POPARIA SCIENCE AND INTERNATIONAL PROPERTY OF STATE OF ST

MMAKER: Multi MedAI Agents with Knowledge-Enhanced Reasoning for Cancer Clinical Decision Support

Sophia Liu

Shanghai American School, China

Abstract

Cancer is one of the leading causes of global mortality. Despite recent advances in cancer detection and treatment, survival rates are still unacceptably low. The current solutions, such as Specialized models and Large Language Models, are limited in ability and often experience limitations such as hallucinations.

To address these challenges, we developed novel MMAKER (Multi MedAI Agents with Knowledge-Enhanced Reasoning), a multi-agent framework that enhances survival prediction accuracy, condition diagnosis, and cancer treatment planning and decision support for doctors. Our system integrates 9 specialized analyst agents to form an AI Multi-Disciplinary Team (MDT): 4 for pathology, 3 for genomics, 1 for clinical analysis, and 1 for survival prediction, integrated by meta-agents and a central reasoning agent. Another key innovation is MedPred, a multimodal survival prediction model that evaluates patients' pathology, genomics, and clinical data of patients to improve prognosis prediction. To reduce the hallucination of AI agents, we developed an LLM-enhanced knowledge base combining vector and graph representations with expert-curated resources that enhance output reliability by 17.2% on MedQA testing dataset. Furthermore, reinforcement learning with chain-of-thought (CoT) optimization ensures the LLM reasoning is transparent. We evaluated our MMAKER on Visual Question Answering (VQA) dataset PathVQA, survival prediction tasks, and oncologist-led response quality assessments. Results demonstrate that our system outperforms state-of-the-art unimodal and multimodal models, achieving up to 72.7% survival prediction accuracy and delivering more reliable natural language responses. These improvements highlights how MMAKER can provide more holistic patient evaluation and clinical decision support for oncologists, especially those in rural areas.

Keywords: Multi-agent systems; artificial intelligence in oncology; cancer survival prediction; multimodal analysis; large language models; chain-of-thought reasoning; personalized medicine; medical knowledge base

Contents

1	Intr	roduction Goal	3			
	1.1	Background Research and Motivation	3			
	1.2	Main Contributions	4			
	1.3	Related Works	4			
	3 (f		1			
2		thodology	6			
	2.1	Overview of Multi MedAI Agent				
	2.2	Multi MedAI Agent Team	7			
		2.2.1 Pathology Analyst Agents	7			
		2.2.1 Pathology Analyst Agents	8			
		2.2.3 Clinical Analyst Agent	9			
		2.2.4 Survival Analyst Agent	9			
		2.2.5 Meta-Agents	10			
		2.2.6 Reasoning Agent	10			
	2.3	MedPred: Specialized Survival Analyst Agent	10			
		2.3.1 Encoders	11			
		2.3.2 Cross-Modal Attention Module	11			
	2.4	2.3.3 Decoders	12			
	2.4	LLM-enhanced Knowledge Base	12			
		2.4.1 Document Processing	13			
		2.4.2 Data Retrieval	14			
		2.4.3 Generation	15			
	2.5	CoT & Reinforcement Learning	15			
3	Res	sults	17			
	3.1					
	3.2	Co-Attention Heatmaps	18			
	3.3	VQA Dataset Evaluation	18			
		3.3.1 Evaluation Dataset	18			
		3.3.2 VQA Evaluation Results	19			
	3.4	Survival Prediction Evaluation	19			
	3.5	Response Quality Evaluation	20			
		/ 0.3				
4			20			
		Ablation Studies	20			
	4.2	Future Work	21			
5 Conclusion						
6	6 Acknowledgements					

1 Introduction Goal

1.1 Background Research and Motivation

Cancer remains a critical global health challenge and a leading cause of mortality worldwide. According to updated reports from the Global Cancer Observatory in 2024, there were approximately 22 million new cancer cases diagnosed globally last year. Furthermore, ongoing respiratory health impacts from the COVID-19 pandemic have contributed to a rise in lung nodule detections, with recent studies indicating that nearly 35% of lung nodules larger than 10mm carry malignant potential. China's vast population accounts for about 24.5% of these global cancer cases. Current data from the American Cancer Society in 2025 identifies lung cancer as the most prevalent cancer type among males, while breast cancer remains the most common among females.

Despite advances in cancer detection and therapy, survival outcomes for most cancers remain low. For example, the global five-year survival rate for pancreatic cancer remains among the lowest at around 10–13%. For lung cancer, survival averages only 10–20% in many regions, reflecting its high lethality despite improvements in targeted and immunotherapies. Even for breast cancer survival in some countries often falls well below 70%. These statistics underscore the urgent need for improved predictive and diagnostic technologies to better manage cancer recurrence risks and survival forecasting worldwide.

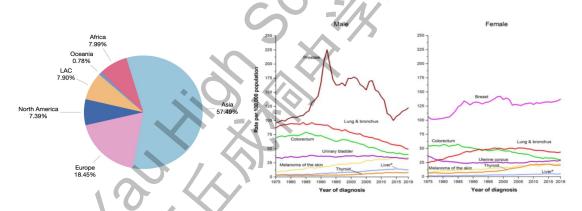


Figure 1: Cancer Statistics [2], [4].

Due to family history with bile duct cancer, we have seen the pain and stress that cancer causes to patients and our family members. Consequently, we wish to innovate a technology that can improve the lives of cancer patients. First, we hope to enable early diagnosis of cancer by leveraging multimodal data analysis, thereby expanding the range of medical treatment options and improving patient outcomes through holistic and personalized cancer care. Second, to accurately calculate risk and survival predictions by analyzing complex factors such as metastasis, recurrence, and treatment responses, addressing key challenges faced by oncologists in formulating effective treatment plans. Through communication with doctors from the Fudan University Cancer Center, we learned that the main challenge doctors currently face when helping cancer patients is with formulating treatment plans based on the patient's metastasis, recurrence, survival prediction for different types of treatment.

Risk prediction refers to analyzing risk factors in order to estimate the probability that a cancer will re-occur in a patient after treatment. Addressing these challenges require a huge amount of multimodal data analysis for accurate prediction, which could be improved with the help of AI. Third, to provide cancer care recommendations and communicate effectively with patients using large language models enriched with professional cancer knowledge and up-to-date information, ensuring high-quality natural language interactions. To understand the queries of the patients and provide high quality conversation to convey analysis results, systems also need professional medical knowledge and natural language processing ability.

1.2 Main Contributions

In order to provide more comprehensive and accurate analysis and interactions with natural language for doctors. My project introduces 4 main contributions:

- 1) We created a novel multi-agents framework, composed of Reasoning agent, Meta agents and Specialized agents, such as clinical record analyst and pathology analyst, genomics analyst etc., forming an AI Multi-Disciplinary Team (MDT) to support cancer clinical decision making.
- 2) We developed a novel specialized tool for cancer survival prediction, which utilizes multimodal data from pathology, genomics, and clinical medical domains to holistically evaluates a patient for accurate prognosis prediction.
- 3) We developed an expert cancer knowledge base with LLM-enhanced optimized intelligent chunking and enhanced categorization to support the Multi Agents analysis for the cancer patient.
- 4) We investigated and implemented oncologist's Chain of Thoughts (COT) when evaluating a cancer patient into our reasoning agents with Reinforcement Learning (RL).

1.3 Related Works

Existing solutions for cancer patient evaluation include Specialized models and Large Language Models. Specialized models are specialized to a particular task within a medical domain; therefore, one Specialized model does not suffice in representing the whole medical domain. Traditional specialized multimodal algorithms require multiple data types for accurate analysis, but they lack matched multimodal patient datasets to support their analysis. Moreover, pathology and genomics data often result in severe data imbalance during fusion. LLMs, on the other hand, lacks professional alignment with real clinical practices and lack cross-domain information when performing crossmodal analysis, leading to more hallucinations.

Specialized Models

For most medical settings, many doctors or professionals choose to use specialized algorithms. These algorithms are specially trained for professional medical tasks such as diagnosing, survival and risk prediction. Some of the most recent include MeTra [7], MCAT [8], SNN [9] etc. The advantages of these algorithms are they include a wide dataset of related medical data and are trained on specialized cancer data.

However, the main issues of specialized algorithms is that they lack adequately matched multimodal patient data. Not all patients have complete data across all medical domains due to variations in testing availability, costs, and clinical practices, resulting in missing and heterogeneous data that complicate data integration. Another critical issue in these models is

the fusion of pathology and genomics data, which often results in severe data imbalance. This imbalance arises because pathology data, such as high-dimensional images, and genomics data, consisting of sparse molecular features, have fundamentally different structures, scales, and noise levels. These challenges are further exacerbated by the scarcity of large, well-annotated matched datasets and difficulties standardizing data collected from diverse sources.

Another significant challenge that traditional models face is that they tend to specialize in very narrow tasks within their medical domain. This occurs because these models are trained on specific datasets representing distinct conditions or features, limiting their ability to recognize or interpret a wider range of abnormalities within the same medical domain. This narrow specialization means that relying on a single domain-specific model is insufficient to fully represent the complexity of a medical domain, as each model focuses on different subsets of features or diseases. Additionally, models trained on limited or homogeneous data may struggle with generalization and robustness when encountering out-of-distribution cases or rare conditions. Consequently, multiple specialized models or a multi-agent framework is often necessary to comprehensively cover the diverse diagnostic tasks within a medical domain while maintaining high accuracy across various pathology types and clinical scenarios.

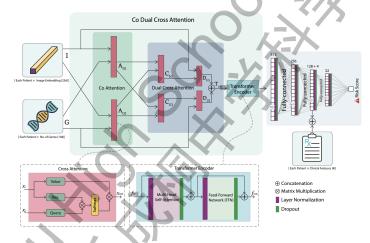


Figure 2: Examples of specialized models (BioFusionNet)

Medical Large Language Models

Medical large language models are algorithms that have the capability to understand natural language inputs and medical answer questions entered by the users. Examples include LLaVA-Med [13], CancerLLM [14] and HuatuoGPT [15]. In the medical field, there has also been an increasing number of medical large language models due to their efficiency in communicating with patients. However, despite their ability to utilize and understand natural language inputs, medical large language models lack specialized training in medical fields such as specific cancer types. This often leads to the generation of incorrect information, or the presence of hallucination in natural language responses when asked a specialized question related to cancer. Additionally, their reasoning processes lack transparency, making it challenging for users to retrace back to the original source of information used for training. Their effectiveness is also heavily dependent on the quality and comprehensiveness of the training data, which may not always reflect the latest advances in cancer research. Another critical challenge is that LLMs lack effective integration of cross-domain information when

performing multimodal analysis. Without properly incorporating diverse data types—such as clinical notes, imaging, pathology, and genomics—LLMs have limited context, increasing the likelihood of hallucinations. This refers to the generation of responses that seem plausible but are in reality factually incorrect and misleading. In a medical context, hallucination can pose a high risk, such as recommending inappropriate medications for treatment or leaving out critical information during patient evaluation.

Thus, to be safely and effectively used in healthcare, LLMs require improved alignment with clinical practices, better incorporation of multimodal data, advanced interpretability, ongoing updates, integration of specialized medical tools.

In summary, the current challenges with specialized algorithms for cancer patient evaluation are the lack of matched multimodal data available for training and the severe imbalance of data type weighing during crossmodal analysis. For the medical large language models, they lack professional cancer knowledge and specialized tools and are also unable to keep up with the most recent publications or research in the field, leading to repeated hallucinations.

2 Methodology

There are three types of agents in our Multi MedAI Agent Framework. Specialized analyst agents are used for analyzing specific areas of the patient data within a medical domain. Meta-agents help consolidate reports from multiple specialized analyst agents within a medical domain. Finally, we utilize a central diagnostic reasoning agent to synthesize genomic, pathology, and clinical domain reports to produce a final, comprehensive and accurate cancer care recommendation.

The Reasoning Agent is also supported with a expert, LLM-enhanced knowledge base to provide specialized oncological knowledge, access recent and relevant papers, and reduce hallucinations. We also include reinforcement learning and Chain Of Thought (COT) from our collaboration with doctors to help improve the Reasoning Agent's thought process.

2.1 Overview of Multi MedAI Agent

Overall, our MedAI MDT includes specialized analyst agents, meta-agents, and a central diagnostic reasoning agent. Categories of analyst agents include pathology, genomic, clinical and a survival analyst agent. As there are multiple specialized analyst agents in the pathology and genomic domains, we utilize a meta-agent to formulate and synthesis a domain report.

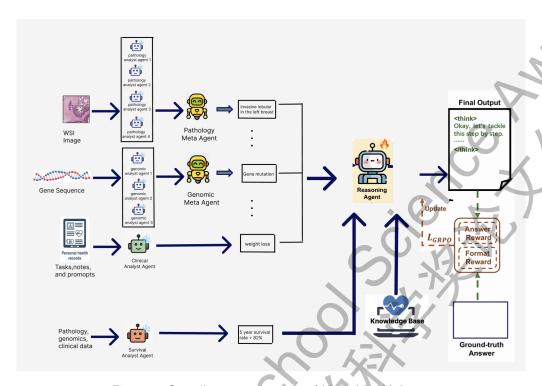


Figure 3: Overall system structure of Multi MedAl Agent

2.2 Multi MedAI Agent Team

In our Multi MedAI Agent Team, there are 9 analyst agents including: 4 pathology analyst agents, consolidated by a pathology meta-agent; 3 genomic analyst agents, consolidated by a genomics meta-agent; a clinical report analyst agent; and a survival analyst agent.

Specialized Analysts Agents are AI agents within the Multi MedAI Agent framework, each agent specializes in processing particular data and completing specific tasks within a medical domain. In our system, we utilize multiple specialized analysts agents for each medical domains (pathology, genomics, and clinical data) and an additional specialized agent for survival prediction. These specialized agents can also call on specialized tools that analyzes data in a certain medical domain to help complete their task.

2.2.1 Pathology Analyst Agents

Currently, this multi-agent team includes 4 pathology analyst agents based on pathology specialized tools (e.g., CHIEF, CLAM, PATHCHAT, LLaVA-Med). We selected specialized pathology tools based on the evaluation shown in the graph below so each have their own advantage and complement each other. There is also a pathology meta-agent responsible for summarizing the pathology analysis report.

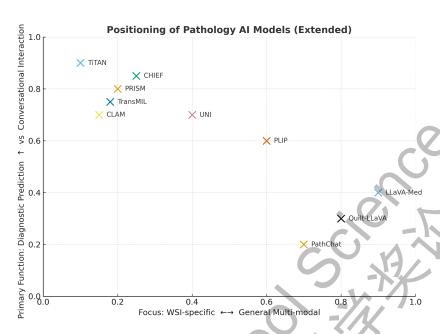


Figure 4: Selection of specialized pathology tool

Pathology analyst agent #1 focuses on identifying diagnostically relevant subregions via attention mechanisms and employs instance-level clustering to refine feature representations. This agent uses function call to access CLAM, which is a weakly supervised deep learning method designed for processing whole-slide images without extensive manual annotations.

Pathology analyst agent #2 aims to extract diverse microscopic pathology representations for cancer cell detection, tumor origin identification, genomic profiling, and prognostic prediction. This agent uses function call to access CHIEF, a weakly supervised machine learning framework for pathology image analysis aimed at systematic cancer evaluation, which relies on the patches and feature extraction that CLAM performs.

Pathology analyst agent #3 aims to enable natural language communication, providing answers, explanations, and insights. Therefore, this agent uses function call to access PathChat, a specialized medical LLM integrated with a vision encoder to facilitate interactive pathology data analysis through natural language.

Pathology analyst agent #4 aims to read medical images, analyze and extract key features, and use natural language to describe and summarize its findings. This agent uses function call to access LLaVA-Med, which is a specialized medical large language and vision assistant that can analyze medical images, including pathology and radiology.

2.2.2 Genomic Analyst Agents

For genomics analysis, genomics-focused specialized agents work together to extract and analyze complex genomics data, providing insights critical for cancer diagnosis and prognosis.

Genomics analyst agent #1 aims to leverage mutation information to classify cancer subtypes, and predicts patient prognosis. This agent uses function call to access SNN, a deep learning network, to extract robust genomic signatures from noisy, high-dimensional gene expression data to support classification and prediction.

- Strength: captures complex nonlinear relationships, no need for manual feature selection.
- Limitation: black-box nature, risk of overfitting, limited interpretability.

Genomics analyst agent #2 performs comparative analysis to other patients begins with quantifying cohort-level gene expression differences and projecting them to individual patients using normalization, differential expression testing, and patient-specific outlier detection to confirm subtype alignment. Next, pathway enrichment methods transform gene-level data into per-patient pathway activity profiles, enabling the identification of therapeutic targets. Protein-protein interaction (PPI) networks overlay expression signals to pinpoint key network hubs and actionable modules. Finally, machine learning selects robust gene panels. After variance filtering, feature selection techniques like LASSO, SVM-RFE, or Random Forest prioritize features. Predictive models then generate patient risk scores that integrate genomics and clinical data for prognosis prediction.

- Strength: interpretability and mechanistic insight.
- Limitation: requires downstream validation and not immediately actionable for patient

Genomics analyst agent #3 follows clinical guidelines, introducing clinical factors to support personalized, guideline-aligned treatment decisions. This comparison delivers actionable biomarkers through immunohistochemistry (IHC) and in situ hybridization(ISH) to assess receptor status, targeted NGS panels to identify key driver mutations, homologous recombination deficiency (HRD) scores to predict response to DNA repair therapies, and PD-L1 assays to guide immunotherapy eligibility. Integrated with clinical parameters, this third clinical analyst agent can provide guideline-backed, personalized treatment decisions, ensuring genomics findings translate effectively into patient care.

- Strength: guideline-backed, evidence-based, and directly usable in the clinic.
- Limitation: limited biomarker scope; cannot fully explain tumor heterogeneity.

Together these three analyst agents form a closed loop, linking discovery, prediction, and clinical application to advance precision oncology. There is also a genomics meta-agent responsible for summarizing the genomics analysis report from the 3 genomic analyst agents and supporting the Reasoning Agent with the genomics report.

2.2.3 Clinical Analyst Agent

For clinical data analysis, the clinical analyst agent aims to process diverse clinical inputs into clinical reports to support the Reasoning Agent's analysis. To perform these tasks, the clinical agent accesses doctor's notes and clinical data tables to analyze and summarize into a clinical report to assist doctors in further analysis. The clinical analyst agent can also function call on Paddle OCR to read clinical table or doctors notes written on paper.

2.2.4 Survival Analyst Agent

The survival analyst agent aims to leverage multimodal clinical data to provide essential insights that assist clinicians in tailoring treatment plans and managing patient expectations.

The survival prediction analyst we used is our novel MedPred which is based on end-to-end multimodal Transformer that holistically evaluates a patient for accurate prognosis prediction. Its main functions include survival, metastasis and recurrence prediction to assist doctors in treatment plan decision.

2.2.5 Meta-Agents

Meta Agents act as the coordinators within the Multi MedAI Agent framework, managing communication and workflow between the central Reasoning Agent and various Specialized Analyst Agents. Their primary role is to manage the analysis results of Specialized Analyst Agents within a specific medical domain and generate comprehensive report or analysis of the information extracted within that medical domain to the Reasoning Agent.

The framework gains from complementary viewpoints, since different medical tools may reveal distinct abnormalities or emphasize different aspects of each domain. For each medical domain, the Meta agents generate a report or summary of observations and initial interpretations. These reports are forwarded and integrated by the Reasoning agent. They act as essential inputs for the concluding diagnostic analysis, ensuring the system's comprehensive evaluation is guided by both thorough multimodal interpretation and specialized survival prediction, improving the accuracy and completeness of cancer care recommendations.

The Meta Agent also balances the medical domains that have unequal data representation. For example, pathology, which often have high spacial resolution, compared to the tabular data of genomics, often leads to over-representation of results extracted from that medical domain.

2.2.6 Reasoning Agent

The Reasoning Agent acts as the central decision-making component of the Multi MedAI Agent system, designed to emulate the clinical reasoning process of expert oncologists by integrating and analyzing diverse multimodal data. It consolidates inputs from genomics, pathology, clinical records, and imaging processed by specialized domain models and the LLM-enhanced knowledge base. This comprehensive synthesis enables the generation of precise, clinically relevant cancer care recommendations that are both accurate and interpretable.

The agent interacts naturally with users through conversational language while maintaining dialogue context across multiple interactions, enabling personalized and continuous communication with patients and healthcare providers. By explicitly modeling the diagnostic reasoning pathway, the Reasoning Agent transforms complex multimodal reports into understandable insights, supporting informed clinical decision-making.

2.3 MedPred: Specialized Survival Analyst Agent

After communicating with doctors, we found that out of all the specialized tasks completed by the specialized models, survival prediction is essential in helping doctors develop a suitable treatment plan. To address this issue, we developed our own survival prediction specialized model for our survival analyst agent to call on. Our survival prediction specialized model is developed to address the question of survival, metastasis and reccurence rate for cancer patients, helping provide personalized cancer care.

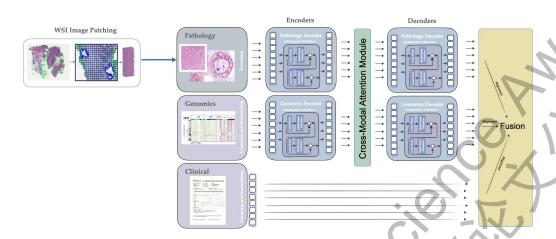


Figure 5: Survival prediction specialized model

MedPred is an end-to-end cross-modal Transformer system that holistically evaluates a patient for accurate prediction. We included a clinical branch on top of pathology and genomics to provide more accurate information for prediction. MedPred simultaneously analyzes three modalities of patient data, including clinical records, genomics data, and pathology images for a more holistic evaluation of the cancer patient. In Med-Pred, the cross-modal attention module identifies and explores connections between all 3 different modalities to share information for better holistic risk and survival prediction. Our novel model also utilizes an encoder-decoder structure and cross-modal attention module for pathology and genomics, performing fusion with the clinical branch as the last step.

2.3.1 Encoders

For the Pathology, Genomics, and Clinical Encoders, they each have a learnable class token that gathers information from the pathological patch features, gene sequences, and clinical texts, respectively [25].

Pathology Encoder: The initial input of the pathology encoder experiences Layer Normalization and Multi-Head Self-attention, then information from the original tokens is added to the result. Then, there is a Pyramid Position Encoding Generator (PPEG) module, which explores the relationships between different patches in the pathology images. Afterwards, there is another layer of Layer Normalization and Multi-Head Self-attention to produce the output.

Genomics Encoder and Clinical Encoder: The initial inputs of the genomics encoder and clinical encoder go through a similar process as the pathology encoder, except without the PPEG module, then produces the output.

The class tokens for the pathology encoder, genomics encoder, and clinical encoder represent the intra-modal characteristics for their corresponding module.

2.3.2 Cross-Modal Attention Module

To simultaneously analyze pathology, genomics, and clinical data for a more holistic evaluation, the relationships between each data modality's tokens must be explored. The instance tokens of the pathology, genomics, and clinical data encoders are each denoted as a series. The cross-modal attention module then creates attention maps for each modality.

With the attention maps, it is now possible to extract genomics and clinical-related information in pathology tokens, pathology and clinical-related information in genomics tokens, and genomics and pathology-related information in clinical tokens. Such extraction is achieved by obtaining the product of the maps with the original instance token series multiplied by a learned parameter. This allows our model to leverage the cross-modal information for learning the complementary characteristics between the multiple modalities.

2.3.3 Decoders

In clinical situations, doctors can estimate gene expressions from pathological images or identify potential pathological phenotypes based on genomic data. The decoders are intended to imitate this process for translating cross-modal information. The pathology decoder has a similar structure to the genomics and clinical encoders, while the genomics and clinical decoders have similar structures as the pathology encoder.

Pathology decoder: The decoder contains two Multi-head Self-attention layers that are applied on the genomics-related and clinical-related information in pathology for information translation.

Genomics Decoder and Clinical Decoder: The process is similar to the pathology decoder, except a PPEG module is applied for information translation.

The class token in the outputs of each decoder represents the cross-modal representation learned from its module. These feature representations go through feature alignment and fusion to produce the final prediction result.

2.4 LLM-enhanced Knowledge Base

Medical LLMs can provide personalized clinical support and guide oncologists in providing cancer care. However, the challenges these medical LLMs face mainly include its lack of knowledge specializing in cancer and lack of access to relevant and most recent papers. These tend to cause hallucinations where medical LLMs generate seemingly plausible but actually unverified or incorrect information.

To address and reduce hallucinations during interactions between oncologists and large language models (LLMs), existing solutions commonly utilize Retrieval-Augmented Generation (RAG) frameworks. However, RAG's current chunking methods, such as fixed or semantic chunking, can often lead to retrieval of incomplete information, which lowers the quality of generated responses.

To ensure the quality of our responses, we developed our novel LLM-enhanced knowledge base. We also obtained recent domain-specific publications during our collaboration with medical experts and building cancer related knowledge into both vector and graph databases. Through these methods, we can minimize the chances of hallucinations occurring because of the external professional knowledge that supplements LLMs with critical information.

Our knowledge base built in collaboration with medical experts empowers our system's agents to generate more accurate and professional results. Since we want efficient retrieval for doctors in clinical environments and comprehensive retrieval for capturing the complex

relationships in cancer care, we combine vector and graph-based representations. By collaborating with professional oncologists from hospitals, we significantly improve the medical context engineering for our agents, supplementing them with the most relevant and professional resources for clinical support.

The LLM-enhanced expert knowledge base method involves three main steps: indexing, retrieval, and generation.

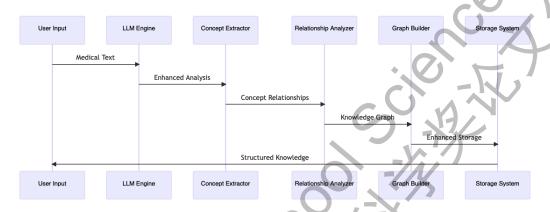


Figure 6: LLM-Enhanced knowledge base processing pipeline

2.4.1 Document Processing

The data indexing contains 4 main steps: text preprocessing, document structure analysis, intelligent chunking, and LLM empowered enhancement. Text preprocessing performs basic operations such as removing extra whitespace, standardizing punctuation, and cleaning the formatting. Next, the document structure analysis provides the basis for intelligent chunking by determining the properties of a given document. The LLM then performs intelligent chunking of the raw context text, such as an entire textbook or published paper, to improve efficiency. The LLM achieves intelligent chunking by integrating multiple methods, including original content chunks which preserves the source material, summary chunks which are AI-generated summaries, concept chunks which are categorized into medical concepts, and key points chunks which are the important information. Chunks are then encoded into vectors through the embedding model and embedded into the vector database. They are also used to construct the graph knowledge base with relationships. For each chunk, we perform LLM-powered enhancement. This involves extracting LLM concepts, analyzing the text difficulty, identifying the target audience, and extracting the key topics.

After the LLM performs chunking and enhancement, we move onto LLM-empowered extraction to construct a knowledge graph. The LLM first extracts concepts from 8 categories, including diseases, symptoms, treatments, medications, examinations, risk factors, prevention, and medical concepts. Afterwards, the LLM identifies 5 main types of relationships to build a graph representation, including causal, treatment, symptom, risk, and prevention. This relational information allows us to create graph nodes (concepts) and graph edges (relationships). The LLM-enhanced data previously mentioned is then added to the nodes and edges. Lastly, the deduplication function finds and merges identical entities and relations to reduce overhead for the graph operations, allowing for improved efficiency

In summary the 6 main steps to constructing the knowledge graph are: (1) extracting all concepts from categories, (2) identifying relationships using an LLM, (3) creating graph nodes (concepts), (4) creating graph edges (relationships), (5) adding metadata to the nodes and edges, and (6) validating the graph structure. This vector and graph-based indexing method allows for comprehensive information understanding and enhanced retrieval performance.

2.4.2 Data Retrieval

The Dual-Level Data Retriever contains two main strategies: high-level retrieval and low-level retrieval. Since we combine graph and vector representations, agents can gain an understanding of the relationships between entities during the retrieval process. The LLM-empowered retrieval process contains 3 steps: Query Keyword Extraction, Keyword Matching, and Incorporating High-Order Relatedness. Given a query, the LLM based retrieval algorithm extracts both local and global query keywords. Then, the LLM uses the efficient vector database to match local query keywords with potential entities and match global query keywords with relationships tied to global keys. Finally, the higher-order relatedness is maintained for the query by identifying neighboring nodes within local subgraphs of the retrieved graph elements. This allows for efficient retrieval with the vector representations and comprehensive retrieval with structured relational knowledge from the constructed graph.



Figure 7: LLM-enhanced cancer expert knowledge base.



Figure 8: Graph-based Medical Dataset.

We included 18 standard medical textbooks used in medical school, covering a wide range of diseases, to provide our agents with professional medical knowledge. We collaborated with oncologists and professors at the Cancer Hospital of Fudan University to locate and analyze hundreds of recent cancer guidelines and published papers from Pubmed, websites, textbooks, and databases. We also leveraged the National Comprehensive Cancer Network (NCCN) database and European Society of Medical Oncology (ESMO) database that frequently updates treatment guidelines for each type of cancer.

When addressing patient queries, LLM-enhanced expert knowledge base retrieves the sources that are most relevant to the query to provide context for the LLM. Then, it generates a natural language response through combining the user's questions and retrieved

chunks. This allows us to empower the agents for more accurate clinical support in cancer. With LLM-empowered knowledge graph understanding and construction, the combination of vector and graph-based representations, and the professional medical resources from doctors. The MedQA dataset contains multiple-choice questions derived from professional medical exams, such as the US Medical Licensing Examination (USMLE) and China's medical licensing exams, covering a wide spectrum of medical knowledge. Under MedQA test, our RAG accuracy improves output accuracy by around 17.2%, from 75.3% with the Regular Knowledge Base to 92.5% with the LLM-Enhanced Knowledge Base.

Metric	Regular Knowledge Base	LLM-Enhanced Knowledge Base
Concept Extraction	3 concepts	14+ concepts
Relationship Recognition	None	5 relationship types
Semantic Understanding	75.3% accuracy	92.5% accuracy
Knowledge Connectivity	Isolated	Networked graph, full connectivity

Figure 9: LLM-enhanced knowledge base accuracy improvements

2.4.3 Generation

With support from the LLM-enhanced expert knowledge base, our Multi MedAI Agent outputs the final response to the user's question in natural language. The user's natural language input at the beginning of the dialogue is synthesized with the most relevant chunks of information from the given sources into a final prompt, which is then given to the LLM. Given this re-invented prompt based on user query and information from the knowledge base, the LLM is tasked with formulating a natural response to the original question.

To ensure the accuracy and effectiveness of communication in our Multi MedAI Agent, the LLM also utilizes conversational history and previous dialogue as context. This allows the patient to have a continued, multi-turn dialogue, optimizing efficiency as the system will have strong familiarity with the users.

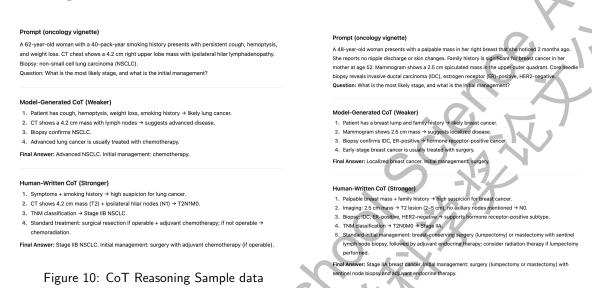
2.5 CoT & Reinforcement Learning

Chain of Thought (CoT) reasoning is a sequential generation process involving multiple reasoning steps (history synthesis, evidence linking, diagnostic hypotheses, etc.) that lead to a final clinical impression or recommendation. The process depends on the context, so models explore multiple reasoning paths before arriving at the final answer. To enhance the CoT reasoning of our agents, we explored three main ways. First is prompt integration, where the user enters the thinking step-by-step process as part of their prompt. Second is using few shots to give patient examples to the LLM to teach the thought process. Lastly, I am working on using Reinforcement Learning (RL) to encourage the simulation of clinician thought process.

The first strategy I used is directly integrating the CoT thought process into the prompt. I collaborated with professional doctors to get 30 CoT reasoning samples to use as the examples given to the agents. We add a description of each step of the chain of thought reasoning process into the prompt to ensure it follows the steps to arrive at an answer. The reasoning process

steps is generalizable to more patients and clinical situations but it aligns with how doctors would actually evaluate a patient as I collaborated with doctors in creating this chain of thought process description.

Sample data:



In addition to the directly integrating the CoT into the prompt to improve our agents' CoT reasoning process, the second strategy I used is few-shot prompting, which is when the user provides the agent with several examples of an ideal response to a similar question to help the agent learn the key characteristics of a successful answer. These reasoning samples clearly depict the steps that doctors take to arrive at their conclusion about a patient given a realistic scenario. For each prompt to the agent, we include around 3-5 sample CoT responses to ensure we maximize the utility of the context window. In addition to directly adding these examples to the prompt, we will also be exploring how to enhance the model's CoT process using these samples for RL training. I asked the clinicians I collaborated with to assess and compare 60 model-generated responses per clinician, allowing me to perform RL training.

Regarding the Reinforcement Learning training process I am working on for the Chain of Thought process, I will be using the Guided Reward Policy Optimization (GRPO) to train the model to follow the correct think-then-answer method. First, the agent explicitly conveys its reasoning process — enclosed within limits for clarity — and then presents its final diagnostic conclusions in a standard, easily extractable format. GRPO provides reward signals that incentivize both accurate predictions and adherence to the logic. This encourages both diagnostic correctness and transparency in the agent's reasoning process. By separating the reasoning and conclusion, I am working to ensure that each diagnostic output is accompanied by a clear, step-by-step explanation. This reasoning-centric approach enhances diagnostic transparency and enables explainable AI by providing clinicians with not only the final diagnostic conclusions but also the logical pathway through which these conclusions were reached.

3 Results

3.1 Co-Attention Heatmaps

For visualizing the genomic-related whole slide image embeddings, we use heatmaps. The coattention weights between genomics and pathology are overlayed with visual assessment from two pathologists for a low and a high risk case in the BRCA dataset [28]. The gradient from blue to red patches indicate the gene's attention weight from low to high. In each whole slide tissue image's heatmap, each patch is assigned a color along the gradient. There are 6 main types of genes we visualize in the following co-attention heatmaps: tumor suppressor genes, oncogenes, protein kinases, cell differentiation markers, transcription factors, and cytokines and growth factors.

Tumor suppressor genes are typical genes that slow down cell division or manage apoptosis. When they malfunction, cells begin to grow out of control which results in cancer [31]. Oncogenes are proto-oncogenes that have mutated and have the potential to cause cancer, leading to uncontrollable cell division [32]. Protein kinases mainly add phosphate groups to proteins, making them the key regulators of cellular processes [33].

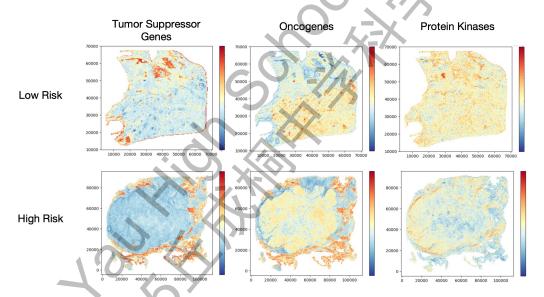


Figure 11: Co-Attention Visualization 1

Cell differentiation markers can indicate tumor behaviors, where less differentiated (immature) tumors are typically more aggressive [34]. Transcription factors ensure that the appropriate genes are expressed at the right location, at the correct time, and to the right degree [35]. Cytokines are signaling proteins that control inflammation, but an exceeding quantity damages tissues and leads to diseases like cancer [36]. Growth factors are groups of proteins that stimulate the growth of tissues; some cancer cells can produce growth factors that increase their own proliferation rate [37], [38].

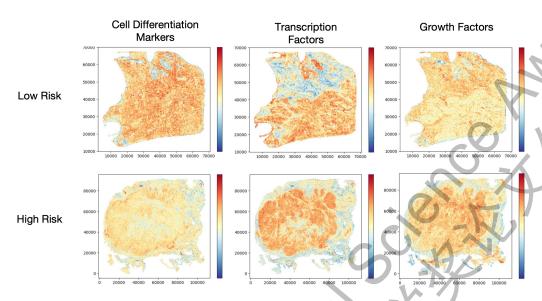


Figure 12: Co-Attention Visualization 2

3.2 Evaluation Metrics

After training the algorithm, we evaluated it using 3 methods. First, we used PathVQA [23] to test the multi-modal abilities of my model, using evaluation metrics recall and accuracy. Second, we tested the performance of our model on specialized tasks for various cancers. Thirdly, we collaborated with professors and oncologists and asked them to be referees in determining the quality of our model's natural language response and outputs in comparison with other models.

3.3 VQA Dataset Evaluation

3.3.1 Evaluation Dataset

The Visual Question Answering (VQA) dataset we used to evaluate our system's multi-modal analysis ability is PathVQA. PathVQA is a pathology evaluation dataset that includes 32799 Q&A pairs and 4998 pathology images.

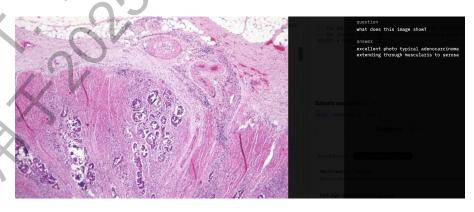


Figure 13: Example of Q&A in VQA datasets used during result evaluation [20]

3.3.2 VQA Evaluation Results

With our first evaluation method, we compared our model's accuracy and recall performance on PathVQA, with and without context, with other existing and state-of-the-art (SOTA) models. We achieved an overall average increase of around 9.7 10.5% accuracy compared to the SOTA LLaVA and LLaVA-Med AI algorithms.

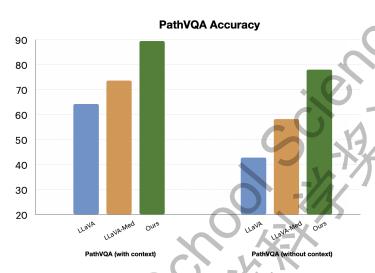


Figure 14: VQA test results compared with LLaVA and LLaVA-Med

3.4 Survival Prediction Evaluation

Second, we evaluated our models accuracy on survival prediction for breast cancer. By utilizing specialized models and LLM-enhanced expert knowledge base in our novel agent framework, our model outperforms existing unimodal and multimodal algorithms. We achieve up to 72.7% accuracy which is 3.6% higher than current state-of-the-art methods, evaluating images and genetic domain data. We also outperform the unimodal specialized models, by 10.3% compared to genetic only models and 6.3% compared to imaging only models.

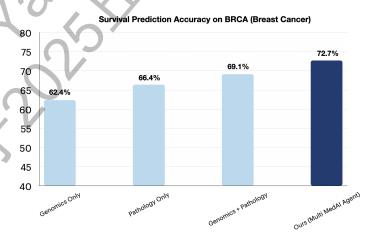


Figure 15: Survival prediction on breast cancer

3.5 Response Quality Evaluation

For our last evaluation method, we collaborated with 4 doctors and received their feedback on ten natural language responses generated by our model under the 3 categories: disease description, condition diagnosis, and medical advice. Disease description refers to the Multi MedAI Agent's ability to accurately extract information from given patient information. Condition diagnosis is the system's accuracy in analyzing the patients medical condition and giving a holistic evaluation and diagnosis. Medical advice refers to medical AI Agents' ability in providing reliable and detailed medical advice based on the patients situation, including treatment plans, recommended diet, activities, and more.

In a blind experiment where the oncologists I collaborated with compared the natural language responses from LLaVA, LLaVA-Med, and our Multi MedAI Agent, our model had overall better quality response ratings, evaluating based on the accuracy and language quality. We received a total of around 240 comparative evaluations from the four doctors we collaborated with.

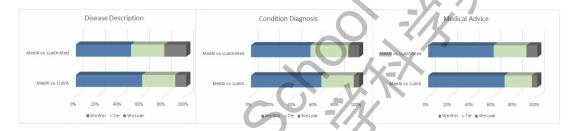


Figure 16: Natural language response quality compared with LLaVA and LLaVA-Med

In the 3 evaluations methods above, our Multi MedAI Agent Team outperformed current existing solution in both accuracy and natural language output quality.

4 Discussion

4.1 Ablation Studies

We also conducted an ablation study with a baseline of pathology and genomics, to determine the impact of early and late fusion of clinical data had on output accuracy. Late fusion refers to the integration of clinical data into the prediction process after the main information from the pathology and genomics data had been extracted and evaluated. Early fusion is the process of prediction where the clinical data is evaluated alongside with the pathology and genomics data. Through experimenting, we found that the late fusion of clinical data resulted in better prediction quality compared to early fusion. The clinical information for each patient is at a low dimensionality, while the features for pathology and genomics are at a much higher dimensionality. This difference results in the early fusion of clinical data posing a less significant role in the prediction outcome, leading to lower survival prediction accuracy.

Table 1: Ablation Study of Survival Prediction Performance with early and late fusion of clinical data.

Medical Domain	Prediction Accuracy(\uparrow)
Pathology + Genomics(Baseline) Pathology + Genomics + Clinical (Early Fusion)	$\begin{array}{c c} 69.1.5_{\pm 2.5} \\ 70.9_{\pm 2.3} \end{array}$
Pathology + Genomics + Clinical (Late Fusion) (OURS)	$72.7_{\pm 1.5}$

4.2 Future Work

The potential applications and enhancements of our Multi MedAI Agent system span various dimensions. A roadmap of these advancements, collaborations, and adaptations in real life is discussed below. When initiating our project, we communicated with doctors to understand their perspective on the main challenges faced in clinical situations for cancer patients, and learned that determining treatment plans and cancer care is one of the most critical challenges. Throughout our project, we maintained constant communication to learn whether our project aligns with clinical needs.

Further Collaborate with Doctors and Expand knowledge base

In the future, we wish to continue collaborating with hospitals to validate our model with more data analysis and assist doctors with selecting treatment plans. Through this communication with doctors and professors, we hope to acquire more domain-specific knowledge and expand the LLM-enhanced expert knowledge base to cover more cancers.

Enhance CoT and Reinforcement Learning in Reasoning Agent

In the future we hope to learn more about doctor's Chain of Thought (COT) to continue improving the thought process of our Reasoning Agent through reinforcement learning.

Implement in Cancer Community

We will provide Multi MedAI Agent to the cancer community to provide more assistance in analyzing and providing advice on cancer in real life scenarios. We also hope to contact users on these platforms to receive authentic feedback from real cancer patients or their family on to further improve.

5 Conclusion

In this research, we developed the Multi MedAI Agent, a novel multi agent framework aimed to provide more comprehensive and accurate analysis and interactions with natural language for doctors.

Our system can provide cancer care through holistic medical analysis combined with effective natural language communications with doctors about specific details of cancer condition.

Collaborating with doctors, we integrated our Multi MedAI Agents with local LLM (Deepseek 70b, GPT-OSS 120b), and we also developed a mobile phone app connected to our Multi MedAI Agent. This allows doctors to access the information and ask domain-specific questions with an accurate response. Data that can be uploaded through the interface in-

cluding genomics, pathology, radiology, clinical data, doctors' notes etc. Doctors can enter patient data to access assistance in deciding the treatment or therapy for the patient later on.

Our MMAKER (Multi MedAI Agent with Knowledge-Enhanced Reasoning) innovatively integrates 9 specialized agents, coordinated by meta-agents and a central reasoning agent. Our technical development of the 9 specialized agents is based on state-of-the-art models with different output information. By examining multiple sources of accurate low-level details from specialized agents, (4 for pathology, 3 for genomics, 1 for clinical records, and 1 for survival prediction), the higher-level analysis outputs can be more holistic and accurate. We also capitalize on the unique advantages of multi-agents compared to multi-modal systems, specifically how each agent can perform the delegated actions while communicating between each other in natural language for oncologists to understand the underlying reasoning behind the outputs.

Our LLM-enhanced Knowledge Base uses LLMs to perform document processing and data retrieval because they can understand longer natural language context, ensuring the construction of a high quality knowledge base. We constructed this knowledge base in collaboration with professional oncologists who uploaded hundreds the most updated cancer-related resources they use. In addition to the knowledge base, we also improve the Chain-of-Thought process with Reinforcement Learning that rewards both logical clarity and final correctness.

References

- Nagai, H., & Kim, Y. H. (2017, March). Cancer prevention from the perspective of global cancer burden patterns. Journal of thoracic disease. https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5394024/
- [2] Sung, H., Ferlay, J., Siegel, R. L., & Laversanne, M. (2021, May). Global cancer statistics 2020: Globocan estimates of incidence and mortality worldwide for 36 cancers in 185 countries. National Center for Biotechnology Information. https://pubmed.ncbi.nlm.nih.gov/33538338/
- [3] Adams, S. J., Stone, E., Baldwin, D. R., Vliegenthart, R., Lee, P., & Fintelmann, F. J. (2022, December 20). Lung cancer screening. Lancet (London, England). https://pubmed.ncbi.nlm.nih.gov/36563698/
- [4] Siegel, R. L., Miller, K. D., Wagle, N. S., & Jamil, A. (2023). Cancer Statistics, 2023. American Cancer Society. https://acsjournals.onlinelibrary.wiley.com/doi/full/10.3322/caac.21763
- [5] Pancreatic Cancer Survival Rates. Pancreatic Cancer Action Network. (2023, November 29). https://pancan.org/facing-pancreatic-cancer/about-pancreatic-cancer/survival-rate/
- [6] Liu, X., Yuan, P., Li, R., & Zhang, D. (2022, May 16). Predicting breast cancer recurrence and metastasis risk by integrating color and texture features of histopathological images and Machine Learning Technologies. Computers in Biology and Medicine. https://www.sciencedirect.com/science/article/abs/pii/S0010482522003614
- [7] Khader, F., Kather, J. N., Müller-Franzes, G., Wang, T., Han, T., Tayebi Arasteh, S., Hamesch, K., Bressem, K., Haarburger, C., Stegmaier, J., Kuhl, C., Nebelung, S., & Truhn, D. (2023, July 1). Medical Transformer for multimodal survival prediction in Intensive Care: Integration of imaging and non-imaging data. Nature News. https://www.nature.com/articles/s41598-023-37835-1
- [8] Chen, R. J., Lu, M. Y., Weng, W.-H., & Chen, T. Y. (2021). Multimodal Co-Attention Transformer for Survival Prediction in Gigapixel Whole Slide Images. Institute of Electrical and Electronics Engineers. https://ieeexplore.ieee.org/document/9710773/
- [9] Klambauer, G., Unterthiner, T., Mayr, A., & Hochreiter, S. (2017, September 7). Self-normalizing neural networks. arXiv.org. https://arxiv.org/abs/1706.02515
- [10] Jiang, A. Q., Sablayrolles, A., Mensch, A., Bamford, C., Chaplot, D. S., Casas, D. de las, Bressand, F., Lengyel, G., Lample, G., Saulnier, L., Lavaud, L. R., Lachaux, M.-A., Stock, P., Scao, T. L., Lavril, T., Wang, T., Lacroix, T., Sayed, W. E. (2023, October 10). Mistral 7B. ArXiv.org. https://doi.org/10.48550/arXiv.2310.06825
- [11] Touvron, Hugo, et al. "LLaMA: Open and Efficient Foundation Language Models." ArXiv:2302.13971 [Cs], 27 Feb. 2023, arxiv.org/abs/2302.13971.
- [12] Liu, H., Li, C., Wu, Q., Lee, Y. J. (2023, April 17). Visual Instruction Tuning, ArXiv.org. https://doi.org/10.48550/arXiv.2304.08485
- [13] Li, Chunyuan, et al. "LLaVA-Med: Training a Large Language-And-Vision Assistant for Biomedicine in One Day." ArXiv.org, 1 June 2023, arxiv.org/abs/2306.00890.
- [14] Li, M., Blaes, A., Johnson, S., Liu, H., Xu, H., Zhang, R. (2024). CancerLLM: A Large Language Model in Cancer Domain. ArXiv.org. https://arxiv.org/abs/2406.10459
- [15] Zhang, H., Chen, J., Jiang, F., Fu, Y., Chen, Z., Chen, G., Li, J., Wu, X., Zhang, Z., Xiao, Q., Xiang, W., Wang, B., Li, H. (2023). HuatuoGPT, Towards Taming Language Model to Be a Doctor. Association for Computational Linguistics. https://doi.org/10.18653/v1/2023.findings-emnlp.725
- [16] Lorica, Ben. "Best Practices in Retrieval Augmented Generation Gradient Flow." Gradient Flow, 19 Oct. 2023, gradientflow.com/best-practices-in-retrieval-augmented-generation/. Accessed 31 Jan. 2025
- [17] Ilse, M., Tomczak, J. M., & Welling, M. (2018, June 28). Attention-based deep multiple instance learning arxiv.org. arxiv. https://arxiv.org/pdf/1802.04712.pdf
- [18] Lu, M. Y., Williamson, D. F. K., Chen, T. Y., Chen, R. J., Barbieri, M., & Mahmood, F. (2020). Data Efficient and Weakly Supervised Computational Pathology on Whole Slide Images. ArXiv:2004.09666 [Cs, Eess, Q-Bio]. https://arxiv.org/abs/2004.09666
- [19] Ding, Y., Yang, F., Han, M., Li, C., Wang, Y., Xu, X., Zhao, M., Zhao, M., Yue, M., Deng, H., Yang, H., Yao, J., & Liu, Y. (2023, July 13). Multi-center study on predicting breast cancer lymph node status from core needle biopsy specimens using multi-modal and multi-instance Deep Learning. Nature News. https://www.nature.com/articles/s41523-023-00562-x
- [20] Lau et al. "Papers with Code VQA-RAD Dataset." Paperswithcode.com, 2022, paperswithcode.com/dataset/vqa-rad. Accessed 31 Jan. 2025.
- [21] Lau, J. J., Soumya Gayen, Demner, D., Asma Ben Abacha. (2018). Visual Question Answering in Radiology (VQA-RAD). OSF. https://doi.org/10.17605/OSF.IO/89KPS
- [22] Liu, B., Zhan, L.-M., Xu, L., Ma, L., Yang, Y., Wu, X.-M. (2021, February 18). SLAKE: A Semantically-Labeled Knowledge-Enhanced Dataset for Medical Visual Question Answering. ArXiv.org. https://doi.org/10.48550/arXiv.2102.09542
- [23] He, X., Zhang, Y., Mou, L., Xing, E., Xie, P. (2020). PathVQA: 30000+ Questions for Medical Visual Question Answering. ArXiv.org. https://arxiv.org/abs/2003.10286
- [24] Li, R., Wu, X., Li, A., & Wang, M. (2022, April 28). HFBSURV: Hierarchical multimodal fusion with factorized bilinear models for cancer survival prediction. National Center for Biotechnology Information. https://pubmed.ncbi.nlm.nih.gov/35188177/

- [26] Johnson, A., Bulgarelli, L., Pollard, T., Horng, S., Celi, L. A., & Mark, R. (2021, March 16). MIMIC-IV. PhysioNet. https://physionet.org/content/mimiciv/1.0/
- [27] Johnson, A. E. W., Pollard, T. J., Greenbaum, N. R., Lungren, M. P., Deng, C., Peng, Y., Lu, Z., Mark, R. G., Berkowitz, S. J., & Horng, S. (2019, November 14). Mimic-CXR-JPG, a large publicly available database of labeled chest radiographs. arXiv.org. https://arxiv.org/abs/1901.07042
- [28] TCGA Dataset. Cancer Imaging Archive. (n.d.). https://nbia.cancerimagingarchive.net
- [29] Pan-Cancer Datasets. pan-cancer-dataset-sources. (n.d.). https://lab-rasool.github.io/pan-cancer-dataset-sources/
- [30] Ding, K., Zhou, M., Wang, H., Gevaert, O., Metaxas, D., & Zhang, S. (2023, April 21). A large-scale synthetic pathological dataset for deep learning-enabled segmentation of breast cancer. Nature News. https://www.nature.com/articles/s41597-023-02125-y
- [31] Oncogenes, Tumor Suppressor Genes, and DNA Repair Genes. American Cancer Society. (2022, August 31). https://www.cancer.org/cancer/understanding-cancer/genes-and-cancer/oncogenes-tumor-suppressor-genes.html
- [32] Liu, P. P. (2024, January 12). Oncogene. Genome.gov. https://www.genome.gov/genetics-glossary/Oncogene
- $[33] \begin{tabular}{ll} Protein Kinases. Cell Signaling Technology. (n.d.). $https://www.cellsignal.com/learn-and-support/protein-kinases. (n.d.). $https://www.cellsignal.com/learn-and-support/prot$
- [34] Jögi, A., Vaapil, M., Johansson, M., & Påhlman, S. (2012, May). Cancer cell differentiation heterogeneity and aggressive behavior in solid tumors. Upsala journal of medical sciences. https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3339553/
- [35] Robertson, S. (2019, February 26). Role of transcription factors. News Medical Life Sciences. https://www.news-medical.net/life-sciences/Role-of-Transcription-Factors.aspx
- $[36] \quad \text{Cytokines. Cleveland Clinic. (2023, January 3). my. cleveland clinic. org/health/body/24585-cytokines.}$
- [37] Growth Factor. Encyclopædia Britannica. (2019, February 28). https://www.britannica.com/science/growth-factor
- [38] National Center for Biotechnology Information. (n.d.). The Development and Causes of Cancer. National Library of Medicine. https://www.ncbi.nlm.nih.gov/books/NBK9963/

6 Acknowledgements

I would like to express my gratitude to Professor WenGuang Chen (uncompensated) from Tsinghua University for his valuable feedback and guidance throughout our project. The relationship was that of mentor and mentee, focused on project direction decision, technology selection, and scientific review. Professor Chen provided guidance on the state-of-the-art advancements in artificial intelligence, suggested system improvements, reviewed the manuscript, and ensured academic rigor, but was not involved in programming, data analysis, or experimental execution. All manuscript writing and technical implementation remained under the sole responsibility of me.

I would also like to thank Dr. Xiang Li, professor at Harvard Medical School and Mass General Brigham, for his advice on my multi-agent framework, strategies to reduce hallucination, and guidance in helping me learn about the latest AI agent applications in medical scenarios.

On the clinical side, I would like to thank Dr. ZhenBing Shen and Dr. Xue Feng (both uncompensated) from the hospital of Fudan University for their continuous support in collaborating with us to develop the research on cancer in medical practices as well as providing feedback on our system's logic ensuring the research remained grounded in real-world medical needs. .

In this project, I developed the initial project idea with guidance from my advisor Professor WenGuang Chen. The data used for training and result evaluation purposes was obtained from open-sourced public cancer dataset, TCGA. I completed the data cleaning from TCGA. For the response quality evaluation section of our result analysis, we received feedback with the doctor's we collaborated with. I wrote this research paper with the support of my advisor.

The entire research process was designed and implemented through a series of distinct, transparent stages:

The research topic was selected after in-depth discussions with oncologists from the Fudan University Cancer Center, combined with the investigator's personal observation of the challenges faced by cancer patients and families. Key barriers identified by clinicians included difficulties in risk stratification, the complexity of integrating multimodal data, and the necessity of personalized prognoses. These insights guided the development of the MMAKER system, emphasizing multimodal analysis and intelligent clinical decision support.

Data for model development and evaluation was entirely sourced from open-access repositories, particularly The Cancer Genome Atlas (TCGA). Pathology images, genomics data, and patient clinical records were curated, cleaned, and pre-processed by me. Additional benchmarking datasets for model validation, such as the PathVQA and MedQA set for visual question answering, were likewise independently acquired and managed.

The project's methodology was built around a multi-agent system architecture featuring nine specialized analyst agents: four for pathology, three for genomics, one for clinical analysis, and one for survival prediction. Each agent performed distinct analytic functions, coordinated through meta-agents and a central reasoning agent emulating multi-disciplinary team decision-making. The specialized survival prediction tool, MedPred, was developed using cutting-edge transformer architectures to synthesize information from pathology, genomics, and clinical records.

The investigator was responsible for all aspects of algorithm design, parameter selection,

data pipeline construction, cross-modal attention module development, coding, and model training. Techniques such as reinforcement learning from human feedback (RLHF) and Chain-of-Thought reasoning modules were designed and implemented based on clinical interviews and direct collaboration with medical experts.

I authored the manuscript, structured results, and synthesized clinical and technical evaluations, with advisors providing review for accuracy, clarity, and scientific rigor.

Every technical, analytical, and written element was performed by me. Advisors and medical experts contributed feedback, guidance, validation, and domain-specific expertise, but were not leading any hands-on coding, experimental operation, or manuscript writing.

The key challenges and solutions in this project include:

- Data Heterogeneity and Imbalance: Successfully addressed by innovating modular fusion and cross-modal attention strategies, ensuring robust synthesis across diverse data types.
- Scarcity of Matched Multimodal Data: Mitigated through open-source data mining, augmentations, and simulated experiments.
- Risk of AI Hallucination: Combated by implementing a knowledge base combining vector and graph-based representations and supplementing results through RL-powered reasoning transparency, with continuous feedback from clinical partners.
- Alignment of Clinical Needs: Maintained through regular expert feedback and validation, resulting in iterative system refinement to address real clinical workflows.
- Generalization and Robustness: Ablation studies and comprehensive evaluation against state-of-the-art models confirmed the system's improved accuracy, efficacy, and reliability in cancer survival prediction and clinical decision support.

Throughout the research, constant communication with medical professionals guided the technical direction, assured real-world relevance, and verified the alignment of results with clinical needs. The MMAKER system demonstrated strong performance in benchmarking and clinical evaluations and provides a scalable foundation for further collaboration, system enhancement, and real-world deployment.