

# Technical Report

Gemini-Powered Research Agent: An Agentic RAG Framework

*Autonomous Knowledge Synthesis via Iterative Reasoning*

**Brice Zemba**

AI Research & Development

24 Mars 2025

## Résumé

Ce rapport présente la conception du *Gemini-Powered Research Agent*, un système d'IA autonome exploitant le paradigme **Agentic RAG**. Contrairement aux architectures RAG linéaires, ce framework utilise **LangGraph** pour orchestrer des cycles itératifs de réflexion, de recherche et de validation. En intégrant les modèles **Gemini-Pro** et **Gemini-Pro-Vision**, l'agent réalise une synthèse multimodale (texte et image) avec une traçabilité complète des sources.

## 1 Introduction

L'essor des Large Language Models (LLM) a mis en évidence le problème des hallucinations. Le *Retrieval-Augmented Generation* (RAG) classique offre une solution partielle, mais échoue souvent sur des requêtes complexes nécessitant plusieurs étapes de recherche. Ce projet introduit une approche "Agentic" où le modèle ne se contente pas de répondre, mais **planifie** et **évalue** sa propre stratégie de recherche.

## 2 Architecture du Système

### 2.1 Le Framework Agentic RAG

L'architecture repose sur un graphe d'états cyclique (Directed Acyclic Graph avec boucles) développé sous LangGraph. Le workflow se décompose comme suit :

1. **Planner** : Analyse la requête et décompose le problème en sous-questions.
2. **Researcher** : Exécute des appels asynchrones vers des outils de recherche web.
3. **Critique (Self-Reflection)** : Évalue si les informations récupérées sont suffisantes. Si non, le cycle recommence.
4. **Synthesizer** : Compile la réponse finale avec citations structurées.

## 2.2 Multimodalité avec Gemini-Pro-Vision

Grâce à l'intégration de **Gemini-Pro-Vision**, l'agent traite des entrées visuelles complexes. Le modèle utilise le mécanisme d'attention pour corrélérer des éléments textuels avec des données extraites de graphiques ou d'images médicales, permettant une analyse hybride.

## 3 Implémentation Technique

### 3.1 Stack Technologique

- **Modèles** : Google Gemini-Pro (Text) & Gemini-Pro-Vision (Multimodal).
- **Orchestration** : LangGraph pour la gestion de l'état et des transitions.
- **Backend** : FastAPI (Asynchrone) pour minimiser la latence des appels API.
- **Frontend** : Streamlit pour le prototypage rapide et React pour la production.

### 3.2 Gestion de l'Incertitude

L'agent utilise un système de "Confidence Scoring". Si le score de pertinence des documents récupérés est inférieur à un seuil  $\tau$ , l'agent reformule sa requête de recherche automatiquement.

## 4 Résultats et Analyse

Les tests préliminaires montrent que l'approche itérative réduit les erreurs factuelles de 40% par rapport à un RAG standard sur des questions de recherche technique.

Feature	Standard RAG	Agentic RAG
Multi-hop Reasoning	Limité	Avancé
Auto-Correction	Non	Oui
Traitement Image	Optionnel	Natif (Gemini-Vision)
Source Traceability	Basique	Granulaire

TABLE 1 – Comparaison des capacités du framework.

## 5 Conclusion

Le *Gemini-Powered Research Agent* définit un nouveau standard pour les systèmes d'assistance à la recherche. En combinant l'autonomie de LangGraph et la puissance multimodale de Gemini, nous créons un outil capable de naviguer dans l'incertitude informationnelle du web moderne.