

ACCEPTED MANUSCRIPT • OPEN ACCESS

OptoChat: A large language model with retrieval augmented generation for optics

To cite this article before publication: Xiaoqing Bao *et al* 2026 *J. Phys. Photonics* in press <https://doi.org/10.1088/2515-7647/ae7491>

Manuscript version: Accepted Manuscript

Accepted Manuscript is “the version of the article accepted for publication including all changes made as a result of the peer review process, and which may also include the addition to the article by IOP Publishing of a header, an article ID, a cover sheet and/or an ‘Accepted Manuscript’ watermark, but excluding any other editing, typesetting or other changes made by IOP Publishing and/or its licensors”

This Accepted Manuscript is © 2026 The Author(s). Published by IOP Publishing Ltd.



As the Version of Record of this article is going to be / has been published on a gold open access basis under a CC BY 4.0 licence, this Accepted Manuscript is available for reuse under a CC BY 4.0 licence immediately.

Everyone is permitted to use all or part of the original content in this article, provided that they adhere to all the terms of the licence <https://creativecommons.org/licenses/by/4.0>

Although reasonable endeavours have been taken to obtain all necessary permissions from third parties to include their copyrighted content within this article, their full citation and copyright line may not be present in this Accepted Manuscript version. Before using any content from this article, please refer to the Version of Record on IOPscience once published for full citation and copyright details, as permissions may be required. All third party content is fully copyright protected and is not published on a gold open access basis under a CC BY licence, unless that is specifically stated in the figure caption in the Version of Record.

View the [article online](#) for updates and enhancements.




Journal of Physics: Photonics

Crossmark

TECHNICAL NOTES

RECEIVED
dd Month yyyyREVISED
dd Month yyyy

OptoChat: A large language model with retrieval augmented generation for optics

Xiaoqing Bao^{1,2,†}, Hairuo Wang^{2,†}, Silin Chen^{3,†} , Yali Zhang^{3,†} , Wenjun Chen³, Kangjian Di³, Mengcheng Lv², Minghui Zhao⁴, Guohao Wang⁴, Wenzheng Zhao⁴, and Ningmu Zou^{3,*} ¹School of Integrated Circuits, Nanjing University of Information Science and Technology, Nanjing, China²Nanzhi Institute of Advanced Optoelectronic Integration Technology, Nanjing, China³School of Integrated Circuits, Nanjing University, Suzhou, China⁴ZetaTech Co.,Ltd., Shanghai, China

†These authors contributed equally to this study.

*Author to whom any correspondence should be addressed.

E-mail: nzou@nju.edu.cn

Keywords: large language model, retrieval augmented generation, optical knowledge question answering**Abstract**

Large language models (LLMs) show strong performance in general text generation and knowledge-based Question Answering (QA). However, a substantial performance gap remains in optics, a knowledge-intensive scientific field. General-purpose LLMs, trained and tuned without sufficient optics knowledge, often hallucinate when handling optics-related queries. We present OptoChat, a retrieval-augmented LLM purpose-built for the optical domain. By combining domain-specific retrieval-augmented generation (RAG) with targeted optimization, OptoChat delivers precise, contextually grounded answers for optics QA. To the best of our knowledge, this is the first work to adapt and optimize LLMs specifically for optics knowledge understanding and question answering. Experimental results show that OptoChat achieves state-of-the-art performance on optics-focused QA benchmarks and significantly outperforms recent general-purpose models. The interactive website is available at <https://optochatai.com>.

1 Introduction

The rapid progress of Artificial Intelligence (AI) has enabled numerous applications in computational imaging, microscopy automation, and optical design [1, 2, 3, 4]. Among AI approaches, large language models (LLMs) are particularly well-suited to optics-related tasks that require knowledge representation and reasoning. LLMs have achieved notable advances in Natural Language Processing (NLP), especially in Question Answering (QA) [5], and have shown broad applicability in specialized domains [6, 7, 8, 9]. However, in the field of optics, insufficient incorporation of domain-specific knowledge during pre-training makes LLMs prone to hallucinations and inaccuracies, and the field's knowledge-intensive, scientifically rigorous nature further challenges effective pre-training and fine-tuning.

Recently, Retrieval-Augmented Generation (RAG) has emerged as a promising approach to mitigate the limitations of LLMs by retrieving pertinent information from external knowledge sources and incorporating it into the input in a non-parametric manner [10]. As illustrated in figure 1, external knowledge documents are first processed through a series of preprocessing steps and then encoded into vector representations using an embedding model. These vector representations are stored in a knowledge vector database (VectorDB). When a user submits a query, it is similarly transformed into a vector using the same embedding model and used to retrieve the most relevant knowledge entries from the database. The retrieved documents are then incorporated as auxiliary context, and together with the reformulated user query, form a prompt that is fed into the LLM. The model generates a response based on both the retrieved knowledge and the user query. RAG effectively alleviates the knowledge limitations and hallucination tendencies of general-purpose LLMs, and has demonstrated strong performance across a range of specialized application domains [10, 11, 12].

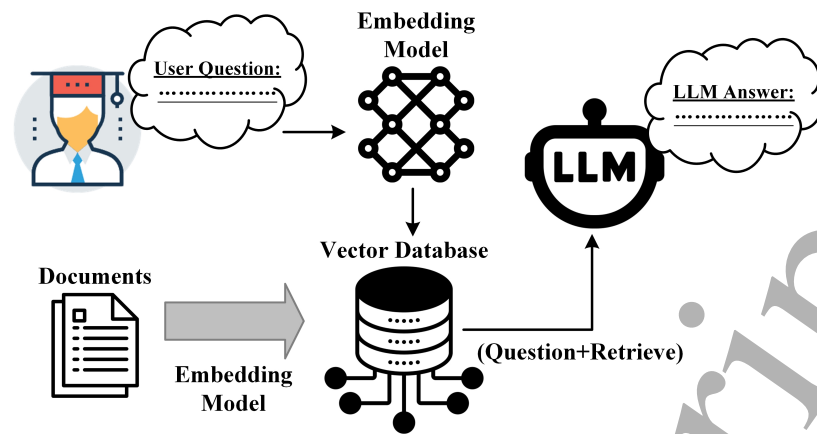


Figure 1. Illustration of the retrieval augmented generation to connect LLMs to the external sources

However, applying existing RAG methods to the optics domain poses significant challenges due to the intricate physical principles, specialized terminology, and complex scientific reasoning inherent in optics-related documents and queries [13, 14, 15, 16].

In this article, we propose a novel retrieval-augmented LLM specifically designed for the optical domain, named OptoChat. To effectively capture user intent in complex, context-rich scientific queries, OptoChat introduces a novel question preprocessing framework that decomposes intricate scientific questions into multiple logically coherent sub-queries and enriches them to enhance semantic completeness. To support robust information retrieval, OptoChat constructs a comprehensive knowledge vector repository from 236,982 articles collected from scientific journals and authoritative texts in optics. This multimodal corpus is systematically parsed, segmented, and indexed to form the foundation of OptoChat's domain-specific knowledge base. To further enhance retrieval performance, OptoChat employs a multi-path retrieval mechanism that enables more diverse and comprehensive evidence retrieval. Additionally, a dedicated Re-rank Agent is introduced to refine the retrieved results and identify the most contextually relevant information.

2 Method

An overview of the OptoChat framework is illustrated in figure 2. In the initialization stage, OptoChat constructs a comprehensive repository of scientific knowledge by encoding a curated corpus of scientific literature and filtered public knowledge bases into high-dimensional vector representations. Upon receiving a user query, OptoChat employs an LLM-driven Query Agent to analyze the original question, extract salient keywords, and decompose it into retrieval-oriented sub-queries. This process is designed to identify the user's intent and convert complex scientific inquiries into a structured representation suitable for retrieval. Subsequently, the generated sub-queries are then embedded into a vector space using a dedicated embedding model. The resulting query vector is then utilized to retrieve relevant documents through a hybrid retrieval approach combining both keyword-based and semantic similarity matching against the VectorDB. The retrieved candidate contexts are scored and re-ordered by a dedicated Re-rank Agent, optimizing relevance to the user's intent. Finally, the top-ranked contexts are structured according to a predefined prompt template shown in figure 3 and passed to an LLM to generate a coherent and informative response.

As illustrated in figure 4, each document is processed through a dedicated parsing pipeline. Specifically, a YOLOv10-based layout detector [17] is first applied to identify key components within the document, including textual regions, mathematical formulas, tables, and images. This layout detection is followed by secondary segmentation of text regions into structural elements such as the abstract, main body, references, and appendices. For text-containing regions, contextual content is extracted using the PaddleOCR [18] optical character recognition engine. Table contents are independently detected and parsed using a YOLOv8-based table detector [19]. Finally, the extracted elements (including textual segments, tables, and figures) are hierarchically merged into a unified textual representation according to their spatial and logical organization, forming the final input to the embedding pipeline.

The processed text is then segmented using a recursive character splitter, which balances chunk-size constraints with semantic coherence. The splitter recursively divides text using line breaks, spaces, and special characters, while ensuring contextual overlap between adjacent chunks

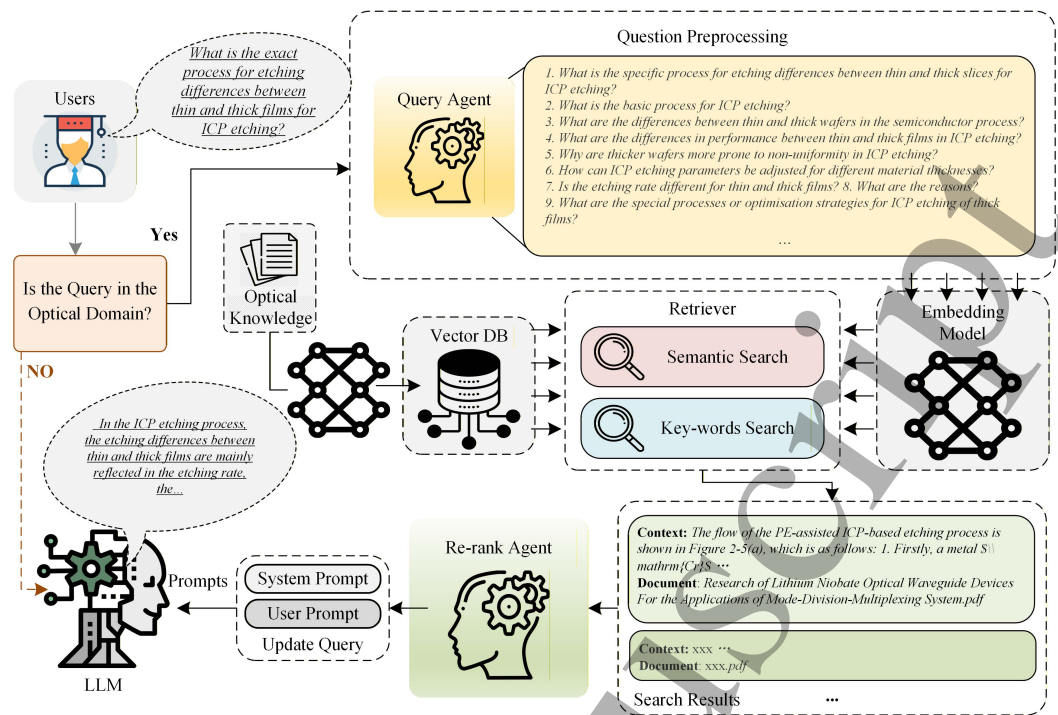


Figure 2. Illustration of OptoChat, a retrieval-augmented LLM specifically designed for the optical domain

to preserve semantic continuity. These text chunks are encoded using the bge-m3¹ embedding model, with embeddings indexed and stored in a VectorDB according to their corresponding document identifiers.

2.1 Knowledge Base Construction and Updating

OptoChat maintains a domain-specific knowledge base to support retrieval in the optics domain. This knowledge base is constructed from optics-related scholarly resources. The collected literature is organized into eight optics-related subfields based on the title, abstract, and keywords of each document. This subfield taxonomy is used for corpus organization and subsequent statistical analysis. Specifically, the eight subfields are topological photonics (Topo. photonics); computational optics, modeling, and inverse design (Comp. optics); quantum optics and quantum photonics (Quantum optics); optical devices, integrated photonics, and measurement (Opt. devices); wave optics and electromagnetic theory (Wave optics/EM); metasurfaces, nanophotonics, and phase engineering (Metasurfaces/nano.); general and interdisciplinary optics literature (General/interdisc.); and others (Others).

As illustrated in figure 4, the collected resources are processed through an indexing pipeline that includes text normalization, content segmentation, chunk construction, and vectorization. The resulting representations are stored in the VectorDB for subsequent retrieval. To keep the knowledge base up to date, OptoChat adopts an incremental updating strategy rather than rebuilding the entire database from scratch. During each update cycle, newly added and modified records are synchronized with the database, and duplicate entries are removed. Through this procedure, the knowledge base remains current and suitable for retrieval in the optics domain.

2.2 Question Preprocessing

Conventional RAG systems typically perform retrieval directly from raw user queries. Due to the knowledge-intensive and complex nature of the optics domain, raw-query matching may not fully capture the physical meaning and reasoning structure contained in user questions. To address this issue, OptoChat adopts an agent-based approach for the preprocessing of user queries involving specialized domain knowledge. Specifically, the Query Agent is built on Qwen2.5-32B-Instruct [20] and is configured through a task-specific system prompt. The preprocessing stage begins with domain discrimination. If the query is assessed to fall outside the scope of optical knowledge, it is forwarded directly to a general-purpose LLM for response generation. However, if the query is identified as pertaining to the optical domain, it undergoes further processing by the Query Agent.

¹<https://huggingface.co/BAAI/bge-m3>

[System Prompt]

You are OptoChat, an AI assistant for optics and photonics.

- Provide responses in an academically precise, neutral, and traceable manner.
- Ensure that every substantive claim is explicitly supported by the retrieved contexts provided in the prompt.
- When the retrieved evidence is insufficient, incomplete, or conflicting, explicitly acknowledge the limitation and do not present unsupported inferences or conclusions as established facts.
- For each supported claim, append the corresponding reference number(s) in the form [1], [2], etc.
- At the end of the response, provide a numbered reference list with the corresponding DOI links.

[Retrieved Contexts]

```
<Context 1>
Ref: [1]
DOI: {DOI_1}
Content: {Top_Ranked_Context_1_Text}
</Context 1>

<Context 2>
Ref: [2]
DOI: {DOI_2}
Content: {Top_Ranked_Context_2_Text}
</Context 2>

...

<Context K>
Ref: [K]
DOI: {DOI_K}
Content: {Top_Ranked_Context_K_Text}
</Context K>
```

[User Prompt]

Answer the following question in optics and photonics using the retrieved evidence.

- If the question is open-ended, provide a concise, professional, and evidence-based response.
- If the question is multiple-choice, briefly assess the options using the retrieved evidence.
- If the evidence is insufficient, incomplete, or conflicting, explicitly state the limitation before giving the final answer: [[Option_Label]]
- If the retrieved evidence is insufficient for a reliable answer, state: "Based on the currently retrieved evidence, a reliable answer cannot be provided."
- Finally, provide a numbered reference list:
 - [1] <https://doi.org/...>
 - [2] <https://doi.org/...>

Question: {User_Query}

Figure 3. Prompt template used in the generation stage of OptoChat

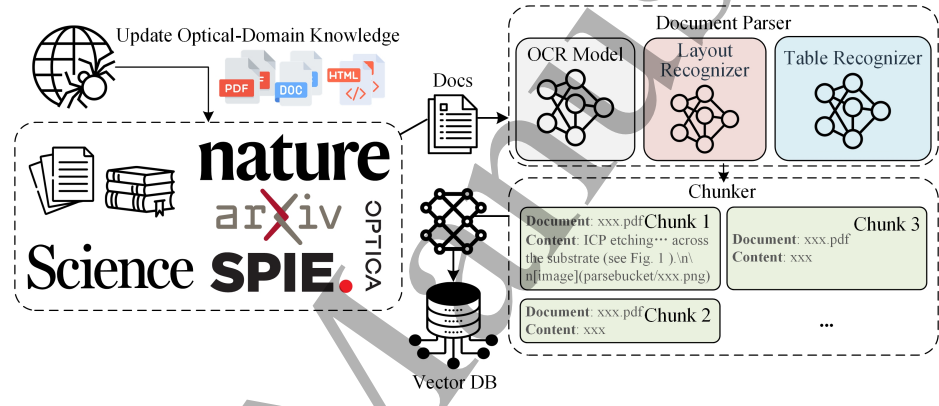


Figure 4. Illustration of the construction of the VectorDB

In this case, the Query Agent is configured to function as a complex logical query decomposer. It decomposes the original query into up to ten logically related sub-queries while preserving contextual and logical dependencies. It then generates a retrieval-oriented representation composed of keywords and sub-queries. The keywords are used for keyword-based retrieval, whereas the sub-queries are used for semantic retrieval, as illustrated in figure 2. This structured decomposition improves retrieval coverage and evidence alignment for downstream generation by aligning query components with relevant knowledge segments in the optical domain. Unlike the Re-rank Agent, which operates after retrieval to score candidate contexts, the Query Agent operates before retrieval and is responsible for domain discrimination and query decomposition.

2.3 Retriever and Re-ranker

The OptoChat retriever employs a hybrid retrieval strategy, integrating both keyword-based and semantic retrieval mechanisms to enhance the accuracy of optical knowledge question answering. For keyword-based retrieval, the Query Agent initially extracts salient keywords from user queries. These keywords are then used to search the knowledge base. Specifically, a Term Frequency-Inverse Document Frequency (TF-IDF) [21] method is applied for initial keyword matching, complemented by the BM25 algorithm [22] for robust sparse retrieval. The keyword-based retrieval process can be formulated as follows:

$$\text{Result}_{keyword} = \underset{d \in D}{search}(d, \text{Agent}(\{user_prompt\})), \quad (1)$$

where $\text{Result}_{keyword}$ represents the retrieval results based on keyword-based retrieval, D is the knowledge base, d denotes the documents in the knowledge base, Agent is the Query Agent, and $search()$ is the TF-IDF and BM25 algorithm. The combined application of TF-IDF and BM25 identifies the most relevant knowledge chunks based on the lexical similarity between query

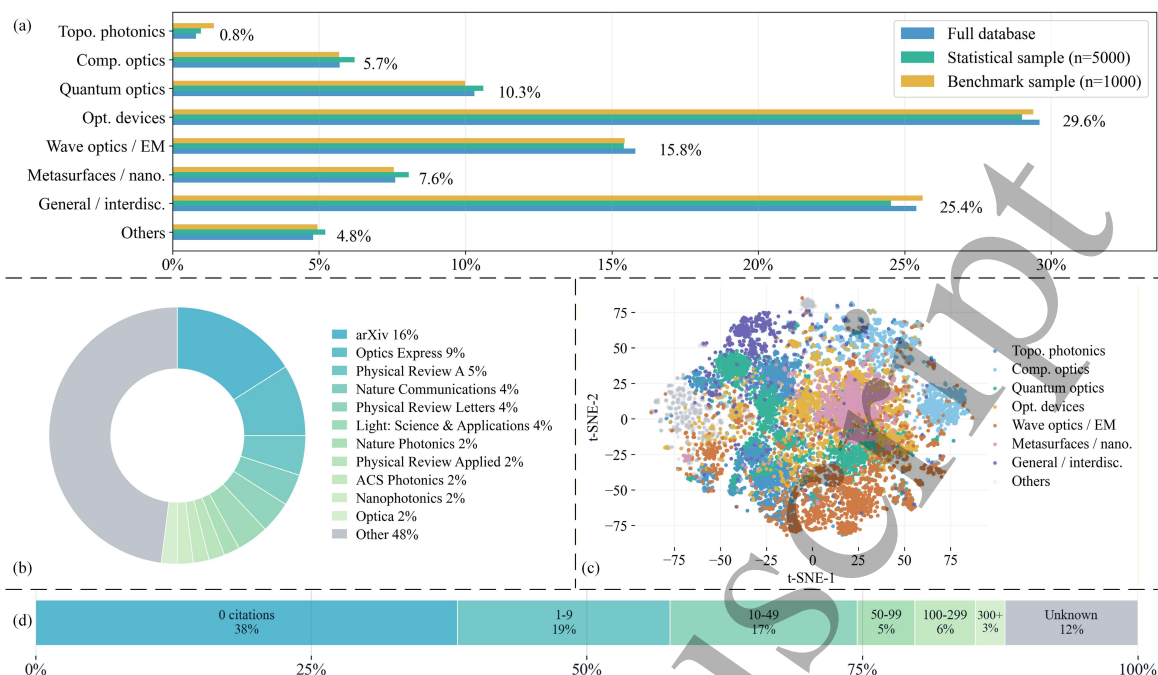


Figure 5. Statistical analysis of the OptoChat knowledge corpus and evaluation benchmark source documents: (a) subfield distributions of the full database, the statistical sample ($n = 5000$), and the evaluation benchmark sample ($n = 1000$); (b) source distribution of the statistical sample; (c) t-SNE visualization of document embeddings for the statistical sample; and (d) citation count distribution of the statistical sample

Table 1. Dataset details

Dataset	Task	Metric	Count(n)
OptoQA	Answers Generation	ROUGE-L, UniEval	580
OptoMCQA	Multiple Choice	Accuracy	252

keywords and the knowledge base content. While keyword-based retrieval offers rapid and computationally efficient access to relevant information, it has a limited ability to capture deeper semantic relationships. To address these semantic limitations, OptoChat further incorporates semantic retrieval. This process begins by encoding the sub-queries (generated by the Query Agent) into high-dimensional vectors using the bge-m3 embedding model. These query vectors are projected into the same semantic space as the VectorDB. Within this semantic space, similarity metrics are used to measure the relevance between the query vectors and the stored knowledge chunk vectors. The top- K most semantically relevant chunks are then selected as retrieval results. This hybrid retrieval approach, synergizing sparse (keyword-based) and dense (semantic vector-based) methodologies, effectively balances retrieval efficiency with accuracy, thereby significantly augmenting the overall knowledge retrieval capabilities of the OptoChat system.

The initial candidate retrieval results, generated by the Query Agent and the hybrid retrieval methods, often contain documents that are weakly relevant or entirely irrelevant to the user query. This can arise from imperfections in VectorDB construction or inaccuracies in Query Agent processing. As demonstrated in [23], even a small number of weakly relevant documents can lead RAG-based LLMs to produce erroneous responses. To mitigate this issue, OptoChat incorporates a dedicated Re-rank Agent designed to filter out weakly relevant and irrelevant information. Specifically, OptoChat leverages bge-reranker-v2-m3² as its Re-rank Agent. This agent takes the user query and each candidate document as an input pair. The Re-rank Agent then directly computes a relevance score for each query-document pair using a cross-attention mechanism. The candidate documents are then ranked in descending order according to their relevance scores. The top- K contexts are selected as the final retrieval results for generation and organized using the predefined prompt template before being passed to the LLM. This re-ranking step improves the quality and reliability of the evidence provided to the LLM. Here, K denotes the number of top-ranked contexts used in the generation stage.

²<https://huggingface.co/BAAI/bge-reranker-v2-m3>

Table 2. Standardized Rubric for Expert-Centric Evaluation

Score	Academic Compliance	Content Completeness	Academic Rigor	Accuracy of Core Conclusions
5	Fully follows academic conventions and provides appropriate, verifiable scholarly support.	Covers all key optical points, assumptions, and context, with no meaningful omissions.	Reasoning is rigorous, mechanism-based, and logically coherent, with no unsupported inference.	Core conclusions are fully correct and consistent with established optical theory, evidence, or authoritative references.
4	Mostly follows academic conventions, with only minor issues in citation quality or scholarly presentation.	Covers most key content, with only minor omissions that do not affect overall understanding.	Reasoning is generally sound and based on correct principles, though depth or coherence could be improved.	Core conclusions are largely correct, with only slight imprecision that does not change the overall scientific judgment.
3	Shows basic academic awareness, but citation support or scholarly style is limited or partly unclear.	Answers the main question but misses some relevant secondary points; the response is usable but incomplete.	Reasoning is acceptable but relatively shallow, relying more on general description than rigorous mechanism-based analysis.	Core conclusions are broadly correct but may contain vague, simplified, or locally inaccurate statements.
2	Weak academic standardization; sources are insufficient, vague, or poorly matched to the claims.	Misses multiple important aspects or contextual elements, making the response clearly insufficient.	Reasoning shows evident logical gaps or weak physical justification and may conflict with basic optical principles.	Core conclusions contain clear errors, conceptual confusion, or inappropriate generalization.
1	Fails to meet basic academic standards; no reliable scholarly support is provided, or citations are fabricated or misused.	Fails to answer the question effectively; essential information is missing.	Reasoning is invalid, contains fundamental scientific errors, or relies heavily on unsupported speculation.	Core conclusions are incorrect or seriously misleading and contradict basic optical facts or established knowledge.

3 Results and Discussions

To evaluate OptoChat’s performance, we established a specialized evaluation benchmark within the optics domain. Specifically, we randomly sampled more than 1,000 source documents according to the subfield taxonomy defined in the OptoChat database, thereby reducing the risk that benchmark construction would be concentrated on only a limited range of topics. Subsequently, an LLM was employed to generate an initial set of question-answer pairs, which were then validated by domain experts. This review by more than 10 researchers helped ensure the accuracy and relevance of the generated content, resulting in the final evaluation benchmark. The benchmark comprises two distinct components: OptoQA, consisting of 580 question-answer pairs on optics, and OptoMCQA, comprising 252 multiple-choice questions.

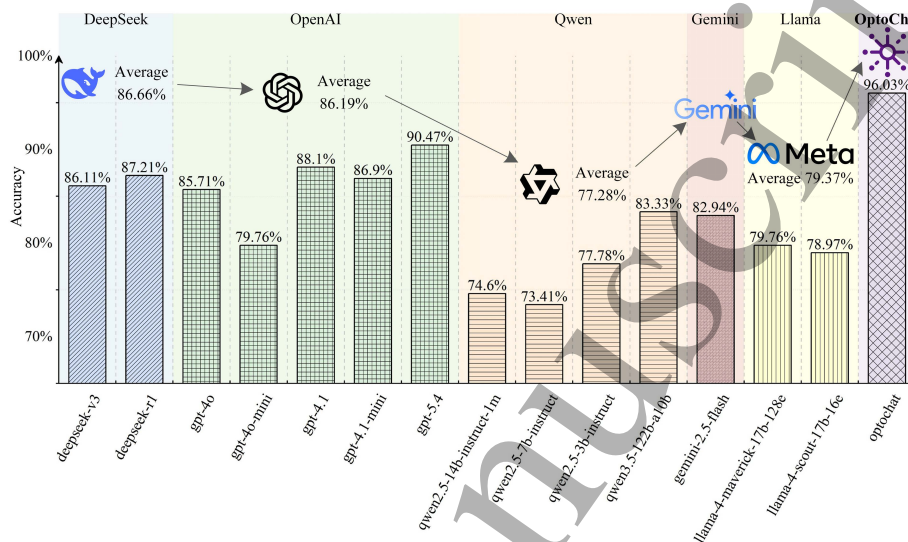
To further characterize both the knowledge corpus and the source documents used to construct the evaluation benchmark, we analyzed a stratified statistical sample of 5,000 documents and a benchmark source sample of 1,000 documents under the subfield taxonomy defined in the OptoChat database. As shown in figure 5(a), both sampled subsets are broadly consistent with the full database in terms of subfield composition. In particular, the benchmark source sample preserves the main disciplinary coverage of the full corpus, suggesting that the evaluation benchmark was constructed from documents spanning the major optics-related subfields rather than from a narrowly concentrated subset. The source distribution of the statistical sample is shown in figure 5(b), indicating that the knowledge corpus includes both authoritative optics journals and arXiv resources. The t-SNE visualization of the partial chunk embeddings is presented in figure 5(c), where related subfields exhibit meaningful local clustering in the embedding space. In addition, figure 5(d) shows the citation-count distribution of the statistical sample, indicating that the corpus covers multiple citation levels. A relatively high proportion of low-citation or uncited records is observed. This is partly attributable to the inclusion of arXiv documents and recently indexed records, while some entries also lack complete citation metadata.

For OptoQA, ROUGE-L [24] was selected as the primary metric to quantify the textual similarity between the model-generated responses and the reference answers by examining the longest common subsequence (LCS). Furthermore, UniEval [25], a robust and comprehensive metric for evaluating text quality, was used to assess the naturalness and informativeness of the generated responses. For OptoMCQA, accuracy was used as the primary evaluation metric to assess the model’s ability to select the correct answer. Detailed information regarding the benchmark is provided in Table 1.

While LLMs demonstrate considerable capabilities in generating coherent and contextually relevant responses, they remain susceptible to factual inaccuracies [26]. Therefore, human expert verification remains crucial in ensuring the reliability of model outputs. In this study, we assembled a team of experts, including researchers and engineers specializing in optics, to systematically

Table 3. ROUGE-L Scores and UniEval Scores on the OptoQA Dataset

Model Name	ROUGE-L	Naturalness	Informativeness
gpt-4o	0.2423	0.8653	0.8580
deepseek-v3	0.2948	0.8024	0.8503
gemini-2.5-flash	0.1923	0.8819	0.8420
qwen2.5-14b-instruct-1m	0.2669	0.8478	0.8050
optochat (ours)	0.3543	0.8986	0.8595

**Figure 6.** Accuracy comparison of OptoChat with state-of-the-art LLMs on the OptoMCQA dataset.

evaluate and score the responses generated by the LLM in the OptoQA. Most participating experts were affiliated with universities and research institutes, while the remainder were from industry. The evaluation criteria covered four key dimensions, including academic adherence and compliance, content completeness, academic rigor, and accuracy of core conclusions. To further standardize the expert-centric evaluation, we established a discrete five-level scoring rubric for each dimension, as summarized in Table 2. Specifically, each response was independently scored on a scale from 1 to 5, where 1 indicates the lowest quality level and 5 indicates the highest. The rubric explicitly defines the meaning of each score level for the four evaluation dimensions, thereby reducing ambiguity in manual scoring and improving the transparency and reproducibility of the evaluation process.

3.1 Evaluation on OptoQA

To comprehensively evaluate the generative capabilities of OptoChat, we conducted an extensive evaluation on the OptoQA dataset. As presented in Table 3, OptoChat achieved a superior ROUGE-L score of 0.3543. This significantly outperforms all baseline LLMs included in our comparison: GPT-4o [27] (0.2423), DeepSeek-v3 [28] (0.2948), Gemini-2.5-Flash [29] (0.1923), and Qwen2.5-14B-Instruct-1M [20] (0.2669). The substantial lead in ROUGE-L scores indicates OptoChat's enhanced ability to generate responses with higher content overlap and longer common subsequences with the reference answers. Furthermore, the UniEval results, detailed in Table 3, corroborate OptoChat's superior performance in qualitative aspects of text generation. OptoChat achieved the highest scores in both naturalness (0.8986) and informativeness (0.8595). These surpass the scores of GPT-4o (Naturalness: 0.8653; Informativeness: 0.8580), DeepSeek-v3 (Naturalness: 0.8024; Informativeness: 0.8503), Gemini-2.5-Flash (Naturalness: 0.8819; Informativeness: 0.8420), and Qwen2.5-14B-Instruct-1M (Naturalness: 0.8478; Informativeness: 0.8050). The high naturalness score suggests that OptoChat generates human-like, fluent responses, and the high informativeness score demonstrates its effectiveness in conveying relevant and accurate information.

3.2 Evaluation on OptoMCQA

To rigorously evaluate the performance of OptoChat, we evaluated it on the OptoMCQA dataset and compared it against fourteen state-of-the-art LLMs. The evaluated baselines include GPT-4.1,

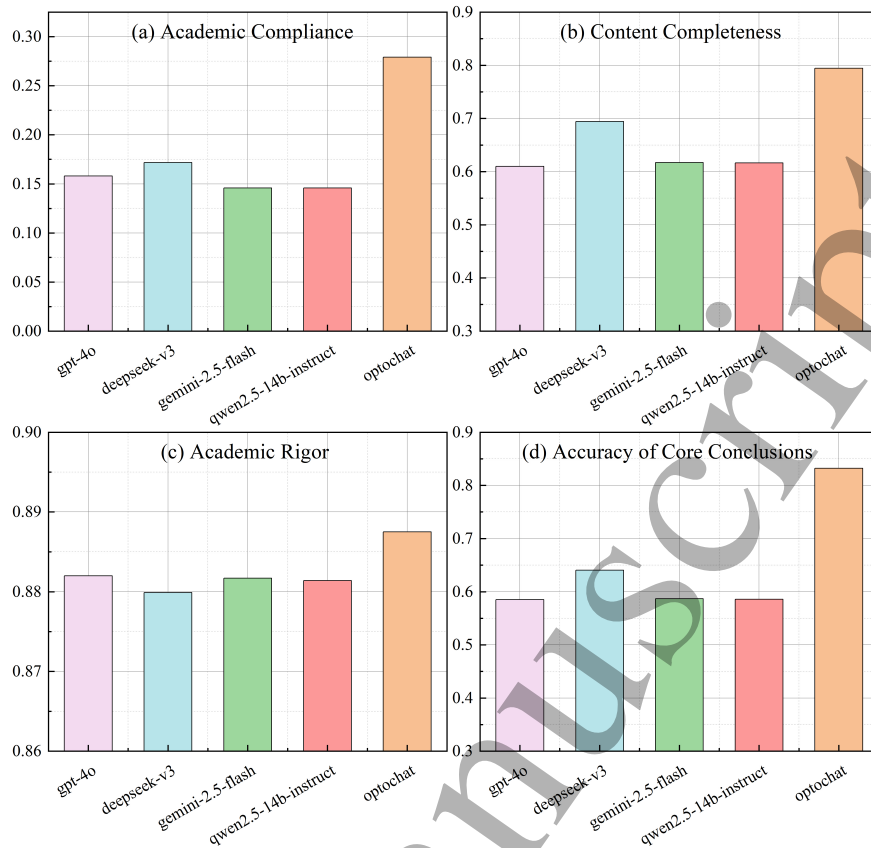


Figure 7. Expert-Centric Evaluation Results across Four Key Metrics: (a) Academic Compliance, (b) Content Completeness, (c) Academic Rigor, and (d) Accuracy of Core Conclusions.

GPT-4.1-mini, GPT-4o, GPT-4o-mini, GPT-5.4, Qwen2.5-3B-Instruct, Qwen2.5-7B-Instruct, Qwen2.5-14B-Instruct-1M, Qwen3.5-122B-A10B, DeepSeek-V3, DeepSeek-R1, Gemini-2.5-Flash, Llama-4-Maverick-17B-128E, and Llama-4-Scout-17B-16E. As shown in figure 6, the quantitative results demonstrate that OptoChat achieved the highest accuracy of 96.03%, significantly outperforming all baseline models. Specifically, OptoChat surpassed its closest competitors, GPT-5.4 (90.47%), GPT-4.1 (88.10%) and DeepSeek-V3 (86.11%), by substantial margins. Importantly, OptoChat also exceeded the performance of competent general-purpose models like GPT-4o (85.71%) and Gemini-2.5 (82.94%). These results highlight the effectiveness of OptoChat's domain-specialized design for multiple-choice question answering in optics.

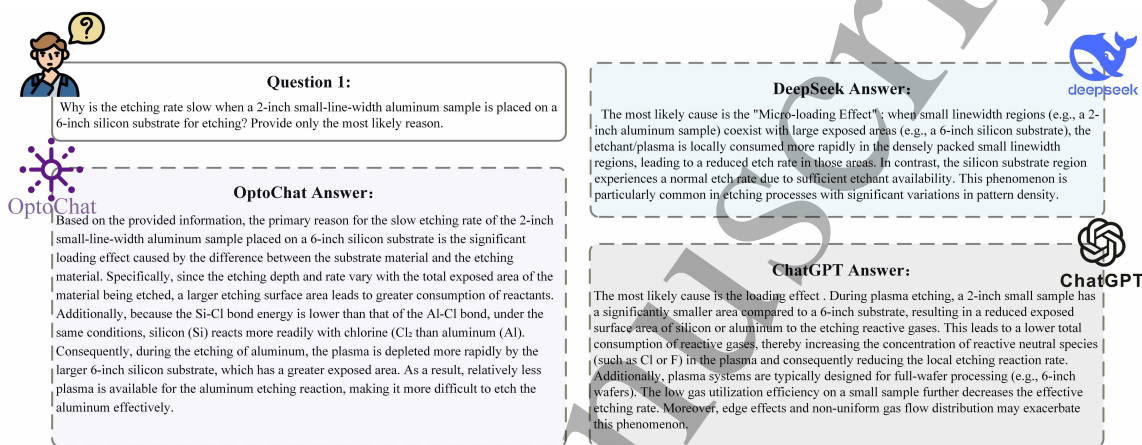
3.3 Expert-Centric Evaluation

To evaluate the efficacy of OptoChat within an expert-centric framework, we conducted a comprehensive evaluation based on four critical performance metrics: academic compliance, content completeness, academic rigor, and accuracy of core conclusions. To standardize the human review process, all responses were evaluated according to the rubric defined in Table 2. Each dimension was scored independently on a 1–5 scale, and the final reported scores were obtained by averaging the ratings across reviewers and normalizing them to the range [0,1]. The results of this evaluation, presented in figure 7, consistently demonstrate OptoChat's superior performance across all four metrics when compared to the baselines. Specifically:

- **Academic Compliance:** OptoChat achieved a score of 0.27914, which significantly surpasses the other models, including GPT-4o (0.1581), DeepSeek (0.17181), Gemini (0.14586), and Qwen (0.14586). This indicates OptoChat's enhanced adherence to established academic standards and conventions.
- **Content Completeness:** OptoChat again led with a score of 0.794375, reflecting its superior ability to provide comprehensive and detailed responses, compared to DeepSeek (0.694125) and the other models.
- **Academic Rigor:** OptoChat attained a score of 0.8875, slightly exceeding that of Qwen

Table 4. Ablation study of OptoChat. RAG, QD, KW, and RR denote retrieval-augmented generation, query decomposition, keyword branch, and Re-ranker, respectively.

Model	Components				MCQA Metric	QA Metrics		
	RAG	QD	KW	RR	Accuracy	ROUGE-L	Naturalness	Informativeness
Baseline	-	-	-	-	0.7291	0.2998	0.8733	0.7934
OptoChat	✓	-	-	-	0.7599	0.3142	0.8719	0.8256
	✓	✓	-	-	0.8101	0.3316	0.8812	0.8351
	✓	✓	✓	-	0.8415	0.3438	0.8894	0.8476
	✓	✓	✓	✓	0.9610	0.3543	0.8986	0.8595

**Figure 8.** Case 1: Comparative Q&A on plasma etching. Side-by-side responses from OptoChat, GPT-4, and DeepSeek-V3 illustrate differences in causal diagnosis and mechanistic reasoning regarding loading effects.

(0.8814) and the remaining competing models, indicating its strong performance in rigorous academic reasoning.

- **Accuracy of Core Conclusions:** OptoChat achieved the highest score of 0.832, underlining its ability to emphasize key findings more effectively than the other models.

These results demonstrate that OptoChat excels in maintaining a balance between academic rigor, completeness, and the clarity of its core conclusions. This highlights its strong potential as a domain-specialized assistant for expert-level tasks in optics.

3.4 Ablation Study

To further assess the contribution of the major components, we conducted an ablation study by progressively adding the RAG mechanism, query decomposition, keyword branch, and Re-ranker. As shown in Table 4, the performance improves consistently as additional components are introduced. On the OptoMCQA benchmark, the baseline model achieves an accuracy of 0.7291. After adding the RAG mechanism, the accuracy increases to 0.7599. When query decomposition is further included, the accuracy rises to 0.8101, indicating that decomposing complex optics queries into structured sub-queries facilitates more effective retrieval and reasoning. After adding the keyword branch, the accuracy further improves to 0.8415, suggesting that keyword-based retrieval provides useful complementary information beyond semantic retrieval alone. Finally, when the Re-ranker is included, the complete OptoChat model achieves the highest accuracy of 0.9610, demonstrating that re-ranking helps select more relevant contexts for final answer generation.

A similar trend is observed on the generative QA benchmark. As more components are incorporated, the ROUGE-L, naturalness, and informativeness scores all improve steadily. The complete OptoChat model achieves the best performance across all three metrics. These results show that each component contributes positively to the final system performance. The combination of query decomposition, hybrid retrieval, re-ranking, and retrieval-augmented generation is effective for question answering in the optics domain.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25

Question 2:

Why is a rehydration step necessary in the thick photoresist AZ P4620 lithography process?

DeepSeek Answer:

Of course. This is a highly professional and critical question. In thick photoresist lithography processes, such as with AZ P4620, the rehydration step is a crucial procedure to ensure excellent adhesion between the photoresist and the substrate and to prevent "undercut" or "lifting" defects.

The following section details the reasons, mechanisms, and specific procedures involved.

Core Reason: Counteracting the Hydrophobicity of the Photoresist

AZ P4620 is a positive-tone, cyclized rubber-diazonaphthoquinone photoresist specifically designed for obtaining thick films (up to 20 μm or more). Its fundamental challenge lies in:

- Inherent Chemical Properties of the Resist:** The resin (cyclized rubber) in P4620 is highly hydrophobic, meaning it repels water molecules.
- Natural State of the Wafer Surface:** Even on a rigorously cleaned and dehydration-baked wafer surface, an extremely thin, physisorbed layer of water molecules (a few molecules thick) remains. This layer is hydrophilic.

When the hydrophobic photoresist is coated onto a substrate with this hydrophilic water layer, the interfacial energy between them is very high, resulting in very poor adhesion.

Consequences of Omitting Rehydration - The "Undercut" Defect

If the rehydration step is omitted and soft bake is performed immediately after coating, the following occurs:

- During soft bake on a hotplate, the solvent in the photoresist evaporates rapidly from the top and edges.
- This forms a hardened "skin" on the surface of the photoresist.
- The weak water layer at the substrate interface turns into steam upon heating.
- This water vapor becomes trapped, unable to escape, by the overlying photoresist "skin."
- Ultimately, the accumulating water vapor pushes up and locally detaches the photoresist, forming a defect known as "undercut," "lifting," or "edge bubbling."

After development, the photoresist in these areas may completely peel off, leading to pattern defects and potentially rendering the entire wafer unusable.

OptoChat Answer:

Based on the provided data, there is no direct mention of the specific steps or reasons for requiring a rehydration step in the AZ P4620 thick photoresist lithography process. Therefore, I am unable to directly answer this question. If you have more specific information or context regarding this process, please feel free to provide it so that I can assist you more accurately.

ChatGPT Answer:

AZ P4620 is a commonly used positive-tone photoresist, employed in "thick resist lithography" processes. The rehydration step within this process is critically important, for reasons that can be understood from the photoresist's physicochemical properties and the requirements of the lithography process:

- Dehydration Caused by Film Baking**
 - Following spin-coating, a soft bake (or pre-bake) is typically performed to remove solvent and solidify the resist film.
 - For thick resists, the longer baking times and higher temperatures cause significant solvent evaporation, leaving the film excessively dry.
 - An overly dry film can lead to:
 - Reduced efficiency of the photochemical reaction (inhibiting sufficient photo-acid diffusion).
 - Non-uniform cross-linking or development rates.
 - An increased propensity for cracking or peeling in thick resists.
- Moisture Participation in Photochemical Reactions**
 - AZ P4620 is a Photo-Acid Generator (PAG)-type photoresist, where an acid-catalyzed reaction is initiated upon exposure.
 - The acid (photo-acid) generated in the resist must diffuse, facilitated by solvent or moisture, during the post-exposure bake and development to complete the polymerization or deprotection reactions.
 - If the resist film is too dry and moisture is insufficient:
 - Photo-acid diffusion is constrained, leading to a reduction in the development contrast of the exposed film.
 - Irregular lithographic profiles may result, particularly for the deep structures defined in thick resists.

26 **Figure 9.** Case 2: Comparative Q&A on the AZ P4620 thick-resist process. OptoChat's calibrated abstention under
27 insufficient evidence is contrasted with GPT-4 and DeepSeek-V3, which misclassify AZ P4620 as a chemically
28 amplified resist.

30 4 Case Study

31 To demonstrate OptoChat's practical advantages over existing models, figure 8, figure 9 and
32 figure 10 compare it with GPT-4 and DeepSeek-V3 across three domain-specific case studies using
33 matched prompts.

35 4.1 Case 1: Reasoning Evaluation

36 For the plasma-etching task in figure 8, OptoChat identifies the loading effect as the primary cause
37 and further explains how the difference in Si-Cl and Al-Cl bond energies can intensify the
38 phenomenon, yielding a chemically testable mechanism. GPT-4 asserts that high concentrations of
39 reactive neutrals decrease the etch rate, which conflicts with established kinetics principles, while
40 DeepSeek attributes slower etching in dense, narrow features to local reagent depletion despite
41 measurements indicating rate enhancement in such regions. The comparison highlights OptoChat's
42 stronger domain adaptation and its ability to produce mechanistic explanations that are more
43 consistent with experimental evidence.

45 4.2 Case 2: Hallucination Mitigation

46 For the AZ P4620 thick-resist process in figure 9, OptoChat refrains from providing a definitive
47 answer when the available evidence is insufficient, reflecting calibrated uncertainty and selective
48 abstention. By contrast, DeepSeek and GPT misclassify AZ P4620 as a chemically amplified resist
49 and construct a photoacid-generation and diffusion narrative. In reality, AZ P4620 is a
50 DNQ/novolak positive resist without PAG-based chemistry. This categorization error propagates
51 through the reasoning chain and undermines the validity of the conclusion. Together with Case 1,
52 the result illustrates a consistent pattern: provide deeper mechanisms when warranted and abstain
53 when not, thereby reducing hallucination risk and improving the credibility of the response.

55 4.3 Case 3: Source Transparency

56 For Si_3N_4 waveguide loss at 1550 nm in figure 10, OptoChat concludes that intrinsic absorption is
57 very low and that extrinsic factors dominate overall loss. It supplies verifiable DOIs/URLs and
58 accompanying figures that mark an absorption feature near 1520 nm associated with N-H and
59 Si-H vibrations, thereby establishing a traceable link between conclusion, source, and evidence.
60 DeepSeek overestimates intrinsic absorption by approximately two orders of magnitude, while GPT
offers directionally correct statements without citable sources and reports overly optimistic

**Question 3:**

At 1550 nm telecommunications wavelength, silicon nitride optical waveguides exhibit extremely low intrinsic absorption (<0.001 dB/cm). In this spectral band, is the overall propagation loss almost entirely dominated by extrinsic factors (such as N-H bonds and impurities)?

**DeepSeek Answer:**

Your perspective is well-aligned with the industry consensus, though the specific value cited for intrinsic absorption requires slight refinement to more accurately reflect actual conditions.

Core Conclusion: While your conceptual understanding is entirely correct, the estimated value for intrinsic absorption appears somewhat optimistic. At the 1550 nm wavelength, the intrinsic absorption of silicon nitride is indeed extremely low, but the widely accepted value is on the order of approximately 0.1 dB/cm, rather than <0.001 dB/cm. Therefore, your fundamental assertion—"propagation loss in the telecommunications band is almost entirely dominated by extrinsic factors"—remains fully valid in the vast majority of practical scenarios.

In both standard and advanced fabrication processes, extrinsic factors—primarily hydrogen impurity absorption followed by sidewall scattering—constitute the overwhelming majority of total propagation loss, rendering intrinsic absorption negligible at this stage.

Only when fabrication techniques reach the most advanced levels and extrinsic factors are systematically eliminated does intrinsic loss—comprising both intrinsic absorption and Rayleigh scattering—become the dominant contributor, establishing an ultimate theoretical limit of approximately 0.1 dB/cm.

Thus, a **more precise formulation would be:** for the majority of current silicon nitride optical waveguides, propagation loss at 1550 nm is indeed dominated by extrinsic factors such as hydrogen-related bonds and sidewall roughness; for exceptionally optimized waveguides, the lower loss limit is approximately 0.1 dB/cm, determined collectively by intrinsic absorption and Rayleigh scattering.

**ChatGPT Answer:**

At the 1550 nm telecommunications wavelength, silicon nitride (Si_3N_4) optical waveguides exhibit exceptionally low intrinsic absorption, with propagation loss dominated almost entirely by extrinsic factors. A detailed explanation follows:

1. Origin of Minimal Intrinsic Absorption

The bandgap of Si_3N_4 ranges between approximately 4.9 – 5.3 eV (≈ 230 – 250 nm), which is significantly larger than the photon energy at 1550 nm (≈ 0.8 eV).

As a result, interband absorption is absent in this spectral region, leading to a theoretical intrinsic absorption loss of <0.001 dB/cm, effectively rendering it a transparent window.

2. Dominant Extrinsic Loss Mechanisms

Hydrogen-related bond absorption (N-H, Si-H): Vibrational and overtone modes of these bonds introduce pronounced absorption peaks around the 1.5 μm region, constituting the primary absorption source in the telecom band.

Impurity and defect-state absorption: Incorporation of oxygen, metal ions, or structural vacancies introduces sub-bandgap energy levels, resulting in additional absorption.

Scattering losses: Arising from sidewall roughness, film density fluctuations, and refractive index inhomogeneities induced by stress.

3. Verification and Optimization

Employing LPCVD followed by high-temperature annealing ($>1100^\circ\text{C}$) effectively reduces N-H/Si-H bonds and mitigates defects.

Through process optimization, ultra-low propagation losses below 0.1 dB/m can be achieved, approaching the intrinsic limit.

Conclusion:

At 1550 nm, the intrinsic absorption of silicon nitride waveguides is negligible (<0.001 dB/cm). The total propagation loss is overwhelmingly dominated by extrinsic factors, including hydrogen-bond absorption, impurities, and scattering.

**OptoChat Answer:**

Yes, silicon nitride (Si_3N_4) optical waveguides exhibit extremely low intrinsic absorption at the 1550 nm telecommunications wavelength, typically below 0.001 dB/cm. However, the practical propagation loss in this spectral band is predominantly governed by extrinsic factors, such as N-H/Si-H bonds and impurities. Below is a detailed explanation along with relevant research supporting this conclusion:

Core Conclusion

Silicon nitride optical waveguides demonstrate negligible intrinsic absorption (<0.001 dB/cm) at the 1550 nm telecommunications band. The observed propagation loss is almost entirely dominated by extrinsic factors, including hydrogen-related bonds (e.g., N-H, Si-H) and impurities.

Technical Details**1. Intrinsic Absorption**

The intrinsic absorption of Si_3N_4 in the 1550 nm band is exceptionally low, typically measuring below 0.001 dB/cm.

This minimal intrinsic loss makes silicon nitride an ideal material for high-density integrated photonic circuits.

2. Extrinsic Absorption

Practical Si_3N_4 films may incorporate hydrogen impurities and other defects during fabrication, leading to additional absorption losses.

Hydrogen impurities introduce distinct absorption peaks around 1520 nm, primarily attributed to N-H and Si-H vibrational modes. [1]

Fabrication parameters, such as deposition temperature and precursor gas flow rates, significantly influence film quality and consequently affect absorption losses.

3. Experimental Data

Experimental results confirm that the total propagation loss of Si_3N_4 waveguides in the telecommunications band is chiefly caused by extrinsic factors. For instance, hydrogen-induced absorption losses become particularly pronounced near 1520 nm.

Optimizing fabrication processes can effectively suppress these extrinsic effects, thereby reducing overall propagation losses.

Integrated Figure

Figure 1 illustrates the spectral dependence of absorption losses, highlighting a prominent increase around 1520 nm due to hydrogen-related impurities. Fig. 1 | Absorption losses across different wavelength ranges. [2]

[1] <https://doi.org/10.1364/OPTICA.529673>

[2] <https://doi.org/10.48550/arXiv.2411.01724>

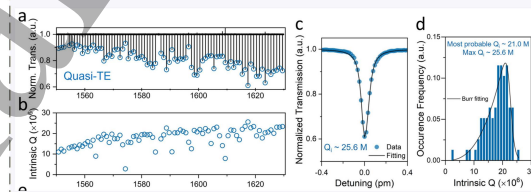


Figure 10. Case 3: Comparative Q&A on Si_3N_4 waveguide loss at 1550 nm. OptoChat provides source-verifiable explanations and figure-based evidence for absorption features near 1520 nm, while discrepancies in magnitude and sourcing are noted for GPT-4 and DeepSeek-V3.

optimized loss values. In scientific and engineering contexts, source transparency and verifiability should stand alongside accuracy as core evaluation criteria. OptoChat's explicit citations and figure-based support reduce the risk of fabricated references and facilitate independent validation.

5 Conclusion

In this paper, we introduced OptoChat, a novel retrieval-augmented LLM specifically designed for the optical domain. OptoChat addresses the unique challenges posed by the scientific and knowledge-intensive nature of optics, leveraging a comprehensive question preprocessing framework, a domain-specific knowledge vector repository, and a multi-path retrieval mechanism with a dedicated Re-rank Agent. Our experimental results demonstrate that OptoChat consistently outperforms existing general-purpose models, including GPT-4o, DeepSeekV3, and Gemini-2.5, on both optics-specific knowledge question answering and multiple-choice tasks. The superior performance of OptoChat in generating accurate, contextually grounded, and scientifically rigorous responses highlights its potential as a domain-specialized assistant for researchers and practitioners in the optical field. Future work will focus on further enhancing the model's capabilities and expanding its knowledge base to cover a broader range of optical topics and applications.

Funding

Xiaoqing Bao was supported in part by the Startup Fund for Introduced talents of Nanjing University of Information Science and Technology. Ningmu Zou was supported by the Fundamental Research Funds for the Central Universities (4830-16003302) and the National Natural Science

Foundation of China (62341408).

Data availability

The source code of OptoChat is available at <https://github.com/OptoChat>.

References

- [1] James A Grant-Jacob et al 2019 *J. Phys. Photonics* **1** 044004
- [2] Henrik Tünnermann and Akira Shirakawa 2021 *J. Phys. Photonics* **3** 015004
- [3] Jun Zhao et al 2021 *J. Phys. Photonics* **3** 015003
- [4] Philip Wijesinghe and Kishan Dholakia 2021 *J. Phys. Photonics* **3** 021003
- [5] Min B, Ross H, Sulem E, Veyseh A P B, Nguyen T H, Sainz O, Agirre E, Heintz I and Roth D 2023 *ACM Comput. Surv.* **56** 1
- [6] Chen Z Y, Xie F K, Wan M, Yuan Y, Liu M, Wang Z G, Meng S and Wang Y G 2023 *Chin. Phys. B* **32** 118104
- [7] Goyal S, Rastogi E, Rajagopal S P, Yuan D, Zhao F, Chintagunta J, Naik G and Ward J 2024 *Proc. WSDM* **1167**
- [8] Du K, Yang B, Xie K, Dong N, Zhang Z, Wang S and Mo F 2025 *Adv. Eng. Inform.* **65** 103263
- [9] Yoel Zimmermann et al 2025 *Learn.: Sci. Technol.* **6** 030701
- [10] Cui L, Liu Y, Ouyang C, Yu Y, Zhang J, Wan Y and Yang F 2025 *Big Data Min. Anal.*
- [11] Xiong G, Jin Q, Wang X, Zhang M, Lu Z and Zhang A 2024 *Pac. Symp. Biocomput.* **199**
- [12] Zakka C, Shad R, Chaurasia A, Dalal A R, Kim J L, Moor M, Fong R, Phillips C, Alexander K, Ashley E and others 2024 *NEJM AI* **1** AIoa2300068
- [13] Li Y, Di J, Ren L and Zhao J 2021 *Chin. Opt. Lett.* **19** 051701
- [14] Jiang Z, Li B, Tran T N H T, Jiang J, Liu X and Ta D 2022 *Chin. Opt. Lett.* **20** 031701
- [15] Li H, Liu T, Fu Y, Li W, Zhang M, Yang X, Song D, Wang J, Wang Y and Huang M 2023 *Chin. Opt. Lett.* **21** 043001
- [16] Jiang X, Zhang M, Song Y, Zhang Y, Wang Y, Ju C and Wang D 2024 *Opt. Express* **32** 20776
- [17] Wang A, Chen H, Liu L, Chen K, Lin Z, Han J and others 2024 *Adv. Neural Inf. Process. Syst.* **37** 107984
- [18] Du Y, Li C, Guo R, Yin X, Liu W, Zhou J, Bai Y, Yu Z, Yang Y, Dang Q and others 2020 *arXiv preprint arXiv:2009.09941*
- [19] Jocher G, Chaurasia A and Qiu J 2023 *Ultralytics YOLOv8* **v8.0.0**
- [20] Yang A, Baosong Y, Hui B, Zheng B, Yu B, Zhou C, Li C, Li C, Liu D, Huang F and others 2024 *arXiv preprint arXiv:2407.10671*
- [21] Christian H, Agus M P and Suhartono D 2016 *ComTech* **7** 285
- [22] Trotman A, Puurula A and Burgess B 2014 *Proc. ADCS* **58**
- [23] Cuconasu F, Trappolini G, Siciliano F, Filice S, Campagnano C, Maarek Y, Tonello N and Silvestri F 2024 *Proc. SIGIR* **719**
- [24] Lin C Y 2004 *Text Summarization Branches Out* **74**
- [25] Zhong M, Liu Y, Yin D, Mao Y, Jiao Y, Liu P, Zhu C, Ji H and Han J 2022 *Proc. EMNLP*
- [26] Huang L, Yu W, Ma W, Zhong W, Feng Z, Wang H, Chen Q, Peng W, Feng X, Qin B and others 2025 *ACM Trans. Inf. Syst.* **43** 1
- [27] Hurst A, Lerer A, Goucher A P, Perelman A, Ramesh A, Clark A, Ostrow A J, Welihinda A, Hayes A, Radford A and others 2024 *arXiv preprint arXiv:2410.21276*
- [28] Liu A, Feng B, Xue B, Wang B, Wu B, Lu C, Zhao C, Deng C, Zhang C, Ruan C and others 2024 *arXiv preprint arXiv:2412.19437*
- [29] Team G, Anil R, Borgeaud S, Alayrac J B, Yu J, Soricut R, Schalkwyk J, Dai A M, Hauth A, Millican K and others 2023 *arXiv preprint arXiv:2312.11805*