

ELEVATING ANALYTICS



Knowledge Graph Infused Natural Language to SQL with Domain Adapted LLM for Patient Data

M⊆KESSON Compile[™]

Bharath Bommakanti, Manager – Data Science Sameer Saraf, Principal Data Scientist

Agenda

• Introduction

 $\,\circ\,$ Current challenges with data democratization

 $\,\circ\,$ Challenges with traditional LLMs / RAG

- McKesson Compile's NL to SQL approach
- NL to SQL showcase
- Results and Metrics
- Conclusion



Current Challenges in Data Democratization

• Expertise limitations

 \odot Analyst required to query data using SQL

- Domain-specific complexities
 - \odot Specialized nomenclature
 - \circ Multiple diverse coding systems
 - $\,\circ\,$ Dilution of context from business user to analyst to SQL query
 - $\,\circ\,$ Diverse healthcare elements hosted on complex database models
- Limitations of traditional LLMs / RAG
 - $\,\circ\,$ Complex schematic understanding
 - $\,\circ\,$ Non-standard table and schema nomenclature
 - $\,\circ\,$ Low-precision retrieval due to lack of semantic knowledge

Challenges with Traditional LLMs / RAG

- Complex schematic understanding
- Non-standard table and schema nomenclature
- Low-precision/ Low-recall retrieval due to lack of semantic knowledge
- Lengthy prompts to guide through the intricacies of healthcare data analytics
 Financial implications
- Ground-truth or the retrieved context do not capture the relationships in the data
- Co-pilots and task specific LLMs (NL to SQL) lack the required context of the domain and data elements





Data | Data Understanding

Schematic data

- Schema table column relationships
- Table table relationships

Semantic data

- Description of the node types
- Description of the node values
- Understanding the data types and data distribution
- Semantic tagging using the diverse node values and aliases
- Business logic



Embedded Knowledge Graph (KG)

- Semantic knowledge embedding with schematic relationships
- Relationship prediction using semantic tagging
 - Useful for relevant sub-graph / community extraction
- High-precision retrieval and generation

 Constrained sub-graph retrieval from KG
 Less guidance required for complex
 - problems through structured prompting
 - Fine-tuned LLM leveraging gisting approach for output generation





Query Synthesis & Continual Learning

- Agentic SQL generation
 - Iterative retrieval agent to ensure completeness in the answer through conditional retrieval
 - $\,\circ\,$ SQL Validation engine
 - $\,\circ\,$ Agent that ensures syntactical accuracy
 - $\,\circ\,$ Contextual relevance and completeness
 - $\odot\,$ Fine-tuned LLM leveraging gisting approach for output generation
- Feedback loop
 - Quick feedback through bucketized granular options
 - $\,\circ\,$ Specific feedback through free flow text
 - Maintaining output feedback pairs for further fine-tuning of the LLM through RLHF



NL to SQL showcase: Example 1

Q: Extract Top 10 HCPs (based on patient volume) that are diagnosing patients with Lung Cancer disease (Diagnosis code: C34) for the year of 2023



Figure: Visualization of SQL query generated by NL to SQL

0000 Pmsa

LEADING EDGE INSIGHTS, COASTAL VIEWS

NL to SQL showcase: Example 1 Output

Q: Extract Top 10 HCPs (based on patient volume) that are diagnosing patients with Lung Cancer disease (Diagnosis code: C34) for the year of 2023

НСР	PATIENT_COUNT
HCP 1	3,123
HCP 2	2,513
НСР 3	2,441
HCP 4	2,158
HCP 5	1,578
НСР 6	1,555
HCP 7	1,503
HCP 8	1,496
НСР 9	1,460
HCP 10	1,431

Note: Data shown is for illustrative purposes only

LEADING EDGE INSIGHTS, COASTAL VIEWS



NL to SQL showcase: Example 2

Q: Provide feasibility estimates (patient counts) by year from 2020 to 2024 in 25 year age group buckets for patients aged ≥18 years with a non-inpatient, non-diagnostic* claim and having an hMPV diagnosis.

*Non-diagnostic claims are claims without any procedures for Radiology, Venipuncture Pathology and Laboratory.



NL to SQL showcase: Example 2 Outputs

METRIC	SERVICE_YEAR	PATIENT_COUNT
Age Group: 18-49 years	2020	149
Age Group: 18-49 years	2021	162
Age Group: 18-49 years	2022	318
Age Group: 18-49 years	2023	443
Age Group: 18-49 years	2024	368
Age Group: 50-74 years	2020	299
Age Group: 50-74 years	2021	169
Age Group: 50-74 years	2022	379
Age Group: 50-74 years	2023	623
Age Group: 50-74 years	2024	479
Age Group: ≥75 years	2020	334
Age Group: ≥75 years	2021	106
Age Group: ≥75 years	2022	295
Age Group: ≥75 years	2023	466
Age Group: ≥75 years	2024	379



Pmsa

Note: Data shown is for illustrative purposes only

LEADING EDGE INSIGHTS, COASTAL VIEWS

Results Demonstrate a Reduction in 3 Key Areas

₽ 95%

Non-executable queries

- Schema misidentification & incorrect syntax
- Attributed to
 - SQL Validation agent
 - KG relationships

↓ 5x

Hallucinations

- Semantic / contextual relevance & completeness
- Attributed to
 - $\circ\,$ Agent orchestrator
 - Conditional retrieval agent

↓15x

Token count

- Attributed to
 - Fine-tuned LLM with gisting
 - Condensed context through KG relationships



Conclusions and Future Directions

- Real-world applications
 - $\circ~$ Clinical decision intelligence
 - Analytical applications like feasibility analysis, cohort analysis, etc.
 - $\circ~$ Advanced provider network analytics
 - Operational improvements
- Improved precision and recall in data democratization
- Cost-efficient GenAl application
- Strong foundation for NL to insights





Thank you for your attention

Questions? Visit McKesson Compile at booth 212