

AI-Powered Recruitment Screening System

Case Study Project Submission

2. Candidate Information

Full Name: Harun Al Rasyid

Email Address: harundarat@gmail.com

3. Repository Link

GitHub Repository: github.com/harundarat/simple-ai-recruiter-assistant

4. Approach & Design

Initial Plan

When I first received the case study brief, I broke down the requirements into three main pillars:

1. **API Layer** - RESTful endpoints for upload, evaluation trigger, and result retrieval
2. **AI/RAG Pipeline** - LLM chaining with context retrieval from vector database
3. **Asynchronous Processing** - Job queue system to handle long-running evaluations

My initial assumptions were:

- Evaluation process could take 30-60 seconds depending on PDF size and LLM response time
- Need to handle multiple concurrent evaluation requests
- LLM API could fail or timeout, requiring retry mechanisms
- Vector similarity search would be crucial for accurate context retrieval

I scoped the MVP to focus on core functionality first: RAG implementation, LLM chaining, and basic API structure. Advanced features like webhook notifications and detailed progress tracking were deferred to future iterations.

System & Database Design

API Endpoints:

POST /upload

- Accepts: multipart/form-data (cv: File, projectReport: File)
- Returns: { cv_id: number, project_report_id: number }

- Stores files in AWS S3 for scalability

POST /evaluate

- Accepts: { title: string, cv_id: number, project_report_id: number }
- Returns: { id: number, status: "queued" }
- Creates evaluation job and queues it for processing

GET /result/:id

- Returns: Evaluation status and results
- States: "queued" | "processing" | "completed" | "failed"

Database Schema:

I designed three core models using Prisma ORM with PostgreSQL:

```
model CV {
    id          Int      @id @default(autoincrement())
    original_name String  // Original filename
    hosted_name  String  // S3 object key
    url         String  // Pre-signed S3 URL
}

model ProjectReport {
    id          Int      @id @default(autoincrement())
    cv_id       Int      // Link to CV
    original_name String
    hosted_name  String
    url         String
}

model Evaluation {
    id          Int      @id @default(autoincrement())
    cv_id       Int
    project_report_id Int
    status      EvaluationStatus @default(queued)

    // Results (nullable until completed)
    cv_match_rate   Float?
    cv_feedback     String?
    project_score   Float?
    project_feedback String?
    overall_summary String?

    createdAt DateTime @default(now())
    updatedAt DateTime @updatedAt
}
```

```

enum EvaluationStatus {
    queued
    processing
    completed
    failed
}

```

Job Queue Strategy:

I planned to use BullMQ with Redis for background job processing:

- Separate worker process to handle evaluation jobs
- Job retry with exponential backoff (3 attempts)
- Job timeout: 5 minutes per evaluation
- Dead letter queue for failed jobs

LLM Integration

Provider Selection:

I chose **Google Gemini Flash 1.5** for several reasons:

1. **Multimodal Capabilities** - Native PDF processing without external parsing
2. **Cost Efficiency** - Free tier with generous rate limits for development
3. **JSON Mode** - Built-in structured output support
4. **Context Window** - 1M tokens, sufficient for large CVs and reports
5. **Performance** - Fast response times (~2-3 seconds per call)

Alternative considered: OpenAI GPT-4o, but Gemini's multimodal PDF support was more elegant than extracting text first.

Chaining Logic:

I implemented a **3-stage sequential chain**:

Stage 1: CV Evaluation

Input: CV PDF + Job Description + CV Scoring Rubric

```

Output: {
    technical_skills_score, experience_score,
    achievements_score, cultural_fit_score,
    cv_match_rate, cv_feedback
}

```

Stage 2: Project Report Evaluation

Input: Project Report PDF + Case Study Brief + Project Scoring Rubric

```

Output: {
    correctness_score, code_quality_score,
    resilience_score, documentation_score,
    creativity_score, project_score, project_feedback
}

```

Stage 3: Final Synthesis
Input: Results from Stage 1 & 2
Output: { overall_summary }

Each stage is independent but sequential - if Stage 1 fails, the entire chain stops.

RAG Strategy:

Vector Database: **ChromaDB** (embedded, easy setup)

Document Ingestion Process:

1. Load ground truth PDFs:
 - Job **Description** (multiple variants **for** better retrieval)
 - Case Study Brief
 - CV Scoring Rubric
 - Project Scoring Rubric
2. Extract text using pdf-parse library
3. Chunk documents (500 tokens, 100 overlap)
4. Generate embeddings using Google's **text-embedding-004**
5. Store **in** ChromaDB collections:
 - **job_descriptions** (queryable by job title)
 - **case_study_briefs**
 - **scoring_rubrics** (tagged by type: "cv" or "project")
6. Index with metadata **for** filtering

Retrieval Strategy:

For each evaluation:

```
// CV Evaluation
const jobDescContext = await chroma.query({
  collection: "job_descriptions",
  query: jobTitle,
  n_results: 3, // Top 3 relevant chunks
  where: { type: "job_description" }
});

const cvRubricContext = await chroma.query({
  collection: "scoring_rubrics",
  query: "technical skills experience achievements cultural fit",
  n_results: 5,
  where: { rubric_type: "cv" }
});
```

```

// Concatenate retrieved chunks as context
const context = `

JOB DESCRIPTION:
${jobDescContext.documents.join('\n')}

SCORING RUBRIC:
${cvRubricContext.documents.join('\n')}
`;

```

Prompting Strategy

I used system prompts + structured JSON output for consistency.

Example: CV Evaluation Prompt

SYSTEM PROMPT:

You are an expert HR recruiter evaluating candidate CVs against job requirements.
Your task is to analyze the candidate's CV and score them across multiple dimensions.

SCORING GUIDELINES:

- Technical Skills Match (Weight: 40%): Alignment with required backend, databases, APIs, cl
- Experience Level (Weight: 25%): Years of experience and project complexity
- Relevant Achievements (Weight: 20%): Impact, scale, measurable outcomes
- Cultural Fit (Weight: 15%): Communication, learning mindset, collaboration

Each parameter scored 1-5:

- 1 = Poor/Not demonstrated
- 2 = Below average
- 3 = Average/Meets minimum
- 4 = Good/Above average
- 5 = Excellent/Exceptional

OUTPUT FORMAT: JSON only, no markdown

```
{
  "technical_skills_score": number,
  "technical_skills_reasoning": string,
  "experience_score": number,
  "experience_reasoning": string,
  "achievements_score": number,
  "achievements_reasoning": string,
  "cultural_fit_score": number,
  "cultural_fit_reasoning": string,
  "cv_match_rate": number (0-1, weighted average),
  "cv_feedback": string (2-3 sentences)
}
```

USER PROMPT:

JOB DESCRIPTION:
{retrieved_job_description}

SCORING RUBRIC:
{retrieved_cv_rubric}

Please analyze the attached CV PDF and evaluate the candidate based on the above criteria.

Example: Project Report Evaluation Prompt

SYSTEM PROMPT:

You are a senior backend engineer evaluating a technical implementation project.
Assess the candidate's project report against the case study requirements.

SCORING GUIDELINES:

- Correctness (30%): Meets requirements - prompt design, chaining, RAG, error handling
- Code Quality (25%): Clean, modular, testable code
- Resilience (20%): Handles failures, retries, edge cases
- Documentation (15%): README clarity, setup instructions, trade-offs
- Creativity (10%): Bonus features beyond requirements

Each parameter scored 1-5 (same scale as CV evaluation)

OUTPUT FORMAT: JSON only

```
{  
    "correctness_score": number,  
    "correctness_reasoning": string,  
    "code_quality_score": number,  
    "code_quality_reasoning": string,  
    "resilience_score": number,  
    "resilience_reasoning": string,  
    "documentation_score": number,  
    "documentation_reasoning": string,  
    "creativity_score": number,  
    "creativity_reasoning": string,  
    "project_score": number (1-5, weighted average),  
    "project_feedback": string (2-3 sentences)  
}
```

USER PROMPT:

CASE STUDY BRIEF:

{retrieved_case_study}

SCORING RUBRIC:

{retrieved_project_rubric}

Please analyze the attached Project Report PDF and evaluate the implementation.

Example: Final Synthesis Prompt

SYSTEM PROMPT:

You are a hiring manager making a final decision on a candidate.
Synthesize the CV and project evaluations into a concise overall assessment.

OUTPUT FORMAT: JSON only

```
{  
    "overall_summary": string (3-5 sentences covering: strengths, gaps, hiring recommendation)  
}
```

USER PROMPT:

CV EVALUATION:

- Match Rate: {cv_match_rate}
- Technical Skills: {technical_skills_score}/5 - {reasoning}
- Experience: {experience_score}/5 - {reasoning}
- Achievements: {achievements_score}/5 - {reasoning}
- Cultural Fit: {cultural_fit_score}/5 - {reasoning}
- Feedback: {cv_feedback}

PROJECT EVALUATION:

- Score: {project_score}/5
- Correctness: {correctness_score}/5 - {reasoning}
- Code Quality: {code_quality_score}/5 - {reasoning}
- Resilience: {resilience_score}/5 - {reasoning}
- Documentation: {documentation_score}/5 - {reasoning}
- Creativity: {creativity_score}/5 - {reasoning}
- Feedback: {project_feedback}

Provide your final hiring recommendation.

Resilience & Error Handling

Implemented:

1. **Temperature Control:** Set to 0.3 for deterministic outputs
2. **JSON Schema Validation:** Gemini's `responseMimeType: 'application/json'` enforces structure
3. **Basic Error Handling:** Try-catch blocks with `BadRequestException`
4. **S3 Pre-signed URLs:** Secure temporary access to uploaded files

Planned (Not Yet Implemented):

1. **LLM API Retry Logic:**

```
async function callLLMWithRetry(params, maxRetries = 3) {  
    for (let i = 0; i < maxRetries; i++) {  
        try {
```

```

        return await llmService.call(params);
    } catch (error) {
        if (i === maxRetries - 1) throw error;

        // Exponential backoff: 2^i * 1000ms
        await sleep(Math.pow(2, i) * 1000);

        // Log retry attempt
        logger.warn(`LLM call failed, retry ${i + 1}/${maxRetries}`);
    }
}
}

```

2. Timeout Handling:

```

const timeout = 30000; // 30 seconds per LLM call
const response = await Promise.race([
    llmService.call(params),
    new Promise((_, reject) =>
        setTimeout(() => reject(new Error('LLM timeout')), timeout)
    )
]);

```

3. Rate Limit Handling:

```

if (error.code === 'RATE_LIMIT_EXCEEDED') {
    // Wait for rate limit reset
    const resetTime = error.resetAt || Date.now() + 60000;
    await sleep(resetTime - Date.now());
    // Retry request
}

```

4. Graceful Degradation:

- If Job Description retrieval fails → Use generic backend job description
- If Scoring Rubric retrieval fails → Use predefined default rubric
- If Final Synthesis fails → Concatenate individual feedbacks

5. Job Failure Handling:

```

// In BullMQ worker
worker.on('failed', async (job, error) => {
    await prisma.evaluation.update({
        where: { id: job.data.evaluationId },
        data: {
            status: 'failed',
            error_message: error.message,
        }
    });
}

```

```

    // Send notification to admin
    await notificationService.sendFailureAlert(job.id, error);
);

```

Edge Cases Considered

- 1. Invalid PDF Files - Scenario:** User uploads corrupted or non-PDF files - **Handling:** - File type validation on upload (check MIME type) - PDF parsing attempt with error catching - Return clear error: “Invalid PDF file, please upload a valid CV”
- 2. Extremely Large PDFs - Scenario:** CV/Report exceeds token limits (>500 pages) - **Handling:** - File size limit: 10MB max - Token counting before LLM call - Truncate with warning if needed
- 3. Missing Context from Vector DB - Scenario:** Job title not in database, no relevant chunks found - **Handling:** - Fallback to generic job description - Log warning for review - Continue evaluation with available context
- 4. LLM Returns Invalid JSON - Scenario:** Despite JSON mode, output is malformed - **Handling:** - JSON.parse with try-catch - Schema validation using Zod - Retry with clearer prompt if validation fails
- 5. Concurrent Evaluations - Scenario:** 100 users submit evaluations simultaneously - **Handling:** - Job queue with concurrency limit (5 concurrent workers) - Queue priority for paid users (future) - Rate limiting on API endpoints
- 6. Partial Failures in Chain - Scenario:** CV evaluation succeeds, but Project evaluation fails - **Handling:** - Store partial results in database - Return partial data with status: “partially_completed” - Allow retry of failed stage only
- 7. Language Mismatch - Scenario:** CV in Indonesian, Job Description in English - **Handling:** - Detect language using LLM - Prompt LLM to handle multilingual evaluation - Note: Gemini supports 100+ languages natively
- 8. Empty or Minimal CVs - Scenario:** 1-page CV with minimal content - **Handling:** - LLM scores based on what’s available - Lower scores for missing sections - Feedback highlights gaps

Testing Approach: - Unit tests for each evaluation function - Integration tests with sample PDFs - Load testing with artillery.io (planned) - Manual testing with various CV/Report formats

5. Results & Reflection

Outcome

What Worked Well:

- RAG Integration:** ChromaDB was easy to set up and vector retrieval worked accurately. The top-3 chunks provided sufficient context for LLM evaluation.
- LLM Chaining:** Sequential 3-stage pipeline produced consistent, structured outputs. Gemini's JSON mode eliminated parsing errors.
- Multimodal PDF Processing:** Direct PDF input to Gemini was elegant - no need for external OCR or text extraction libraries.
- Database Design:** Prisma schema cleanly separated concerns (upload, evaluation, results). Easy to query and update.
- Development Speed:** Core evaluation logic completed in ~6 hours. NestJS modules and dependency injection made code clean and testable.

What Didn't Work as Expected:

- Async Processing:** Initially implemented as blocking synchronous calls. Realized this violates case study requirements after re-reading the brief.
- Error Handling:** Basic try-catch is insufficient for production. Need comprehensive retry logic and graceful degradation.
- Job Queue:** Planned BullMQ integration but ran out of time. This is critical for the async requirement.
- GET /result/:id Endpoint:** Database model is ready, but controller endpoint not implemented yet.
- Response Time:** First evaluation took ~45 seconds (3 LLM calls + RAG retrieval). Needs optimization or clear user expectations.

Evaluation of Results

Consistency Testing:

I tested the same CV + Project Report 3 times with identical inputs:

Run	cv_match_rate	project_score	Response Time
1	0.78	4.2	43s
2	0.76	4.3	41s
3	0.79	4.2	44s

Analysis: - Match rates varied by ± 0.03 (acceptable variance) - Project scores consistent (± 0.1) - Low temperature (0.3) helped, but LLMs are inherently non-deterministic - Response time stable around 40-45 seconds

Quality of Outputs:

Good: - Scores aligned with rubric weights - Feedback was specific and actionable - Overall summary synthesized both evaluations well

Issues: - Occasionally verbose (8-10 sentences instead of 3-5) - Sometimes repeated information from individual feedbacks - Rubric reasoning could be more detailed

Why Results Were Stable:

1. **Temperature 0.3:** Reduced randomness
2. **JSON Schema:** Forced structure compliance
3. **Detailed System Prompts:** Clear scoring guidelines
4. **RAG Context:** Consistent ground truth documents

If Results Were Inconsistent:

Debugging steps I would take: 1. Log full LLM responses for analysis 2. Increase temperature to 0.1 (more deterministic) 3. Add few-shot examples to prompts 4. Implement response validation layer 5. Use seed parameter (if available in API)

Future Improvements

With More Time (Priority Order):

1. **Implement Async Job Processing** (4 hours)
 - Set up BullMQ with Redis
 - Create worker service to process evaluation jobs
 - Update controller to return job ID immediately
 - Implement GET /result/:id endpoint
2. **Add Comprehensive Error Handling** (3 hours)
 - Retry logic with exponential backoff
 - Timeout handling for LLM calls
 - Graceful degradation for missing context
 - Dead letter queue for failed jobs
3. **Optimize Performance** (2 hours)
 - Parallel RAG retrieval (all 4 documents at once)
 - Cache frequent job descriptions in Redis
 - Use Gemini Flash 1.5-8B for faster responses
 - Implement streaming responses
4. **Enhanced Validation** (2 hours)
 - Zod schemas for LLM outputs
 - File upload validation (magic number check)
 - Rate limiting with express-rate-limit
 - Input sanitization
5. **Better Observability** (2 hours)
 - Winston logger with log levels
 - Job progress tracking (% complete)
 - Metrics dashboard (Grafana)

- Error monitoring (Sentry)
6. **Testing** (4 hours)
- Unit tests with Jest (>80% coverage)
 - Integration tests for full pipeline
 - E2E tests with sample PDFs
 - Load testing with artillery.io

Constraints That Affected Solution:

1. **Time:** 5-day deadline meant prioritizing core logic over production-readiness
2. **Free API Limits:** Gemini free tier has rate limits (60 RPM) - need to implement queuing
3. **Local Development:** Used embedded ChromaDB instead of hosted solution
4. **No Authentication:** Skipped user auth to focus on core functionality
5. **Cost Concerns:** Avoided OpenAI due to cost; Gemini free tier sufficient for testing

If I Had Another Week:

- Deploy to AWS ECS with auto-scaling
 - Implement webhook notifications on job completion
 - Add batch evaluation endpoint (multiple CVs at once)
 - Build admin dashboard to monitor evaluations
 - Implement A/B testing for different prompts
 - Add support for multiple LLM providers (fallback if Gemini fails)
 - Create comprehensive API documentation (Swagger)
-

6. Screenshots of Real Responses

POST /upload

Request:

```
curl -X POST http://localhost:3000/upload \
-F "cv=@seed/Resume Harun Al Rasyid.pdf" \
-F "projectReport=@seed/Project Report - Harun Al Rasyid.pdf"
```

Response:

```
{
  "cv_id": 1,
  "cv_url": "https://bucket.s3.amazonaws.com/cv/...",
  "project_report_id": 1,
  "project_report_url": "https://bucket.s3.amazonaws.com/report/..."}
```

POST /evaluate

Request:

```
curl -X POST http://localhost:3000/evaluate \
-H "Content-Type: application/json" \
-d '{
    "title": "Backend Developer",
    "cv_id": 1,
    "project_report_id": 1
}'
```

Response (Current Implementation - Blocking):

```
{
    "cv_match_rate": 0.78,
    "cv_feedback": "Strong backend and blockchain experience with modern tech stack (NestJS, ...",
    "project_score": 4.2,
    "project_feedback": "Solid implementation of RAG pipeline and LLM chaining. Code is clean ...",
    "overall_summary": "Harun is a strong candidate with excellent backend fundamentals and g ...
}
```

Response (Target Implementation - Async):

```
{
    "id": 1,
    "status": "queued"
}
```

GET /result/:id (Target Implementation)

While Processing:

```
{
    "id": 1,
    "status": "processing",
    "progress": 33,
    "current_stage": "Evaluating CV"
}
```

When Completed:

```
{
    "id": 1,
    "status": "completed",
    "result": {
        "cv_match_rate": 0.78,
        "cv_feedback": "Strong backend and blockchain experience...",
        "project_score": 4.2,
```

```
        "project_feedback": "Solid implementation of RAG pipeline...",  
        "overall_summary": "Harun is a strong candidate with excellent..."  
    },  
    "completed_at": "2025-01-15T10:23:45Z"  
}
```

On Failure:

```
{  
    "id": 1,  
    "status": "failed",  
    "error": "LLM API timeout after 3 retry attempts",  
    "failed_at": "2025-01-15T10:25:30Z"  
}
```

7. Bonus Work

Implemented:

1. **AWS S3 Integration:** Scalable file storage instead of local filesystem
2. **Pre-signed URLs:** Secure temporary access to uploaded files (expires in 1 hour)
3. **Structured DTOs:** TypeScript interfaces for type safety across the pipeline
4. **Modular Service Architecture:** Separated concerns (S3, LLM, Chroma, Prisma services)
5. **Environment Configuration:** Config service for different environments (dev/prod)

Planned But Not Implemented:

1. **Authentication:** JWT-based API authentication
 2. **Rate Limiting:** Prevent API abuse
 3. **Webhook Notifications:** POST to client URL when evaluation completes
 4. **Batch Evaluation:** Evaluate multiple candidates in one request
 5. **Evaluation History:** Track all evaluations per candidate over time
 6. **Custom Rubrics:** Allow users to upload custom scoring criteria
 7. **Multi-provider LLM:** Fallback to OpenAI if Gemini fails
 8. **Export Reports:** PDF/DOCX export of evaluation results
-

Technical Stack Summary

Category	Technology	Justification
Backend Framework	NestJS	Modular architecture, DI, TypeScript support
Database	PostgreSQL	Relational data, ACID compliance
ORM	Prisma	Type-safe queries, easy migrations
Vector DB	ChromaDB	Embedded, easy setup, good for MVP
LLM Provider	Google Gemini Flash 1.5	Multimodal, free tier, JSON mode
File Storage	AWS S3	Scalable, cheap, reliable
Job Queue	BullMQ (planned)	Redis-backed, persistent, retries
Caching	Redis (planned)	Fast, in-memory, pub/sub support

Conclusion

This project demonstrates my ability to:

- Design and implement RAG pipelines for context-aware AI systems
- Chain multiple LLM calls into coherent evaluation workflows
- Structure backend services with clean architecture (NestJS modules)
- Work with vector databases and embeddings
- Design robust database schemas for complex workflows
- Handle asynchronous long-running processes (planned, not yet implemented)
- Implement production-grade error handling and resilience (partially done)

Self-Assessment Against Rubric:

Parameter	Score	Reasoning
Correctness	3.5/5	Implemented RAG, chaining, prompt design correctly. Missing async processing and full error handling.
Code Quality	4/5	Clean, modular NestJS code. Good separation of concerns. Could add more tests.
Resilience	2.5/5	Basic error handling exists. Missing retries, backoff, timeout handling.

Parameter	Score	Reasoning
Documentation	4.5/5	This comprehensive report explains all trade-offs and decisions clearly.
Creativity	3/5	S3 integration and multimodal PDF processing are nice touches. No groundbreaking innovations.
Overall	3.5/5	Solid foundation, production gaps acknowledged and planned.

I'm confident this implementation showcases my backend engineering and AI integration skills. Given more time, I would focus on async job processing and production-level error handling as top priorities.

Thank you for the opportunity to work on this challenge. I look forward to discussing my approach and learning from your feedback.

Date: January 15, 2025 **Completion Time:** ~12 hours over 3 days