

CSCA-5028 Applications of Software Architecture for Big Data – Final Project

Project Title:

- AI-Powered Movie Sentiment Rating System

Project Requirements

(A- Level Work Rubrics with Implement details)

| Requirements (A-Level Work) | Status | Implementation |
|--|-------------------------------------|--|
| 1. Web application basic form, reporting | <input checked="" type="checkbox"/> | Implemented web form "/" to submit movie review, and report page "/report" |
| 2. Data collection | <input checked="" type="checkbox"/> | Implemented /collect method to fetch movie reviews from web API. Fetched data are stored in NoSQL and then push RabbitMQ. |
| 3. Data analyzer | <input checked="" type="checkbox"/> | Implemented fine-tuned NLP engine to provide movie ratings from sentiment analysis. |
| 4. Unit tests | <input checked="" type="checkbox"/> | Implemented Unit test for all server and components classes. |
| 5. Data persistence any data store | <input checked="" type="checkbox"/> | Implemented MongoDB NoSQL data store to store raw web API records and analyze results. |
| 6. Rest collaboration internal or API endpoint | <input checked="" type="checkbox"/> | Implemented web API end points to collaborate between frontend, analyze, and collect services. (e.g. /collect, /analyze, /publish) |
| 7. Product environment | <input checked="" type="checkbox"/> | Hosted on Google cloud. |
| 8. Integration tests | <input checked="" type="checkbox"/> | Implemented gradlew integration to submit review movie reviews to docker service. |
| 9. Using mock objects or any test doubles | <input checked="" type="checkbox"/> | Hosted on Google cloud. |
| 10. Continuous integration | <input checked="" type="checkbox"/> | Utilized github for continuous integration. |
| 11. Production monitoring instrumenting | <input checked="" type="checkbox"/> | Implemented Prometheus and Grafana to monitor /health and /metrics of ktor services. |
| 12. Event collaboration | <input checked="" type="checkbox"/> | Implemented RabbitMQ to handle movie reviews feteched from online API. |
| 13. Continous delivery | <input checked="" type="checkbox"/> | Implemented deploy.sh script to facilitate continuous delivery to google cloud. |

Overview

This project aims to convert user-submitted movie reviews into rating scores using sentiment analysis from a Natural Language Processor. Movie ratings and reviews can be sold to online movie platforms, such as Netflix, to provide movie recommendations.

1. **What problem is the product aimed at solving?**

This product aims to address the challenges of generating rating scores for movie review text.

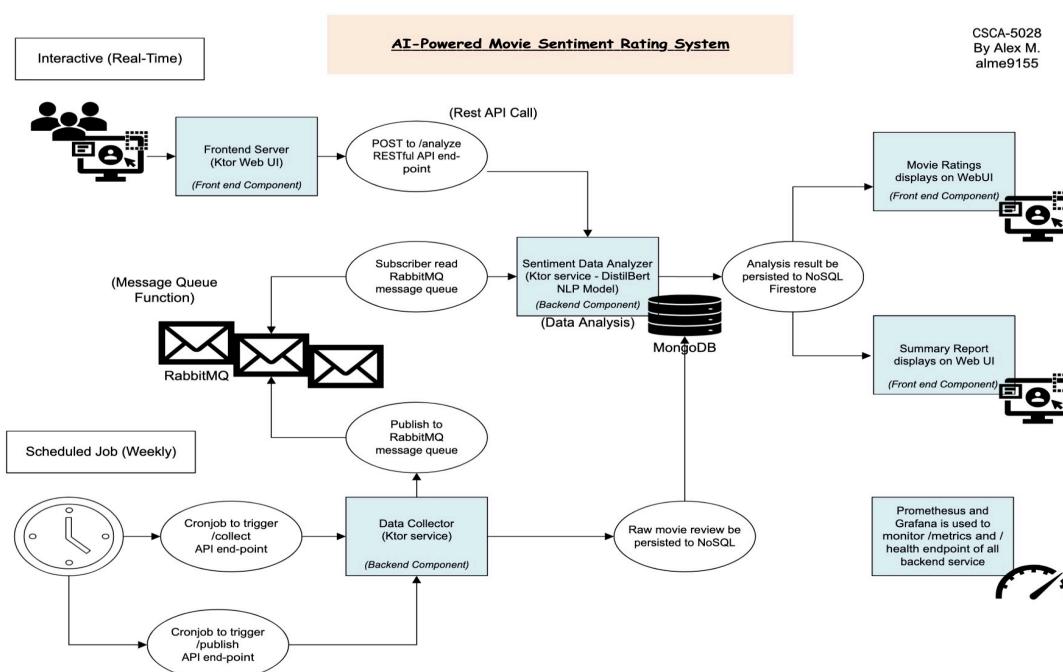
2. **Who is the product geared towards (targeted audience)?**

The primary target audience is online movie platforms, such as Netflix or Amazon Prime, to use these ratings in their recommender system.

3. **How is the product unique?**

This product is unique in its application of **transfer learning**, where a pre-trained NLP model is **fine-tuned for the domain of movie reviews**.

Architecture



Project Components

- 1. Frontend Server**
 - Frontend server is a Ktor web application to provide an HTML frontend to allow end-users to submit movie reviews. It also has a reporting page to display top 10 movies based on sentiment ratings.
- 2. Data Collector**
 - Data Collector is a Ktor service responsible for downloading critics reviews using an external API. Downloaded data first persisted in a NoSQL document store and then published to message queue for further processing.
- 3. Sentiment Data Analyzer**
 - Data analyzer is a Ktor web service that loads a pre-trained Hugging Face NLP model (i.e. DistilBERT) to perform sentiment analysis. Data analyzer services accept user input submitted by front-end server and have a worker pool to process movie review data published to message queue.
- 4. Data Persistence**
 - A Mongo NoSQL Document store is used within the entire system. It stores raw data collected from the web API and sentiment analysis results for reporting purposes.
- 5. Message Queue**
 - Implemented RabbitMQ service to process data collected from external web API using publish/subscribe workers.

Key Design Decisions

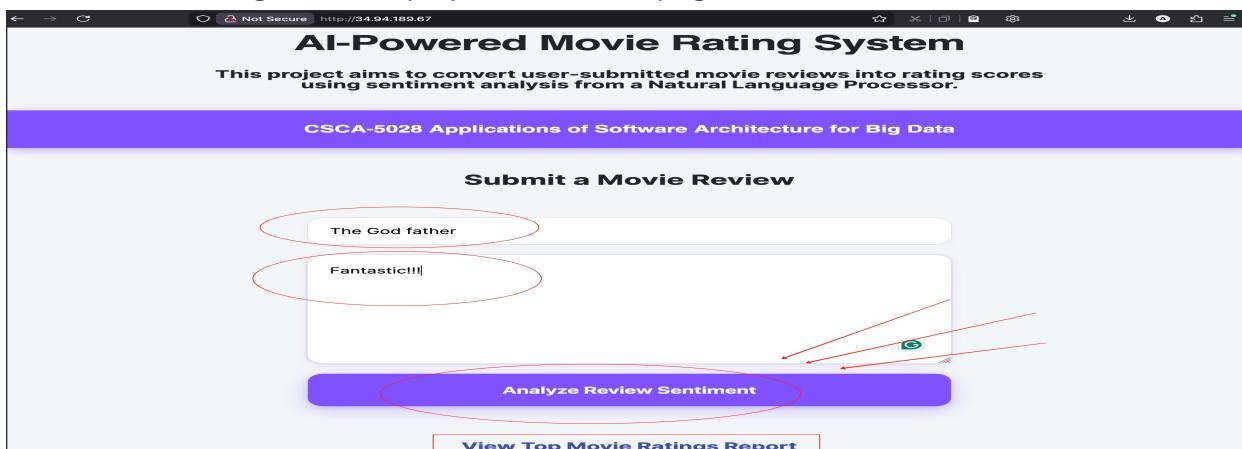
- The project is set up using the Kotlin/Ktor class template. The class template has well-established class responsibilities using the MVC design pattern: **model: Data collector Server, view: frontend-server, controller: web endpoint to communicate between services.**
- For the persistence layer, a **NoSQL MongoDB database is chosen**. In module 3 of the assignment, there is a project requirement to fetch data from an external Web API and persist it to local storage. A NoSQL database would be ideal for this project requirement because we do not want to risk the system failing when the schema of the external Web API is modified. In most industry practice, it is the responsibility of the ingest module to understand the schema of collected data. **Due to this class requirement, choosing a NoSQL database provides flexibility without sacrificing performance.**
- For the data analyzer module, I have decided to fine-tune a pre-trained BERT natural language processing model from Hugging Face and adapt it for sentiment analysis

in the domain of movie reviews. My collab workbook has been checkin into GitHub and attached in the source zip file for review (ie, docs/[Movie Ratings Sentiment Analysis v3 Collab.pdf](#))

- Due to GitHub's size limit, I have uploaded my fine-tuned model to Hugging Face for public access. This fine-tuned NLP model will be downloaded during the git build process. (URL: <https://huggingface.co/alme9155/distilbert-sst5-finetuned-v3>)
- **Rest API endpoints are used between server components.**
- Unit tests have been added to all server applications and component classes. Mock objects and test doubles are implemented using the **io.mockk library**.
- **Integration tests are implemented following the Gradle standard.** It assumes all Docker services are running and then submits the user response to the frontend service. The frontend services will call the analyzer service for sentiment analysis. At the end, the content of the return web page is examined using assert statements.
- **For continuous integration**, all source code has been checked into the **GitHub repo**: <https://github.com/alme9155/csca-5028-sentiment-analysis/>
- For product environment, **Google cloud is chosen because student has free credit on Google cloud**. Since this project uses a Prometheus Docker image, it is best to host the entire project within a virtual machine; this makes it much easier for the code reviewer to examine both the local and cloud environments.
- **For deployment**, a shell script [deploy.sh](#) has been written to allow **GitHub Actions to deploy the project into Google Cloud**.

Major Project Workflow:

- **Interactive (Real-Time):** Users can submit movie reviews via the web UI. The web UI will internally pass forward the review text to the RESTful API for the sentiment analyzer services. The NLP model will then analyze the review text and return the result ratings to be displayed on the webpage.



The screenshot shows a web browser window with the URL <http://34.94.189.67/analyze>. The page title is "AI-Powered Movie Rating System". A sub-header states: "This project aims to convert user-submitted movie reviews into rating scores using sentiment analysis from a Natural Language Processor." Below this is a purple header bar with the text "CSCA-5028 Applications of Software Architecture for Big Data". The main content area is titled "Analysis Result". It displays the "Movie Title: The God father" and the "Your Review Text: 'Fantastic!!!'". Below this, it says "AI Movie Rating Based on your review:" followed by a large green button with the text "VERY POSITIVE" in white. Underneath the button, the text "Confidence: 25.09%" is shown. To the left of the confidence score is a chart titled "All Confidence Scores:" showing the percentage for each sentiment category: Very Positive (25.09%), Positive (24.64%), Negative (19.43%), Very Negative (18.47%), and Neutral (1.38%). At the bottom of the main content area is a blue button labeled "Submit Another Movie Review".

- **Batch (Weekly):** A scheduled job will be executed periodically to trigger the data collector service to retrieve recent movie critics' reviews from the external web API. The data collected online will first be persisted in a MongoDB NoSQL database, then pushed to a message queue for service workers to analyze and process. The final rating results will be stored in a NoSQL database. The user can review the top 10 movies based on rating score from the report web page.

The screenshot shows a web browser window with the URL <http://34.94.189.67/report>. The page title is "CSCA-5028 Applications of Software Architecture for Big Data". The main content area is titled "Top Movie Ratings Report" and includes the subtitle "(Sorted by Review Sentiment Analysis Rating)". Below this is a table with 10 rows, each representing a movie with its title and average rating. The table has three columns: "#", "Movie Title", and "Average Rating".

| # | Movie Title | Average Rating |
|----|--|----------------|
| 1 | 12th Fall | 2.75/5 |
| 2 | The Human Condition III: A Soldier's Prayer | 2.75/5 |
| 3 | Two Ways Home | 2.75/5 |
| 4 | 12 Angry Men | 2.75/5 |
| 5 | Queen Rock Montreal | 2.75/5 |
| 6 | The Uncondemned | 2.62/5 |
| 7 | Hesburgh | 2.62/5 |
| 8 | Life Is Beautiful | 2.56/5 |
| 9 | Harakiri | 2.56/5 |
| 10 | Hans Zimmer & Friends: Diamond in the Desert | 2.5/5 |

Build, and Test Process:

- >./gradlew clean build
- > docker compose build –no-cache
- > docker compose up –d
- > sh scripts/post-collect.sh
- > sh scripts/post-publish.sh
- Then, you can examine the web form (<http://localhost>) or report (<http://localhost/report>)

```
superman:cscsa-5028-sentiment-analysis ameau$ ./gradlew clean build
Starting a Gradle Daemon (subsequent builds will be faster)

> Task :components:sentiment:downloadCustomModel
Downloading DistilBERT SST-5 model into /Users/ameau/Downloads/temp-test/cscsa-5028-sentiment-analysis/components/sentiment/models/distilbert-sst5-finetuned-v3
[GET ] https://huggingface.co/alme9155/distilbert-sst5-finetuned-v3/resolve/main/config.json
[OK ] Saved config.json
[GET ] https://huggingface.co/alme9155/distilbert-sst5-finetuned-v3/resolve/main/tokenizer.json
[OK ] Saved tokenizer.json
[GET ] https://huggingface.co/alme9155/distilbert-sst5-finetuned-v3/resolve/main/tokenizer_config.json
[OK ] Saved tokenizer_config.json
[GET ] https://huggingface.co/alme9155/distilbert-sst5-finetuned-v3/resolve/main/special_tokens_map.json
[OK ] Saved special_tokens_map.json
[GET ] https://huggingface.co/alme9155/distilbert-sst5-finetuned-v3/resolve/main/vocab.txt
[OK ] Saved vocab.txt
[GET ] https://huggingface.co/alme9155/distilbert-sst5-finetuned-v3/resolve/main/model.safetensors
[OK ] Saved model.safetensors
[GET ] https://huggingface.co/alme9155/distilbert-sst5-finetuned-v3/resolve/main/training_args.bin
[OK ] Saved training_args.bin

> Task :components:messaging:compileKotlin
w: file:///Users/ameau/Downloads/temp-test/cscsa-5028-sentiment-analysis/components/messaging/src/main/kotlin/cu/cscsa5028/alme9155/messaging/ReviewsDataHandler.kt:207:21 This is a delicate API and its use requires care. Make sure you fully read and understand documentation of the declaration that is marked as a delicate API.
Java HotSpot(TM) 64-Bit Server VM warning: Sharing is only supported for boot loader classes because bootstrap classpath has been appended
Java HotSpot(TM) 64-Bit Server VM warning: Sharing is only supported for boot loader classes because bootstrap classpath has been appended
Java HotSpot(TM) 64-Bit Server VM warning: Sharing is only supported for boot loader classes because bootstrap classpath has been appended

Deprecated Gradle features were used in this build, making it incompatible with Gradle 9.0.

You can use '--warning-mode all' to show the individual deprecation warnings and determine if they come from your own scripts or plugins.

For more on this, please refer to https://docs.gradle.org/8.7/userguide/command\_line\_interface.html#sec:command\_line\_warnings in the Gradle documentation.

BUILD SUCCESSFUL in 2m 9s
65 actionable tasks: 52 executed, 13 up-to-date
```

```

superman:cscsa-5028-sentiment-analysis amea$ docker compose build --no-cache
[+] Building 19.0s (27/27) FINISHED
=> [data-collector internal] load build definition from Dockerfile
=> => transferring dockerfile: 1.90kB
=> [frontend-server internal] load build definition from Dockerfile
=> => transferring dockerfile: 1.90kB
=> [data-analyzer internal] load build definition from Dockerfile
=> => transferring dockerfile: 1.90kB
=> [data-collector internal] load metadata for docker.io/library/alpine:3.20
=> [data-analyzer internal] load metadata for docker.io/library/eclipse-temurin:21-alpine
=> [data-collector auth] library/eclipse-temurin:pull token for registry-1.docker.io
=> [data-collector auth] library/alpine:pull token for registry-1.docker.io
=> [frontend-server internal] load .dockerignore
=> => transferring context: 2B
=> [data-analyzer internal] load .dockerignore
=> => transferring context: 2B
=> [data-collector internal] load .dockerignore
=> => transferring context: 2B
=> CACHED [data-analyzer stage-1 1/7] FROM docker.io/library/alpine:3.20@sha256:765942a4039992336de8dd5db68058
=> CACHED [data-analyzer jre-build 1/2] FROM docker.io/library/eclipse-temurin:21-alpine@sha256:c4799f335a65b1
=> [frontend-server internal] load build context
=> => transferring context: 51.99MB
=> [data-collector stage-1 2/7] RUN apk add --no-cache ca-certificates tini
=> [data-collector internal] load build context
=> => transferring context: 51.85MB
=> [frontend-server jre-build 2/2] RUN /opt/java/openjdk/bin/jlink --add-modules java.base,java.logging,j
=> [data-analyzer internal] load build context
=> => transferring context: 51.83MB
=> [data-collector stage-1 3/7] RUN addgroup -g 1001 kotlin && adduser -S -D -u 1001 -G kotlin kotlin
=> [data-analyzer stage-1 4/7] COPY --from=jre-build /javaruntime /opt/java/openjdk
=> [frontend-server stage-1 5/7] WORKDIR /app
=> [data-analyzer stage-1 6/7] RUN mkdir -p /opt/applications && chown kotlin:kotlin /opt/applications
=> [data-analyzer stage-1 7/7] COPY applications/data-analyzer-server/build/libs/*.jar /opt/applications/app.j
=> [frontend-server stage-1 7/7] COPY applications/frontend-server/build/libs/*.jar /opt/applications/app.jar
=> [data-collector stage-1 7/7] COPY applications/data-collector-server/build/libs/*.jar /opt/applications/app
=> [data-analyzer] exporting to image
=> => exporting layers
=> => writing image sha256:58b57abd6854f0d6d7adce9dda81c7d3e53b519a0feff708a8a9a9975b0e6c71
=> => naming to docker.io/library/cscsa-5028-sentiment-analysis-data-analyzer
=> [data-collector] exporting to image
=> => exporting layers
=> => writing image sha256:7b309095b1f765b0d4fe8e49e2112086dd18ccfd9cbb55f9b93ed7b66e42cbce
=> => naming to docker.io/library/cscsa-5028-sentiment-analysis-data-collector
=> [frontend-server] exporting to image
=> => exporting layers
=> => writing image sha256:51be3fb795cccec14366baff2c86480cac8b3dc10e4cd842125a17be57e9765
=> => naming to docker.io/library/cscsa-5028-sentiment-analysis-frontend-server
superman:cscsa-5028-sentiment-analysis amea$ docker compose up -d
[+] Running 7/8
  . Network cscsa-5028-sentiment-analysis_default  Created          24.6s
  ✓ Container rabbitmq                         Healthy           13.7s
  ✓ Container mongodb                          Healthy           24.2s
  ✓ Container prometheus                      Started          1.6s
  ✓ Container grafana                          Started          2.2s
  ✓ Container sentiment-data-collector       Started          2.3s
  ✓ Container sentiment-analyzer-frontend    Started          2.4s
  ✓ Container sentiment-analyzer             Started          24.2s
superman:cscsa-5028-sentiment-analysis amea$ sh scripts/post-collect.sh
{"fetchedCount":100,"upsertedCount":100}superman:cscsa-5028-sentiment-analysis amea$ 
superman:cscsa-5028-sentiment-analysis amea$ sh scripts/post-publish.sh
{"dbCount":100,"msgCount":245}superman:cscsa-5028-sentiment-analysis amea$ 

```

Unit Test and Integration Test

- Unit Test are executed as part of build process, or using the command.

> ./gradlew test

- Assuming source code are built and all the docker containers are running, then execute the gradlew integration test

> ./gradlew :integration-test:workflowEndToEndTest

```
superman:cscsa-5028-sentiment-analysis amea$ ./gradlew test
[...]
Deprecated Gradle features were used in this build, making it incompatible with Gradle 9.0.
You can use '--warning-mode all' to show the individual deprecation warnings and determine if they come from your own
scripts or plugins.

For more on this, please refer to https://docs.gradle.org/8.7/userguide/command_line_interface.html#sec:command_line_w
arnings in the Gradle documentation.

BUILD SUCCESSFUL in 4s
43 actionable tasks: 43 up-to-date
superman:cscsa-5028-sentiment-analysis amea$ ./gradlew :integration-test:workflowEndToEndTest
[...]
Deprecated Gradle features were used in this build, making it incompatible with Gradle 9.0.
You can use '--warning-mode all' to show the individual deprecation warnings and determine if they come from your own
scripts or plugins.

For more on this, please refer to https://docs.gradle.org/8.7/userguide/command_line_interface.html#sec:command_line_w
arnings in the Gradle documentation.

BUILD SUCCESSFUL in 5s
25 actionable tasks: 1 executed, 24 up-to-date
superman:cscsa-5028-sentiment-analysis amea$ [...]
```

Product environment

- Project has been deployed to Google cloud.
- Web form and report can be accessed <http://34.94.189.67/>

(Note: ONLY HTTP protocol is supported)

Monitoring

- Prometheus can be used to track service health: <http://localhost:9090/targets>
- Grafana can be login using (admin / password): <http://localhost:3000>

Alexander Meau (alme9155)

The screenshot shows the Prometheus web interface at <http://34.94.189.67:9090/targets>. The top navigation bar includes links for Query, Alerts, and Status > Target health. The main content area displays three service entries:

- data-analyzer**: Endpoint <http://data-analyzer:8080/metrics>, Labels: instance="data-analyzer:8080", job="data-analyzer".
- data-collector**: Endpoint <http://data-collector:8080/metrics>, Labels: instance="data-collector:8080", job="data-collector".
- frontend-server**: Endpoint <http://frontend-server:8080/metrics>, Labels: instance="frontend-server:8080", job="frontend-server".

