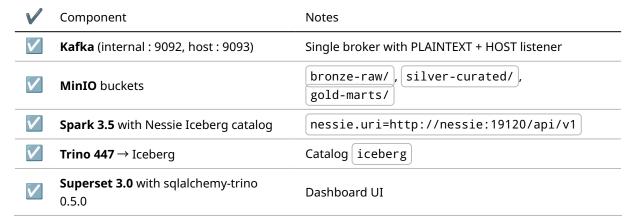
Social-Media Analytics on the Medallion Platform

A high-level delivery plan that turns your working infra (Kafka \rightarrow Spark \rightarrow Iceberg \rightarrow Trino \rightarrow Superset) into a production-grade pipeline for X / Instagram / TikTok data.

0 Foundations (already done 🔽



No more infra stories—everything below is data & application.

1 Product Vision

"Marketing wants minute-level engagement KPIs and hashtag sentiment trend lines for all active campaigns."

- Freshness SLA: < 2 min end-to-dashboard.
- History: Keep raw forever; Silver + Gold for 400 days.
- Scalability target: 10 M social events / hr (\approx 2 K msg/s).

gh repo clone uzair14137/Medallion-Data-Platform-on-MinIO2 Roadmap & Sprints (2-week cadence)

Enrichment & ClickHouse metrics | • User lookup enrichment (followers, verified)• ClickHouse MV likes_last_5m | SELECT returns correct counts | | 4 | Gold marts + Superset dashboards | • Iceberg view gold.social.campaign_perf • Superset dashboard "Real-Time Engagement" with 6 charts | Stakeholder sign-off | | 5 | Hardening & observability | • Prom/ Grafana alerts: lag, DAG failures• Retry/ dead-letter topic | All dashboards stay green for 72 h |

3 Detailed Backlog (next two sprints)

Sprint 1 – Ingestion MVP

- 1. **Story S1-1** Twitter API connector (search/recent by hashtag)
- 2. Use tweepy v4, polling every 15 s.
- 3. **\$1-2** Dockerise collector
- 4. Build collector-social:0.1 image; env vars for API keys, Kafka brokers.
- 5. **\$1-3** Kafka topic setup script
- 6. kafka-topics --create ... --configs
 cleanup.policy=compact,retention.ms=604800000
- 7. **\$1-4** Schema Registry setup (optional)
- 8. **S1-5** Airflow DAG ingest_raw.py
- 9. KubernetesPodOperator or DockerOperator.

Sprint 2 - Bronze → Silver

- 1. **S2-1** Spark job skeleton
- 2. Read from kafka , write to nessie.social.posts_raw . Demo run on cluster mode.
- 3. **S2-2** Parse + basic transformations
- 4. created_at to UTC, extract hashtags array, language.
- 5. **S2-3** Unit tests with spark-testing-base.
- 6. **S2-4** Airflow DAG bronze_to_silver.py (Streaming Add Jar operator).

4 Data Model (Silver → Gold)

Layer	Table / View	Partitions	Notes
Silver	social.posts	ds_hour	Clean, enriched posts
Silver	social.hashtags	ds_hour, hashtag	Exploded array
Gold	social.campaign_perf	daily	campaign_id, metrics

Layer	Table / View	Partitions	Notes	
Gold	social.sentiment_hourly	hourly	hashtag	integrate HuggingFace sentiment

5 Orchestration Graph (Airflow)

```
ingest_raw → bronze_to_silver →
enricolickhouse
gold_marts → superset_refresh
```

Use Airflow-Nessie plugin (optional) to tag/merge after successful DAG run.

6 Monitoring & Alerting

- Kafka lag → JMX exporter + Prometheus rule lag_seconds > 30 for 5m.
- **Spark streaming** → StructuredStreaming metrics sink.
- Airflow → failure Slack alert via Webhook.
- MinIO capacity → MinIO Prometheus exporter, 80 % threshold.

7 Security Notes

- Store API keys in Docker secrets and Airflow Connections .
- Consider SASL/SSL for Kafka (listeners | SASL_HOST | on 9094) before prod.
- Enable Trino password auth if Superset is exposed outside dev PC.

Appendix A Bronze-layer Ingestion Guide

Goal: move every JSON event that lands in Kafka topic ` into the immutable bucket s3a://bronze-raw/social/`within seconds.

1 Create / verify the topic

```
# idempotent - safe to run at deploy time or Airflow task
kafka-topics \
   --bootstrap-server kafka:9092 \
   --create \
```

```
--topic social_raw \
--partitions 3 \
--replication-factor 1 \
--config cleanup.policy=delete \
--config retention.ms=604800000 # 7 days for replay
```

2 Producer skeleton (Python + FastAPI)

```
# src/collector/twitter_producer.py
import os, json, asyncio, tweepy, aiokafka
BOOTSTRAP = os.getenv("KAFKA_BROKERS", "kafka:9092")
           = os.getenv("TOPIC", "social_raw")
async def main():
   producer = aiokafka.AIOKafkaProducer(bootstrap_servers=BOOTSTRAP)
   await producer.start()
   try:
        for tweet in tweepy.Cursor(api.search tweets, g="#marketing",
tweet mode="extended").items():
            await producer.send and wait(TOPIC,
json.dumps(tweet._json).encode())
    finally:
        await producer.stop()
if __name__ == "__main__":
    asyncio.run(main())
```

Build as image collector-social:0.1 and run via Airflow DockerOperator.

3 Spark Structured-Streaming job (Bronze writer)

Package with zipapp or --py-files, submit via

```
spark-submit --packages org.apache.iceberg:iceberg-spark-runtime-3.5_2.12:1.5.2
\
bronze_to_minio.py
```

4 Airflow DAG snippet

```
from airflow import DAG
from airflow.providers.apache.spark.operators.spark_submit import
SparkSubmitOperator
from airflow.utils.dates import days_ago
dag = DAG(
    "bronze_to_minio", start_date=days_ago(1), schedule_interval="@once",
catchup=False)
write bronze = SparkSubmitOperator(
    task_id="write_bronze", dag=dag,
    application="/opt/airflow/jobs/bronze_to_minio.py",
    name="bronze-social",
    conf={
        "spark.hadoop.fs.s3a.endpoint": "http://minio:9000",
        "spark.hadoop.fs.s3a.access.key": "minio",
        "spark.hadoop.fs.s3a.secret.key": "minio123",
        "spark.hadoop.fs.s3a.path.style.access": "true",
    },
    packages="org.apache.iceberg:iceberg-spark-runtime-3.5_2.12:1.5.2")
```

5 Validation

Check	Command	
Kafka lag	kafka-consumer-groupsbootstrap-server kafka:9092 describegroup <spark-group></spark-group>	
Files in bucket	`mc ls minio/bronze-raw/social/	head`
Query raw in Spark	Spark.read.1Son("S3a://pronze-raw/Social/").count()	

(this week)

- 1. Approve sprint scope
- 2. Generate Twitter dev keys (marketing)
- 3. Merge Compose changes (HOST listener, spark-defaults) to main
- 4. Start **Sprint 1-Day 1** stand-up Monday 09:00