

Projektarbeit: OpenWeatherMap-Streaming

Datum: 26.05.2025

Autor: Philippe Christen

Weiterbildung: CAS Data Engineering an der FHNW

Der vollständige Quellcode inklusive Docker-Konfiguration und einer Beispieldatei «.env.example» ist im zugehörigen Git-Repository verfügbar. Sensiblen Konfigurationswerte (z. B. API-Schlüssel) sind darin bewusst nicht enthalten und müssen vom Benutzer lokal in einer .env-Datei ergänzt werden. Die Umgebungsvariable für den Zugriff auf die OpenWeatherMap-API ist dabei wie folgt zu definieren: `API_KEY="DEIN_API_KEY"`

Problemstellung (Ist-Zustand)

Im Kontext moderner IoT- und Smart-City-Anwendungen besteht ein wachsender Bedarf, Wetterdaten in nahezu Echtzeit zu analysieren. Die OpenWeatherMap-API erlaubt das periodische Abfragen von Wetterdaten. Es fehlt jedoch an standardisierten, containerisierten Pipelines zur robusten Verarbeitung und Speicherung dieser Datenströme. Herausforderungen ergeben sich insbesondere durch das asynchrone Eintreffen der Daten (verschiedene Eventzeiten je Stadt), das Handling unvollständiger Datensätze sowie die Notwendigkeit, Datenverluste und Duplikate zu vermeiden. Zudem besteht ein Bedarf, meteorologische Auffälligkeiten automatisiert zu erkennen und diese für spätere Analysen bereitzustellen.

Ein zusätzlicher Aspekt betraf die Limitierung und Pflege der externen Datenquelle: OpenWeatherMap definiert im kostenfreien API-Zugang (Free Access Plan) klare Nutzungsbeschränkungen. Neben einem Monatslimit von 1.000.000 API-Aufrufen dürfen:

- pro Standort maximal ein API-Aufruf alle 10 Minuten erfolgen (entspricht der Aktualisierungsfrequenz des OpenWeather-Modells),
- systemweit maximal 60 API-Aufrufe pro Minute abgesetzt werden (Rate-Limit pro API-Key).

Diese Begrenzungen dienen der Lastverteilung und dem Schutz der Infrastruktur von OpenWeatherMap.

Quellen: [OpenWeatherMap Pricing](#), [API Guidelines](#)

Daher wurde in der Umsetzung ein Abrufintervall von 600 Sekunden (10 Minuten) pro Stadt konfiguriert, um regelkonform zu bleiben und einen stabilen Betrieb sicherzustellen.

New Products	Services	API keys	Billing plans	Payments	Block logs	My orders	My profile	Ask a question
Name	Description	Price plan	Limits	Details				
Weather	Current weather and forecast	Free plan	Hourly forecast: unavailable Daily forecast: unavailable Calls per minute: 60 3 hour forecast: 5 days	view				

Abb.1: Verwendeter OpenWeatherMap Price-Plan, Quelle: Philippe Christen

Zielarchitektur (Soll-Zustand)

Ziel ist der Aufbau einer containerisierten Streaming-Architektur zur Erfassung, Anreicherung, Speicherung und Visualisierung von Wetterdaten in nahezu Echtzeit. Die Architektur basiert auf einem modularen Aufbau:

- OpenWeatherMap API: Quelle der Sensordaten (Polling)
- Kafka: Messaging Layer zur Entkopplung der Verarbeitung
- Apache Beam: Verarbeitungspipeline mit Filterung, Standardisierung, Anomalieerkennung und Persistenz
- PostgreSQL: Zentrale relationale Datenhaltung für strukturierte Auswertungen
- Grafana: Visualisierung mit Dashboards und Zeitreihenanalysen

Die gesamte Architektur ist mit Docker Compose orchestriert. Zeitstempel (event_time, received_time, processing_time) werden zur Nachvollziehbarkeit aller Verarbeitungsschritte in der Datenbank gespeichert. Duplikate werden durch ein kombiniertes Schlüsselattribut (UNIQUE(city, timestamp)) und ON CONFLICT DO NOTHING verhindert.

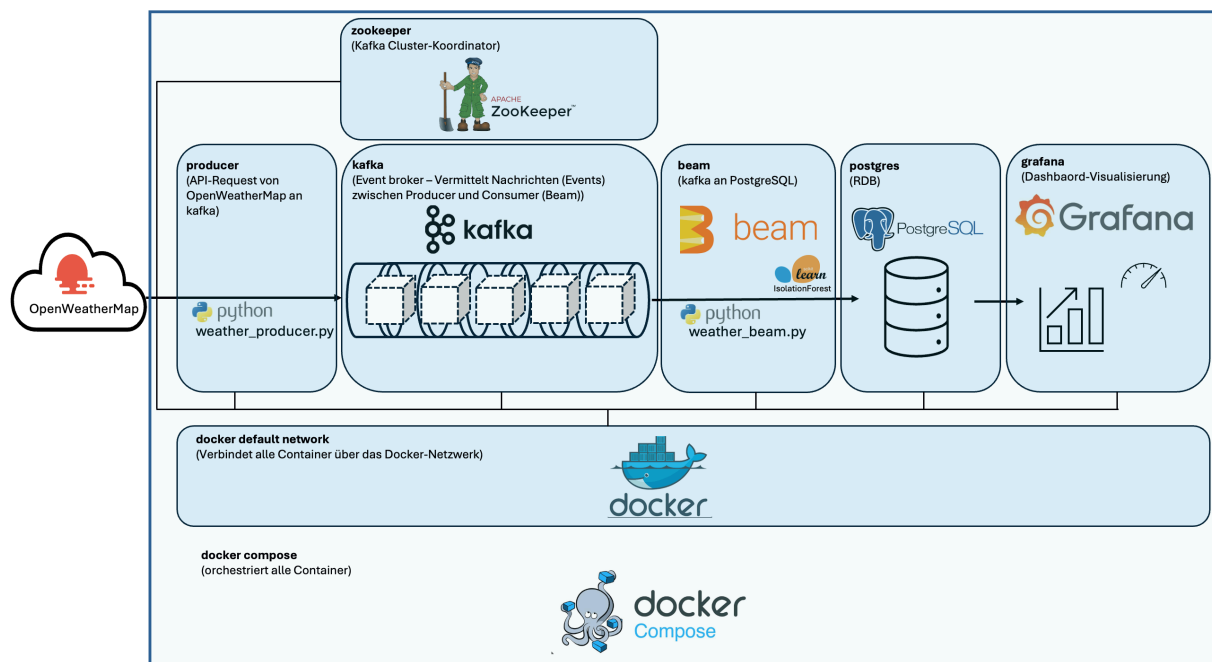


Abb.2: Zielarchitektur, Quelle: Philippe Christen

Technologie-Entscheidungen

Apache Beam wurde gewählt, um eine erweiterbare, portierbare und testbare Streaming-verarbeitung zu ermöglichen. Beam trennt klar zwischen Pipeline-Definition und Ausführung (Runner), was lokale Tests mit dem DirectRunner und später Cloud-Migrationen (z. B. Google Cloud DataFlow) erleichtert.

Kafka dient als Puffer und Messaging-System. Durch die Verwendung von Topics wird eine lose Kopplung zwischen Producer und Beam-Consumer erreicht. Nachrichten können mehrfach gelesen oder reprocessed werden.

PostgreSQL bietet eine robuste, relationale Speicherung, ist mit Grafana integrierbar und erlaubt komplexe Zeitreihenanalysen. Funktionen wie `date_trunc`, Aggregationen und Filter unterstützen flexible Auswertungen und Heatmaps.

Grafana wurde als Visualisierungsplattform gewählt, da es eine direkte Anbindung an PostgreSQL erlaubt, interaktive Dashboards unterstützt und für Monitoring- und Analysezwecke im DevOps-Umfeld etabliert ist. Zudem ermöglicht der JSON-Export ein schnelles Reproduzieren der Visualisierungen in anderen Umgebungen.

Python wurde als Sprache aufgrund der Bibliotheken `kafka-python`, `psycpg2`, `apache-beam`, `joblib`, `pandas` und `dotenv` verwendet. Alle Services werden via Docker containerisiert und gemeinsam orchestriert.

```
[+] Running 7/7
  ✓ Network openweathermap-bigdata-project_default          Created
  ✓ Container openweathermap-bigdata-project-kafka-1        Started
  ✓ Container openweathermap-bigdata-project-grafana-1       Started
  ✓ Container openweathermap-bigdata-project-postgres-1      Started
  ✓ Container openweathermap-bigdata-project-zookeeper-1     Started
  ✓ Container openweathermap-bigdata-project-producer-1      Started
  ✓ Container openweathermap-bigdata-project-beam-1          Started
  ○ (base) philippe-christen:openweathermap-bigdata-project philippe
  ○ (base) philippe-christen:openweathermap-bigdata-project philippe
  ● (base) philippe-christen:openweathermap-bigdata-project philippe docker ps
```

CONTAINER ID	IMAGE	COWARD	CREATED	STATUS	PORTS	NAMES
093d7b92bb9b	confluentinc/cp-kafka:7.5.0		53 minutes ago	Up 53 minutes	9092/tcp	cool_heyrovsky
afdc2d08d1d9	openweathermap-bigdata-project-beam	"python weather_beam..."	53 minutes ago	Up 53 minutes		openweathermap-bigdata-project-beam-1
ad857d62943f	openweathermap-bigdata-project-producer	"python weather_prod..."	53 minutes ago	Up 53 minutes		openweathermap-bigdata-project-producer-1
9358f4ebcf9f2	postgres:14	"docker-entrypoint.su..."	53 minutes ago	Up 53 minutes	0.0.0.0:5433→5433/tcp	openweathermap-bigdata-project-postgres-1
afc96543da35	confluentinc/cp-zookeeper:7.5.3	"etc/confluent/dock..."	53 minutes ago	Up 53 minutes	2888/tcp, 0.0.0.0:2181→2181/tcp, 3888/tcp	openweathermap-bigdata-project-zookeeper-1
5b95769dec1	confluentinc/cp-kafka:7.5.3	"etc/confluent/dock..."	53 minutes ago	Up 53 minutes	0.0.0.0:9092→9092/tcp	openweathermap-bigdata-project-kafka-1
9e77d017d12b	grafana/grafana:10.0.3	"/run.sh"	53 minutes ago	Up 53 minutes	0.0.0.0:3000→3000/tcp	openweathermap-bigdata-project-grafana-1

Abb.3: Container-Umgebung, Quelle: Philippe Christen

```
{ "event_time": 1748179955, "received_time": 1748179955 }
INFO:kafka.com-drokersconnection client_id=kafka-python-producer-2, node_id=1 host=kafka:9092 connected [IPv4 ('172.21.0.2', 9092)]: connecting to kafka:9092 [IPv4 ('172.21.0.2', 9092)]: Connection complete.
INFO:kafka.com-drokersconnection client_id=kafka-python-producer-2, node_id=bootstrap-kafka:9092 connected [IPv4 ('172.21.0.2', 9092)]: closing connection.
INFO:_main__Send to Kafka (city: 'Leipzig', temperature: 18.74, humidity: 71, pressure: 1016, wind_speed: 5.36, cloud_coverage: 100, weather_main: 'Clouds', weather_description: 'Bedeckt', 'timestamp': 1748179956, 'lon': 8.175, 'lat': 47.385, 'sys_country': 'C
INFO:_main__Send to Kafka (city: 'Aarau', temperature: 18.88, humidity: 69, pressure: 1016, wind_speed: 3.88, cloud_coverage: 100, weather_main: 'Clouds', weather_description: 'Bedeckt', 'timestamp': 1748179734, 'lon': 8.042, 'lat': 47.3925, 'sys_country': 'C
INFO:_main__Send to Kafka (city: 'Hausen AG', temperature: 19.43, humidity: 74, pressure: 1016, wind_speed: 3.48, cloud_coverage: 100, weather_main: 'Clouds', weather_description: 'Bedeckt', 'timestamp': 1748179958, 'lon': 8.2099, 'lat': 47.464, 'sys_country': 'C
INFO:_main__Send to Kafka (city: 'Zürich', temperature: 19.77, humidity: 69, pressure: 1016, wind_speed: 3.55, cloud_coverage: 100, weather_main: 'Clouds', weather_description: 'Bedeckt', 'timestamp': 1748179959, 'lon': 8.2184, 'lat': 47.479, 'sys_country': 'C
INFO:_main__Send to Kafka (city: 'Zürich', temperature: 18.8, humidity: 66, pressure: 1016, wind_speed: 5.14, cloud_coverage: 75, weather_main: 'Clouds', weather_description: 'Überwiegend bewölkt', 'timestamp': 1748179858, 'lon': 8.55, 'lat': 47.3667, 'sys_coun
INFO:_main__Send to Kafka (city: 'Zürich', temperature: 13.81, humidity: 69, pressure: 1017, wind_speed: 3.89, cloud_coverage: 75, weather_main: 'Rain', weather_description: 'Leichter Regen', 'timestamp': 1748179674, 'lon': 8.4813, 'lat': 46.8211, 'sys_coun
INFO:_main__Send to Kafka (city: 'Sitten', temperature: 20.85, humidity: 40, pressure: 1016, wind_speed: 2.86, cloud_coverage: 75, weather_main: 'Clouds', weather_description: 'Überwiegend bewölkt', 'timestamp': 1748179780, 'lon': 7.3594, 'lat': 46.2291, 'sys
INFO:_main__Send to Kafka (city: 'Sitten', temperature: 18.84, humidity: 69, pressure: 1016, wind_speed: 1.34, cloud_coverage: 100, weather_main: 'Clouds', weather_description: 'Bedeckt', 'timestamp': 1748179963, 'lon': 8.3421, 'lat': 47.3511, 'sys_country': 'C
INFO:_main__Send to Kafka (city: 'Lugano', temperature: 22.04, humidity: 58, pressure: 1015, wind_speed: 3.6, cloud_coverage: 20, weather_main: 'Clouds', weather_description: 'Ein paar Wolken', 'timestamp': 1748179806, 'lon': 8.96, 'lat': 46.8101, 'sys_country': 'C
INFO:_main__Send to Kafka (city: 'Winterthur', temperature: 18.83, humidity: 62, pressure: 1016, wind_speed: 5.14, cloud_coverage: 75, weather_main: 'Clouds', weather_description: 'Überwiegend bewölkt', 'timestamp': 1748179719, 'lon': 8.75, 'lat': 47.5, 'sys_c
INFO:_main__Send to Kafka (city: 'Biel', temperature: 18.32, humidity: 67, pressure: 1016, wind_speed: 3.6, cloud_coverage: 75, weather_main: 'Clouds', weather_description: 'Überwiegend bewölkt', 'timestamp': 1748179967, 'lon': 7.2441, 'lat': 47.1324, 'sys_cou
INFO:_main__Send to Kafka (city: 'Zürich', temperature: 18.74, humidity: 71, pressure: 1016, wind_speed: 5.36, cloud_coverage: 100, weather_main: 'Clouds', weather_description: 'Bedeckt', 'timestamp': 1748179956, 'lon': 8.175, 'lat': 47.385, 'sys_country': 'C
INFO:_main__Send to Kafka (city: 'Aarau', temperature: 18.88, humidity: 69, pressure: 1016, wind_speed: 3.88, cloud_coverage: 100, weather_main: 'Clouds', weather_description: 'Bedeckt', 'timestamp': 1748179734, 'lon': 8.042, 'lat': 47.3925, 'sys_country': 'C
INFO:_main__Send to Kafka (city: 'Hausen AG', temperature: 19.43, humidity: 74, pressure: 1016, wind_speed: 3.48, cloud_coverage: 100, weather_main: 'Clouds', weather_description: 'Bedeckt', 'timestamp': 1748179958, 'lon': 8.2099, 'lat': 47.464, 'sys_country': 'C
INFO:_main__Send to Kafka (city: 'Zürich', temperature: 19.77, humidity: 69, pressure: 1016, wind_speed: 3.55, cloud_coverage: 100, weather_main: 'Clouds', weather_description: 'Bedeckt', 'timestamp': 1748179959, 'lon': 8.2184, 'lat': 47.479, 'sys_country': 'C
INFO:_main__Send to Kafka (city: 'Zürich', temperature: 18.8, humidity: 66, pressure: 1016, wind_speed: 5.14, cloud_coverage: 75, weather_main: 'Clouds', weather_description: 'Überwiegend bewölkt', 'timestamp': 1748179858, 'lon': 8.55, 'lat': 47.3667, 'sys_coun
INFO:_main__Send to Kafka (city: 'Zürich', temperature: 13.81, humidity: 69, pressure: 1017, wind_speed: 3.89, cloud_coverage: 75, weather_main: 'Rain', weather_description: 'Leichter Regen', 'timestamp': 1748179674, 'lon': 8.4813, 'lat': 46.8211, 'sys_coun
INFO:_main__Send to Kafka (city: 'Sitten', temperature: 20.85, humidity: 40, pressure: 1016, wind_speed: 2.86, cloud_coverage: 75, weather_main: 'Clouds', weather_description: 'Überwiegend bewölkt', 'timestamp': 1748179780, 'lon': 7.3594, 'lat': 46.2291, 'sys
INFO:_main__Send to Kafka (city: 'Sitten', temperature: 18.84, humidity: 69, pressure: 1016, wind_speed: 1.34, cloud_coverage: 100, weather_main: 'Clouds', weather_description: 'Bedeckt', 'timestamp': 1748179963, 'lon': 8.3421, 'lat': 47.3511, 'sys_country': 'C
INFO:_main__Send to Kafka (city: 'Lugano', temperature: 22.04, humidity: 58, pressure: 1015, wind_speed: 3.6, cloud_coverage: 20, weather_main: 'Clouds', weather_description: 'Ein paar Wolken', 'timestamp': 1748179806, 'lon': 8.96, 'lat': 46.8101, 'sys_country': 'C
INFO:_main__Send to Kafka (city: 'Winterthur', temperature: 18.83, humidity: 62, pressure: 1016, wind_speed: 5.14, cloud_coverage: 75, weather_main: 'Clouds', weather_description: 'Überwiegend bewölkt', 'timestamp': 1748179719, 'lon': 8.75, 'lat': 47.5, 'sys_country': 'C
INFO:_main__Send to Kafka (city: 'Biel', temperature: 18.32, humidity: 67, pressure: 1016, wind_speed: 3.6, cloud_coverage: 75, weather_main: 'Clouds', weather_description: 'Überwiegend bewölkt', 'timestamp': 1748179967, 'lon': 7.2441, 'lat': 47.1324, 'sys_cou
```

Abb.4: Producer-Log, Quelle: Philippe Christen

```
(base) philippe-christen:openweathermap-bigdata-project philippe docker run -it --network=openweathermap-bigdata-project_default confluentinc/cp-kafka:7.5.0 kafka-console-consumer --bootstrap-server kafka:9092 --topic weather_data --from-beginning
{"city": "Bregenz", "temperature": 18.89, "humidity": 69, "pressure": 1016, "wind_speed": 3.48, "cloud_coverage": 100, "weather_main": "Clouds", "weather_description": "Bedeckt", "timestamp": 1748179955, "lon": 8.2087, "lat": 47.481, "sys_country": "CH", "event_time": 1748179955, "received_time": 1748179955 }
{"city": "Leipzig", "temperature": 18.74, "humidity": 71, "pressure": 1016, "wind_speed": 5.36, "cloud_coverage": 100, "weather_main": "Clouds", "weather_description": "Bedeckt", "timestamp": 1748179956, "lon": 8.175, "lat": 47.385, "sys_country": "CH", "event_time": 1748179956, "received_time": 1748179956 }
{"city": "Aarau", "temperature": 18.88, "humidity": 69, "pressure": 1016, "wind_speed": 3.88, "cloud_coverage": 100, "weather_main": "Clouds", "weather_description": "Bedeckt", "timestamp": 1748179734, "lon": 8.042, "lat": 47.3925, "sys_country": "CH", "event_time": 1748179734, "received_time": 1748179734 }
{"city": "Hausen AG", "temperature": 19.43, "humidity": 74, "pressure": 1016, "wind_speed": 3.48, "cloud_coverage": 100, "weather_main": "Clouds", "weather_description": "Bedeckt", "timestamp": 1748179958, "lon": 8.2099, "lat": 47.464, "sys_country": "CH", "event_time": 1748179958, "received_time": 1748179958 }
{"city": "Zürich", "temperature": 19.77, "humidity": 69, "pressure": 1016, "wind_speed": 3.55, "cloud_coverage": 100, "weather_main": "Clouds", "weather_description": "Bedeckt", "timestamp": 1748179959, "lon": 8.2184, "lat": 47.479, "sys_country": "CH", "event_time": 1748179959, "received_time": 1748179959 }
{"city": "Zürich", "temperature": 18.8, "humidity": 66, "pressure": 1016, "wind_speed": 5.14, "cloud_coverage": 75, "weather_main": "Clouds", "weather_description": "Überwiegend bewölkt", "timestamp": 1748179858, "lon": 8.55, "lat": 47.3667, "sys_country": "CH", "event_time": 1748179858, "received_time": 1748179858 }
{"city": "Zürich", "temperature": 13.81, "humidity": 69, "pressure": 1017, "wind_speed": 3.89, "cloud_coverage": 75, "weather_main": "Rain", "weather_description": "Leichter Regen", "timestamp": 1748179674, "lon": 8.4813, "lat": 46.8211, "sys_country": "CH", "event_time": 1748179674, "received_time": 1748179674 }
{"city": "Sitten", "temperature": 20.85, "humidity": 40, "pressure": 1016, "wind_speed": 2.86, "cloud_coverage": 75, "weather_main": "Clouds", "weather_description": "Überwiegend bewölkt", "timestamp": 1748179780, "lon": 7.3594, "lat": 46.2291, "sys_country": "CH", "event_time": 1748179780, "received_time": 1748179780 }
{"city": "Sitten", "temperature": 18.84, "humidity": 69, "pressure": 1016, "wind_speed": 1.34, "cloud_coverage": 100, "weather_main": "Clouds", "weather_description": "Bedeckt", "timestamp": 1748179963, "lon": 8.3421, "lat": 47.3511, "sys_country": "CH", "event_time": 1748179963, "received_time": 1748179963 }
{"city": "Lugano", "temperature": 22.04, "humidity": 58, "pressure": 1015, "wind_speed": 3.6, "cloud_coverage": 20, "weather_main": "Clouds", "weather_description": "Ein paar Wolken", "timestamp": 1748179806, "lon": 8.96, "lat": 46.8101, "sys_country": "CH", "event_time": 1748179806, "received_time": 1748179806 }
{"city": "Winterthur", "temperature": 18.83, "humidity": 62, "pressure": 1016, "wind_speed": 5.14, "cloud_coverage": 75, "weather_main": "Clouds", "weather_description": "Überwiegend bewölkt", "timestamp": 1748179719, "lon": 8.75, "lat": 47.5, "sys_country": "CH", "event_time": 1748179719, "received_time": 1748179719 }
{"city": "Biel", "temperature": 18.32, "humidity": 67, "pressure": 1016, "wind_speed": 3.6, "cloud_coverage": 75, "weather_main": "Clouds", "weather_description": "Überwiegend bewölkt", "timestamp": 1748179967, "lon": 7.2441, "lat": 47.1324, "sys_country": "CH", "event_time": 1748179967, "received_time": 1748179967 }
```

Abb.5: Kafka Live-Stream, Quelle: Philippe Christen

```
INFO:kafka.consumer.subscription_state:Updating subscribed topics to: ('weather_data',)
INFO:_main__Start Batch-Verarbeitung in über-Schritten...
INFO:kafka.consumer.subscription_state:Updated partition assignment: [TopicPartition(topic='weather_data', partition=0)]
INFO:kafka.com-drokersconnection client_id=kafka-python-2.2.10, node_id=1 host=kafka:9092 connected [IPv4 ('172.21.0.2', 9092)]: connecting to kafka:9092 [IPv4 ('172.21.0.2', 9092)]: Connection complete.
INFO:kafka.com-drokersconnection client_id=kafka-python-2.2.10, node_id=bootstrap-kafka:9092 connected [IPv4 ('172.21.0.2', 9092)]: closing connection.
INFO:kafka.consumer.fetcher:Resetting offset for partition TopicPartition(topic='weather_data', partition=0) to offset 0.
INFO:_main__Empfang von Kafka (city: 'Bregenz', temperature: 18.89, humidity: 69, pressure: 1016, wind_speed: 3.48, cloud_coverage: 100, weather_main: 'Clouds', weather_description: 'Bedeckt', 'timestamp': 1748179955, 'lon': 8.2087, 'lat': 47.481, 'sys_country': 'CH', 'event_time': 1748179955, 'received_time': 1748179955 }
INFO:_main__Empfang von Kafka (city: 'Leipzig', temperature: 18.74, humidity: 71, pressure: 1016, wind_speed: 5.36, cloud_coverage: 100, weather_main: 'Clouds', weather_description: 'Bedeckt', 'timestamp': 1748179956, 'lon': 8.175, 'lat': 47.385, 'sys_country': 'CH', 'event_time': 1748179956, 'received_time': 1748179956 }
INFO:_main__Empfang von Kafka (city: 'Aarau', temperature: 18.88, humidity: 69, pressure: 1016, wind_speed: 3.88, cloud_coverage: 100, weather_main: 'Clouds', weather_description: 'Bedeckt', 'timestamp': 1748179734, 'lon': 8.042, 'lat': 47.3925, 'sys_country': 'CH', 'event_time': 1748179734, 'received_time': 1748179734 }
INFO:_main__Empfang von Kafka (city: 'Hausen AG', temperature: 19.43, humidity: 74, pressure: 1016, wind_speed: 3.48, cloud_coverage: 100, weather_main: 'Clouds', weather_description: 'Bedeckt', 'timestamp': 1748179958, 'lon': 8.2099, 'lat': 47.464, 'sys_country': 'CH', 'event_time': 1748179958, 'received_time': 1748179958 }
INFO:_main__Empfang von Kafka (city: 'Zürich', temperature: 19.77, humidity: 69, pressure: 1016, wind_speed: 3.55, cloud_coverage: 100, weather_main: 'Clouds', weather_description: 'Bedeckt', 'timestamp': 1748179959, 'lon': 8.2184, 'lat': 47.479, 'sys_country': 'CH', 'event_time': 1748179959, 'received_time': 1748179959 }
INFO:_main__Empfang von Kafka (city: 'Zürich', temperature: 18.8, humidity: 66, pressure: 1016, wind_speed: 5.14, cloud_coverage: 75, weather_main: 'Clouds', weather_description: 'Überwiegend bewölkt', 'timestamp': 1748179858, 'lon': 8.55, 'lat': 47.3667, 'sys_country': 'CH', 'event_time': 1748179858, 'received_time': 1748179858 }
INFO:_main__Empfang von Kafka (city: 'Zürich', temperature: 13.81, humidity: 69, pressure: 1017, wind_speed: 3.89, cloud_coverage: 75, weather_main: 'Rain', weather_description: 'Leichter Regen', 'timestamp': 1748179674, 'lon': 8.4813, 'lat': 46.8211, 'sys_country': 'CH', 'event_time': 1748179674, 'received_time': 1748179674 }
INFO:_main__Empfang von Kafka (city: 'Sitten', temperature: 20.85, humidity: 40, pressure: 1016, wind_speed: 2.86, cloud_coverage: 75, weather_main: 'Clouds', weather_description: 'Überwiegend bewölkt', 'timestamp': 1748179780, 'lon': 7.3594, 'lat': 46.2291, 'sys_country': 'CH', 'event_time': 1748179780, 'received_time': 1748179780 }
INFO:_main__Empfang von Kafka (city: 'Sitten', temperature: 18.84, humidity: 69, pressure: 1016, wind_speed: 1.34, cloud_coverage: 100, weather_main: 'Clouds', weather_description: 'Bedeckt', 'timestamp': 1748179963, 'lon': 8.3421, 'lat': 47.3511, 'sys_country': 'CH', 'event_time': 1748179963, 'received_time': 1748179963 }
INFO:_main__Empfang von Kafka (city: 'Lugano', temperature: 22.04, humidity: 58, pressure: 1015, wind_speed: 3.6, cloud_coverage: 20, weather_main: 'Clouds', weather_description: 'Ein paar Wolken', 'timestamp': 1748179806, 'lon': 8.96, 'lat': 46.8101, 'sys_country': 'CH', 'event_time': 1748179806, 'received_time': 1748179806 }
INFO:_main__Empfang von Kafka (city: 'Winterthur', temperature: 18.83, humidity: 62, pressure: 1016, wind_speed: 5.14, cloud_coverage: 75, weather_main: 'Clouds', weather_description: 'Überwiegend bewölkt', 'timestamp': 1748179719, 'lon': 8.75, 'lat': 47.5, 'sys_country': 'CH', 'event_time': 1748179719, 'received_time': 1748179719 }
INFO:_main__Empfang von Kafka (city: 'Biel', temperature: 18.32, humidity: 67, pressure: 1016, wind_speed: 3.6, cloud_coverage: 75, weather_main: 'Clouds', weather_description: 'Überwiegend bewölkt', 'timestamp': 1748179967, 'lon': 7.2441, 'lat': 47.1324, 'sys_country': 'CH', 'event_time': 1748179967, 'received_time': 1748179967 }
INFO:_main__Verarbeiten Batch mit 8 Nachrichten via Beam.
INFO:apache_beam.runners.runner.statecache:Creating state cache with size 1000000000
INFO:_main__Empfang von Kafka (city: 'Zürich', temperature: 18.8, humidity: 66, pressure: 1016, wind_speed: 5.14, cloud_coverage: 75, weather_main: 'Clouds', weather_description: 'Überwiegend bewölkt', 'timestamp': 1748179858, 'lon': 8.55, 'lat': 47.3667, 'sys_country': 'CH', 'event_time': 1748179858, 'received_time': 1748179858 }
INFO:_main__Empfang von Kafka (city: 'Zürich', temperature: 13.81, humidity: 69, pressure: 1017, wind_speed: 3.89, cloud_coverage: 75, weather_main: 'Rain', weather_description: 'Leichter Regen', 'timestamp': 1748179674, 'lon': 8.4813, 'lat': 46.8211, 'sys_country': 'CH', 'event_time': 1748179674, 'received_time': 1748179674 }
INFO:_main__Empfang von Kafka (city: 'Sitten', temperature: 20.85, humidity: 40, pressure: 1016, wind_speed: 2.86, cloud_coverage: 75, weather_main: 'Clouds', weather_description: 'Überwiegend bewölkt', 'timestamp': 1748179780, 'lon': 7.3594, 'lat': 46.2291, 'sys_country': 'CH', 'event_time': 1748179780, 'received_time': 1748179780 }
INFO:_main__Empfang von Kafka (city: 'Sitten', temperature: 18.84, humidity: 69, pressure: 1016, wind_speed: 1.34, cloud_coverage: 100, weather_main: 'Clouds', weather_description: 'Bedeckt', 'timestamp': 1748179963, 'lon': 8.3421, 'lat': 47.3511, 'sys_country': 'CH', 'event_time': 1748179963, 'received_time': 1748179963 }
INFO:_main__Empfang von Kafka (city: 'Lugano', temperature: 22.04, humidity: 58, pressure: 1015, wind_speed: 3.6, cloud_coverage: 20, weather_main: 'Clouds', weather_description: 'Ein paar Wolken', 'timestamp': 1748179806, 'lon': 8.96, 'lat': 46.8101, 'sys_country': 'CH', 'event_time': 1748179806, 'received_time': 1748179806 }
INFO:_main__Empfang von Kafka (city: 'Winterthur', temperature: 18.83, humidity: 62, pressure: 1016, wind_speed: 5.14, cloud_coverage: 75, weather_main: 'Clouds', weather_description: 'Überwiegend bewölkt', 'timestamp': 1748179719, 'lon': 8.75, 'lat': 47.5, 'sys_country': 'CH', 'event_time': 1748179719, 'received_time': 1748179719 }
INFO:_main__Empfang von Kafka (city: 'Biel', temperature: 18.32, humidity: 67, pressure: 1016, wind_speed: 3.6, cloud_coverage: 75, weather_main: 'Clouds', weather_description: 'Überwiegend bewölkt', 'timestamp': 1748179967, 'lon': 7.2441, 'lat': 47.1324, 'sys_country': 'CH', 'event_time': 1748179967, 'received_time': 1748179967 }
```

Abb.6: Beam-Log, Quelle: Philippe Christen

city	temperature	humidity	pressure	wind_speed	cloud_coverage	weather_main	weather_description	timestamp	lon	lat	sys_country	event_time	received_time	processing_time	anomaly
Yverdon-les-Bains	14.26	68	1017	14.7	93	Clouds	Bedeckt	1748180001	6.6412	46.7785	CH	1748180001	1748180001	1748180002	f
Aigle VD	22.1	63	1016	12.9	98	Clouds	Bedeckt	1748180000	6.9646	46.3181	CH	1748180000	1748180000	1748180002	f
Sarnen	16.67	66	1017	11.1	75	Clouds	Überwiegend bewölkt	1748179999	8.2587	46.8985	CH	1748179999	1748179999	1748180002	f
Herisau	15.21	76	1017	7.7	100	Clouds	Bedeckt	1748179998	9.2702	47.3862	CH	1748179998	1748179998	1748180002	f
Altstätten	15.84	76	1017	3.6	100	Clouds	Bedeckt	1748179997	8.6444	46.8884	CH	1748179997	1748179997	1748180002	f
Appenzell	14.14	76	1017	5.3	100	Clouds	Bedeckt	1748179994	9.41	47.1331	CH	1748179994	1748179994	1748179996	f
Delémont	15.4	65	1016	14	100	Rain	Mässiger Regen	1748179993	7.6217	46.7912	CH	1748179993	1748179993	1748179996	f
Wil SG	17.23	75	1016	14.5	100	Clouds	Bedeckt	1748179988	9.0455	47.4615	CH	1748179988	1748179988	1748179991	f
Burgdorf BE	16.64	76	1016	13	75	Rain	Leichter Regen	1748179984	7.6279	47.059	CH	1748179984	1748179984	1748179986	f
Thun	16.24	62	1016	9.6	99	Rain	Leichter Regen	1748179983	7.2445	47.3649	CH	1748179983	1748179983	1748179986	f
Locarno	23.28	61	1014	6.4	77	Clouds	Überwiegend bewölkt	1748179980	8.7995	46.1709	CH	1748179980	1748179980	1748179986	f
Samedan	11.89	58	1017	24.1	75	Clouds	Überwiegend bewölkt	1748179979	9.8712	46.5342	CH	1748179979	1748179979	1748179986	t
Olten	16.24	76	1017	7	96	Clouds	Bedeckt	1748179978	9.5329	46.8499	CH	1748179978	1748179978	1748179986	f
St. Gallen	15.74	63	1017	11.1	75	Clouds	Überwiegend bewölkt	1748179977	9.3748	47.4239	CH	1748179977	1748179977	1748179986	f
Kanton Neuenburg	21.41	46	1017	16.7	75	Clouds	Überwiegend bewölkt	1748179975	6.8333	46.9167	CH	1748179975	1748179975	1748179976	f
Freiburg	16.94	56	1017	16.7	75	Clouds	Überwiegend bewölkt	1748179974	7.1513	46.8024	CH	1748179974	1748179974	1748179976	f
Biel	19.32	67	1016	13	75	Clouds	Überwiegend bewölkt	1748179967	7.2441	47.1324	CH	1748179967	1748179966	1748179970	f
Bern	16.99	70	1016	14.5	75	Clouds	Überwiegend bewölkt	1748179966	7.4474	46.9481	CH	1748179966	1748179966	1748179970	f
Bregarten	16.94	69	1016	4.8	100	Clouds	Bedeckt	1748179963	8.2421	47.3511	CH	1748179963	1748179963	1748179965	f
Basel	21.23	49	1016	22.2	0	Clear	Klarer Himmel	1748179959	7.5733	47.5584	CH	1748179959	1748179959	1748179973	f
Windisch	15.77	69	1016	12.8	100	Clouds	Bedeckt	1748179959	8.2184	47.479	CH	1748179959	1748179959	1748179959	f
Hausen AG	19.43	74	1016	12.5	100	Clouds	Bedeckt	1748179958	8.2899	47.464	CH	1748179958	1748179958	1748179959	f
Lenzburg	16.74	71	1016	19.3	100	Clouds	Bedeckt	1748179956	8.175	47.3885	CH	1748179956	1748179956	1748179959	f
Brugg	15.89	69	1016	12.5	100	Clouds	Bedeckt	1748179955	8.2087	47.481	CH	1748179955	1748179955	1748179959	f
Bulle FR	16.94	58	1017	17.7	92	Rain	Leichter Regen	1748179859	7.8567	46.6195	CH	1748179859	1748179991	1748179991	f
Zürich	18.8	66	1016	18.5	75	Clouds	Überwiegend bewölkt	1748179850	8.55	47.3667	CH	1748179850	1748179960	1748179965	f
Olten	16.87	66	1016	11.6	100	Clouds	Bedeckt	1748179847	7.9833	47.35	CH	1748179847	1748179985	1748179986	f

Abb.7: *psql-Abfrage, Quelle: Philippe Christen*

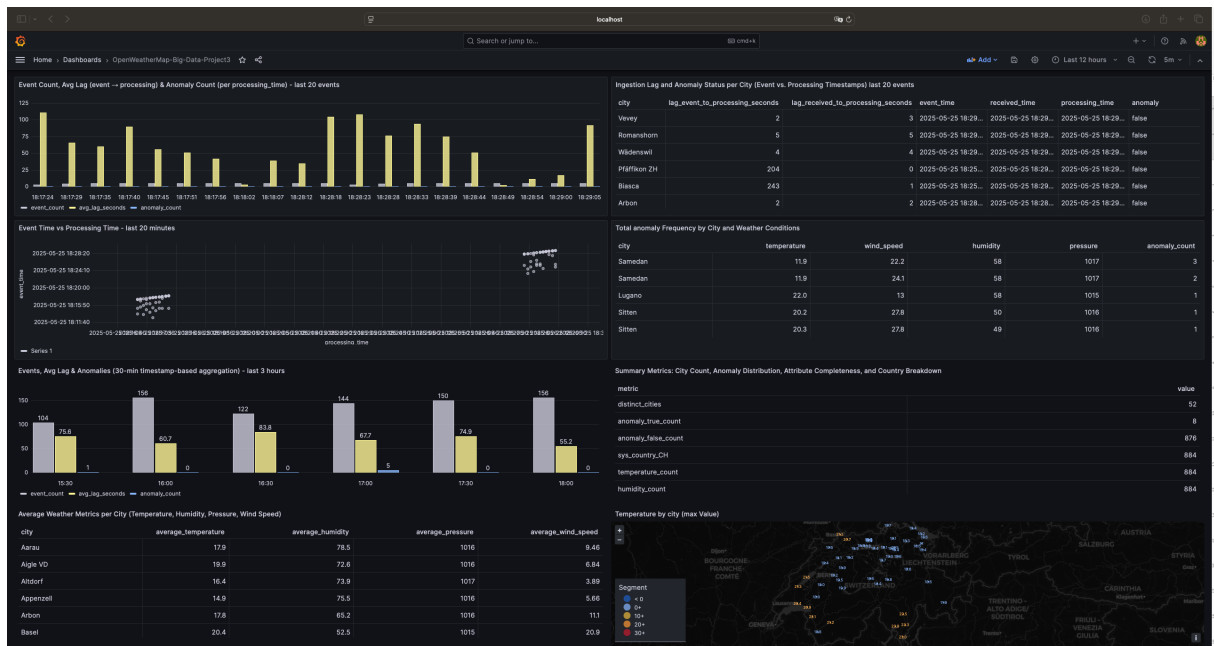


Abb.8: *Grafana Dashboard, Quelle: Philippe Christen*

Erweiterung: Anomalieerkennung mit Machine Learning

Zur qualitativen Erweiterung wurde ein auf IsolationForest basierendes ML-Modell integriert, das meteorologische Anomalien (z. B. unplausible Kombinationen von Wind, Temperatur, Luftdruck) erkennt. Das Modell wurde offline trainiert und via joblib serialisiert. In der Beam-Pipeline wird es mit einem standardisierten Feature-Vektor (6 Merkmale: Temperatur, Wind, Luftdruck, Luftfeuchtigkeit, Längengrad, Breitengrad) verwendet. Die Vorverarbeitung erfolgt mittels StandardScaler. Die Features werden in einem pandas.DataFrame organisiert, um eine skalierbare Übergabe an den StandardScaler zu ermöglichen. Der boolesche Prädiktor «anomaly» wird in der Datenbank gespeichert und zur Analyse in Grafana genutzt. Fehlerquellen – wie inkonsistente Typen oder fehlende Spaltennamen – wurden durch explizite Konvertierung der Anomalievariable zu bool und die Übergabe benannter Spalten im DataFrame behoben.

Fazit und Ausblick

Die entwickelte Pipeline zeigt, wie sich mit Apache Beam und Kafka ein modular aufgebautes, robustes Streaming-System realisieren lässt. Durch die Anomalieerkennung wurde die analytische Tiefe der Daten erheblich erhöht. Die Visualisierung in Grafana erlaubt Echtzeit-Monitoring und Auswertung historischer Latenzen (Lag).

Zukünftig denkbare Erweiterungen

- Einführung von Windowing und Triggering in Beam für dynamische Zeitfenster
- Verwendung von Apache Flink als Runner für produktive Latenzanforderungen
- Modell-Drift-Überwachung und Nachtraining des ML-Modells
- Echtzeit-Alerting bei Anomalien via Webhook/Slack/E-Mail
- Die Architektur ist flexibel, lokal lauffähig, nachvollziehbar und auf zukünftige Anforderungen erweiterbar.
- Deployment in Cloud-Umgebungen wie Google Cloud Dataflow oder AWS MSK zur weiteren Skalierung
- Orchestrierung mit Kubernetes, um Hochverfügbarkeit, Auto-Scaling und Ausfallsicherheit zu gewährleisten
- Data Governance & Monitoring: Implementierung eines zentralen Monitorings (z. B. mit Prometheus und Alertmanager), um Logs automatisiert auszuwerten und bei API-Ausfällen oder Fehlern im Streaming-Prozess frühzeitig Benachrichtigungen auszulösen. Dies erhöht die Betriebssicherheit und unterstützt Data Governance durch transparente Überwachung der Pipeline