



TEXT-TO-SQL





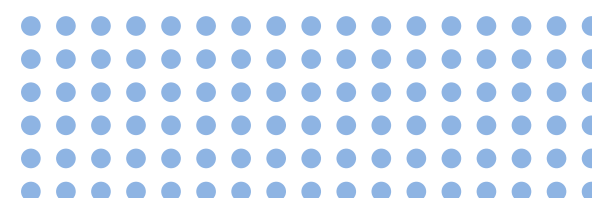
Khái niệm

Text-to-SQL là một công nghệ cho phép chuyển đổi câu hỏi hoặc yêu cầu bằng ngôn ngữ tự nhiên (tiếng Anh, tiếng Việt,...) thành câu lệnh SQL (Structured Query Language) để tương tác với cơ sở dữ liệu (CSDL).

Ví dụ:

- Câu hỏi tự nhiên: "Hiển thị danh sách nhân viên có lương cao hơn 10 triệu"
- SQL tương ứng:

```
SELECT * FROM employees WHERE salary > 10000000;
```



Tại sao cần Text-to-SQL ?



1. Giảm rào cản kỹ thuật:

- Người dùng không cần biết SQL vẫn có thể truy vấn dữ liệu bằng ngôn ngữ tự nhiên.
- Hữu ích cho người dùng phổ thông, nhà quản lý, hoặc nhân viên không chuyên về công nghệ.

2. Tăng tốc độ phát triển:

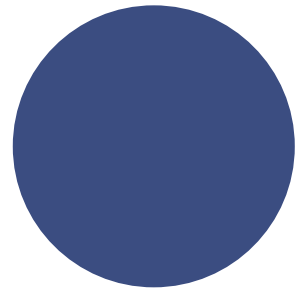
- Lập trình viên tiết kiệm thời gian viết SQL thủ công, đặc biệt với các truy vấn phức tạp.

3. Ứng dụng trong AI và chatbot:

- Hệ thống hỏi đáp tự động (ví dụ: chatbot hỗ trợ khách hàng) có thể trả lời câu hỏi liên quan đến CSDL mà không cần can thiệp thủ công.

4. Hỗ trợ phân tích dữ liệu nhanh:

- Chuyển đổi ngay lập tức yêu cầu phân tích (ví dụ: "Doanh thu tháng này so với tháng trước?") thành SQL để lấy kết quả.



Công nghệ đằng sau

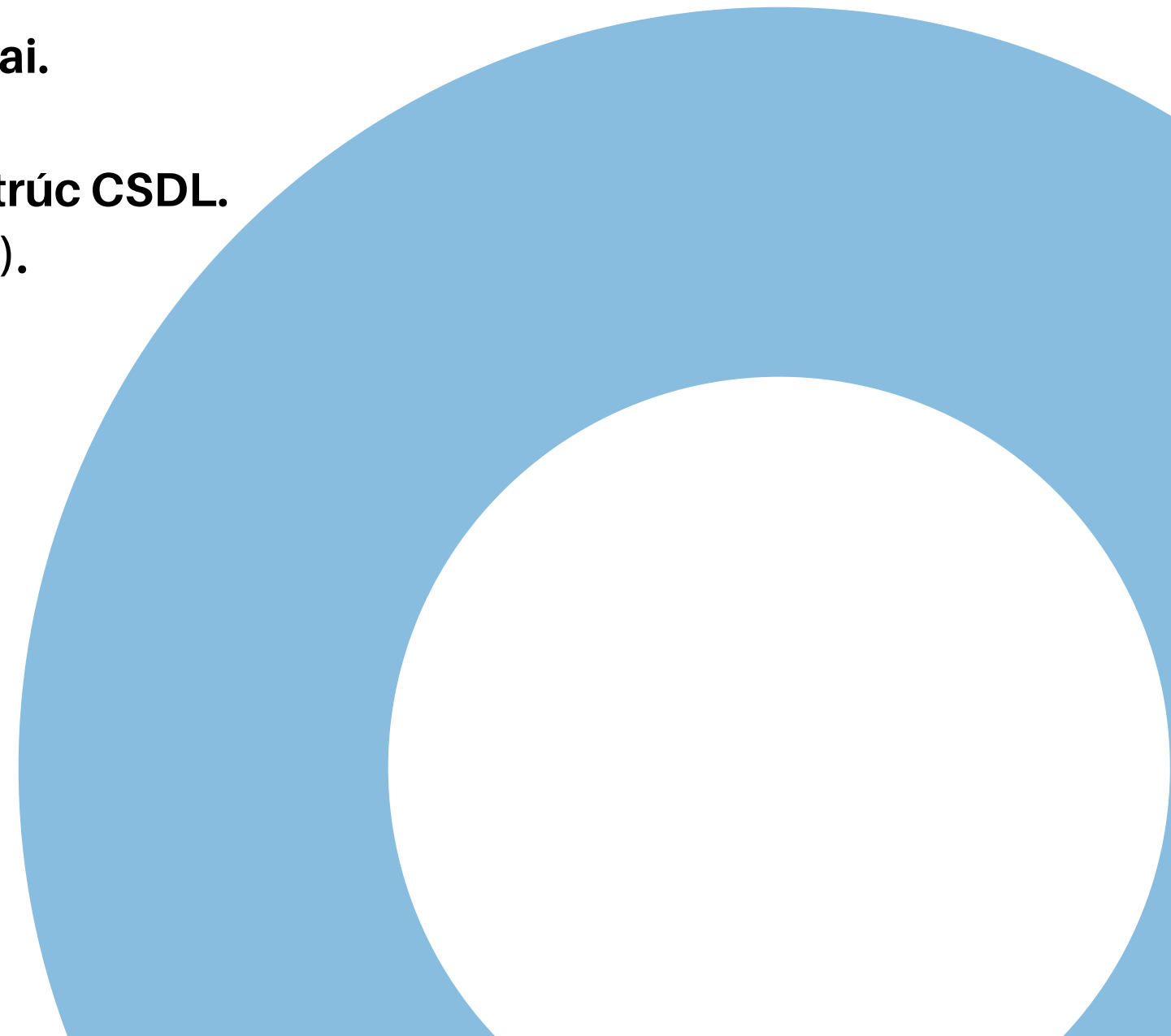
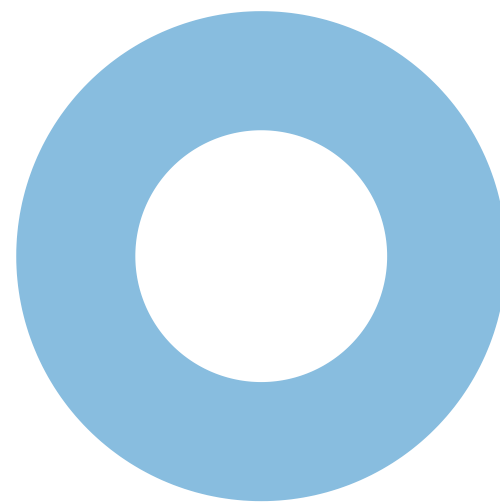
Text-to-SQL thường sử dụng AI/Machine Learning (như mô hình NLP: BERT, GPT, T5) để hiểu ngữ nghĩa câu hỏi và ánh xạ vào cấu trúc CSDL (bảng, cột, khóa...).

Ví dụ các công cụ:

- **Google's TAPAS, Microsoft's Semantic Machines**
- **Các thư viện Python như LangChain, SQLGlue hỗ trợ triển khai.**

Hạn chế

- **Độ chính xác phụ thuộc vào độ phức tạp của câu hỏi và cấu trúc CSDL.**
- **Cần huấn luyện mô hình trên dữ liệu đặc thù (schema cụ thể).**





Ứng dụng thực tế

a. Trong Doanh Nghiệp

- Báo cáo tự động: Nhân viên marketing hỏi "Tỉ lệ chuyển đổi từ quảng cáo Facebook trong Q1?" → Hệ thống tự sinh SQL và trả kết quả.
- Hỗ trợ khách hàng: Chatbot trả lời câu hỏi như "Sản phẩm nào còn hàng ở chi nhánh Hà Nội?" bằng cách truy vấn CSDL.

b. Trong Phân Tích Dữ Liệu (BI)

- Tableau/Power BI: Dùng voice query hoặc nhập text để vẽ biểu đồ thay vì viết SQL thủ công.
- Google BigQuery/NLP Integration: Google cho phép hỏi tự nhiên để phân tích dữ liệu lớn.

c. Hệ Thống Nội Bộ

- Quản lý kho: "Hiển thị mặt hàng sắp hết hạn trong 30 ngày tới" → Tự động query từ CSDL kho.
- HR: "Liệt kê nhân viên nghỉ việc năm 2023" → Truy vấn từ database nhân sự.

Cách Text-to-SQL hoạt động

Text-to-SQL kết hợp Xử lý ngôn ngữ tự nhiên (NLP) và Hiểu cấu trúc cơ sở dữ liệu (Database Schema) để sinh ra câu lệnh chính xác. Quy trình gồm các bước:

1. Phân tích cú pháp (Parsing): Tách từ khóa, xác định ý định (ví dụ: "tìm", "lọc", "tính tổng").
2. Ánh xạ ngữ nghĩa (Semantic Mapping): Liên kết từ ngữ tự nhiên với tên bảng/cột trong CSDL (ví dụ: "nhân viên" → bảng employees).
3. Sinh SQL (SQL Generation): Tạo câu lệnh hoàn chỉnh, kiểm tra tính hợp lệ.



Ví dụ:

- Câu hỏi: "Cho tôi top 3 nhân viên bán hàng có doanh số cao nhất tháng 5"
SQL sinh ra:

```
SELECT employee_name, sales
FROM sales_records
WHERE MONTH(date) = 5
ORDER BY sales DESC
LIMIT 3;
```

Ưu nhược điểm của Text-to-SQL

1. Ưu điểm

a. Truy cập dữ liệu dễ dàng

- Không cần biết SQL: Người dùng chỉ cần nhập yêu cầu bằng ngôn ngữ tự nhiên (tiếng Anh, tiếng Việt) thay vì viết câu lệnh phức tạp.
- Phù hợp với người không chuyên (nhà quản lý, nhân viên kinh doanh, khách hàng).

b. Tiết kiệm thời gian

- Tự động hóa truy vấn: Giảm thời gian viết SQL thủ công, đặc biệt với CSDL lớn hoặc yêu cầu phức tạp.
- Tốc độ phân tích nhanh: Ví dụ, thay vì mất 10 phút viết SQL, hệ thống trả kết quả trong vài giây.

c. Ứng dụng linh hoạt

- Tích hợp vào chatbot (hỗ trợ khách hàng, tra cứu thông tin).
- Hỗ trợ Business Intelligence (BI): Dùng trong Tableau, Power BI để tạo báo cáo bằng giọng nói hoặc văn bản.
- Tương tác với nhiều loại CSDL (MySQL, PostgreSQL, BigQuery...).

d. Giảm lỗi cú pháp SQL

- Mô hình AI tự động sinh câu lệnh chuẩn, tránh sai sót do viết tay (ví dụ: thiếu WHERE, nhầm tên bảng).

2. Nhược điểm

a. Phụ thuộc vào chất lượng dữ liệu huấn luyện

- Nếu mô hình chưa được huấn luyện trên schema (cấu trúc bảng) cụ thể, nó có thể sinh SQL sai.
- Ví dụ: Câu hỏi "Tìm sản phẩm bán chạy" → Cần biết bảng nào chứa doanh số (sales hay orders?).

b. Khó xử lý câu hỏi phức tạp

- Câu hỏi đa tầng:
 - "So sánh doanh thu tháng này với tháng trước, rồi tính phần trăm thay đổi" → Cần nhiều bước SQL (JOIN, GROUP BY, subquery).
- Ngữ nghĩa mơ hồ:
 - "Hiển thị đơn hàng gần đây" → "Gần đây" là 7 ngày, 1 tháng hay 1 năm?

c. Yêu cầu hiểu biết về CSDL

- Người dùng cần biết từ khóa quan trọng trong CSDL (ví dụ: tên bảng users thay vì "khách hàng").
- Schema thay đổi → Mô hình có thể hoạt động sai.

d. Hiệu suất không ổn định

- Mô hình nhỏ (ví dụ: SQLNet) dễ mắc lỗi logic.
- Mô hình lớn (GPT-4, TAPAS) tốn tài nguyên và chi phí.

e. Bảo mật và quyền truy cập

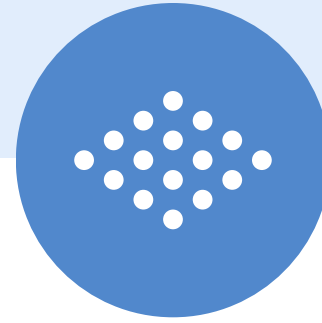
- Nếu không kiểm soát, người dùng có thể truy vấn dữ liệu nhạy cảm (ví dụ: lương nhân viên).
- Cần tích hợp RBAC (Role-Based Access Control) để giới hạn quyền.

>>> Giải pháp khắc phục nhược điểm

1. Fine-tuning mô hình trên CSDL riêng để tăng độ chính xác.
2. Kết hợp RAG (Retrieval-Augmented Generation): Truy xuất thông tin schema trước khi sinh SQL.
3. Kiểm tra SQL tự động bằng thư viện như SQLGlot hoặc sqlparse.
4. Giới hạn quyền truy cập theo vai trò người dùng.

● Kết luận:

- Dùng Text-to-SQL khi: Cần truy vấn nhanh, người dùng không rành SQL, ứng dụng chatbot hoặc BI.
- Hạn chế dùng khi: Truy vấn quá phức tạp hoặc yêu cầu độ chính xác tuyệt đối.



**THANK
YOU!**

