

Legalreach: A Multi-Tier AI-Augmented Platform for Democratizing Legal Consultation Access in India

Kushal Khune¹, Aditya Akankar², Shruti Mundhe³, Sakshi Kasurkar⁴, Anjali Shegokar⁵, Prof. Amit Patil⁶
^{1,2,3,4,5}*Student, Dept. of Information Technology Engineering, Mauli Group of Institution's College of Engineering and Technology, Shegaon, India*

⁶*Guide, Mauli Group of Institution's College of Engineering and Technology, Shegaon, India*

Abstract—Access to legal counsel in India is severely constrained by geographic barriers, opaque pricing, and the absence of a unified digital intermediary between citizens and legal professionals. Existing platforms address access and discovery but provide no integrated consultation workflow, document security, or AI-assisted guidance. This paper presents Legal Reach, a multi-tier web platform that unifies domain-constrained Retrieval-Augmented Generation (RAG) legal guidance, freemium-gated real-time chat, WebRTC peer-to-peer video consultation, and presigned cloud document storage within a single role-differentiated system. The RAG pipeline, built on FAISS vector search and a domain-limited system prompt, achieves 85% top-3 retrieval accuracy across four Indian statutes. All REST API operations except the RAG inference step satisfy a 300 ms P95 latency target. Security controls correctly reject 100% of tampered payment signatures and expire presigned URLs at the 15-minute boundary.

Index Terms—FAISS vector search, Legal Tech, RAG, real-time consultation, Retrieval-Augmented Generation, Socket.io, WebRTC, legal access, AI-assisted legal guidance

I. INTRODUCTION

Access to legal counsel is widely recognized as a prerequisite for the equitable enforcement of civil rights. Yet in India, structural barriers prevent the majority of citizens from obtaining timely, affordable, and competent legal advice. A 2022 NLSIU survey [1] found that over 75% of respondents in Tier-2 and Tier-3 cities had never consulted a lawyer, citing cost opacity, geographic inaccessibility, and unfamiliarity with the legal system as primary deterrents. The ratio of lawyers to population in India stands at approximately 1:1,000, among the lowest of any major

jurisdiction, and is concentrated disproportionately in metropolitan centres.

The problem is compounded on the supply side. Independent practitioners in Tier-2 and Tier-3 cities lack purpose-built digital infrastructure for practice management. Appointment scheduling, client communication, document exchange, and payment collection occur across fragmented, insecure channels—email threads, WhatsApp groups, and physical couriers—introducing inefficiency, privacy risk, and revenue leakage.

Existing platforms address this problem only partially. General-purpose directories (Justdial, Sulekha) provide contact information but offer no consultation workflow, document security, real-time communication, or AI assistance. Specialised services either restrict access by language, jurisdiction, or premium subscription, or deploy unrestricted general-purpose large language models (LLMs) that hallucinate statutory references—a critical failure mode in a domain where incorrect legal information can cause direct harm to users.

This paper makes the following contributions: (1) a multi-tier web architecture that integrates AI-driven legal guidance, real-time chat, and peer-to-peer video consultation within a single role-differentiated platform; (2) a domain-constrained RAG pipeline that grounds every AI response in a retrievable statutory corpus, preventing parametric hallucination; (3) a freemium payment gate implemented via HMAC-SHA256 webhook verification that enables sustainable platform monetisation while preserving basic access; (4) a presigned URL document storage model that provides time-limited, audit-logged access to sensitive client documents without exposing cloud credentials.

The remainder of this paper is organized as follows. Section II reviews related work. Section III describes the proposed system architecture. Section IV details the implementation of key novel components. Section V presents evaluation results. Section VI concludes with directions for future work.

II. RELATED WORK

A. Legal Tech Platforms and Their Limitations

Susskind [2] argued that legal services are undergoing AI-driven disruption, predicting a shift from bespoke human counsel to commoditized, online-delivered legal products. This prediction is validated by GPT-4's reported bar exam passage [3], which demonstrated that transformer-based models now encode sufficient legal knowledge for credentialled performance. However, the clinical deployment of such models for citizen-facing advice remains problematic: unconstrained LLMs produce fluent but ungrounded statutory citations, and no existing commercial platform has published a methodology for constraining LLM outputs to a verified legal corpus in the Indian jurisdiction.

B. Retrieval-Augmented Generation for Legal Domains

Lewis et al. [4] introduced Retrieval-Augmented Generation as a methodology that augments LLM generation with dense retrieval over external corpora, substantially improving factual accuracy for knowledge-intensive tasks. Johnson et al. [5] demonstrated that FAISS IndexFlatL2 structures support billion-scale approximate nearest-neighbour search with sub-millisecond query latency on CPU, making them suitable for deployment in resource-constrained environments. No prior work has applied RAG with domain-restriction prompting specifically to Indian statutory text for citizen-facing legal guidance.

C. Real-Time Communication in Professional Service Platforms

Pimentel et al. [6] formally specified the WebSocket protocol, enabling the persistent full-duplex channels that underlie Socket. Io's room-based event architecture. Holmberg [7] evaluated WebRTC for clinical telemedicine, demonstrating acceptable video latency over consumer broadband but noting the

critical dependency on ICE/STUN infrastructure for NAT traversal. Neither source addresses the integration of payment gating with real-time consultation initiation, which is the novel architectural contribution of Legal Reach's chat module.

D. Security in Cloud-Based Professional Platforms

Zisis and Lekkas [8] identified presigned URL-based access control as the recommended approach for cloud document storage in regulated professional contexts, avoiding the exposure of direct bucket credentials to client applications. NIST SP 800-63B [14] mandates HMAC-based webhook authentication for financial transaction verification. LegalReach implements both recommendations and adds a 15-minute expiry window on presigned URLs to limit the blast radius of credential interception.

The gap in existing literature is clear: no published work describes a unified, role-differentiated platform combining domain-constrained RAG, real-time consultation, WebRTC video, and freemium payment gating within a single deployable system targeting the Indian legal services market.

III. PROPOSED SYSTEM ARCHITECTURE

A. Architectural Overview

LegalReach is designed as a loosely coupled, multi-tier system comprising four distinct bounded contexts: (1) a Presentation Layer, (2) an Application and API Layer, (3) an AI Microservice Layer, and (4) a Data and Storage Layer. Each layer communicates through well-defined interfaces—REST over HTTPS, WebSocket over WSS, and internal HTTP between the Node.js gateway and the FastAPI microservice—ensuring that each bounded context can be independently scaled, replaced, or audited.

The Presentation Layer is a React.js 18 single-page application using React Router DOM v6 for client-side routing, Tailwind CSS 3.4 for responsive layout, and the React Context API for global authentication state. Role-specific routing guards redirect unauthenticated users and enforce the Client/Lawyer role boundary at the UI level, providing defence-in-depth alongside the server-side JWT middleware.

B. Communication Protocol Design

TLS 1.3 is enforced at the Nginx boundary with HSTS headers preventing protocol downgrade. REST API

calls and WebSocket connections (WSS) share the same origin, eliminating cross-origin preflight overhead. The Socket.io server is mounted as a sub-path on the same Express application, enabling a single Nginx upstream block to handle both HTTP and WebSocket upgrade requests with minimal configuration surface.

C. Role-Based Access Control Model

The system defines two primary roles: Client and Lawyer. Role assignment occurs at registration and is encoded in the JWT payload. An AUTH Middleware function extracts the Bearer token from the Authorization header, verifies the signature against the server-side secret, and attaches the decoded payload to the Express request object. Route-level role guards then compare req.user.role against the required role, returning HTTP 403 for violations. This model follows the principle of least privilege: Lawyer-only routes such as appointment management and profile publication are inaccessible to Client-role tokens even if the client-side routing guard is bypassed.

IV. IMPLEMENTATION OF KEY COMPONENTS

A. Domain-Constrained RAG Pipeline

The AI Legal Assistant is the most technically novel component of the system and the primary differentiator from general-purpose LLM deployments. The ingestion phase uses PyPDF2 to extract raw text from four Indian statutes: the Consumer Protection Act 2019, the Information Technology Act 2000, the Indian Contract Act 1872, and the Right to Information Act 2005. Text is segmented into 512-token chunks with a 50-token overlap using a sliding window algorithm implemented in FastAPI's startup event handler.

Each chunk is encoded into a 384-dimensional dense vector using the all-MiniLM-L6-v2 Sentence Transformer model [9] and indexed in a FAISS IndexFlatL2 structure. The index and associated chunk metadata are persisted to disk and loaded into memory on service startup, eliminating per-request re-indexing overhead. At query time, the user's question is encoded with the same model, and the top-3 nearest neighbours are retrieved by Euclidean distance. The retrieved chunks are concatenated into a context block and injected into a structured prompt delivered to a locally-hosted Ollama Llama 3 instance.

This directive enforces the critical constraint that distinguishes LegalReach's assistant from a generic LLM: parametric knowledge is suppressed, and every answer is grounded in retrievable, citable statutory text. If no relevant chunk is retrieved above a similarity threshold, the system returns a standardised disclaimer rather than generating an unsupported answer.

B. Freemium-Gated Real-Time Chat System

The chat system is a deliberate hybrid of REST and WebSocket protocols. Session initiation and state mutations (create, accept, close) are handled via REST to maintain durable, auditable state in MongoDB. Real-time message delivery uses Socket.io rooms keyed on the session ObjectId, allowing both participants to receive messages with sub-100 ms round-trip latency on local network benchmarks.

After the client completes Razorpay payment, the backend verifies the transaction's authenticity by computing HMAC-SHA256 over the concatenation of razorpayOrderId and razorpayPaymentId using the merchant secret key. The computed digest is compared to the razorpaySignature field in the webhook payload using a constant-time comparison to prevent timing attacks. Only sessions with verified payment status are eligible for Socket.io room admission, enforced by the server-side session state check on the connection event handler.

C. WebRTC Peer-to-Peer Video Consultation

Video consultation uses a signalling-server-assisted WebRTC architecture. The Socket.io server acts purely as a signalling relay: it forwards SDP offers, SDP answers, and ICE candidates between the two participants but carries no media. Once the WebRTC peer connection is established, all audio and video flows directly between the two browsers over DTLS-SRTP, with the signalling server completely bypassed. This architecture scales media bandwidth linearly with the number of active consultations rather than requiring server-side media processing.

V. EVALUATION AND RESULTS

A. Evaluation Methodology

Evaluation was conducted across three dimensions: (1) RAG retrieval quality, assessed by querying the pipeline with 40 domain-specific questions spanning

the four statutes in the corpus and evaluating whether the correct statutory basis appeared in the top-3 retrieved chunks; (2) system latency, measured by instrumenting each API endpoint with server-side timing middleware and recording median and 95th-percentile response times over 200 synthetic requests per endpoint; (3) security control correctness, verified by injecting tampered payment payloads and expired presigned URLs to confirm rejection behaviour.

B. RAG Retrieval Accuracy

The top-3 FAISS retrieval returned chunks containing the correct statutory basis for 34 of 40 test queries, yielding a retrieval accuracy of 85%. The 6 misses occurred on queries that required synthesis across multiple sections of different statutes—a known limitation of single-index flat retrieval. Analysis of the miss cases indicates that per-statute FAISS indices with cross-index re-ranking would address 4 of the 6 failures; the remaining 2 required temporal reasoning about amendment history not present in the current corpus.

C. Latency Benchmarks

Table I summarises latency measurements across all key system operations. All API operations satisfy the 300 ms P95 target with significant margin, except the RAG pipeline which—at 2,800 ms P95—reflects the sequential nature of CPU-bound Ollama inference. The Socket.io message round-trip of 18 ms median meets real-time communication requirements. WebRTC connection setup at 3,500 ms median is consistent with published benchmarks for STUN-assisted ICE negotiation over consumer broadband.

TABLE I System Latency Benchmarks (Median / P95 in milliseconds)

Operation	Median (ms)	P95 (ms)
User Auth (Login + JWT)	45	110
Lawyer Search Query	62	190
Socket.io Msg Round-Trip	18	55
AI RAG Full Pipeline	1,200	2,800
PDF Upload to S3 (1 MB)	820	N/A
Razorpay Order Creation	310	N/A
React App Initial Load	1,800	2,900
WebRTC Connection Setup	3,500	5,000

D. Security Verification

All HMAC-SHA256 signature verification tests correctly rejected 100% of payloads with tampered payment identifiers. Presigned S3 URLs were confirmed to return HTTP 403 from AWS at exactly the 15-minute expiry boundary. These results confirm that the payment and document security controls operate correctly under adversarial conditions consistent with the OWASP Top 10 threat model [10].

VI. CONCLUSION AND FUTURE WORK

This paper presented LegalReach, a multi-tier AI-augmented web platform designed to address the documented access deficit in the Indian legal services market. The platform's primary technical contributions are a domain-constrained RAG pipeline that grounds AI legal guidance in retrievable statutory text, a freemium-gated real-time consultation system with cryptographically verified payment, and a WebRTC peer-to-peer video architecture that scales media bandwidth without server-side media processing.

Evaluation demonstrates that all API operations except the RAG inference pipeline meet the 300 ms P95 latency target, that the security controls correctly reject tampered payment signatures and expired presigned URLs, and that the RAG pipeline achieves 85% top-3 retrieval accuracy on Indian statutory queries—a result that, while not production-grade, establishes a viable baseline for domain-specific legal AI in resource-constrained deployment environments. Future work includes: (i) GPU-accelerated Ollama inference to reduce RAG pipeline P95 below 800 ms; (ii) per-statute FAISS indices with cross-index re-ranking to address cross-statute query failures; (iii) expansion of the statutory corpus to include state-level legislation and subordinate regulations; (iv) a multilingual embedding model to support regional language queries; (v) a formal security audit against OWASP ASVS Level 2 criteria prior to public deployment.

ACKNOWLEDGMENT

The authors thank their project guide and the Department of Computer Engineering for guidance throughout this work. The authors acknowledge the open-source communities behind React.js, Socket.io,

FastAPI, FAISS, and Ollama, whose freely available tools made this research feasible.

REFERENCES

- [1] National Law School of India University, "Access to Justice Survey Report," NLSIU, Bengaluru, India, 2022.
- [2] R. Susskind, *Tomorrow's Lawyers: An Introduction to Your Future*, 2nd ed. Oxford University Press, 2017.
- [3] D. M. Katz, M. J. Bommarito, S. Gao, and P. Arredondo, "GPT-4 Passes the Bar Exam," *Philosophical Transactions of the Royal Society A*, vol. 382, no. 2270, 2024.
- [4] P. Lewis, E. Perez, A. Piktus et al., "Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks," *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 33, pp. 9459–9474, 2020.
- [5] J. Johnson, M. Douze, and H. Jégou, "Billion-Scale Similarity Search with GPUs," *IEEE Transactions on Big Data*, vol. 7, no. 3, pp. 535–547, 2019.
- [6] V. Pimentel, B. Nickerson, and M. Schipper, "The WebSocket Protocol," RFC 6455, IETF, 2012.
- [7] A. Holmberg, "WebRTC for Healthcare: Evaluating Peer-to-Peer Video Architectures in Clinical Settings," *Journal of Medical Internet Research*, vol. 20, no. 4, e10219, 2018.
- [8] D. Zissis and D. Lekkas, "Addressing Cloud Computing Security Issues," *Future Generation Computer Systems*, vol. 28, no. 3, pp. 583–592, 2012.
- [9] N. Reimers and I. Gurevych, "Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks," in *Proc. EMNLP-IJCNLP*, 2019, pp. 3982–3992.
- [10] OWASP Foundation, "OWASP Top 10: 2021," 2021. [Online]. Available: <https://owasp.org/Top10/>
- [11] Ollama Team, "Ollama: Run Large Language Models Locally," 2024. [Online]. Available: <https://ollama.com>
- [12] MongoDB Inc., "MongoDB Architecture Guide," MongoDB Documentation, 2024. [Online]. Available: <https://www.mongodb.com/docs/>
- [13] Razorpay, "Payment Gateway Integration and Webhook Signature Verification," Razorpay Developer Documentation, 2024. [Online]. Available: <https://razorpay.com/docs/>
- [14] NIST, "Special Publication 800-63B: Digital Identity Guidelines," U.S. Dept. of Commerce, 2020.
- [15] Amazon Web Services, "S3 Security Best Practices: Presigned URLs," AWS Documentation, 2024. [Online]. Available: <https://docs.aws.amazon.com/s3/>