

# **A Generative AI and Knowledge Graph Framework for TRIZ-Based Innovation**

## **Part 1: Foundational Technology Stack Analysis**

The creation of a generative AI-powered application for TRIZ (Theory of Inventive Problem Solving) requires a carefully architected technology stack. Each component must not only be best-in-class for its specific function but also integrate seamlessly to form a cohesive, intelligent system. This initial analysis addresses the critical technology choices that will serve as the foundation for the entire project, focusing on the selection of an open-source Retrieval-Augmented Generation (RAG) engine, and articulating the synergistic roles of the chosen Large Language Model (LLM) and graph database.

### **Section 1.1: Selecting the Optimal Open-Source RAG Engine**

The request for a "RAG Database" points to a central component in modern generative AI systems. A Retrieval-Augmented Generation (RAG) system is fundamentally a process that enhances an LLM's ability to generate accurate, context-aware responses by first retrieving relevant information from an external knowledge source.<sup>1</sup> This process involves a retriever, typically a vector database for searching unstructured data, and a generator, the LLM itself. The market offers a spectrum of solutions, from standalone vector databases that act as the retrieval component to comprehensive RAG engines that orchestrate the entire data ingestion and retrieval pipeline. For a complex, end-to-end application like the proposed TRIZ platform, the choice between a simple component and a full engine is a critical strategic decision that impacts development velocity, maintainability, and overall system capability.

### Subsection 1.1.1: The RAG Landscape for Local Development and Production Scaling

The selection of a RAG technology must satisfy two primary constraints: it must be deployable on a developer's local machine (Windows/Mac) for prototyping and development, and it must have a clear, viable path to scale for production use. The following analysis evaluates leading open-source options against these criteria, as well as their suitability for the specific demands of a TRIZ application, which involves processing dense, complex documents like patents and academic papers.

#### Vector Databases (The "Retrieval" Component)

These databases specialize in storing and searching high-dimensional vector embeddings, which are numerical representations of data like text or images.<sup>2</sup> They are the core of the "retrieval" step in RAG.

- **ChromaDB:** Positioned as "the open-source embedding database for AI applications," ChromaDB is renowned for its simplicity and ease of use.<sup>3</sup> It can be installed with a single pip command and run in-memory or persistently on a local machine, making it an excellent choice for rapid prototyping.<sup>5</sup> It supports both Python and JavaScript clients and offers features like metadata filtering and full-text search.<sup>3</sup> However, while ideal for small-to-medium workloads, its storage efficiency and performance are considered less robust for massive, enterprise-scale datasets when compared to more specialized databases like Milvus.<sup>8</sup> Its primary advantage is its low barrier to entry.
- **Milvus:** Milvus is an open-source vector database explicitly designed for "massive-scale vector data" and is a popular choice for enterprise-grade RAG applications.<sup>8</sup> It boasts high performance through features like GPU acceleration, distributed querying, and a variety of efficient indexing methods (e.g., HNSW, IVF) that allow for a trade-off between search speed and accuracy.<sup>2</sup> Critically for this project, Milvus offers **Milvus Lite**, a lightweight, embedded version that runs on macOS and Ubuntu (and thus on Windows via WSL) and can be installed directly via pip.<sup>11</sup> This provides a seamless development-to-production path, allowing developers to work locally with the same API they would use for a large-scale, distributed Milvus cluster in production.

- **Weaviate:** Weaviate is a versatile vector database that distinguishes itself with a strong focus on hybrid search (combining keyword and vector search) and a distributed architecture designed for scalability and efficiency.<sup>8</sup> It has a user-friendly, API-first design and offers modularity and vector compression to enhance storage efficiency.<sup>8</sup> The creators of Weaviate also developed Verba, an open-source, end-to-end RAG chatbot application, which serves as a powerful demonstration of the database's capabilities in a complete system.<sup>14</sup>
- **Qdrant:** Qdrant is a vector similarity search engine and database recognized for its high performance and advanced filtering capabilities.<sup>2</sup> A key feature is its ability to store arbitrary JSON payloads alongside vectors, enabling complex filtering conditions during search queries.<sup>2</sup> It also supports hybrid search using sparse vectors, which can significantly improve retrieval accuracy for queries with specific keywords that dense vector models might miss.

## RAG Frameworks & Engines (The "Orchestration" Layer)

These solutions go beyond simple vector storage and provide tools or full platforms to manage the entire RAG pipeline.

- **LangChain & LlamaIndex:** These are the dominant open-source *frameworks* for building LLM applications, including RAG systems.<sup>10</sup> They are not standalone servers but rather extensive Python/TypeScript libraries that provide modular, composable components for every step of the RAG process: data loaders, text splitters, embedding models, vector store integrations, and retrievers.<sup>17</sup> They offer maximum flexibility and control, integrating with virtually every popular LLM and vector database. The trade-off is that the developer is responsible for writing the code to assemble, orchestrate, and manage this entire pipeline.
- **RAGFlow:** RAGFlow stands out as a complete, open-source RAG *engine*.<sup>1</sup> It is a self-hosted platform that provides a streamlined, visual workflow for building and managing RAG systems.<sup>14</sup> Its core philosophy is "deep document understanding," focusing on intelligently parsing complex, unstructured data formats like PDF, DOCX, and even tables within documents—a feature that is exceptionally relevant for processing the TRIZ knowledge base of patents and scientific literature.<sup>1</sup> RAGFlow runs via Docker, ensuring compatibility with Windows, macOS, and Linux for local development.<sup>18</sup> Furthermore, its documented capability for "Graph-Enhanced RAG" or "GraphRAG" suggests a native understanding of graph structures, presenting a compelling opportunity for future synergy with the project's Neo4j database.<sup>1</sup>

The choice between a simple vector database, a development framework, or a

full-fledged engine is a pivotal architectural decision. The project already involves significant complexity in building a custom Flutter UI, a Gemini-powered reasoning core, and a novel TRIZ knowledge graph in Neo4j. The RAG component's role, while critical, is to perform the standardized (though complex) task of ingesting and retrieving context from a document corpus. Building this pipeline from scratch using a framework like LangChain would consume considerable development resources that could be better allocated to the unique, core-value aspects of the application—namely, the AI-driven TRIZ logic and the Neo4j integration.

An engine like RAGFlow abstracts away the undifferentiated heavy lifting of the RAG pipeline. It provides a pre-built, configurable server for data ingestion, chunking, embedding, and retrieval, complete with a user interface for managing knowledge bases.<sup>14</sup> This approach de-risks and accelerates the project by allowing the development team to focus on the application's novel features rather than reinventing the RAG wheel. RAGFlow's specific features, such as its deep document understanding and visual workflow editor, are tailored to solve the exact data preparation challenges this project will face when dealing with dense TRIZ source materials. Therefore, selecting a RAG engine over a more granular framework or a simple vector database represents a strategic move to optimize development effort and leverage a specialized, purpose-built tool.

**Table 1: Open-Source RAG/Vector DB Comparison Matrix**

Tool	Type	Key Features	Local Dev Setup (Win/Mac)	Scalability Path	Python Integration	Best For (Use Case)
<b>ChromaDB</b>	Vector DB	Simple API, in-memory & persistent modes, metadata filtering, full-text search. <sup>3</sup>	Very Easy (pip install). <sup>6</sup>	Good for small-medium scale; cloud-hosted version in preview. <sup>3</sup>	Excellent (native Python client). <sup>6</sup>	Rapid prototyping, small to medium applications, developers prioritizing simplicity. <sup>8</sup>
<b>Milvus</b>	Vector DB	Massive scalability, GPU acceleration	Easy (Milvus Lite via pip). <sup>11</sup>	Excellent (Designed for distribute	Excellent (PyMilvus SDK). <sup>11</sup>	Enterprise-grade, large-scale vector

		on, distributed querying, rich indexing options (HNSW, etc.). <sup>2</sup>		d, web-scale deployment). <sup>8</sup>		search applications requiring high performance. <sup>8</sup>
<b>Weaviate</b>	Vector DB	Hybrid search, distributed architecture, vector compression, API-first design. <sup>8</sup>	Moderate (Docker recommended).	Excellent (Distributed architecture). <sup>8</sup>	Good (Python client available).	Large-scale or enterprise deployments needing flexible hybrid search. <sup>8</sup>
<b>Qdrant</b>	Vector DB	High performance, advanced filtering with JSON payloads, sparse vector support. <sup>2</sup>	Moderate (Docker recommended).	Good (Production-ready service). <sup>2</sup>	Good (Python client available).	Applications requiring precise filtering and hybrid search accuracy.
<b>LangChain / LlamaIndex</b>	RAG Framework	Modular components, vast integrations, maximum flexibility, data connectors. <sup>1</sup>	Easy (Python library).	Depends on chosen vector DB and deployment architecture.	Native (Core libraries are Python/JS). <sup>16</sup>	Custom RAG applications where fine-grained control over the entire pipeline is required. <sup>15</sup>
<b>RAGFlow</b>	RAG Engine	End-to-end system, deep	Moderate (Docker Compose)	Good (Designed for	Good (Provides Python	Document-heavy applicatio

		document understanding, visual workflow editor, graph-RAG. <sup>1</sup>	. <sup>18</sup>	enterprise scalability). <sup>1</sup>	API). <sup>18</sup>	ns needing to process complex formats and accelerate development. <sup>10</sup>
--	--	---	-----------------	---------------------------------------	---------------------	---

### Subsection 1.1.2: Deep Dive and Getting Started Guides

To provide a practical understanding of the setup process for the top contenders, the following are condensed, step-by-step guides for local installation on a Windows or macOS machine.

#### Guide 1: Getting Started with ChromaDB (The "Quick Start" Option)

ChromaDB is designed for an extremely low-friction startup experience, running directly on your machine with a simple package installation.<sup>5</sup>

1. **Installation:** Open a terminal or command prompt and install the ChromaDB Python package using pip. This command installs everything needed to run Chroma locally.<sup>6</sup>

Bash

```
pip install chromadb
```

2. **Create a Client:** In a Python script, import the library and create a client. An ephemeral client runs in-memory and data is lost on termination. For persistence, a persistent client can be used which saves data to disk.<sup>6</sup>

Python

```
import chromadb
```

```
# For an in-memory instance
```

```
# client = chromadb.Client()
```

```
# For a persistent instance that saves to disk
```

```
client = chromadb.PersistentClient(path="/path/to/save/db")
```

3. **Create a Collection:** Collections are analogous to tables and are used to store

embeddings, documents, and metadata.<sup>7</sup>

Python

```
collection = client.get_or_create_collection(name="triz_documents")
```

4. **Add Documents:** Add text documents to the collection. Chroma automatically handles the embedding process using a default model (all-MiniLM-L6-v2), though this can be customized.<sup>5</sup> Each document requires a unique ID.

Python

```
collection.add(
    documents=,
    metadatas=[{"source": "triz_manual_p1"}, {"source": "triz_manual_p2"}],
    ids=["principle_1", "principle_2"]
)
```

5. **Query the Collection:** Query the collection using natural language text to find the most similar documents.<sup>6</sup>

Python

```
results = collection.query(
    query_texts=["how to break down a system?"],
    n_results=1
)
print(results)
```

## Guide 2: Getting Started with Milvus Lite (The "Scalable" Option)

Milvus Lite provides the core functionality of the powerful Milvus database in a lightweight package that can be embedded directly in a Python application, making it ideal for local development.<sup>13</sup>

1. **Prerequisites:** Ensure your system meets the requirements: macOS >= 11.0 or Ubuntu >= 20.04 (runnable on Windows via WSL2) and Python 3.8+.<sup>12</sup>
2. **Installation:** Install the pymilvus package, which includes Milvus Lite. The -U flag ensures you get the latest version.<sup>11</sup>

Bash

```
pip install -U pymilvus
```

3. **Set Up and Connect:** In your Python script, instantiate the MilvusClient. By providing a local file path as the URI, you instruct the client to use Milvus Lite and create a self-contained database file.<sup>11</sup>

Python

```
from pymilvus import MilvusClient
```

```
# This will create a file "milvus_triz.db" in the current directory  
client = MilvusClient("milvus_triz.db")
```

4. **Create a Collection:** Define a schema for your collection, specifying the vector dimension and any other fields. For simplicity, Milvus can also create a schema automatically if you enable the dynamic field option.<sup>11</sup>

Python

```
# Minimum required schema: name and vector dimension
```

```
client.create_collection(  
    collection_name="triz_patents",  
    dimension=768 # Example dimension for a common embedding model  
)
```

5. **Insert and Search Data:** Data insertion and searching follow a similar pattern to other vector databases, using the created client to interact with the collection. The API is consistent with the full distributed version of Milvus, ensuring a smooth transition to production.<sup>12</sup>

### Guide 3: Getting Started with RAGFlow (The "Engine" Option)

RAGFlow is a complete server application and is best run using Docker and Docker Compose, which standardizes the deployment across different operating systems.<sup>18</sup>

1. **Prerequisites:** Install Docker and Docker Compose on your system (Docker Desktop for Windows/Mac handles both). Ensure your machine meets the resource requirements ( $\geq 4$  cores,  $\geq 16$  GB RAM).<sup>18</sup>
2. **Configure System Settings:** RAGFlow's underlying Elasticsearch component requires a higher `vm.max_map_count`.
  - **On macOS:** Run `docker run --rm --privileged --pid=host alpine sysctl -w vm.max_map_count=262144`. For persistence, a LaunchDaemon plist file must be created as detailed in the official documentation.<sup>18</sup>
  - **On Windows with WSL2:** Run `wsl -d docker-desktop -u root sysctl -w vm.max_map_count=262144`. To make this change persistent, add `kernelCommandLine = "sysctl.vm.max_map_count=262144"` to the `[wsl2]` section of your `%USERPROFILE%\wslconfig` file and restart Docker Desktop.<sup>18</sup>
3. **Clone and Configure:** Clone the RAGFlow repository from GitHub and navigate into the docker directory.

Bash

```
git clone https://github.com/infiniflow/ragflow.git
```



```
cd ragflow/docker
```

4. **Start the Server:** Use Docker Compose to pull the pre-built images and start all the necessary services (RAGFlow server, MinIO, Elasticsearch, Redis, MySQL).<sup>18</sup>

```
Bash
```

```
# For CPU-based operation
```

```
docker compose -f docker-compose.yml up -d
```

5. **Access the UI:** Once all containers are running, access the RAGFlow web interface by navigating to <http://localhost> (or your machine's IP address) in a web browser. The default login credentials can be found in the documentation.<sup>19</sup>
6. **Usage:** From the UI, you can configure LLM providers (like Gemini), create a new "Knowledge Base," upload documents (PDFs, DOCX, etc.), and monitor the parsing and chunking process visually.<sup>18</sup> The platform provides a full end-to-end environment for the "Retrieval" and "Augmentation" parts of your system.

### Subsection 1.1.3: Definitive Recommendation and Justification

After a thorough analysis of the available open-source options, the definitive recommendation for this project is **RAGFlow**.

While ChromaDB offers unparalleled simplicity for getting started and Milvus provides a robust, direct path to massive scale, RAGFlow presents the most strategically advantageous balance of features, ease of use, and alignment with the specific goals of building a generative AI-powered TRIZ application. The justification for this recommendation rests on five key pillars:

1. **Significant Reduction in Development Overhead:** RAGFlow is an end-to-end engine, not just a library or a single-purpose database. It provides a fully operational, self-hosted service that manages the entire RAG pipeline: document ingestion, parsing, chunking, embedding, storage, and retrieval.<sup>14</sup> This is a substantial advantage, as it abstracts away a vast amount of complex, undifferentiated engineering work. By using RAGFlow, the development team is freed from building and maintaining this pipeline and can instead invest its resources in the novel and high-value components of the application: the sophisticated TRIZ reasoning logic, the unique Neo4j knowledge graph integration, and the polished Flutter user experience.
2. **Optimized for Complex Document Processing:** The knowledge base for a

powerful TRIZ application will be built upon dense, complex, and variably formatted documents, including patent filings, scientific papers, and technical manuals. RAGFlow is explicitly designed for this challenge, with "deep document understanding" as a core feature.<sup>1</sup> It includes capabilities like template-based chunking and visual inspection of parsing results, which are crucial for extracting high-quality, meaningful information from these sources and ensuring the LLM receives clean, relevant context.<sup>14</sup>

3. **Seamless Cross-Platform Local Deployment:** The requirement to support development on both Windows and Mac laptops is fully met by RAGFlow's Docker-based deployment model.<sup>18</sup> Docker provides a consistent, isolated environment, ensuring that the application behaves identically regardless of the host operating system. While the initial setup is more involved than a simple pip install, it is a one-time cost that yields a stable and reproducible development environment for the entire team.
4. **Future-Proofing with Graph-Aware Capabilities:** The mention of "Graph-Enhanced RAG" and "GraphRAG" in RAGFlow's feature set is a powerful indicator of its forward-looking architecture.<sup>1</sup> While the initial integration will involve the backend orchestrator querying Neo4j and RAGFlow separately, this suggests a potential future path for a more direct, native integration between the two. A RAG engine that is already thinking in terms of graphs is the ideal partner for a system where a knowledge graph is a first-class citizen.
5. **Fully Open-Source and Scalable:** RAGFlow is released under the permissive Apache 2.0 license, satisfying the core project requirement.<sup>20</sup> It is also designed with enterprise scalability in mind, ensuring that the architecture chosen for local development can be scaled up for a production environment without a complete re-platforming effort.<sup>1</sup>

In conclusion, RAGFlow is the most strategically sound choice. It offloads the generic RAG pipeline complexity, is purpose-built for the type of documents the project will use, works across development platforms, and is architecturally aligned with the project's graph-centric approach.

## Section 1.2: The Synergy of Gemini 2.5 and Neo4j

The selection of Google Gemini 2.5 as the LLM and Neo4j as the graph database is an astute architectural decision. These two components are not redundant but form a powerful, complementary partnership that functions as the application's "brain."

Gemini provides the fluid, generative reasoning capabilities, while Neo4j provides the structured, verifiable knowledge base. Their synergy is what will elevate the application from a simple information retrieval tool to a genuine ideation partner.

### Gemini 2.5 as the "Generative Reasoning Core"

The role of the Gemini 2.5 LLM extends far beyond simple text generation. It will act as the application's central cognitive processor, performing complex reasoning tasks that bridge the gap between ambiguous, human-language problem descriptions and the highly structured, logical world of the TRIZ methodology. Its key responsibilities will include:

- **Problem Abstraction:** Translating a user's free-text problem statement (e.g., "I need to make our product stronger, but it's becoming too heavy") into the formal language of TRIZ. This involves identifying the improving feature ("Strength") and the worsening feature ("Weight of stationary object") from the list of 39 standard engineering parameters.<sup>21</sup>
- **Contradiction Identification:** Formally identifying the core conflict as either a **Technical Contradiction** (where improving one parameter degrades another) or a **Physical Contradiction** (where a single parameter must have two opposite states, e.g., a coffee cup must be hot to drink but cool to hold).<sup>21</sup> This step is foundational to applying the correct TRIZ tools.
- **Solution Synthesis:** This is the most critical task. The LLM will receive structured inputs from the knowledge tier—the *what* from Neo4j (e.g., "Principles 1, 15, and 35 are recommended") and the *how* from RAGFlow (e.g., "Here are three patent snippets showing how Principle 15, 'Dynamics,' was applied in the automotive industry"). Gemini's role is to weave these disparate pieces of information into a coherent, creative, and domain-specific set of solution concepts for the user.
- **Creative Ideation:** Using the retrieved principles and examples as a conceptual launchpad, Gemini will brainstorm novel applications of those principles to the user's specific problem, moving beyond mere summarization to active, generative ideation.

### Neo4j as the "Structured Knowledge Memory"

If Gemini is the reasoning core, Neo4j is the system's high-fidelity, long-term memory. Its purpose is to store the explicit, proven, and highly interconnected knowledge of the TRIZ methodology itself. A graph is the natural way to represent this domain, as TRIZ is fundamentally about the relationships between problems, parameters, and principles.<sup>23</sup> Neo4j's role is to ground the LLM's creativity in established fact, ensuring the system's suggestions are both innovative and sound. Its key responsibilities

include:

- **Storing the TRIZ Canon:** The graph database will serve as the definitive, persistent model of the core TRIZ tools. This includes nodes for each of the 40 Inventive Principles, the 39 Engineering Parameters, and the relationships that form the Contradiction Matrix.<sup>23</sup> This creates a computable version of the foundational TRIZ knowledge.
- **Enabling Pattern Discovery:** The power of a graph database lies in its ability to traverse relationships and uncover non-obvious patterns.<sup>26</sup> The system will be able to execute complex Cypher queries that go beyond what a relational database could easily handle. For example, a query could find "inventive principles that are frequently used to solve contradictions involving 'Reliability' and have also been successfully applied in the 'Medical Devices' industry, and are related to the principle of 'Segmentation'." This ability to find distant but relevant connections is the essence of computer-aided innovation.
- **Grounding the LLM:** A primary weakness of LLMs is their potential to "hallucinate" or generate plausible but incorrect information. The Neo4j knowledge graph acts as a powerful grounding mechanism.<sup>28</sup> The LLM's reasoning is constrained and informed by the factual, verifiable relationships stored in the graph. When the system suggests a principle, it's not because the LLM "thinks" it's a good idea in isolation; it's because there is a concrete path in the knowledge graph, derived from decades of patent analysis, that links the user's specific type of contradiction to that principle. This provides a scaffold of logic that makes the AI's output trustworthy and explainable.

## Part 2: Architecting the AI-Powered TRIZ Knowledge Core

With the foundational technologies selected, the next critical step is to translate the abstract concepts of the TRIZ methodology into a concrete, machine-readable format. This involves a process of knowledge engineering: deconstructing TRIZ into its computable components and then designing a robust graph schema in Neo4j to represent this knowledge. This "Knowledge Core" will be the intellectual heart of the application, transforming it from a generic AI tool into a specialized innovation engine.

### Section 2.1: Deconstructing TRIZ for AI Implementation

Before any code can be written, the human-centric TRIZ framework must be analyzed and broken down into a set of structured data entities, relationships, and logical workflows that a computer system can understand and execute. This process ensures that the AI's operations are firmly rooted in the established theory.

### Subsection 2.1.1: Modeling the 40 Inventive Principles and Contradiction Matrix

The core of classical TRIZ, as developed by Genrich Altshuller through the analysis of hundreds of thousands of patents, is its highly structured knowledge base.<sup>30</sup> This inherent structure makes it exceptionally well-suited for representation within a database.

- **The 40 Inventive Principles:** Each of the 40 principles represents a generalized strategy for resolving an inventive problem. These principles, such as Principle 1: Segmentation, Principle 2: Taking Out (or Extraction), and Principle 4: Asymmetry, will be modeled as distinct entities within the system.<sup>31</sup> For each principle, the database must store its canonical number, its name, a detailed description of the strategy, and a collection of illustrative examples that show how the principle has been applied in practice.<sup>31</sup>
- **The 39 Engineering Parameters:** These parameters are the standardized characteristics of a technical system that can be improved or that might degrade as a result of an improvement. The list includes attributes like Weight of a moving object, Speed, Strength, Temperature, and Reliability.<sup>21</sup> These 39 parameters form the axes of the classical Contradiction Matrix and are the vocabulary used to define a technical conflict. Each will be a distinct entity in our model.
- **The Contradiction Matrix:** The matrix is the primary tool for resolving **Technical Contradictions**, which occur when an attempt to improve one parameter leads to the worsening of another.<sup>21</sup> For example, increasing the Power (improving parameter) of an engine often increases its Weight of moving object (worsening parameter).<sup>36</sup> The matrix itself is not a single data object but rather a dense network of relationships. It maps a specific pair of conflicting parameters to a small, rank-ordered set of Inventive Principles that have been historically successful at resolving that exact type of conflict.<sup>36</sup> This structure is perfectly suited for a graph model. The relationship can be expressed as:  
(:Parameter {name: 'Power'})-->(:Contradiction)<--(:Parameter {name: 'Weight of

moving object'})), where the :Contradiction node is then linked to suggested principles: (:Contradiction)-->(:InventivePrinciple). This approach allows the system to computationally replicate the core function of the matrix: given a conflict, find the most promising solution pathways.

### Subsection 2.1.2: The ARIZ Algorithm as a Computational Workflow

While the Contradiction Matrix is powerful for simpler problems, the **Algorithm for Inventive Problem Solving (ARIZ)** is the central analytical tool in TRIZ for tackling complex, non-standard problems.<sup>38</sup> ARIZ is not a static piece of data; it is a rigorous, step-by-step logical process designed to guide a problem-solver from a vague initial problem to a concrete, inventive solution.<sup>41</sup> For this application, ARIZ will be modeled as the primary, stateful workflow orchestrated by the backend system.

The application's main user journey should not be a static dashboard of tools but rather a guided, interactive "AI-Assisted ARIZ Session." This transforms the user experience from one of "looking up information" to one of "collaborative problem-solving." A novice user, who would be overwhelmed by a raw list of 40 principles, can be effectively guided through the innovation process. Each of the nine major parts of the ARIZ-85C version will correspond to a distinct phase in the application's workflow, managed by the backend and presented through a series of screens or interactions in the Flutter UI.<sup>39</sup>

An example mapping of ARIZ Part 1 (Analysis of the Problem) to a computational workflow demonstrates this concept:

1. **User Input (Start of Session):** The user enters a free-text description of their problem into the Flutter app.
2. **LLM Task (Corresponds to ARIZ Step 1.1):** The Python backend sends the user's text to the Gemini 2.5 API. The LLM is prompted to perform an initial analysis: identify the system's basic function, its key components, and the primary useful and harmful functions involved.<sup>42</sup>
3. **UI Interaction & Confirmation (Corresponds to ARIZ Step 1.2):** The backend processes the LLM's analysis and presents a structured summary to the user in the Flutter app. This includes a proposed "mini-problem" statement for confirmation, such as: "The goal is to achieve [useful function] without causing [harmful function], using minimal changes to the existing system." The user

validates or refines this statement. This human-in-the-loop step is crucial for ensuring accuracy.

4. **LLM Task & Neo4j Query (Corresponds to ARIZ Step 1.3):** Based on the confirmed problem statement, the backend tasks the LLM to formulate the core **Technical Contradiction**. The LLM is specifically prompted to identify the improving and worsening parameters from the canonical list of 39 TRIZ parameters. Once these parameters are identified (e.g., Strength and Weight), the backend immediately queries the Neo4j knowledge graph to find the inventive principles historically suggested for this specific conflict.
5. **Advance Workflow:** With the contradiction identified and potential solution paths retrieved from the knowledge graph, the system is ready to advance the user to the next part of the ARIZ workflow, where these principles will be explored in detail.

This approach operationalizes ARIZ, turning a complex human-centric algorithm into an executable sequence of AI-driven analysis, database lookups, and user interactions.

## Section 2.2: Designing the TRIZ Knowledge Graph in Neo4j

The Neo4j knowledge graph is the structured memory of the application. A well-designed graph data model is paramount, as it determines query performance, data integrity, and the system's ability to uncover the insightful, non-obvious connections that drive innovation.<sup>43</sup> The Labeled Property Graph (LPG) model, native to Neo4j, is used, which consists of nodes (entities), relationships (connections), and properties (attributes on nodes and relationships).<sup>25</sup>

### Subsection 2.2.1: Core Data Model: Nodes, Relationships, and Properties

The following schema defines the core entities and connections for the TRIZ knowledge graph. This model is designed to be both a faithful representation of TRIZ theory and an extensible framework for capturing new knowledge over time.

**Table 2: TRIZ Knowledge Graph Schema Definition**



Element Type	Name	Properties	Source Node	Target Node	Description/ Purpose
Node Label	Problem	id: string, description: text, domain: string, status: string, createdAt: datetime	-	-	Represents a specific problem instance submitted by a user.
Node Label	Parameter	id: int, name: string, description: text	-	-	Represents one of the 39 standard TRIZ engineering parameters (e.g., Speed, Strength).
Node Label	Inventive Principle	id: int, name: string, description: text, summary: text	-	-	Represents one of the 40 TRIZ inventive principles (e.g., Segmentation, Dynamics).
Node Label	Contradiction	id: string, type: string ('Technical' or 'Physical')	-	-	An intermediate node representing a specific conflict between two parameters.
Node Label	Example	id: string, source_document: string, text_snippet: text, embedding_i	-	-	A concrete, real-world example of a principle, ingested from the



		d: string			RAG corpus (e.g., a patent abstract).
Node Label	Industry	name: string	-	-	A classification node for problems, examples, and solutions (e.g., 'Aerospace', 'Medical').
Relationship Type	HAS_CONTRADICTION	is_primary: boolean	Problem	Contradiction	Connects a user's problem to the identified contradiction.
Relationship Type	IMPROVES		Contradiction	Parameter	Identifies the parameter to be improved in a technical contradiction.
Relationship Type	WORSENS		Contradiction	Parameter	Identifies the parameter that degrades in a technical contradiction.
Relationship Type	SUGGESTS	rank: int	Contradiction	Inventive Principle	<b>Encodes the Contradiction Matrix.</b> Links a conflict to a set of ranked,

					suggested principles.
Relationship Type	APPLIES_PRINCIPLE	context: text	Example	InventivePrinciple	Connects a real-world example to the specific principle it demonstrates.
Relationship Type	FROM_INDUSTRY		Example	Industry	Classifies an example by its industry of origin.
Relationship Type	RELATED_TO	weight: float	InventivePrinciple	InventivePrinciple	A heuristic link capturing conceptual similarity between principles (e.g., 'Segmentation' is related to 'Fragmentation').

This schema provides a rich, interconnected structure. It not only captures the classical TRIZ matrix (Contradiction -> SUGGESTS -> InventivePrinciple) but also enriches it with real-world evidence (Example nodes from the RAG system) and semantic context (Industry and RELATED\_TO links).

### Subsection 2.2.2: Querying for Innovation: Cypher Patterns for Problem Solving

The value of the graph model is realized through queries. Neo4j's Cypher query language is a declarative, pattern-matching language that is exceptionally expressive for exploring graph structures.<sup>44</sup> The following examples demonstrate how the schema can be queried to support the innovation process.

### Example Query 1: Classic Contradiction Matrix Lookup

This query replicates the fundamental function of the TRIZ Contradiction Matrix. Given an improving parameter and a worsening parameter, it finds the rank-ordered list of suggested principles.

Cypher

```
// Given: improvingParam = "Power", worseningParam = "Weight of moving object"
MATCH (p_improve:Parameter {name: $improvingParam})
MATCH (p_worsen:Parameter {name: $worseningParam})

MATCH (c:Contradiction)-->(p_improve)
MATCH (c)-->(p_worsen)
MATCH (c)-->(ip:InventivePrinciple)

RETURN ip.id, ip.name, ip.summary, r.rank
ORDER BY r.rank ASC
```

This query efficiently finds the specific Contradiction node that links the two conflicting parameters and then follows the SUGGESTS relationships to retrieve the recommended principles, ordered by their historical effectiveness.

### Example Query 2: Finding Cross-Domain Inspiration

A core tenet of TRIZ is that solutions are repeated across industries.<sup>30</sup> This query helps a user find inspiration from outside their immediate domain. Given a principle suggested for their problem, it finds concrete examples of that principle from a different industry.

Cypher

```
// Given: principleName = "Dynamics", userIndustry = "Automotive"
MATCH (ip:InventivePrinciple {name: $principleName})
MATCH (ip)<--(ex:Example)-->(ind:Industry)
```

```
WHERE ind.name <> $userIndustry
```

```
RETURN ex.text_snippet, ex.source_document, ind.name AS inspiration_industry  
LIMIT 10
```

This query demonstrates the power of the enriched graph. It starts at a specific InventivePrinciple, traverses to its real-world Example nodes, and filters them to show applications from industries other than the user's own, directly stimulating cross-domain thinking.

### Example Query 3: Discovering Related Solution Strategies

This query uses the heuristic RELATED\_TO links to suggest alternative or complementary solution pathways. If a user is exploring one principle, this can show them other, conceptually similar principles.

Cypher

```
// Given: principleName = "Segmentation"  
MATCH (p1:InventivePrinciple {name: $principleName})--(p2:InventivePrinciple)  
RETURN p2.name, p2.summary, r.weight  
ORDER BY r.weight DESC
```

By traversing these relationships, the application can proactively suggest, "Users who found 'Segmentation' helpful also explored 'Fragmentation' and 'Local Quality'," guiding the user toward a more comprehensive set of potential solutions. These queries illustrate how a well-designed knowledge graph moves beyond being a simple data repository to become an active participant in the discovery and innovation process.

## Part 3: System Architecture and Implementation Blueprint

This section synthesizes the foundational technology selections and the TRIZ knowledge model into a cohesive, end-to-end system architecture. It provides a

comprehensive blueprint for the application, detailing the structure of the backend engine, the design of the frontend user portal, and the communication protocols that bind them together into a functional whole.

### Section 3.1: High-Level System Architecture

The proposed architecture is a modern, multi-tiered system designed for scalability, maintainability, and a clear separation of concerns. It consists of a client tier, a backend tier responsible for business logic and AI orchestration, and a knowledge tier for data persistence and retrieval.

A textual description of the architectural diagram is as follows:

- **Client Tier:** This is the user-facing layer, comprising the **Flutter Application**. It is designed as a single codebase that can be deployed across multiple platforms (Web, Windows, macOS, and potentially mobile).<sup>46</sup> The Flutter app is responsible for all UI rendering and user interaction. It communicates with the backend exclusively through a well-defined API.
- **Backend Tier (Python):** This is the central nervous system of the application. It is built using a modern Python web framework like **FastAPI** or **Flask**.<sup>48</sup> This tier hosts several key services:
  - **API Gateway:** A single, secure entry point for all requests from the Flutter client. It handles authentication, rate limiting, and routing requests to the appropriate internal service.
  - **TRIZ Orchestrator Service:** This is the core business logic module. It implements the state machine for the ARIZ algorithm, managing a user's progress through a problem-solving session. It is responsible for coordinating calls to the other services.
  - **Generative AI Service:** This is a dedicated wrapper around the **Google Gemini 2.5 API**. It manages prompt engineering, API key security, and the formatting of requests and responses to and from the LLM.
  - **Hybrid Retrieval Service:** This service is responsible for fetching information from the knowledge tier. It contains the logic to query both the structured data in Neo4j and the unstructured data in RAGFlow.
- **Knowledge Tier:** This layer is responsible for the storage and retrieval of all TRIZ-related knowledge. It consists of two independent, open-source databases:
  - **RAGFlow Server:** This RAG engine manages the knowledge base of

unstructured documents (e.g., patents, academic papers, technical articles). It handles the entire pipeline of ingestion, chunking, embedding, and vector search, exposing its functionality via an API.

- **Neo4j Database:** This graph database stores the highly structured TRIZ Knowledge Graph, containing the principles, parameters, contradictions, and their intricate relationships.

### End-to-End Data Flow for a Problem-Solving Session:

The following sequence illustrates the flow of data and logic during a typical user interaction:

1. A user inputs a problem description (e.g., "My drone's battery runs out too quickly") into the **Flutter App**.
2. The Flutter app packages this text into a JSON payload and sends a POST request to the `/api/v1/session/analyze-problem` endpoint on the **API Gateway**.
3. The API Gateway authenticates the request and forwards it to the **TRIZ Orchestrator Service**.
4. The Orchestrator initiates an ARIZ workflow. It calls the **Generative AI Service**, passing the problem description. The Generative AI Service uses a specialized prompt to instruct Gemini 2.5 to identify the technical contradiction (e.g., Improving: Duration of action, Worsening: Weight of moving object).
5. The Orchestrator receives the identified contradiction parameters. It then calls the **Hybrid Retrieval Service**.
6. The Hybrid Retrieval Service executes two parallel queries:
  - a. A Cypher query to the Neo4j Database to retrieve the list of InventivePrinciple nodes suggested for this specific contradiction.
  - b. A vector search query to the RAGFlow Server to find real-world examples and documents related to the identified parameters (e.g., patents about lightweight, long-duration batteries).
7. The Hybrid Retrieval Service returns the structured principles from Neo4j and the unstructured text snippets from RAGFlow to the Orchestrator.
8. The Orchestrator sends this combined, multi-modal context back to the **Generative AI Service**. Gemini 2.5 is prompted to synthesize this information into a set of human-readable, creative solution ideas tailored to the user's drone problem.
9. The final, synthesized response is passed back through the API Gateway to the **Flutter App**, where it is displayed to the user in a structured, interactive format.

This architecture effectively decouples the UI, the core logic, and the data sources,

allowing for independent development, testing, and scaling of each component.

Section 3.2: The Backend Engine: Python, Gemini, and RAG Orchestration

The backend is the engine that drives the application's intelligence. Using Python is an excellent choice due to its extensive ecosystem of AI/ML libraries and robust web frameworks.<sup>48</sup> FastAPI is recommended for its high performance, asynchronous capabilities (crucial for handling concurrent I/O operations to the LLM and databases), and automatic API documentation generation.

Subsection 3.2.1: API Design for Frontend-Backend Communication

A clear, well-documented RESTful API is the contract between the frontend and backend, enabling parallel development and ensuring maintainability. All data will be exchanged in JSON format.<sup>48</sup>

Table 3: Core API Endpoint Definitions

Endpoint	Method	Request Body (JSON)	Success Response (200 OK, JSON)	Description
/api/v1/auth/register	POST	{ "email": "...", "password": "..." }	{ "user_id": "...", "token": "..." }	Registers a new user account.
/api/v1/auth/login	POST	{ "email": "...", "password": "..." }	{ "user_id": "...", "token": "..." }	Authenticates a user and returns a session token.
/api/v1/sessions	POST	{ "problem_title": "...", "initial_description": "..." }	{ "session_id": "...", "status": "created" }	Creates a new TRIZ problem-solving session.

/api/v1/sessions/{session_id}/analysis	GET	(None)	{ "contradiction": {...}, "parameters": [...]}	Retrieves the AI's analysis of the problem, including the identified contradiction.
/api/v1/sessions/{session_id}/principles	GET	(None)	{ "principles": [{ "id":..., "name":..., "summary":... }] }	Retrieves the list of inventive principles suggested for the session's contradiction.
/api/v1/sessions/{session_id}/examples	GET	?principle_id=1	{ "examples": [{ "snippet":..., "source":... }] }	Retrieves real-world examples for a specific inventive principle.
/api/v1/sessions/{session_id}/synthesis	POST	{ "selected_principles": }	{ "ideas": [{ "title":..., "description":... }] }	Triggers the final synthesis of ideas based on selected principles and returns the generated concepts.
/api/v1/sessions/{session_id}	GET	(None)	{ "session_id":..., "problem_title":..., "history": [...] }	Retrieves the entire state and history of a problem-solving session.

### Subsection 3.2.2: The Generative Reasoning Core

This module is the interface to the Gemini 2.5 API and is responsible for a critical task: prompt engineering. The quality of the AI's output is directly proportional to the quality of the prompts it receives. For a specialized domain like TRIZ, simple, open-ended prompts are insufficient. The system must employ sophisticated prompting strategies to guide the LLM's reasoning process and ensure accurate,



reliable results.

The approach should be to break down complex reasoning into a sequence of simpler, verifiable steps, a technique analogous to Chain-of-Thought reasoning. Instead of asking the LLM one large, complex question, the system will ask a series of smaller, targeted questions. For instance, to identify a contradiction, the prompt should not be "What is the contradiction here?". A more robust, multi-step prompt would be:

"You are an expert in the TRIZ methodology. Analyze the following user-submitted problem statement.

1. **Identify the primary useful function** the user is trying to achieve.
2. **Identify the primary harmful effect** or negative consequence that occurs.
3. **Map the desired improvement** to one of the following 39 standard TRIZ Engineering Parameters: [List of 39 parameters].
4. **Map the harmful effect** to one of the same 39 parameters.
5. **State the final Technical Contradiction** in the format: 'Improving [Parameter A] worsens'."

This structured prompt forces the model to externalize its reasoning process, which improves accuracy and allows the system to potentially validate each step.

Furthermore, this can be enhanced with **Few-Shot Prompting**, where the prompt includes 2-3 complete, high-quality examples of correctly analyzed problems before presenting the new user problem. This technique primes the model, showing it exactly the format and quality of output that is expected, which has been shown to be highly effective for adapting general-purpose LLMs to specialized tasks.<sup>49</sup>

### Subsection 3.2.3: The Hybrid Retrieval Pipeline

The "TRIZ Orchestrator" is the heart of the backend. It acts as a state machine for the ARIZ workflow and implements the hybrid retrieval pipeline. At each step of the workflow, it determines what knowledge is required and orchestrates queries to the appropriate data sources.

This hybrid approach is essential. A purely graph-based retrieval would miss the rich, unstructured context found in patents and articles. A purely vector-based RAG would lack the logical, proven structure of the TRIZ methodology. The orchestrator combines

the strengths of both:

- **When a contradiction is identified:** It queries **Neo4j** for the structured, canonical answer: the list of suggested principles. This is a query for *deductive truth*.
- **When a principle is being explored:** It queries **RAGFlow** for unstructured, real-world context: examples of how that principle has been applied. This is a query for *inductive evidence*.

The orchestrator then aggregates these two distinct types of information and passes them to the Gemini core for the final synthesis step. This fusion of structured logical reasoning (from the graph) and contextual pattern recognition (from RAG) is a hallmark of advanced, knowledge-grounded AI systems.

### Section 3.3: The User Portal: Flutter/Dart for Interactive Ideation

The frontend, built with Flutter and Dart, is the user's window into this complex system. Its primary goal is to make the powerful backend accessible and intuitive through a guided, interactive experience. Flutter's cross-platform capabilities are a significant asset, enabling the creation of a single application that runs natively on Windows, macOS, and in the browser.<sup>46</sup>

#### Subsection 3.3.1: UI/UX Flow for a Guided TRIZ Session

The user interface should guide the user through the ARIZ problem-solving process, abstracting away the complexity of the underlying AI and database queries. The flow should be designed as a collaborative conversation between the user and the AI.

1. **Screen 1: Problem Definition:** A clean, minimalist screen with a large text area where the user can describe their problem in natural language.
2. **Screen 2: AI Analysis & Confirmation:** After the user submits the problem, this screen presents the AI's initial analysis in a structured format. For example: "I understand your problem. It seems you are facing a **Technical Contradiction**: You want to improve **Parameter X**, but this is causing **Parameter Y** to get worse. Is this correct?" This "human-in-the-loop" confirmation step is vital. It gives the

user agency, builds trust, and allows for correction, dramatically improving the overall accuracy of the session.

3. **Screen 3: Principle Explorer:** Once the contradiction is confirmed, this screen displays the inventive principles suggested by the Neo4j knowledge graph, perhaps as a series of interactive cards. Each card would show the principle's name, number, and a brief summary.
4. **Screen 4: Example Browser:** Clicking on a principle card navigates the user to this screen. It displays a list of real-world examples for that principle, retrieved by RAGFlow. The examples should be sourced, linking back to the original patent or article, enhancing transparency and credibility.<sup>18</sup> The UI could allow filtering these examples by industry.
5. **Screen 5: Ideation Workbench:** This is the final workspace. It presents the AI-synthesized solution concepts generated by Gemini. This should be an editable space where the user can refine the AI's ideas, add their own notes, and ultimately export a final innovation report.

### Subsection 3.3.2: State Management and API Integration in Flutter

Building a multi-screen application with a complex, asynchronous data flow requires a robust state management strategy.

- **State Management:** For an application of this scale, a predictable state management library like **Bloc** or **Riverpod** is highly recommended over simpler approaches. These libraries help manage the application state in a structured way, separating business logic from the UI and making the app easier to test and maintain as it grows. The state would include the current session ID, the problem description, the identified contradiction, retrieved principles, examples, and generated ideas.
- **API Integration:** The Flutter http package, or a more advanced wrapper like the dio package, will be used to handle all communication with the Python backend's REST API.<sup>48</sup> For each API endpoint, a corresponding data model (a Dart class) should be created to represent the expected JSON structure. This allows for type-safe handling of data received from the backend. All API calls are asynchronous and should be handled gracefully in the UI, with loading indicators to show when the app is "thinking" (i.e., waiting for a response from the backend) and clear error messages if a request fails.

## Part 4: Phased Development Roadmap and Strategic Insights

This final part of the report provides a pragmatic, step-by-step plan for building the application and situates the project within the context of current academic and industry research. This ensures the development process is manageable and the final product is informed by the latest advancements in the field.

### Section 4.1: A Phased Development Plan

A complex project like this should not be built monolithically. An iterative, phased approach is essential to manage risk, gather feedback, and deliver value incrementally. The following four-phase plan breaks the project into manageable milestones, from a foundational MVP to a full-featured, scalable application.

Table 4: Phased Development Roadmap

Phase	Key Objectives	Core Tasks	Deliverable/ Milestone
<b>Phase 1: Foundation &amp; Core Knowledge Ingestion (MVP)</b>	Establish the core infrastructure and populate the knowledge bases with the canonical TRIZ data.	Backend: Deploy Neo4j and RAGFlow servers (e.g., via Docker Compose). Set up a basic Python/Fast API project. Data: Write Python scripts to parse and ingest the 40	A functional backend with populated Neo4j and RAGFlow instances, verifiable via direct database queries and the RAGFlow UI.

		<p>Principles, 39 Parameters, and the full Contradiction Matrix into the Neo4j graph.</p> <p>RAG: Use the RAGFlow UI to create an initial knowledge base and upload a seed corpus of 50-100 high-quality TRIZ articles and patents.</p>	
<p><b>Phase 2: Implementing the AI Reasoning &amp; Hybrid Retrieval Pipeline</b></p>	<p>Develop the core AI logic and the system's ability to respond to a problem statement.</p>	<p>Backend: Build the Generative AI Service to interface with the Gemini API. Implement the TRIZ Orchestrator to manage the initial steps of the ARIZ workflow.</p> <p>API: Implement the key API endpoints for problem analysis (/sessions/{id}/analysis) and principle retrieval (/sessions/{id}/principles).</p>	<p>A functional backend API that can accept a problem description via a tool like Postman or cURL and return a structured JSON response containing the identified contradiction and suggested principles.</p>

		<p>d}/principles)</p> <ul style="list-style-type: none"> <li>Logic: Write and thoroughly test the prompt engineering strategies and the hybrid queries to Neo4j and RAGFlow.</li> </ul>	
<b>Phase 3: Building the Interactive Flutter UI &amp; User Portal</b>	<p>Create the user-facing application and connect it to the backend services.</p>	<p>Frontend: Develop the Flutter application screens based on the defined UI/UX flow (Problem Definition, AI Analysis, Principle Explorer). Integration: Implement the API client in Flutter to call the backend endpoints and handle the responses. State Management : Set up a robust state management solution (e.g., Bloc, Riverpod) to manage the</p>	<p>A working, interactive prototype of the application. A user can log in, start a new session, enter a problem, see the AI's analysis, and view the list of suggested principles.</p>

		<p>session state across different screens.</p> <p>Features:</p> <p>Implement basic user authentication and session management.</p>			
<p><b>Phase 4: Scaling, Evaluation, &amp; Advanced Features</b></p>	<p>Enhance the application for production use, implement evaluation, and add advanced capabilities.</p>	<p>Data:</p> <p>Massively expand the RAG corpus with thousands of patents and articles to improve the quality and breadth of examples.</p> <p>Evaluation:</p> <p>Implement an evaluation framework (e.g., using a library like RAGAS 10) to quantitatively measure the quality, relevance, and novelty of the AI's suggestions. Incorporate a user feedback mechanism (e.g., rating generated</p>	<p>Features:</p> <p>Explore a multi-agent architecture where specialized AIs handle different parts of the ARIZ process.<sup>51</sup></p> <p>Automate the analysis of new patents using RAGFlow to continuously enrich the</p>	<p>Example nodes in the Neo4j graph.</p> <p>Ops: Prepare for production deployment with robust logging, monitoring, and scaling strategies.</p>	<p>A production-ready, scalable, and continuously improving application with a rich knowledge base and mechanisms for quality assurance.</p>

		ideas).			
--	--	---------	--	--	--

## Section 4.2: Insights from Academic and Industry Research

Connecting this project to the broader scientific conversation provides validation for its core concepts and offers strategic direction for its evolution. Recent research at the intersection of AI and TRIZ is highly relevant.

### Subsection 4.2.1: Summary of Key Papers and Their Architectural Takeaways

Several recent academic papers directly support and inform the proposed architecture.

- **AutoTRIZ: Artificial Ideation with TRIZ and Large Language Models** <sup>52</sup>:  
This paper from early 2024 demonstrates the fundamental feasibility of using LLMs to automate the core reasoning process of TRIZ. The researchers developed a tool that takes a problem statement and automatically generates a solution report following the TRIZ workflow.
  - **Architectural Takeaway:** This provides strong academic validation for the project's central premise. The challenges noted in the paper—namely, the complexity of TRIZ concepts and the reliance on user expertise—are precisely the problems this application is designed to solve through its guided, AI-assisted workflow. Their focus on automating the "Method of Inventive Principles" aligns perfectly with the proposed plan to model the Contradiction Matrix in Neo4j.
- **TRIZ Agents: A Multi-Agent LLM Approach for TRIZ-Based Innovation** <sup>51</sup>:  
This 2025 conference paper proposes moving beyond a single LLM to a multi-agent system for TRIZ-based innovation. It simulates a team of specialized AI agents working together.
  - **Architectural Takeaway:** This provides a clear and compelling roadmap for the "Advanced Features" in Phase 4. While the initial architecture can be monolithic for simplicity, it can be evolved into a more sophisticated system where a `Problem_Analyst_Agent` identifies the contradiction, a `Solution_Generator_Agent` brainstorms ideas based on principles, and a



Critique\_Agent evaluates the generated solutions for feasibility. This modular approach can lead to higher-quality outputs.

- Evaluating the Effectiveness of Generative AI in TRIZ: A Comparative Case Study<sup>49</sup>.

This study compared the results of a real-world innovation project using traditional TRIZ tools versus a generative AI-assisted approach. The key finding was that the combination of GenAI and TRIZ produces "feasible, cross-domain preliminary conceptual directions with satisfactory scientific substantiation."

- **Architectural Takeaway:** This provides empirical evidence that the output of the proposed application will have tangible, real-world value. It validates the core value proposition: saving innovators significant time during the crucial early stages of conceptual exploration and idea generation.
- Extracting TRIZ Contradictions from Patents using RAG<sup>53</sup>:  
A recent paper highlights a method using a RAG approach with GPT-4 to accurately extract TRIZ contradictions and principles directly from patent text.
  - **Architectural Takeaway:** This directly validates the architectural decision to include a powerful RAG component (RAGFlow). It confirms that RAG is a suitable and effective technology for processing the project's primary source documents (patents). This also suggests an advanced feature for Phase 4: using the RAG pipeline not just for retrieving examples, but for automatically analyzing new patents to discover and populate new Example and Contradiction nodes in the Neo4j graph, creating a self-enriching system.

#### Subsection 4.2.2: Addressing Key Challenges: Prompt Engineering, AI Evaluation, and Fostering User Trust

The research also highlights critical challenges that must be addressed architecturally.

- **Prompt Engineering:** The variability of LLM outputs based on subjective prompts is a known issue.<sup>49</sup> As detailed in Section 3.2.2, the solution is to move away from simple prompts and implement highly structured, few-shot, chain-of-thought prompts that guide the LLM's reasoning process. This is a key engineering task for the backend team.
- **AI Evaluation:** A significant challenge is determining the "quality" of the AI's generated ideas. This is inherently subjective. The architecture must include mechanisms for evaluation. This includes implementing an automated evaluation

pipeline using frameworks like RAGAS, which can measure metrics like faithfulness (how well the answer is supported by the context) and relevance.<sup>10</sup> Crucially, it also means building a human feedback loop into the Flutter UI, allowing users to rate the novelty, feasibility, and usefulness of the generated solutions. This feedback is invaluable data for future fine-tuning and system improvement.

- **Fostering User Trust:** Users will not rely on a tool they perceive as an untrustworthy "black box." The proposed architecture addresses this directly through three key design principles:
  1. **Grounding:** The Neo4j knowledge graph ensures that all suggestions are grounded in the established, verifiable logic of the TRIZ methodology. The AI is not inventing principles; it is applying them.
  2. **Transparency and Citations:** The RAGFlow component is capable of providing citations, linking retrieved examples back to their source documents.<sup>17</sup> The Flutter UI must surface these links, allowing users to "click through" and verify the information for themselves.
  3. **Human-in-the-Loop:** As designed in the UI/UX flow, the user is an active participant, not a passive recipient. By asking the user to confirm the AI's analysis at key stages, the system becomes a collaborative partner, which builds confidence and ensures the final output is aligned with the user's understanding of the problem.

### Section 4.3: Concluding Recommendations and Future Outlook

This report has laid out a comprehensive blueprint for creating a novel, open-source, generative AI-powered TRIZ application. The analysis leads to a set of concluding strategic recommendations to guide the project's development and ensure its success.

#### Core Strategic Recommendations:

1. **Adopt a RAG Engine to Accelerate Development:** The most critical early decision is to leverage a full-featured RAG engine like RAGFlow. This will abstract away the significant complexity of building a data ingestion and retrieval pipeline, allowing the development team to focus on the application's unique value proposition.
2. **Structure the User Journey Around the ARIZ Algorithm:** Frame the

application's core feature as an "AI-Assisted ARIZ Session." This transforms the product from a static knowledge base into an interactive, guided problem-solving partner, which is a far more compelling and valuable user experience.

3. **Implement a Hybrid Knowledge Architecture:** The true power of this system comes from the synergy between the two knowledge tiers. The backend must be architected to leverage both the structured, logical reasoning of the Neo4j graph and the unstructured, contextual evidence from the RAGFlow engine. This hybrid retrieval pipeline is central to generating innovative yet grounded solutions.
4. **Follow a Phased, Iterative Development Plan:** The provided four-phase roadmap should be followed to manage complexity, mitigate risks, and build the application incrementally. Starting with a foundational MVP and progressively adding layers of intelligence and functionality is the most pragmatic path to success.

### Future Outlook:

The architecture described in this report provides a powerful foundation, but it also opens the door to a fascinating future outlook. The ultimate vision for this application should be to create a **self-improving innovation system**. This aligns with one of the core ideas of TRIZ itself: the "laws of technical systems evolution," which state that systems evolve toward increased ideality and dynamism.<sup>30</sup>

By capturing user feedback on the quality of generated solutions and logging which solution paths (contradiction -> principle -> example) lead to highly-rated outcomes, the system can begin to learn. This data can be used to update the weights on relationships within the Neo4j graph (e.g., strengthening the link between a specific contradiction and a principle that proves highly effective in a certain domain). The RAG engine can be used to automatically scan new patent databases, identify emerging technologies, and use the LLM to map them back to TRIZ principles, continuously enriching the graph with new, relevant examples.

In this future state, the application would not just be a tool for applying TRIZ; it would become a living embodiment of it, a system that actively participates in and learns from the process of innovation, evolving its own knowledge and becoming a more effective creative partner over time.

### Works cited

1. Top 10 Open-Source RAG Frameworks you need!! - DEV Community, accessed July 1, 2025,  
[https://dev.to/rohan\\_sharma/top-10-open-source-rag-frameworks-you-need-3fh](https://dev.to/rohan_sharma/top-10-open-source-rag-frameworks-you-need-3fh)

e

2. Top 10 open source vector databases - Instaclustr, accessed July 1, 2025, <https://www.instaclustr.com/education/vector-database/top-10-open-source-vector-databases/>
3. Chroma, accessed July 1, 2025, <https://www.trychroma.com/>
4. chromadb - PyPI, accessed July 1, 2025, <https://pypi.org/project/chromadb/>
5. Learn How to Use Chroma DB: A Step-by-Step Guide | DataCamp, accessed July 1, 2025, <https://www.datacamp.com/tutorial/chromadb-tutorial-step-by-step-guide>
6. Getting Started - Chroma Docs, accessed July 1, 2025, <https://docs.trychroma.com/getting-started>
7. Getting Started - Chroma Docs, accessed July 1, 2025, <https://docs.trychroma.com/docs/overview/getting-started?lang=typescript>
8. Top 9 Vector Databases as of June 2025 | Shakudo, accessed July 1, 2025, <https://www.shakudo.io/blog/top-9-vector-databases>
9. Top 5 Open Source Vector Databases in 2024 - LLM & Langchain Blogs, accessed July 1, 2025, <https://www.langchain.ca/blog/top-5-open-source-vector-databases-2024/>
10. 15 Best Open-Source RAG Frameworks in 2025 - Firecrawl, accessed July 1, 2025, <https://www.firecrawl.dev/blog/best-open-source-rag-frameworks>
11. Quickstart | Milvus Documentation, accessed July 1, 2025, <https://milvus.io/docs/quickstart.md>
12. How to Get Started with Milvus, accessed July 1, 2025, <https://milvus.io/blog/how-to-get-started-with-milvus.md>
13. Run Milvus Lite Locally | Milvus Documentation, accessed July 1, 2025, [https://milvus.io/docs/milvus\\_lite.md](https://milvus.io/docs/milvus_lite.md)
14. 15 Best Open-Source RAG Frameworks in 2025 - Apidog, accessed July 1, 2025, <https://apidog.com/blog/best-open-source-rag-frameworks/>
15. Compare the Top 7 RAG Frameworks in 2025 - Pathway, accessed July 1, 2025, <https://pathway.com/rag-frameworks/>
16. RAG Frameworks You Should Know: Open-Source Tools for Smarter AI | DataCamp, accessed July 1, 2025, <https://www.datacamp.com/blog/rag-framework>
17. 25+ Best Open Source RAG Frameworks in 2025 - Signity Solutions, accessed July 1, 2025, <https://www.signitysolutions.com/blog/best-open-source-rag-frameworks>
18. Get started | RAGFlow, accessed July 1, 2025, <https://ragflow.io/docs/dev/>
19. How to Use RAGFlow(Open Source RAG Engine): A Complete Guide - Apidog, accessed July 1, 2025, <https://apidog.com/blog/ragflow/>
20. RAGFlow with Local LLMs and Ollama: Step-by-Step Guide for Free Retrieval-Augmented Generation - YouTube, accessed July 1, 2025, <https://www.youtube.com/watch?v=kAuVcCwqgwQ&pp=0gcJCfwAo7VqN5tD>
21. TRIZ Contradictions - We ask and you answer! The best answer wins! - Benchmark Six Sigma Forum, accessed July 1, 2025, <https://www.benchmarksixsigma.com/forum/topic/39396-triz-contradictions/>

22. Unlocking Innovative Solutions with TRIZ: A Powerful Problem-Solving Methodology - SixSigma.us, accessed July 1, 2025, <https://www.6sigma.us/six-sigma-in-focus/triz-inventive-problem-solving-methodology/>
23. Neo4j-based knowledge storage and visualization. - ResearchGate, accessed July 1, 2025, [https://www.researchgate.net/figure/Design-knowledge-visualization-search-example-display\\_fig1\\_373561926](https://www.researchgate.net/figure/Design-knowledge-visualization-search-example-display_fig1_373561926)
24. The Instantiation Process from RDF to Neo4j. | Download Scientific Diagram - ResearchGate, accessed July 1, 2025, [https://www.researchgate.net/figure/The-Instantiation-Process-from-RDF-to-Neo4j\\_fig5\\_354553567](https://www.researchgate.net/figure/The-Instantiation-Process-from-RDF-to-Neo4j_fig5_354553567)
25. Graph database concepts - Getting Started - Neo4j, accessed July 1, 2025, <https://neo4j.com/docs/getting-started/appendix/graphdb-concepts/>
26. How Do You Know If a Graph Database Solves the Problem? - Neo4j, accessed July 1, 2025, <https://neo4j.com/blog/developer/how-do-you-know-if-a-graph-database-solves-the-problem/>
27. Does anyone have any specific examples of graph DB (neo4j) solutions to business problems. - Reddit, accessed July 1, 2025, [https://www.reddit.com/r/datascience/comments/11umvz6/does\\_anyone\\_have\\_any\\_specific\\_examples\\_of\\_graph/](https://www.reddit.com/r/datascience/comments/11umvz6/does_anyone_have_any_specific_examples_of_graph/)
28. Knowledge Graph Use Cases Driving Innovation in 2025 - PingCAP, accessed July 1, 2025, <https://www.pingcap.com/article/knowledge-graph-use-cases-2025/>
29. Knowledge graphs | The Alan Turing Institute, accessed July 1, 2025, <https://www.turing.ac.uk/research/interest-groups/knowledge-graphs>
30. TRIZ - Wikipedia, accessed July 1, 2025, <https://en.wikipedia.org/wiki/TRIZ>
31. TRIZ 40 Design Principles - ipface.org, accessed July 1, 2025, [https://www.ipface.org/pdfs/reading/TRIZ\\_Principles.pdf](https://www.ipface.org/pdfs/reading/TRIZ_Principles.pdf)
32. Guide to the 40 TRIZ Principles (Table Format) - Quality Gurus, accessed July 1, 2025, <https://www.qualitygurus.com/guide-to-the-40-triz-principles-table-format/>
33. 40 TRIZ Principles, accessed July 1, 2025, [https://www.triz40.com/aff\\_Principles\\_TRIZ.php](https://www.triz40.com/aff_Principles_TRIZ.php)
34. 40 Inventive Principles for Business - The Triz Journal, accessed July 1, 2025, <https://the-trizjournal.com/40-inventive-business-principles-examples/>
35. What is TRIZ Matrix and How to Use it ? - Management Study Guide, accessed July 1, 2025, <https://www.managementstudyguide.com/triz-matrix.htm>
36. TRIZ Contradiction Analysis - CreatingMinds.org, accessed July 1, 2025, [http://creatingminds.org/tools/triz/triz\\_contradiction\\_analysis.htm](http://creatingminds.org/tools/triz/triz_contradiction_analysis.htm)
37. 40 principles of invention - Wikipedia, accessed July 1, 2025, [https://en.wikipedia.org/wiki/40\\_principles\\_of\\_invention](https://en.wikipedia.org/wiki/40_principles_of_invention)
38. ARIZ - Innovation Wiki by verrocchio Institute, accessed July 1, 2025, <https://www.innovation.wiki/en/method/ariz/>
39. Introduction to ARIZ - ResearchGate, accessed July 1, 2025,

[https://www.researchgate.net/profile/Umakant-Mishra/publication/235742388\\_An\\_Introduction\\_to\\_ARIZ\\_-The\\_Algorithm\\_of\\_Inventive\\_Problem\\_Solving/links/09e4151306cccc7248000000/An-Introduction-to-ARIZ-The-Algorithm-of-Inventive-Problem-Solving.pdf](https://www.researchgate.net/profile/Umakant-Mishra/publication/235742388_An_Introduction_to_ARIZ_-The_Algorithm_of_Inventive_Problem_Solving/links/09e4151306cccc7248000000/An-Introduction-to-ARIZ-The-Algorithm-of-Inventive-Problem-Solving.pdf)

40. ARIZ - Technical Innovation Center, Inc., accessed July 1, 2025, <https://triz.org/ariz/>
41. Ariz Explored: A Step-by-Step Guide to Ariz, the Algorithm for Solving Inventive-图书简介, accessed July 1, 2025, <http://www.chinaims.org/index.php?c=content&a=show&id=69>
42. ARIZ : The Algorithm for Inventive Problem Solving, accessed July 1, 2025, <https://www.metodolog.ru/triz-journal/archives/1998/04/d/index.htm>
43. What is graph data modeling? - Getting Started - Neo4j, accessed July 1, 2025, <https://neo4j.com/docs/getting-started/data-modeling/>
44. Using the Neo4j Graph Database and Cypher To Solve This Brain Teaser. Why Argue?, accessed July 1, 2025, <https://medium.com/neo4j/using-the-neo4j-graph-database-and-cypher-to-solve-this-brain-teaser-why-argue-350fde86da14>
45. What is TRIZ - The Theory of Inventive Problem Solving?, accessed July 1, 2025, <https://www.sixsigmadaily.com/triz-theory-of-inventive-problem-solving/>
46. Flutter architectural overview, accessed July 1, 2025, <https://docs.flutter.dev/resources/architectural-overview>
47. Flet: Build multi-platform apps in Python powered by Flutter, accessed July 1, 2025, <https://flet.dev/>
48. Can We Use Python as a Backend for a Flutter App?, accessed July 1, 2025, <https://shivlab.com/blog/use-python-as-a-backend-for-flutter-app/>
49. Evaluating the Effectiveness of Generative AI in TRIZ: A Comparative Case Study, accessed July 1, 2025, [https://www.researchgate.net/publication/385348808\\_Evaluating\\_the\\_Effectiveness\\_of\\_Generative\\_AI\\_in\\_TRIZ\\_A\\_Comparative\\_Case\\_Study](https://www.researchgate.net/publication/385348808_Evaluating_the_Effectiveness_of_Generative_AI_in_TRIZ_A_Comparative_Case_Study)
50. linking between backend on python code and flutter - Stack Overflow, accessed July 1, 2025, <https://stackoverflow.com/questions/59121076/linking-between-backend-on-python-code-and-flutter>
51. A Multi-Agent LLM Approach for TRIZ-Based Innovation - SciTePress, accessed July 1, 2025, <https://www.scitepress.org/Papers/2025/133219/133219.pdf>
52. arxiv.org, accessed July 1, 2025, <https://arxiv.org/html/2403.13002v2>
53. Automotive innovation landscaping using LLM - arXiv, accessed July 1, 2025, <https://arxiv.org/html/2409.14436v1>