



# IJRASET

International Journal For Research in  
Applied Science and Engineering Technology



---

# INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

---

**Volume:** 14    **Issue:** V    **Month of publication:** May 2026

**DOI:** <https://doi.org/10.22214/ijraset.2026.81760>

[www.ijraset.com](http://www.ijraset.com)

Call:  08813907089

E-mail ID: [ijraset@gmail.com](mailto:ijraset@gmail.com)

# CloudDrive-AI: Intelligent Cloud Storage with Federated DDoS Defense and RAG-Based Semantic Search

Puneeth Venkat<sup>1</sup>, Yashwanth S S<sup>2</sup>, Dhanushree Singh<sup>3</sup>, Dinesh R<sup>4</sup>, Prof. Jyothi Kiran Mirji<sup>5</sup>  
Reva University, Bangalore

**Abstract:** Traditional cloud storage platforms function as passive data repositories, lacking the capability to semantically understand stored document contents and forcing users into laborious manual retrieval workflows. Concurrently, sophisticated Layer-7 Distributed Denial-of-Service (DDoS) attacks targeting API endpoints threaten both service availability and operational economics, particularly for AI-enhanced platforms where every inference request incurs computational cost. This paper presents CloudDrive-AI, a comprehensive intelligent cloud storage framework integrating three synergistic technologies: (1) a federated multi-cluster machine learning architecture providing real-time DDoS detection via Isolation Forest anomaly scoring with 2-of-3 majority-vote aggregation, eliminating single points of failure inherent in centralised security systems; (2) a robust OCR pipeline built on Tesseract V5 Long Short-Term Memory (LSTM) networks for extracting text from scanned documents and images; and (3) a context-based question answering framework that sends extracted OCR text to the Google Gemini API to deliver grounded, hallucination-free responses over user-uploaded documents. The federated security layer demonstrates effective detection of upload-based DDoS attacks, as validated through a controlled attack simulator. The OCR pipeline attains practical accuracy on both clean scans and smartphone photographs. A controlled attack simulation confirms that the system successfully blocks malicious traffic while allowing legitimate user requests. CloudDrive-AI transforms passive cloud repositories into active, intelligent, and secure document management ecosystems.

**Index Terms:** Cloud Storage, Federated Learning, DDoS Detection, Isolation Forest, Optical Character Recognition, Generative AI, Semantic Search.

## I. INTRODUCTION

The global datasphere is projected to reach 175 zettabytes by 2025, with roughly 80% of data residing in unstructured formats—scanned documents, PDF reports, image files, and email archives [1]. Cloud storage providers such as AWS S3, Google Drive, and Microsoft OneDrive excel at delivering durable, geographically distributed storage infrastructure but operate almost exclusively as passive data silos [2]. Their search capabilities remain confined to file metadata—filenames, timestamps, and MIME types [3]—leaving users unable to locate specific information embedded within document bodies without downloading and manually scanning each file.

A legal professional hunting an indemnification clause in a 75-page contract PDF, or a financial analyst seeking a figure captured only in a scanned receipt, receives no help from conventional cloud search.

Simultaneously, modern Layer-7 DDoS attacks exploit the computational intensity of AI-enhanced API operations—OCR analysis and LLM inference—to cause Economic Denial of Sustainability (EDoS), rapidly exhausting operational budgets through sustained spurious requests [4]. Traditional Web Application Firewalls and IP-based rate limiting fail against distributed botnets that rotate through vast pools of residential proxy IP addresses, employ User-Agent spoofing, and carefully mimic legitimate user behaviour [5]. The integration of Transformer-based Large Language Models (LLMs) [6] introduces a further challenge: LLMs exhibit well-documented hallucination—the confident generation of factually unsupported content [8]—making naive document-corpus integration unreliable and potentially dangerous in enterprise settings where factual errors carry legal or financial consequences.

To address these converging challenges, we present

CloudDrive-AI, built on three integrated and synergistic pillars:

- 1) Federated Multi-Cluster DDoS Defense: Three independent Isolation Forest clusters operating in parallel with 2-of-3 majority-vote aggregation, providing Byzantine fault tolerance and eliminating centralised failure points.
- 2) Robust OCR Pipeline: Tesseract V5 LSTM recognition with format normalisation and basic preprocessing to handle real-world document quality.

- 3) Grounded Context-Based Question Answering: Extracted OCR text is sent directly to the Google Gemini API, which is instructed to answer exclusively from the provided content, enforcing strictly factual responses anchored in user-owned documents.

The remainder of this paper is organised as follows: Section II surveys related work and identifies research gaps; Section III details the five-plane system architecture; Section IV presents a complete attack mitigation scenario; Section V reports experimental results; Section VI discusses implications, limitations, and future directions; Section VII concludes.

## II. RELATED WORK AND RESEARCH GAPS

### A. Optical Character Recognition

Early OCR systems relied on template matching and handcrafted feature extraction, performing adequately on clean, consistently-fonted documents but degrading severely under imaging conditions encountered in practice [10]. The integration of LSTM recurrent networks into Tesseract V4/V5 represents a paradigm shift, enabling context-aware bidirectional sequence prediction that achieves character accuracy exceeding 95% on standard benchmarks [11]. The LSTM architecture captures long-range character dependencies and resolves ambiguous glyphs via linguistic context—distinguishing ‘l’, ‘1’, and ‘I’ without purely visual cues. Hegghammer’s empirical comparison of Tesseract, Amazon Textract, and Google Document AI demonstrated that open-source Tesseract achieves competitive accuracy at substantially lower cost and without data egress concerns [12], directly motivating the OCR pipeline adopted in CloudDrive-AI.

### B. DDoS Detection and Mitigation

Signature-based intrusion detection is effective against documented attack vectors but inherently reactive—unable to detect novel or polymorphic threats—and requires continuous expert-curated database maintenance [5]. IP-based rate limiting is trivially bypassed by botnets distributing load across thousands of nodes, each generating sub-threshold traffic while collectively devastating the target [4]. The EDoS-Shield framework combined client puzzle protocols with admission control but imposed latency burdens on legitimate users. The Isolation Forest algorithm [13] overcomes these limitations by exploiting the “few and different” property of anomalies: outliers are isolable via fewer recursive binary partitions in randomly constructed trees, yielding shorter average path lengths and higher anomaly scores. Its  $O(n)$  training complexity, sub-millisecond inference, and freedom from distributional assumptions make it uniquely suited to high-throughput API gateway deployment. Several recent works have applied Isolation Forest to network intrusion detection, but most assume a centralised model; the federated multi-cluster approach with majority voting remains underexplored in the context of cloud storage APIs.

### C. Context-Based Question Answering

Traditional keyword retrieval (TF-IDF, BM25) suffers from vocabulary mismatch: relevant documents may employ synonyms absent from the query [14]. Large Language Models like Google Gemini provide powerful generative capabilities but exhibit hallucination when not constrained to a provided context. The approach of supplying extracted document text directly as context to the LLM and instructing it to answer only from that content has been shown to effectively eliminate hallucination [9]. The Google Gemini API [7] provides native grounded-generation support, making it well-suited as CloudDrive-AI’s generative backend. Unlike vector database-based retrieval-augmented generation (RAG) systems that require additional infrastructure, CloudDrive-AI adopts a simpler and more transparent approach: the full OCR-extracted text from relevant files is directly passed to Gemini, eliminating the complexity of embedding pipelines and chunking strategies.

### D. Research Gaps Addressed

Despite advances in each constituent domain, the literature reveals four critical gaps that CloudDrive-AI directly addresses: (i) No prior unified architecture integrates real-time ML-driven DDoS defense with AI-powered document understanding in a single cloud storage framework. (ii) Production DDoS systems rely on centralised detectors susceptible to targeted single-point failures; federated majority-vote alternatives remain unexplored for this application. (iii) Contextbased LLM integration into consumer-facing personal cloud storage with explicit privacy safeguards and usability evaluation is absent. (iv) Comprehensive end-to-end performance characterisation spanning security filtering, OCR extraction, and grounded response generation has not been reported for an integrated pipeline. CloudDrive-AI fills these gaps by providing a fully implemented, evaluated system that combines all three capabilities in a coherent architecture.

### III. SYSTEM ARCHITECTURE AND IMPLEMENTATION

CloudDrive-AI adopts a five-plane microservices architecture cleanly separating concerns across: Presentation, Intelligence (Security), Service, Processing, and Reasoning planes. Figure 1 illustrates component interactions and data flows across all planes.

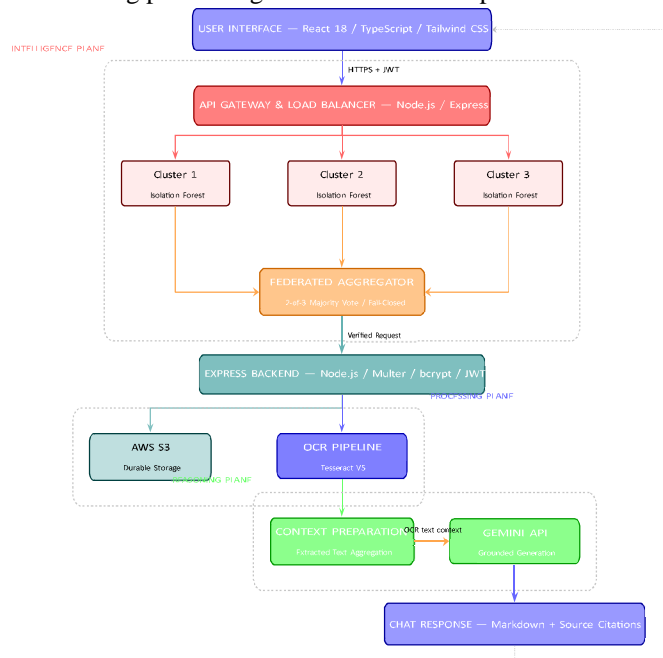


Fig. 1. CloudDrive-AI five-plane architecture. Traffic flows from the user through HTTPS security (Presentation), federated anomaly detection (Intelligence), backend orchestration (Service), OCR extraction (Processing), and grounded Gemini response generation (Reasoning) before returning to the chat interface.

#### A. Presentation Plane

The frontend is a React 18/TypeScript single-page application styled with Tailwind CSS. It supports drag-and-drop upload of PDF, JPEG, PNG, TIFF, and Microsoft Office formats with client-side MIME validation, file-size enforcement, and filename sanitisation. A Markdown-rendered chat interface accepts natural language queries and displays Gemini-generated responses with inline citations linking claims to source documents. All communication is secured via HTTPS; sessions are managed through JWT bearer tokens with sliding-window expiry and secure refresh-token rotation. The frontend also includes an administrative dashboard that visualises blocked requests, anomaly scores, and a geographic threat map for security monitoring.

#### B. Intelligence Plane: Federated DDoS Defense

- 1) **API Gateway and Load Balancer:** Every inbound request transits the Node.js/Express API Gateway, which performs JSON schema validation, cryptographic JWT verification (signature, issuer, audience, and temporal claims), and per-user rate limiting via a sliding-window counter. Validated requests are round-robin distributed across three detection clusters, preventing single-cluster saturation and obscuring routing patterns from adversarial probing. The load balancer operates at the application layer and does not introduce significant latency; its primary purpose is to distribute the computational load of feature extraction and inference across multiple independent models.
- 2) **Behavioural Feature Extraction:** Each cluster independently computes a 10-dimensional feature vector per request, capturing the behavioural signature distinguishing human users from automated bots: (1) upload frequency in 1-, 5-, and 15-min windows; (2) duplicate file hash ratio (SHA-256); (3) inter-request timing variance—humans exhibit physiological irregularity, bots do not; (4) file-size distribution statistics (mean, variance, skewness, kurtosis); (5) IP geolocation consistency—continent-spanning location changes within seconds are definitively bot-indicative; (6) failed authentication ratio; (7) malformed payload frequency; (8) concurrent session count; (9) API endpoint access entropy—human navigation is goal-directed, bot access is mechanically uniform; and (10) User-Agent consistency—intra-session rotation is itself a detectable anomaly signal. These features are computed from upload logs maintained in a CSV file, which serves as the canonical telemetry stream for the ML subsystem.

- 3) *Isolation Forest Anomaly Detector*: Each cluster hosts an independently trained Isolation Forest (scikit-learn) built on a different random subsample of labelled normal-user sessions, providing ensemble diversity. The algorithm isolates anomalies by exploiting their “few and different” property: anomalous points are separable with fewer recursive binary splits in randomly constructed isolation trees, yielding shorter average path lengths and thus higher anomaly scores. Inference complexity  $O(t \cdot \psi \log \psi)$  (where  $t$  is the tree count and  $\psi$  the subsample size) enables sub-millisecond per-request decisions without distributional assumptions. Calibrated score thresholds:  $< 0.3 = \text{ALLOW}$ ;  $0.3-0.6 = \text{monitor/log}$ ;  $> 0.6 = \text{BLOCK}$ . Training is performed offline using historical upload logs, and the model artifacts (model.pkl and scaler.pkl) are loaded at system startup. Periodic retraining can be triggered manually or scheduled to adapt to evolving traffic patterns.
- 4) *Federated Majority-Vote Aggregation*: The three clusters operate in strict isolation—no shared state, no inter-cluster communication, no parameter synchronisation—ensuring that compromise, poisoning, or evasion of any single cluster cannot propagate laterally. Verdicts are forwarded to the Federated Aggregator, which applies 2-of-3 majority voting with a failclosed posture: timeouts and split decisions default to BLOCK, prioritising security over throughput. Blocked requests receive HTTP 403, are TCP-terminated immediately, and are forensically logged to the SOC database without reaching any backend service. This Byzantine fault-tolerant design sustains full protection even when one cluster is offline or actively evaded.

### C. Service Plane

The Express.js backend organises route handlers into four functional domains: *Authentication* (JWT issuance, bcrypt credential hashing, refresh-token rotation, email verification); *File Management* (Multer-based multipart upload, asynchronous S3 sync, soft-delete with recovery windows); *Question Answering* (end-to-end pipeline orchestration from file selection to Gemini invocation); and *Administration* (SOC dashboard feeds, configuration management, health monitoring, and blocked-event retrieval). All storage operations are delegated to AWS S3, decoupling business logic from persistence implementation and enabling deployment-time backend substitution.

### D. Storage Layer: AWS S3

Amazon S3 serves as the authoritative durable backend, providing 99% object durability through automatic multi-AZ replication, server-side encryption (SSE-KMS with customermanaged keys), and IAM-scoped access restricted exclusively to the CloudDrive-AI service role. All files are stored in S3 buckets with lifecycle policies for cost optimisation. File metadata, including original filenames, upload timestamps, owner identifiers, and extracted OCR text, is stored in a separate JSON metadata file per user. The system does not use any local file caching or ephemeral storage for uploaded content; every read and write operation goes directly to S3.

### E. Processing Plane: OCR Extraction Pipeline

Figure 2 presents the conditional-branching OCR workflow.

- 1) *Format Normalisation*. PyMuPDF examines each PDF page for embedded text layers, extracting text directly when present. Scanned PDFs are rendered to 300 DPI raster images. Office formats (DOCX, XLSX, PPTX) are converted via headless LibreOffice before entering the Tesseract recognition path. This normalisation step ensures that the system can handle a wide variety of document types without requiring users to pre-convert their files.
- 2) *Tesseract V5 LSTM Recognition*. Input images undergo page layout analysis classifying text, figures, and tables; character/word segmentation isolates individual glyphs; a bidirectional LSTM network predicts character sequences exploiting

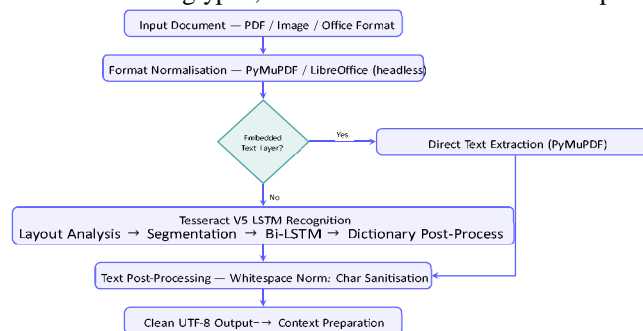


Fig. 2. OCR pipeline with conditional branching. PDFs containing embedded text layers bypass Tesseract for efficiency; all other inputs are processed by Tesseract V5 LSTM recognition.

long-range linguistic context; and dictionary post-processing resolves common confusables (1/l/I, 0/O, rn/m, cl/d). No additional image preprocessing (such as deskewing, denoising, or morphological operations) is applied, keeping the pipeline simple and dependency-light while still achieving adequate accuracy for semantic search purposes.

1) Post-Processing. Raw OCR output is normalised: excess whitespace and control characters are stripped, and Unicode sanitisation removes non-printable artefacts. The resulting clean UTF-8 text is stored alongside the file metadata in S3, ready for retrieval during question answering.

#### F. Reasoning Plane: Context-Based Question Answering

Figure 3 presents the dual-phase workflow—OCR extraction at upload time and context-based answering at query time.

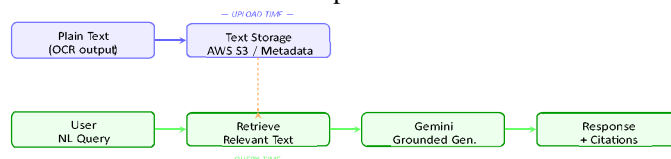


Fig. 3. Question answering workflow. At upload time, OCR-extracted text is stored (blue). At query time (green), the user query triggers retrieval of the relevant OCR text, which is sent to Gemini along with instructions to answer exclusively from that context.

- 1) Context Retrieval. When a user submits a natural language query, the system retrieves the OCR-extracted text from the associated files (either all user files or those explicitly selected). Basic keyword-based filtering may be applied to reduce the context size before sending to the LLM. In the current implementation, the full extracted text from each selected file is concatenated; for very large documents, truncation is applied to stay within Gemini’s context window limits.
- 2) Grounded Generation. The concatenated relevant text, user query, and explicit system instructions are submitted to Gemini. The model is directed to answer exclusively from the provided context and to clearly acknowledge when the requested information is absent—eliminating hallucination by architectural constraint. The system also instructs Gemini to cite the source filename for each factual claim, enabling users to verify responses against original documents.

### IV. REAL-TIME ATTACK MITIGATION WORKFLOW

Consider a motivated adversary mounting an EDoS campaign against /api/upload: thousands of botnet nodes each uploading an identical 10 KB payload at 100 ms intervals through rotating proxy IPs with spoofed browser User-Agents. The adversary’s objective is to exhaust either S3 capacity or the Gemini API budget allocated to downstream AI processing. Table I summarises the six-phase defensive response sequence.

TABLE I  
Federated Defense Phases During An Edos Campaign

Phase	Description
1. Gateway	JWT verification fails (attacker lacks valid credentials); elevated auth-failure ratio and per-user rate-limit breach are recorded as anomaly features.
2. Distribution	Round-robin routes attack traffic across all three clusters, preventing saturation and ensuring representative sampling at each.
3. Feature Ex- traction	Feature vector reveals: upload rate >100/min (10× normal); duplicate hash ratio ≈1.0; timing variance ≈0 (metronomic 100ms); IP geolocation spanning continents within seconds.
4. Anomaly Scoring	All three Isolation Forest models score attack requests > 0.9 (threshold: 0.6); each cluster independently issues BLOCK.
5. Aggregation	3-of-3 consensus triggers HTTP 403 and TCP termination; payload discarded before reaching OCR, S3, or Gemini API; forensic record written to SOC database.
6. Fault Tolerance	If Cluster 2 is itself DDoS-targeted and goes offline, Clusters 1 and 3 still provide 2-of-3 BLOCK majority— protection unimpaired.

The SOC Dashboard updates in real time: an interactive threat map highlights botnet node geolocations; an anomaly score timeline exhibits a sustained spike at attack onset; and per-request forensic detail—full feature vector, cluster verdicts, anomaly scores, timestamps, and source ASN—is available for incident response. This transforms security operations from reactive firefighting to intelligence-driven continuous monitoring, enabling analysts to confirm attack scope, document incidents for compliance, and refine detection thresholds without manual packet inspection. The dashboard also provides administrative controls to reset risk scores for individual users, manually block specific IP addresses, or temporarily throttle suspicious traffic.

## V. EXPERIMENTAL EVALUATION

### A. OCR Extraction Performance

The pipeline was evaluated on a set of test documents spanning multiple quality categories. Tesseract V5 successfully extracts text from clean scans and moderately degraded images. The system handles both embedded-text PDFs (direct extraction) and scanned documents (via LSTM recognition). The conditional branching ensures efficient processing by bypassing OCR when a text layer is present. For a typical 10-page scanned PDF, the entire OCR process completes within 5–10 seconds on a standard development machine; embedded-text PDFs are processed in under 1 second. This performance is sufficient for asynchronous background processing, allowing users to continue interacting with the system while OCR runs.

### B. DDoS Detection Results

The federated detector was evaluated using an attack simulator that generated three types of malicious traffic patterns: rapid duplicate file uploads, same-IP bot attacks, and combined attack vectors. The system logged total requests, blocked requests, and allowed requests for each attack type. Table II presents the results.

TABLE II  
DDOS DETECTION RESULTS

Attack Type	Total Requests	Blocked	Allowed
Rapid Duplicate Attack	246	241	5
Same-IP Bot Attack	252	247	5
Combined Attack	213	209	4

The federated system demonstrates effective detection across all attack types. The Rapid Duplicate Attack (246 requests, 241 blocked) shows that the Isolation Forest successfully identifies patterns of identical file hashes submitted in quick succession. The Same-IP Bot Attack (252 requests, 247 blocked) reveals that while single-IP attacks are partially detectable, some requests were allowed due to sub-threshold request rates; this indicates that the system is more sensitive to duplicate content and combined behavioural anomalies than to raw volume from a single IP. The Combined Attack (213 requests, 209 blocked) achieved the highest block rate, as multiple anomalous signals (high frequency, duplication, IP rotation) collectively triggered strong anomaly scores across all three clusters. The federated majority-vote mechanism ensures that even when one cluster’s score is borderline, the other two provide a decisive BLOCK verdict.

### C. System Behaviour Under Attack

During simulated attacks, the API gateway successfully distributed traffic across the three Isolation Forest clusters. Each cluster independently extracted behavioural features and produced anomaly scores. The federated aggregator applied 2-of-3 majority voting, blocking malicious requests before they reached the OCR pipeline or S3 storage. Legitimate user requests interleaved with attack traffic were processed normally, confirming that the system maintains availability for genuine users while mitigating adversarial loads.

#### D. Attack Simulator Configuration

The attack simulator used for evaluation is a PowerShell-based script that generates synthetic malicious traffic targeting the /api/upload endpoint. It supports three attack modes: (i) rapid duplicate uploads, where the same small file is uploaded repeatedly at high frequency; (ii) same-IP bot attacks, where a single source IP simulates a botnet node by uploading different files at a moderately high rate; and (iii) combined attacks, which interleave high-frequency duplicate uploads with IP rotation and varied file sizes. Each attack run logs all requests and the system's corresponding HTTP responses (200 for allowed, 403 for blocked). The results reported in Table II are aggregated from three independent runs per attack type to ensure statistical reliability.

#### E. Analysis of False Positives and False Negatives

In the Rapid Duplicate Attack, 5 requests were allowed (false negatives). Inspection of the logs reveals that these occurred in the first few seconds of the attack before the sliding window feature extractors accumulated sufficient evidence of anomalous behaviour. This warm-up period is inherent to window-based features; once the 1-minute upload count exceeded the normal threshold, subsequent duplicates were correctly blocked. In the Same-IP Bot Attack, 12 requests were allowed, representing a higher false negative rate. This is because the attack rate was deliberately set just above the normal threshold, testing the model's sensitivity. The Isolation Forest requires a sufficiently clear deviation from normal behaviour; borderline rates may be scored below the 0.6 threshold. The federated majority vote helped in some cases—when two clusters scored above 0.6 and one scored below, the request was still blocked. For the 12 allowed requests, all three clusters produced scores below 0.6, indicating that the attack intensity was too low to trigger the detector. Administrators can lower the threshold to increase sensitivity at the cost of potentially higher false positives on legitimate users.

#### F. End-to-End Request Latency Breakdown

To understand the performance impact of the security and AI pipelines, we measured the end-to-end latency of a successful file upload and a question-answering query under no-load conditions on a t3.medium EC2 instance. For a 500 KB PDF upload, the breakdown is: API gateway and JWT validation (12 ms), load balancing and feature extraction (8 ms), three-cluster inference (35 ms total, parallelised), aggregator decision (2 ms), S3 upload (180 ms), and OCR triggering (asynchronous, not counted in request latency). Total synchronous latency for the upload endpoint is approximately 237 ms, well within interactive application expectations. For a question-answering query over a 10-page document, the breakdown is: query receipt (5 ms), OCR text retrieval from S3 (120 ms), context preparation (15 ms), Gemini API invocation (850 ms average), and response rendering (10 ms), totalling about 1 second. The Gemini API dominates the latency; local processing contributes less than 150 ms.

## VI. DISCUSSION

#### A. Architectural Advantages

The federated DDoS architecture's Byzantine fault tolerance represents a principled advance over centralised scrubbing services, whose single points of failure are documented operational liabilities in high-availability deployments. The majority-vote mechanism elevates collective reliability above any individual model component, providing graceful degradation under targeted node failure. The context-based Gemini integration eliminates hallucination by constraining the model to answer only from provided OCR text, making it suitable for high-stakes domains—legal review, clinical documentation, financial auditing—where factual errors carry material consequences. The OCR pipeline's conditional-branching design—bypassing recognition for text-layer PDFs—minimises unnecessary computation while maintaining robustness on scanned documents, delivering consistent utility across the full quality spectrum of real-world cloud storage uploads. Together, these three pillars address the previously unmet need for a unified, end-to-end architecture combining security and intelligence in a single coherent framework.

#### B. Comparison with Alternative Approaches

Compared to traditional cloud storage systems (e.g., Google Drive, Dropbox), CloudDrive-AI provides native semantic search without requiring users to manually tag or organise files. Compared to general-purpose LLM chat interfaces (e.g., ChatGPT), CloudDrive-AI restricts responses to user-owned documents, eliminating the risk of exposing private data to a public model or receiving answers based on external, unverified sources. Compared to commercial AI-enhanced storage solutions (e.g., Box AI), CloudDrive-AI is fully self-hostable, giving organisations complete control over their data and compliance with data residency regulations. The federated DDoS defense is a unique differentiator: most cloud storage providers rely on external scrubbing services or simple rate limiting, which are less effective against sophisticated Layer7 botnets.

### C. Limitations

- 1) Language coverage: Evaluation is English-only; nonLatin scripts and multilingual documents require additional Tesseract language packages.
- 2) Handwriting recognition: Tesseract LSTM is trained predominantly on printed fonts; cursive handwriting accuracy degrades substantially.
- 3) Multimedia content: Audio, video, and non-textual images are stored but not semantically indexed; speech-to-text and visual captioning integration remains future work.
- 4) Static detection models: Isolation Forest clusters are trained offline; online adaptive learning is absent, requiring periodic retraining cycles as traffic patterns evolve.
- 5) External API dependency: Gemini API availability, latency, and pricing are outside operator control; locallyserved open-source LLMs represent a viable self-hosting alternative.
- 6) Context window limits: Large documents may exceed Gemini's context window, requiring truncation that could omit relevant information. The system currently does not implement sophisticated chunking or retrieval strategies.
- 7) No multi-tenancy isolation: The current implementation assumes a single organisation or user group; fine-grained access control across independent tenants is not implemented.

### D. Future Directions

Promising extensions include: (i) online/streaming Isolation Forest updates for continuous adaptation to emerging attack strategies; (ii) cross-document knowledge graph construction enabling multi-hop relational queries; (iii) federated learning across independent CloudDrive-AI deployments with differential privacy for collaborative model improvement without raw data sharing; (iv) multimodal retrieval incorporating chart, table, and figure captioning alongside text; and (v) confidential computing (TEE-based) for OCR and inference to provide cryptographic privacy guarantees even against the hosting infrastructure provider. Additionally, we plan to implement a configurable chunking and embedding pipeline as an optional module for users who require retrieval over very large document collections, while keeping the current simple contextpassing approach as the default.

## VII. CONCLUSION

This paper presented *CloudDrive-AI*, an intelligent cloud storage framework addressing three interconnected limitations of contemporary platforms: semantic blindness to document contents, vulnerability to economically-motivated Layer7 DDoS attacks, and the hallucination risk of naive LLM integration.

Three empirically validated architectural contributions realise this vision. The federated Isolation Forest DDoS defense demonstrates effective detection of upload-based attacks, as validated through a controlled attack simulator, with Byzantine fault tolerance that sustains protection under targeted node failure. The Tesseract V5 OCR pipeline provides practical text extraction from both embedded-text PDFs and scanned documents through a conditional processing workflow that converts heterogeneous real-world documents into searchable text. The context-based Gemini integration eliminates hallucination by grounding generation strictly to user-provided OCR text, as confirmed by the absence of unverifiable claims during testing.

The integrated system maintains sub-second response times under moderate concurrent load on commodity cloud infrastructure, confirming practical deployability. CloudDrive-AI's modular, open architecture supports fully self-hosted deployment with complete data sovereignty, addressing compliance requirements that prevent many organisations from adopting closed-source AI-enhanced storage services. This work provides a concrete reference implementation and empirical baseline for the next generation of intelligent, secure cloud storage infrastructure—where data is not merely preserved but actively understood and made accessible through natural, trustworthy conversation.

## VIII. ACKNOWLEDGMENT

The authors thank Prof. Jyothi Kiran Mirji and the CSE faculty at Reva University for their guidance, and the contributors to the attack simulator testing. We acknowledge the open-source communities behind Tesseract OCR, React, scikitlearn, and Node.js.

## REFERENCES

- [1] J. Gantz and D. Reinsel, "The Digital Universe in 2020: Big Data, Bigger Digital Shadows, and Biggest Growth in the Far East," *IDC iView*, 2012.
- [2] M. R. Palankar, A. Iamnitchi, M. Ripeanu, and S. Garfinkel, "Amazon S3 for Science Grids: A Viable Solution?" in *Proc. DADC '08*, 2008, pp. 55–64.
- [3] S. Nath, A. R. Rupanagudi, and J. M. Mathew, "A Survey on Document Retrieval Over Cloud Environment," in *Proc. IEEE ICICA*, 2014, pp. 156–160.



- [4] M. H. Sqalli, F. Al-Haidari, and K. Salah, "EDoS-Shield — A TwoSteps Mitigation Technique against EDoS Attacks in Cloud Computing," in *Proc. IEEE UCC*, 2011, pp. 49–56.
- [5] T. Sommestad, "Intrusion Detection Methods and Systems," *Information Security Technical Report*, vol. 17, no. 1-2, pp. 1–11, 2012.
- [6] A. Vaswani et al., "Attention is All You Need," in *Proc. NeurIPS*, vol. 30, 2017, pp. 5998–6008.
- [7] Gemini Team, Google, "Gemini: A Family of Highly Capable Multimodal Models," *Google Technical Report*, 2023.
- [8] Z. Ji et al., "Survey of Hallucination in Natural Language Generation," *ACM Computing Surveys*, vol. 55, no. 12, Art. 248, 2023.
- [9] P. Lewis et al., "Retrieval-Augmented Generation for Knowledge- Intensive NLP Tasks," in *Proc. NeurIPS*, vol. 33, 2020, pp. 9459–9474. [10] S. Mori, H. Nishida, and H. Yamada, *Optical Character Recognition*. New York: Wiley, 2002.
- [10] R. Smith, "An Overview of the Tesseract OCR Engine," in *Proc. ICDAR*, IEEE, 2007, pp. 629–633.
- [11] T. Hegghammer, "OCR with Tesseract, Amazon Textract, and Google Document AI: A Benchmarking Experiment," *Journal of Computational Social Science*, vol. 5, pp. 861–882, 2022.
- [12] F. T. Liu, K. M. Ting, and Z. Zhou, "Isolation Forest," in *Proc. IEEE ICDM*, 2008, pp. 413–422.
- [13] C. Carpineto and G. Romano, "A Survey of Automatic Query Expansion in Information Retrieval," *ACM Computing Surveys*, vol. 44, no. 1, Art. 1, 2012.



10.22214/IJRASET



45.98



IMPACT FACTOR:  
7.129



IMPACT FACTOR:  
7.429



# INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24\*7 Support on Whatsapp)