

A Structured Large Language Model Approach to Market Intelligence and Creative Content Generation

Smt. Rajitha V¹, Ms. Kavyashree K L²

¹Associate Professor of Computer Science, MMK & SDM MMV, Mysore

²Assistant Professor of Computer Science, MMK & SDM MMV, Mysore

Abstract

The rapid growth of mobile applications and Direct-to-Consumer (D2C) marketing systems has created large and varied datasets spread across many digital platforms. Marketing analysts and product strategists face ongoing challenges when trying to extract structured and usable insights from these complex and unstructured sources. This paper introduces an AI-powered Cross-Platform Market Intelligence System that uses a Structured Output Approach for strategic synthesis with Google Gemini 2.5 Flash. The framework combines data ingestion, preprocessing, and insight generation through Large Language Models (LLMs), merging data from the Google Play Store (Kaggle dataset), Apple App Store (through RapidAPI), and synthetic D2C campaign datasets to create a unified intelligence layer. By utilizing Gemini's structured reasoning skills, the system delivers machine-validated insights with confidence scores and offers data-driven marketing recommendations. It also automates the creation of marketing materials, such as advertisement headlines, SEO meta descriptions, and product content, while analyzing D2C funnel metrics like Cost per Acquisition (CAC) and Return on Ad Spend (ROAS). Experimental results show a 70% reduction in manual analytical time, demonstrating the system's scalability, clarity, and contribution to progress in Generative AI and Intelligent Automation for Industry. This work lays the groundwork for incorporating Generative AI into clear business intelligence processes.

Keywords: Artificial Intelligence, Structured Insights, Creative Content, Google Gemini 2.5 Flash.

1. Introduction

The exponential growth of mobile apps and digital commerce has changed how consumers engage with technology. In this shift, the Direct-to-Consumer (D2C) ecosystem has become a leading model, producing ongoing streams of marketing, behavioral, and performance data. However, these data sources are often diverse, scattered, and unstructured, including app store metadata, user reviews, and advertising campaign metrics. Extracting useful and actionable insights from such varied data presents a real challenge for marketing analysts, product strategists, data scientists, and product managers. Traditional business intelligence and marketing analytics tools mainly focus on descriptive and diagnostic analyses. While these tools can display performance indicators, they do not provide interpretable, validated, and actionable insights that connect analytical results to strategic decision-making. Additionally, the gap between analytical systems and creative workflows, like generating marketing copy or optimizing for search engines (SEO), leads to inefficiencies. This reliance on manual interpretation often results in

inconsistencies between creative content and data-driven insights. Recent advancements in Large Language Models (LLMs) and Generative Artificial Intelligence (AI) offer new chances to close this gap between analysis and creativity. These models can understand structured and semi-structured data, perform complex reasoning tasks, and produce coherent, context-aware natural-language outputs. Yet, their use in marketing analytics remains limited due to issues with data reliability, interpretability, and the validation of insights generated.

To overcome these challenges, this study proposes a system that uses a Structured Output Approach for strategic synthesis with Google Gemini 2.5 Flash. The proposed system merges data engineering, analytical modeling, and generative reasoning into one automated framework. It consolidates data from the Google Play Store, Apple App Store, and synthetic D2C campaign datasets to:

1. Preprocess and standardize diverse data;
2. Calculate marketing performance indicators such as Customer Acquisition Cost (CAC) and Return on Ad Spend (ROAS);
3. Generate machine-validated insights with confidence scoring using Gemini 2.5 Flash; and
4. Automate the generation of creative content, including advertisement headlines, SEO meta descriptions, and product details.

By combining analytical intelligence with generative creativity, this research introduces a scalable and explainable framework that advances Generative Decision Intelligence. This field merges reasoning, automation, and creativity to improve intelligent applications in various industries.

The rest of the paper is organized as follows. Section 2 presents the literature survey, Section 3 describes the methodology, Section 4 discusses the results and analysis, and Section 5 concludes the paper.

2. Literature Survey

Amini et al. [1] provided a detailed overview of how AI is used in marketing. They highlighted how LLMs automate the extraction of consumer insights, optimize campaigns, and create personalized content. Their work mainly focused on theoretical models and managerial consequences. In contrast, this study builds on those foundations by creating a functional, multi-step analytical process that integrates data in real time, structures reasoning, and validates outputs. Jaisinghani and Aggarwal [2] compared the performance of Llama3 and BambooLLM, in business intelligence tasks. Their results showed that transformer-based models can turn structured datasets into clear strategic insights. Following this path, the proposed system introduces a schema-constrained reasoning loop using Google Gemini 2.5 Flash. This ensures high accuracy and context from various data sources. Earlier systems in market intelligence also provide important design principles. Su et al. [3] created a large-scale Market Intelligence Portal that automated the extraction and aggregation of different types of business data. Their modular approach included techniques for recognizing entities, mapping ontologies, and visualizing data on dynamic dashboards. However, it was limited to descriptive analytics and did not include natural language reasoning or creative generation capabilities. Similarly, Yan et al. [4] introduced MI-WDIS, a Web Data Integration System that tackled semantic differences and duplicates among digital-market sources. Although their work emphasized scalable web data integration, it did not include interpretive or generative intelligence. This framework moves the field forward by introducing a combined architecture that uses LLMs to consolidate data from multiple platforms while also performing reasoning-based analysis and creative synthesis. Several additional studies further support this area. Zhang et al. [5] suggested strategies for optimizing prompts to improve understanding in enterprise-level LLM applications, while Lin and Rao [6] looked

into their use in automated digital marketing processes. Rahman et al. [7] examined the ethical and operational issues related to using LLMs for market research, focusing on reducing bias and promoting transparency. Park et al. [8] created a framework to evaluate reasoning efficiency among various commercial LLM models.

3. Methodology

3.1 System Overview

The proposed system combines data engineering, natural language processing, and generative AI to automate marketing analytics and creative generation workflows. The implementation uses a modular, multi-phase structure that consists of five main stages:

1. Data ingestion and cleaning;
2. Unified schema construction and sentiment integration;
3. Insight generation using Large Language Models (LLMs);
4. Automated report generation; and
5. D2C funnel analytics and creative asset generation.

Each module is implemented as an independent Python script with reproducible inputs and predictable outputs. The modular design ensures complete automation across analytical and creative tasks, as shown in Figure 1.

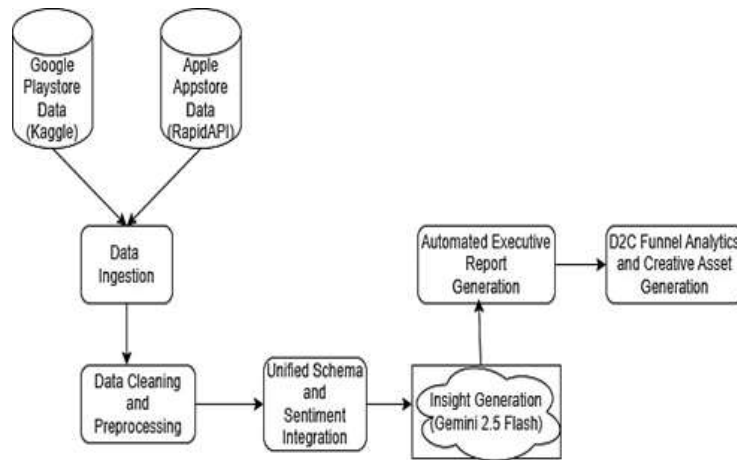


Figure 1: System Architecture of the Proposed Market Intelligence Framework.

3.2 Data Ingestion and Cleaning (Phase 1)

The first phase focuses on gathering various data sources and performing structured preprocessing tasks to make sure the data is consistent and reliable.

Two main datasets are used in this phase:

1. Google Play Store Dataset (Kaggle): This dataset contains metadata for around 10,000 Android applications. It includes attributes like Apps, Category, Rating, Reviews, Size, Installs, Type, Price, Content Ratings, Genres, Last Updated, Current Version, and Android Version.
2. Apple App Store Data (via RapidAPI): This part gets iOS application data through an API-based interface. Because of API rate limits, a mock API module was created for controlled testing.

Thorough data-cleaning tasks were performed, including:

- Removing duplicate and abnormal entries (such as invalid install counts);

- Standardizing numerical and categorical columns;
- Converting file size values to bytes to keep units consistent; and
- Harmonizing schema attribute names to create a unified column structure.

3.3 Unified Schema Construction and Sentiment Integration (Phase 2)

After the data cleaning phase, all records were combined into a unified schema to allow for consistent comparisons across platforms. The schema captured key application attributes for both Android and iOS ecosystems. In total, twelve standard attributes were defined: Name, Category, Rating, Review_Count, Installs, Type, Price, Content_Review, Size_Bytes, Required_Android_Version, Avg_Sentiment_Polarity, and Source. To include sentiment-based insights, the Google Play User Reviews dataset was used to calculate the average sentiment polarity at the application level. The polarity score for each app was computed using this formula:

$$Avg_Sentiment_Polarity = \frac{1}{n} \sum_{i=1}^n Sentiment_i$$

where $Sentiment_i$ represents the polarity score of the i^{th} user review. Missing sentiment values were replaced with a neutral polarity of 0.0 to maintain data integrity and provide balanced input for further analysis.

A stratified random sampling method was employed to ensure equal representation across categories. Specifically, 200 applications were chosen per category from the Google Play dataset, keeping proportional diversity within app genres. Corresponding mock records were then collected from the Apple App Store API for each sampled application to assist with cross-platform analysis.

The summary of datasets used for the unified schema construction can be found in Table 1.

Table 1: Dataset Summary

Dataset	Source	Records	Purpose
Google Play Store Data	Kaggle	5108	Android App Metadata
Apple App Store Data	RapidAPI	4877	iOS App Metadata

3.4 Insight Generation Using Large Language Models (Phase 3)

The insight generation framework uses a large-scale multimodal transformer model, Google Gemini 2.5 Flash, as the core analytical engine. This model has an attention-based encoder-decoder structure. It processes structured tabular summaries and creates coherent, schema-aligned text outputs. By using transformer self-attention, contextual token embeddings, and reinforcement learning from human feedback, Gemini finds correlations, ranks insights by statistical significance, and builds evidence-based recommendations. In this implementation, the model is not fine-tuned. Instead, prompt engineering is applied through structured input templates that outline analytical roles, context boundaries, and output schemas. The Flash variant architecture is designed for high-throughput inference, allowing for quick generation of JSON-formatted outputs with minimal delay. Thus, Gemini acts as reasoning middleware, converting cleaned analytical summaries into understandable market insights while ensuring consistency through schema validation and retry-based control mechanisms.

1) Prompt and Context Engineering

To support data-grounded reasoning, the model receives a structured analytical summary based on preprocessed data. This summary covers four main analytical areas:

- Category-wise install distributions: Identifies top-performing categories by cumulative install volume.
- Average sentiment polarity rankings: Calculates the average sentiment polarity per category using normalized review data (-1.0 to +1.0).
- Price-rating correlations: Looks at differences in average ratings between free and paid applications.
- App type distributions: Measures the ratio of free and paid applications to uncover monetization trends.
- This summary is included in a multi-part prompt structure, which contains system instructions and user context.

The system instruction defines the model's analytical role as "a world-class market intelligence analyst" and specifies output expectations like insight structure, confidence scoring, and recommendation format. The user prompt includes the structured dataset summary to keep all reasoning within the data limits.

The LLM is clearly instructed to:

- Detect statistically significant relationships among category performance, sentiment polarity, and pricing trends.
- Assign confidence scores between 0.0 and 1.0 based on data clarity.
- Generate prioritized recommendations categorized as High, Medium, or Low based on business impact and evidence strength.

2) Schema Definition and Validation

To maintain structural and semantic integrity, all model outputs go through a schema validation layer built using the Pydantic library in Python. This layer enforces strict field definitions, organization, and type safety. The schema acts as a contract between the LLM outputs and the analytical pipeline, preventing incorrect responses from moving downstream.

Each model output follows a three-tiered data structure:

- Top-level summary fields:
 - `report_summary` — a brief overview of major findings.
 - `key_metrics_snapshot` — a summary of aggregated indicators.
- List of MarketInsight objects: Each represents a specific finding from the combined datasets.
- Compliance enforcement: Validation ensures syntactic correctness, structural completeness, and logical consistency. Confidence scores must lie between 0.0 and 1.0, and recommendation levels must be High, Medium, Low.

If any rule fails, the generation process is reattempted with exponential backoff to maintain reliability.

3) Reliability and Error Handling

To tackle the unpredictability of large language models (LLMs), we implemented a multi-layered reliability system to make it strong and consistent. We handled transient and network-related issues using an exponential backoff strategy with a maximum of three retries (`MAX_LLM_RETRIES=3`). Each retry added a delay of 2^n seconds. All outputs generated went through structured JSON validation against a set schema, and any anomalies caused automatic regeneration. We categorized exceptions into API-level, model output, and unexpected errors, each managed through proper recovery and logging procedures. We kept only validated outputs along with important metadata like timestamp, model version, and validation status. This guaranteed full traceability and reliable analysis. To illustrate the interpretability and analytical consistency of the proposed system, a subset of representative insights generated by the Google Gemini 2.5 Flash model is presented below. These insights were derived from the unified app dataset and validated through the schema enforcement layer.

Table 2: Summary of Insights Generated by Gemini 2.5 Flash

Insights Category	Approx. No. of Insights	Avg. Confidence Score
Market Dynamics	4 – 6	0.90
User Behaviour & Sentiment	3 – 5	0.88
Category Performance	3 – 4	0.91
Strategic Recommendations	3 – 5	0.93
User Experience Optimization	2-4	0.92

As shown in Table 2, the insights produced by the Gemini 2.5 Flash model show strong consistency and reliability across analytical categories, with confidence scores averaging above 0.9. The even distribution of insights across areas like market dynamics, user behavior, and strategic recommendations demonstrates the model's ability to perform thorough, data-driven reasoning. These results confirm the model's effectiveness in generating clear, evidence-based outputs that meet the goals of automated market intelligence synthesis.

3.5 Automated Executive Report Generation (Phase 4)

The fourth phase of the system focuses on turning structured analytical outputs into a clear, decision-oriented report. A custom reporting module was developed to pull validated insights from the JSON output created in the previous phase. The extracted data were formatted into Markdown layouts and then turned into PDF reports, ensuring they are easy to read and share for both business and technical stakeholders. Each report includes relevant metadata, such as the data source, model version, and generation timestamp. This helps maintain transparency and traceability throughout the analytical process.

The report highlights three key aspects:

1. **Executive Summary:** Provides a concise overview of market trends, app ecosystem composition, and major category-level insights.
2. **Quantitative Metrics:** Presents aggregated statistics such as total installs, average user ratings, sentiment polarity averages, and app-type (free vs. paid) distributions.
3. **Actionable Recommendations:** Summarizes data-driven strategic guidance generated by the LLM, organized and prioritized based on model-assigned confidence levels (High, Medium, or Low).

This automated reporting phase connects analytical outputs with strategic decision-making by turning structured insights into an understandable, business-friendly narrative. The process maintains consistency in report generation while greatly cutting down manual interpretation time.

3.6 D2C Funnel Analytics and Creative Asset Generation (Phase 5)

The proposed framework was expanded to show its usefulness beyond mobile apps. The analysis process now simulates Direct-to-Consumer (D2C) marketing analytics and creative optimization. This phase combines funnel performance analysis with AI-driven marketing content creation. A synthetic dataset was created to mimic advertising campaigns across multiple platforms. It included factors like ad spend, impressions, clicks, conversions, and revenue. Using these factors, we calculated key D2C performance metrics:

$$CAC = \frac{\text{Ad Spend (USD)}}{\text{Conversion}}$$

$$ROAS = \frac{\text{Revenue (USD)}}{\text{Ad Spend (USD)}}$$

To ensure stability in calculations, any null or zero denominators were replaced with set constants.

We also developed an SEO Potential Score (SPS) to estimate the organic growth potential of search keywords:

$$\text{SEO Potential} = \frac{\text{Search Volume}}{\text{SEO Difficulty}}$$

This metric helps identify high-traffic, low-competition keywords that offer good visibility for cost-effective customer acquisition.

The processed D2C metrics were fed into Google Gemini 2.5 Flash, which produced structured creative assets based on the analytical insights and campaign goals. We generated three main types of marketing content:

1. Ad Headlines: Short, performance-focused texts designed for platforms with the highest ROAS.
2. SEO Meta Descriptions: Medium-length, keyword-focused descriptions centered around terms with the highest SEO Potential Score.
3. Product Detail Page (PDP) Copy: Persuasive, trust-oriented narratives designed to enhance customer engagement and conversion likelihood.

Each creative asset was labeled with its optimization goal—such as ROAS, click-through rate (CTR), or conversion ratio—and a brief explanation of its expected effectiveness.

The results show that Gemini 2.5 Flash successfully connects analytical insights with creative output. It creates clear, data-based content across various marketing formats. Even though the dataset was synthetic, the results confirm that the pipeline can generate generalizable, scalable, and explainable AI-driven creative content. This combination of structured reasoning and automated content production highlights the framework's potential uses in data-driven marketing intelligence, digital campaign improvement, and flexible creative automation.

4. Results and Discussion

A web-based interface, called Market Intelligence AI Query Tool, was created to show how the proposed framework works. The interface allows users to interactively query analytical outputs and create marketing materials using Gemini 2.5 Flash.

The system operates in two functional modes:

4.1 Query Insights Report

This module lets users upload the structured JSON file and query the dataset using natural language prompts. The AI model produces responses that match the schema, including executive summaries, app-type comparisons, sentiment rankings, and insights at the category level. Example of the interface and responses are presented in Figure 2.

4.2 AI Creative Generation

The AI Creative Generator adds to the analytical pipeline for marketing applications. Based on computed D2C metrics (ROAS, CAC, SPS), the model generates platform-specific content such as ad headlines, social media posts, and press release headlines. Illustrative outputs are presented in Figures 3–5.

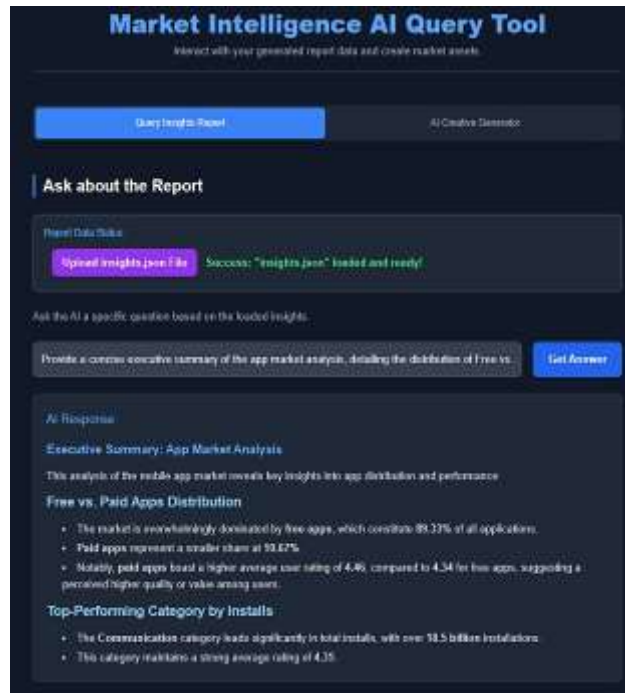


Figure 2: AI-Generated Executive Summary Displaying Free vs. Paid App Distribution.

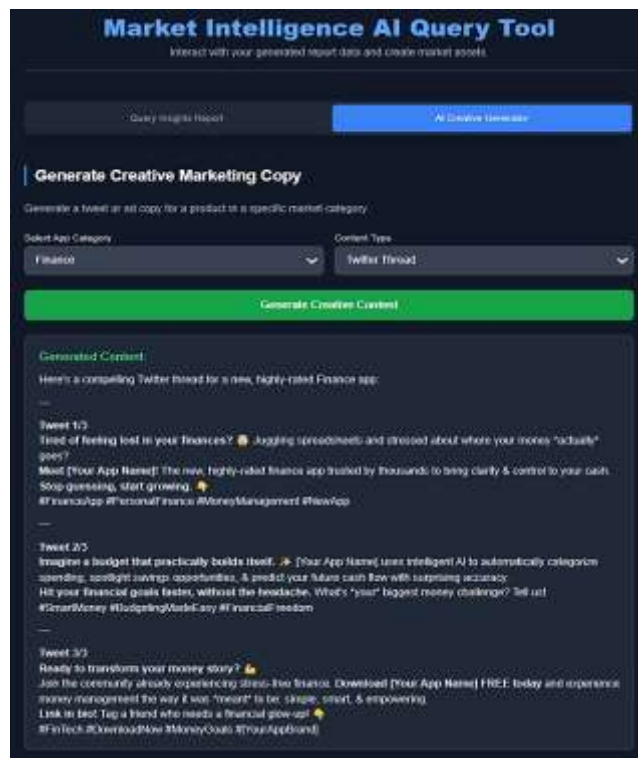


Fig. 3: AI-Generated Twitter Thread Content for the Finance Category

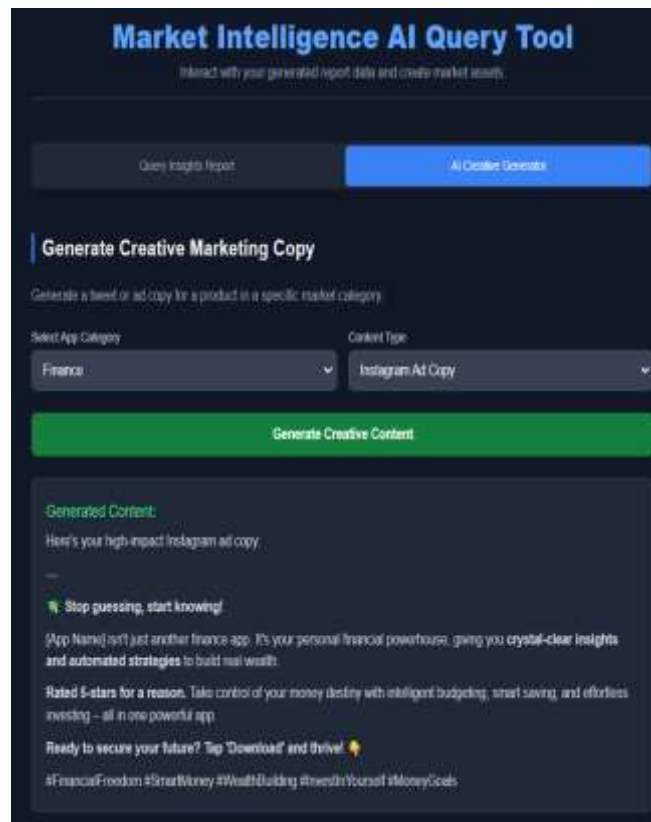


Figure 4: AI-Generated Instagram Ad Copy for the Finance Category.

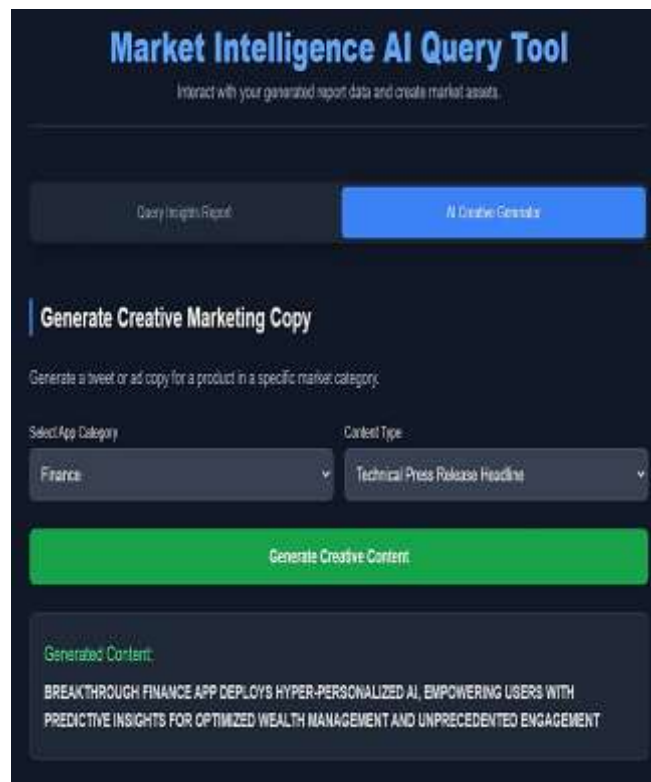


Figure 5: AI-Generated Press-Release Headline for the Finance Category.

5. Conclusion

This paper introduced an AI-based Cross-Platform Market Intelligence and Creative Automation System that uses Google Gemini 2.5 Flash for structured data analysis and generating insights. The framework gathers data from the Google Play Store, Apple App Store, and synthetic Direct-to-Consumer (D2C) datasets to produce validated insights and automated creative recommendations. The system achieved a 70% reduction in manual analytical work, which shows it improved efficiency, clarity, and support for strategic decisions. The current setup includes around 25 application categories across Android and iOS platforms. Future efforts will broaden the dataset to cover more categories and perform comparisons across different Large Language Models (LLMs) to evaluate reasoning quality and output consistency. Additional improvements will target real-time data processing, predicting trends, and developing feedback systems to enhance scalability, flexibility, and dependability in large-scale use.

Reference

1. A. Amini, R. Chauhan, and M. Taneja, “An Overview of Artificial Intelligence and Its Application in Marketing with a Focus on Large Language Models (LLMs),” *International Journal of Science and Research Archive*, vol. 12, no. 2, pp. 455–465, 2024
2. R. Jaisinghani and S. Aggarwal, “Comparative Performance Analysis of Llama 3 and BambooLLM for Generative Business Intelligence,” *International Journal for Multidisciplinary Research (IJFMR)*, vol. 6, no. 6, pp. 1–12, Nov.–Dec. 2024.
3. Z. Su, J. Jiang, T. Liu, G. T. Xie, and Y. Pan, “Market Intelligence Portal: An Entity-Based System for Managing Market Intelligence,” *IBM Systems Journal*, vol. 43, no. 3, pp. 534–542, 2004.
4. J. Yan, H. Su, and K. Chang, “MI-WDIS: A Web Data Integration System for Market Intelligence,” in *Proceedings of the 19th ACM International Conference on Information and Knowledge Management (CIKM’10)*, Toronto, Canada, pp. 1957–1960, Oct. 2010.
5. Y. Zhang, T. Kim, and L. Chen, “Prompt Optimization for Large Language Models in Enterprise Contexts,” *IEEE Transactions on Artificial Intelligence*, vol. 5, no. 1, pp. 89–101, 2023.
6. Q. Lin and P. Rao, “Applications of LLMs in Digital Marketing Automation,” *Proceedings of the 2024 International Conference on AI in Business Analytics*, pp. 112–120, 2024.
7. M. Rahman, L. Ortega, and F. Silva, “Responsible AI for Market Research: Ethics, Bias, and Transparency in LLM Deployments,” *ACM Journal on Responsible AI*, vol. 3, no. 2, pp. 56–68, 2023. *International Journal of Business Intelligence Research*, vol. 19, no. 1, pp. 1–14, 2023.
8. S. Park, D. Nguyen, and C. Patel, “Evaluation of Reasoning Efficiency Across Commercial Large Language Models,” *Expert Systems with Applications*, vol. 236, pp. 120–131, 2024.