

## **DECLARATION**

We hereby declare that the Capstone Project Phase - 1 entitled "**A Cost-Effective, Proactive Hallucination Routing System for LLMs**" has been carried out by us under the guidance of **Dr. Ravi Gorripati, Associate Professor** and submitted in partial fulfilment of the course requirements for the award of the degree of **Bachelor of Technology in Computer Science and Engineering (Artificial Intelligence and Machine Learning)** of **PES University, Bengaluru** during the academic semester Aug – Dec 2025. The matter embodied in this report has not been submitted to any other university or institution for the award of any degree.

**PES1UG23AM313**  
**PES1UG23AM917**  
**PES1UG23AM314**  
**PES1UG23AM315**

**Sourabh S M**  
**Chandan R**  
**Sreephaneesha K**  
**Sri Charan D A**

---

---

---

---

## **ACKNOWLEDGEMENT**

I would like to express my gratitude to **Dr. Ravi Gorripati**, Department of Computer Science & Engineering (Artificial Intelligence and Machine Learning), PES University, for his continuous guidance, assistance, and encouragement throughout the development of UE23AM320A- Capstone Project Phase – 1.

I am grateful to Capstone Project Coordinators, **Dr. Chetana Srinivas** for organizing, managing, and helping with the entire process.

I take this opportunity to thank **Dr. Jayashree R**, Professor & Chairperson, Department of Computer Science and Engineering (Artificial Intelligence and Machine Learning), PES University, for all the knowledge and support I have received from the department. I would like to thank **Dr. K S Sridhar** Dean of Faculty and Registrar, PES University for his help.

I am deeply grateful to **Prof. Jawahar Doreswamy, Chancellor, PES University**, **Dr. Suryaprasad J, Vice-Chancellor, PES University**, for providing me with various opportunities and enlightenment every step of the way. Finally, Phase 1 of the project could not have been completed without the continual support and encouragement I have received from my family and friends.

## ABSTRACT

Large Language Models are powerful, but they also love to make things up—confidently. Most existing solutions try to fix these hallucinations after the model has already produced a bad answer, which wastes compute and doesn't stop the damage. This paper takes a different route: instead of reacting after the fact, it judges before the model even begins generating whether the answer is likely to be trustworthy.

The proposed system estimates the model's confidence using three signals:

- How closely its internal representations match reliable reference embeddings,
- How steadily its reasoning progresses through layers, and
- A learned predictor trained directly on activation patterns.

These signals are combined into one confidence score. Based on that score, the system routes the question to one of several paths letting the small model answer if it's confident, using retrieval when certainty dips, escalating to a bigger model when needed, and falling back to a human only when the model is truly lost.

Across standard knowledge-heavy QA benchmarks, this proactive approach catches hallucinations far better than older methods while using significantly less compute.

Instead of treating hallucinations like an afterthought, this framework brings early awareness into LLM pipelines—making them faster, safer, and more reliable for real-world use

## **TABLE OF CONTENTS**

<b>Chapter No.</b>	<b>Title</b>	<b>Page No.</b>
1.	<b>INTRODUCTION</b>	<b>01</b>
2.	<b>PROBLEM DEFINITION</b>	<b>05</b>
3.	<b>LITERATURE SURVEY</b>	<b>06</b>
	3.1 Foundations of Hallucination in LLMs	
	3.2 Reactive Hallucination Detection Approaches	
	3.3 Retrieval-Augmented Generation (RAG)	
	3.4 Confidence-Aware and Proactive Methods	
	3.5 Limitations in Existing Literature	
4.	<b>RESEARCH / TECHNOLOGY GAPS AND CHALLENGES</b>	<b>10</b>
5.	<b>OBJECTIVES AND PROJECT SCOPE</b>	<b>13</b>
6.	<b>CONCLUSION OF CAPSTONE PROJECT PHASE – 1</b>	<b>18</b>
7.	<b>PLAN OF WORK FOR CAPSTONE PROJECT PHASE – 2</b>	<b>19</b>

## **REFERENCES/BIBLIOGRAPHY**

## **APPENDIX A DEFINITIONS, ACRONYMS, AND ABBREVIATIONS**

## **LIST OF TABLES**

<b>Table No.</b>	<b>Title</b>	<b>Page No.</b>
<b>Table 1.1</b>	<b>Reactive vs Proactive</b>	<b>03</b>
<b>Table 3.1</b>	<b>Summary of Literature Survey</b>	<b>09</b>
<b>Table 5.1</b>	<b>Types of Signals</b>	<b>14</b>
<b>Table 7.1</b>	<b>Plan of work</b>	<b>21</b>

## **LIST OF FIGURES**

<b>Figure No.</b>	<b>Title</b>	<b>Page No.</b>
<b>Figure 1.1</b>	<b>Traditional Reactive Hallucination Mitigation Pipeline</b>	<b>04</b>
<b>Figure 4.1</b>	<b>Conditional Hybrid Confidence System (3-Interval)</b>	<b>17</b>