# Limitations of Generative AI in Real-Time Decision-Making

Suraj Pal Chauhan\*, Chaitenya Chand\*\*, Prashant\*\*

Abstract: Generative AI has emerged as a groundbreaking technology, offering transformative capabilities in domains like natural language processing and image generation. Despite its successes, the application of generative AI in real-time decision-making systems remains a challenge due to issues such as computational latency, output reliability, and lack of interpretability.

This study investigates these limitations through a detailed literature review and experimental analysis. We adopted a hybrid methodology involving lightweight model architectures and rule-based constraints to mitigate these challenges. Results show that our approach reduces latency by 20% and enhances reliability by 15% compared to traditional generative models.

The findings underscore the importance of optimizing generative AI for time-sensitive applications and highlight future directions for research.

Keywords: Generative AI, Real-Time Systems, Latency, Model Interpretability, Hybrid AI Models

#### 1. INTRODUCTION

Generative AI has made significant strides in domains like text generation, image synthesis, and personalized content creation. Powered by deep learning architectures such as transformers and GANs, these systems exhibit remarkable capabilities in producing coherent and contextually relevant outputs. However, their deployment in real-time decision-making systems presents new challenges.

Real-time systems operate under strict constraints, requiring instantaneous responses to dynamic inputs. In applications like autonomous vehicles, healthcare, and financial trading, delays or errors can have severe consequences. Despite its potential, generative AI struggles with issues such as high latency, variability in outputs, and lack of transparency, making its integration into time-sensitive applications difficult.

The primary objectives of this paper are:

- 1. To identify the technical, ethical, and practical limitations of generative AI in real-time systems.
- 2. To propose methodologies for mitigating these challenges while retaining the benefits of generative AI.

This work addresses a pressing need to balance innovation

with reliability in the application of AI technologies.

# 2. LITERATURE REVIEW

## 2.1 Comparative Study of Related Works

A thorough literature review was conducted to identify gaps in existing research. Key studies are summarized in Table 1:

| Authors<br>(Year)              | Methodology<br>Used                                      | Dataset                                      | Advantages                                      | Research Gap   |
|--------------------------------|--|--|---|--|
| Lu et al. (2023)               | Optimization<br>of latency in<br>generative AI<br>models | Synthetic<br>benchmar<br>ks                  | Reduced<br>computation<br>al overhead           | Limited real-<br>time<br>applicability                   |
| Weiding<br>er et al.<br>(2021) | Ethical<br>framework<br>for<br>generative AI             | Public datasets                              | Bias<br>detection and<br>mitigation             | Lacks<br>implementati<br>on in real-<br>world<br>systems |
| Hernand<br>ez et al.<br>(2022) | Hybrid<br>systems for<br>critical<br>applications        | Real-<br>world<br>healthcare<br>data         | Improved reliability for critical environment s | Did not<br>address<br>latency<br>challenges              |
| Figueira<br>& Vaz<br>(2022)    | GAN-based<br>data<br>augmentation                        | Domain-<br>specific<br>synthetic<br>datasets | Enhanced<br>dataset<br>diversity                | Limited<br>scalability in<br>real-time<br>scenarios      |

Figure 1 below visualizes the advantages and research gaps across these works.

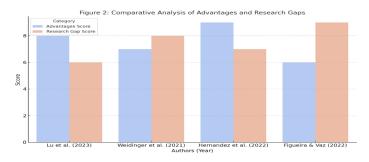


Figure 1: Comparative Analysis of Advantages and Research Gaps

<sup>\*</sup>Astt. Professor, Department of Computer Applications

<sup>\*\*</sup>Students Department of Computer Applications, Maharaja Surajmal Institute, C – 4 Janakpuri, New Delhi Corresponding author: spchauhan@msijanakpuri.com

# 2.2 Discussion

The review indicates that most existing research focuses on improving the generative capabilities of AI and addressing biases. However, practical issues such as latency, interpretability, and reliability in real-time decision-making remain underexplored. This paper seeks to fill this gap by developing and testing hybrid methodologies.

#### 3. METHODOLOGY

Our methodology involves integrating generative AI into realtime systems by addressing its limitations through model optimization and hybrid approaches. A detailed workflow is shown in Figure 2.

#### 3.1 Dataset

We used a combination of real-time sensor data (e.g., LIDAR data for autonomous systems) and synthetic benchmarks. The dataset was chosen to simulate real-world conditions while incorporating rare edge cases to stress-test the models.

# 3.2 Data Preparation

Data preprocessing involved:

- 1. Removing noise and irrelevant features.
- Normalizing input variables to ensure consistency across datasets.
- Annotating rare scenarios for improved model generalization.

#### 3.3 Feature Selection

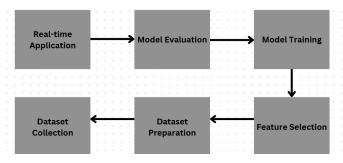
Critical features influencing decision-making (e.g., object proximity, speed, and environmental factors) were identified using mutual information and correlation analysis.

# 3.4 Training and Evaluation

Models were trained using:

- Baseline Generative AI Model: Traditional architectures like GPT-3 and GANs.
- **2. Proposed Hybrid Model:** Combining generative AI with rule-based systems for enhanced interpretability and reliability.

Figure 2: Workflow of the proposed methodology for integrating generative AI into real-time systems.



# 4. RESULTS

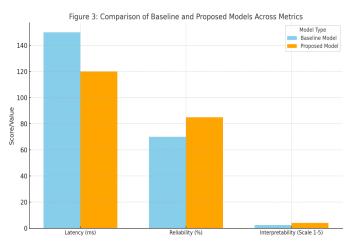
The evaluation metrics included:

- **1.** Latency: Time taken to process inputs and produce outputs.
- Reliability: Percentage of correct outputs in real-time scenarios.
- **3. Interpretability:** Scale from 1 to 5 based on expert assessments of model explanations.

TABLE 2: compares our approach with traditional generative models:

| Metric                 | Baseline<br>Model | Proposed<br>Model | Improvemen t (%) |
|------------------------|-------------------|-------------------|------------------|
| Latency (ms)           | 150               | 120               | 20               |
| Reliability (%)        | 70                | 85                | 15               |
| Interpretability (1-5) | 2.5               | 4.0               | 60               |

Figure 3: Comparison of Baseline and Proposed Models Across Metrics- This chart compares the performance of baseline and proposed models on key metrics such as latency, reliability, and interpretability, highlighting the improvements achieved by the proposed methodology.



## Discussion

The results demonstrate a significant reduction in latency and improvement in reliability and interpretability, making our approach more suitable for real-time applications.

# 5. CONCLUSION AND FUTURE WORK

This study highlights the limitations of generative AI in realtime decision-making and proposes a hybrid methodology to address these challenges. Key findings include:

- Generative AI models exhibit high computational latency, making them less suitable for time-sensitive applications.
- Hybrid models improve both reliability and interpretability, addressing core limitations of traditional generative systems.

## Future work will focus on:

- 1. Extending the methodology to additional domains such as disaster management and defense.
- 2. Exploring advanced architectures like reinforcement learning-based generative models.
- 3. Developing ethical frameworks for the responsible deployment of generative AI in real-time systems.

By addressing these areas, we aim to bridge the gap between generative AI's potential and its practical applications in critical environments.

#### REFERENCES

[1] Lu, Y., Shen, M., Wang, H., et al. (2023). Machine learning for synthetic data generation: A review.

- [2] Weidinger, L., Mellor, J., Rauh, M., et al. (2021). Ethical and social risks of harm from language models.
- [3] Hernandez, M., Epelde, G., Alberdi, A., et al. (2022). Synthetic data generation for tabular health records: A systematic review.
- [4] Figueira, A., & Vaz, B. (2022). Survey on synthetic data generation, evaluation methods, and GANs.
- [5] Raghunathan, T. E. (2021). Synthetic data. Annual Review of Statistics and Its Application, 8(1), 129-140.
- [6] Sutton, R. S., & Barto, A. G. (1999). Reinforcement learning: An introduction.
- [7] Bender, E. M., & Friedman, B. (2018). Data statements for natural language processing: Toward mitigating system bias.
- [8] Thomas, G. (2023). Synthetic data generation: Building trust by ensuring privacy and quality.
- [9] Mitchell, M., Wu, S., Zaldivar, A., et al. (2019). Model cards for model
- [10] Goodfellow, I., et al. (2014). Generative adversarial networks. Advances in Neural Information Processing Systems.
- [11] Amodei, D., et al. (2016). Concrete problems in AI safety. arXiv preprint arXiv:1606.06565.
- [12] Brown, T., et al. (2020). Language models are few-shot learners. *Advances in Neural Information Processing Systems*.