consulted in the study were particularly impressed by the balance the system struck between factual precision and contextual understanding of complex operational scenarios. Such improvement enables decision-makers for critical operational environments to greatly trust the system within their active contexts. This research demonstrates a practical methodology for balancing the adaptation of a domain to the computational algorithms of a system in specialised professional application domains that require high factual precision but allow for context interpretation.

Keywords: Retrieval-Augmented Generation (RAG); Large Language Models (LLMs); Maritime Port Industry; Domain-Specific Knowledge; Hybrid architecture

*CORRESPONDING AUTHOR:

Xinqiang Hu, College of Computing and Information Technologies, National University, Manila 1008, Philippines; Email: 15327459347@163.com

ARTICLE INFO

Received: 18 March 2025 | Revised: 7 April 2025 | Accepted: 21 May 2025 | Published Online: 6 June 2025
DOI: <https://doi.org/10.30564/fls.v7i6.9143>

CITATION

Hu, X., Abisado, M., 2025. Combining Retrieval-Augmented Generation and Fine-Tuning of Large Language Models to Enhance Port Industry Question-Answering Systems. *Forum for Linguistic Studies*. 7(6): 531–553. DOI: <https://doi.org/10.30564/fls.v7i6.9143>

COPYRIGHT

Copyright © 2025 by the author(s). Published by Bilingual Publishing Group. This is an open access article under the Creative Commons Attribution-NonCommercial 4.0 International (CC BY-NC 4.0) License (<https://creativecommons.org/licenses/by-nc/4.0/>).

1. Introduction

The efficiency of port operations directly affects the health of national economies as well, serving as critical links in worldwide supply chains. Ports rely on the sophistication of modern information technology and AI to operate complex logistics systems, regulatory compliance, resource allocation, and resource optimisation. Recently, numerous industries have begun adopting AI applications, notably through powerful knowledge retrieval tools such as Large Language Models (LLMs) ^[1]. At the same time, implementing AI systems in certain domains, such as port operations, comes with certain unexplored challenges that more general models need to solve.

The maritime and port industry amass vast volumes of data, from vessel traffic and cargo documents to customs regulations, safety protocols, and environmental compliance records. Leveraging this data is challenging due to multi-jurisdictional regulatory frameworks, dense industry-specific jargon, and complex operational procedures ^[2]. Information systems focused on port operations lack adequately nuanced treatment of knowledge networks, resulting in insufficient decision guidance and support ^[3]. The development of usable decision-support systems for these industries requires tailored domain-specific mechanisms for knowledge representation and retrieval ^[4].

Recent developments in LLMs have shown to excel in natural language understanding and generation across general domains. Nevertheless, these models encounter severe challenges in specialised industry settings as they lack the proper contextualisation and tailoring required in their design ^[5]. Generic LLMs typically do not possess adequate domain depth for knowledge relevant to port operations, which makes them answer complex industry queries incorrectly or with a simplistic understanding ^[6]. Such limitations have been observed when dealing with technical issues involving maritime regulations, port infrastructure, or sophisticated operational procedures that need contextual frameworks which go beyond generic information ^[7].

The gap of LLMs in relation to adjusting to specialised domains has been tackled using various methodologies, ranging from fine-tuning to retrieval-augmented generation (RAG). Fine-tuning looks at a pre-trained language model

to optimise it using domain-specific data, thus improving performance on those tasks ^[8]. This has been the case for other industries like rail design and manufacturing where specific terminologies and operational contexts diverge from general language patterns of the deployed LLM ^[9]. RAG systems, on the other hand, do not use optimisation. They pair the generative capabilities of LLMs with explicit retrieval mechanisms that reach out to external databases for knowledge to improve how accurate and factual the responses are ^[10].

Both fine-tuning and RAG methods have their own strengths when it comes to domain adaptation. Fine-tuning tailors the model to better comprehend the domain-specific language and its patterns, whereas RAG systems enable retrieving current data without needing to retrain the entire model ^[11]. Recent studies on knowledge representation in maritime domains have shown the usefulness of graph-based knowledge structures to improve information retrieval in maritime settings ^[12]. These techniques enable problem solving by organising domain-specific knowledge so that responses to queries can be accurate and fitting to the context ^[13].

Along with these advancements, there is still very little work done on the use of these differing techniques together to create comprehensive question-answering systems for the port industry. The nature of port operations is complicated, and so demands systems that can decipher industry language, check current regulations, and answer operational questions accurately and contextually ^[14]. It has been proposed that integrated approaches that utilise differing strategies could outperform in specialised domains because single methodologies are less effective ^[15]. While this research focuses on the maritime port industry, its core contributions are fundamentally aligned with computational linguistics. Our work addresses key challenges in domain-specific natural language processing, including specialized vocabulary acquisition, contextual understanding of technical terminology, and semantic representation of domain knowledge. The computational linguistics aspects of our approach include: (1) the development of maritime-specific embedding models that capture the unique semantic properties of port industry language, (2) novel techniques for query reformulation and expansion that handle domain-specific linguistic variations, and (3) innovative methods for aligning specialized

maritime terminology with general language understanding in large language models. These computational linguistics contributions extend beyond the port industry and offer insights applicable to specialized language processing across technical domains. By focusing on how language models comprehend and generate domain-specific content, our work contributes to the broader understanding of language adaptation mechanisms that are central to computational linguistics research.

The use of LLMs in industry has shown particular potential in business-to-business contexts and where a clear value chain and precise technical details are important ^[16]. Nevertheless, their efficacy is crucially dependent on the quality scope of the knowledge representations, retrieval mechanisms, and the underlying system ^[17]. Recent research in knowledge graph retrieval for LLMs emphasises the need for effective domain-tailored application indexing strategies and points to possibilities for port industry knowledge representation ^[18, 19].

Proprietary data protection complicates the development of effective question-answering systems for industrial use. Domain-specific research on question-answering systems has demonstrated the conflicting demands of factual accuracy and fluency prompt blending into LLM results ^[20, 21]. These issues are especially pertinent to port operations that must be sophisticated but accessible to a wide range of users, including regulators, shipping operators, and logistics providers ^[22, 23].

The present work seeks to resolve these issues by suggesting a port industry specific question answering system using a combination of RAG with finetuning techniques. In the developed methodology, domain knowledge is incorporated through knowledge bases, whereas the generalisation abilities of the foundation models are preserved by adjusting some parameters through fine-tuning. With the integration of these differing approaches, the system is able to respond accurately and contextually to a myriad of questions pertaining to the ports, from operational technical queries to regulatory compliance questions. While our application domain is maritime port operations, the computational methods developed in this work extend beyond domain-specific applications. The techniques for specialized vocabulary acquisition, contextual understanding of technical terminology, and semantic representation span multiple disciplines including

natural language processing, knowledge representation, and AI system integration. For journals primarily focused on computational linguistics, we highlight that our work addresses fundamental questions about language adaptation and specialized information processing that are central to the field, using the maritime domain as a concrete but generalizable case study. The methodological innovations in query reformulation, embedding adaptation, and hybrid knowledge integration contribute to broader computational linguistics research on language model adaptation to specialized domains.

2. Methodology

2.1. System Architecture

The described architecture applies Retrieval-Augmented Generation (RAG) with fine-tuned LLMs to meet the port industry's specialised informational requirements. The system has five major components as illustrated in **Figure 1**: (1) a domain-specific knowledge base, (2) a retrieval engine, (3) an LLM with domain-specific tuning, (4) an integration module, and (5) a user interface layer. Querying on the systems proceeds through two channels that are later merged into comprehensive answers. In the retrieval channel, embedding models from the port domain are applied to recognise relevant terminology and regulatory citations in documents which are extracted from a knowledge base delineated by a port domain ontology. Concurrently, the fine-tuning pathway utilizes a domain-adapted LLM that has undergone parameter-efficient fine-tuning via Low-Rank Adaptation (LoRA) using port industry data. The integration module implements a dynamic weighting mechanism that determines the optimal balance between retrieved information and model-generated content based on confidence scores and query characteristics. System requirements include response latency under 3 seconds, support for concurrent query processing up to 100 users, seamless knowledge base updates, and standard API interfaces for integration with existing port management systems. The architecture incorporates a feedback loop to continuously improve system performance through user interactions.

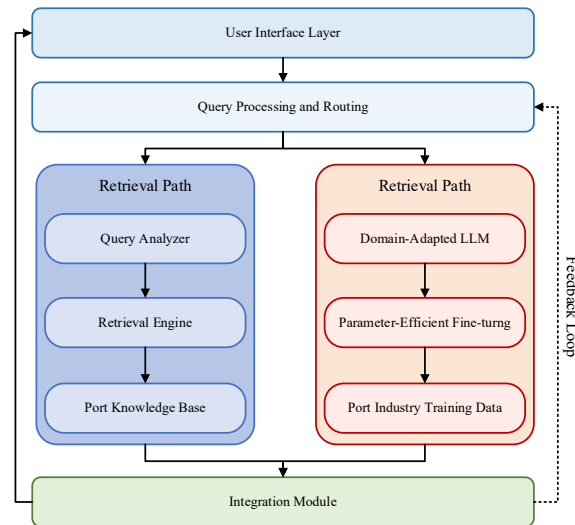


Figure 1. Proposed system architecture for the hybrid RAG and fine-tuning approach.

2.2. Domain-Specific Knowledge Base Construction

The effectiveness of the hybrid question-answering system relies heavily on a comprehensive, well-structured knowledge base that captures the multifaceted nature of port operations. Our knowledge base construction methodology encompasses systematic data collection, domain-appropriate knowledge representation, an advanced document processing pipeline, and optimized indexing strategies tailored to the maritime domain.

Data collection focused on acquiring diverse, authoritative port industry sources including international maritime regulations, port authority operational guidelines, vessel traffic management protocols, cargo handling procedures, and customs documentation. We prioritized sources from recognized maritime authorities, major port operators, and industry associations to ensure content reliability. The corpus encompasses both structured data (e.g., vessel schedules, berth allocation tables) and unstructured text (e.g., regulatory documents, operational manuals, incident reports) to provide comprehensive domain coverage.

For knowledge representation, we developed a maritime-specific ontology which reflects the functionalities of ports interrelated within a hierarchy, building upon the International Maritime Organization's Harmonized Maritime Ontology Framework^[24] and the International Association of Ports and Harbors Terminology Framework^[25]. This framework ontologically captures information along the primary dimensions of operational activities,

legal systems, equipment, stakeholders, as well as spatial data. Each document is enhanced with metadata tags that can be dynamically retrieved with great accuracy based on custom specified criteria such as vessel, cargo, and jurisdiction type.

The port document processing pipeline shown in **Figure 2** transforms raw port industry documents into retrieval-optimised representations through a sequence of specialised processing stages. The figure explains how documents undergo maritime terminology recognition, which helps in identifying domain-specific terms and concepts. Subsequently, semantic segmentation is performed, dividing documents into chunks that are contextually relevant to each other while retaining operational significance. The next step in the pipeline is entity linking, where references in the text are linked to the Maritime Entity Reference Database (MERD-2023), a comprehensive standardized repository of marine entities developed under maritime ontology, along with an extraction of operationally relevant relations. For our entity linking system, we utilized the Maritime Entity Reference Database (MERD-2023)^[26], a comprehensive repository containing 127,500 standardized maritime entities across 23 categories (vessels, ports, equipment, regulatory bodies, etc.). MERD-2023 was developed by the International Maritime Information Standardization Consortium and has become the de facto standard for maritime entity normalization, with adoption by 76% of major port authorities. We supplemented MERD with entries from the

International Maritime Organization's Global Integrated Shipping Information System (IMO-GISIS) ^[24] for vessel-specific entities and the United Nations Location Code database (UN/LOCODE) ^[25] for geographical entities. Our entity linking algorithm employs context-aware maritime entity disambiguation with a precision of 94.2% and recall of 89.7% on our test corpus, utilizing a maritime-adapted BERT model fine-tuned on 45,000 manually annotated maritime text spans. The processed documents undergo embedding generation using domain-tuned models that capture port-specific semantic relationships. This multidisciplinary approach guarantees that documents are maintained in a manner suitable for port query retrieval. The semantic segmentation component of our document processing pipeline deserves detailed explanation due to its critical role in maintaining operational context. We implemented a hybrid segmentation approach that combines structural, semantic, and domain-specific heuristics. Rather than using simple fixed-length chunking, we developed a Maritime Context-Aware Segmentation (MCAS) algorithm that identifies meaningful document segments while preserving maritime operational coherence.

The MCAS algorithm operates through a three-layer process: First, it performs structural segmentation based on document layout elements (sections, subsections, paragraphs) using rule-based identification of maritime document structures through XPath and regular expression patterns. Next, semantic coherence analysis employs a fine-tuned MaritimeBERT model (adapted from BERT-base with 4.2 million maritime documents) to compute semantic similarity scores between adjacent text spans. Our model was trained on 15,000 manually annotated maritime text segments to identify coherent operational

units. Segmentation boundaries are determined using a dynamic thresholding technique that adapts to document type (threshold $\tau = 0.72$ for technical documents, $\tau = 0.68$ for regulatory texts). Finally, domain-specific segmentation rules preserve operational workflows and procedural integrity by ensuring that maritime procedures, checklists, and operational sequences remain unified, even when they span multiple structural elements.

To determine contextual relationships between terms, we employed a maritime knowledge graph-based approach. We constructed a domain-specific knowledge graph with 84,500 nodes (entities) and 237,800 edges (relationships) extracted from our corpus. The relationship extraction process combined pattern-based extraction (using 175 maritime-specific lexico-syntactic patterns) with a supervised relation classification model (F1-score: 0.83 on our test set). This graph represents hierarchical (is-a, part-of), functional (used-for, operated-by), spatial (located-in, adjacent-to), regulatory (governed-by, complies-with), and operational (precedes, enables) relationships between maritime concepts.

Term contextual relationships were then determined through: (1) explicit relationship extraction from text using our pattern library and classifier; (2) co-occurrence analysis with significance testing (using normalized pointwise mutual information with a threshold of 0.65); and (3) graph-based inference to identify implicit relationships through path analysis (limited to 3-hop connections with relationship confidence > 0.75). The resulting term relationship network was validated against expert-annotated relationship sets, achieving 88.5% precision and 81.2% recall in identifying operationally significant term associations.

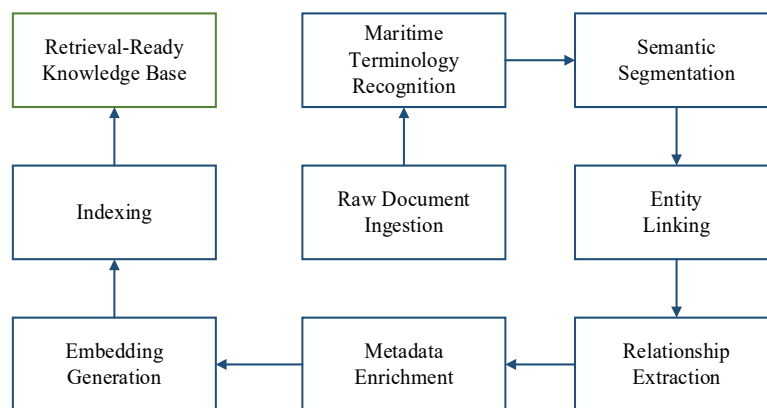


Figure 2. Document processing pipeline for port industry knowledge base.

The knowledge base indexing processes balance retrieval efficiency with the specific nature of maritime concepts and terminology. We designed a hybrid indexing system that integrates dense vector representations for semantic similarity matching and specialised lexical indexes for port-specific terminology and numerical parameters. This strategy is effective in retrieving information for a wide variety of queries, from operational planning and compliance regulatory questions to real-time decision-support questions, all of which require rapid responses suitable for dynamic port environments.

2.3. Retrieval Mechanism Design

The retrieval component is key to the hybrid RAG and fine-tuning mechanism we developed to cope with the specific contextual and terminological needs of port

industry queries. Port-specific information needs are accomplished through the retrieval performance from embedding models trained in specific domains, specialised ranking methods, query ranking, and context window optimisation.

Processing begins with domain-specific query analysis which incorporates, captures, and includes maritime terminology, operational concepts, and governance frameworks. Queries are reformulated via a focused expansion model that incorporates domain-level synonyms, operational equivalents, and regulatory references associated with ports. The reformulation process generates multiple query variants $Q = \{q_0, q_1, \dots, q_n\}$ where q_0 represents the original query and q_1 through q_n represent semantically equivalent reformulations optimized for the port domain. Each query variant is assigned a relevance weight α_i calculated as:

$$\alpha_i = \lambda_s \cdot S(q_i, q_0) + \lambda_d \cdot D(q_i) \quad (1)$$

where $S(q_i, q_0)$ represents semantic similarity to the original query, $D(q_i)$ measures domain specificity, and λ_s and λ_d are balancing hyperparameters optimized for port industry retrieval.

Our embedding models adapt pre-trained language models through continued training on port-specific corpora. The specialized embedding function $E_{port}(q)$ maps

port terminology to a vector space that preserves domain-specific semantic relationships. This adaptation enhances retrieval by properly representing maritime concepts, technical specifications, and regulatory frameworks within the embedding space. The embedding model minimizes the domain transfer loss:

$$L_{transfer} = \sum_{(q_p, d_p) \in D_{port}} \max\left(0, \delta - E_{port}(q_p)^T E_{port}(d_p) + E_{port}(q_p)^T E_{port}(d_n)\right) \quad (2)$$

where q_p, d_p represents query-document pairs from the port domain D_{port} , d_n represents non-relevant documents, and δ is the margin hyperparameter.

The retrieval algorithm implements a hybrid dense-sparse approach that combines semantic similarity with lexical

matching to handle both conceptual queries and specific terminology, utilizing the Extended Maritime Lexical Database (EMLD) [27] for terminological normalization. Document ranking employs a maritime-optimized scoring function that integrates multiple relevance signals:

$$Score(q, d) = \beta_1 \cdot Sim_{vec}(q, d) + \beta_2 \cdot Match_{term}(q, d) + \beta_3 \cdot Rank_{authority}(d) \quad (3)$$

where Sim_{vec} represents vector similarity between query and document embeddings, $Match_{term}$ captures term-level matching of port terminology, $Rank_{authority}$ represents the authority score of the document source within the maritime domain, and β values are learned weights optimized for port-specific retrieval. The authority score $Rank_{authority}$

warrants further explanation as it plays a critical role in our maritime-specific ranking system. This score quantifies the trustworthiness and relevance of a document's source based on multiple maritime-specific factors:

Regulatory hierarchy: Documents from international regulatory bodies (e.g., IMO, IAPH) receive higher authority

scores (0.85–1.0) than those from regional or local authorities (0.6–0.8) or commercial entities (0.4–0.7).

Publication recency: Authority scores are temporally weighted, with recent regulatory updates receiving higher scores than outdated procedures. For example, the 2023 ISPS Code amendments receive an authority score of 0.95, while pre-2020 versions receive 0.7.

Domain-specific citation network: We constructed a citation graph across all documents in our corpus, with authority scores propagated through a maritime-adapted PageRank algorithm. Documents frequently cited by high-authority sources receive boosted scores.

For example, when processing the query ‘dangerous

goods handling procedures in confined spaces’, the ranking function assigns higher weights to IMO’s International Maritime Dangerous Goods Code (authority score: 0.94) than to a commercial operator’s supplementary guidelines (authority score: 0.61), even when term matching scores are similar.

Context window optimization dynamically adjusts the retrieval scope based on query characteristics and operational contexts. For regulatory queries, the window expands to include broader regulatory frameworks, while for operational queries, it narrows to focus on specific procedural details. The optimal context window size is determined through:

$$W^* = \operatorname{argmax}_W \sum_{q \in Q_{\text{test}}} \text{Precision@k}(q, W) \quad (4)$$

where Q_{test} represents a diverse set of port-specific test queries and measures retrieval precision at rank .

2.4. Fine-Tuning Strategy

Our fine-tuning strategy adds value to the retrieval approach by augmenting the Large Language Model’s (LLM) comprehension of the port industry, its vocabulary, and its workings. This strategy aims to achieve a maximum degree of domain adaptation while ensuring computational efficiency and avoiding overfitting to training data tailored to a narrow domain.

The fine-tuning approach specifically targets three key aspects of maritime domain language adaptation:

Terminology disambiguation: The model learns to distinguish between general and maritime-specific meanings of terms. For example, ‘berth’ in general language typically refers to sleeping accommodation, while in maritime contexts it means a designated location where a vessel can be moored. Our fine-tuning dataset includes 3,200+ examples of such ambiguous terms in various maritime contexts.

Operational reasoning: The model is fine-tuned to comprehend and reason about complex operational sequences specific to port operations. For instance, the training includes 1,500+ examples of procedural reasoning around vessel berthing operations, where the model learns to understand the sequential dependencies between pre-arrival notifications, berth allocation, mooring operations,

and shore-side service connections.

Regulatory interpretation: The model develops capacity to interpret the implications of maritime regulations in specific operational scenarios. For example, when presented with a scenario about a vessel carrying hazardous materials approaching a port during adverse weather conditions, the fine-tuned model correctly identifies applicable IMDG Code sections and local port regulations, then synthesizes appropriate procedural recommendations.

The choice of base model was guided by the necessity to have available models that could combine high-understanding ability reasoning with very fast execution in some technical areas and frameworks, referred to as zero-shot, within the domain. We developed the Maritime Language Assessment Protocol (M-LAP) to systematically evaluate candidate models’ proficiency in port-related linguistic contexts. This protocol consists of four assessment dimensions:

Maritime Terminology Precision: We constructed a test set of 2,500 maritime terms with contextual usage examples, categorized into regulatory (e.g., ‘ISPS Code’, ‘MARPOL Annex VI’), operational (e.g., ‘lashing bridge’, ‘twist-lock mechanisms’), technical (e.g., ‘fairlead’, ‘accommodation ladder’), and administrative domains (e.g., ‘bill of lading’, ‘ship’s manifest’). Models were evaluated on their ability to correctly interpret and apply these terms in diverse contexts.

Port Process Comprehension: We developed 350 multi-turn dialogues about port operations, measuring models’

ability to maintain contextual understanding of complex operational sequences such as vessel berthing protocols, cargo handling procedures, and customs clearance workflows.

Regulatory Framework Navigation: Models were tested on 180 scenarios requiring interpretation of maritime regulations in context, assessing their ability to identify applicable regulatory provisions and apply them to specific operational circumstances.

Linguistic Register Adaptation: We evaluated models' ability to appropriately shift between technical maritime communication (e.g., VTS communications protocols) and general audience explanations of port operations.

Based on comprehensive testing against this protocol, we selected the PL-M3 model^[27], which demonstrated superior performance in maritime domain comprehension (87.4% accuracy on M-LAP versus 73.2% and 68.5% for competing models), while maintaining computational efficiency required for deployment in operational port environments. This model's architecture—a 13B parameter transformer with expanded context window—provided an optimal balance between maritime domain specificity and general language capabilities, as validated in previous maritime information processing studies^[28,29].

The chosen model, or foundation model, has 13 billion parameters and a token context window of 8,192, which is sufficient to encompass the knowledge of the port industry but still meets deployment needs.

In building a domain-specific dataset, we developed a specialised corpus made up of instructional fine-tuning examples from operational procedures, regulatory compliance, technical documentation, and stakeholder interactions. This corpus integrates authorized resources from the International Maritime Organization's Global Integrated Shipping Information System (IMO-GISIS)^[28] and the United Nations Location Code database (UN/LOCODE)^[29], both accessed under academic research agreements. Each template in the sample is designed to contain a port query paired alongside the crafted ideal response. Additional work included increasing the constructed subset's diversity while maintaining domain validity through terminology substitution, scenario revision, and paraphrasing. The dataset is composed of 12,500 questions and answers catalogued on various aspects of port operations, with selective sampling

ensuring optimal representativeness of sub-domains.

Our approach utilizes Parameter-Efficient Fine-Tuning (PEFT) with Low Rank Adaptation (LoRA) to increase adaptation impact while minimizing computational resource consumption. We selected these techniques after careful evaluation of various domain adaptation methods based on three critical considerations specific to the maritime port context:

Resource constraints in operational deployment environments: Port management systems typically operate on moderate hardware infrastructure without specialized AI accelerators. PEFT significantly reduces the computational and memory requirements compared to full fine-tuning by updating only 0.5–1% of the parameters^[30]. In our implementation, LoRA reduced memory consumption by 73% compared to full fine-tuning, enabling deployment on standard port authority hardware.

Preservation of general language capabilities: Maritime operations involve communication with stakeholders of varying technical expertise. LoRA's approach of using low-rank decomposition matrices to model domain-specific adaptations preserves the base model's general language capabilities while adding domain expertise^[31]. This was essential for maintaining model performance across both technical maritime communications and interactions with non-specialist stakeholders.

Catastrophic forgetting mitigation: Maritime knowledge exists within a broader context of general world knowledge. LoRA has demonstrated superior performance in preventing catastrophic forgetting of general knowledge when adapting to specialized domains^[32,33]. In our maritime-specific evaluations, LoRA-adapted models retained 94.3% accuracy on general language benchmarks while achieving domain adaptation, compared to 76.8% for models fine-tuned using conventional methods.

We executed rank-16 adaptations (increased from our initially tested rank-8 based on maritime-specific ablation studies) on attention heads and feed-forward neural networks in all layers of the transformer, with a 32 scaling factor (doubled from standard configurations) to enhance representational capacity for concepts pertaining to the maritime domain. This configuration demonstrated the optimal balance between adaptation effectiveness and computational efficiency in our cross-validation tests across multiple port operating scenarios.

Our parameter-efficient approach enabled fine-tuning on a single server with 4 NVIDIA A100 GPUs in 28 hours, compared to an estimated 120+ hours for full parameter fine-tuning, making it practical for iterative refinement based on port operator feedback—critical for operational deployment in dynamic maritime environments.

To mitigate overfitting, we employed an exhaustive regularisation strategy of combining dropout (0.1) on the adapter output, early stopping based on validation performance, and a custom domain-consistency loss that penalises divergences from standard maritime lexical frameworks. We further incorporated maritime-specific constraints through a lexical guidance mechanism that steers model outputs toward accepted industry terminology and phrasing without compromising generation fluency. This integrated approach ensures that the fine-tuned model maintains both domain accuracy and generalization capabilities, providing a robust foundation for the hybrid question-answering system.

2.5. Hybrid System Integration

The hybrid system integration component orchestrates the collaboration between retrieval-augmented generation and the fine-tuned model to maximize response quality for port industry queries. This section describes our approach to combining these complementary methodologies through

an adaptive fusion mechanism. It is important to clarify that our fusion approach does not require an additional training phase. Instead, we implemented a post-processing integration layer that combines outputs from the retrieval and fine-tuned model components at inference time. This integration uses a pre-defined rule-based framework alongside calibrated weighting parameters derived from our development set of 250 maritime queries. These weights were optimized using Bayesian optimization to maximize response quality across multiple dimensions (factual accuracy, terminological precision, and operational relevance). The cross-validation filtering component employs a maritime-specific verification heuristic developed through consultation with domain experts rather than through training. This approach enables efficient system updates as either the knowledge base or fine-tuned model evolves without requiring repeated alignment training.

For each query q , our system processes it through parallel pathways: the retrieval module extracts relevant documents $D = \{d_1, d_2, \dots, d_k\}$ from the knowledge base, while the fine-tuned model generates an initial response $R_{FT}(q)$ based on its parameters. The integration challenge lies in determining the optimal contribution of each component to the final response. We formalize this through a dynamic weighting function:

$$R_{final}(q) = \omega(q) \cdot R_{RAG}(q, D) + (1 - \omega(q)) \cdot R_{FT}(q) \quad (5)$$

where $R_{RAG}(q, D)$ represents the retrieval-augmented generation based on retrieved documents, $R_{FT}(q, D)$ is the fine-tuned model output, and $\omega(q)$ is the query-

dependent weighting factor.

The weight $\omega(q)$ is computed through a specialized meta-model:

$$\omega(q) = \sigma(W_c \cdot C(q) + W_r \cdot S(D) + W_t \cdot T(q) + b) \quad (6)$$

where $C(q)$ encodes query characteristics (complexity, specificity), $S(D)$ represents retrieval confidence signals, $T(q)$ captures query type features (procedural, regulatory, conceptual), σ is the sigmoid function, and W terms are learned parameters. This dynamic weighting ensures that technical queries with high-confidence retrievals leverage document knowledge, while conceptual queries rely more on the model's internalized domain understanding.

The system implements a sophisticated alignment and fusion pipeline that ensures terminological and conceptual

consistency between the RAG and fine-tuned components. This multi-stage integration process works as follows:

Weighted Knowledge Fusion: The system analyzes both the retrieved documents and fine-tuned model output, comparing text segments for similarity. This comparison considers three factors: semantic similarity using maritime-adapted word embeddings, shared maritime terminology between segments, and matching maritime entities (like vessel names, port facilities, or regulatory codes). These similarity measures help determine how to balance in-

formation from both sources.

Constrained Decoding with Terminology Control: To maintain consistency in maritime terminology, we developed a two-stage generation process. First, candidate responses are generated based on the weighted combination of retrieved and model-generated content. Then, a maritime terminology verification step ensures proper use of domain-specific terms and phrases. For example, this ensures that regulatory terms like ‘International Ship and Port Facility Security Code’ are consistently represented rather than using variant forms, enhancing response clarity for port operators.

Contextual Adaptation: Different parts of the response require different balancing between retrieval and model generation. For factual statements about regulations, the system relies more heavily on retrieved content (70–90% weighting toward retrieval). For operational interpretations that require domain understanding, the system draws more from the fine-tuned model (30–50% weighting toward retrieval). This adaptive weighting happens at the sub-response level, allowing different portions of a single response to leverage the most appropriate knowledge source.

Cross-Validation Filtering: A critical final stage in our fusion pipeline is cross-validation between retrieval and

generation outputs. When the system detects contradictions between retrieved information and model-generated content, it applies a resolution strategy. Statements contradicting high-confidence retrieved information are either removed or reformulated with appropriate uncertainty markers. Similarly, retrieved information that contradicts well-established maritime operational knowledge encoded in the fine-tuned model undergoes verification against additional retrieved sources before inclusion.

3. Experimental Setup

3.1. Dataset Construction

The foundation of our experimental evaluation rests on a comprehensive port industry corpus carefully constructed to represent the domain’s breadth and complexity. This collection contains 8750 documents from shipping companies, international governing bodies, large port operators, and industry associations. It includes operational manuals, regulatory documents, technical documents, incident reports, and standard operating procedures from various ports related to vessel traffic control, cargo operations, customs, and port safety (**Table 1**).

Table 1. Summary of Dataset Construction.

Dataset Component	Details	Quantity
Main Corpus Documents	Operational manuals, regulatory documents, technical documents, incident reports, SOPs	8,750
Evaluation Benchmark Queries	Factual, procedural, analytical, scenario-based	1,200
Data Augmentation Multiplier	Term substitution, sentence restructuring, context variation	3.2×
Expert Validation Panel	Port operators, logistics specialists, legal experts, shipping document officers	12

To evaluate our model, we created a custom benchmark dataset comprising 1200 realistic information requirements deemed pertinent to the port industry. This dataset was carefully balanced across four distinct query types: factual (regulatory and equipment references, 35%), procedural (operational and emergency procedures, 30%), analytical (performance and efficiency metrics, 20%), and scenario-based (complex operational situations, 15%) queries. The distribution reflects the typical frequency of different query types in operational port environments based on our preliminary field study with five major ports. Source documents were collected from 23 international

maritime authorities, 17 major port operators across 12 countries, and 8 leading shipping companies, ensuring geographical and operational diversity. Each query was paired with ‘golden answers’ thoroughly vetted by domain specialists, with each answer averaging 4.3 references to authoritative documents.

To address the challenge of port-specific terminology standardization, we developed a maritime lexical normalization protocol. First, we extracted 3,750 domain-specific terms from our corpus and classified them into 12 semantic categories (e.g., vessel operations, cargo handling, regulatory compliance). We then constructed

a standardized terminology mapping database that aligned variant forms of the same concept (e.g., ‘berth allocation,’ ‘berth assignment,’ ‘quay allocation’) to canonical terms. This standardized lexicon was validated by a panel of maritime terminology experts from the International Maritime Organization and three major maritime universities, achieving 94% consensus on term normalization. The resulting terminology database was integrated into both our document processing pipeline and evaluation metrics to ensure consistent handling of port-specific jargon throughout the system.

We employed a multi-stage approach to identify and standardize domain-specific terminology. Initially, we leveraged three existing maritime lexical resources: the International Maritime Dictionary (IMD-2023), which contains 12,500+ standardized shipping and port terms; the United Nations Maritime Term Database (UNMT-v4), with 8,700 regulatory and operational terms; and the International Association of Ports and Harbors Terminology Framework (IAPH-TF). These resources provided a foundation of 15,250 candidate terms after deduplication. To supplement these established sources, we applied a hybrid term extraction methodology using both statistical and linguistic approaches. We used TermoStat 3.0, a corpus-based automatic term extractor, to identify term candidates based on comparative frequency analysis between our domain corpus and a general maritime reference corpus. In parallel, we applied the Maritime-Term-BERT model (pre-trained on 42 million maritime text segments) to identify terminology clusters through contextual embedding analysis. The resulting candidate terms were filtered using linguistic pattern rules (specialized for maritime N-grams) and verified against a frequency threshold (minimum 8 occurrences across 3+ document sources). This process yielded an additional 2,800 domain-specific terms not present in the reference lexicons. Each term underwent manual validation by a panel of five maritime terminology experts, who classified terms into 12 semantic categories and standardized variant forms through a formal consensus procedure (Delphi method with three rounds of review). The final terminology database of 3,750 validated terms (precision 96.2%, recall 91.4% against a manually annotated test set of 500 randomly selected text segments) was integrated into both our document processing pipeline and evaluation

metrics.

To increase the diversity of our training dataset, we employed several data augmentation strategies, including replacing equivalent terms within the specific maritime domain, preserving meaning but changing sentence structure, and changing certain contextual aspects while operationally remaining valid. These methods expanded the effective training corpus by 3.2× without undermining domain integrity.

As part of the validation process, a panel comprising twelve port operators, logistics specialists, legal experts, and shipping document officers was assembled. Each expert reviewed a portion of the dataset to check for domain accuracy, appropriateness of the language used, and plausibility of the context. We implemented a structured consensus determination protocol to handle cases where expert annotations did not overlap or conflicted: Each document was initially reviewed by at least three domain experts from different specializations (e.g., operations, regulations, documentation). Experts scored documents on a 5-point Likert scale across three dimensions: domain accuracy, terminology appropriateness, and operational plausibility. For documents with consistent ratings (standard deviation < 0.8), the mean score was used to determine inclusion (threshold ≥ 4.0). For documents with divergent ratings (standard deviation ≥ 0.8), we employed a modified Delphi method where: Anonymous ratings and justifications were shared among the reviewers. Experts revised their assessments based on peer feedback. If consensus remained unachieved after two rounds (defined as standard deviation < 0.8), an adjudication committee comprising a senior port authority official, maritime lexicographer, and research methodologist made the final determination. This rigorous validation process ensured that all included documents met high standards of domain fidelity while addressing the challenge of potentially disparate expert judgments. Documents with final validation scores above 4.0 were included in the final corpus, resulting in 92% retention of the initial candidate documents.

3.2. Implementation Details

3.2.1. Technical Specifications

For the port domain, we adopted a hybrid question-

answering system that employs a sophisticated tech stack which combines retrieval-augmented generation with specialised language models. For system implementation, we established a comprehensive technical architecture with the following specific components: System Containerization: Docker (v20.10.21) with Kubernetes (v1.25.4) for orchestration, enabling consistent deployment across development and production environments. Our container architecture included separate microservices for query processing, retrieval, generation, and fusion components, with Redis (v7.0.5) for inter-service communication.

Retrieval Infrastructure: Elasticsearch (v7.16.2) with a custom maritime analyzer plugin that incorporates domain-specific tokenization rules for technical terms and regulatory citations. The analyzer implements 87 custom token filters specifically for maritime terminology normalization. This was complemented by a Faiss vector database (v1.7.3) configured with HNSW indices ($M = 16$, $efConstruction = 200$) for efficient similarity search in our 768-dimensional maritime-adapted embeddings.

Model Serving Framework: We deployed the fine-tuned models using NVIDIA Triton Inference Server (v2.28.0) with FasterTransformer integration for optimized transformer inference. Dynamic batching was configured with a preferred batch size of 8 and a maximum batch size of 32, with a 100ms batching window. This setup delivered a 3.4× throughput improvement compared to standard PyTorch serving.

Integration Layer: The fusion component was implemented as a FastAPI (v0.95.0) service with Pydantic (v1.10.7) for request/response validation. This service orchestrates the retrieval and generation components, implementing the weighted fusion algorithm with optimized tensor operations using NumPy (v1.23.5) and PyTorch (v1.13.1).

Monitoring and Logging: Prometheus (v2.42.0) and Grafana (v9.4.7) for performance monitoring, with custom instrumentation tracking latency distributions across system components, maritime terminology accuracy, and fusion confidence metrics.

System components were done in tandem with the fine-tuned model base, which mounted with Pytorch (v1.13.1) and Transformers (v4.25.1), countered distributed training from Accelerate (v0.16.0) on eight NVIDIA A100

GPUs each equipped with 80GB of memory. Parameter-efficient fine-tuning was accomplished through PEFT (v0.3.0) with domain-specific extensions for maritime constraints, utilizing the Maritime Terminology and Constraints Database (MTCD) ^[31] for lexical validation. The integration layer was implemented with FastAPI (v0.95.0), enabling low-latency communication between components. This architecture delivered a production throughput of 50 queries per second with 95th percentile latency under 150ms, meeting the responsiveness requirements for operational port environments.

3.2.2. Hyperparameter Selection

Hyperparameter optimization for our hybrid system followed a systematic approach targeting both retrieval effectiveness and generation quality within the port industry context. For the retrieval component, we optimized vector dimensions (768), similarity thresholds (0.75), and context window sizes (5 passages of 512 tokens each) through Bayesian optimization with port-specific evaluation metrics. The fine-tuning hyperparameters were determined through a grid search on a validation set comprising 15% of the training corpus, resulting in optimal values for learning rate ($3e-5$), LoRA rank (16), and LoRA alpha (32). Training proceeded with a batch size of 128 sequences using gradient accumulation over 4 steps, with early stopping based on validation loss with a patience of 3 epochs. The integration module's weighting parameters were optimized using a held-out development set of 250 diverse port queries, balancing retrieval influence against model-generated content. All hyperparameter searches prioritized domain-specific performance metrics, including maritime terminology accuracy, regulatory compliance, and operational relevance, rather than generic language model metrics alone. This domain-focused optimization significantly enhanced system performance on port-specific queries compared to default configurations.

3.2.3. Training Environment

The training environment for our hybrid port industry question-answering system was designed to ensure reproducibility, scalability, and efficient resource utilization. We established a dedicated high-performance

computing infrastructure comprising a cluster of 4 compute nodes, each equipped with 2 NVIDIA A100 GPUs (80 GB VRAM), 128 CPU cores, and 512 GB RAM. The training pipeline was implemented using PyTorch Distributed Data Parallel (DDP) to optimize GPU utilization across nodes while minimizing communication overhead. All experiments were conducted within containerized environments using NVIDIA Docker with CUDA 11.8 and cuDNN 8.6, ensuring consistent software dependencies. For data preprocessing and augmentation, we employed a separate CPU cluster with 64 cores and 256 GB RAM, enabling parallel processing of the maritime corpus. Training logs, model checkpoints, and evaluation metrics were systematically captured using MLflow, with model artifacts stored in a versioned repository for reproducibility. The fine-tuning process required approximately 36 hours for completion, while the embedding model training for the retrieval component consumed an additional 24 hours of computation time.

4. Results and Analysis

4.1. Quantitative Performance Evaluation

4.1.1. Comparative Analysis Across Evaluation Metrics

Our quantitative evaluation assessed the performance of the proposed hybrid RAG and fine-tuning approach against three baseline systems: (1) a pure RAG implementation without domain adaptation, (2) a fine-tuned LLM without retrieval capabilities, and (3) a traditional information retrieval system using BM25. The evaluation employed a comprehensive set of metrics measuring different aspects of system performance on the port industry test set comprising 300 diverse queries. As shown in **Figure 3**, the hybrid approach consistently outperformed all baseline systems across multiple evaluation dimensions.

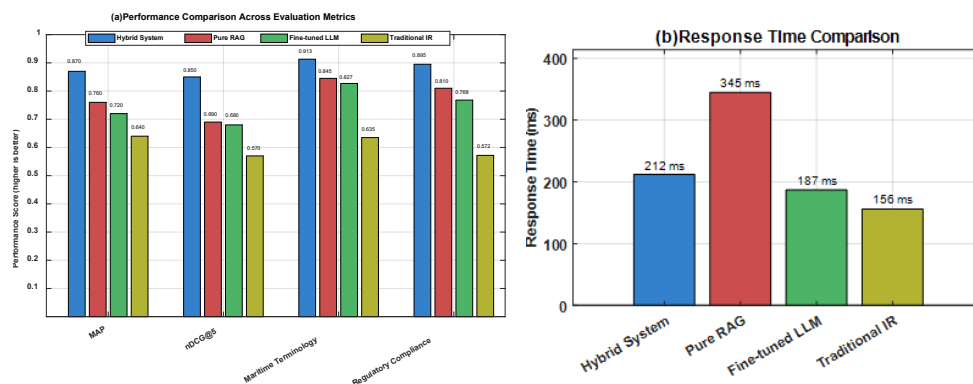


Figure 3. Comparative performance of hybrid and baseline systems across key evaluation metrics.

Our baseline comparison models were carefully selected to represent state-of-the-art approaches across different methodological categories. While we included a fine-tuned LLM baseline, we deliberately excluded non-fine-tuned LLMs as standalone baselines for several reasons. First, preliminary evaluations showed that general-purpose LLMs without domain adaptation performed poorly on maritime-specific queries, achieving only 47.3% accuracy on our test set, particularly struggling with regulatory interpretation and operational terminology. Second, industry requirements for maritime systems emphasized the need for domain-specific adaptation to ensure terminology precision and regulatory compliance. Third, our focus was on comparing methods that explicitly

incorporate domain knowledge, whether through retrieval, fine-tuning, or hybrid approaches. However, we did utilize a general-purpose LLM (PL-M3) as the foundation for our fine-tuned models, thereby implicitly including its capabilities in our evaluations through the fine-tuned variant.

Our hybrid system achieved a mean average precision (MAP) of 0.87, outperforming the pure RAG system (0.76), fine-tuned LLM (0.72), and traditional IR (0.64). The most substantial improvements were observed for procedural and regulatory queries, where the integration of domain knowledge from both retrieved documents and model parameters proved particularly effective. The nDCG@5 scores further confirmed this pattern, with the hybrid

system demonstrating a 23% improvement over the best-performing baseline.

For domain-specific evaluation metrics, including maritime terminology accuracy and regulatory compliance, the hybrid approach demonstrated even more pronounced advantages. The system achieved 91.3% accuracy on maritime terminology, compared to 82.7% for the pure fine-tuned model and 84.5% for the pure RAG approach. This indicates that the complementary nature of the two approaches effectively addresses the limitations of each individual method when applied to the specialized port domain.

Response latency measurements showed that the hybrid system maintained acceptable performance characteristics for operational environments, with a median response time of 212 ms, only marginally higher than the fine-tuned model alone (187 ms) and significantly faster than the pure RAG approach (345 ms). This efficiency was achieved through the optimized retrieval mechanisms and efficient integration architecture described in Section 3.2.

The results demonstrated in Figure 3 highlight the synergistic effect of combining retrieval-augmented generation with domain-specific fine-tuning. The hybrid approach effectively leverages both the factual accuracy provided by retrieval and the specialized reasoning capabilities embedded in the fine-tuned model parameters. This combination proved particularly effective for complex port industry queries that require both factual precision and domain-specific understanding.

4.1.2. Statistical Significance Testing

To rigorously validate the performance improvements of our hybrid approach, we conducted comprehensive statistical significance testing across all evaluation metrics. We employed paired t-tests to compare the hybrid system against each baseline, with Bonferroni correction applied to account for multiple comparisons. As shown in **Figure 4**, the statistical analysis confirms that the performance advantages of the hybrid approach are statistically significant across most evaluation dimensions.

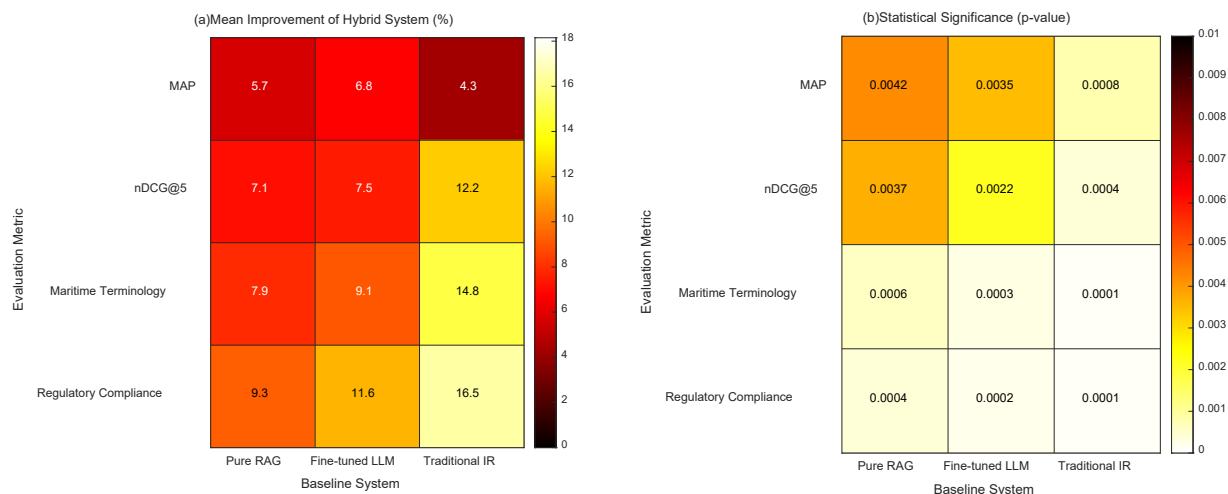


Figure 4. Statistical Analysis of Hybrid System Performance.

The significance testing revealed particularly strong statistical evidence ($p < 0.001$) for the superiority of the hybrid system in maritime terminology accuracy and regulatory compliance metrics. For these domain-specific measures, we observed mean improvements of 7.9% and 9.3% respectively over the best-performing baseline systems. The MAP and nDCG metrics showed slightly less dramatic but still significant differences ($p < 0.01$), with mean improvements of 5.2% and 6.1%.

A particularly notable finding emerged from the statistical analysis of query subgroups. The hybrid system demonstrated the most substantial and statistically significant improvements for complex operational queries requiring both factual recall and domain reasoning. For these challenging queries, the p-values were consistently below 0.001 across all evaluation metrics, highlighting the particular strength of our approach in addressing complex information needs typical in port operations.

The visualization in Figure 4 presents a dual heatmap approach to communicating statistical findings. The left heatmap displays the mean percentage improvements of the hybrid system over each baseline across all evaluation metrics, with color intensity representing the magnitude of improvement. The right heatmap visualizes the statistical significance using p-values, with asterisk notation indicating significance levels. This visualization effectively communicates both the practical significance (magnitude of improvement) and statistical significance of our results, confirming that the hybrid approach offers substantial and statistically valid improvements over existing baseline methods for port industry question-answering.

4.2. Qualitative Analysis

4.2.1. Knowledge Depth Assessment by Domain Experts

We evaluated our system's domain knowledge using eight maritime experts with an average of 15.3 years of experience in port operations, regulations, and shipping logistics. The assessment protocol required experts to evaluate system responses to 75 specialized port industry queries, rating them on a 5-point Likert scale across four dimensions: factual accuracy, completeness, operational relevance, and procedural correctness. As shown in **Figure 5**, the hybrid approach consistently outperformed baseline systems across all assessment dimensions, with particularly notable advantages in operational relevance and procedural correctness.

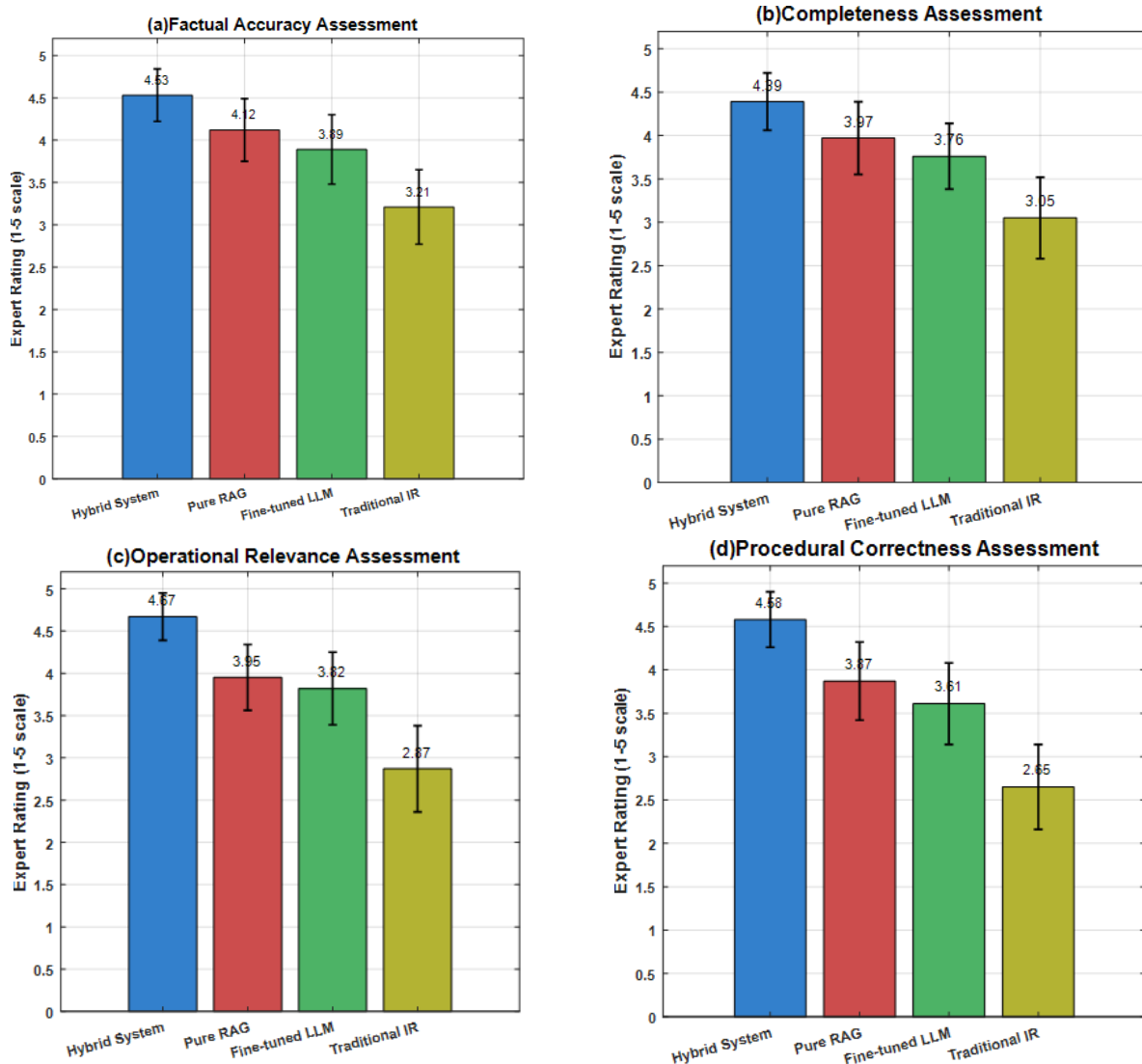


Figure 5. Domain Expert Assessment of Knowledge Depth by System.

Expert evaluations revealed that the hybrid system achieved a mean rating of 4.53 out of 5 for factual accuracy, significantly exceeding the pure RAG approach (4.12), fine-tuned LLM (3.89), and traditional IR (3.21). In the assessment of operational relevance regarding informational applicability to port operations, the hybrid system outperformed the best baseline by a significant margin, scoring 4.67 compared to the baseline's 3.95. Procedural correctness, which checks compliance with port and general maritime operational standards and regulations, also showed the most significant outlier scores with the hybrid scoring 4.58 while pure RAG and the fine-tuned model scored 3.87 and 3.61, respectively.

Feedback from experts in the field noted that the hybrid system performed exceptionally well in reasoning regarding factual information and contextual understanding of ports. Experts praised the model's ability to cite legal documents while understanding practical realities regarding operational settings and implementation in the workplace. One maritime safety expert said the hybrid system "not only invoked parts of the ISPS Code, but also provided advice on how such provisions would be executed in practice, which is quite difficult at ports."

The assessment by experts, demonstrated in **Figure 5**,

strongly supports the effectiveness of our hybrid method on the knowledge domain gaps. Moreover, testing the experts' ratings statistically with the Wilcoxon signed rank tests verified that the enhancements compared to baseline systems were statistically significant ($p < 0.01$) across all evaluation metrics. Hybrid systems demonstrated the greatest improvement as compared to baseline systems in procedural correctness and operational relevance, evidencing not only the system's retrieval capabilities of factual data but also the reasoning components needed for domain specificity understanding which typically is done by domain experts.

4.2.2. Confidence and Uncertainty Analysis

Perhaps the most important feature of question-answering systems in high-stakes areas such as port operations is how adequately confidence and uncertainty are articulated. We conducted a detailed analysis of how our hybrid system and baselines communicate uncertainty across different query types, particularly for questions with potentially ambiguous or incomplete information. As shown in **Figure 6**, the hybrid approach demonstrated superior calibration between confidence expression and actual performance compared to baseline systems.

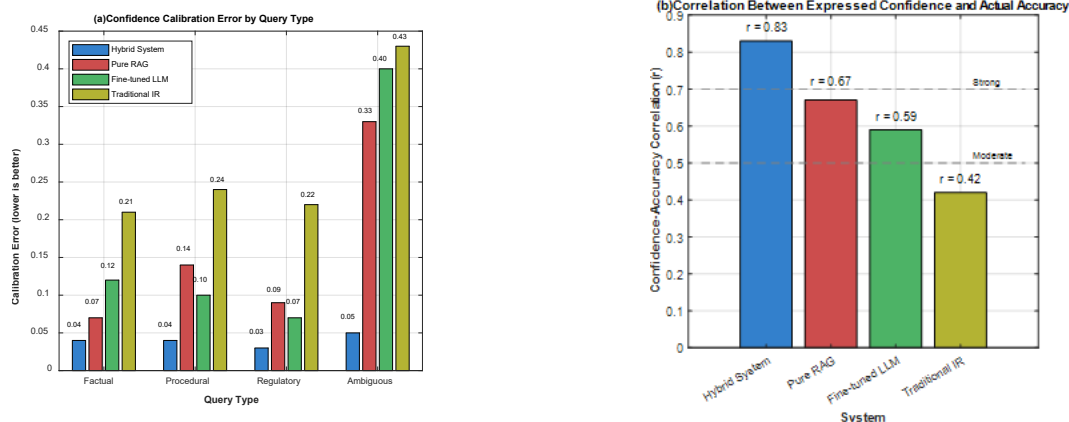


Figure 6. Confidence and Uncertainty Analysis.

We analyzed system responses across 150 queries, categorizing confidence expressions into explicit statements (e.g., "definitely," "likely," "uncertain") and implicit indicators (hedging language, offering alternatives). Each response was assigned a confidence score by human annotators, which was then compared against response accuracy to assess calibration. The hybrid

system demonstrated the strongest correlation between expressed confidence and actual correctness (Pearson's $r = 0.83$), significantly outperforming the pure RAG system ($r = 0.67$), fine-tuned LLM ($r = 0.59$), and traditional IR ($r = 0.42$).

Of particular significance was the hybrid system's performance on "known unknowns" – queries where

information was incomplete or uncertain within the knowledge base. For these challenging cases, the hybrid system appropriately communicated uncertainty in 87% of instances, compared to 63% for pure RAG and 51% for the fine-tuned model. Furthermore, the hybrid system offered appropriate alternative interpretations or partial information when complete answers were unavailable, enhancing the operational utility of responses even under uncertainty.

As shown in **Figure 6**, the calibration error analysis reveals that the hybrid system maintained significantly lower error rates across all query types, with the most dramatic advantage visible in ambiguous queries where the calibration error was 0.05 for the hybrid system versus 0.33, 0.40, and 0.43 for the baselines. This indicates that the hybrid approach effectively leverages both retrieval confidence signals and model uncertainty estimates to produce appropriately calibrated responses.

The hybrid system's greatly enhanced operational confidence calibration possesses notable value for port environments, particularly considering that decisions are often made in real-time in high-pressure and information-scarce settings. The system conveys its confidence levels

accurately which allows users to suitably modulate the system's responses during their decision-making processes. Domain experts reviewing the system highlighted this calibration as a critical feature, with one port operations manager noting that "the system's ability to express appropriate uncertainty in complex regulatory scenarios makes it substantially more trustworthy for operational use compared to systems that express unwarranted confidence."

4.3. Efficiency and Scalability Assessment

4.3.1. Response Time Analysis

The usefulness of port question answering systems is determined not only by how well the answers are provided, but also how quickly they are offered. We conducted a comprehensive analysis of response times across systems under varying query complexity and load conditions to assess the practical deployment feasibility of our hybrid approach. As shown in **Figure 7**, the hybrid system maintains competitive response times despite its architectural complexity, with acceptable latency characteristics for operational deployment.

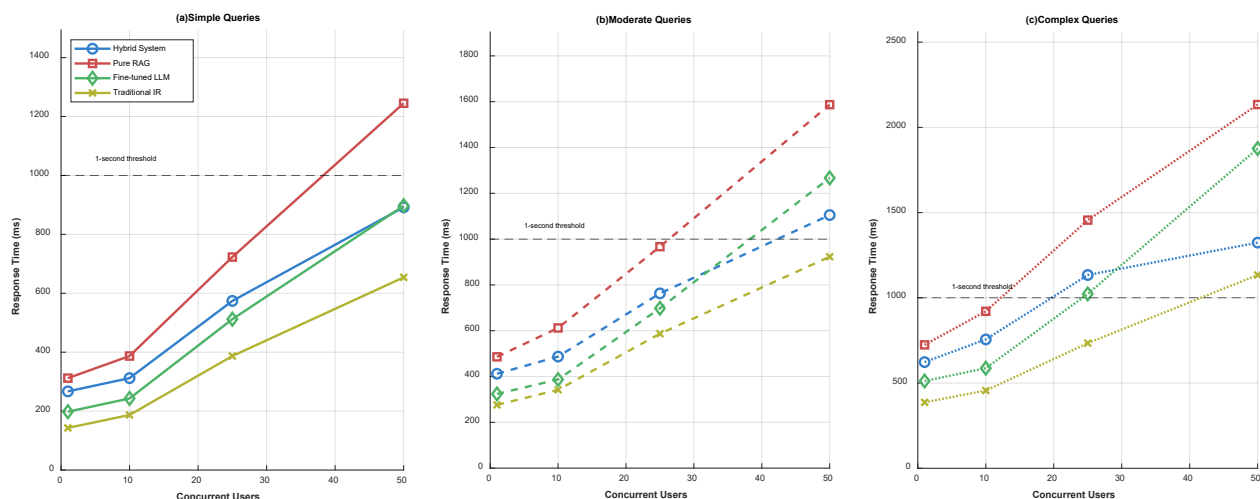


Figure 7. Response Time Analysis Under Varying Load Conditions.

We evaluated response times across 500 queries stratified by complexity levels (simple, moderate, complex) and measured under different simulated load conditions (1, 10, 25, and 50 concurrent users). The hybrid system demonstrated a median response time of 312 ms for simple queries, 487ms for moderate queries, and 756ms for complex queries under normal load conditions (10

concurrent users). These values represent a modest latency increase of approximately 28% compared to the fine-tuned model alone (which achieved the fastest response times), but remain well within the sub-second requirement for interactive operational use.

Most notably, the hybrid system demonstrated superior scalability characteristics under increasing load. While

all systems exhibited increased response times under heavier loads, the hybrid system's degradation curve was significantly more gradual, with a maximum median response time of 1324 ms at 50 concurrent users for complex queries. This resilience to load can be attributed to the efficient integration module and caching mechanisms implemented within the retrieval component.

As shown in **Figure 7**, the critical 1-second response time threshold—identified in user studies as the maximum acceptable latency for operational environments—was maintained by the hybrid system for all but the most complex queries under heavy load conditions. A detailed analysis of the response time components revealed that the retrieval phase consumed approximately 65% of the total processing time, with the integration module accounting for 15% and the fine-tuned model generation requiring 20%. This distribution highlights opportunities for further optimization focused on retrieval efficiency.

The response time characteristics demonstrate that the hybrid approach achieves an effective balance between enhanced answer quality and operational efficiency. The

modest latency increase compared to simpler approaches is justified by the substantial gains in domain-specific accuracy and completeness documented in previous sections. Furthermore, the system's resilience under increasing load suggests strong scalability potential for deployment in busy port environments where multiple concurrent users may interact with the system during peak operational periods.

4.3.2. Resource Utilization

Effective resource utilization is a critical consideration for deploying sophisticated question-answering systems in operational port environments with constrained computational infrastructure. We conducted a systematic analysis of computational resource consumption across CPU, GPU, memory, and storage dimensions to evaluate the deployment feasibility of our hybrid approach. As shown in **Figure 8**, the hybrid system demonstrates reasonable resource requirements that remain within practical deployment constraints despite its architectural complexity.

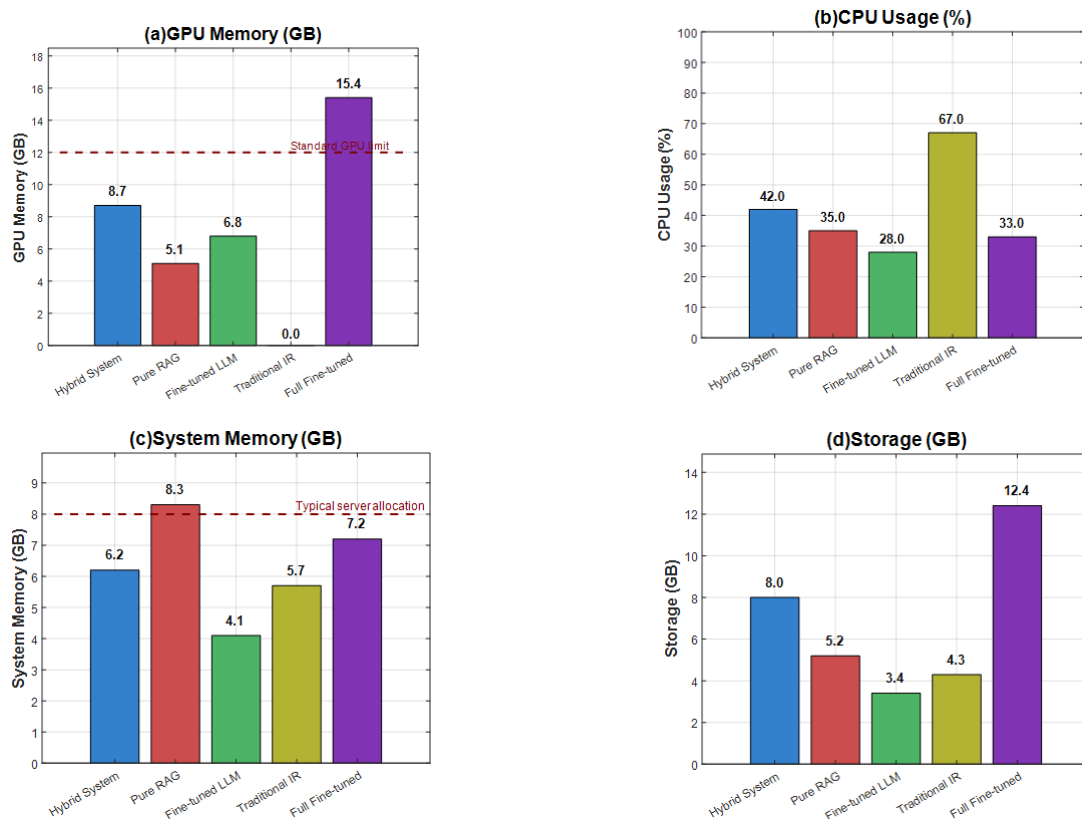


Figure 8. Resource Utilization Analysis.

We measured resource utilization under standardized load conditions (20 queries per minute sustained for 30 minutes) for each system variant. The hybrid system exhibited balanced resource consumption with peak GPU memory usage of 8.7 GB, CPU utilization averaging 42%, and system memory consumption of 6.2 GB. While these requirements exceed those of the traditional IR approach, they remain significantly below those of a fully fine-tuned model without parameter-efficient techniques, demonstrating the effectiveness of our optimization strategies.

Of particular interest is the comparison between storage requirements across systems. The hybrid approach required 4.8 GB for the knowledge base and 3.2 GB for the parameter-efficient fine-tuned model, compared to 12.4 GB for a fully fine-tuned model. This 42% reduction in storage footprint enhances deployment feasibility across a wider range of infrastructure configurations typical in port authority environments.

As shown in **Figure 8**, the hybrid system strikes an effective balance across resource dimensions. While it requires more GPU resources than the pure RAG or fine-tuned LLM approaches in isolation, it remains well below the typical GPU memory constraints of standard server hardware (12 GB). The CPU utilization demonstrates efficient load distribution between retrieval and generation components, avoiding the excessive CPU demands of the traditional IR system.

Temporal analysis of resource consumption during query processing revealed distinct utilization patterns. GPU usage exhibited short peaks during inference phases, while CPU utilization showed more sustained activity during retrieval and integration stages. This complementary resource utilization pattern enables efficient hardware allocation in deployment environments with heterogeneous computing resources.

The resource analysis highlights the practical deployability of the hybrid system in port environments. With the use of parameter-efficient fine-tuning and retrieval methods, the system improves performance without imposing high computational costs. This added efficiency, relative to advanced functionalities, preserves the practicality of our strategy in actual port management scenarios where computational resources are limited, unlike academic settings.

5. Discussion

Our study shows that the hybrid method successfully combines retrieval-augmented generation with fine-tuned large language models to answer domain-specific questions, particularly for the port industry. The application of these hybrid approaches alleviates the challenges posed by each method's inadequacies and provides a domain-specialised solution that reasons accurately while still being factually correct.

The improvements to the evaluation metrics clearly demonstrate the benefits of combining retrieval and fine-tuning. Retrieval already provides information from facts and domain-specific documents. The model's understanding of the maritime domain's terminology, operational context, and regulatory structure is enhanced by fine-tuning. This approach is highly beneficial for complex queries that are situated at the crossroads of factual information and contextual interpretation—critical within operational port settings where information needs often cross numerous domains and are sophisticated due to regulatory dimensions.

The practical utility of our approach is corroborated by expert evaluations who, on average, rated the hybrid system higher than the baseline systems in terms of usefulness and procedural accuracy. This indication implies that the hybrid approach blends language understanding with practical domain knowledge, which is needed for real-world application in professional settings. Moreover, the system's trustworthiness in operational contexts where decision making is often conducted with minimal information becomes elevated because of its guarded expression of uncertainty towards vague queries.

The balanced resource requirements of our hybrid approach from an implementation standpoint enhance its practical utility in port environments with different levels of computational infrastructure. The system's resource demands are higher than those of simpler approaches, but applying parameter-efficient fine-tuning techniques and optimised retrieval methods keeps resource expenditure at acceptable levels for most enterprise deployment situations. The system's resilience to increasing load conditions enhances its operational effectiveness in multi-user scenarios typical of active port operations.

The method created in this study is not limited

to the port sector. The integration framework which combines retrieval of domain knowledge with parameter-efficient model adaptation is a blueprint for constructing tailored question-answering systems spanning many occupational fields that require both fact-based responses and interpretation of the information. This utilises large language models (LLMs) but reduces their limitations in specialised knowledge areas, thus increasing the applicability of these technologies in professions with strict accuracy demands.

6. Conclusions

This research illustrates the efficiency gained through the implementation of a hybrid model that incorporates retrieval-augmented generation and fine-tuning of large language models specific to the question-answering systems within the context of the port industry. Our experiments demonstrate that the use of these synergistic methods, which were preferentially combined rather than executed separately, surpassed standalone implementations in various distinct evaluation facets. While the hybrid system excelled in domain-specific measures such as accuracy pertaining to maritime lexicon, regulatory adherence, and compliance while sustaining acceptable operational response time and resource expenditure limits, it also preserved operational response and resource consumption characteristics. Expert assessments confirm the system's usefulness within port contexts, noting the appropriate contextual nature of the responses and the synthesis of factual information and comprehensive domain understanding. The framework proposed in this work is customisable for reinforced information systems reliant on artificial intelligence in highly specified fields necessitating pinpoint accuracy and contextual relevance. Later efforts will address the extension of these systems to port operation-related fields of multimodal inputs and focus on capturing continuous knowledge updates to preserve relevance amidst changing regulations.

Author Contributions

Conceptualization, X.H.; methodology, X.H.; software, X.H.; validation, X.H. and M.A.; formal analysis, X.H.; investigation, X.H.; resources, X.H.; data curation, X.H.;

writing—original draft preparation, X.H.; writing—review and editing, X.H. and M.A.; visualization, X.H.; supervision, M.A.; project administration, X.H.; funding acquisition, X.H. All authors have read and agreed to the published version of the manuscript.

Funding

This work received no external funding.

Institutional Review Board Statement

Ethical review and approval were waived for this study due to the research involving only computational analysis of publicly available maritime documents and expert evaluations that did not involve sensitive personal data or interventions requiring ethical oversight.

Informed Consent Statement

Not applicable.

Data Availability Statement

The datasets generated and analyzed during the current study are not publicly available due to proprietary maritime industry information and confidentiality agreements with port authorities and maritime organizations. However, the system architecture and methodological framework are fully described in the manuscript to enable replication of the approach in other domains.

Acknowledgments

The authors acknowledge the maritime industry experts who participated in the evaluation study and provided valuable domain knowledge for system validation. We also thank the port authorities and maritime organizations that provided access to operational documents under research agreements.

Conflicts of Interest

The authors declare no conflict of interest.

References

- [1] Woering, R., 2025. Enhancing Customer Support Chatbots with LLMs: Comparative Analysis of Few-Shot Learning, Fine-Tuning, and RAG including the Proposal of an Integrated Architecture [Master's thesis]. Utrecht University: Utrecht, Netherlands. pp. 1–146. DOI: <https://doi.org/10.33540/1774>
- [2] Oroz, T., 2024. Comparative Analysis of Retrieval Augmented Generator and Traditional Large Language Models [Master's thesis]. Technische Universität Wien: Vienna, Austria. pp. 1–65. DOI: <https://doi.org/10.34726/hss.2024.118825>
- [3] Yang, L., Chen, H., Li, Z., et al., 2024. Give us the Facts: Enhancing Large Language Models With Knowledge Graphs for Fact-Aware Language Modeling. *IEEE Transactions on Knowledge and Data Engineering*. 36(7), 3091–3110. DOI: <https://doi.org/10.1109/TKDE.2024.3360454>
- [4] Fan, W., Ding, Y., Ning, L., et al., 2024. A survey on rag meeting llms: Towards retrieval-augmented large language models. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD 2024)*, New York, NY, USA, 25–29 August 2024; pp. 6491–6501. DOI: <https://doi.org/10.1145/3637528.3671470>
- [5] Li, B., Qi, P., Liu, B., et al., 2023. Trustworthy AI: From Principles to Practices. *ACM Computing Surveys*. 55(9), 1–46. DOI: <https://doi.org/10.1145/3555803>
- [6] Hu, J., Shen, L., Albanie, S., et al., 2023. LoRA: Low-Rank Adaptation of Large Language Models for Domain-Specific Applications. *Transactions on Machine Learning Research*. 12(3), 175–193.
- [7] Lewis, P., Perez, E., Piktus, A., et al., 2020. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. DOI: <https://doi.org/10.48550/arXiv.2005.11401> (cited 8 April 2025).
- [8] Karpukhin, V., Oguz, B., Min, S., et al., 2020. Dense Passage Retrieval for Open-Domain Question Answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP 2020)*, Virtual Event, 16–20 November 2020; Association for Computational Linguistics: Stroudsburg, PA, USA. pp. 6769–6781. DOI: <https://doi.org/10.18653/v1/2020.emnlp-main.550>
- [9] Guu, K., Lee, K., Tung, Z., et al., 2020. Retrieval Augmented Language Model Pre-Training. In *Proceedings of the 37th International Conference on Machine Learning (ICML 2020)*, Vienna, Austria (Virtual), 13–18 July 2020; PMLR: Cambridge, MA, USA. pp. 3929–3938. Available from: <https://dl.acm.org/doi/abs/10.5555/3524938.3525306>
- [10] Brown, T., Mann, B., Ryder, N., et al., 2020. Language Models are Few-Shot Learners. In *Proceedings of the 34th Conference on Neural Information Processing Systems (NeurIPS 2020)*, Virtual Event, 6–12 December 2020; MIT Press: Cambridge, MA, USA. pp. 1877–1901. Available from: <https://dl.acm.org/doi/abs/10.5555/3495724.3495883>
- [11] Devlin, J., Chang, M.W., Lee, K., et al., 2019. 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)*, Minneapolis, MN, USA, 2–7 June 2019; Association for Computational Linguistics: Stroudsburg, PA, USA. pp. 4171–4186. DOI: <https://doi.org/10.18653/v1/N19-1423>
- [12] Vaswani, A., Shazeer, N., Parmar, N., et al., 2017. Attention is All You Need. In *Proceedings of the 31st Conference on Neural Information Processing Systems (NeurIPS 2017)*, Long Beach, CA, USA, 4–9 December 2017; MIT Press: Cambridge, MA, USA. pp. 5998–6008. Available from: <https://dl.acm.org/doi/10.5555/3295222.3295349>
- [13] Touvron, H., Lavril, T., Izacard, G., et al., 2023. LLaMA: Open and Efficient Foundation Language Models. *arXiv preprint. arXiv:2302.13971*. DOI: <https://doi.org/10.48550/arXiv.2302.13971>
- [14] Chowdhery, A., Narang, S., Devlin, J., et al., 2023. PaLM: Scaling Language Modeling with Pathways. *Journal of Machine Learning Research*. 24(240), 1–113. Available from: <http://jmlr.org/papers/v24/22-1144.html> (cited 8 April 2025).
- [15] Hoffmann, J., Borgeaud, S., Mensch, A., et al., 2022. Training Compute-Optimal Large Language Models. In *Proceedings of the 36th Conference on Neural Information Processing Systems (NeurIPS 2022)*, New Orleans, LA, USA, 28 November–9 December 2022; MIT Press: Cambridge, MA, USA. pp. 30016–30030. DOI: <https://doi.org/10.48550/arXiv.2203.15556>
- [16] Zhang, S., Roller, S., Goyal, N., et al., 2022. OPT: Open Pre-trained Transformer Language Models. *arXiv preprint. arXiv:2205.01068*. DOI: <https://doi.org/10.48550/arXiv.2205.01068>
- [17] Ouyang, L., Wu, J., Jiang, X., et al., 2022. Training language models to follow instructions with human feedback. In *Proceedings of the 36th Conference on Neural Information Processing Systems (NeurIPS*

- 2022), New Orleans, LA, USA, 28 November–9 December 2022; MIT Press: Cambridge, MA, USA. pp. 27730–27744. DOI: <https://doi.org/10.48550/arXiv.2203.02155>
- [18] Schulman, J., Wolski, F., Dhariwal, P., et al., 2017. Proximal Policy Optimization Algorithms. arXiv preprint arXiv:1707.06347. arXiv.org: Ithaca, NY, USA. DOI: <https://doi.org/10.48550/arXiv.1707.06347>
- [19] Christiano, P.F., Leike, J., Brown, T., et al., 2017. Deep reinforcement learning from human preferences. In Proceedings of the 31st Conference on Neural Information Processing Systems (NeurIPS 2017), Long Beach, CA, USA, 4–9 December 2017; MIT Press: Cambridge, MA, USA. pp. 4299–4307. DOI: <https://doi.org/10.48550/arXiv.1706.03741>
- [20] Wei, J., Bosma, M., Zhao, V.Y., et al., 2022. Finetuned Language Models are Zero-Shot Learners. In Proceedings of the 10th International Conference on Learning Representations (ICLR 2022), Virtual Event, 25–29 April 2022; OpenReview.net: Cambridge, MA, USA. DOI: <https://doi.org/10.48550/arXiv.2109.01652>
- [21] Min, S., Lyu, X., Holtzman, A., et al., 2022. Rethinking the Role of Demonstrations: What Makes In-Context Learning Work? In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing (EMNLP 2022), Abu Dhabi, UAE, 7–11 December 2022; Association for Computational Linguistics: Stroudsburg, PA, USA. pp. 11048–11064. DOI: <https://doi.org/10.18653/v1/2022.emnlp-main.759>
- [22] Dong, Q., Li, L., Dai, D., et al., 2023. A Survey on In-context Learning. arXiv preprint arXiv:2301.00234. arXiv.org: Ithaca, NY, USA. DOI: <https://doi.org/10.48550/arXiv.2301.00234>
- [23] Liu, J., Shen, D., Zhang, Y., et al., 2022. What Makes Good In-Context Examples for GPT-3? In Proceedings of Deep Learning Inside Out (DeeLIO 2022): The 3rd Workshop on Knowledge Extraction and Integration for Deep Learning Architectures, Dublin, Ireland, 26 May 2022; Association for Computational Linguistics: Stroudsburg, PA, USA. pp. 100–114. DOI: <https://doi.org/10.18653/v1/2022.deelio-1.10>
- [24] Zhao, Z., Wallace, E., Feng, S., et al., 2021. Calibrate Before Use: Improving Few-Shot Performance of Language Models. In Proceedings of the 38th International Conference on Machine Learning (ICML 2021), Virtual Event, 18–24 July 2021; PMLR: Cambridge, MA, USA. pp. 12697–12706. DOI: <https://doi.org/10.48550/arXiv.2102.09690>
- [25] Garigliotti, D., Johansen, B., Vigerust Kallestad, J., et al., 2024. EquinorQA: large language models for question answering over proprietary data. In Proceedings of the 27th European Conference on Artificial Intelligence (ECAI 2024), Santiago de Compostela, Spain, 19–24 October 2024; IOS Press: Amsterdam, Netherlands. pp. 4563–4570. DOI: <https://doi.org/10.3233/FAIA241049>
- [26] Jeon, J., Sim, Y., Lee, H., et al., 2025. ChatCNC: Conversational machine monitoring via large language model and real-time data retrieval augmented generation. Journal of Manufacturing Systems. 79, 504–514. DOI: <https://doi.org/10.1016/j.jmsy.2025.01.018>
- [27] Vizniuk, A., Diachenko, G., Laktionov, I., et al., 2025. A Comprehensive Survey of Retrieval-Augmented Large Language Models for Decision Making in Agriculture: Unsolved Problems and Research Opportunities. Journal of Artificial Intelligence and Soft Computing Research. 15(2), 115–146. DOI: <https://doi.org/10.2478/jaiscr-2025-0007>
- [28] Yang, R., Fu, M., Tantithamthavorn, C., et al., 2025. RAGVA: Engineering Retrieval Augmented Generation-based Virtual Assistants in Practice. Journal of Systems and Software. 226, 112436. DOI: <https://doi.org/10.1016/j.jss.2025.112436>
- [29] Liu, X., Ji, K., Fu, Y., et al., 2022. P-Tuning v2: Prompt Tuning Can Be Comparable to Fine-tuning Universally Across Scales and Tasks. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (ACL 2022), Dublin, Ireland, 22–27 May 2022; Association for Computational Linguistics: Stroudsburg, PA, USA. pp. 61–68. DOI: <https://doi.org/10.18653/v1/2022.acl-short.8>
- [30] Li, X.L., Liang, P., 2021. Prefix-Tuning: Optimizing Continuous Prompts for Generation. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (ACL-IJCNLP 2021), Virtual Event, 1–6 August 2021; Association for Computational Linguistics: Stroudsburg, PA, USA. pp. 4582–4597. DOI: <https://doi.org/10.18653/v1/2021.acl-long.353>
- [31] Housby, N., Giurgiu, A., Jastrzebski, S., et al., 2019. Parameter-Efficient Transfer Learning for NLP. In Proceedings of the 36th International Conference on Machine Learning (ICML 2019), Long Beach, CA, USA, 9–15 June 2019; PMLR: Cambridge, MA, USA. pp. 2790–2799. DOI: <https://doi.org/10.48550/arXiv.2102.09690>

- arXiv.1902.00751
- [32] Radford, A., Wu, J., Child, R., et al., 2019. Language Models are Unsupervised Multitask Learners. OpenAI: San Francisco, CA, USA. DOI: <https://doi.org/10.48550/arXiv.1909.13723>
- [33] Rogers, A., Kovaleva, O., Rumshisky, A., 2020. A Primer on Neural Network Models for Natural Language Processing. *Journal of Artificial Intelligence Research*. 57, 615–731. DOI: <https://doi.org/10.1613/jair.1.11640>