Semantic Modeling and Text-to-SQL Pipelines in Snowflake for AI Retrieval

Raman Krishnaswami

Chief Information Officer

ABSTRACT: The paper will discuss the creation of Text-to-SQL pipelines and a semantic model in Snowflake to simplify the acquisition of AI-based data retrieval by businesses. The system can convert natural language queries into SQL queries using Snowflake Cortex, enabling users to access data without technical skills. Semantic modeling, hybrid search, and RAG are adopted together in the approach in a bid to enhance accuracy and relevance. We refer to pipeline design, workflow, performance, and governance, and demonstrate how the queries based on AI can be run quickly, securely, and in compliance. Real-life examples of financial, operational, and customer analytics indicate that timely decision-making and reduced manual labor can be achieved. This article puts Snowflake into focus as a powerful solution with respect to AI-based enterprise analytics.

KEYWORDS: AI, SQL, Semantic Modeling, Pipeline, Snowflake

I. INTRODUCTION

Artificial intelligence (AI) has transformed how business is conducted and perceived over the past few years. Previously, business intelligence (BI) tools were employed to prepare reports and dashboards. However, these tools often required specialists to write complex SQL queries, and most business users struggled to extract insights from large databases.

To address this issue, scholars invented the Text-to-SQL systems, which are able to convert easy natural language queries into SQL queries automatically. As an example, a user can query and display the total revenue for the past month, and the system will provide the appropriate SQL code to retrieve the result. Research indicates that these systems have improved due to large training data sets like Spider, WikiSQL, and CoSQL [2]. These data have provided AI models with the ability to understand natural language and translate it into accurate SQL queries. Consequently, Text-to-SQL systems have increasingly become an important component of AI tools, making data more accessible to everyone, not just technical users.

Although the Text-to-SQL translation has been enhanced significantly, it remains difficult to apply in large enterprises. Big organizations normally have databases that are very large and intricate, containing numerous tables and columns. They also have their data distributed in various systems, and the queries may have complex joins and computations. Besides, no two database platforms share the same SQL format entirely. Indicatively, Snowflake, BigQuery, and Databricks have their special SQL features.

It makes it difficult to have one Text-to-SQL model that would perform well with all of them. The second problem is data governance and security. Businesses should regulate access to vulnerable data and ensure that personal information (PII) is secure. These security rules are not observed in many open Text-to-SQL systems. Due to this, organizations require an effective and robust system with the ability to integrate both Text-to-SQL translation and good governance, scale, and regulatory measures.

According to the latest studies, cloud data appliances like Amazon Redshift, Google BigQuery, Databricks, and Snowflake are beginning to directly execute enormous language models (LLMs) within SQL systems [10]. This means users can communicate with the database in plain English and receive accurate answers without writing any code.

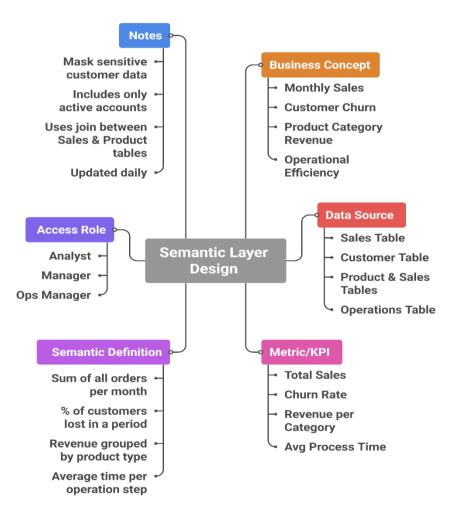
Snowflake is a unique product, due to its Cortex AI capabilities. Cortex enables seamless integration of AI with Snowflake's structured and unstructured data. This makes Text-to-SQL pipelines easier and more secure to use. Snowflake has the potential to become a leading platform for conversational analytics and AI-driven data retrieval in business.

II. SEMANTIC MODELING IN SNOWFLAKE

Semantic modeling is a very important part of building smart data systems in Snowflake. It helps to create a clear and semantic model, which is regarded as a highly significant aspect of constructing intelligent sets of data in Snowflake. It also aids in developing a unified and uniform method of explaining business information to ensure that it can be interpreted easily by humans and AI tools. Most companies keep data in various tables and databases with complex names and formats. The lack of an appropriate model makes the data unclear as to what the data represents and in what ways it should be interpreted.

This issue is addressed using a semantic model that creates a layer defining key business terms, metrics, and relationships. For example, instead of writing complex SQL joins, a user can simply type a business term such as 'monthly sales.' In Snowflake, developers can build this structure using metadata, catalogs, and data dictionaries that describe each dataset. This approach simplifies data usage, improves accuracy, and reduces cross-team confusion.

Semantic Layer Design Components



Semantic modeling is also beneficial in AI-based contexts. AI agents or Text-to-SQL systems rely on the meaning of words to answer questions. An appropriately developed semantic model provides AI systems with background on the relationship between data and the meaning of each metric. For example, when an AI system is asked to provide revenue growth over the last quarter, it must answer questions like: 'Where is the data on revenue?', 'What timeframe constitutes a quarter?', and 'How do you compute growth?'

These definitions are offered in the semantic model in Snowflake. It connects natural language queries to structured queries, making responses more precise and meaningful. It also supports governance by standardizing definitions across the company, so everyone shares the same understanding of terms such as customer churn or profit margin.

Semantic-aware chunking, metadata use, and optimized indexing have proven instrumental in improving the retrieval process of Retrieval-Augmented Generation (RAG) systems [7]. The same principles apply to Snowflake. When constructing a semantic layer, data can be subdivided into small, meaningful chunks with distinct metadata tags. This enables Cortex Search and vector-based retrieval in Snowflake to quickly find the required information.

Metadata helps prioritize results based on importance and context, improving AI responses and saving time in accessing the appropriate data. Semantic modeling in Snowflake thus enhances AI retrieval, making analytics faster, more accurate, and easier for business users to handle.

III. TEXT-TO-SOL PIPELINES

The primary component of the natural language—structured data connection is the text-to-SQL pipeline. These pipelines accept a simple text query from a user and convert it into an SQL query, which can be executed in an SQL database such as Snowflake. This aims to ensure that the data is accessible to all, even to people who do not know how to write SQL. A typical workflow would begin with a user posing a query like, "Get total revenue last year."

The system will then interpret the meaning of the question, identify where the data is stored in tables and columns, generate an SQL statement, and provide the output. The advantage of creating such a pipeline within Snowflake is that it provides an opportunity to connect AI models directly to the company's real-time data. Users can also obtain conversational responses using automatic SQL generation via Snowflake Cortex Analyst, without leaving Snowflake.

It is not a smooth process, though. Large enterprise systems often contain databases with hundreds or even thousands of tables that have long and complicated names. It may even be difficult for a human to find the correct table or column to use. Research such as ReFoRCE suggests that Text-to-SQL systems should account for schema complexity and various SQL versions [1]

Text-to-SQL Pipeline



ReFoRCE applies intelligent techniques such as pattern-based table grouping and self-refinement to ensure that SQL queries remain valid even when using more than one platform, such as Snowflake or BigQuery. These methods can influence how Cortex Analyst enhances the precision of large and complex databases. The model understands the structure and logic of the database instead of making errors such as combining inappropriate tables or applying excessive filtration.

Another useful concept is the work of DAIL-SQL, which aims to increase the efficiency and performance of tokens [3]. This involves ensuring that the model consumes fewer tokens (or words) when comprehending and constructing queries, making it faster and cheaper to run. Snowflake Cortex Analyst has the opportunity to use the same concept to optimize its input query processing and internal context. This helps manage multiple users and large workloads without complications.

According to research surveys, Text-to-SQL models can be improved through prompt engineering and fine-tuning strategies [4]. This includes example-based and feedback-driven teaching to enable the model to respond to various types of questions. Through these methods, the Text-to-SQL pipeline in Snowflake improves over time. It observes user behavior, expands its understanding of business terms, and eliminates the need for manual SQL writing. This is why AI-driven analytics in Snowflake is easy, fast, and precise for all users.

IV. PIPELINE ARCHITECTURE AND WORKFLOW

Snowflake text to SQL platform is connected through the pipeline, where all components of the Text-to-SQL and AI retrieval process are attached. It enables the movement of data across ingestion and query implementation. The design processes start with data collection and end with the users receiving answers in natural language. The first process is the ingestion of data, in which the data from various sources is loaded into Snowflake.

Which search configuration should be used?



This may comprise application data, file data, or API data. The data can be brought in real-time automatically with the help of such tools as Snowpipe and Step Functions. The data can then be cleaned after being ingested, structured, and enriched with semantic data to ensure that AI models can comprehend it more effectively. When the data is prepared, it is connected to Snowflake Cortex, where natural language processing and Text-to-SQL translation take place. Users are then able to pose questions in plain language, which is then converted into an SQL query by the pipeline, and the queries are executed, and the results are instantly returned to the users.

The main characteristic of this architecture is that it can support hybrid search. This involves matching keyword search (which searches for an exact match of the words being searched) and vector search (which searches for a meaning and

context). Research indicates that applying the two techniques in combination with each other enhances the clarity and applicability of search results [7].

An illustration is given by asking, "Show revenue by customer type." A keyword search might include tables containing both revenue and customer in their names, whereas a vector search would interpret customer type as a specific category or field within the data. The hybrid search in Snowflake Cortex assists AI models in finding the most appropriate information for Retrieval-Augmented Generation (RAG) searches. This ensures that responses are more accurate and reduces confusion in complex questions.

The system must be evaluated and monitored to ensure that it works well and delivers appropriate results. The AutoNuggetizer model is one way to check RAG pipeline performance automatically [6]. It uses bits of information to verify the presence of correct facts in answers. This concept can be applied to analyze a hybrid pipeline offered by Snowflake to detect areas of weakness, giving developers an opportunity to improve data accuracy over time.

Incremental updates are also supported in the pipeline, so new data or modifications can be made without recreating everything. This keeps AI responses current. By assembling these steps—data ingestion, semantic enrichment, Text-to-SQL translation, hybrid search, and automatic evaluation—Snowflake provides a complete and effective workflow for AI-driven data retrieval and decision-making.

V. GOVERNANCE AND SECURITY

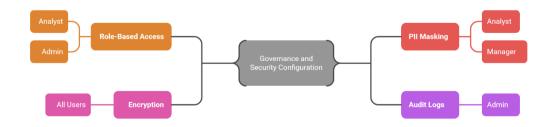
Any data and AI system is highly sensitive; therefore, governance and security are of utmost importance when handling sensitive business information. At Snowflake, governance means ensuring that the correct users can review and use the correct data, while security ensures there is no misuse or leakage of data. Both guarantee that the system complies with laws and company policies.

Role-based access control (RBAC) helps Snowflake achieve this by allowing access based on job roles. For example, a data analyst can only read data, whereas a database administrator can update tables or manage permissions. This hierarchy eliminates unauthorized access and provides each user with sufficient permission to perform their job.

Snowflake can also prevent the theft of personally identifiable information (PII), such as names, contact details, or financial records. Sensitive columns may be masked so that authorized users can access actual data, while others see hidden or partial values. For instance, an email address may appear as xxxxx@gmail.com instead of the full text.

Snowflake also supports tokenization and encryption; therefore, even if someone has access to the database, they cannot read protected information without the appropriate keys. These features are particularly useful in constructing Text-to-SQL pipelines, as user prompts or query outputs may contain sensitive information. Such exposure can be prevented with proper data masking and encryption during AI processing.

Governance and Security Configuration



Data lineage and auditing are also part of governance, monitoring how data flows and evolves over time. Snowflake automatically records activities such as role changes, queries, and data uploads. These logs help companies comply with standards like GDPR, HIPAA, and SOC 2. This is especially important when handling enterprise-level AI systems, which

require explanations of how data was used or modified. Snowflake's secure sharing feature enables teams to share data and information across departments or even organizations without creating copies. Shared data is always subject to the same access and masking principles.

Snowflake follows strong governance and security practices that ensure AI and Text-to-SQL workflows are trustworthy. These practices safeguard information, protect user privacy, and support compliance with industry regulations. Such thoughtful design allows organizations to use AI products like Snowflake Cortex Analyst without worrying about data safety, compliance, or manageability.

VI. PERFORMANCE AND SCALABILITY

Performance and scalability are highly important to ensure that AI and Text-to-SQL systems in Snowflake work concurrently, even when the data is extremely large or when many users are working simultaneously. Good performance means the system can produce answers quickly and accurately. Scalability refers to its ability to process additional information and traffic without slowing down.

This is achieved because storage and compute layers are separated in Snowflake, allowing users to scale one without scaling the other. For example, if an organization has large data tables but few users, it can maintain low costs. Snowflake can instantly add more computing power when demand is high, keeping the system fast.

Studies such as Arctic-Text2SQL-R1 indicate that reinforcement learning can improve SQL execution performance in AI models [9]. The model learns from past query results and gradually improves its ability to select the most efficient SQL statements. This concept could be applied in Snowflake Cortex to enhance the speed and accuracy of Text-to-SQL generation.

For instance, when a query is very slow or returns excessive data, the system can learn to reformulate a better query in the future. This personalized improvement strategy reduces costs and improves response time. It also prevents common SQL errors, such as unnecessary joins or missing filters, which slow down query performance.

Chain-of-Thought (CoT) fine-tuning is another concept used in ExCoT [8]. It enables the model to show its reasoning in steps before providing the final SQL output. This stepwise approach helps identify and correct errors early. Snowflake Cortex can combine this with semantic refinement, where the model checks the meaning of both the query and the data schema twice before forming the SQL query.

This minimizes query failures and ensures the SQL aligns with user intent. If an error still occurs, the system can trigger error-repair functions, such as re-running the query with modified parameters or providing a helpful message to the user.

These concepts make Snowflake and its Text-to-SQL pipeline successful and effective. Chain-of-Thought fine-tuning offers better reasoning and accuracy, while reinforcement learning provides speed and efficiency. Combined with Snowflake's scalable infrastructure, these methods enable companies to process large and complex datasets efficiently, respond quickly, and maintain high performance even during peak demand.

VII. ENTERPRISE USE CASES

The Text-to-SQL and AI-powered data systems offered by Snowflake can be highly useful in most enterprise cases. They simplify the process by allowing business users to pose questions in simple language and receive immediate responses without requiring extensive technical knowledge. Studies indicate similar systems are applied in areas such as healthcare, finance, and education, where fast access to accurate data enhances decision-making and effectiveness [2]. These ideas are applicable in enterprise domains, including financial analysis, business operations, and customer insights.

Text-to-SQL provides rapid information retrieval in finance, assisting analysts and managers in locating details on revenues, expenses, and performance trends. Instead of writing long SQL queries, users simply ask questions, and the system automatically generates and executes the appropriate SQL—for example, "Show the profit growth by region in the last quarter."

This saves time and reduces errors in manual query writing. It also helps non-technical teams take action using data. Albased insights can further identify suspicious transactions or risk patterns, enabling finance teams to improve compliance and detect fraud.

Snowflake Cortex can support monitoring and reporting activities across various business operations. For instance, a supply chain manager might request "supply shortage by supplier this month" or "average movement time by region." The Text-to-SQL engine translates these questions into queries within seconds and provides answers in an understandable format. This helps operations teams respond quickly to changes and streamline processes. AI models can also uncover hidden patterns, such as delays or cost inefficiencies, that human analysts might miss.

To gain customer insights, AI-assisted SQL generation enables marketing and customer care teams to understand customer behavior without technical assistance. They can pose natural questions such as, "Which group of customers was most satisfied during the last quarter?" or "What products have the strongest demand among new customers?" The system integrates and relates data from multiple sources to deliver comprehensive and accurate solutions. This helps teams enhance services, personalize customer experiences, and improve predictions.

These enterprise use cases show that Snowflake's Text-to-SQL and AI integration expands access to data, making it available to everyone. It eliminates the technological gap between information and decision-making. Raw data can now deliver faster, more precise, and intelligent responses for organizations, regardless of their size or business function.

VIII. FUTURE DIRECTIONS

Text-to-SQL and AI systems at Snowflake have a highly promising future, with numerous opportunities for development and advancement. The existing models work well with single questions, though they still struggle with multi-turn conversations and generalization across domains [4]. Multi-turn dialogue means the system can retain and comprehend a sequence of related queries rather than treat each query as separate.

For example, when a user asks, "Show total sales in 2023," and then says, "Now show only in Asia," the system must recognize that the second query depends on the first. Such conversational interaction would make Snowflake Cortex more natural and easier to use. Domain generalization means the system can adapt to various types of data and sectors, such as finance, healthcare, or manufacturing, without extensive retraining. Enhancing these two aspects will make AI-based SQL tools more adaptable and useful in most business settings.

Emerging research also emphasizes advanced techniques such as multihop retrieval, metadata utilization, and agentic RAG systems [7]. Multihop retrieval allows the model to correlate information across multiple sources or tables to answer complex questions. For instance, a user might ask, "What are the suppliers who caused project delays in the highest cost-overrun projects?"

Answering this requires data from several tables and a series of logical steps. Snowflake Cortex can execute such deep, connected queries more accurately by enabling multihop retrieval. Metadata utilization expands information about tables, field definitions, and data relationships to improve query comprehension. This helps the model select appropriate data sources and generate additional correct SQL queries.

Agentic RAG systems further enhance this by giving the AI model greater autonomy. Instead of merely retrieving and answering, the system acts like an intelligent agent, planning the best way to locate, combine, and validate information before presenting a response. This approach could enable Snowflake Cortex to develop fully self-regulating semantic AI pipelines that adapt based on user feedback and system performance.

The future of Snowflake's AI systems lies in enabling more natural dialogue, smarter retrieval techniques, and self-learning pipelines. These enhancements will help organizations make better use of their data, reduce manual effort, and simplify access to business insights through natural language.

IX. CONCLUSION

Snowflake Text-to-SQL pipelines with semantic modeling offer the simplest and most effective method for an enterprise to use AI to retrieve data. The system provides answers to non-technical users quickly and accurately by converting natural

language into SQL. Hybrid search and RAG improve relevance, while governance and security features ensure control over sensitive data.

Performance and scalability aspects ensure that large workloads run efficiently. As seen in real-life applications in finance, operations, and customer analytics, reduced manual work and shorter decision-making times are achieved. Snowflake Cortex, with its semantic-aware pipelines, provides a strong foundation for AI-driven enterprise analytics and future developments, such as multi-turn dialogues and advanced retrieval.

REFERENCES

- [1] Deng, M., Ramachandran, A., Xu, C., Hu, L., Yao, Z., Datta, A., & Zhang, H. (2025). ReFoRCE: A Text-to-SQL Agent with Self-Refinement, Format Restriction, and Column Exploration. arXiv (Cornell University). https://doi.org/10.48550/arxiv.2502.00675
- [2] Singh, A., Shetty, A., Ehtesham, A., Kumar, S., & Khoei, T. T. (2024). A survey of Large Language Model-Based Generative AI for Text-to-SQL: Benchmarks, applications, use cases, and challenges. arXiv (Cornell University). https://doi.org/10.48550/arxiv.2412.05208
- [3] Gao, D., Wang, H., Li, Y., Sun, X., Qian, Y., Ding, B., & Zhou, J. (2024). Text-to-SQL empowered by large language models: a benchmark evaluation. Proceedings of the VLDB Endowment, 17(5), 1132–1145. https://doi.org/10.14778/3641204.3641221
- [4] Shi, L., Tang, Z., School of Computer Science, Peking University, China, SINGDATA CLOUD PTE. LTD, USA, SINGDATA CLOUD PTE. LTD, China, Yang, Z., & School of Computer Science, Peking University, China. (2025). A survey on employing large language models for Text-to-SQL tasks. In ACM Comput. Surv. (Vols. 1–1, Issue 1, pp. 1–36) [Journal-article]. https://doi.org/10.1145/3737873
- [5] Singh, A., Shetty, A., Ehtesham, A., Kumar, S., & Khoei, T. T. (2025). A Survey of Large Language Model-Based Generative AI for Text-to-SQL: Benchmarks, Applications, Use Cases, and Challenges. A Survey of Large Language Model-Based Generative AI for Text-to-SQL: Benchmarks, Applications, Use Cases, and Challenges. https://arxiv.org/pdf/2412.05208
- [6] Pradeep, R., University of Waterloo, Thakur, N., University of Waterloo, Upadhyay, S., University of Waterloo, Campos, D., Snowflake, Craswell, N., Microsoft, Lin, J., & University of Waterloo. (2025). [The Great Nugget Recall: Automating Fact Extraction and RAG Evaluation with Large Language Models]. The Great Nugget Recall: Automating Fact Extraction and RAG Evaluation With Large Language Models. https://arxiv.org/html/2504.15068v1
- [7] James, A., Trovati, M., & Bolton, S. (2025). Retrieval-Augmented generation to generate knowledge assets and creation of action drivers. Applied Sciences, 15(11), 6247. https://doi.org/10.3390/app15116247
- [8] Zhai, B., Xu, C., He, Y., Yao, Z., & Snowflake Inc. (2025). ExCoT: Optimizing Reasoning for Text-to-SQL with Execution Feedback. In arXiv. https://arxiv.org/abs/2503.19988v1
- [9] Yao, Z., Sun, G., Borchmann, L., Shen, Z., Deng, M., Zhai, B., Zhang, H., Li, A., Snowflake AI Research, University of Maryland, College Park, & University of California, San Diego. (2025). Arctic-Text2SQL-R1: Simple rewards, Strong reasoning in Text-to-SQL [Report]. https://huggingface.co/Snowflake/Arctic-Text2SQL-R1-7B
- [10] Akillioglu, K., Chakraborty, A., Voruganti, S., & Özsu, M. T. (2025). Research challenges in relational database management systems for LLM queries. VLDB 2025 Workshop: Applied AI for Database Systems and Applications (AIDB 2025). https://arxiv.org/html/2508.20912v1