금융 데이터 수집을 통한 데이터 분석환경 구축



김수민 김형준 이수영 조주혜 한기호

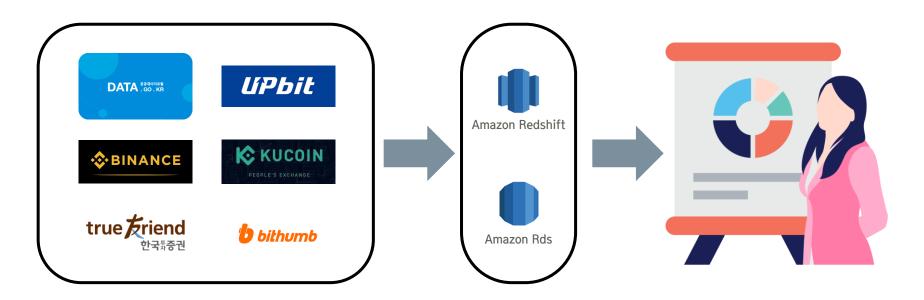
CONTENTS

01 주제 선정배경/목표

02 사용기술

- 03 구축과정
 - 실시간처리
 - 배치처리
- 04 최종 시각화

○1 프로젝트 목표



금융 데이터 통합

편리한 데이터 분석 가능

02 사용기술

스케쥴러



데이터 수집 도구





실시간 데이터 처리





Amazon EC2

데이터 전처리



DB



Amazon S3 (DataLake)



Amazon Redshift Amazon Rds (DataWarehouse)



(MySQL)

시각화 도구



Amazon Quicksight

기타도구









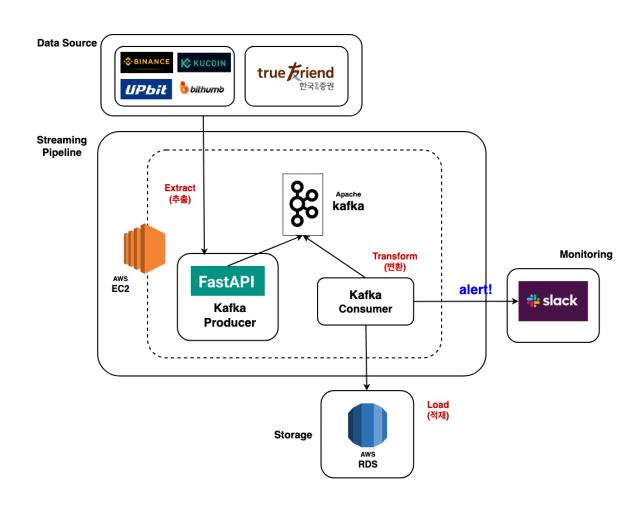






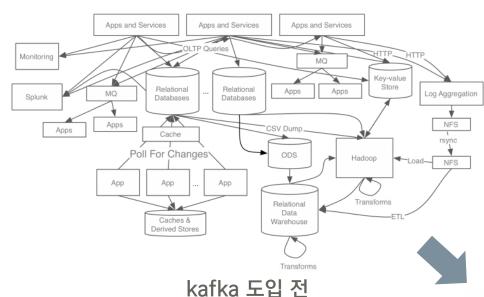
- 실시간처리

실시간처리 아키텍처



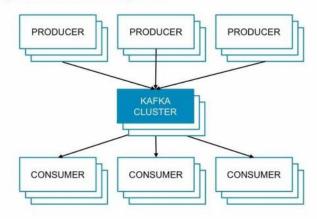
- 실시간처리

kafka 도입 이유



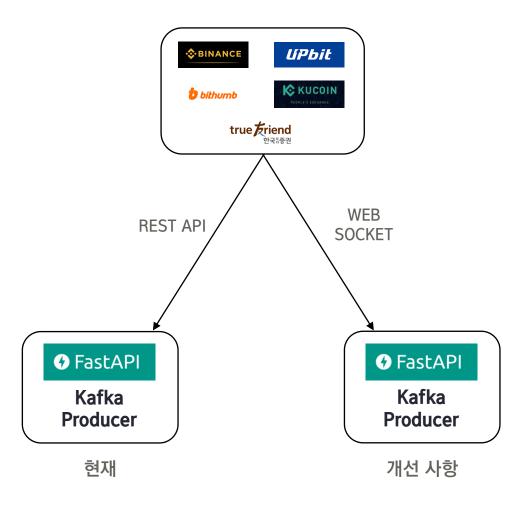
kafka 도입 후

Kafka Inputs and Outputs

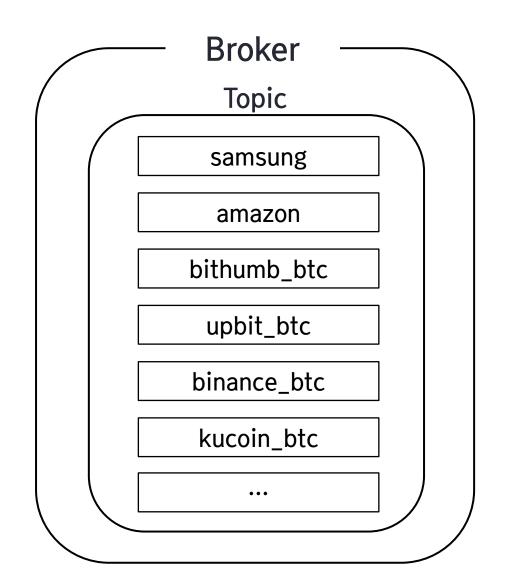


- 실시간처리

Producer

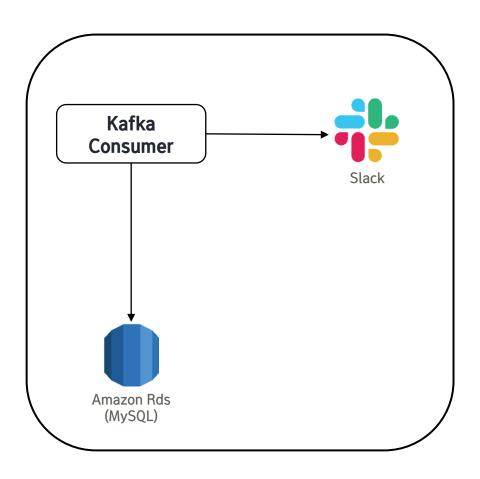


- 실시간처리



- 실시간처리

Consumer

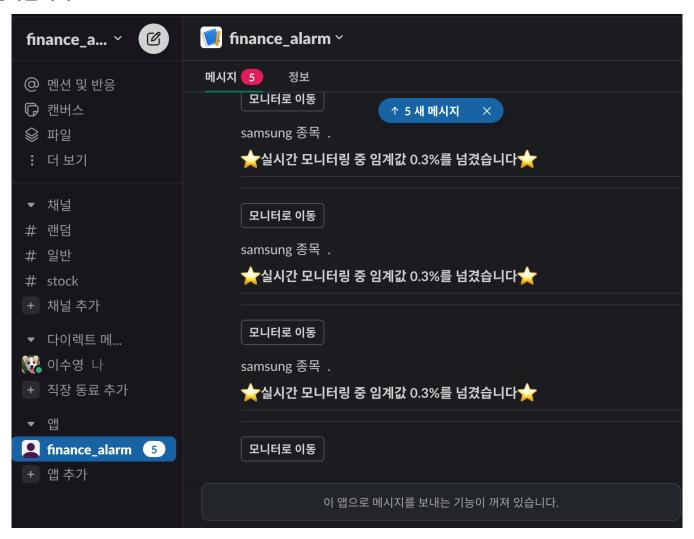


- 실시간처리

	id 🔻	<pre>Ø get_data_timestamp</pre>	RBC ticker 🔻	RBC stock_name T	123 market_cap	123 current_price This is a second of the	123 change
60	60	2023-09-04 21:24:48	005930	삼성전자	4,256,460	71,300	0.14
61	61	2023-09-04 21:24:49	005930	삼성전자	4,250,480	71,200	0
62	62	2023-09-04 21:24:24	005930	삼성전자	4,256,460	71,300	0.14
63	63	2023-09-04 21:24:26	005930	삼성전자	4,250,480	71,200	0
64	64	2023-09-04 21:24:27	005930	삼성전자	4,256,460	71,300	0.14
65	65	2023-09-04 21:24:38	005930	삼성전자	4,250,480	71,200	0
66	66	2023-09-04 21:24:50	005930	삼성전자	4,256,460	71,300	0.14
67	67	2023-09-04 21:24:54	005930	삼성전자	4,250,480	71,200	0
68	68	2023-09-04 21:24:55	005930	삼성전자	4,256,460	71,300	0.14
69	69	2023-09-04 21:24:58	005930	삼성전자	4,250,480	71,200	0
70	70	2023-09-04 21:24:17	005930	삼성전자	4,250,480	71,200	0
71	71	2023-09-04 21:24:21	005930	삼성전자	4,250,480	71,200	0
72	72	2023-09-04 21:24:23	005930	삼성전자	4,250,480	71,200	C
73	73	2023-09-04 21:24:32	005930	삼성전자	4,256,460	71,300	0.14
74	74	2023-09-04 21:24:33	005930	삼성전자	4,256,460	71,300	0.14
75	75	2023-09-04 21:24:36	005930	삼성전자	4,250,480	71,200	0
76	76	2023-09-04 21:24:39	005930	삼성전자	4,250,480	71,200	0
77	77	2023-09-04 21:24:42	005930	삼성전자	4,256,460	71,300	0.14
78	78	2023-09-04 21:24:43	005930	삼성전자	4,256,460	71,300	0.14
79	79	2023-09-04 21:24:47	005930	삼성전자	4,250,480	71,200	0
30	80	2023-09-04 21:24:51	005930	삼성전자	4,250,480	71,200	0
31	81	2023-09-04 21:24:53	005930	삼성전자	4,250,480	71,200	0
32	82	2023-09-04 21:24:22	005930	삼성전자	4,256,460	71,300	0.14
33	83	2023-09-04 21:24:28	005930	삼성전자	4,256,460	71,300	0.14
84	84	2023-09-04 21:24:37	005930	삼성전자	4,256,460	71,300	0.14
35	85	2023-09-04 21:24:40	005930	삼성전자	4,256,460	71,300	0.14
36	86	2023-09-04 21:26:04	005930	삼성전자	4,256,460	71,300	0
37	87	2023-09-04 21:26:00	005930	삼성전자	4,250,480	71,200	-0.14
38	88	2023-09-04 21:26:01	005930	삼성전자	4,250,480	71,200	-0.14
89	89	2023-09-04 21:26:03	005930	삼성전자	4,250,480	71,200	-0.14

삼성전자 테이블

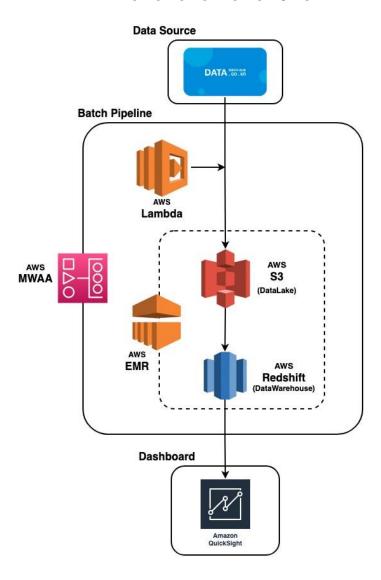
- 실시간처리



Slack에 변동률 알림

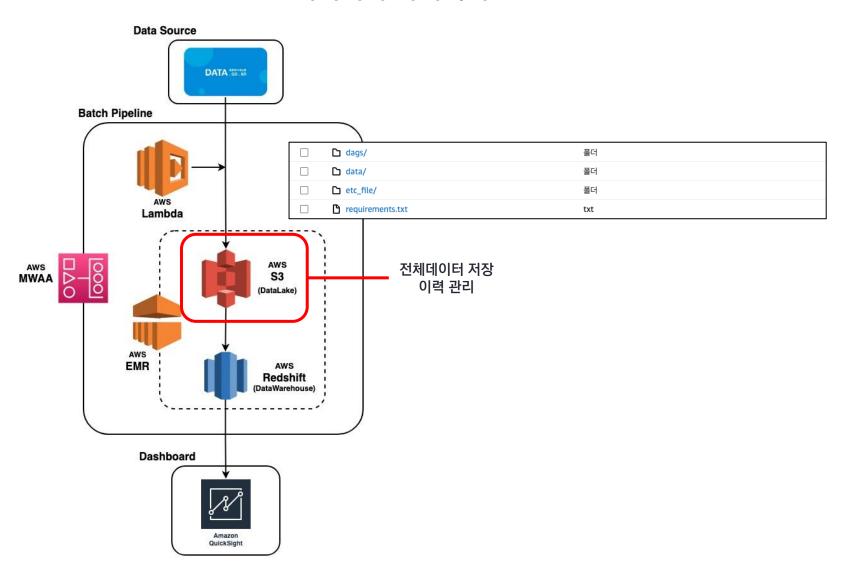
- 배치처리

배치처리 아키텍처



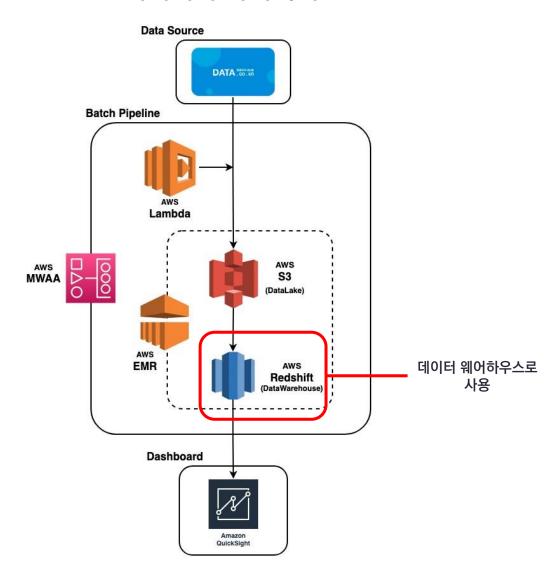
- 배치처리

배치처리 아키텍처



- 배치처리

배치처리 아키텍처



- 배치처리

스케쥴러



장점

- 비용 효율성
- 유연성

단점

- 초기세팅 필요
- 관리 어려움



장점

- 뛰어난 보안성
- 높은 이용률

단점

- 러닝커브가 있을 수 있음



장점

- 편리한 관리 기능

단점

- DAG 수 제한

- 리소스 대비 많은 비용 부과

- 배치처리



장점

- 비용 효율성
- 유연성

단점

- 초기세팅 필요
- 관리 어려움



장점

- 뛰어난 보안성
- 높은 이용률

단점

- 러닝커브가 있을 수 있음

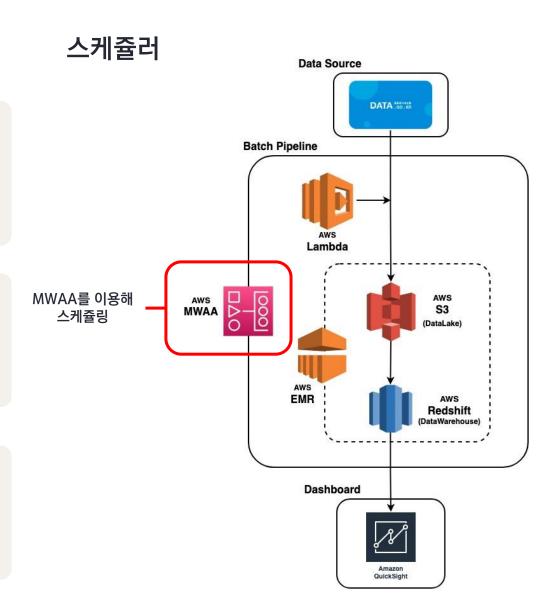


장점

- 편리한 관리 기능

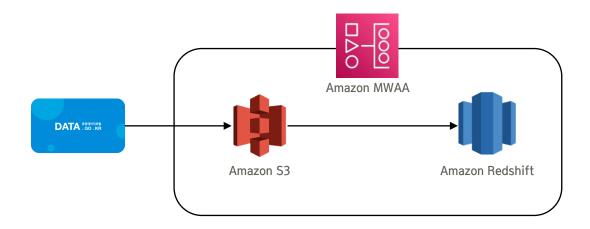
단점

- DAG 수 제한
- 리소스 대비 많은 비용 부과



- 배치처리

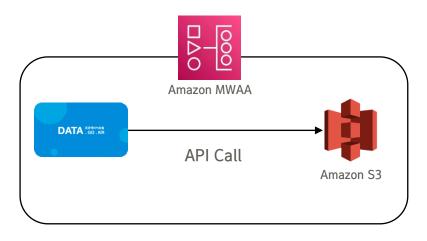
초기 데이터 적재 방법



매일 full refresh

- 배치처리

데이터 적재

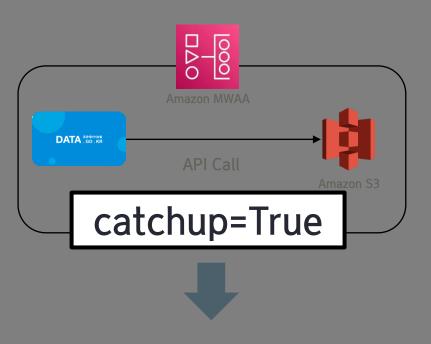




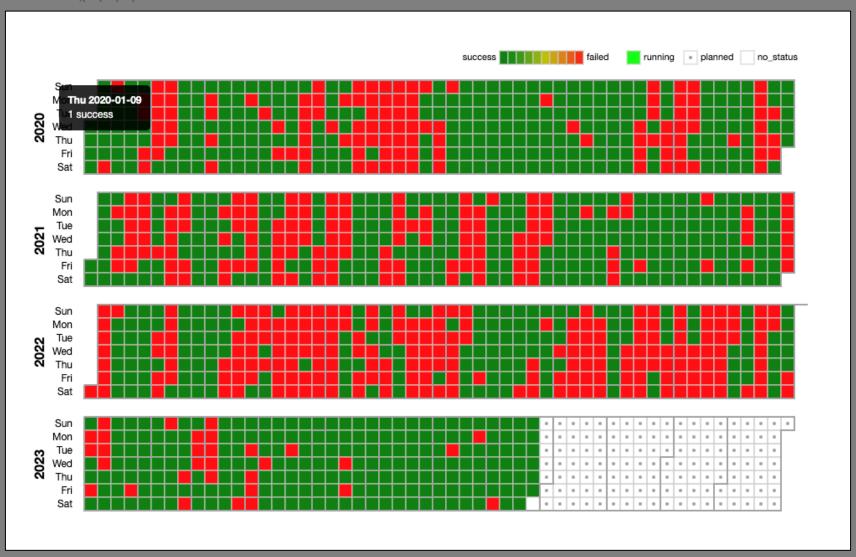
데이터가 많은 경우 문제가 생기게 됨!

- 배치처리

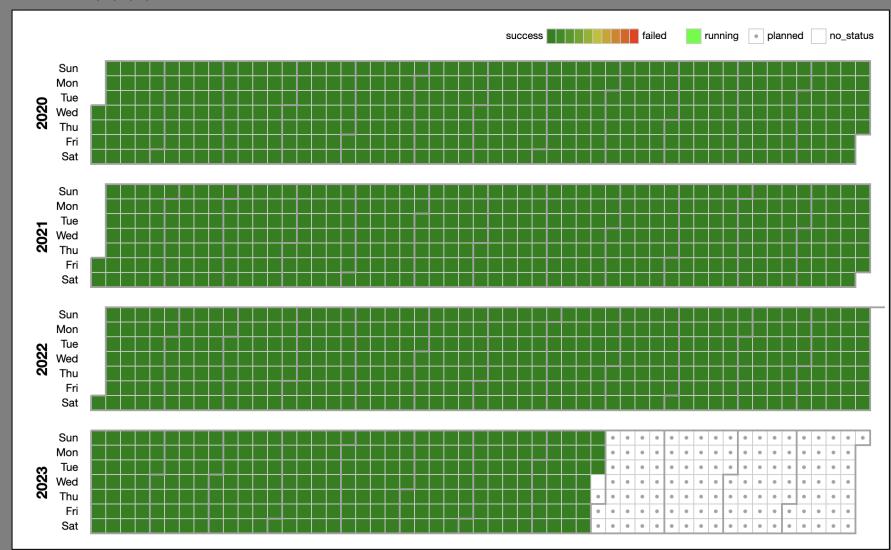
데이터 적재



데이터가 많은 경우 문제가 생기게 됨!

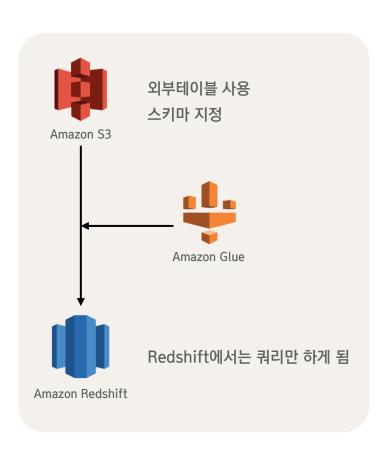






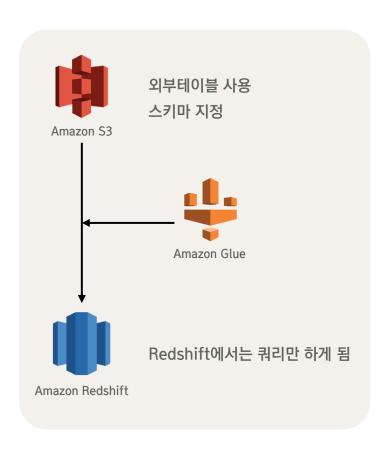
- 배치처리

외부테이블 사용과 파티셔닝



- 배치처리

외부테이블 사용과 파티셔닝



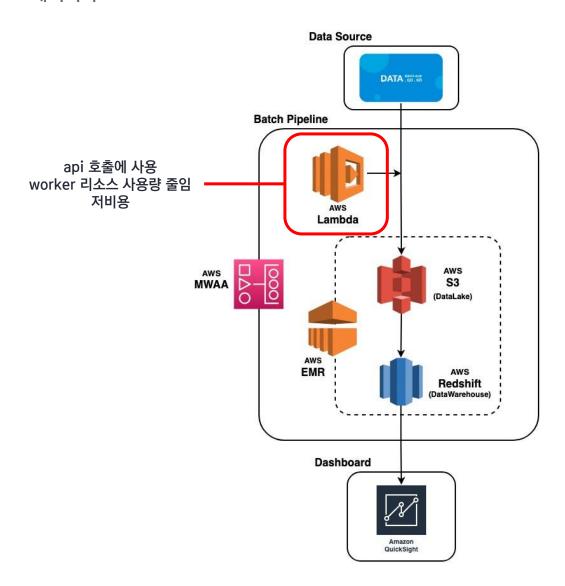
☐ dt=2022-11/
dt=2022-12/
dt=2023-01/
dt=2023-02/
☐ dt=2023-03/
dt=2023-04/
dt=2023-05/
dt=2023-06/
dt=2023-07/
☐ dt=2023-08/

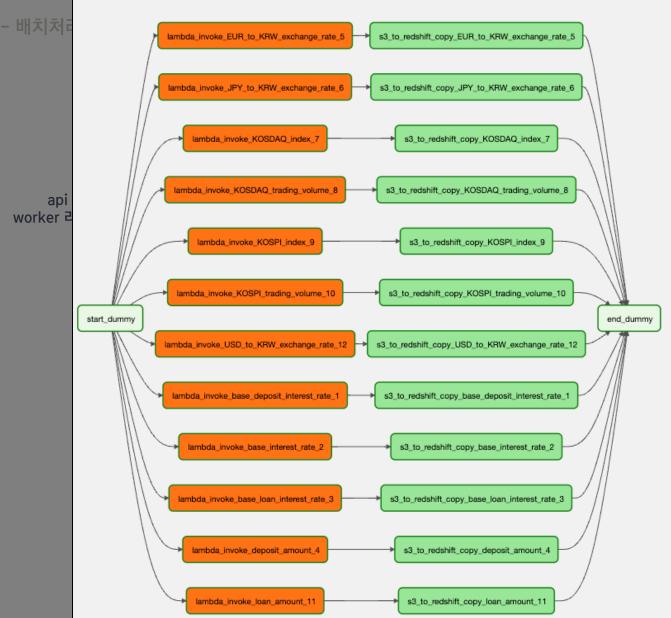
날짜 파티셔닝 → 최적화

- 배치처리

MWAA 리소스 제한

환경 클래스 정보										
각 Amazon MWAA 환경에는 스케줄러, 웹 서버 및 1 작업자가 포함됩니다. 작업자는 시스템 로드에 따라 확장 및 축소됩니다. 환경의 로드 를 모니터링하고 클래스를 언제든지 수정할 수 있습니다.										
ı	DAG 용량*	스케줄러 CPU	작업자 CPU	웹 서버 CPU						
mw1.small	최대 50	1 vCPU	1 vCPU	0.5 vCPU						
O mw1.medium	최대 250	2 vCPU	2 vCPU	1 vCPU						
O mw1.large	최대 1000	4 vCPU	4 vCPU	2 vCPU						
최대 작업자 수 환경에서 확장할 수 있도록 8	허용된 최대 작업자 수입니다.			*일반적인 사용량에서						
10										
1~25여야 합니다. 최소 작업자 수 환경에 항상 존재하는 최소 작	작업자 수입니다.									
1										
최대 작업자보다 작거나 같0	l야 합니다. 최소 작업자 1명									
스케줄러 수 환경에서 사용할 스케줄러 수	-입니다.									
2										
2~5여야 합니다.										





- 배치처리

전처리 과정



특징

- 대규모 데이터 처리에 적합
- Cluster 관리필요
- SparkStreaming 지원 0



특징

- 스크립트 실행에 중점을 둠
- Cluster관리 불필요
- SparkStreaming 지원 X

- 배치처리

전처리 과정



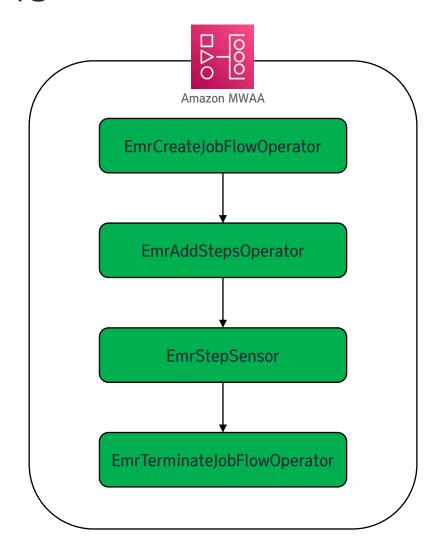
특징

- 대규모 데이터 처리에 적합
- Cluster 관리필요
- SparkStreaming 지원 0



특징

- 스크립트 실행에 중점을 둠
- Cluster관리 불필요
- SparkStreaming 지원 X



- 배치처리

전처리 과정



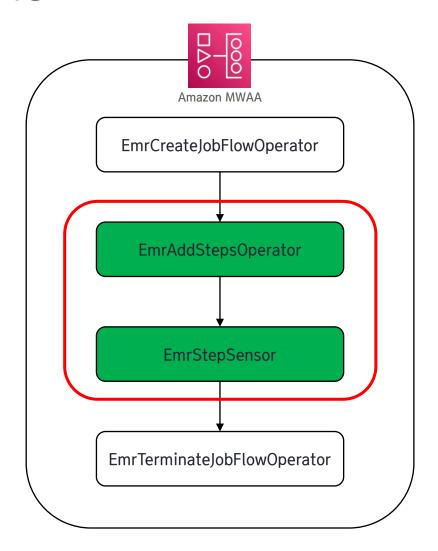
특징

- 대규모 데이터 처리에 적합
- Cluster 관리필요
- SparkStreaming 지원 0



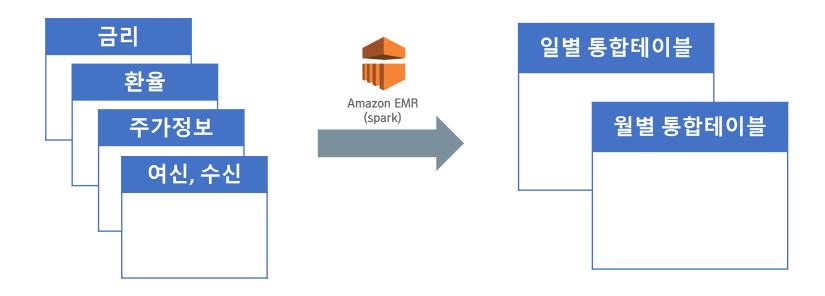
특징

- 스크립트 실행에 중점을 둠
- Cluster관리 불필요
- SparkStreaming 지원 X



- 배치처리

전처리 과정



04 최종 시각화

부동산 데이터 시각화



04 최종 시각화

지표 데이터 시각화



Q & A