

Laporan Proyek: Backend Service untuk Screening Kandidat Berbasis AI

Informasi Kandidat

- **Nama Lengkap:** Rafli Ramadhani
- **Alamat Email:** raflirama7446@gmail.com

Link Repositori

- <https://github.com/zendParadox/backend-ai-screener>

1. Pendekatan & Desain Sistem

Rencana Awal

Tantangan ini saya pecah menjadi beberapa komponen utama:

1. **API Layer:** Membuat RESTful API yang jelas untuk interaksi dari luar.
2. **Proses Asynchronous:** Mengidentifikasi bahwa evaluasi AI akan memakan waktu, sehingga perlu dipisahkan dari alur request-response HTTP utama.
3. **RAG (Retrieval-Augmented Generation):** Merancang alur untuk menyimpan dan mengambil dokumen referensi (ground truth) untuk memberikan konteks yang akurat kepada LLM.
4. **Integrasi LLM:** Menghubungkan sistem dengan layanan AI generatif untuk melakukan analisis akhir.

Desain Sistem dan Database

API Endpoints: Sistem ini diekspos melalui tiga endpoint utama:

- **POST /upload:** Menerima file multipart/form-data (CV dan Laporan), menyimpannya di server, dan mengembalikan nama file sebagai ID unik.
- **POST /evaluate:** Menerima ID file dan detail pekerjaan, lalu menambahkan tugas evaluasi ke dalam antrian (queue) dan segera mengembalikan job_id.
- **GET /result/{id}:** Memungkinkan klien untuk memeriksa status pekerjaan (queued, processing, completed, failed) dan mengambil hasil akhir jika sudah selesai.

Antrian Pekerjaan (Job Queue): Untuk menangani proses evaluasi yang berjalan lama, saya menggunakan **BullMQ** dengan **Redis** sebagai backend. Pendekatan ini memastikan bahwa endpoint /evaluate dapat memberikan respons dengan cepat tanpa harus menunggu panggilan API ke LLM selesai. Ini juga membangun pondasi yang kuat untuk skalabilitas dan penanganan kegagalan di masa depan.

Database Vektor: Saya memilih **Qdrant** sebagai vector database. Dokumen-dokumen referensi seperti Deskripsi Pekerjaan, Case Study Brief, dan Rubrik Penilaian di-ingest ke dalam sebuah *collection* di Qdrant. Ini memungkinkan sistem untuk melakukan pencarian semantik yang efisien guna menemukan konteks yang paling relevan untuk setiap evaluasi.

Integrasi LLM

Pilihan LLM: Saya menggunakan model **gemini-2.5-flash** dari Google AI. Model ini dipilih karena menawarkan keseimbangan yang baik antara kecepatan, kemampuan mengikuti instruksi kompleks (seperti menghasilkan format JSON yang ketat), dan efisiensi biaya.

Strategi RAG (Retrieval-Augmented Generation): Alur RAG diimplementasikan sebagai berikut:

1. Saat evaluasi dimulai, konten dari CV kandidat diubah menjadi vektor embedding menggunakan model Xenova/all-MiniLM-L6-v2 yang berjalan secara lokal.
2. Vektor ini digunakan untuk melakukan pencarian di Qdrant untuk mengambil dokumen referensi yang paling relevan.
3. Kumpulan dokumen yang relevan ini (konteks) kemudian disuntikkan ke dalam *prompt* yang dikirim ke Gemini.

Strategi Prompting: Prompt dirancang dengan sangat hati-hati untuk memastikan output yang konsisten dan terstruktur:

1. **Persona:** AI diberi peran sebagai "Asisten HR Ahli".
2. **Instruksi Ketat:** Memberikan instruksi yang jelas untuk menggunakan dokumen konteks sebagai satu-satunya sumber kebenaran.
3. **Struktur Output:** Mendefinisikan skema JSON yang wajib diikuti.
4. **Penyisipan Data:** Menyediakan konteks dari RAG dan data dari CV/Laporan kandidat secara terpisah dengan penanda yang jelas.

2. Hasil & Refleksi

Hasil: Arsitektur asinkron dengan BullMQ bekerja dengan sangat baik. Sistem mampu menerima pekerjaan, memprosesnya di latar belakang, dan menyajikan hasilnya tanpa memblokir klien. Alur RAG juga berhasil mengambil konteks yang relevan, yang secara signifikan meningkatkan kualitas output dari Gemini dibandingkan tanpa konteks.

Refleksi: Tantangan terbesar adalah memastikan kompatibilitas versi antara semua komponen (server Qdrant, klien Node.js, dll.). Selain itu, kualitas embedding sangat menentukan keberhasilan RAG. Menggunakan model embedding yang berjalan secara lokal memberikan keuntungan dari segi privasi dan biaya, meskipun memerlukan waktu untuk inisialisasi pada proses pertama.

3. Peningkatan di Masa Depan

Dengan lebih banyak waktu, beberapa area yang dapat ditingkatkan adalah:

- **Error Handling:** Mengimplementasikan logika *retry* dengan *exponential backoff* yang lebih canggih untuk panggilan API ke Gemini.
- **Testing:** Menambahkan unit test dan integration test untuk memastikan keandalan setiap komponen.
- **Frontend:** Membangun antarmuka pengguna sederhana untuk memudahkan proses upload dan pengecekan hasil.

4. Screenshots of Real Responses

