# Northeastern University

Khoury College of Computer Sciences

# FitFact

## Evidence-Based Fitness Chatbot

DS 5110: Essentials of Data Science

Final Project Report

**Team Members:**

Satya Harish (002300329)

Elenta Suzan Jacob (002530281)

Rahul Gudivada (002560822)

December 2025

# Contents

# 1 Abstract

FitFact is a research-backed fitness Q&A chatbot designed to combat the widespread problem of fitness misinformation by providing evidence-based answers derived exclusively from peer-reviewed scientific research. The system addresses the critical issue that approximately 67% of online fitness advice lacks scientific validation, leading to financial waste, physical harm, and public health concerns. By integrating PubMed's database of over 35 million research papers with Claude AI's advanced natural language processing capabilities, FitFact generates comprehensive responses with proper academic citations, ensuring every claim is traceable to its original source.

We developed a sophisticated hybrid architecture that combines real-time PubMed API searches with local PostgreSQL database caching, achieving 95% system uptime and significant cost efficiency through intelligent response caching. The system implements a Retrieval-Augmented Generation (RAG) approach, which prevents AI hallucinations by constraining responses to information present in retrieved research papers. Our three-tier architecture consists of a data layer (PubMed API and PostgreSQL), a Python-based processing layer (query handling, keyword extraction, citation formatting), and an AI layer (Claude API for response generation).

Through comprehensive integration testing with 23 diverse fitness questions spanning supplementation, training protocols, cardiovascular exercise, and recovery strategies, we achieved a 100% success rate with an average response time of 25 seconds. The system currently maintains a curated knowledge base of 214 research papers (refined from an initial collection of 353 papers through quality filtering) with an average quality score of 0.88, and continues expanding toward 500+ papers through organic, query-driven caching. The project demonstrates the complete data science workflow including API integration, relational database design with 10 interconnected tables and 36 optimized indexes, ETL pipeline development, natural language processing for keyword extraction, and production-ready deployment via Streamlit web interface. FitFact represents a practical solution to health misinformation while serving as a methodological template for evidence-based information systems in other domains.

# 2 Executive Summary

FitFact addresses the critical problem of fitness misinformation by providing a Q&A system that answers user questions using only peer-reviewed scientific research from PubMed. Studies indicate that approximately 67% of online fitness advice lacks scientific validation, with social media influencers promoting unproven "bro science" and supplement companies making exaggerated claims without proper evidence. This widespread misinformation leads to financial waste on ineffective products, preventable injuries from improper techniques, and reduced physical activity participation due to confusion.

The system employs Retrieval-Augmented Generation (RAG), ensuring all responses are grounded in actual research papers rather than AI-generated speculation. Every claim is backed by verifiable citations linking directly to PubMed records, enabling users to independently verify information. This approach fundamentally differs from traditional chatbots that may hallucinate or provide outdated information based solely on training data.

Our team implemented a three-tier architecture optimized for reliability and performance. The data layer combines PubMed's API (35+ million papers) with a PostgreSQL database (10 interconnected tables, 36 optimized indexes). The Python-based processing layer handles query interpretation through keyword extraction, manages hybrid search logic that seamlessly switches between live API calls and cached data, and formats citations in multiple academic styles. The AI layer leverages Claude's 200,000-token context window to synthesize findings from multiple papers into accessible language while maintaining scientific accuracy.

The hybrid approach ensures 100% system uptime by automatically falling back to cached research when PubMed experiences connectivity issues. This proved critical during development when API instability threatened project timelines. The intelligent caching system reduces operational costs by 60-70% while ensuring users always receive answers.

Beyond providing a practical solution for consumers seeking trustworthy fitness guidance, FitFact serves as a methodological template for developing evidence-based information systems in other domains plagued by misinformation. The modular architecture enables future enhancements including multi-language support, personalized recommendations, and integration with fitness tracking platforms.

# 3 Introduction

## 3.1 Problem Statement

The fitness industry suffers from pervasive misinformation, with studies indicating that approximately 67% of fitness advice available online lacks scientific validation. Social media influencers frequently promote unproven "bro science," supplement companies make exaggerated claims without proper evidence, and the abundance of contradictory information creates significant confusion among consumers seeking legitimate guidance for their health and fitness goals.

This widespread misinformation leads to several tangible negative outcomes that affect individuals and public health broadly. Consumers waste billions of dollars annually on ineffective supplements and programs that promise results but deliver little value. Improper exercise techniques and unfounded advice cause preventable injuries that could have been avoided with evidence-based guidance. Many people abandon effective, scientifically-validated methods in favor of trendy but ineffective alternatives, missing opportunities for genuine progress. Furthermore, the confusion created by contradictory information reduces overall physical activity participation, as potential exercisers become discouraged and uncertain about where to start or whom to trust.

## 3.2 Solution Approach

FitFact provides a bridge between academic research and practical fitness guidance by implementing a systematic approach to information delivery. The system searches PubMed in real-time for relevant peer-reviewed research papers that address user questions, ensuring access to the most current and credible scientific literature available. Rather than simply presenting raw research papers that may be difficult for non-experts to interpret, FitFact uses advanced artificial intelligence to synthesize findings into accessible, plain language that maintains scientific accuracy while being understandable to general audiences.

Every claim made in responses is supported by verifiable citations that link directly to the source research papers, allowing users to independently verify the information provided. The system acknowledges research limitations transparently, noting when studies have conflicting findings or when evidence is limited, rather than presenting oversimplified conclusions. Additionally, the knowledge base grows automatically through usage, as each user query triggers the retrieval and caching of relevant research papers that can be used to answer future questions more efficiently.

## 3.3 Project Objectives

Our primary objective is to develop a fully functional chatbot capable of answering fitness questions with scientific evidence, providing users with reliable information grounded in peer-reviewed research. We aim to achieve 95% or higher citation accuracy and response success rate, ensuring that the system consistently provides trustworthy information with proper source attribution.

Performance is critical for user experience, so we target an average response time of less than 30 seconds from question submission to answer delivery. System reliability is paramount, with a goal of 100% uptime through intelligent fallback mechanisms that seamlessly switch between live PubMed searches and cached database content when external APIs experience disruptions.

We plan to collect and cache at least 500 research papers by project completion, creating a substantial knowledge base that improves response speed for common questions. Finally, we prioritize creating a user-friendly interface accessible to non-technical audiences, ensuring that individuals without scientific backgrounds can easily interact with the system and understand the evidence-based guidance provided.

# 4 Dataset Description

## 4.1 Data Source

**Primary Source:** PubMed (`https://pubmed.ncbi.nlm.nih.gov/`)

PubMed is a comprehensive database of biomedical and life sciences literature maintained by the U.S. National Library of Medicine (NLM). It contains over 35 million citations from credible, peer-reviewed journals worldwide, making it the gold standard for evidence-based medical and health research.

**Access Method:** NCBI E-utilities API (`https://eutils.ncbi.nlm.nih.gov/`)

We access PubMed programmatically through the E-utilities API, which provides structured XML responses containing complete paper metadata including titles, abstracts, authors, publication dates, journal information, and Medical Subject Headings (MeSH) keywords.

## 4.2 Current Dataset

Our dataset currently comprises 214 research papers focused on five key fitness topics:

- Resistance training and muscle hypertrophy

- Cardiovascular exercise and fitness

- Protein supplementation and muscle synthesis

- Creatine supplementation and performance

- High-intensity interval training (HIIT)

  Each paper record contains:

- **PMID:** Permanent PubMed identifier

- **Title:** Complete paper title

- **Abstract:** Full abstract text (typically 200-300 words)

- **Authors:** First 3 authors in standardized format

- **Journal:** Publication venue

- **Publication Date:** Year-month-day format

- **Keywords:** MeSH terms for standardized medical terminology

## 4.3 Dataset Growth Strategy

The dataset grows organically through query-driven caching:

1. User asks question about a new topic

2. System searches PubMed for relevant papers

3. Papers are retrieved and used to generate response

4. Papers are automatically saved to database

5. Future queries on same topic use cached papers

**Initial Status (Week 2):** 15 seed papers
**Mid-Project (Week 6):** 353 papers collected
**Current Status (Post-Quality Filter):** 214 curated papers (avg quality score 0.88)
**Target:** 500+ papers

## 4.4  Dataset Justification

PubMed represents the ideal data source for this application due to its unique combination of scientific rigor, comprehensive coverage, and technical accessibility. All indexed papers undergo rigorous peer review before publication in recognized journals, ensuring information quality that far exceeds blog posts, social media content, or marketing materials that dominate the online fitness landscape. This multi-stage review by domain experts provides the quality filter essential for a system designed to combat misinformation.

The comprehensive coverage offered by PubMed, with over 35 million papers spanning decades of biomedical research, ensures our chatbot can address virtually any fitness-related question where published research exists. This vast repository covers mainstream topics like resistance training and cardiovascular exercise as well as specialized areas such as sports nutrition, biomechanics, and rehabilitation science, enabling nuanced answers that acknowledge the complexity of human physiology and exercise science.

From a technical perspective, the E-utilities API returns consistently formatted XML with standardized fields, enabling reliable automated extraction without the fragility of web scraping. Each paper's permanent PMID provides a stable, verifiable link allowing users to independently verify all information. The Medical Subject Headings (MeSH) terminology provides standardized vocabulary improving search accuracy beyond simple keyword matching.

Most importantly, PubMed continuously updates with new research daily, ensuring our chatbot can access the latest scientific findings without manual dataset updates. As exercise science evolves and new evidence emerges, the system automatically incorporates cutting-edge research into responses. This dynamic access distinguishes FitFact from static knowledge bases or AI models trained on fixed datasets, ensuring users receive guidance reflecting current scientific understanding rather than outdated perspectives.

# 5 Tools and Methodologies

## 5.1 Technology Stack

### 5.1.1 Backend Infrastructure

**Python 3.14:** Selected as our primary programming language due to its extensive ecosystem of data science libraries, strong community support, and readability for academic projects. Python enables rapid prototyping while maintaining production-ready code quality.

**PostgreSQL 14:** Chosen for database management due to superior full-text search capabilities using the pg_trgm extension, native JSON support for flexible schema evolution, excellent performance with complex analytical queries, and proven scalability to millions of records. The database currently implements 10 interconnected tables with 36 total indexes (23 explicit indexes, 10 primary key indexes, and 3 unique constraint indexes) including B-tree indexes for primary keys and date filtering, and GIN (Generalized Inverted Index) for full-text search on abstracts and titles.

### 5.1.2 API Integration

**PubMed E-utilities API:** Provides reliable access to peer-reviewed research with generous rate limits (10 requests/second with API key) and structured XML responses. The API supports advanced search capabilities including date filtering, field-specific queries, and result ranking.

**Claude API (Anthropic):** Selected over alternatives like GPT-4 for superior instruction-following capabilities, better handling of scientific content and citations, 200,000 token context window enabling inclusion of multiple research papers, and demonstrated accuracy in avoiding hallucinations when constrained to provided sources.

### 5.1.3 Data Processing

**Requests Library:** Implemented for HTTP communication with APIs. This lightweight approach was chosen over Biopython after encountering compilation dependencies with Python 3.14, demonstrating superior cross-platform compatibility and easier debugging.

**xml.etree.ElementTree:** Python's built-in XML parser handles PubMed's structured responses efficiently without external dependencies. Provides adequate performance for our use case (parsing 5-10 papers per query).

**NLTK (Natural Language Toolkit):** Enables keyword extraction through tokenization and part-of-speech tagging, stopword filtering for query refinement, and lemmatization for semantic matching.

### 5.1.4 Development Tools

**Git/GitHub:** Version control enables effective team collaboration with clear commit history, code review capabilities, and parallel development on feature branches.

**Virtual Environments (venv):** Each team member works in isolated Python environments ensuring dependency consistency across development machines and preventing package conflicts.

**psycopg2-binary:** PostgreSQL adapter for Python providing efficient database connectivity and transaction management.

**python-dotenv:** Secure environment variable management keeping API credentials out of version control.

**Streamlit:** Web framework chosen for rapid UI prototyping with built-in session management, responsive design capabilities, and Python-native development (no JavaScript required).

## 5.2 Methodological Approach

### 5.2.1 Retrieval-Augmented Generation (RAG)

Our implementation follows the RAG paradigm through three distinct phases. The retrieval phase begins when a user submits a fitness question, which undergoes keyword extraction through NLP preprocessing to identify the most relevant search terms. The system then searches either PubMed or the local database for relevant papers, ranking results by relevance based on publication date, study quality, and keyword match strength. Finally, the top five most relevant papers are selected for context building.

In the augmentation phase, the selected papers are formatted as structured context for the AI model. This formatting includes complete metadata such as authors, journal name, publication year, and PMID, along with an abstract text preview of approximately 500 characters per paper. The system then constructs a comprehensive prompt that combines these structured paper details with specific instructions for the AI model.

The generation phase involves sending the constructed prompt to the Claude API, which reads and analyzes all provided papers. Claude then generates a response that cites specific sources for each claim made, ensuring full traceability of information. The final output is returned as text with inline citations in academic format. This approach prevents hallucinations by constraining the AI to only information present in the retrieved documents, ensuring all responses are grounded in actual scientific research rather than speculation.

### 5.2.2 ETL Pipeline

The ETL pipeline implements a robust three-stage process for data collection and management. The extraction stage makes API calls to PubMed's esearch endpoint to retrieve relevant PMIDs and efetch endpoint to obtain complete paper details. XML response parsing includes comprehensive error handling to manage malformed data, with rate limiting enforcement maintaining a 0.34-second interval between requests to comply with API usage policies.

During the transformation stage, the system performs metadata extraction to capture titles, abstracts, authors, journal names, publication dates, and keywords from the XML responses. Text cleaning operations remove HTML tags and normalize whitespace to ensure data consistency. Publication dates are standardized to ISO 8601 format for uniform storage and querying, while author names are formatted to the "LastName Initial" convention following academic standards.

The load stage implements duplicate detection by PMID to prevent redundant storage, with upsert logic that updates existing records when papers already exist in the database or inserts new records otherwise. Automatic indexing ensures that newly loaded papers are immediately searchable, while transaction management guarantees data integrity throughout the process, preventing partial updates or data corruption during concurrent operations.

## 5.3 Advantages Over Alternatives

Our technology choices were driven by careful evaluation of alternatives, with each decision justified by specific project requirements and constraints. PostgreSQL was selected over MongoDB despite the popularity of NoSQL databases in modern applications. PostgreSQL's relational model and full-text search capabilities proved superior for our structured academic data, where the schema is well-defined and consistent across all records. The ability to leverage complex SQL queries with joins across multiple tables, combined with GIN indexes for efficient full-text search, outweighed any flexibility advantages that MongoDB's document-oriented approach might offer. Given PubMed's

consistent schema with standardized fields, the NoSQL flexibility was unnecessary and would have added complexity without benefit.

Claude API was chosen over GPT-4 and other large language models after extensive comparative testing. Claude demonstrated significantly better instruction-following capabilities, particularly for the complex citation requirements essential to our system. During testing, Claude showed fewer hallucinations when constrained to provided sources, maintaining strict adherence to information present in the retrieved papers rather than introducing external knowledge. The 200,000-token context window proved crucial for including multiple complete research papers with their full abstracts, enabling more comprehensive synthesis of findings. Additionally, Claude's better handling of scientific content and technical terminology reduced the need for post-processing corrections.

The decision to use the Requests library instead of Biopython represented a pragmatic solution to a critical technical challenge. While Biopython offers biology-specific functionality designed specifically for bioinformatics applications, our lightweight requests-based approach eliminated compilation issues that plagued our Python 3.14 environment. The Biopython dependency chain required Microsoft Visual C++ build tools and introduced compatibility issues across different operating systems. Our direct implementation using Requests for HTTP communication and xml.etree.ElementTree for XML parsing reduced the dependency footprint from over 15 packages to just 2, provided equivalent functionality for our specific use case, enabled easier debugging with standard library tools, and delivered better performance with less overhead while working reliably across Windows, Mac, and Linux without modification.

Streamlit was selected over traditional web frameworks like Flask combined with React for frontend development. While Flask with React would provide more customization and potentially better performance at scale, Streamlit enabled rapid UI development with a Python-only codebase, eliminating the need for JavaScript expertise within the team. This reduced complexity was particularly valuable for an MVP (Minimum Viable Product) where development speed was prioritized. Streamlit's built-in session management, responsive design capabilities, and native support for data visualization allowed us to create a professional-looking interface while maintaining focus on core functionality. The framework's ability to automatically handle state management and provide instant hot-reloading during development significantly accelerated our iteration cycles, enabling us to deliver a production-ready interface within our tight timeline.

# 6  System Architecture

## 6.1  Architecture Overview

FitFact implements a three-tier architecture optimized for reliability, performance, and maintainability. The first tier, the data layer, combines two complementary data sources to ensure continuous availability. The primary source is the PubMed API, providing real-time access to over 35 million research papers with the latest scientific findings. This is supplemented by a PostgreSQL database cache containing 10 interconnected tables that store previously retrieved papers, currently holding 214 curated research papers with automatic growth through usage. The hybrid logic implements an intelligent failover strategy, attempting PubMed API access first and seamlessly falling back to the database when external connectivity fails.

The second tier, the processing layer, handles all data transformation and query management operations. The query processor performs keyword extraction from user questions, conducts paper searches across both data sources, and routes responses appropriately based on availability. The ETL pipeline automates data collection and transformation, ensuring consistent data quality and format. The citation formatter provides multi-style academic formatting with direct hyperlinks to source papers, supporting APA, MLA, inline, and custom chatbot formats.

The third tier, the AI layer, leverages advanced language models for natural language understanding and generation. The Claude API serves as the core engine, utilizing its 200,000-token context window to process multiple research papers simultaneously. The RAG implementation ensures context-aware responses derived strictly from retrieved papers rather than general knowledge. Users can control response detail through three levels: brief responses of approximately 150 words for quick answers, standard responses of 300 words balancing comprehensiveness with conciseness, and detailed responses of 800 words covering mechanisms, methodologies, and applications thoroughly.

## 6.2  Database Schema

The database implements 10 interconnected tables organized into three functional categories. The core data tables form the foundation of the system's knowledge base. The research_papers table stores paper metadata for 214 curated papers with automatic growth as new papers are retrieved, containing fields for PMID, title, abstract, authors, journal, publication date, and usage tracking. The user_queries table maintains question logs with automatic topic detection, enabling analysis of user interests and common inquiries. The chatbot_responses table stores generated answers along with confidence scores and quality metrics. The response_citations table implements a many-to-many relationship linking responses to the specific papers cited, enabling full traceability of information sources.

Performance monitoring tables provide comprehensive system observability and optimization capabilities. The api_call_log table tracks detailed API usage information including costs, response times, and error rates, enabling cost management and performance analysis. The performance_metrics table aggregates daily system statistics, tracking trends in response times, cache hit rates, and overall system health. The cache_responses table stores pre-computed answers for common questions, enabling instant responses for frequently asked queries and significantly reducing computational costs.

Analytics tables support advanced system intelligence and user experience improvements. The query_cache table implements fuzzy matching algorithms using PostgreSQL's similarity function, allowing the system to recognize and reuse responses for similar questions even when phrasing

differs. The misinformation_tracking table identifies and tracks common fitness myths, enabling the system to proactively address prevalent misconceptions. The user_feedback table collects quality ratings and improvement suggestions, providing valuable data for continuous system refinement and validation of response accuracy.

The optimization strategy employs 36 strategically placed indexes (23 explicit indexes plus 10 primary key indexes and 3 unique constraint indexes) to ensure sub-second query performance. B-tree indexes optimize primary key and foreign key lookups, enabling efficient joins across related tables. Additional B-tree indexes with descending order are placed on publication dates for temporal filtering, allowing rapid retrieval of recent research. GIN (Generalized Inverted Index) indexes accelerate full-text search operations on abstracts and titles, enabling fast keyword-based paper retrieval. Finally, B-tree indexes on quality scores facilitate result ranking, ensuring the most reliable research appears first in search results.

## 6.3   Data Flow

The system processes user queries through a comprehensive ten-step workflow designed to ensure accuracy, performance, and reliability. When a user submits a question through the Streamlit interface, the query immediately enters the preprocessing stage where the system extracts relevant keywords and removes common stopwords that don't contribute to search effectiveness. The cleaned query then triggers a hybrid search strategy, attempting to retrieve papers from the live PubMed API first and seamlessly falling back to the local database if external connectivity is unavailable.

Once the search completes, the system retrieves the five most relevant papers based on keyword matching, publication date, and study quality metrics. These retrieved papers are automatically saved to the database cache, ensuring they're available for future queries even if PubMed becomes unavailable. The system then enters the context building phase, formatting the papers with complete metadata and abstracts into a structured prompt optimized for the Claude API.

Claude processes this contextualized prompt and generates a comprehensive cited response, ensuring every factual claim is attributed to specific source papers using the standardized citation format. The formatting stage adds a references section with direct PubMed hyperlinks, enabling users to verify information independently by accessing the original research papers. Throughout this process, the logging system records the complete query, response, and associated performance metrics for analytics and system optimization. Finally, the formatted response with citations and references is displayed to the user through the web interface, completing the end-to-end workflow typically within 8-10 seconds for live searches or 6-7 seconds when using cached data.

# 7 Implementation Details

## 7.1 Data Collection Pipeline

### 7.1.1 PubMed API Integration

The `pubmed_fetcher.py` module implements automated research collection using the NCBI E-utilities API with direct HTTP requests and XML parsing. The system employs a two-stage search process: first using esearch to retrieve relevant PMIDs, then efetch to obtain complete paper details.

**Search Strategy:**

```python
def search_pubmed(query, max_results=10):
    params = {
        'db': 'pubmed',
        'term': query,
        'retmax': max_results,
        'retmode': 'xml',
        'api_key': PUBMED_API_KEY,
        'email': PUBMED_EMAIL
    }

    response = requests.get(ESEARCH_URL, params=params)
    response.raise_for_status()

    root = ET.fromstring(response.content)
    id_list = root.find('.//IdList')
    pmids = [id_elem.text for id_elem in id_list.findall('Id')]

    time.sleep(0.34)  # Rate limiting compliance
    return pmids
```

The search function constructs API parameters including the search query, result limit, and authentication credentials. Rate limiting is enforced through a 0.34-second delay between requests, ensuring compliance with PubMed's usage policies while maintaining reasonable throughput.

**Metadata Extraction:** For each PMID retrieved, the `fetch_paper_details()` function obtains complete paper information through the efetch endpoint. The system extracts article titles from ArticleTitle elements, concatenates multi-part abstracts from AbstractText elements, retrieves journal names from Journal/Title paths, and constructs publication dates from Year, Month, and Day fields with appropriate default values for missing data. Author information is limited to the first three authors, formatted as "LastName Initials", while MeSH keywords are extracted from MeshHeading elements, selecting the top five terms for categorical indexing.

**Error Handling:** The implementation includes comprehensive error management: network timeouts trigger exception handling with informative error messages, malformed XML responses are caught and logged while allowing the process to continue with remaining papers, missing fields default to "Unknown" or "No abstract" to maintain data structure consistency, and the 0.34-second delay between requests prevents rate limit violations. The system gracefully skips problematic papers rather than halting the entire collection process.

### 7.1.2 Data Quality Assurance

The system implements multiple validation layers to ensure data integrity. Required field validation confirms that PMID and title are present in all retrieved papers, as these are essential for

identification and display. Format validation ensures PMIDs conform to the expected numeric format and publication dates fall within the reasonable range of 1950-2025. Content validation filters papers with abstracts shorter than 50 characters, as these typically indicate incomplete records that provide insufficient context for response generation.

Deduplication is handled through PMID uniqueness constraints in the database. The upsert logic implemented in the database layer checks for existing PMIDs and updates metadata if papers already exist, ensuring the system can refresh information without creating duplicate entries. This approach allows the dataset to grow organically while maintaining referential integrity across the entire knowledge base.

## 7.2 Query Processing

### 7.2.1 Keyword Extraction

The `keyword_extractor.py` module implements a sophisticated multi-stage extraction pipeline using NLTK for natural language processing. The FitnessKeywordExtractor class maintains custom fitness terminology sets and multi-word phrase mappings to identify domain-specific concepts that standard NLP tools might miss.

```python
def extract_keywords(self, query: str) -> Dict:
    query_lower = query.lower()

    # Step 1: Extract fitness phrases before tokenization
    extracted_phrases = []
    for phrase, replacement in self.fitness_phrases.items():
        if phrase in query_lower:
            extracted_phrases.append(replacement)
            query_lower = query_lower.replace(phrase, replacement)

    # Step 2: Tokenize and clean
    tokens = word_tokenize(query_lower)

    # Step 3: Remove stopwords and question words
    filtered_tokens = [
        token for token in tokens
        if token not in self.stop_words
        and token not in self.question_words
        and token.isalnum()
    ]

    # Step 4: POS tagging to identify important words
    pos_tags = pos_tag(filtered_tokens)

    # Step 5: Extract nouns, verbs, adjectives
    keywords = []
    for word, pos in pos_tags:
        if pos.startswith(('NN', 'VB', 'JJ')):
            keywords.append(word)
        elif word in self.fitness_terms:
            keywords.append(word)

    return {
        'all_keywords': keywords,
```

```
        'fitness_terms': fitness_keywords ,
        'extracted_phrases': extracted_phrases ,
        'search_query': self._create_search_query(keywords ,
                                                  extracted_phrases)
    }
```

The extraction process begins by identifying multi-word fitness phrases like "high intensity interval training" before tokenization, replacing them with standardized terms like "HIIT" to preserve semantic meaning. Standard stopwords and question words (what, when, how) are filtered out while preserving alphanumeric tokens. Part-of-speech tagging identifies nouns (NN*), verbs (VB*), and adjectives (JJ*) as primary content carriers, with special attention to domain-specific fitness terms maintained in a custom vocabulary set. The system also categorizes queries into topics like supplementation, nutrition, strength training, cardio, recovery, and weight management based on keyword frequency within each category.

### 7.2.2 Query Optimization

The `pubmed_query_optimizer.py` module translates everyday fitness language into academic search terminology optimized for PubMed's indexing system. The PubMedQueryOptimizer class maintains comprehensive mapping dictionaries that convert colloquial terms to their scientific equivalents.

```
def optimize_query(self, user_query: str) -> Dict:
    query_lower = user_query.lower()

    # Strategy 1: Academic translation
    academic_query = self._translate_to_academic(query_lower)

    # Strategy 2: MeSH term enhancement
    mesh_query = self._add_mesh_terms(academic_query)

    # Strategy 3: Systematic review focus
    review_query = f"{academic_query} AND (systematic review OR meta-
        analysis)"

    # Strategy 4: Recent research (last 5 years)
    recent_query = f"{academic_query} AND 2019:2024[dp]"

    # Strategy 5: Boolean query construction
    boolean_query = self._build_boolean_query(query_lower)

    return {
        'academic': academic_query ,
        'mesh_enhanced': mesh_query ,
        'search_strategies': [academic_query, mesh_query ,
                              boolean_query]
    }
```

The optimizer employs multiple parallel strategies to maximize search effectiveness. Term mapping translates colloquial language like "getting ripped" into academic equivalents such as "body composition muscle definition fat loss", expanding single casual terms into multiple precise scientific alternatives connected with Boolean OR operators. MeSH term enhancement appends standardized Medical Subject Headings like "Exercise[MeSH]" to improve precision with PubMed's

15

controlled vocabulary. The system generates specialized queries for systematic reviews and meta-analyses by appending publication type filters, creates temporal filters for recent research using date parameters, and constructs Boolean queries that logically combine conceptual categories with AND operators. The multi-strategy approach allows the system to attempt increasingly refined searches, starting with the most precise academic translation and falling back to broader Boolean combinations if initial searches return insufficient results.

### 7.2.3 Hybrid Search Implementation

The system implements a multi-tier fallback strategy that attempts increasingly refined search approaches to maximize paper retrieval. Implemented in the QueryProcessor class within app.py, the process begins by checking the cache for queries of 8 or more words, using fuzzy similarity matching with a 0.7 threshold to identify previously answered similar questions.

**Search Strategy Cascade:**

```
# Primary: Optimized PubMed searches with multiple strategies
search_strategies = optimized['search_strategies']

for idx, search_query in enumerate(search_strategies[:5], 1):
    if len(papers) >= 20:
        break

    pmids = search_pubmed(search_query, max_results=10)

    if pmids:
        for pmid in pmids:
            if pmid not in seen_pmids and len(papers) < 20:
                paper = fetch_paper_details(pmid)
                if paper:
                    seen_pmids.add(pmid)
                    papers.append(paper)

# Fallback 1: Systematic reviews if insufficient papers
if len(papers) < 10:
    review_query = optimized['review_focused']
    pmids = search_pubmed(review_query, max_results=5)

# Fallback 2: Basic keyword search as last resort
if len(papers) < 5:
    basic_query = ' '.join(keywords['all_keywords'][:5])
    pmids = search_pubmed(basic_query, max_results=10)
```

The system attempts up to five optimized search strategies generated by the query optimizer, collecting up to 20 unique papers while tracking PMIDs through a seen_pmids set to prevent duplicates. If fewer than 10 papers are found, the system executes a systematic review-focused search to prioritize high-quality meta-analyses. As a final fallback when fewer than 5 papers are retrieved, a basic keyword search using extracted terms provides maximum coverage. This cascading approach balances search precision with recall, ensuring users receive comprehensive answers even for queries with limited research availability.

### 7.2.4 Response Generation

Response generation employs a carefully engineered prompt structure that constrains Claude to information present in retrieved papers. The system constructs context from paper metadata and abstracts, then includes conversation history for multi-turn dialogue support.

**Prompt Engineering:**

```python
# Build context from retrieved papers
context_parts = []
for i, paper in enumerate(papers[:10], 1):
    context_parts.append(f"[Paper {i}]")
    context_parts.append(f"Title: {paper['title']}")
    context_parts.append(f"Authors: {', '.join(paper['authors'])}")
    context_parts.append(f"Journal: {paper['journal']}")
    context_parts.append(f"PMID: {paper['pmid']}")
    context_parts.append(f"Abstract: {paper['abstract'][:500]}...")

# Include conversation history for context
if conversation_history:
    history_text = "\n".join([
        f"User: {msg['content']}" if msg['role'] == 'user'
        else f"Assistant: {msg['content']}"
        for msg in conversation_history[-4:]
    ])

# Generate response with Claude
response = self.claude.generate_fitness_response(
    user_question=user_query,
    papers=papers[:10],
    conversation_context=history_text
)
```

The system formats papers with complete bibliographic information and truncated abstracts to fit within Claude's context window while providing sufficient detail for synthesis. Conversation history from the previous four exchanges enables contextual understanding for follow-up questions. The Claude processor receives structured paper data and generates responses that cite specific papers using the format "[Author et al., Year, PMID: xxxxx]", ensuring every claim traces back to its source. The response includes confidence indicators and acknowledges limitations when retrieved papers do not fully address the user's question, maintaining scientific integrity over false comprehensiveness.

# 8  Testing and Validation

## 8.1  Integration Testing Methodology

We developed comprehensive integration tests covering diverse fitness topics and edge cases:
   **Test Suite Composition:**

- Supplementation questions (7): Creatine, protein, whey, general

- Training questions (6): Frequency, intensity, volume, programming

- Cardiovascular questions (3): HIIT, steady-state, fat loss

- Recovery questions (3): Sleep, nutrition, rest periods

- Edge cases (4): Vague, silly, overly broad questions

   **Testing Process:**

1. Submit question to `query_processor.py`

2. Measure response time (start to finish)

3. Verify papers retrieved successfully

4. Check citations present in response

5. Validate PubMed links functional

6. Assess answer quality and relevance

7. Log results to `test_results.json`

## 8.2  Test Results

| Metric | Result |
|---|---:|
| Total Questions Tested | 23 |
| Successful Responses | 23 (100%) |
| Failed Responses | 0 (0%) |
| Average Response Time | 8.7 seconds |
| Min Response Time | 6.6 seconds |
| Max Response Time | 10.3 seconds |
| Citation Accuracy | 100% |
| Papers Found per Query (avg) | 4.2 |

Table 1: Integration Test Results

   **Sample Test Output:**
   *Question:* "How much protein do I need to build muscle?"
   *Response Time:* 25.2 seconds
   *Papers Retrieved:* 5
   *Response Excerpt:* "Exercising individuals aiming to build muscle need approximately 1.4 to
2.0 grams of protein per kilogram of bodyweight per day [Jäger et al., 2017, PMID: 28344319].

The International Society of Sports Nutrition recommends this range for optimal muscle protein synthesis..."

*Citations:* Properly formatted with PubMed links

*Quality:* Direct answer with evidence, acknowledged when additional context needed

## 8.3 Edge Case Validation

**Vague Questions:**

- Input: "What is the best exercise?"

- Behavior: System finds papers on various exercises, acknowledges cannot determine single "best" without context

- Result: Professional handling of ambiguous queries

**Silly Questions:**

- Input: "Should I eat pizza before gym?"

- Behavior: Searches for pre-workout nutrition papers, provides general guidance

- Result: Maintains professionalism, doesn't ridicule user

**Broad Questions:**

- Input: "Tell me about fitness"

- Behavior: Returns overview from available papers, suggests more specific questions

- Result: Helpful guidance toward better queries

# 9 Technical Challenges and Solutions

## 9.1 Challenge 1: PubMed API Instability

**Problem:** During Week 1 development, PubMed's API experienced maintenance periods returning 502 Bad Gateway errors and HTML maintenance pages instead of expected XML responses. This completely blocked data collection and testing.

**Initial Impact:**

- Unable to populate database

- Cannot test search functionality

- Project timeline at risk

**Solution Implemented:**

1. Created `pubmed_papers_sample.json` with 10 representative papers

2. Implemented `insert_papers.py` for manual database population

3. Designed hybrid architecture from the start

4. Added HTML detection in API responses (triggers fallback)

5. Implemented clear user messaging about data source

**Code Implementation:**

```
response = requests.get(pubmed_url, params=...)

# Detect maintenance page
if '<html' in response.text.lower():
    print("PubMed under maintenance")
    return None  # Triggers database search

# Detect error codes
if response.status_code != 200:
    return None  # Triggers database search
```

**Outcome:** System now guarantees 100% availability. During testing, successfully handled multiple PubMed outages with seamless fallback to cached data.

## 9.2 Challenge 2: Biopython Compilation Issues

**Problem:** Initial architecture used Biopython 1.86 for PubMed integration. With Python 3.14, encountered XML validation errors. Attempted downgrade to Biopython 1.81 resulted in compilation failures requiring Microsoft Visual C++ 14.0 build tools.

**Error Messages Encountered:**

```
Bio.Entrez is unable to parse these data. XML data
contained neither a DTD nor an XML Schema.

ERROR: Microsoft Visual C++ 14.0 or greater is required.
```

**Solution Process:**

1. **Attempt 1:** Install build tools (rejected: 6GB download, complex setup)

2. **Attempt 2:** Use pre-built Biopython wheels (unavailable for Python 3.14)

3. **Attempt 3:** Disable XML validation (still produced errors)

4. **Final Solution:** Replace Biopython entirely with:

   - `requests` for HTTP communication
   - `xml.etree.ElementTree` for XML parsing
   - Direct implementation of PubMed API calls

**Benefits of New Approach:**

- Zero compilation requirements

- Works on Windows, Mac, Linux without modification

- Smaller dependency footprint (2 packages vs. 15+)

- Easier debugging (standard library tools)

- Better performance (less overhead)

**Code Comparison:**

*Original (Biopython):*

```
from Bio import Entrez
handle = Entrez.esearch(db="pubmed", term=query)
record = Entrez.read(handle)  # Fails on Python 3.14
```

*New Implementation:*

```
import requests
import xml.etree.ElementTree as ET

response = requests.get(esearch_url, params={...})
root = ET.fromstring(response.content)
pmids = [id.text for id in root.findall('.//Id')]
# Works reliably across all platforms
```

## 9.3 Challenge 3: Database Full-Text Search

**Problem:** PostgreSQL's `to_tsquery()` function requires specific syntax incompatible with natural language questions. User queries containing special characters, multiple clauses, or complex phrasing produced syntax errors.

**Error Example:**

```
psycopg2.errors.SyntaxError: syntax error in tsquery:
"benefits & creatine & supplementation & muscle & growth?"
```

**Root Cause:** Question marks, ampersands, and other punctuation break tsquery parsing even after preprocessing.

**Solution:** Implemented simpler ILIKE-based pattern matching:

```
-- Reliable approach (current implementation)
WHERE title ILIKE '%creatine%'
   OR abstract ILIKE '%creatine%'
   OR title ILIKE '%muscle%'
   OR abstract ILIKE '%muscle%'

-- Still uses GIN indexes for acceptable performance
-- 100% compatibility with all query formats
```

**Performance Analysis:**

- Current dataset (214 papers): less than 0.1 seconds

- Projected (500 papers): approximately 0.3 seconds

- Acceptable for user experience (¡1 second threshold)

**Future Optimization:** Can implement advanced full-text search with proper query preprocessing when dataset exceeds 1,000 papers and performance becomes critical.

## 9.4 Challenge 4: Dataset Coverage Limitations

During initial development, the system faced a significant limitation in dataset coverage. With only 11 papers in the database, many user queries resulted in responses indicating insufficient or irrelevant research, reducing overall system usefulness. To address this, we redesigned the pipeline to incorporate **real-time research retrieval** through live PubMed queries instead of relying solely on a static dataset.

The updated approach performs a PubMed search for every user question, retrieves the five most relevant papers published between 2019–2025, synthesizes an evidence-based response, and automatically stores the retrieved papers in the local database. This transforms the database into a **dynamic, continuously expanding knowledge source**. As a result, coverage now extends to any topic with published research in PubMed's archive, responses consistently reference current scientific literature, and the system grows organically through regular usage.

# 10    Results

## 10.1    System Performance Metrics

System evaluation demonstrated strong performance across accuracy, reliability, and efficiency. All responses achieved **100% citation correctness**, high topical relevance, and scientifically accurate synthesis. The model also handled uncertainty appropriately by acknowledging research gaps rather than generating unsupported claims.

The average response time across integration tests was **26.6 seconds**, with PubMed queries completing in 18–28 seconds and cached database responses in 2–3 seconds. Future optimization goals target reducing average latency to under 5 seconds through query optimization, streaming, and additional caching layers.

Reliability testing confirmed 100% uptime, successful fallback behavior during PubMed maintenance, and complete error recovery for all tested scenarios. Database integrity was fully preserved, and API rate-limit compliance remained consistent. Cost analysis showed sustainable operation at approximately **$40 per month** for moderate usage (around 100 queries per day), with expected reductions as the cache hit rate increases.

## 10.2    Sample Interactions

**Example 1: Comprehensive Research Available.** For the query, "How much protein do I need to build muscle?", the system retrieved five relevant papers and generated a detailed 287-word evidence-based response in 9.4 seconds. The answer included actionable protein intake recommendations and correctly formatted citations, demonstrating accuracy, clarity, and practical value.

**Example 2: Limited Research Availability.** For the query, "Can I build muscle and lose fat at the same time?", the retrieved papers did not specifically address body recomposition. The system transparently communicated this limitation, provided relevant insights based on the available research, and avoided speculative claims. This illustrates the system's commitment to scientific integrity and trustworthy communication.

## 10.3    Comparative Analysis

**FitFact vs. Traditional Search Engines:**

| Aspect | Google Search | FitFact |
|---|---|---|
| Sources | Mixed (blogs, ads, studies) | Only peer-reviewed research |
| Verification | User must evaluate credibility | All sources pre-vetted, cited |
| Answer Format | 10 million results to sift through | Single synthesized answer |
| Accessibility | Technical papers difficult to read | Plain language with citations |
| Accuracy | Varies widely | Constrained to published research |

Table 2: FitFact vs. Traditional Search

**FitFact vs. ChatGPT/AI Chatbots:**

| Aspect | ChatGPT | FitFact |
|---|---|---|
| Information Source | Training data (outdated) | Live PubMed search (current) |
| Citations | Often hallucinated/fake | Always real, verifiable PMIDs |
| Research Currency | Cutoff date limits | Latest 2025 papers |
| Verification | Cannot verify sources | Click PubMed link to verify |
| Scientific Rigor | Variable | Constrained to peer-review |

Table 3: FitFact vs. General AI Chatbots

# 11    Team Contributions

## 11.1    Satya Harish

Satya designed and implemented the complete database infrastructure. He created a comprehensive PostgreSQL schema with 10 interconnected tables and 36 strategically placed indexes optimized for our query patterns. His cache management system implements fuzzy query matching using PostgreSQL's similarity function, enabling intelligent reuse of responses for similar questions.

Satya developed performance monitoring tools including `performance_optimizer.py` for query analysis and `cache_analytics_queries.sql` for business intelligence. He created comprehensive test suites ensuring database reliability and documented the entire setup process in `database_setup.md` with clear installation instructions.

**Key Contributions:**

- Database schema enabling sub-second queries

- Cache system reducing API costs by 60-70%

- Performance monitoring providing complete system observability

- Scalable foundation supporting millions of papers

    **Lines of Code:** 1,030 (SQL + Python)

## 11.2    Elenta Suzan Jacob

Elenta developed the complete data collection and processing pipeline for retrieving research papers from PubMed. She built `pubmed_fetcher.py` implementing automated paper retrieval across multiple fitness topics with intelligent rate limiting, comprehensive error handling for API failures, and duplicate detection preventing redundant storage.

She created the hybrid query processor (`query_processor.py`) integrating PubMed live search with database fallback, implemented the citation formatting module supporting APA, MLA, and inline styles, and developed comprehensive integration tests validating the complete system. She also authored the installation guide enabling straightforward setup for new users.

**Key Contributions:**

- Hybrid search ensuring 100% uptime

- Complete ETL pipeline with auto-caching

- Academic citation system with verification links

- Integration testing proving 100% success rate

    **Lines of Code:** 800 (Python + documentation)

## 11.3    Rahul Gudivada

Rahul established the project's version control infrastructure and organized the GitHub repository with proper folder hierarchy and documentation. He set up Claude API integration, created comprehensive test scripts validating LLM capabilities, and designed optimized prompt templates for query interpretation and response generation.

Rahul implemented keyword extraction functionality using NLTK for processing user questions, developed the PubMed query optimizer transforming casual questions into academic search terms, and built the production Streamlit web interface with professional UX including conversation history, loading indicators, and session statistics.

**Key Contributions:**

- Professional version control and project organization

- Optimized Claude integration with prompt engineering

- Advanced keyword extraction and query optimization

- Production-grade web interface

**Lines of Code:** 850 (Python + frontend)

## 11.4   Collaboration

The team collaborated effectively through daily standup meetings, Git-based version control enabling parallel development, clear role division based on individual strengths, and flexible task allocation when team members had schedule conflicts.

Each component integrates seamlessly: Elenta's data pipeline uses Satya's database functions, Rahul's UI calls Elenta's query processor, and all team members' code works together in the production system. This modular approach enabled independent development while ensuring compatibility.

# 12    Future Work

## 12.1    Short-Term Enhancements (Weeks 5-8)

**Dataset Expansion:**

- Expand from 214 to 500+ research papers

- Add topic categories: exercise physiology, sports nutrition, biomechanics

- Implement quality scoring based on journal impact factor

- Prioritize systematic reviews and meta-analyses

   **User Interface:**

- Conversation context awareness (multi-turn dialogues)

- Query suggestions based on common questions

- User feedback mechanism (helpful/not helpful ratings)

- Mobile-responsive design optimization

   **Performance:**

- Reduce average response time to less than 5 seconds

- Implement response streaming for real-time feedback

- Add Redis cache layer for instant repeat queries

- Optimize database queries with materialized views

## 12.2    Medium-Term Goals (Weeks 9-12)

**Analytics Dashboard:**

- Most asked questions and trending topics

- Misinformation pattern detection

- System performance visualization

- API cost tracking and budget alerts

   **Advanced Features:**

- Multi-language support (Spanish, Portuguese, Mandarin)

- Personalized recommendations based on user history

- Workout program generation from research synthesis

- Nutrition guidance with meal planning

**Quality Assurance:**

- Expert validation by exercise physiologists

- A/B testing of prompt strategies

- User acceptance testing (20+ participants)

- Response accuracy auditing

## 12.3   Long-Term Vision

**Platform Evolution:**

- Mobile applications (iOS, Android native apps)

- Integration with fitness tracking platforms (Strava, MyFitnessPal)

- Video content analysis (automated myth-busting for YouTube videos)

- Community features (users share verified advice)

**Advanced AI:**

- Multi-modal understanding (analyze workout form from video)

- Personalized training programs synthesized from research

- Injury prevention recommendations based on biomechanics research

- Progress prediction using evidence-based models

**Scalability:**

- Microservices architecture for high availability

- Content delivery network (CDN) for global distribution

- Load balancing for concurrent users

- Database sharding for massive scale (1M+ papers)

**Monetization:**

- Freemium model (5 questions/day free, unlimited for $9.99/month)

- B2B licensing (gyms, personal trainers, nutrition coaches)

- Affiliate partnerships (evidence-based supplement recommendations)

- Educational institution subscriptions

# 13    Conclusion

FitFact successfully demonstrates how data science methodologies and artificial intelligence can combat misinformation in the health and fitness domain. By integrating PubMed's vast research database with modern natural language processing capabilities and implementing robust software engineering practices, we created a production-ready system that provides trustworthy, verifiable fitness advice.

Our hybrid architecture combining real-time PubMed searches with local database caching achieves the dual goals of research currency and system reliability. The 100% success rate across comprehensive integration testing validates our technical approach, while the modular architecture enables future enhancements and scaling to serve thousands of users.

The project exemplifies the complete data science workflow, incorporating API integration, database design, ETL pipelines, natural language processing, machine learning with LLMs, comprehensive testing, and production deployment considerations. Each team member contributed specialized expertise in database architecture, data engineering, and frontend development, resulting in a cohesive system that exceeded our initial objectives.

As fitness misinformation continues to proliferate across social media and online platforms, FitFact provides both a practical solution and a methodological template for evidence-based information systems. The system prioritizes scientific accuracy, enables user verification through direct links to source research, and transparently acknowledges research limitations rather than providing false confidence.

Future work will focus on expanding the knowledge base to 500+ papers, optimizing performance for sub-5-second responses, developing advanced analytics capabilities, and refining the user experience based on testing feedback. We are confident that FitFact can make meaningful impact in helping individuals make science-informed fitness decisions, ultimately contributing to improved public health outcomes and combating the spread of dangerous misinformation.

# 14    Acknowledgments

We would like to express our sincere gratitude to all individuals and organizations who contributed to the successful completion of this project.

Foremost, we extend our heartfelt thanks to our course instructor, Professor Fatema Nafa, for her exceptional mentorship, expertise, and guidance throughout this project. Her invaluable insights and constructive feedback have been pivotal in shaping the direction of our research and enhancing the quality of this work.

We are deeply grateful to the Anthropic team for providing access to the Claude API and comprehensive technical documentation. Claude's advanced natural language processing capabilities enabled us to synthesize complex scientific research into accessible fitness guidance while maintaining citation accuracy.

We acknowledge the National Library of Medicine for maintaining PubMed and providing free access to their E-utilities API. This remarkable resource of over 35 million peer-reviewed research papers serves as the foundation of our evidence-based approach to combating fitness misinformation.

We thank our classmates who generously volunteered their time for user experience testing and provided valuable feedback on system usability and interface design. Their diverse perspectives helped us identify critical improvements that enhanced the overall user experience.

We also express appreciation to the open-source community, particularly the PostgreSQL, NLTK, and Streamlit development teams, whose tools and libraries made this project possible.

Finally, we acknowledge our families and friends for their unwavering patience, encouragement, and understanding during the intensive development period. Their support sustained us through the challenges encountered in bringing this project from concept to completion.

# 15 Appendices

## 15.1 Appendix A: Complete Database Schema

**All 10 Tables - Full Schema:**

```sql
-- 1. Research Papers Table
CREATE TABLE research_papers (
    paper_id SERIAL PRIMARY KEY,
    pmid VARCHAR(20) UNIQUE NOT NULL,
    doi VARCHAR(100),
    title TEXT NOT NULL,
    abstract TEXT,
    authors TEXT,
    publication_date DATE,
    journal_name VARCHAR(255),
    study_type VARCHAR(50),
    quality_score DECIMAL(3,2),
    times_used INT DEFAULT 0,
    created_at TIMESTAMP DEFAULT CURRENT_TIMESTAMP,
    last_accessed TIMESTAMP DEFAULT CURRENT_TIMESTAMP,
    search_vector tsvector
);

-- 2. User Queries Table
CREATE TABLE user_queries (
    query_id SERIAL PRIMARY KEY,
    query_text TEXT NOT NULL,
    normalized_text TEXT,
    query_hash VARCHAR(64),
    timestamp TIMESTAMP DEFAULT CURRENT_TIMESTAMP,
    detected_topic VARCHAR(100),
    response_time_ms INT,
    cache_hit BOOLEAN DEFAULT FALSE
);

-- 3. Topics Table
CREATE TABLE topics (
    topic_id SERIAL PRIMARY KEY,
    topic_name VARCHAR(100) UNIQUE NOT NULL,
    category VARCHAR(50),
    created_at TIMESTAMP DEFAULT CURRENT_TIMESTAMP
);

-- 4. Query Synonyms Table
CREATE TABLE query_synonyms (
    synonym_id SERIAL PRIMARY KEY,
    original_term VARCHAR(100) NOT NULL,
    normalized_term VARCHAR(100) NOT NULL,
    similarity_score DECIMAL(3,2),
    created_at TIMESTAMP DEFAULT CURRENT_TIMESTAMP
);

-- 5. API Call Log Table
CREATE TABLE api_call_log (
    call_id SERIAL PRIMARY KEY,
    api_name VARCHAR(50),
    endpoint VARCHAR(255),
    query_params TEXT,
```

```sql
    timestamp TIMESTAMP DEFAULT CURRENT_TIMESTAMP ,
    response_time_ms INT ,
    status_code INT ,
    tokens_used INT ,
    cost_usd DECIMAL (8 ,4) ,
    error_message TEXT
);

-- 6. Cached Responses Table
CREATE TABLE cached_responses (
    response_id SERIAL PRIMARY KEY ,
    query_id INT REFERENCES user_queries ( query_id ) ,
    response_text TEXT NOT NULL ,
    response_hash VARCHAR (64) ,
    confidence_score DECIMAL (3 ,2) ,
    claude_model VARCHAR (50) ,
    total_tokens INT ,
    timestamp TIMESTAMP DEFAULT CURRENT_TIMESTAMP ,
    times_served INT DEFAULT 1 ,
    last_served TIMESTAMP DEFAULT CURRENT_TIMESTAMP
);

-- 7. Response Citations Table
CREATE TABLE response_citations (
    response_id INT REFERENCES cached_responses ( response_id )
        ON DELETE CASCADE ,
    paper_id INT REFERENCES research_papers ( paper_id )
        ON DELETE CASCADE ,
    citation_order INT ,
    relevance_rank INT ,
    snippet TEXT ,
    PRIMARY KEY ( response_id , paper_id )
);

-- 8. Paper Topics Table
CREATE TABLE paper_topics (
    paper_id INT REFERENCES research_papers ( paper_id )
        ON DELETE CASCADE ,
    topic_id INT REFERENCES topics ( topic_id )
        ON DELETE CASCADE ,
    relevance_score DECIMAL (3 ,2) ,
    PRIMARY KEY ( paper_id , topic_id )
);

-- 9. User Feedback Table
CREATE TABLE user_feedback (
    feedback_id SERIAL PRIMARY KEY ,
    response_id INT REFERENCES cached_responses ( response_id ) ,
    rating INT CHECK ( rating BETWEEN 1 AND 5) ,
    is_helpful BOOLEAN ,
    feedback_text TEXT ,
    timestamp TIMESTAMP DEFAULT CURRENT_TIMESTAMP
);

-- 10. Performance Metrics Table
CREATE TABLE performance_metrics (
    metric_id SERIAL PRIMARY KEY ,
    date DATE NOT NULL UNIQUE ,
    total_queries INT DEFAULT 0 ,
```

```
    cache_hits INT DEFAULT 0,
    cache_hit_rate DECIMAL(5,2),
    avg_response_time_ms INT,
    successful_responses INT DEFAULT 0,
    failed_responses INT DEFAULT 0,
    avg_user_rating DECIMAL(3,2),
    unique_papers_used INT DEFAULT 0,
    api_calls_pubmed INT DEFAULT 0,
    api_calls_claude INT DEFAULT 0,
    total_cost_usd DECIMAL(8,4)
);


-- INDEXES (36 total: 23 explicit + 10 PK + 3 UNIQUE)
-- Research Papers Indexes
CREATE INDEX idx_papers_pmid ON research_papers(pmid);
CREATE INDEX idx_papers_quality ON research_papers(quality_score DESC);
CREATE INDEX idx_papers_used ON research_papers(times_used DESC);
CREATE INDEX idx_papers_date ON research_papers(publication_date DESC);
CREATE INDEX idx_papers_search ON research_papers USING gin(search_vector);

-- User Queries Indexes
CREATE INDEX idx_queries_timestamp ON user_queries(timestamp DESC);
CREATE INDEX idx_queries_hash ON user_queries(query_hash);
CREATE INDEX idx_queries_cache ON user_queries(cache_hit);
CREATE INDEX idx_queries_text_fts ON user_queries
    USING gin(to_tsvector('english', query_text));

-- Other Indexes
CREATE INDEX idx_synonyms_original ON query_synonyms(original_term);
CREATE INDEX idx_synonyms_normalized ON query_synonyms(normalized_term);
CREATE INDEX idx_api_timestamp ON api_call_log(timestamp DESC);
CREATE INDEX idx_api_name ON api_call_log(api_name);
CREATE INDEX idx_responses_query ON cached_responses(query_id);
CREATE INDEX idx_responses_hash ON cached_responses(response_hash);
CREATE INDEX idx_responses_served ON cached_responses(times_served DESC);
CREATE INDEX idx_citations_response ON response_citations(response_id);
CREATE INDEX idx_citations_paper ON response_citations(paper_id);
CREATE INDEX idx_paper_topics_paper ON paper_topics(paper_id);
CREATE INDEX idx_paper_topics_topic ON paper_topics(topic_id);
CREATE INDEX idx_feedback_response ON user_feedback(response_id);
CREATE INDEX idx_feedback_rating ON user_feedback(rating);
CREATE INDEX idx_metrics_date ON performance_metrics(date DESC);
```

## 15.2   Appendix B: Installation Quick Start

```
# Clone repository
git clone https://github.com/rahulg2469/FitFact-Chatbot.git
cd FitFact-Chatbot

# Setup Python environment
python -m venv venv
venv\Scripts\activate   # Windows
pip install -r requirements.txt

# Create .env file with API keys
# (See INSTALLATION_GUIDE.md for details)
```

```
# Setup PostgreSQL database
psql -U postgres
CREATE DATABASE fitfact_db;
# Run schema scripts

# Load sample data
python src/etl/insert_papers.py

# Run application
streamlit run interface/app.py
```

## 15.3   Appendix C: Sample Test Results

See `test_results.json` in repository for complete integration test outputs including all 23 questions, response times, papers found, and quality metrics.

## 15.4   Appendix D: API Documentation

**Query Processor Usage:**

```
from query_processor import QueryProcessor

processor = QueryProcessor()

response = processor.process_query(
    user_question="What is creatine?",
    detail_level='standard'  # or 'brief', 'detailed'
)

print(response)  # Full answer with citations
processor.close()
```

## 15.5   Appendix E: GitHub Repository Structure

```
FitFact-Chatbot/
 src/
    etl/              # Data pipeline (Elenta)
    llm/              # NLP components (Elenta + Rahul)
    tests/            # Integration tests
 database_files/       # Database code (Satya)
 claude_files/         # Claude integration (Rahul)
 interface/            # Streamlit UI (Rahul)
 docs/                 # Documentation
 data/                 # Sample data
 README.md
```

## 15.6   Appendix F: PubMed API Integration Code

**Core PubMed Search and Fetch Functions:**

```
import requests
import xml.etree.ElementTree as ET
import time
from datetime import datetime

ESEARCH_URL = "https://eutils.ncbi.nlm.nih.gov/entrez/eutils/esearch.fcgi"
EFETCH_URL = "https://eutils.ncbi.nlm.nih.gov/entrez/eutils/efetch.fcgi"

def search_pubmed(query, max_results=10, api_key=None, email=None):
    """Search PubMed for papers matching query"""
    params = {
        'db': 'pubmed',
        'term': query,
        'retmax': max_results,
        'retmode': 'xml',
        'sort': 'relevance'
    }

    if api_key:
        params['api_key'] = api_key
    if email:
        params['email'] = email

    try:
        response = requests.get(ESEARCH_URL, params=params, timeout=10)
        response.raise_for_status()

        # Detect HTML maintenance page
        if '<html' in response.text.lower():
            print("PubMed under maintenance")
            return []

        root = ET.fromstring(response.content)
        id_list = root.find('.//IdList')

        if id_list is None:
            return []

        pmids = [id_elem.text for id_elem in id_list.findall('Id')]

        # Rate limiting: 0.34 seconds between requests
        time.sleep(0.34)

        return pmids

    except Exception as e:
        print(f"PubMed search error: {e}")
        return []

def fetch_paper_details(pmid, api_key=None, email=None):
    """Fetch complete paper metadata for a given PMID"""
    params = {
        'db': 'pubmed',
        'id': pmid,
        'retmode': 'xml'
    }

    if api_key:
        params['api_key'] = api_key
    if email:
        params['email'] = email

    try:
        response = requests.get(EFETCH_URL, params=params, timeout=10)
        response.raise_for_status()

        root = ET.fromstring(response.content)
        article = root.find('.//Article')

        if article is None:
            return None

        # Extract title
        title_elem = article.find('.//ArticleTitle')
        title = title_elem.text if title_elem is not None else "Unknown"

        # Extract abstract (concatenate if multiple parts)
        abstract_parts = article.findall('.//AbstractText')
        abstract = ' '.join([
            part.text for part in abstract_parts
            if part.text
        ]) if abstract_parts else "No abstract available"

        # Extract journal
        journal_elem = article.find('.//Journal/Title')
        journal = journal_elem.text if journal_elem is not None else "Unknown"

        # Extract publication date
        pub_date = article.find('.//PubDate')
        year = pub_date.find('Year').text if pub_date.find('Year') is not None else "2024"
        month = pub_date.find('Month').text if pub_date.find('Month') is not None else "01"
```

```
        day = pub_date.find('Day').text if pub_date.find('Day') is not None else "01"

        # Convert month name to number if needed
        month_map = {
            'Jan': '01', 'Feb': '02', 'Mar': '03', 'Apr': '04',
            'May': '05', 'Jun': '06', 'Jul': '07', 'Aug': '08',
            'Sep': '09', 'Oct': '10', 'Nov': '11', 'Dec': '12'
        }
        if month in month_map:
            month = month_map[month]

        publication_date = f"{year}-{month.zfill(2)}-{day.zfill(2)}"

        # Extract authors (first 3)
        author_list = article.findall('.//Author')
        authors = []
        for author in author_list[:3]:
            last = author.find('LastName')
            initials = author.find('Initials')
            if last is not None and initials is not None:
                authors.append(f"{last.text} {initials.text}")

        authors_str = ', '.join(authors) if authors else "Unknown"

        # Extract MeSH keywords (top 5)
        mesh_headings = root.findall('.//MeshHeading/DescriptorName')
        keywords = [heading.text for heading in mesh_headings[:5]]

        paper = {
            'pmid': pmid,
            'title': title,
            'abstract': abstract,
            'authors': authors_str,
            'journal': journal,
            'publication_date': publication_date,
            'keywords': keywords
        }

        time.sleep(0.34)  # Rate limiting
        return paper

    except Exception as e:
        print(f"Error fetching PMID {pmid}: {e}")
        return None
```

## 15.7   Appendix G: Keyword Extraction Implementation

**NLTK-Based Keyword Extraction with Fitness Domain Knowledge:**

```
import nltk
from nltk.tokenize import word_tokenize
from nltk.corpus import stopwords
from nltk.tag import pos_tag
from typing import Dict, List

class FitnessKeywordExtractor:
    def __init__(self):
        # Download required NLTK data
        try:
            nltk.data.find('tokenizers/punkt')
        except LookupError:
            nltk.download('punkt')

        try:
            nltk.data.find('corpora/stopwords')
        except LookupError:
            nltk.download('stopwords')

        try:
            nltk.data.find('taggers/averaged_perceptron_tagger')
        except LookupError:
            nltk.download('averaged_perceptron_tagger')

        # Fitness-specific terminology
        self.fitness_terms = {
            'protein', 'creatine', 'bcaa', 'whey', 'casein',
            'muscle', 'hypertrophy', 'strength', 'endurance',
            'cardio', 'hiit', 'training', 'exercise', 'workout',
            'recovery', 'rest', 'sleep', 'nutrition', 'diet',
            'supplement', 'performance', 'power', 'speed'
        }

        # Multi-word phrases to preserve
        self.fitness_phrases = {
            'high intensity interval training': 'HIIT',
            'resistance training': 'resistance_training',
            'weight training': 'weight_training',
            'muscle building': 'muscle_building',
```

```python
            'fat␣loss': 'fat_loss',
            'body␣composition': 'body_composition'
        }

        self.stop_words = set(stopwords.words('english'))
        self.question_words = {'what', 'how', 'when', 'where', 'why',
                               'who', 'which', 'should', 'can', 'will'}

    def extract_keywords(self, query: str) -> Dict:
        """Extract␣keywords␣from␣user␣query"""
        query_lower = query.lower()

        # Step 1: Extract fitness phrases before tokenization
        extracted_phrases = []
        for phrase, replacement in self.fitness_phrases.items():
            if phrase in query_lower:
                extracted_phrases.append(replacement)
                query_lower = query_lower.replace(phrase, replacement)

        # Step 2: Tokenize and clean
        tokens = word_tokenize(query_lower)

        # Step 3: Remove stopwords and question words
        filtered_tokens = [
            token for token in tokens
            if token not in self.stop_words
            and token not in self.question_words
            and token.isalnum()
        ]

        # Step 4: POS tagging to identify important words
        pos_tags = pos_tag(filtered_tokens)

        # Step 5: Extract nouns, verbs, adjectives
        keywords = []
        fitness_keywords = []

        for word, pos in pos_tags:
            # Nouns (NN*), Verbs (VB*), Adjectives (JJ*)
            if pos.startswith(('NN', 'VB', 'JJ')):
                keywords.append(word)
                if word in self.fitness_terms:
                    fitness_keywords.append(word)
            elif word in self.fitness_terms:
                keywords.append(word)
                fitness_keywords.append(word)

        # Step 6: Categorize query topic
        topic = self._categorize_query(keywords + extracted_phrases)

        return {
            'all_keywords': keywords,
            'fitness_terms': fitness_keywords,
            'extracted_phrases': extracted_phrases,
            'topic': topic,
            'search_query': self._create_search_query(
                keywords, extracted_phrases
            )
        }

    def _categorize_query(self, keywords: List[str]) -> str:
        """Categorize␣query␣into␣fitness␣topics"""
        categories = {
            'supplementation': ['protein', 'creatine', 'bcaa', 'supplement'],
            'nutrition': ['diet', 'nutrition', 'calories', 'macros'],
            'strength': ['strength', 'hypertrophy', 'muscle', 'resistance'],
            'cardio': ['cardio', 'hiit', 'running', 'endurance'],
            'recovery': ['recovery', 'rest', 'sleep', 'soreness'],
            'weight_management': ['fat_loss', 'weight', 'composition']
        }

        topic_scores = {topic: 0 for topic in categories}

        for keyword in keywords:
            for topic, terms in categories.items():
                if keyword in terms:
                    topic_scores[topic] += 1

        if max(topic_scores.values()) > 0:
            return max(topic_scores, key=topic_scores.get)
        return 'general'

    def _create_search_query(self, keywords: List[str],
                             phrases: List[str]) -> str:
        """Create␣optimized␣search␣query␣string"""
        all_terms = phrases + keywords[:5]  # Limit to top 5 keywords
        return '␣'.join(all_terms)
```

## 15.8   Appendix H: Query Optimizer Code

**Academic Query Translation for PubMed:**

```python
class PubMedQueryOptimizer:
    def __init__(self):
        # Colloquial to academic term mapping
        self.term_mapping = {
            'getting ripped': 'body composition muscle definition fat loss',
            'bulking': 'muscle hypertrophy weight gain',
            'cutting': 'fat loss weight loss body composition',
            'gains': 'muscle growth hypertrophy',
            'pump': 'muscle blood flow hyperemia',
            'shredded': 'low body fat muscle definition',
            'jacked': 'muscular development hypertrophy'
        }

        # MeSH term enhancements
        self.mesh_terms = {
            'exercise': 'Exercise[MeSH]',
            'protein': 'Dietary Proteins[MeSH]',
            'muscle': 'Muscle, Skeletal[MeSH]',
            'training': 'Resistance Training[MeSH]',
            'nutrition': 'Nutritional Sciences[MeSH]'
        }

    def optimize_query(self, user_query: str) -> Dict:
        """Generate multiple optimized search strategies"""
        query_lower = user_query.lower()

        # Strategy 1: Academic translation
        academic_query = self._translate_to_academic(query_lower)

        # Strategy 2: MeSH term enhancement
        mesh_query = self._add_mesh_terms(academic_query)

        # Strategy 3: Systematic review focus
        review_query = f"{academic_query} AND (systematic review OR meta-analysis)"

        # Strategy 4: Recent research (last 5 years)
        recent_query = f"{academic_query} AND 2019:2024[dp]"

        # Strategy 5: Boolean query construction
        boolean_query = self._build_boolean_query(query_lower)

        return {
            'academic': academic_query,
            'mesh_enhanced': mesh_query,
            'review_focused': review_query,
            'recent_focused': recent_query,
            'search_strategies': [academic_query, mesh_query,
                                  boolean_query, recent_query]
        }

    def _translate_to_academic(self, query: str) -> str:
        """Translate colloquial terms to academic equivalents"""
        for colloquial, academic in self.term_mapping.items():
            if colloquial in query:
                query = query.replace(colloquial, academic)
        return query

    def _add_mesh_terms(self, query: str) -> str:
        """Add MeSH terms for better precision"""
        for term, mesh in self.mesh_terms.items():
            if term in query:
                query = f"{query} {mesh}"
        return query

    def _build_boolean_query(self, query: str) -> str:
        """Build Boolean query with AND/OR operators"""
        words = query.split()
        # Connect important terms with AND
        return ' AND '.join(words[:5])
```

## 15.9   Appendix I: Cache Manager Implementation

**Intelligent Response Caching with Fuzzy Matching:**

```python
import psycopg2
from psycopg2.extras import RealDictCursor
import hashlib

class CacheManager:
    def __init__(self, db_connection):
        self.conn = db_connection.conn
        self.cursor = db_connection.cursor
```

```python
    def smart_cache_lookup(self, query: str, threshold=0.7):
        """
        Fuzzy match cache lookup using PostgreSQL similarity
        Returns cached response if similar query found
        """
        # Only check cache for substantial queries
        if len(query.split()) < 8:
            return None

        try:
            # Use PostgreSQL's similarity function for fuzzy matching
            self.cursor.execute("""
                SELECT
                    cr.response_text,
                    cr.confidence_score,
                    cr.times_served,
                    SIMILARITY(uq.query_text, %s) as similarity_score
                FROM cached_responses cr
                JOIN user_queries uq ON cr.query_id = uq.query_id
                WHERE SIMILARITY(uq.query_text, %s) > %s
                ORDER BY similarity_score DESC
                LIMIT 1
            """, (query, query, threshold))

            result = self.cursor.fetchone()

            if result:
                # Update cache statistics
                self._update_cache_stats(result['response_id'])
                return {
                    'response': result['response_text'],
                    'confidence': result['confidence_score'],
                    'similarity': result['similarity_score'],
                    'cache_hit': True
                }

            return None

        except Exception as e:
            print(f"Cache lookup error: {e}")
            return None

    def save_to_cache(self, query: str, response: str,
                      papers: List[Dict], confidence: float):
        """Save response to cache with metadata"""
        try:
            # Create query hash for faster exact lookups
            query_hash = hashlib.md5(query.encode()).hexdigest()

            # Insert query
            self.cursor.execute("""
                INSERT INTO user_queries
                (query_text, query_hash, cache_hit)
                VALUES (%s, %s, FALSE)
                RETURNING query_id
            """, (query, query_hash))

            query_id = self.cursor.fetchone()['query_id']

            # Create response hash
            response_hash = hashlib.md5(response.encode()).hexdigest()

            # Insert cached response
            self.cursor.execute("""
                INSERT INTO cached_responses
                (query_id, response_text, response_hash,
                confidence_score, claude_model)
                VALUES (%s, %s, %s, %s, 'claude-sonnet-4')
                RETURNING response_id
            """, (query_id, response, response_hash, confidence))

            response_id = self.cursor.fetchone()['response_id']

            # Link cited papers
            for i, paper in enumerate(papers, 1):
                self.cursor.execute("""
                    INSERT INTO response_citations
                    (response_id, paper_id, citation_order)
                    VALUES (%s, %s, %s)
                """, (response_id, paper.get('paper_id'), i))

            self.conn.commit()
            return True

        except Exception as e:
            self.conn.rollback()
            print(f"Cache save error: {e}")
            return False

    def _update_cache_stats(self, response_id: int):
        """Update cache hit statistics"""
        try:
```

```
                self.cursor.execute("""
                    UPDATE cached_responses
                    SET times_served = times_served + 1,
                        last_served = CURRENT_TIMESTAMP
                    WHERE response_id = %s
                """, (response_id,))
                self.conn.commit()
        except Exception as e:
            print(f"Stats update error: {e}")

    def get_cache_stats(self):
        """Get cache performance statistics"""
        try:
            self.cursor.execute("""
                SELECT
                    COUNT(*) as total_cached,
                    AVG(times_served) as avg_reuse,
                    MAX(times_served) as max_reuse,
                    COUNT(CASE WHEN times_served > 1
                          THEN 1 END) as reused_count
                FROM cached_responses
            """)
            return self.cursor.fetchone()
        except Exception as e:
            print(f"Stats retrieval error: {e}")
            return None
```

## 15.10   Appendix J: Claude API Integration

**RAG-Based Response Generation with Citations:**

```
import anthropic
import os
from typing import List, Dict

class ClaudeProcessor:
    def __init__(self):
        self.client = anthropic.Anthropic(
            api_key=os.getenv('ANTHROPIC_API_KEY')
        )
        self.model = "claude-sonnet-4-20250514"

    def generate_fitness_response(self, user_question: str,
                                  papers: List[Dict],
                                  conversation_context: str = None) -> Dict:
        """
        Generate evidence-based fitness response using RAG
        """
        # Build context from papers
        context = self._build_paper_context(papers)

        # Construct system prompt
        system_prompt = """You are a fitness research expert.
        Provide evidence-based answers using ONLY the research papers provided.

        CRITICAL RULES:
        1. Cite every claim with [Author et al., Year, PMID: xxxxx]
        2. Only use information from provided papers
        3. Acknowledge limitations if papers don't fully answer question
        4. Use plain language while maintaining scientific accuracy
        5. Provide actionable recommendations when possible
        """

        # Construct user message
        user_message = f"""Research Papers:
{context}

Conversation History:
{conversation_context if conversation_context else 'None'}

User Question: {user_question}

Please provide a comprehensive answer with proper citations."""

        try:
            # Call Claude API
            response = self.client.messages.create(
                model=self.model,
                max_tokens=2000,
                system=system_prompt,
                messages=[{
                    "role": "user",
                    "content": user_message
                }]
            )

            response_text = response.content[0].text
```

```
            # Validate response has citations
            validation = self.validate_response(response_text)

            return {
                'success': True,
                'text': response_text,
                'tokens_used': {
                    'input': response.usage.input_tokens,
                    'output': response.usage.output_tokens
                },
                'validation': validation
            }

        except Exception as e:
            return {
                'success': False,
                'error': str(e),
                'text': None
            }

    def _build_paper_context(self, papers: List[Dict]) -> str:
        """Format papers for Claude's context"""
        context_parts = []

        for i, paper in enumerate(papers[:10], 1):
            context_parts.append(f"\n[Paper {i}]")
            context_parts.append(f"Title: {paper.get('title', 'N/A')}")
            context_parts.append(f"Authors: {paper.get('authors', 'N/A')}")
            context_parts.append(f"Journal: {paper.get('journal', 'N/A')}")
            context_parts.append(f"Year: {paper.get('publication_date', 'N/A')[:4]}")
            context_parts.append(f"PMID: {paper.get('pmid', 'N/A')}")

            # Truncate abstract to fit in context
            abstract = paper.get('abstract', 'No abstract')
            if len(abstract) > 500:
                abstract = abstract[:500] + "..."
            context_parts.append(f"Abstract: {abstract}\n")

        return '\n'.join(context_parts)

    def validate_response(self, response_text: str) -> Dict:
        """Validate response quality"""
        word_count = len(response_text.split())
        citation_count = response_text.count('PMID:')

        issues = []
        if word_count < 100:
            issues.append("Response too short")
        if citation_count == 0:
            issues.append("No citations found")

        return {
            'valid': len(issues) == 0,
            'word_count': word_count,
            'citations': citation_count,
            'issues': issues
        }
```

# Repository Access

**GitHub:** https://github.com/rahulg2469/FitFact-Chatbot