

Reddiment: Reddit Sentiment-Analyse

Technical Report

Tobias Bauer
t.bauer@oth-aw.de

Fabian Beer
f.beer1@oth-aw.de

Daniel Holl
d.holl1@oth-aw.de

Ardian Imeraj
a.imeraj@oth-aw.de

Konrad Schweiger
k.schweiger@oth-aw.de

Philipp Stangl
p.stangl1@oth-aw.de

Wolfgang Weigl
w.weigl@oth-aw.de

Zusammenfassung—Dieser Technical Report beschreibt die Architektur von Reddiment – ein webbasiertes Dashboard zur Sentiment-Analyse von Subreddits.

I. EINFÜHRUNG UND ZIELE

In dem Subreddit `r/wallstreetbets`, auch bekannt als WallStreetBets oder WSB, wird über Aktien- und Optionshandel spekuliert. Er ist bekannt für seine profane Art und die Vorwürfe, dass Nutzer/innen Wertpapiere manipulieren und volatile Kursbewegungen auslösen. Anhand von Sentiment-Analyse sollen nun Subreddits in Bezug auf Aktienkursverläufe analysiert werden. Dazu soll ein webbasiertes Dashboard entwickelt werden, welches den zeitlichen Verlauf von Sentiment und Erwähnungen benutzerdefinierter Schlüsselwörter in ausgewählten Subreddits dem Aktienverlauf gegenüberstellt.

In den weiteren Abschnitten des Technical Reports wird zuerst die Bausteinsicht des Gesamtsystems in Abschnitt II eingegangen. Im Nächsten Abschnitt III wird die Verteilungssicht der Anwendung beschrieben. In Abschnitt IV werden die angewandten Werkzeuge zur Entwicklung der Anwendung vorgestellt. Abschließend wird kurz das Sicherheitskonzept in Abschnitt V vorgestellt und ein Fazit in Abschnitt VI gegeben.

II. BAUSTEINSICHT

Diese Sicht zeigt die statische Zerlegung des Systems in Bausteine sowie deren Beziehungen.

A. Gesamtsystem

Reddiment bezieht Daten aus mehreren Quellen und stellt diese dem Benutzer aggregiert bereit. Die folgende Abbildung 1 zeigt die Interaktionen des Systems mit Fremdsystemen und dem Benutzer. Die beiden Datenquellen für Reddiment sind

- Reddit-API für Subreddit-Daten und
- Yahoo Finance für Aktienmarkt-Daten.

B. Backend

Das Backend wurde als ein unter „Node.js“ [1] laufender Server realisiert. Hierfür kam das Web-Framework „Express“ [2] zum Einsatz. Darauf aufbauend wurde eine GraphQL-Schnittstelle [3] mithilfe des Frameworks „Apollo Server“ [4] implementiert. Die GraphQL-API ist unter der Route `/graphql` verwendbar.

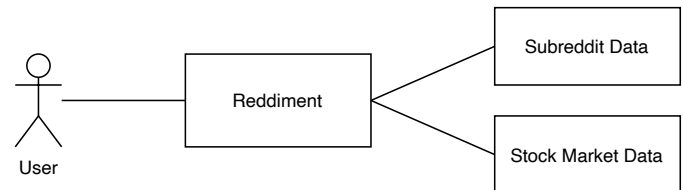


Abbildung 1. Kontextabgrenzung des Reddiment Gesamtsystems

Das Backend stellt die zentrale Kommunikations- und Verwaltungseinheit für die einzelnen Komponenten des Gesamtsystems Reddiment dar. So liefern die proaktiven Crawler (siehe Abschnitt II-D) die Rohdaten über eine GraphQL-Mutation an das Backend. Diese Rohdaten werden anschließend weiterverarbeitet und im Falle von Reddit-Kommentaren durch Aufruf des Sentiment-Services um den Sentiment-Wert erweitert. Alle so erhaltenen Daten werden über die Anbindung zur Datenbank (siehe Abschnitt II-C) persistiert.

Stellt nun das Frontend eine Anfrage an das Backend, werden diese aus der Datenbank geladen und aufbereitet. Die Gruppierung und Aggregation von beispielsweise Reddit-Kommentaren sowie die Interpretation des Sentiment-Werts erfolgt in der Suchanfrage an die Datenbank. Damit wird die rechenaufwendige Aggregation direkt an der Quelle der Daten ausgeführt und gleichzeitig Bandbreite eingespart. Die Konfiguration im Backend hingegen ermöglicht es uns, zukünftig eine feinere oder auch gröbere Zusammenfassung von Sentiment-Daten ohne Datenverlust oder aufwendigen Transformationen vorzunehmen.

Aufgrund des Einsatzes von GraphQL als Schnittstelle können Frontend und ggf. andere (externe) Akteure die benötigten Daten in einer von ihnen gewählten Struktur abfragen. Dies erlaubt eine flexible Entwicklung des Frontends ohne Kommunikations-Overhead. Weiterhin sind serverseitig die sogenannten Resolver modular aufgebaut, sodass verschiedene GraphQL-Felder in getrennten Quellcodedateien behandelt werden können. Alle Quellcodedateien im Zusammenhang mit GraphQL befinden sich im Verzeichnis `backend/src/graphql`.

Das Backend ist aufgrund der zentralen Position im System Reddiment auf die anderen Bestandteile angewiesen.

Um dennoch Unit-Tests für ein isoliertes Backend schreiben und ausführen zu können, können die Datenbank und das Sentiment-Modul dynamisch durch Mocks ersetzt werden. Beim Start des Backends kann über die Umgebungsvariable `PRODUCTION` festgelegt werden, ob die realen Module (`true`) oder die Mock-Module (`false`) verwendet werden. Die entsprechenden Quellcode-dateien liegen im Ordner `backend/src/services`. Da die Crawler proaktiv sind, müssen für diese keine Mocks erstellt werden – vielmehr simulieren die Unit-Tests die Crawler und prüfen, ob die gewünschten Daten von der GraphQL-Schnittstelle zurückgeliefert werden.

Im Ordner `backend/src/util` befinden sich verschiedene Quellcode-dateien mit Hilfsfunktionen. Gestartet wird das Backend durch Ausführen der Datei `backend/server.ts`, die wiederum alle benötigten Module lädt und den Apollo Server mitsamt der GraphQL-Schnittstelle startet.

C. Datenbank

In der NoSQL-Datenbank „Elasticsearch“ [5] werden Daten in Form von Dokumenten gespeichert. Für dieses Projekt werden Daten, welche an das Backend gesendet werden, als Dokumente persistent in der Datenbank gespeichert. Die Reddit-Daten bestehen aus den neun Feldern in Tabelle I.

Tabelle I
REDDIT-DATEN

Feld	Beschreibung
<code>subreddit</code>	Name des Subreddits
<code>text</code>	Kommentar als String
<code>timestamp</code>	Zeitstempel des Kommentars
<code>commentId</code>	Reddit-Kommentar-ID
<code>userId</code>	Reddit-User-ID
<code>articleId</code>	Reddit-Post-ID
<code>upvotes</code>	Anzahl Positiv-Bewertung in Reddit
<code>downvotes</code>	Anzahl Negativ-Bewertung in Reddit
<code>sentiment</code>	Ausgabe der Sentiment-Analyse

Die Aktien-Daten setzen sich aus den acht Feldern in Tabelle II zusammen.

Tabelle II
REDDIT-DATEN

Feld	Beschreibung
<code>stock</code>	Name der Aktie
<code>timestamp</code>	Zeitstempel der Werte
<code>open</code>	Wert bei Tagesbeginn
<code>high</code>	Maximaler Tages-Wert
<code>low</code>	Minimaler Tages-Wert
<code>close</code>	Wert bei Tagesabschluss
<code>adjClose</code>	Dividenden-adjustierter Tagesabschluss-Wert
<code>volume</code>	Marktvolumen der gehandelten Aktie

Um die Daten in der Datenbank zu organisieren, arbeitet Elasticsearch mit Indizes. Dies bedeutet, dass ähnlich strukturierte Dokumente unter demselben Index gespeichert werden. Pro Subreddit gibt es einen Index, dessen Name mit `r_` präfigiert ist, und pro Aktie gibt es einen Index mit Präfix `f_`. Um einen zeitlichen Verlauf einer Aktie oder

des Sentiments zu erzeugen, müssen alle Dokumente eines Index abgerufen werden. Dieses Vorgehen ermöglicht es, dass bei einer Suchanfrage nicht alle Dokumente aller Indizes durchsucht werden müssen, sondern eine Vorauswahl auf eine Untermenge getroffen werden kann.

Elasticsearch verwendet eine global eindeutige ID für jedes Dokument. Bei Dokumenten, welche Reddit-Daten enthalten, wird die ID des Dokumentes gleich der eindeutigen `commentId` gesetzt. Den Dokumenten mit Aktien-Daten wird eine ID nach folgendem Schema zugewiesen: `Aktiename_Zeitstempel`. Dokumente können auch anhand ihrer ID gesucht und ausgegeben werden. Dies wird beispielsweise bei der Aktualisierung von Dokumenten angewendet.

Im Backend wird das Paket „Elasticsearch Node.js client“ [6] verwendet. Das Paket enthält einen Client, der eine Verbindung zur Datenbank herstellt. In der TypeScript-Datei `backend/src/services/database.ts` befindet sich der Quellcode zur Kommunikation mit der Datenbank.

D. Dienste

Dieser Abschnitt beschreibt die Dienste (engl. Services) des Gesamtsystems Reddiment. Es gibt drei Dienste: Sentiment, Reddit-Crawler und Stock-Market-Crawler.

1) *Sentiment*: Dieser Dienst ermittelt für einen Text das Sentiment. Es wurde eine REST-API mit dem Endpunkt `/sentiment` implementiert. An diesen Endpunkt kann mittels einer POST-Anfrage ein Text übertragen werden. Der Dienst ermittelt anschließend die Stimmungslage des Textes und gibt das Ergebnis im JSON-Format zurück. Die Ermittlung des Sentiments erfolgt durch zwei regelbasierte Verfahren, um eine höhere Aussagekraft zu haben. Das erste regelbasierte Verfahren ist „vader“ [7] (Valence Aware Dictionary and Sentiment Reasoner) der Python-Bibliothek *nltk*. Das zweite Verfahren ist „TextBlob“ [8]. Sind sich die beiden Verfahren in der Auswertung einig, wird der Sentiment-Wert von vader zurückgegeben. Andernfalls wird 0 (neutral) zurückgegeben.

2) *Reddit-Crawler*: Der Reddit-Crawler verwendet das Paket „snoowrap“ [9], das JavaScript-Funktionen für den Zugriff auf die Reddit-API bereitstellt. „snoowrap“ unterliegt den API-Regeln von Reddit, worin u.a. das Rate-Limit bei 60 Anfragen pro Minute liegt. Um die Reddit-API nutzen zu können, muss ein API-Key über das Benutzerkonto¹ angefordert werden. Ein API-Key besteht aus den zwei Teilen: einer *Client-ID* und einem *Client-Secret*.

Mit dem API-Key und den Anmeldedaten für das zugehörige Reddit-Konto kann ein Objekt der Klasse `Snoowrap` erstellt werden. Mit diesem Objekt können Reddit-API-Aufrufe durchgeführt werden.

Der Reddit-Crawler stellt in einem zeitlichen Intervall eine Anfrage an das Backend und erhält eine Liste mit Subreddit-Namen. Kommentare dieser Subreddits sollen von Reddit abgefragt und an das Backend gesendet werden. Das Sammeln

¹<https://www.reddit.com/prefs/apps>

der Kommentare erfolgt zyklisch, um das Rate-Limit nicht zu überschreiten. Pro Sammelzyklus werden die Kommentare an das Backend übermittelt.

3) *Stock-Market-Crawler*: Um den Verlauf des Sentiments in der Vergangenheit auch mit der tatsächlichen Marktlage vergleichen zu können, gibt es einen reaktiven Stock-Market-Crawler mit einer REST-API für Marktdaten. Dieser hat den Endpunkt `/post` und erwartet ein valides *Ticker-symbol* (Aktienkürzel) einer Aktie. Der zugehörige Aktienkurs wird von *Yahoo-Finance* abgefragt und zurückgegeben.

E. Frontend

Dieser Abschnitt beschreibt die client-seitige Frontend-Architektur. Das Frontend wird unter Zuhilfenahme des Frontend-Frameworks „SvelteKit“ [10] realisiert. Das Frontend besteht aus zwei Bausteinen: 1) den Routen zur Navigation, und 2) der *Library* für Komponenten und weitere Module.

1) *Routen*: Die Routen zur Navigation sind durch Svelte-Komponenten in `src/routes` festgelegt. Es gibt Elemente, die auf jeder Seite sichtbar sind, z. B. die Navigationsleiste oder eine Fußzeile. Anstatt diese für jede Seite neu zu definieren, wird eine Layout-Komponente namens `src/routes/__layout.svelte` verwendet.

2) *Library*: Die *Library* (`src/lib`) ist eine Sammlung von Frontend-Komponenten, mit denen die eigentliche Anzeige im Webbrowser realisiert wird. Die Komponenten behandeln alle Eingaben und kommunizieren bei Bedarf mit dem Backend über die GraphQL API. Für die Kommunikation mit der GraphQL API wird im Frontend „KitQL“ [11] verwendet, das einen client-seitigen GraphQL Client bereitstellt.

III. VERTEILUNGSSICHT

Zentraler Bestandteil der Verteilungsstruktur sind Docker-Container [12]. Jeder Baustein von Reddiment ist für sich isoliert in einem eigenen Docker-Container untergebracht. Damit das gesamte System mit allen verteilten Komponenten in der korrekten Reihenfolge gestartet wird, werden die Docker-Container mit Docker Compose orchestriert. Docker Compose übernimmt die Netzwerkkonfiguration und die Vergabe von Host-Namen an die jeweiligen Docker-Container. Darüber hinaus können Umgebungsvariablen verwaltet werden. Sensible Daten, wie Zugangspasswörter, werden durch spezielle Mechanismen (siehe Abschnitt V) sicher zur Verfügung gestellt.

IV. ENTWICKLUNGSWERKZEUGE

In diesem Abschnitt wird auf die verwendeten Entwicklungswerkzeuge genauer eingegangen.

A. Paketverwaltung

Die Verwaltung der Abhängigkeiten erfolgt mit „npm“ [13] für auf „Node.js“ basierende Bausteine. Für den Sentiment-Baustein wird „Pipenv“ [14] verwendet.

B. Linting

Im Frontend wird „eslint“ [15] in Verbindung mit „prettier“ [16] verwendet, um die Einhaltung der Codierichtlinien zu gewährleisten. Die Konfigurationen sind jeweils in den Dateien `eslinttrc.js` und `prettierrc.js` hinterlegt.

C. Build-Tools

1) *Backend*: Sowohl das Backend als auch die Crawler werden mit dem TypeScript Compiler „tsc“ [17] in ein für „Node.js“ ausführbares Format überführt. Die Konfiguration befindet sich dabei in der Datei `tsconfig.json`.

2) *Frontend*: Im Frontend ist Vite dafür zuständig, die Anwendung aus dem Quellcode zu erstellen. Dabei gibt es zwei Varianten: Für Entwicklungszwecke wird ein Vite-Dev-Server (mit Reload-Funktionalität) zum Bereitstellen der Anwendung verwendet. Für den Produktiveinsatz werden nur die benötigten Zielformate unter Verwendung des `Static adapter` erstellt, die dann mit einer beliebigen Server-Software ausgeliefert werden können.

D. Unit-Tests

Im Backend werden Unit-Tests mit „mocha“ [18] ausgeführt. Dabei wird „istanbul“ [19] für die Erzeugung der Test-Abdeckung verwendet.

Im Frontend wird „Vitest“ [20] verwendet. Für die Frontend-Komponenten wird zusätzlich die „svelte-testing-library“ [21] eingesetzt. Diese ermöglicht es, die Komponenten zu *rendern*, um Details über die verschiedenen Elemente innerhalb der Komponente zu erhalten.

V. SECRETS-VERWALTUNG

Die Anforderung, sensible Daten jeder Anwendung individuell und zusätzlich einfach konfigurierbar zur Verfügung zu stellen, lässt schnell auf den Einsatz von Umgebungsvariablen schließen. Ein solches Vorgehen birgt jedoch einige sicherheitsrelevante Schwachstellen. Aus diesem Grund werden die komponentenspezifischen Zugangsdaten mittels dem Schlüsselwort `secrets` in der Docker-Compose-Konfigurationsdatei den Docker-Containern zur Verfügung gestellt. Dabei werden im Gegensatz zu Umgebungsvariablen die Passwörter in einer Text-Datei gespeichert und in dem jeweiligen Docker-Container eingebunden.

Ein weiteres zu beachtendes Kriterium in diesem Sicherheitskonzept ist der Zugriff auf die jeweiligen API der einzelnen Komponenten. Beispielsweise sollen nur die Crawler die GraphQL-Mutations des Backends verwenden dürfen. Dafür wurde für jede Anwendung ein individueller „Access-Key“ erstellt. Der Zugriff wird durch Prüfen des Schlüssels bei jeder Anfrage gewährt bzw. unterbunden. Unberechtigte Dritte können dadurch die API nicht verwenden und ein Missbrauch wird dadurch erschwert. Aufgrund zeitlicher Einschränkungen konnte das Sicherheitskonzept nicht vollständig implementiert werden.

VI. FAZIT UND AUSBLICK

Reddiment ist ein webbasiertes Dashboard, welches den zeitlichen Verlauf von Sentiment und Erwähnungen bestimmter Schlüsselwörter in ausgewählten Subreddits dem Aktienverlauf gegenüberstellt. Durch die Auswahl eines Subreddits und eines oder mehrere Schlüsselwörter werden die Diagramme automatisiert mit den vorhandenen Daten befüllt. Durch die strikte Trennung von Front- und Backend besteht

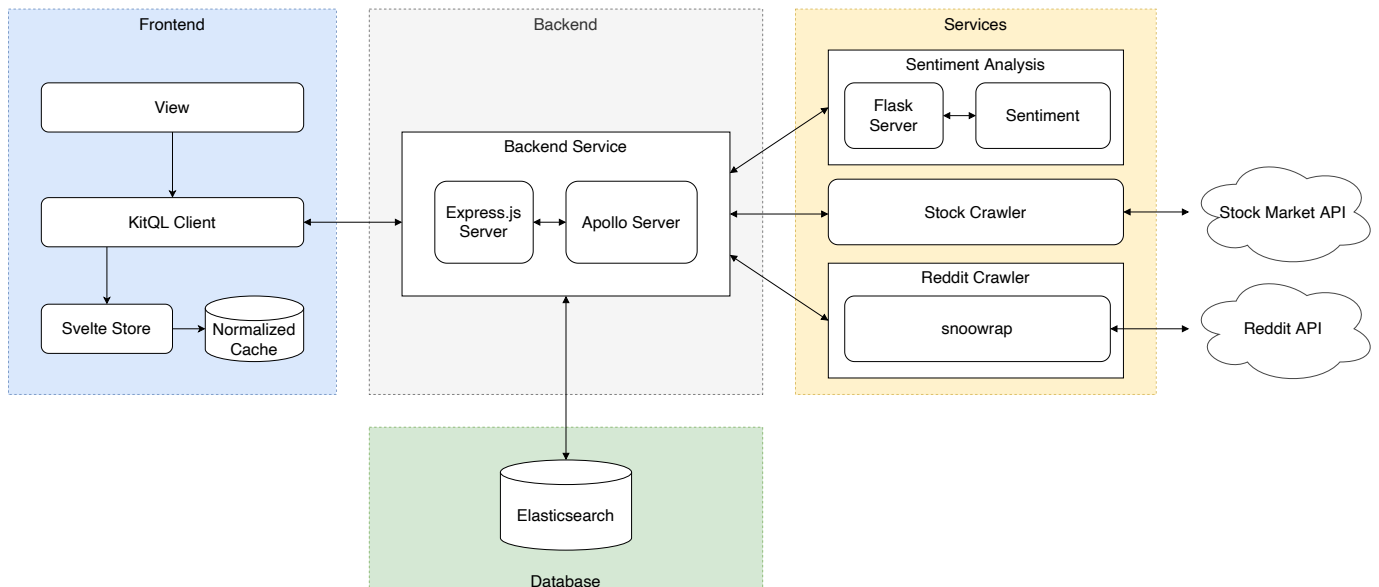


Abbildung 2. Überblick über die Architektur von Reddiment. Die Architektur besteht aus vier Teilen: Das Frontend bietet dem Nutzer ein graphisches Dashboard zur Visualisierung der Metriken (Abschnitt II-E), das Backend stellt (Abschnitt II-B) die GraphQL API bereit, die Datenbank (Abschnitt II-C) speichert persistent Reddit und Aktienmarkt Daten, und den Diensten (Abschnitt II-D) zur Sentiment Analyse als auch zum crawlen der Daten von APIs.

die Möglichkeit, verwendete APIs zu ersetzen oder Weitere einzubinden.

LITERATUR

- [1] Node. "Node.js." [Online]. (2022), Adresse: <https://www.nodejs.org/>.
- [2] Express. "Express." [Online]. (2022), Adresse: <https://expressjs.com/>.
- [3] The GraphQL Foundation. "GraphQL: A query language for your API." [Online]. (2022), Adresse: <https://graphql.org/>.
- [4] Apollo Graph Inc. "Introduction to Apollo Server – Apollo GraphQL Docs." [Online]. (2022), Adresse: <https://www.apollographql.com/docs/apollo-server/>.
- [5] Elastic. "Elasticsearch." [Online]. (2022), Adresse: <https://www.elastic.co/de/elasticsearch/>.
- [6] Elastic. "Elasticsearch Node.js client." [Online]. (2022), Adresse: <https://www.npmjs.com/package/@elastic/elasticsearch>.
- [7] C. Hutto und E. Gilbert, "Vader: A parsimonious rule-based model for sentiment analysis of social media text," in *Proceedings of the international AAAI conference on web and social media*, Bd. 8, 2014, S. 216–225.
- [8] TextBlob. "TextBlob: Simplified Text Processing." [Online]. (2022), Adresse: <https://textblob.readthedocs.io/en/dev/quickstart.html#sentiment-analysis>.
- [9] T. Katz. "snoowrap: A fully-featured JavaScript wrapper for the reddit API." [Online]. (2022), Adresse: <https://github.com/not-an-aardvark/snoowrap>.
- [10] Svelte. "Svelte Kit: The fastest way to build Svelte Apps." [Online]. (2022), Adresse: <https://kit.svelte.dev/>.
- [11] The Guild. "KitQL: GraphQL Framework for Svelte-Kit." [Online]. (2022), Adresse: <https://www.kitql.dev/>.
- [12] Docker. "Docker." [Online]. (2022), Adresse: <https://www.docker.com/>.
- [13] npm. "npm." [Online]. (2022), Adresse: <https://www.npmjs.com/>.
- [14] Python Packaging Authority. "Pipenv: Python Dev Workflow for Humans." [Online]. (2022), Adresse: <https://pipenv.pypa.io/en/latest/>.
- [15] N. Zakas. "ESLint: Pluggable JavaScript linter." [Online]. (2022), Adresse: <https://eslint.org/>.
- [16] Prettier. "Prettier: Opinionated Code Formatter." [Online]. (2022), Adresse: <https://prettier.io/>.
- [17] Microsoft. "TypeScript: JavaScript With Syntax For Types." [Online]. (2022), Adresse: <https://www.typescriptlang.org/>.
- [18] Mochajs. "Mocha: JavaScript test framework running on Node.js and in the browser." [Online]. (2022), Adresse: <https://mochajs.org/>.
- [19] Istanbul. "Istanbul: JavaScript test coverage made simple." [Online]. (2022), Adresse: <https://istanbul.js.org/>.
- [20] Vitest. "Vitest: A Vite-native unit test framework." [Online]. (2022), Adresse: <https://vitest.dev/>.
- [21] Testing Library. "Svelte Testing Library." [Online]. (2022), Adresse: <https://github.com/testing-library/svelte-testing-library>.