

# BAB I

## PENDAHULUAN

### 1.1 Latar Belakang

Pengolahan Data atau *Data Management* saat ini menjadi salah satu aspek utama yang krusial dalam perkembangan teknologi. Dalam dunia digital yang juga kerap berkembang, aktivitas manusia juga menghasilkan jumlah data yang besar sehingga kebutuhan untuk melakukan pengolahan data menjadi kunci untuk kesuksesan dalam hasil yang akan disampaikan kepada *end-user* [1]. Data yang terus berkembang ini disebut sebagai *Big Data* yang merupakan sebuah wadah untuk menyatukan ratusan bahkan ribuan sumber data dan menjadi kesatuan yang dapat dianalisis[2].

Salah satu tantangan besar dalam proses pengolahan *big data* yang jumlahnya sangat besar tersebut adalah bagaimana cara mengolah dan mendapatkan informasi dari banyaknya data yang ada. Dalam hal mendapatkan data, salah satu bahasa pemrograman utama yang digunakan untuk mengolah data tersebut adalah *Structured Query Language* atau yang dikenal dengan SQL. Pada awalnya, SQL hanya digunakan untuk mengolah *Relational Database Management Systems* atau RDBMS. Seiring berjalannya waktu dengan perkembangan *database* sendiri, telah terbukti bahwa SQL dapat menjadi *query language* yang efektif untuk mengolah jumlah data yang sangat banyak dari sistem *Big Data* tersebut[3].

SQL, dengan kemampuannya dalam merumuskan pertanyaan terstruktur, memainkan peran sentral dalam menggali informasi berharga dari Big Data. Melalui perintah SQL yang tepat, analisis data dapat menyaring, menggabungkan, dan mengagregasi data dari berbagai sumber, mengungkap pola, tren, dan korelasi yang tersembunyi[4]. Informasi yang dihasilkan dari analisis SQL ini kemudian dapat digunakan untuk mendukung pengambilan keputusan strategis, meningkatkan efisiensi operasional, dan mendorong inovasi dalam berbagai sektor, mulai dari bisnis hingga penelitian ilmiah[5].

Namun, penggunaan SQL secara efektif seringkali membutuhkan pemahaman mendalam tentang sintaks dan struktur database. Hal ini dapat menjadi hambatan bagi individu yang tidak memiliki latar belakang teknis yang kuat. Disinilah Natural Language Processing (NLP) hadir sebagai salah satu Solusi untuk membantu awam[6]. Model bahasa besar (Large Language Models/LLM), seperti Llama 2, memiliki kemampuan untuk memahami dan memproses bahasa manusia[7]. Dengan memanfaatkan teknik fine-tuning dan metode seperti QLORA, LLM dapat dilatih untuk menerjemahkan perintah bahasa alami menjadi query SQL yang valid[8].

Penerapan LLM dalam konteks text-to-SQL ini membuka peluang baru dalam kemudahan akses terhadap Big Data. Pengguna non-teknis dapat berinteraksi dengan data menggunakan bahasa sehari-hari, tanpa perlu mempelajari sintaks SQL yang rumit. Hal ini dapat mempercepat proses analisis data, meningkatkan produktivitas, dan mendorong kolaborasi lintas disiplin ilmu. Penelitian mengenai fine-tuning model bahasa Llama 2 untuk text-to-SQL menggunakan metode QLORA menjadi langkah penting dalam mewujudkan potensi LLM dalam transformasi cara kita berinteraksi dengan Big Data[9].

Salah satu manfaat penggunaan LLM kedalam *chatbot* terhadap pengolahan data dan membaca data adalah dengan membuat sebuah *chatbot* yang mampu memberikan informasi mendetail mengenai *dataset* yang dimiliki sehingga tanpa adanya keahlian dalam pengolahan data, *user* tetap dapat mendapatkan informasi mengenai data yang dimiliki[10]. Hal ini dapat dilakukan dengan sebuah *chatbot* yang dapat menginterpretasikan sebuah kalimat menjadi sebuah *query sql*. Dengan *chatbot* tersebut, *user* dapat dengan mudah menyampaikan sebuah gambaran visualisasi yang diinginkan dan *chatbot* dapat menginterpretasikannya ke dalam *query sql* yang nantinya dapat diolah menjadi sebuah visualisasi yang diinginkan.

Penelitian terdahulu dalam bidang text-to-SQL telah menghasilkan kemajuan signifikan, terutama dengan munculnya model-model seperti SQL-PaLm[11] dan RAT-SQL [12]. Model-model ini telah menunjukkan

kemampuan untuk menerjemahkan pertanyaan bahasa alami menjadi query SQL yang kompleks, bahkan pada dataset yang menantang seperti Spider[9]. Penelitian terdahulu telah menunjukkan potensi Large Language Model (LLM) dalam tugas text-to-SQL, namun masih terdapat beberapa keterbatasan. Beberapa penelitian tidak menjelaskan secara spesifik model LLM yang digunakan, sehingga sulit untuk memahami bagaimana hasil penelitian dapat direplikasi. Selain itu, metode fine-tuning yang digunakan dalam penelitian sebelumnya seringkali membutuhkan resources yang besar, seperti yang terlihat pada SQL-PaLM. Hal ini dapat menjadi kendala bagi peneliti lain yang ingin mengimplementasikan atau mengembangkan penelitian tersebut. Beberapa penelitian lain, seperti ChatDoctor[13] , Goat[14], Tamil-Llama [15], dan Llama Guard [16], berfokus pada fine-tuning Llama 2 untuk tugas-tugas spesifik seperti domain medis, aritmatika, penerjemahan bahasa Tamil, dan filterisasi kata-kata tidak pantas, namun tidak membahas penerapannya pada tugas text-to-SQL. Penelitian mengenai LLaMA-Adapter[17] dan LoRA fine-tuning[18] memperkenalkan metode fine-tuning yang efisien, tetapi tidak membahas penerapannya pada model Llama 2 atau dampaknya terhadap tugas text-to-SQL. Penelitian ini bertujuan untuk mengatasi keterbatasan-keterbatasan tersebut dengan mengimplementasikan fine-tuning model Llama 2 menggunakan metode QLoRA yang efisien, serta melakukan evaluasi komprehensif terhadap kemampuan model dalam menghasilkan query SQL yang akurat dan relevan. Oleh karena itu, penelitian ini mengusulkan penggunaan model Llama 2, sebuah model open-source yang mudah diakses dan memiliki performa yang baik dalam berbagai tugas NLP. Selain itu, penelitian ini akan menggunakan metode QLoRA untuk fine-tuning Llama 2, sebuah metode yang efisien dan dapat mengurangi penggunaan resources komputasi secara signifikan.

Pembuatan model ini didasarkan pada kebutuhan mengenai pengembangan *Artificial Intelligence* dalam penggunaan sehari-hari[19]. Penggunaan AI dalam berbagai bidang memberikan banyak manfaat terutama dalam membantu seluruh pihak untuk memberikan *output* yang maksimal. Salah satu *case* yang

dapat disampaikan adalah adanya keinginan untuk melakukan implementasi AI dari beberapa perusahaan komunikasi yang ada di Indonesia dalam beberapa kebutuhan perusahaan tersebut. Beberapa *use case* yang diinginkan dari perusahaan-perusahaan tersebut di antara lain adalah membangun AI yang mampu membantu tim HR untuk melakukan filterisasi terhadap calon pegawai yang mendaftar. *Use case* lain yang diinginkan adalah kemampuan untuk menciptakan sebuah AI yang dapat menjadi *co-pilot* atau *co-assistant* dalam proses pengolahan data lewat *text-to-sql* dengan model-model LLM seperti Llama-2 dan GPT yang saat ini digunakan oleh banyak industri.

Skripsi ini berupaya mengatasi keterbatasan tersebut dengan memanfaatkan model bahasa Llama 2 yang telah menunjukkan performa mengesankan dalam berbagai tugas NLP. Dengan menggunakan metode[8], model Llama 2 akan *fine-tuning* secara efisien untuk tugas *text-to-SQL*, dengan harapan dapat mencapai hasil yang lebih baik dibandingkan model-model sebelumnya, terutama dalam hal akurasi dan generalisasi.

Dengan demikian, skripsi ini diharapkan dapat memberikan kontribusi signifikan terhadap pengembangan model *text-to-SQL* yang lebih akurat, efisien, dan mudah digunakan.

## 1.2 Rumusan Masalah

Dengan latar belakang yang sudah dirumuskan, dapat ditarik beberapa rumusan masalah bagi penelitian ini, antara lain:

1. Bagaimana implementasi *fine-tuning* terhadap model Llama-2 dalam pembuatan model *text-to-sql* dapat dilakukan?
2. Bagaimana hasil performa yang diberikan oleh model Llama-2 yang sudah di-*fine tuning* terhadap hasil *query sql* yang diberikan?

## 1.3 Batasan Masalah

Ada pula beberapa batasan masalah yang terdapat pada penelitian ini, antara lain:

1. *Dataset* yang digunakan untuk proses *fine-tuning* merupakan *dataset sql-create-context* yang berisikan *question* dan *answer* mengenai *text-to-sql*
2. Model Llama-02 yang digunakan merupakan model dengan parameter terkecil yaitu Llama-2-7b-chat-hf dengan total 7 Miliar Parameter.
3. Bahasa yang dapat ditranslasikan ke dalam bentuk SQL adalah bahasa inggris.
4. Hasil akhir merupakan model *fine-tuned* Llama-2 yang dapat diintegrasikan ke sebuah sistem *chatbot*

## 1.4 Tujuan dan Manfaat Penelitian

### 1.4.1 Tujuan Penelitian

Tujuan dari adanya penelitian ini adalah sebagai berikut :

1. Mengimplementasi model Llama-02 yang sudah di-*fine tune* untuk membantu translasi *natural language* menjadi *sql query*
2. Menghasilkan model yang dapat mentranslasikan pertanyaan manusia menjadi sebuah *sql query* untuk pengumpulan dan pengolahan data

### 1.4.2 Manfaat Penelitian

Ada pula manfaat pada penelitian ini yang diharapkan sebagai berikut:

1. Membantu penyampaian informasi mengenai *database* yang dimiliki.
2. Menghasilkan model Llama-02 yang dapat digunakan untuk keperluan *text-to-sql* untuk nantinya dapat digunakan oleh khalayak umum

## 1.5 Sistematika Penulisan

Sistematika Penulisan dimuat dalam penelitian ini untuk memberikan gambaran terhadap kerangka laporan secara garis besar dengan menjelaskan beberapa bab, antara lain :

### BAB 1 : PENDAHULUAN

Bab ini menjelaskan perihal latar belakang dan juga masalah yang ingin dijelaskan dan diangkat dalam penelitian ini. Bab ini terdiri atas latar belakang, rumusan masalah, Batasan masalah, tujuan dan manfaat penelitian, dan sistematika penulisan.

## **BAB 2 : LANDASAN TEORI**

Bab ini berisi tentang landasan teori dari penelitian yang dilakukan serta penjelasan mengenai metode apa yang digunakan. Bab ini juga menjelaskan mengenai penelitian terdahulu sebagai dasar dan acuan untuk teori-teori yang berhubungan dengan penelitian.

## **BAB 3 : METODOLOGI PENELITIAN**

Bab ini menjelaskan mengenai penerapan dan proses metode yang dijelaskan pada bab 2. Pada bab ini juga dijelaskan mengenai proses pengumpulan data dan Langkah-langkah penelitian yang diperlukan untuk mencapai tujuan.

## **BAB 4 : ANALISIS DAN HASIL PENELITIAN**

Pada bab 4 ini akan menjelaskan proses CRISP-ML yang dilakukan dengan menjelaskan setiap bagian dan menampilkan hasil masing-masing tahap. Serta menampilkan hasil akhir dari penelitian yang dilakukan

## **BAB 5 : KESIMPULAN DAN SARAN**

Bab 5 Kesimpulan dan Saran berisi mengenai Kesimpulan yang didapatkan dari penelitian dan saran yang dapat disampaikan mengenai hasil penelitian.

UMMN

UNIVERSITAS  
MULTIMEDIA  
NUSANTARA