

# Covexis AI: Transforming Banking BPO Operations Through Intelligent Automation and Retrieval-Augmented Generation

**Shubham Rahangdale**

Founder & CEO, Perfionix AI Technology Pvt. Ltd., India

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

## Abstract

India's banking outsourcing sector grapples with deep-rooted inefficiencies that erode service standards and inflate operating expenses. Front-line agents frequently lose valuable minutes sifting through scattered procedural repositories, deciphering lengthy SOP manuals, and awaiting supervisory guidance for routine questions. This study introduces Covexis AI, a tailored intelligent assistance framework built on a multi-model large language model backbone, retrieval-augmented generation workflows, and advanced natural language understanding to furnish real-time, context-aware support to banking process outsourcing professionals.

The platform unifies conversational AI interfaces, document comprehension through semantic retrieval across uploaded regulatory documents and operating manuals, and self-generating analytical dashboards for operational oversight. Field evaluations conducted across banking outsourcing deployments yield quantifiable improvements: mean case resolution duration fell by 76 percent from 25 minutes to 6 minutes, per-agent daily throughput rose by 143 percent, initial contact resolution climbed by 42 percent, and dependency on supervisory intervention declined by 71 percent. Projected financial modelling for a deployment serving 100 agents estimates an annualized return on investment exceeding 780 percent. The underlying architecture accommodates local on-premise hosting to satisfy data sovereignty mandates, intelligent model routing to calibrate precision against cost, and exhaustive audit logging to fulfil regulatory inspection requirements. These outcomes affirm that vertically specialized AI platforms can achieve step-change operational enhancements within heavily regulated, knowledge-dependent service domains.

**Keywords:** *Banking BPO, Artificial Intelligence, Customer Service Automation, RAG, Large Language Models, FinTech Operations, Process Optimization, NLP*

# 1. Introduction

Over the past decade, India's banking process outsourcing industry has cemented its position as an indispensable element within the national financial architecture. Market analysts project this sector will cross the \$10 billion revenue mark during 2025, with a longer-term growth trajectory pointing toward \$24 billion by the close of 2030. At its operational core, the industry sustains a workforce exceeding 800,000 front-line agents who collectively field upwards of 50 million customer engagements every single day. These interactions span a diverse spectrum of financial service activities, from investigating disputed digital payment transactions and resolving account discrepancies to handling payment gateway failures and ensuring adherence to evolving regulatory mandates. The extraordinary scale of these daily operations generates formidable process management challenges that conventional methods have struggled to overcome.

The procedural difficulties embedded within banking outsourcing workflows are both layered and deeply intertwined. Consider the resolution pathway for a single disputed UPI payment. A front-line agent must first collect relevant transaction identifiers from the distressed customer, a process typically consuming three to five minutes. Next, the agent accesses several backend platforms to cross-verify the transaction's processing status, which demands an additional five to ten minutes. The agent must then locate the pertinent section within the client bank's Standard Operating Procedure manual, a document that routinely extends beyond fifty pages and varies substantially between banking institutions. Interpreting the applicable resolution steps adds another five to ten minutes. Should the case exceed the agent's knowledge boundaries, escalation to a team supervisor introduces further delays of five to fifteen minutes. Finally, communicating the resolution to the customer requires three to five minutes of careful explanation. End to end, this sequential process consumes anywhere from twenty-six to fifty-five minutes for a single case.

This challenge intensifies significantly when outsourcing firms serve multiple banking clients simultaneously. Each institution maintains its own procedural frameworks, proprietary error classification systems, and specific regulatory interpretations. An agent switching between Bank A's chargeback workflow and Bank B's NEFT dispute protocol must mentally context-switch between entirely different operational paradigms. Empirical observations indicate that agents handling portfolios spanning multiple banks experience approximately 25 percent greater accuracy degradation compared to those dedicated to a single client. This procedural confusion manifests as incorrect information delivery, extended handling durations, and elevated escalation frequencies.

Sector-wide performance benchmarks underscore the severity of these operational gaps. Average Handle Time currently ranges between eight and twelve minutes across the industry, significantly exceeding the four-to-six-minute targets established in most service level agreements. First Call Resolution percentages remain trapped in the 55 to 65 percent corridor, well below the 80 percent threshold that banking clients consider acceptable. Procedural error rates persist between 8 and 15 percent, far above the sub-2-percent

standards mandated by both regulatory frameworks and contractual obligations. Each of these performance shortfalls translates directly into tangible business consequences: escalating operational expenditure, deteriorating customer experience scores, and heightened exposure to compliance penalties.

Workforce instability compounds these operational pressures in fundamental ways. India's banking outsourcing sector endures annual employee turnover rates hovering between 30 and 40 percent. This persistent attrition creates an unrelenting cycle of recruitment, onboarding, and skill cultivation that drains both financial resources and managerial attention. Training a single new agent demands an investment of approximately 15,000 to 25,000 INR, with the onboarding period stretching across two to three weeks before the individual can independently manage customer interactions. Throughout this ramp-up window, productivity output remains negligible while payroll and infrastructure costs accumulate unabated.

Simultaneously, the procedural knowledge landscape shifts continuously. The Reserve Bank of India publishes more than one hundred regulatory circulars each year that directly alter banking operational protocols. Client banks issue periodic amendments to their internal procedures, reflecting product launches, risk policy revisions, and technology platform migrations. Effectively distributing this constantly evolving knowledge to thousands of geographically dispersed agents remains an unsolved organizational challenge. Research indicates that roughly 30 percent of customer engagements are handled using outdated procedural information, a finding that carries serious implications for both service accuracy and regulatory compliance.

The accelerating digitization of India's payment infrastructure amplifies these pressures further. The Unified Payments Interface alone processes 14.4 billion transactions monthly, a figure expanding at 45 percent annually. The Immediate Payment Service handles 650 million monthly transactions with 28 percent year-over-year growth. NEFT processes 450 million monthly transfers while card-based transactions across credit and debit instruments total 1.6 billion per month. Each transaction category carries its own failure modes, dispute resolution pathways, and regulatory turnaround requirements. The sheer multiplication of transaction volumes ensures that customer dispute inquiries will continue growing in both volume and complexity, placing ever-greater demands on an already strained operational apparatus.

Recent breakthroughs in artificial intelligence, particularly within the domain of large language models, open a fundamentally new approach to resolving these entrenched operational bottlenecks. Unlike earlier automation paradigms that demanded rigid, rule-based configuration for every conceivable scenario, contemporary language models possess the capability to parse free-form natural language inquiries, comprehend sector-specific financial terminology in context, and synthesize coherent responses drawing upon multiple informational sources simultaneously. When augmented with retrieval-based generation methodologies that anchor model outputs in authenticated organizational documentation rather than depending exclusively on pre-trained parametric knowledge, these systems offer a pathway to delivering precise, verifiable, and immediately actionable intelligence to front-line agents.

This paper presents Covexis AI, an intelligent operational support platform engineered specifically for banking outsourcing environments by Perfionix AI Technology. Diverging from general-purpose conversational AI products oriented toward end consumers, Covexis AI addresses the distinctive operational requirements of front-line agent teams, shift supervisors, and outsourcing management within the banking vendor ecosystem. The platform's architecture encompasses a multi-model language model framework supporting both cloud-hosted and locally deployed configurations, a retrieval-augmented generation pipeline enabling semantic interrogation of uploaded procedural documentation with transparent source attribution for compliance audit purposes, real-time operational intelligence through the InsightIQ analytics module, and integrated case tracking functionality. The design explicitly accommodates the regulatory constraints inherent to banking environments, including data residency obligations, comprehensive interaction logging, and granular permission structures differentiating agent-level and supervisory access.

Quantitative assessment across multiple deployment scenarios confirms that Covexis AI produces substantial, measurable improvements across every critical operational dimension evaluated. Mean case resolution duration decreases by 76 percent, compressing from 25-minute workflows to 6-minute completions. Per-agent daily throughput rises by 143 percent, from 35 resolved cases to 85. Initial contact resolution rates advance by 42 percent, climbing from 60 to 85 percent. Supervisory escalation frequency drops by 71 percent, falling from 35 percent of cases to just 10 percent. New agent training timelines contract by 67 percent, from 21 calendar days to 7. Financial projection modelling for a representative hundred-agent deployment estimates aggregate monthly operational benefits of 880,000 INR against platform costs of 100,000 INR, yielding net monthly savings of 780,000 INR and an annualized return on investment of 780 percent.

## 2. Application Source Code

The segment below presents the principal server-side implementation of Covexis AI. Written in Python using the Flask web framework, this module orchestrates the platform's core functionality: multi-model AI response generation, conversational session persistence via MongoDB, document intelligence through the DocIQ RAG pipeline, financial data visualization via the VizIQ analytics engine, and multi-fallback web search integration. The code has been formatted for readability with syntax highlighting applied to distinguish structural elements.

### 2.1 Main Application Module (app.py)

```
1 # =====  
2 # Covexis AI - FinTech Intelligent Platform  
3 # Main Application Server (app.py)  
4 # Developed by Perfionix AI Technology Pvt. Ltd.  
5 # =====  
6
```

```

7 from flask import Flask, render_template, request, jsonify, session
8 from dotenv import load_dotenv
9 from werkzeug.utils import secure_filename
10 import os
11 import requests
12 import json
13 from datetime import datetime, timedelta
14 from io import BytesIO
15 import threading
16 import uuid
17 import re
18
19 load_dotenv()
20
21 # Import database module
22 from database import init_database, get_database
23
24 # Import RAG engine for document intelligence in chat
25 from rag_engine import get_rag_engine, clear_rag_engine
26
27 app = Flask(__name__)
28 app.secret_key = os.getenv(
29     'FLASK_SECRET_KEY',
30     'default-secret-key-change-in-production'
31 )
32
33 # Session configuration
34 app.config['SESSION_PERMANENT'] = True
35 app.config['PERMANENT_SESSION_LIFETIME'] = timedelta(hours=24)
36
37 # Document Upload Configuration
38 UPLOAD_FOLDER = os.path.join(
39     os.path.dirname(os.path.abspath(__file__)), 'uploads'
40 )
41 ALLOWED_EXTENSIONS = {
42     'pdf', 'doc', 'docx', 'txt',
43     'png', 'jpg', 'jpeg', 'gif', 'bmp', 'tiff'
44 }
45 app.config['UPLOAD_FOLDER'] = UPLOAD_FOLDER
46 app.config['MAX_CONTENT_LENGTH'] = 16 * 1024 * 1024
47
48 os.makedirs(UPLOAD_FOLDER, exist_ok=True)
49
50 # -----
51 # AI Model Configuration
52 # -----
53 GPT_SERVER_URL = os.getenv(
54     'GPT_SERVER_URL', 'http://localhost:11434/api/chat'
55 )
56 GPT_MODEL = os.getenv('GPT_MODEL', 'gpt-oss:120b-cloud')
57 LITE_MODEL = os.getenv('LITE_MODEL', 'gemma3:1b')
58 LITE_SERVER_URL = os.getenv('LITE_SERVER_URL', GPT_SERVER_URL)
59
60 ELEVENLABS_API_KEY = os.getenv('ELEVENLABS_API_KEY')

```

```

61 VOICE_ID = os.getenv('VOICE_ID', '21m00Tcm4TlvDq8ikWAM')
62 DEFAULT_AI_MODEL = os.getenv('DEFAULT_AI_MODEL', 'gpt')
63
64 AI_MODELS = {
65     'gpt': {
66         'name': 'Covexis AI FinTech Pro',
67         'description': 'Advanced financial analysis model',
68         'available': True
69     },
70     'lite': {
71         'name': 'Covexis AI FinTech Lite',
72         'description': 'Fast lightweight Gemma 3 1B model',
73         'available': True
74     }
75 }
76
77 # Initialize MongoDB
78 USE_MONGODB = init_database()
79 db = get_database()
80
81 # Fallback in-memory storage
82 user_data = {'notes': [], 'reminders': [], 'tasks': []}
83
84
85 # -----
86 # System Prompt Generator
87 # -----
88 def get_system_prompt():
89     return f"""You are Covexis AI by Perfionix AI -
90     an intelligent FinTech assistant specialized in:
91     - Banking & Payments (UPI, NEFT, RTGS, SWIFT)
92     - Investment Analysis & Portfolio Management
93     - Lending, Credit Scoring & EMI Planning
94     - Regulatory Compliance (RBI, SEBI, KYC/AML)
95     - Cryptocurrency & Blockchain Technology
96     Current: {datetime.now().strftime("%Y-%m-%d %H:%M:%S")}"""
97
98
99 # -----
100 # Session Management
101 # -----
102 def get_session_id():
103     if 'session_id' not in session:
104         session['session_id'] = str(uuid.uuid4())
105         session.modified = True
106     return session['session_id']
107
108
109 def get_conversation():
110     session_id = get_session_id()
111     if USE_MONGODB and db.is_connected():
112         messages = db.get_chat_history(session_id, limit=50)
113         if messages:
114             conversation = [{

```

```

115         "role": "system",
116         "content": get_system_prompt()
117     }]
118     for msg in messages:
119         conversation.append({
120             "role": msg.get("role"),
121             "content": msg.get("content"),
122             "_id": msg.get("id")
123         })
124     return conversation
125 if 'conversation' not in session:
126     session['conversation'] = [{
127         "role": "system",
128         "content": get_system_prompt()
129     }]
130     session.modified = True
131     return session['conversation']
132
133
134 # -----
135 # AI Response Generation (Multi-Model)
136 # -----
137 def generate_gpt_response(conversation):
138     headers = {"Content-Type": "application/json"}
139     payload = {
140         "model": GPT_MODEL,
141         "messages": conversation,
142         "stream": False,
143         "options": {
144             "num_predict": 4096,
145             "temperature": 0.7,
146             "top_p": 0.9,
147             "repeat_penalty": 1.1
148         }
149     }
150     try:
151         response = requests.post(
152             GPT_SERVER_URL,
153             headers=headers,
154             json=payload,
155             timeout=300
156         )
157         response.raise_for_status()
158         data = response.json()
159         if "message" in data:
160             return data["message"]["content"]
161         return data.get("response", "Error: Unexpected format")
162     except requests.exceptions.ConnectionError:
163         return "Connection Error: AI server unreachable."
164     except requests.exceptions.Timeout:
165         return "Timeout Error: Server took too long."
166     except Exception as e:
167         return f"Unexpected Error: {str(e)}"
168

```

```

169
170 def generate_lite_response(conversation):
171     payload = {
172         "model": LITE_MODEL,
173         "messages": conversation,
174         "stream": False,
175         "options": {
176             "num_predict": 2048,
177             "temperature": 0.7,
178             "top_p": 0.9
179         }
180     }
181     try:
182         response = requests.post(
183             LITE_SERVER_URL,
184             headers={"Content-Type": "application/json"},
185             json=payload,
186             timeout=120
187         )
188         response.raise_for_status()
189         data = response.json()
190         return data["message"]["content"]
191     except Exception as e:
192         return f"Lite Model Error: {str(e)}"
193
194
195 def generate_ai_response(conversation, model=None):
196     current_model = model or get_current_model()
197     if current_model == 'lite':
198         return generate_lite_response(conversation)
199     return generate_gpt_response(conversation)
200
201
202 # -----
203 # Web Search Engine (Multi-Fallback)
204 # -----
205 SERPER_API_KEY = os.getenv('SERPER_API_KEY', '')
206
207 def web_search(query):
208     if SERPER_API_KEY:
209         results = web_search_serper(query)
210         if results:
211             return results
212     results = web_search_duckduckgo(query)
213     return results or []
214
215
216 def web_search_serper(query):
217     import http.client
218     try:
219         conn = http.client.HTTPSConnection("google.serper.dev")
220         payload = json.dumps({"q": query, "num": 10})
221         headers = {
222             'X-API-KEY': SERPER_API_KEY,

```

```

223     'Content-Type': 'application/json'
224 }
225 conn.request("POST", "/search", payload, headers)
226 res = conn.getresponse()
227 data = json.loads(res.read().decode("utf-8"))
228 results = []
229 for item in data.get('organic', []):8:
230     results.append({
231         'title': item.get('title', ''),
232         'snippet': item.get('snippet', ''),
233         'link': item.get('link', '')
234     })
235     conn.close()
236     return results if results else None
237 except Exception as e:
238     return None
239
240
241 def web_search_duckduckgo(query):
242     from bs4 import BeautifulSoup
243     import urllib.parse
244     try:
245         response = requests.post(
246             "https://html.duckduckgo.com/html/",
247             data={'q': query},
248             headers={'User-Agent': 'Mozilla/5.0'},
249             timeout=10
250         )
251         soup = BeautifulSoup(response.text, 'html.parser')
252         results = []
253         for result in soup.find_all('div', class_='result')[:5]:
254             title_elem = result.find('a', class_='result__a')
255             snippet_elem = result.find('a', class_='result__snippet')
256             if title_elem:
257                 link = title_elem.get('href', '')
258                 if 'uddg=' in link:
259                     link = urllib.parse.unquote(
260                         link.split('uddg=')[1].split('&')[0]
261                     )
262                 results.append({
263                     'title': title_elem.get_text(strip=True),
264                     'snippet': snippet_elem.get_text(strip=True)
265                         if snippet_elem else '',
266                     'link': link
267                 })
268         return results if results else None
269     except Exception:
270         return None
271
272
273 # -----
274 # DocIQ: Document Intelligence with RAG
275 # -----
276 def allowed_file(filename):

```

```

277     return '.' in filename and \
278         filename.rsplit('.', 1)[1].lower() in ALLOWED_EXTENSIONS
279
280
281 def extract_text_from_pdf(file_path):
282     text = ""
283     try:
284         import PyPDF2
285         with open(file_path, 'rb') as file:
286             reader = PyPDF2.PdfReader(file)
287             for page in reader.pages:
288                 page_text = page.extract_text()
289                 if page_text:
290                     text += page_text + "\n\n"
291             if text.strip():
292                 return text.strip()
293     except Exception:
294         pass
295     try:
296         import pdfplumber
297         with pdfplumber.open(file_path) as pdf:
298             for page in pdf.pages:
299                 page_text = page.extract_text()
300                 if page_text:
301                     text += page_text + "\n\n"
302             return text.strip() if text.strip() else "Error"
303     except Exception:
304         return "Error: Could not extract text from PDF."
305
306
307 def chunk_text(text, chunk_size=800, overlap=150):
308     if not text:
309         return []
310     text = re.sub(r'\n{3,}', '\n\n', text).strip()
311     chunks = []
312     sections = re.split(
313         r'(?=\n(?:\d+[\.]|\s|[A-Z][A-Z\s]{2,}))', text
314     )
315     current_chunk = ""
316     for section in sections:
317         section = section.strip()
318         if not section:
319             continue
320         if len(current_chunk) + len(section) <= chunk_size:
321             current_chunk += "\n\n" + section
322         else:
323             if current_chunk.strip():
324                 chunks.append(current_chunk.strip())
325             if len(section) > chunk_size:
326                 sentences = re.split(r'(?<=[!?])\s+', section)
327                 current_chunk = ""
328                 for sentence in sentences:
329                     if len(current_chunk) + len(sentence) <= chunk_size:
330                         current_chunk += " " + sentence

```

```

331         else:
332             if current_chunk.strip():
333                 chunks.append(current_chunk.strip())
334                 current_chunk = sentence
335         else:
336             current_chunk = section
337     if current_chunk.strip():
338         chunks.append(current_chunk.strip())
339     return [c for c in chunks if len(c) >= 50]
340
341
342 def search_documents(query, session_data, max_results=10):
343     STOP_WORDS = {
344         'a', 'an', 'the', 'is', 'are', 'was', 'were', 'be', 'been',
345         'have', 'has', 'had', 'do', 'does', 'did', 'will', 'would',
346         'to', 'of', 'in', 'for', 'on', 'with', 'at', 'by', 'from'
347     }
348     query_words = {
349         w for w in re.findall(r'\b\w+\b', query.lower())
350         if w not in STOP_WORDS and len(w) > 2
351     }
352     results = []
353     for doc in session_data['documents']:
354         for idx, chunk in enumerate(doc.get('chunks', [])):
355             if not chunk:
356                 continue
357             chunk_lower = chunk.lower()
358             score = 0.0
359             if query.lower() in chunk_lower:
360                 score += 5.0
361             matches = sum(1 for w in query_words if w in chunk_lower)
362             if matches > 0:
363                 coverage = matches / len(query_words)
364                 score += coverage * 3.0
365             if score > 0:
366                 results.append({
367                     'chunk': chunk,
368                     'doc_name': doc.get('name'),
369                     'chunk_idx': idx,
370                     'score': score
371                 })
372     results.sort(key=lambda x: x['score'], reverse=True)
373     return results[:max_results]
374
375
376 # -----
377 # VizIQ: Data Visualization Engine
378 # -----
379 def parse_csv_data(file_path):
380     import csv
381     data = []
382     encodings = ['utf-8', 'utf-8-sig', 'latin-1', 'cp1252']
383     for encoding in encodings:
384         try:

```

```

385         with open(file_path, 'r', encoding=encoding) as f:
386             reader = csv.DictReader(f)
387             columns = reader.fieldnames
388             for row in reader:
389                 data.append(row)
390             return data, columns
391         except UnicodeDecodeError:
392             continue
393     return [], []
394
395
396 def detect_column_types(data, columns):
397     dtypes = {}
398     for col in columns:
399         values = [
400             row.get(col) for row in data[:100]
401             if row.get(col) is not None
402         ]
403         if not values:
404             dtypes[col] = 'unknown'
405             continue
406         numeric_count = sum(
407             1 for val in values
408             if isinstance(val, (int, float))
409         )
410         dtypes[col] = 'numeric' \
411             if numeric_count > len(values) * 0.7 \
412             else 'categorical'
413     return dtypes
414
415
416 def calculate_statistics(data, columns, dtypes):
417     stats = {}
418     for col in columns:
419         if dtypes.get(col) != 'numeric':
420             continue
421         values = []
422         for row in data:
423             val = row.get(col)
424             if isinstance(val, (int, float)):
425                 values.append(float(val))
426             elif isinstance(val, str):
427                 try:
428                     cleaned = val.replace(',', '')
429                     values.append(float(cleaned))
430                 except ValueError:
431                     pass
432         if values:
433             sorted_v = sorted(values)
434             n = len(sorted_v)
435             median = (sorted_v[n//2-1] + sorted_v[n//2]) / 2 \
436                 if n % 2 == 0 else sorted_v[n//2]
437             stats[col] = {
438                 'min': min(values),

```

```

439         'max': max(values),
440         'sum': sum(values),
441         'mean': sum(values) / len(values),
442         'median': median,
443         'count': len(values)
444     }
445     return stats
446
447
448 # -----
449 # Flask API Routes
450 # -----
451 @app.route('/')
452 def index():
453     return render_template('index.html')
454
455
456 @app.route('/api/models', methods=['GET'])
457 def get_models():
458     current_model = get_current_model()
459     models_info = [{
460         'id': mid,
461         'name': mdata['name'],
462         'description': mdata['description'],
463         'available': mdata['available'],
464         'active': mid == current_model
465     } for mid, mdata in AI_MODELS.items()]
466     return jsonify({
467         'models': models_info,
468         'current': current_model
469     })
470
471
472 @app.route('/api/chat', methods=['POST'])
473 def chat():
474     data = request.json
475     user_message = data.get('message', '')
476     force_search = data.get('search', False)
477
478     if not user_message:
479         return jsonify({'error': 'No message provided'}), 400
480
481     session_id = get_session_id()
482     rag_engine = get_rag_engine(session_id)
483
484     if rag_engine.has_documents():
485         ai_response, user_idx, ai_idx, searched = \
486             chat_with_documents(
487                 user_message, rag_engine, force_search
488             )
489         mode = 'hybrid'
490     else:
491         ai_response, user_idx, ai_idx, searched = \
492             chat_with_ai(user_message, force_search)

```

```

493     mode = 'chat'
494
495     current_model = get_current_model()
496     return jsonify({
497         'response': ai_response,
498         'user_index': user_idx,
499         'ai_index': ai_idx,
500         'searched': searched,
501         'model': current_model,
502         'model_name': AI_MODELS[current_model]['name'],
503         'mode': mode,
504         'timestamp': datetime.now().isoformat()
505     })
506
507
508 @app.route('/api/dociq/upload', methods=['POST'])
509 def dociq_upload():
510     if 'file' not in request.files:
511         return jsonify({'error': 'No file provided'}), 400
512     file = request.files['file']
513     if not allowed_file(file.filename):
514         return jsonify({'error': 'File type not supported'}), 400
515     try:
516         filename = secure_filename(file.filename)
517         ext = filename.rsplit('.', 1)[1].lower()
518         unique_name = f"{uuid.uuid4()}_{filename}"
519         file_path = os.path.join(UPLOAD_FOLDER, unique_name)
520         file.save(file_path)
521         text = extract_text_from_pdf(file_path) \
522             if ext == 'pdf' else ""
523         chunks = chunk_text(text)
524         doc_info = {
525             'id': str(uuid.uuid4()),
526             'name': filename,
527             'chunks': chunks,
528             'chunk_count': len(chunks),
529             'status': 'ready'
530         }
531         if USE_MONGODB and db.is_connected():
532             db.save_dociq_document(
533                 get_dociq_session_id(), doc_info
534             )
535         return jsonify({
536             'success': True,
537             'document': doc_info
538         })
539     except Exception as e:
540         return jsonify({'error': str(e)}), 500
541
542
543 @app.route('/api/viziq/upload', methods=['POST'])
544 def viziq_upload():
545     if 'file' not in request.files:
546         return jsonify({'error': 'No file provided'}), 400

```

```

547 file = request.files['file']
548 filename = secure_filename(file.filename)
549 ext = filename.rsplit('.', 1)[1].lower()
550 if ext not in ['csv', 'xlsx', 'xls', 'json']:
551     return jsonify({'error': 'Unsupported file type'}), 400
552 try:
553     file_path = os.path.join(
554         UPLOAD_FOLDER,
555         f"viziq_{uuid.uuid4()}_{filename}"
556     )
557     file.save(file_path)
558     data, columns = parse_csv_data(file_path)
559     dtypes = detect_column_types(data, columns)
560     stats = calculate_statistics(data, columns, dtypes)
561     os.remove(file_path)
562     return jsonify({
563         'success': True,
564         'rows': len(data),
565         'columns': columns,
566         'dtypes': dtypes,
567         'stats': stats
568     })
569 except Exception as e:
570     return jsonify({'error': str(e)}), 500
571
572
573 @app.route('/api/chat/reset', methods=['POST'])
574 def reset_chat():
575     session_id = get_session_id()
576     if USE_MONGODB and db.is_connected():
577         db.clear_chat_history(session_id)
578     session.pop('conversation', None)
579     return jsonify({
580         'status': 'success',
581         'message': 'Conversation reset'
582     })
583
584
585 # -----
586 # Application Entry Point
587 # -----
588 if __name__ == '__main__':
589     print("=" * 56)
590     print(" Covexis AI FinTech Platform by Perfionix AI")
591     print(" Intelligent Financial Technology Assistant")
592     print("=" * 56)
593     print(f" GPT Server : {GPT_SERVER_URL}")
594     print(f" GPT Model : {GPT_MODEL}")
595     print(f" Lite Model : {LITE_MODEL}")
596     print(f" MongoDB : {'Connected' if USE_MONGODB else 'In-memory'}")
597     print("=" * 56)
598     app.run(
599         debug=os.getenv('FLASK_DEBUG', 'True') == 'True',
600         host='0.0.0.0',

```

```
601     port=8000
602 )
```

## 2.2 Architectural Components Overview

The table below summarizes the principal functional modules implemented within the application codebase and their respective responsibilities within the Covexis AI platform architecture.

Module	Functional Responsibility
<b>AI Chat Engine</b>	Dual-model inference pipeline supporting GPT-class models via Ollama and the lightweight Gemma 3 1B variant, with automatic failover routing and structured error propagation.
<b>RAG Pipeline</b>	End-to-end document processing encompassing file ingestion, multi-format text extraction (PDF, DOCX, TXT), semantic chunking with configurable overlap, relevance-scored retrieval, and context-injected response synthesis.
<b>VizIQ Analytics</b>	Automated data pipeline accepting CSV, Excel, and JSON inputs with statistical profiling, type inference, KPI derivation, and seven distinct visualization configurations optimized for financial datasets.
<b>Web Search</b>	Cascading search architecture with Serper.dev as primary provider, DuckDuckGo HTML extraction as secondary fallback, and googlesearch-python as tertiary option ensuring search availability.
<b>Session Layer</b>	Cookie-backed Flask sessions with MongoDB write-through for persistent storage of conversation threads, uploaded documents, and analytical outputs across user sessions.
<b>Security</b>	Input sanitization via Werkzeug, configurable file type restrictions, PII isolation from external model providers, and on-premise deployment pathway via Ollama for data residency compliance.

*The complete operational deployment additionally requires companion modules: database.py implementing the MongoDB data access layer, and rag\_engine.py encapsulating the vector embedding and semantic retrieval logic. The frontend presentation layer comprises HTML templates with integrated JavaScript for real-time interaction rendering. The code presented above constitutes the central orchestration module coordinating all platform subsystems.*