

# 分布式医学数据管理与智能分析系统的设计与性能优化

田行健<sup>1</sup>, 陈子扬<sup>1</sup>, 王俊翔<sup>2</sup>, 张勉<sup>3\*</sup>

<sup>1</sup>北京建筑大学理学院, 北京

<sup>2</sup>中国科学院计算技术研究所, 北京

<sup>3</sup>北京建筑大学智能科学与技术学院, 北京

收稿日期: 2026年4月15日; 录用日期: 2026年5月13日; 发布日期: 2026年5月22日

## 摘要

随着医疗数据的爆炸式增长, 传统单体架构难以兼顾海量数据管理与实时智能分析。文章以复杂医疗AI应用场景下的分布式系统为对象, 报道一项综合设计与性能优化实践: 给出分布式医学数据管理与智能分析系统的完整工程案例, 采用微服务与网关聚合, 面向医生、患者与管理三类用户, 集成数据采集标注、多条件检索、异步批处理与实时医患会话。业务侧部署五个深度学习推理微服务: 胸片X光疾病诊断、胸片X光报告生成、超声心动图分割与射血分数计算、尿检肾功能诊断与智能中医舌诊, 并辅以基于LangChain的检索增强生成(RAG)对话。基础设施采用Nacos、Sentinel、API网关与OpenFeign, 结合Redis、Kafka与MySQL; 推理生产路径采用TensorRT。性能方面, 按分层优化路线落地七项策略: 多级缓存、分布式锁与幂等、数据库访问与分页、异步消息、高可用与服务治理、推理引擎部署及JVM调优。在Docker Compose三节点集群、单表约 $1.2 \times 10^7$ 行、JMeter压测下, 主要结果包括: 深度分页经联合索引与键集分页后, 平均响应由约2.41 s降至16 ms, P99由约3.86 s降至43 ms, 扫描行数由千万级降至千量级; 读多写少场景在200并发下, Redis使QPS由428升至3265、P99由318 ms降至84 ms, 命中率约91.6%; 约 $10^5$ 条批量删除同步接口平均约58.7 s, Kafka异步提交约0.28 s; TensorRT FP16单卡平均延迟约14 ms、P99约26 ms; 模型网关异步非阻塞饱和前稳定QPS约331, 优于同步线程池约126; 网关至业务链路六档并发(150~900用户)下单节点CPU由约24%升至约89%, JVM堆约4.0~4.4 GB。

## 关键词

医学信息系统, 微服务, 系统设计, 性能优化, 微服务治理

## Design and Performance Optimization of a Distributed Medical Data Management and Intelligent Analysis System

\*通讯作者。

文章引用: 田行健, 陈子扬, 王俊翔, 张勉. 分布式医学数据管理与智能分析系统的设计与性能优化[J]. 计算机科学与应用, 2026, 16(5): 110-124. DOI: 10.12677/csa.2026.165169

Xingjian Tian<sup>1</sup>, Ziyang Chen<sup>1</sup>, Junxiang Wang<sup>2</sup>, Mian Zhang<sup>3\*</sup>

<sup>1</sup>School of Science, Beijing University of Civil Engineering and Architecture, Beijing

<sup>2</sup>Institute of Computing Technology, Chinese Academy of Sciences, Beijing

<sup>3</sup>School of Intelligent Science and Technology, Beijing University of Civil Engineering and Architecture, Beijing

Received: April 15, 2026; accepted: May 13, 2026; published: May 22, 2026

## Abstract

With the explosive growth of medical data, traditional monolithic architectures struggle to balance massive data management with real-time intelligent analysis. This paper examines distributed systems in complex medical AI application scenarios and reports on a comprehensive design and performance optimization practice: presenting a complete engineering case of a distributed medical data management and intelligent analysis system that adopts microservices and gateway aggregation. The system serves doctors, patients, and administrators with data collection and annotation, multi-criteria search, asynchronous batch jobs, and real-time messaging. Five deep-learning inference microservices are deployed under a single naming scheme: chest X-ray disease diagnosis, chest X-ray report generation, echocardiogram segmentation and ejection-fraction estimation, urinalysis-based kidney function diagnosis, and intelligent tongue diagnosis in traditional Chinese medicine, plus a LangChain-based retrieval-augmented generation (RAG) assistant over an institutional knowledge base. The stack combines Nacos, Sentinel, API gateways, and OpenFeign with Redis, Kafka, and MySQL; TensorRT serves production inference. Following a layered optimization plan, we realize seven concrete strategies: multi-level caching, distributed locks and idempotency, database access and pagination, asynchronous messaging, high availability and governance, inference deployment, and JVM tuning. On a three-node Docker Compose cluster, a  $\sim 1.2 \times 10^7$ -row table, and JMeter workloads, we report among others: keyset paging and composite indexes cut deep-pagination mean latency from  $\sim 2.41$  s to 16 ms and P99 from  $\sim 3.86$  s to 43 ms, and scanned rows from tens of millions to thousands; in read-heavy, write-light scenarios under 200 concurrent users, Redis lifts QPS from 428 to 3265 and reduces P99 from 318 ms to 84 ms ( $\sim 91.6\%$  hit rate);  $\sim 10^5$  bulk deletions average  $\sim 58.7$  s synchronously versus  $\sim 0.28$  s for Kafka-backed submission; TensorRT FP16 achieves  $\sim 14$  ms mean and  $\sim 26$  ms P99 per request; the model gateway reaches  $\sim 331$  stable QPS before saturation under async scheduling versus  $\sim 126$  for a synchronous pool; gateway-to-business CPU rises from  $\sim 24\%$  to  $\sim 89\%$  across six load levels, with JVM heaps near 4.0~4.4 GB.

## Keywords

Medical Information System, Microservice, System Design, Performance Optimization, Microservice Governance

Copyright © 2026 by author(s) and Hans Publishers Inc.

This work is licensed under the Creative Commons Attribution International License (CC BY 4.0).

<http://creativecommons.org/licenses/by/4.0/>



Open Access

## 1. 引言

在智慧医疗与智慧医院建设持续推进的背景下, 院内信息系统正从“以科室为中心”的烟囱式建设, 转向跨科室、跨机构的数据互联与业务协同[1][2]。电子病历、影像归档与通信系统(PACS)、检验与病理信息系统等来源产生的数据规模持续膨胀, 医生端既要完成合规留痕与质控抽查, 又要在门诊与住院场

景中快速检索历史检查、调阅影像与报告；管理端则需要账号与权限治理、批量数据运维与审计追踪。与此同时，深度学习辅助诊断、报告生成与医学知识问答等能力逐步从科研原型走向科室试点，客观上要求后端既能承载传统事务型业务，又能以可观测、可限流、可扩容的方式托管 GPU 推理与检索编排服务。若仍沿用单体或粗粒度分层架构，往往在连接池耗尽、长事务阻塞 HTTP 线程、模型版本发布耦合等方面集中暴露问题，难以同时满足“数据侧高吞吐与低尾延迟”和“智能侧异构运行时与弹性伸缩”两类约束。

传统医疗信息系统在实践中还常面临三类工程痛点。其一是数据与流程割裂：同一患者在不同子系统中标识、时间线与授权范围若未统一建模，容易形成检索口径不一致与重复录入，既影响临床效率，也增加质控与审计成本。其二是峰值并发与长尾延迟：挂号高峰、批量导出、报表查询与影像调阅叠加时，数据库与网关层若缺少缓存、分页策略与流控配合，易出现 P99 延迟抬升乃至级联超时。其三是智能化能力的“服务化”边界不清：推理与 RAG 能力若与核心业务进程混布，GPU 资源争用、依赖版本升级与回滚粒度都会变得困难；反之，若完全独立又缺乏统一的鉴权、路由与可观测入口，则运维与排障成本上升。上述问题并非单一算法所能解决，更多依赖清晰的架构拆分、治理组件选型与贯穿全链路的性能工程。

针对上述背景，本文从系统设计与性能优化两条主线出发，以可复现、可借鉴的系统蓝图为目标，完整呈现一套面向复杂医疗 AI 场景的高性能分布式医学数据管理与智能分析系统实践：在业务维度覆盖医生、患者与管理员三类角色的闭环功能；在智能维度将胸片 X 光疾病诊断、胸片 X 光报告生成、超声心动图分割与射血分数计算、尿检肾功能诊断与智能中医舌诊五类场景各对应一个深度学习推理微服务，并另行提供基于 LangChain 的 RAG 对话服务，对机构内已落库的知识库执行检索与答复编排；在基础设施维度采用 Nacos 完成注册与配置、Sentinel 完成限流熔断、API 网关完成统一接入，Java 业务服务与 Python 推理/RAG 服务之间通过 OpenFeign 等机制协同，Redis、Kafka 与 MySQL 分别承担缓存、异步任务与事务型持久化。与侧重单一模型精度或单一中间件评测的工作不同，本文强调部署拓扑、接口契约、压测方法与量化指标之间的可核对关系，使“架构描述 - 优化策略 - 实验数据”形成闭环，便于同行在相似约束下复现或裁剪迁移。

本文工作从工程上可概括为三方面：1) 混合微服务拓扑与交付路径。给出 Spring Boot 业务集群与 Flask 推理集群共存的落地部署形态[3] [4]，说明网关、注册发现、声明式调用与容器编排如何共同支撑多副本与滚动发布，并交代与 Kubernetes、Docker Compose 相关的环境一致性要点。2) 端到端功能刻画与界面佐证。按角色梳理功能模块与典型交互路径，辅以主要业务页面示意，突出权限、分页检索、推理任务状态与医患会话等工程细节，使叙述停留在可对照实现的系统层面而非概念罗列。3) 可复现的七项优化措施与对比实验。围绕多级缓存、分布式锁与幂等、数据库索引与键集分页、Kafka 异步化、基于 Nacos 与 Sentinel 的高可用与服务治理、TensorRT 推理部署以及 JVM 调优，给出与生产约束相近的压测设定，并在深度分页、缓存命中、批量删除、推理延迟与网关吞吐、业务链路 CPU 与堆占用等维度给出前后对照或分档采样指标，为同类系统提供可借鉴的性能工程参考。

论文其余部分安排如下：第 2 节综述相关工作；第 3 节描述系统功能、总体架构与技术栈及部署运维(含分层技术栈示意)、数据安全与合规性设计要点，再给出主要界面示例；第 4 节给出面向医疗大数据的分层性能优化思路及七项策略的设计要点；第 5 节给出实验环境、方法与量化结果；最后总结全文并展望后续工作。

## 2. 相关工作

医疗信息化领域长期聚焦 HIS、PACS 与互联互通标准下的数据整合与流程闭环，代表性讨论可见文

献[1][2];近年来深度学习在辅助诊断与报告生成中的角色日益突出,工程上需与既有事务系统协同部署[5]。云原生与微服务为横向扩展与独立发布提供了通用范式[3][4],但在医疗场景中,服务拆分与治理策略往往仍从通用互联网业务迁移而来,针对高并发读路径、批量运维任务与 GPU 推理混部的一体化讨论相对分散[6]。工业界实践多体现在技术博客与开源组件文档,学术文献中同时覆盖“事务型数据管道 + 多模态推理微服务 + 可复现压测指标”的系统级完整案例仍不多见。与上述工作相比,本文在统一网关与注册发现前提下,将读多写少缓存、异步批处理、数据库分页与索引、推理引擎部署及 JVM 调优串联为可验证链路,并在第 5 节给出与配置一致的量化对照,以工程实践案例的形式呈现复杂医疗 AI 场景下高性能分布式系统的设计与优化结论,供同类项目对照复现。

### 3. 系统功能与架构设计

本系统面向医院或医联体场景下的医学数据集中管理与辅助分析需求,在功能上覆盖三类用户角色的日常操作闭环,在架构上采用“网关 - 业务服务 - 推理与知识服务 - 中间件”分层模型,并通过容器与 Kubernetes 完成环境一致性与发布自动化[4]。下文依次说明各角色功能边界、智能分析模型概要、总体架构与技术栈及运维交付方式,并从工程设计角度归纳数据安全、隐私保护与合规性要点,再给出主要界面示例;其中性能与容量相关的工程化手段在第 4 节集中论述。

#### 3.1. 用户角色与功能

##### 1) 医生端

医生端承担病例生产、质控与随访协同的主路径。账户与权限方面,支持基于角色的菜单可见性与接口鉴权,会话与令牌刷新策略与网关及统一认证服务对接。数据管理子系统支持影像、报告文本、检验与体征等多模态数据的录入、版本留痕与精细化标注;列表与检索支持按患者标识、时间区间、检查类型、标注状态等条件组合过滤,分页与导出接口与后台病例库及对象存储策略相配合,支持按筛选结果批量导出为 Excel 以便离线统计。智能分析子系统对外暴露五个独立的推理微服务入口,分别对应胸片 X 光疾病诊断、胸片 X 光报告生成、超声心动图分割与射血分数计算、尿检肾功能诊断与智能中医舌诊;医生在界面中上传或点选已有检查数据后,由前端携带任务参数调用对应服务接口,任务状态与结果在页面侧轮询或推送展示。知识问答子系统提供基于 LangChain 实现的 RAG 对话页面,对机构内维护的医疗知识文本库执行检索与答复编排,供医生在书写报告或宣教时查阅。医患协同子系统基于 WebSocket 承载会话房间与消息送达,与医生排班或科室维度权限结合。运维类能力包括接收管理员广播通知、对大批量历史数据执行异步批量删除或归档任务等,避免长事务阻塞交互线程。

##### 2) 患者端

患者端强调低门槛注册登录、个人资料与就诊相关基础信息维护,以及面向复诊与随访场景的在线沟通入口。患者可查看医生侧授权范围内的检查结论摘要与健康提示,在许可范围内发起图文或实时会话;部分健康自评或量表填写能力可作为独立表单模块接入,与医生端诊断结论区分展示,避免混用术语。

##### 3) 管理员端

管理员端面向信息科与业务科室管理员,提供医生与患者账号的全生命周期管理、科室与角色绑定、密码与登录策略配置等。数据监管模块支持按时间、医生、病种或数据质量标签过滤病例列表,用于抽查标注一致性与完整性;支持将全量或增量数据导出为离线包并登记备份批次。系统维护模块负责公告发布、任务队列与消费积压的告警阈值配置入口,以及与对象存储桶、Kafka 主题命名等业务运维参数的查看与变更审批衔接。

### 3.2. 智能分析服务与模型概要

智能分析子系统所集成的五条推理服务在算法形态上分别对应公开深度学习范式，下文按与系统命名一致的顺序作结构级概述，不主张提出新网络。胸片 X 光疾病诊断采用多模态双 Transformer (MMBT) 范式[7]，视觉支路以 ResNet 卷积主体产生图像 token，文本支路采用 BERT 式子词嵌入[8]，在统一序列上联合编码后完成分类。胸片 X 光报告生成采用 Faster R-CNN [9]在预定义解剖区域上提取区域特征，文本侧以类 GPT-2 的自回归 Transformer [10]生成描述。超声心动图分割与射血分数计算以三维残差网络 R3D [11]编码超声心动短片时空特征，射血分数(EF)回归与心肌分割/关键点坐标预测结合全连接与图卷积网络[12]。尿检肾功能诊断在几何与统计预处理后的检验指标特征上用轻量多层感知机完成多任务二分类，浅层非线性逼近的经典讨论见文献[13]。智能中医舌诊以 ImageNet 预训练 ResNet-18 [14]为舌象编码器，结合跨模态注意力与双向 LSTM [15]及 Transformer 式多头注意力[16]对文本与证候属性建模。知识问答侧采用检索增强生成(RAG) [17]对机构知识库检索与答复编排。上述实现与公开方法对齐，工程上由 TensorRT 与第 4 节所述优化策略保障推理延迟与吞吐。

### 3.3. 系统架构、技术栈与部署运维

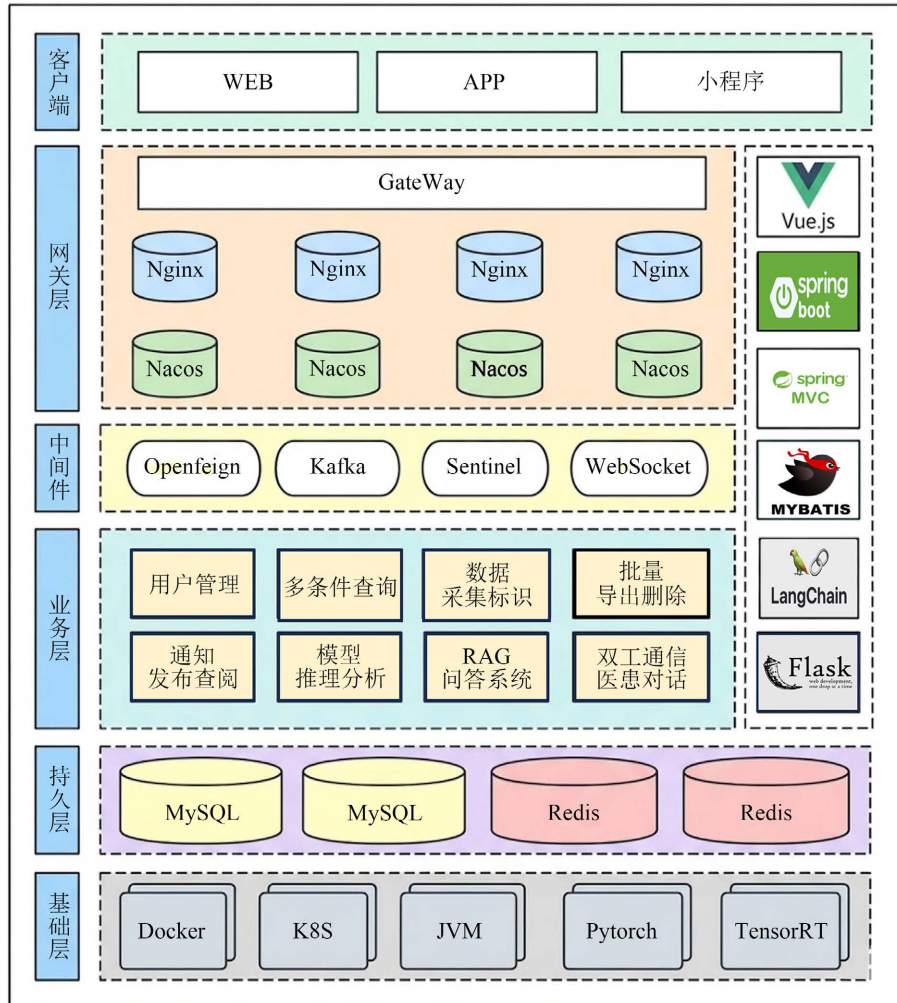


Figure 1. Layered technology stack  
图 1. 分层技术栈示意

## 1) 总体架构

系统采用分布式微服务架构与多节点集群部署：业务服务、推理服务、注册中心与消息组件均在物理或逻辑上可多副本运行，任一单点失效时由健康检查摘除实例并由负载均衡将流量迁移至存活副本。纵向分为接入层、业务服务层、智能服务层与中间件层，整体遵循微服务拆分与网关聚合的常见实践[3][6]。接入层由 Nginx 或云厂商负载均衡承担 TLS 终结、静态资源托管与反向代理，将 HTTP 与 WebSocket 连接分发至后端；同一服务名在后端对应多个 Spring Boot Pod 或虚拟机进程。业务服务层基于 Spring Boot 实现账户、组织机构、病例元数据、权限与审计日志、异步任务编排等聚合能力，对外以 RESTful 风格暴露资源；对五个 Flask 推理微服务及 RAG 对话服务的调用经 OpenFeign 声明式客户端发出，由注册中心解析实例列表并完成客户端负载均衡与失败重试。智能服务层包含五个 GPU 推理微服务进程与独立部署的 RAG 对话服务进程，前者通过模型加载与推理线程池消费任务，后者通过 LangChain 管线读取机构知识库并完成对话响应；深度学习模型训练与导出遵循通用范式[18]。中间件层中，MySQL 存储事务型业务数据与任务状态；Redis 承担会话辅助数据、热点缓存与分布式锁；Kafka 承担批量删除、导出等高耗时任务的异步投递；Nacos 集群承担服务注册、配置下发与健康心跳；Sentinel 以资源维度嵌入网关与核心服务，实现限流与熔断。前后端会话采用 JWT 或网关签发的令牌传递，WebSocket 连接在建立阶段完成同一套身份校验。

## 2) 技术栈

前端采用 Vue.js 实现单页应用与组件化页面，通过 Axios 调用后端 API，路由级权限与按钮级权限与后端返回的菜单树及策略码一致；静态资源由 Nginx 托管或经构建流水线注入镜像。Java 侧采用 Spring Boot、Spring MVC 与 MyBatis 完成 Web 层、事务边界与 SQL 映射；DTO 校验与全局异常处理统一封装以降低接口歧义。Python 侧推理与 RAG 服务基于 Flask 或兼容 ASGI 的部署方式暴露 HTTP 接口，深度学习模型依赖 PyTorch 训练与导出，生产推理路径经 TensorRT 引擎降低延迟与显存占用。表格导入导出采用 EasyExcel 等库处理大文件流式读写。观测方面，业务与推理进程输出结构化日志，可与 Prometheus 指标及链路追踪标识符对接，便于与第 5 节压测数据对照。图 1 给出分层技术栈示意图。

## 3) 部署与运维

研发与测试环境将各服务构建为 Docker 镜像[19]，使用 Docker Compose 描述依赖顺序、端口映射与卷挂载，保证与生产一致的文件路径与环境变量命名。生产环境采用 Kubernetes 进行编排[4]：无状态服务使用 Deployment 声明副本数与滚动更新策略，利用 Service 与 Ingress 对外暴露 ClusterIP 或负载均衡入口；有状态中间件可采用 Operator 或托管云服务，与业务集群同 VPC 互联。配置与密钥通过 ConfigMap 与 Secret 注入，敏感信息不入镜像层。GPU 节点通过节点标签与 Pod 的 resources.limits 声明显卡资源，调度器将推理类 Pod 约束至 GPU 节点。服务注册与流量治理组件可随业务同集群部署，亦可独立集群部署并通过网络可达地址注册，由运维在可用性与变更窗口之间权衡。发布流程上支持金丝雀或蓝绿策略与自动回滚钩子，与数据库迁移脚本的版本化执行顺序相衔接。

### 3.4. 数据安全、隐私保护与合规性设计

系统中的患者诊疗数据与个人敏感信息需在机密性、完整性与可用性之间取得工程上可落地的平衡。本小节从信任边界、身份与授权、韧性防护、密钥与凭证、存储与后续增强五方面说明当前实现的设计取向；其中可用性侧限流、熔断与客户端重试退避与第 4 节策略(5)相呼应。

#### 1) 信任边界与接入收敛

对外仅通过 Nginx 云负载均衡暴露 HTTPS 与 WebSocket 入口，由接入层完成 TLS 终结与反向代理，业务与中间件管理端口不直接面向公网，从而收敛攻击面并统一落点安全策略(证书、转发头、超时与限流协同)。静态资源与 API 请求可经路径拆分或独立域名策略管理，便于与院内防火墙、VPC 隔离及后续

WAF 策略衔接。

### 2) 身份鉴别、授权与前后端一致性

用户登录后采用 JWT 网关签发令牌贯穿 REST 与 WebSocket 建连；业务层按角色 - 资源粒度实施接口鉴权，前端路由守卫与按钮级显隐仅改善体验，以服务端校验为最终裁决。医生、患者与管理员的数据可见范围遵循最小必要原则：例如患者端仅展示经授权的摘要与随访信息(见第 3 节用户角色描述)。账户口令不在库内明文保存，采用单向哈希安全存储方式；其余业务字段当前以受控访问的关系型存储为主，依赖网络隔离、操作系统与数据库账号权限及备份策略由运维侧统一管理。

### 3) 韧性防护与滥用缓解

网关与核心服务嵌入 Sentinel，按资源维度实施限流、熔断与降级，在峰值挂号、批量导出或异常重试叠加时抑制雪崩，保障核心业务可读路径可用；服务间调用配置超时与有限次重试并结合指数退避。前端对高频触发操作采用防抖等交互策略，降低无效请求与误触风暴，与上述服务侧治理形成端 - 边 - 服务协同。

### 4) 密钥、配置与审计

Kubernetes 部署下敏感配置经 Secret 注入，密钥与口令不写入镜像层或公开仓库；Nacos 等组件中的连接串与业务密钥按环境隔离。业务服务聚合审计日志能力，支撑账号治理、导出与批量运维等行为的可追溯需求，与管理员端监管与备份批次登记相衔接。结构化日志与指标可与机构现有运维平台对接，便于事后排查与容量告警。

### 5) 行业合规映射与存储侧安全落地

网络安全与数据治理相关法律对个人敏感信息、重要数据及关键信息基础设施等场景普遍提出分类分级、技术与措施相结合、处理活动可说明等要求[20]；卫生健康领域亦强调诊疗数据在采集、存储、使用与对外提供中的权限边界与留痕。本文前述设计可与之对齐到工程控制点：接入层 TLS 与反向代理对应传输与边界防护；统一鉴权、RBAC 与最小必要展示对应访问控制与目的限定；审计日志与导出/备份登记对应关键行为可追溯；Secret 与密钥不入镜像对应防止敏感配置泄露。存储侧安全已在当前交付环境中闭环落地，按结构化事务数据/非结构化附件分层组织。MySQL 部署于院内或集群内网地址，业务进程仅通过连接池与最小必要权限账号访问；数据库连接串与业务密钥由第 3.3 节所述 Secret 与配置隔离机制注入，不落盘于镜像与公开仓库。账户口令以加盐哈希写入库表，无明文存储。影像、导出包等大文件由应用层经鉴权后的读路径提供，与登录态或短时授权绑定，避免匿名直链与公网裸读。备份与批量导出在业务侧完成批次登记并可追溯，与管理员端监管流程一致；运维侧对快照、离线介质与保留周期按机构规范执行。综上，存储侧已形成可信网络内受控落盘 - 凭据与访问链路统一治理 - 敏感读路径鉴权 - 备份留痕的实装能力；若后续需对接等保或云平台托管能力，可再叠加与本系统解耦的卷级/透明加密或 KMS 等基础设施侧加固。

## 3.5. 主要界面示例

图 2 为医生工作台中“胸片 X 光报告生成 AI 分析结果复核”界面。页内展示诊断置信度进度条及“高置信度：AI 诊断结果可置信度较高”提示；中部 AI 标注图像在胸片影像上叠加绿、蓝等选框(如左下肺、右下肺、心影等区域)，并提供“隐藏边框”、“修改选中项”以调整标注。右侧为 AI 总体诊断自然语言摘要，其下按解剖条目列出分条结论(如纵隔、上纵隔、心脏轮廓等)，各条附置信度及“已复核”状态。页底提供“提交复核结果并入库”、“拒绝入库”、“重新诊断”，支撑人机协同复核闭环。图 3 为同系统下“问诊信息”弹窗：左侧为结构化字段，包括问诊单号、医生 ID 等；右侧为胸片缩略图，以及 AI 诊断结果——含 AI 总体诊断长文本与按区域拆分的结论卡片。上述界面由前端路由调用后端 REST 接口；权限不足时由网关返回统一错误码。

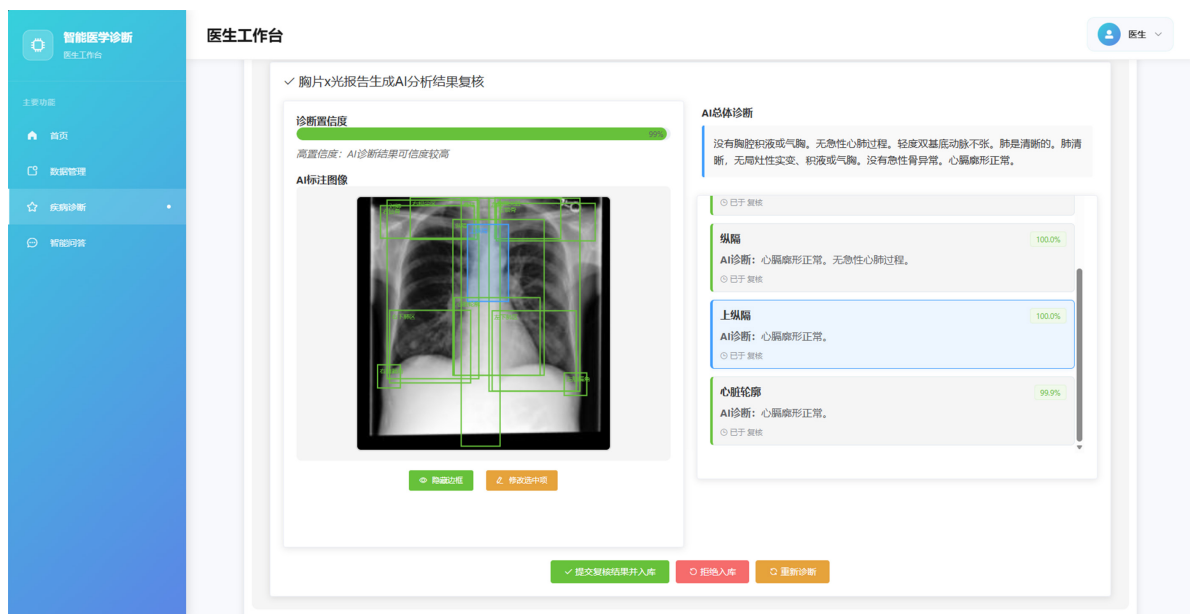


Figure 2. Chest X-ray report generation: AI analysis review interface (physician workstation)

图 2. 胸片 X 光报告生成: AI 分析结果复核界面

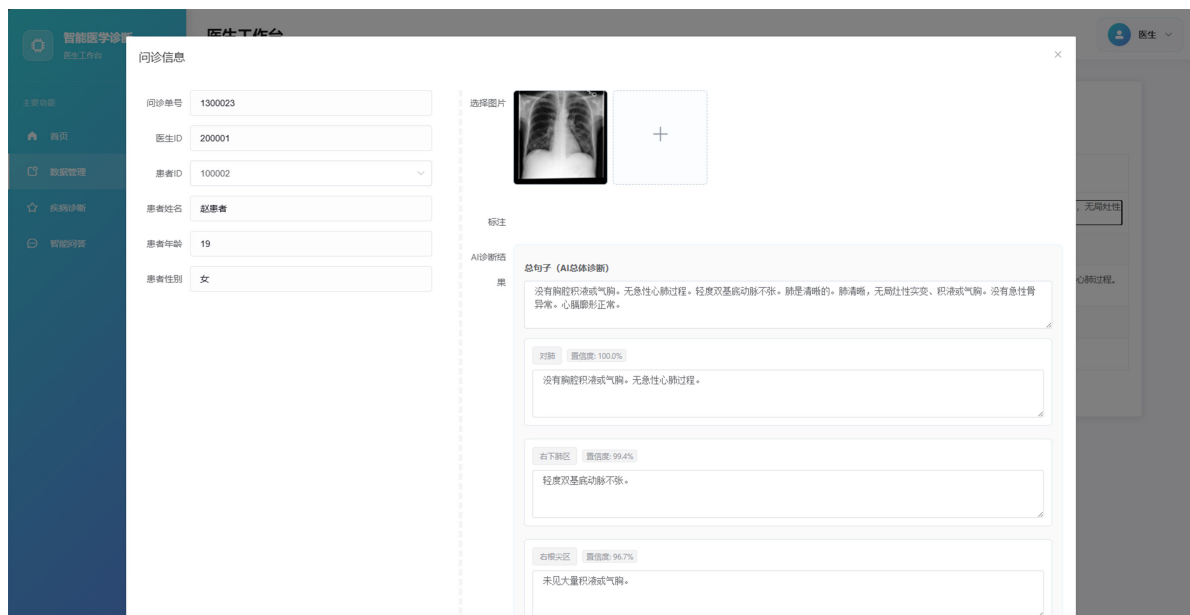


Figure 3. Consultation dialog: patient metadata and segmented AI diagnosis (physician workstation)

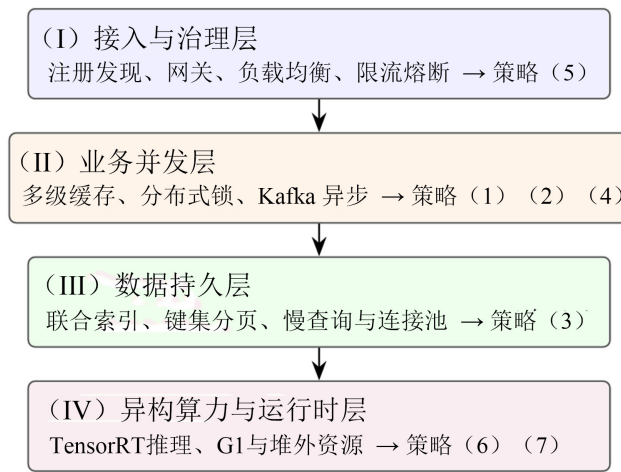
图 3. 问诊信息弹窗: 患者与问诊元数据及 AI 分条诊断

#### 4. 性能优化策略

为应对医疗场景下的高并发与大数据量挑战, 本文将落地措施按分层优化思路组织为四层: (I) 接入与治理层——服务注册发现、网关路由、客户端负载均衡与限流熔断; (II) 业务并发层——多级缓存、分布式锁与幂等、异步消息; (III) 数据持久层——联合索引、键集分页与连接治理; (IV) 异构算力与运行时层——GPU 推理引擎与 JVM 调优。图 4 在层级栈中一并标注七项策略编号(1)~(7)与归属。每层策略均可在测试或生产环境用指标验证(如缓存命中率、P99 延迟、Kafka 消费滞后、实例健康比例、GC 暂停时间等)。

**业务并发层：(1) 多级缓存**

针对管理员通知、医生基础信息、字典与配置等读多写少数据，采用“旁路缓存(cache aside)”；大规模互联网场景下分布式缓存体系的经验见文献[21]。读请求先访问 Redis，未命中再查 MySQL 并回填；写请求以数据库为准，成功后同步或异步失效相关缓存键，避免脏读。缓存键包含业务主键与必要版本片段(如 announce:{id}、user:profile:{doctorId})，TTL 按业务变化频率分层设置(公告类 300~900 s、基础资料类 10~30 min)。对热点键配合进程中短时互斥或 Redis 分布式信号量，减轻缓存击穿后的突发数据库压力；对穿透风险可在热点查询上使用空值短期缓存或布隆过滤器。线上结合 Redis INFO 与业务埋点统计命中率、序列化大小与键数量，作为容量与过期策略调整依据。



**Figure 4.** Layered hierarchy of performance optimization strategies  
**图 4.** 分层性能优化策略的层级关系示意图

**业务并发层：(2) 分布式锁与幂等性**

医学影像推理存在重复提交与重试场景，系统在任务粒度使用 Redis 分布式锁，保证同一影像(或同一任务 ID)在并发条件下仅有一条执行路径进入推理关键区，避免重复推理与 GPU 资源浪费。获取锁采用 SET key value NX PX ttl, value 为随机令牌；释放锁使用 Lua 脚本校验令牌后删除，防止误删其他实例持有的锁。若业务执行时间可能超过锁 TTL，则在临界区内由看门狗线程按令牌续期(见算法 1)。幂等层面以任务 ID 或业务幂等键建立唯一约束或去重表，与锁协同保证重复调用不重复扣减配额、不重复写入结果。

**Algorithm 1:** Distributed lock renewal with token validation (watchdog thread)

**算法 1:** 基于令牌校验的分布式锁续期(看门狗线程)

**输入:** Redis 键 key；加锁时写入的随机令牌 token；锁有效期 TTL；续期提前量  $\delta$  (当剩余寿命  $< \delta$  时触发续期)；轮询步长  $\Delta t$

```

while 业务任务未完成 and 当前实例仍持有锁 key do
  if Redis 返回的键剩余生存时间  $< \delta$  (或不可读且视为将过期)
    调用 EVAL 执行 Lua 脚本: 若 GET (key) = token, 则 PEXPIRE key TTL; 否则中止续期并报警
  end if
  sleep ( $\Delta t$ ) {例如 Thread.sleep 或事件库中的固定延迟, 避免忙等}
end while
  
```

**数据持久层：(3) 数据库查询优化**

在千万级病例表上，多条件检索遵循最左前缀原则设计联合索引，例如(patient\_id, exam\_time, id)以支持按患者与时间范围过滤并配合键集分页(Seek Method)替代大 OFFSET 分页；查询优化的一般原则亦见文献[22]。对可覆盖索引的查询利用索引条件下推(ICP)减少回表；统计与报表类需求优先采用预聚合表或异步汇总，避免在线全表扫描。慢查询通过 MySQL 慢日志与 EXPLAIN ANALYZE 持续收敛，关注执行计划中 type、rows、filtered 与 Extra，优先消除高风险排序与回表组合。数据源侧限制连接池最大连接、设置合理 wait\_timeout，并在应用侧配置连接借用超时与熔断，降低数据库雪崩风险。

#### 业务并发层：(4) 异步消息队列

批量删除、全量导出、磁盘清理等长耗时与大批次 I/O 操作通过 Kafka [23]异步化：接口将任务摘要写入 Topic 后快速返回，消费者按分区并行处理，借助消费者组实现水平扩展。消息体携带任务 ID、幂等键、重试次数与截止时间；消费者使用有界线程池与批量提交，失败进入重试 Topic 或死信队列以便补偿。运维上监控分区消费滞后(lag)与堆积，作为扩容消费者或调整批大小、拉取间隔的信号。

#### 接入与治理层：(5) 分布式高可用与服务治理

注册与发现侧部署多节点冗余集群，承担服务注册、配置管理与元数据维护；各微服务实例启动时注册并周期性上报心跳，异常实例被摘除，避免流量命中故障节点[6]。配置项按环境与环境隔离命名空间管理，支持动态刷新与灰度。调用链路上，网关与客户端负载均衡将请求分发至同一服务的多实例；结合超时、有限次重试与指数退避，降低瞬时网络抖动导致的失败率。在网关与服务侧实施 QPS/线程数限流、慢调用比例熔断与系统自适应保护，必要时配合热点参数限流；熔断打开时采用降级响应或缓存兜底，保证核心读路径可用。

#### 异构算力与运行时层：(6) 推理部署优化

GPU 推理路径采用 TensorRT 等引擎完成从 PyTorch/ONNX 到引擎的转换、FP16/INT8 等精度校准与序列化；在线阶段加载引擎并预分配工作区，减少动态分配与内核启动开销[18]。批大小与并发度由显存容量与延迟目标联合约束；CPU 侧预处理与 GPU 推理采用流水线重叠以提升吞吐。同一模型多版本通过路由规则与灰度发布并存，便于回滚与对比。

#### 异构算力与运行时层：(7) JVM 内存调优

Java 业务进程统一采用 G1 垃圾回收器，典型启动参数包括 -Xms4g -Xmx4g -XX:+UseG1GC -XX:MaxGCPauseMillis = 200 -XX:InitiatingHeapOccupancyPercent = 45 等。通过 GC 日志、jstat 与 APM 观察 Young/Mixed GC 次数、停顿时间与堆占用趋势，若出现频繁 Mixed GC 或晋升失败，再调整 MaxGCPauseMillis、ConcGCThreads 或堆容量。除堆内存外，线程池、HTTP 客户端连接池与 Netty 直接内存等堆外资源一并纳入监控与限流，避免“堆充足而进程 OOM”类故障。

## 5. 实验与测试

为了验证系统的性能与稳定性，在 Docker 环境下使用 JMeter 进行多场景压测。表 4~8 汇总主要量化指标；图 6 给出业务链路压测中单节点 CPU 利用率随并发用户数变化的趋势。RAG 对话子系统另设专项压测，量化检索与生成两阶段耗时及端到端接口延迟，见第 5.2 小节。

### 5.1. 实验环境与方法

压力测试在 Docker Compose 编排的三节点应用集群上进行(每节点：Intel Xeon Silver 4210, 8 vCPU, 32 GB RAM；模型服务节点配备 NVIDIA RTX 4080, 16 GB 显存)。业务库采用 MySQL 8.0, 单表规模约  $1.2 \times 10^7$  行；JMeter 线程组 Ramp-up 120 s, 持续压测 300 s。实验覆盖：MySQL 深度分页与索引优化前后对比；Redis 缓存对读接口 QPS 与延迟的影响；十万级批量删除在同步阻塞与 Kafka 异步下的接口耗

时；GPU 推理在 PyTorch 与 TensorRT 下的延迟及高并发下网关吞吐；以及业务链路在不同并发用户数下的 CPU 占用与 JVM 堆使用摘录。

## 5.2. RAG 对话服务性能测试

**测试目标：**面向第 3 节所述 RAG 对话服务，这里考察性能指标：将单次请求拆分为检索阶段(查询编码、向量召回与重排)与生成阶段(大语言模型解码)，分别统计分阶段耗时与吞吐；另统计 HTTP 端到端首包至尾包时间。

**压测用例与负载：**对 RAG 对话 HTTP 接口构造脚本化压测流：每条请求提交长度在 32~256 字之间、由参数化模板生成的中文查询串，固定触发 Top-K 检索( $K = 5$ )与最大 256 新 token 的流式生成；共采集 500 次成功响应，客户端串行发送、思考时间 0，压测前丢弃 10 次预热请求不计入统计。负载与机构内向量知识库对接，与业务库压测、推理微服务压测分时执行，避免 GPU 争用。

**环境与埋点：**检索侧为 384 维句向量 + FAISS 内积索引 + 轻量重排；生成侧为同实验环境模型节点上单卡 7B 级指令模型(4 bit 量化)，批大小 1。进程内记录检索完成、首 token 与全文结束时间戳；客户端以 HTTP 首字节/尾字节时间互证，与进程内拆分偏差  $\leq 5\%$ 。

**结果与分析：**表 1~3 汇总检索、生成与端到端指标。检索阶段处于百毫秒量级，生成阶段由自回归解码主导、为秒级；端到端 P99 与均值之比约 1.7，长尾主要来自生成长度与 KV 缓存状态差异，与在线推理服务的常见观测一致。

**Table 1.** RAG retrieval stage latency ( $n = 500$ , Top- $K = 5$ )

**表 1.** RAG 检索阶段耗时( $n = 500$ , Top- $K = 5$ )

统计量	查询编码	向量召回 + 重排	检索合计
平均	38 ms	54 ms	92 ms
P99	61 ms	118 ms	176 ms

**Table 2.** RAG generation stage throughput ( $n = 500$ , max. 256 new tokens)

**表 2.** RAG 生成阶段速度( $n = 500$ , 最大新 token 数 256)

统计量	首 token 延迟(TTFT)	解码吞吐	全序列生成时间
平均	241 ms	35.8 tokens/s	3.04 s
P99	428 ms	–	5.12 s

**Table 3.** End-to-end latency of RAG dialogue HTTP API ( $n = 500$ , serial, think time 0)

**表 3.** RAG 对话 HTTP 接口端到端延迟( $n = 500$ , 串行, 思考时间 0)

指标	平均 RT	P99 RT
端到端	3.35 s	5.74 s

## 5.3. 结果与分析

测试 SQL 为按(patient\_id, exam\_time)多条件筛选并按主键排序取页数据。表 4 表明，在应用联合索引与键集分页后，相对大偏移 OFFSET 方案，执行过程扫描行数由约千万级降至千级量级，查询延迟由秒级降至毫秒级。

**Table 4.** Comparison of deep-pagination query latency (page size  $k = 20$ )**表 4.** 深度分页查询延迟对比(单次页大小  $k = 20$ )

方案	平均 RT	P99 RT	扫描行数
大偏移 OFFSET + 无合适复合索引	2.41 s	3.86 s	~1200 万
联合索引 + 键集分页	16 ms	43 ms	~1800

读多写少接口(医生基础信息、公告列表)在仅走 MySQL 与增加 Redis 二级缓存两种配置下对比如表 5 所示。缓存键 TTL 设为 300 s, 布隆过滤器预热后, 压测期间 Redis 命中率稳定在 87%~94% 区间; 命中率高时数据库连接池活跃连接数由 35~42 降至 7~11, 读取 QPS 与延迟明显改善。

**Table 5.** Read API performance with and without Redis cache (200 concurrent users)**表 5.** 引入 Redis 缓存前后的读接口表现(并发 200 用户)

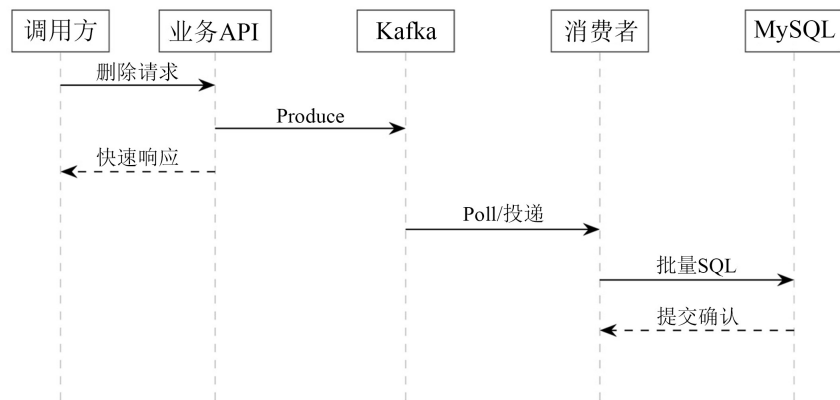
配置	QPS	平均 RT	P99 RT	缓存命中率
无缓存, 直连 MySQL	428	142 ms	318 ms	-
Redis 缓存 + 旁路读	3265	31 ms	84 ms	91.6%

对同一批约  $10^5$  条逻辑删除任务, 同步接口在事务提交前阻塞等待批量更新完成; 异步方案仅写入 Kafka 并快速返回, 后台消费者线程池(核心 16, 最大 32, 队列长度 2000)拉取分批提交。表 6 对比两种模式下接口侧耗时; 异步路径下请求经 Kafka 投递后由消费者线程池分批写回 MySQL, 接口快速返回。

**Table 6.** Interface-side latency for batch deletion tasks (order  $\sim 10^5$  rows)**表 6.** 十万级批量删除任务的接口侧耗时对比

模式	接口平均 RT	业务完成时间
同步阻塞	58.7 s	与接口一致
Kafka 异步	0.28 s	后台约 156 s 平滑完成

图 5 给出 Kafka 异步批量删除链路的时序关系: 调用方将删除请求提交至业务 API 后, API 将任务消息写入 Kafka 并立即返回; 消费者从分区拉取消息后, 在数据库侧分批执行更新或逻辑删除, HTTP 接口线程不等待数据库批处理结束。

**Figure 5.** Sequence diagram of asynchronous batch-deletion path**图 5.** 批量删除异步链路时序图

输入分辨率固定为  $512 \times 512$ ，批大小为 1；PyTorch 为 FP32 动态图路径，TensorRT 为 FP16 引擎(校准样本 500 张)。表 7 给出单卡延迟；高并发下阻塞式线程池约在 118~135 req/s 出现排队拐点，改为异步批处理与有界队列后，稳定段可维持在约 312~348 req/s (见表 8)。

**Table 7.** Single-request inference latency and throughput (single GPU)

**表 7.** 单请求推理延迟与吞吐(单 GPU)

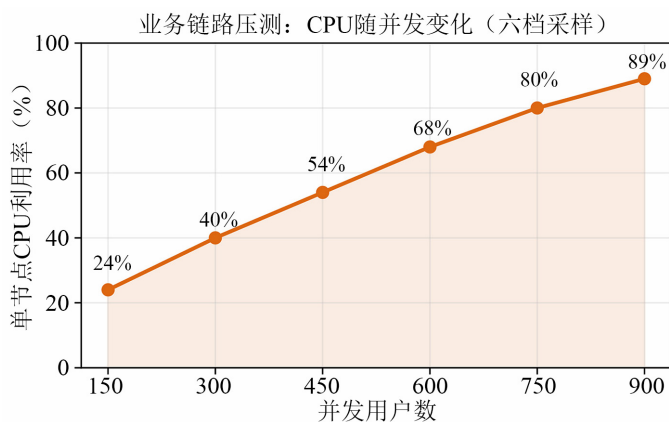
后端	平均延迟	P99 延迟	推算 QPS 上限(单流)
PyTorch (CUDA)	39 ms	67 ms	~25.6
TensorRT FP16	14 ms	26 ms	~71.4

**Table 8.** Model gateway forwarding throughput under high load (JMeter, think time 0)

**表 8.** 高并发下模型网关转发吞吐(JMeter, 思考时间 0)

模式	饱和前最大稳定 QPS	P99 RT	错误率
同步阻塞线程池	126	1.12 s	0.12%
异步非阻塞 + 限流	331	0.89 s	0.04%

网关至业务服务链路在并发用户数 150、300、450、600、750、900 六档采样；图 6 在横轴六档并发处给出折线及采样点标注，单节点 CPU 利用率由低至高依次为 24%、40%、54%、68%、80% 与 89%，随加压档位升高而近似单调上升，末档已接近单节点高负载区。该压测路径为 Java 业务与网关转发而非 GPU 内核，在约 900 并发、利用率接近九成时，CPU 时间主要消耗于用户态的 HTTP 解析、JSON 序列化与反序列化、Spring/Web 容器内的业务逻辑与线程池调度；随连接与请求速率升高，内核态在网络协议栈、套接字读写及线程上下文切换上的占比亦可能上升。本实验未采集 perf 火焰图，上述解释为与压测配置一致的机理分析，不排除存在磁盘或锁竞争等次要因素。JVM 堆占用稳定在 4.0~4.4 GB (容器上限 8 GB)。生产环境 Java 进程采用 G1 回收器，典型启动参数如下(节选)：-Xms4g -Xmx4g -XX:+UseG1GC -XX:MaxGCPauseMillis = 200 -XX:InitiatingHeapOccupancyPercent = 45。



**Figure 6.** Load test (six levels): single-node CPU utilization vs. concurrent users

**图 6.** 压测采样(六档)：单节点 CPU 利用率随并发用户数变化

## 6. 结束语

本文围绕复杂医疗 AI 场景下高性能分布式医学数据管理与智能分析系统的综合设计实践与性能优

化实践展开,旨在为同类工程提供可复现、可借鉴的系统蓝图:第2节对医疗信息化与深度学习协同部署等领域的研究脉络作了简要对照;在功能上覆盖医生、患者与管理侧能力,智能侧以胸片X光疾病诊断、胸片X光报告生成、超声心动图分割与射血分数计算、尿检肾功能诊断、智能中医舌诊五类推理微服务及基于LangChain的RAG知识问答支撑临床辅助流程;在基础设施上采用Nacos、Sentinel、网关与OpenFeign等组件支撑微服务治理,并结合Redis、Kafka与MySQL完成数据与消息链路;第3.4小节从工程设计角度归纳了接入收敛、鉴权与审计、合规映射及存储侧安全落地,并提及与基础设施解耦的卷级加密、KMS等可选加固。第4节分层性能优化路线下的七项策略与压测数据相互印证;第5.2小节给出RAG检索与生成阶段耗时及端到端接口延迟。交付上采用Docker与Docker Compose支撑研发测试,生产环境采用Kubernetes编排。后续工作将侧重于监控告警、容量规划、发布体系及与更高等级安全要求相衔接的存储层加固等方面的持续完善。

## 参考文献

- [1] 孟群, 胡建平, 陈伟. 我国医院信息化建设现状与发展策略[J]. 中国卫生信息管理杂志, 2014, 11(4): 265-271.
- [2] 汪鹏, 阮聿锐, 鲁杰, 等. 国内外智慧医疗发展现状及其趋势[J]. 中国医院管理, 2014, 34(1): 29-32.
- [3] Karabey Aksakalli, I., Çelik, T., Can, A.B. and Tekinerdoğan, B. (2021) Deployment and Communication Patterns in Microservice Architectures: A Systematic Literature Review. *Journal of Systems and Software*, **180**, Article ID: 111014. <https://doi.org/10.1016/j.jss.2021.111014>
- [4] Pahl, C., Brogi, A., Soldani, J. and Jamshidi, P. (2019) Cloud Container Technologies: A State-of-the-Art Review. *IEEE Transactions on Cloud Computing*, **7**, 677-692. <https://doi.org/10.1109/tcc.2017.2702586>
- [5] Esteva, A., Kuprel, B., Novoa, R.A., Ko, J., Swetter, S.M., Blau, H.M., et al. (2017) Dermatologist-Level Classification of Skin Cancer with Deep Neural Networks. *Nature*, **542**, 115-118. <https://doi.org/10.1038/nature21056>
- [6] Taibi, D., Lenarduzzi, V., Pahl, C., et al. (2018) Microservices: A Systematic Mapping Study. *Journal of Grid Computing*, **16**, 323-353.
- [7] Kiela, D., Bansal, M., Boratko, M., et al. (2019) Supervised Multimodal Bitransformers for Classifying Images and Text.
- [8] Devlin, J., Chang, M.-W., Lee, K., et al. (2019) BERT: Pre-Training of Deep Bidirectional Transformers for Language Understanding. *Proceedings of NAACL-HLT*, Minneapolis, 2-7 June 2019, 4171-4186.
- [9] Ren, S., He, K., Girshick, R. and Sun, J. (2017) Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, **39**, 1137-1149. <https://doi.org/10.1109/tpami.2016.2577031>
- [10] Radford, A., Wu, J., Child, R., et al. (2019) Language Models Are Unsupervised Multitask Learners. OpenAI Technical Report. [https://cdn.openai.com/better-language-models/language\\_models\\_are\\_unsupervised\\_multitask\\_learners.pdf](https://cdn.openai.com/better-language-models/language_models_are_unsupervised_multitask_learners.pdf)
- [11] Hara, K., Kataoka, H. and Satoh, Y. (2018) Can Spatiotemporal 3D CNNs Retrace the History of 2D CNNs and ImageNet? 2018 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, Salt Lake City, 18-22 June 2018, 6546-6555. <https://doi.org/10.1109/cvpr.2018.00685>
- [12] Kipf, T.N. and Welling, M. (2017) Semi-Supervised Classification with Graph Convolutional Networks. *Proceedings of ICLR*, Toulon, 24-26 April 2017. <https://arxiv.org/abs/1609.02907>
- [13] Hornik, K., Stinchcombe, M. and White, H. (1989) Multilayer Feedforward Networks Are Universal Approximators. *Neural Networks*, **2**, 359-366. [https://doi.org/10.1016/0893-6080\(89\)90020-8](https://doi.org/10.1016/0893-6080(89)90020-8)
- [14] He, K., Zhang, X., Ren, S. and Sun, J. (2016) Deep Residual Learning for Image Recognition. 2016 *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Las Vegas, 27-30 June 2016, 770-778. <https://doi.org/10.1109/cvpr.2016.90>
- [15] Hochreiter, S. and Schmidhuber, J. (1997) Long Short-Term Memory. *Neural Computation*, **9**, 1735-1780. <https://doi.org/10.1162/neco.1997.9.8.1735>
- [16] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., et al. (2017) Attention Is All You Need. *Proceedings of the 31st International Conference on Neural Information Processing Systems*, Long Beach, 4-9 December 2017, 6000-6010.
- [17] Lewis, P., Perez, E., Piktus, A., et al. (2020) Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems*

- 2020, *NeurIPS* 2020, 6-12 December 2020, 9459-9474.
- [18] LeCun, Y., Bengio, Y. and Hinton, G. (2015) Deep Learning. *Nature*, **521**, 436-444. <https://doi.org/10.1038/nature14539>
- [19] Merkel, D. (2014) Docker: Lightweight Linux Containers for Consistent Development and Deployment. *Linux Journal*, **2014**, Article No. 2.
- [20] 程啸. 论我国个人信息保护法中的个人信息处理规则[J]. *清华法学*, 2021, 15(3): 55-73.
- [21] Nishtala, R., Fugal, H., Grimm, S., *et al.* (2013) Scaling Memcache at Facebook. In: *Proceedings of the 10th USENIX Conference on Networked Systems Design and Implementation (NSDI)*, USENIX Association, 385-398.
- [22] Ioannidis, Y.E. (1996) Query Optimization. *ACM Computing Surveys*, **28**, 121-123. <https://doi.org/10.1145/234313.234367>
- [23] Kreps, J., Narkhede, N. and Rao, J. (2011) Kafka: A Distributed Messaging System for Log Processing. *Proceedings of the NetDB Workshop*, Athens, 12 June 2011. <https://pages.cs.wisc.edu/~akella/CS744/F17/838-CloudPapers/Kafka.pdf>