

Work Experience

2023.02.13 ~ 현재	SOOP (구 AfreecaTV)	데이터기술팀 / 팀원 ML Engineer
2022.06.20 ~ 2022.12.31	Undefined	개발팀 / 팀원 ML Engineer
2019.12.26 ~ 2020.02.29	Kakao	추천팀 / 인턴 Data Scientist
2018.11.05 ~ 2019.04.22	한빛소프트	인공지능 파트 / 팀원 Research Scientist

Education

2020.03.01 ~ 2022.02.25	한양대학교 대학원 컴퓨터소프트웨어학과	ML System Lab. / 석사 DL Model Optimization
2012.03.01 ~ 2018.08.31	강원대학교 산업공학과	Management of Tech. Lab. / 학부 Gamification

Publications

2022 BIB Journal (Briefings in Bioinformatics)	<i>RAMP: Response-Aware Multi-task Learning with Contrastive Regularization for Cancer Drug Response Prediction</i> ( <a href="#">link</a> )
2022 ICEIC (International Conference on Electronics, Information, and Communication)	<i>Quantization training with two-level bit width</i> ( <a href="#">link</a> )

SOOP (구 AfreecaTV)	[1] Global SOOP Reco. Pipeline	(진행중)
	[2] MLOps	1 개월
	[3] Clip & Short-form VOD Reco.	3 + 3 개월
	[4] Streamer Exploration	2 개월
	[5] Streamer Representative VOD	2 개월
	[6] VOD View Valuation	2 개월
	[7] Offline Simulation	2 개월
	[8] LIVE Broadcast Reco.	2 주

# PROJECTS

Undefined	[1] Match Result Recorder	1 개월
	[2] Competition Rule Recommendation	2 개월
	[3] FAQ Chatbot	3 개월
ML System Lab.	[1] Network Embedding Generation	2 년 2 개월
	[2] DNN Model Quantization - 1	2 년
	[3] DNN Model Quantization - 2	3 개월
	[4] Artificial Intelligence Assistant	2 개월
Kakao	[1] Automobile Video Recommendation	2 개월
	[2] Comics Recommendation	2 주
HanbitSoft	[1] (KR) Multi-speaker Speech Synthesis Model	4 개월
	[2] (EN) Text Chatbot	2 개월

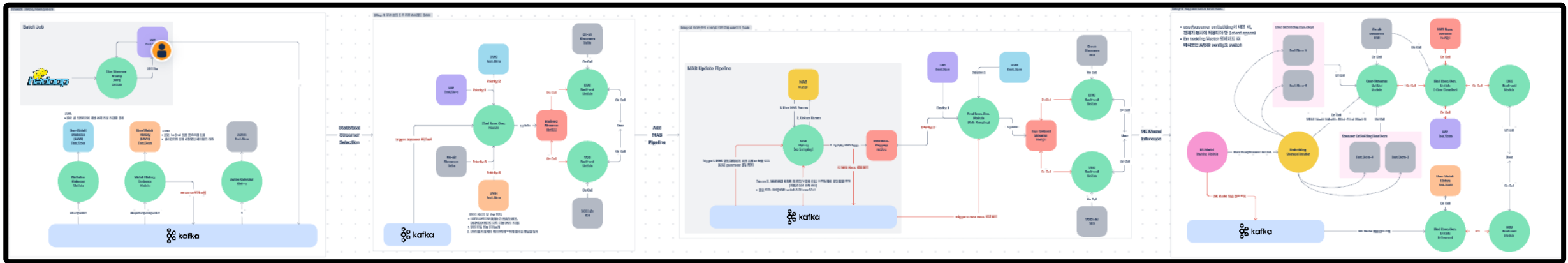
## SOOP (구 AfreecaTV)

- [1] Global SOOP Reco. Pipeline
- [2] MLOps
- [3] Clip & Short-form VOD Reco.
- [4] Streamer Exploration
- [5] Streamer Representative VOD
- [6] VOD View Valuation
- [7] Offline Simulation
- [8] LIVE Broadcast Reco.

## [신규 서비스] 추천 파이프라인 개발

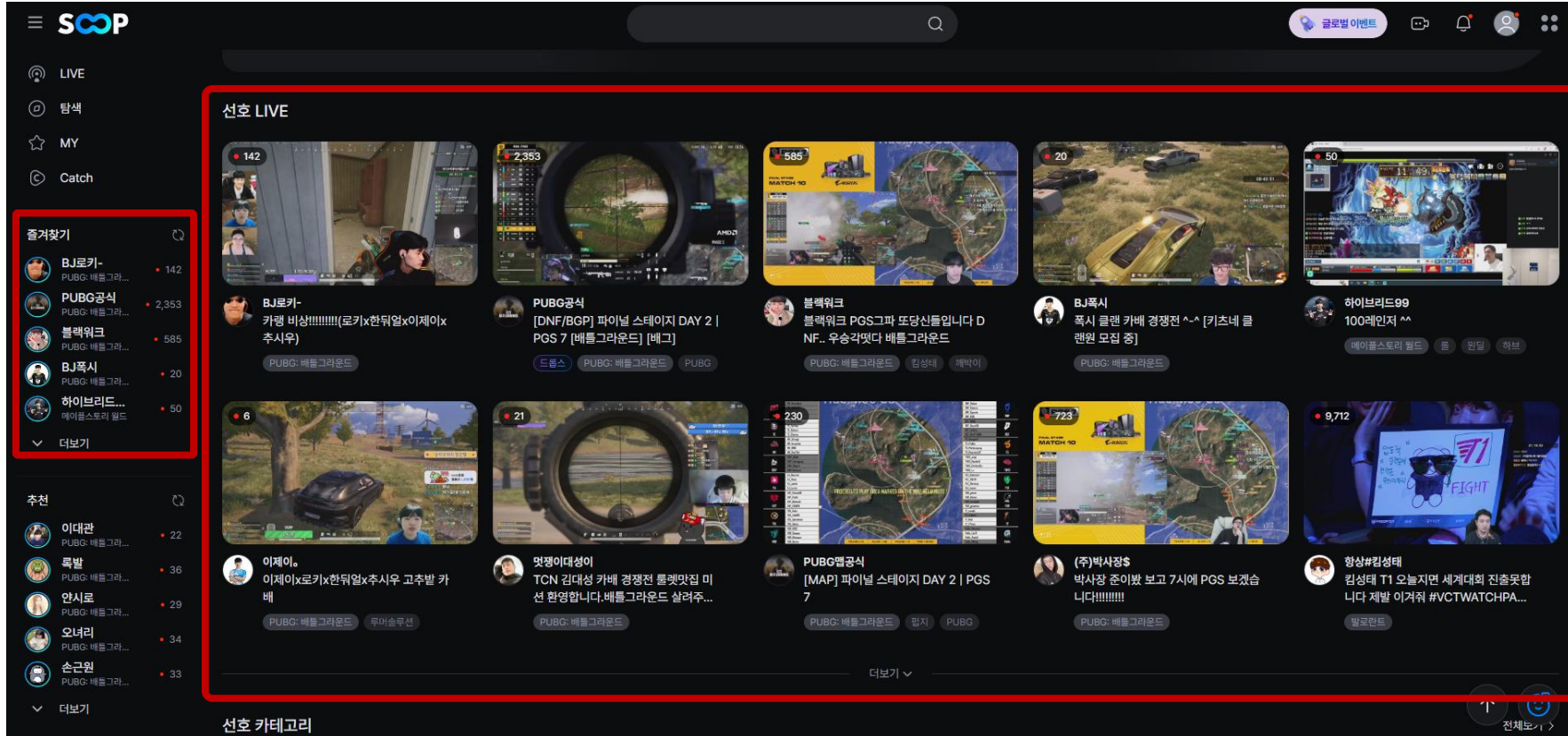
### 추천 파이프라인 개발 4 단계

1. 데이터 수집 및 전처리
2. 장기 선호 스트리머 추천
3. 최근 관심 스트리머 추천 (MAB)
4. 연관 스트리머 추천 (DL)



<p>완료</p>	<ul style="list-style-type: none"> <li>- 신규 서비스 추천 파이프라인 구상 및 계획 수립</li> <li>- '장기 선호 스트리머' 국내 서비스의 메인 페이지, Live 추천 배포</li> </ul>
<p>예정</p>	<ol style="list-style-type: none"> <li>1. MAB를 이용한 '최근 관심 스트리머' 추천</li> <li>2. 유저 시청 데이터 수집 및 Feature Store 구축 (Batch &amp; Stream)</li> <li>3. DL 모델을 이용하여 유저와 스트리머('연관 스트리머')/VOD의 연산 및 추천</li> </ol>
<p>기술</p>	<p>AWS, k8s, Kafka, Airflow, Feast, Hive, Spark, MAB, PyTorch, ElasticSearch, Redis</p>

## [메인 페이지 > Live 추천] 장기 선호 스트리머 데이터



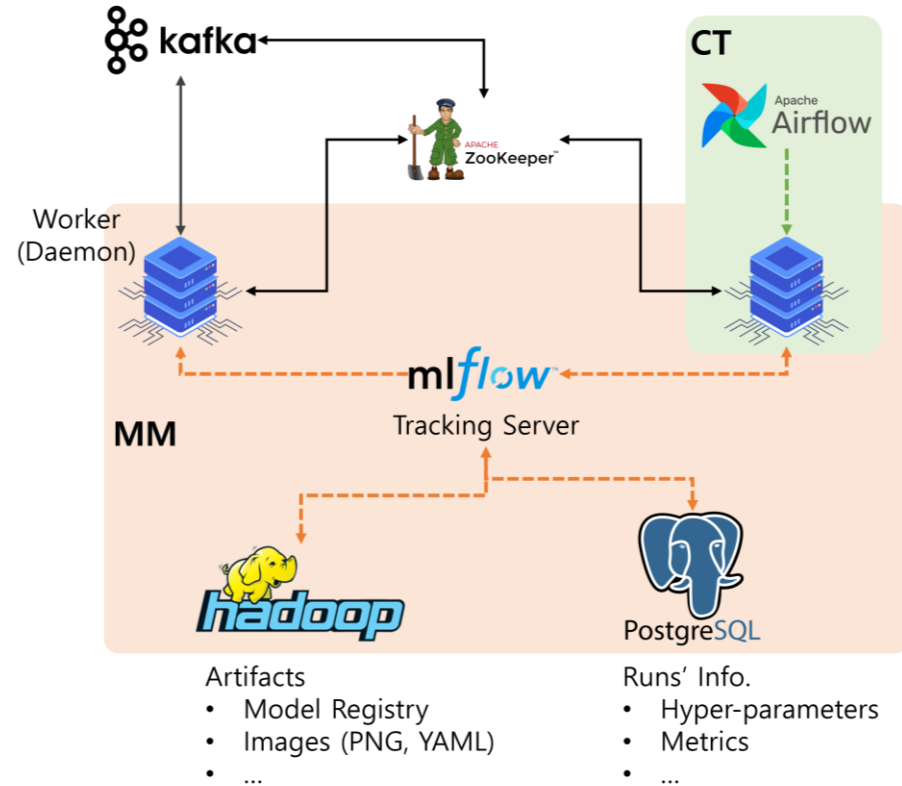
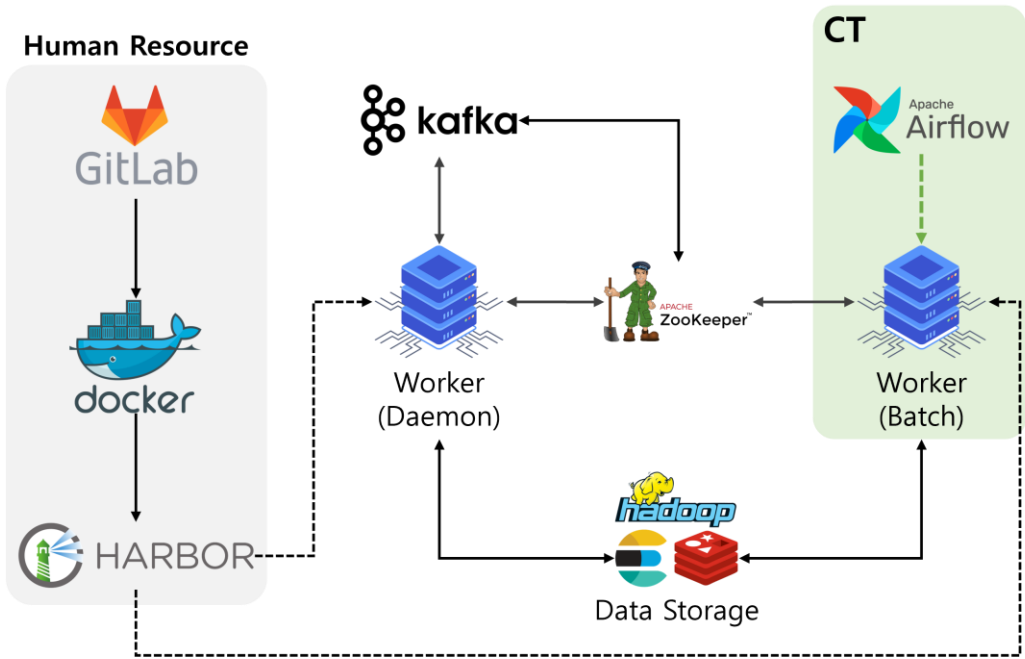
여러 지표 별 랭킹 수립 후 앙상블

- Beta Distribution
- Bayesian Average
- Log(impression) CTR
- ...

Online A/B Test

- PC & Mobile 나눠서 수행
- 평가 지표:
  - CTR
  - MAP
- 결과:  
기존 알고리즘 대비 CTR 상승률 10%

## MLflow 도입



요약	MLflow 도입 및 이중화
근거	ML/DL 연구 내용 자산화 및 모델 관리 필요
기술	MLflow, HDFS, PostgreSQL, nginx

# 메인&즐거찾기 페이지 Clip&Short-form 동영상 추천

The screenshot displays the SOOP website interface. On the left, there is a navigation menu with 'MY' highlighted in a red box. The main content area is divided into several sections:

- 유저클립 (User Clips):** A row of video thumbnails with titles such as '메수. [클립]배민정 vs 코코양 로키님 반응', '플라즈마단 [클립]민고 다영이 잠들어서 지터가 깨있는 게임 나는 상을 뽐내면 이게 / 다영) 나 최근에 했었는데 / 민고)', and '볼고기피자. [클립] 갓오브아레나 피해자 사정 (W백프로)'. This section is highlighted with a red border.
- VOD (Video On Demand):** A row of video thumbnails with titles like 'BJ백서 폭사 카페 경향전 악동 스위트 출동 [키즈네 클렌진 모집 중]', '간,김동하 간) 배그 1.1 갓 오브 아레나 우승 ^^', and '태민98 갓오브아레나 1타1결승'. This section is not highlighted.
- Catch (Catch-up):** A row of video thumbnails with titles such as '대체 뭐가 보이는데?', '[캐시]김민고 열광전 가성비대배', and '연상길러 DG98의 더블킬'. This section is highlighted with a red border.

At the bottom of the page, there is a URL: <https://vod.sooplive.co.kr/player/321333/catchstory?urllocation=my&StoryId=321333-321293-321271-321269-321073-321019-320479-320>

# 메인 & 즐겨찾기 페이지 Clip & Short-form 동영상 추천

## Model Architecture

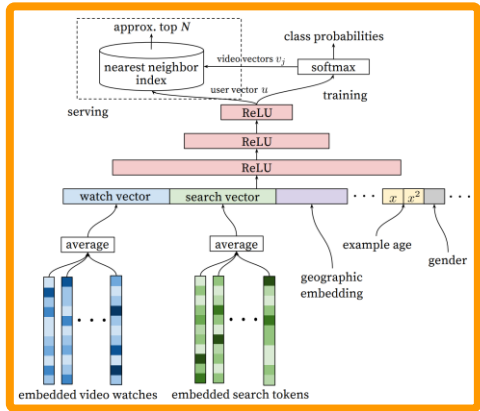
## Multi-task Loss (without task layers)

## Batch Softmax

## Post-processing (preference ranking)

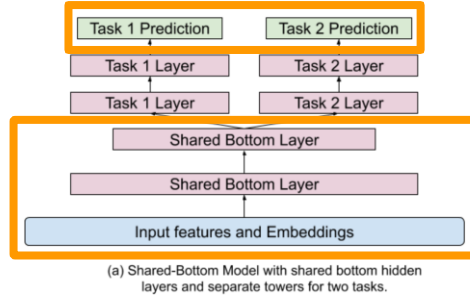
### Deep Neural Networks for YouTube Recommendations

Paul Covington, Jay Adams, Emre Sargin  
Google  
Mountain View, CA  
{pcovington, jka, msargin}@google.com



### Recommending What Video to Watch Next: A Multitask Ranking System

Zhe Zhao, Lichan Hong, Li Wei, Jilin Chen, Aniruddh Nath, Shawn Andrews, Aditee Kumthekar, Maheswaran Sathiamoorthy, Xinyang Yi, Ed Chi  
Google, Inc.  
{zhezhaol, lichan, liwei, jilinc, aniruddhnath, shawnandrews, aditeek, nlogn, xinyang, edchi}@google.com



### Sampling-Bias-Corrected Neural Modeling for Large Corpus Item Recommendations

Xinyang Yi, Ji Yang, Lichan Hong, Derek Zhiyuan Cheng, Lukasz Heldt, Aditee Kumthekar, Zhe Zhao, Li Wei, Ed Chi  
Google, Inc.  
{xinyang, jiyangyi, lichan, zcheng, heldt, aditeek, zhezhaol, liwei, edchi}@google.com

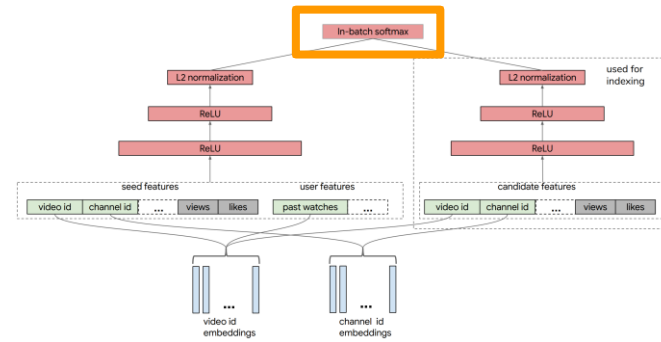


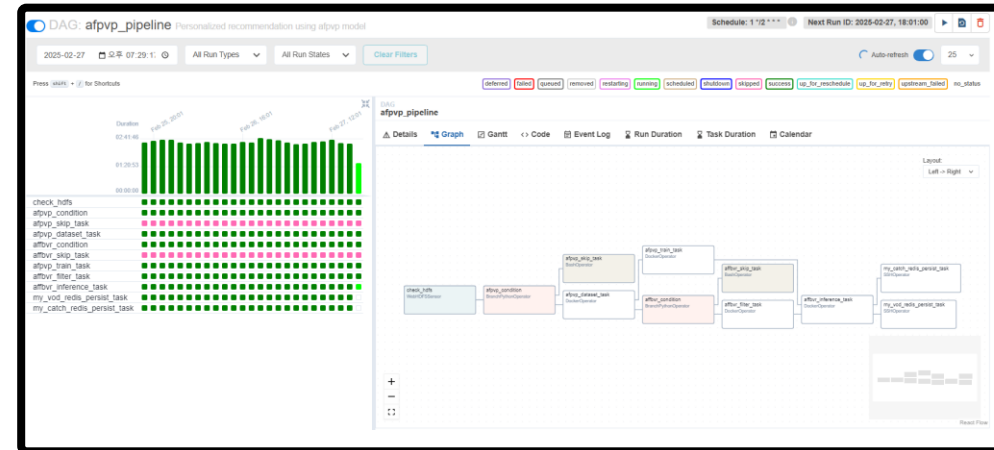
Figure 2: Illustration of the Neural Retrieval Model for YouTube.



Streamer Selection  
1. Short-term  
2. Long-term



## Airflow Dag Graph



요약

1. 서비스 메인 스포츠 동영상 추천 (3개월) > CTR 7% 상승
2. 즐겨찾기 페이지 클립 & 스포츠 동영상 추천 (3개월) > 재생 20% 상승

근거

늘 보는 스트리머 위주로 시청하는 유저에게 추천 할 데이터 필요

기술

Airflow, PyTorch, Hive

## 스트리머 탐험 로직 개발

요약	숏폼 VOD 연속 재생 횟수 증대를 위해, 기존 개인화 추천 데이터에 유저가 시청할 만한 다른 스트리머들의 최신/인기 VOD 추가
근거	다른 VOD 유형 대비, 숏폼을 시청할 때 비교적 다양한 스트리머의 VOD를 시청
방법	1. 네 가지 탐험 유형을 설계 (e.g., 함께 시청되는 스트리머) 2. 로그에 탐험 데이터 유형을 기록하여 평가
기술	Airflow, PySpark, Hive

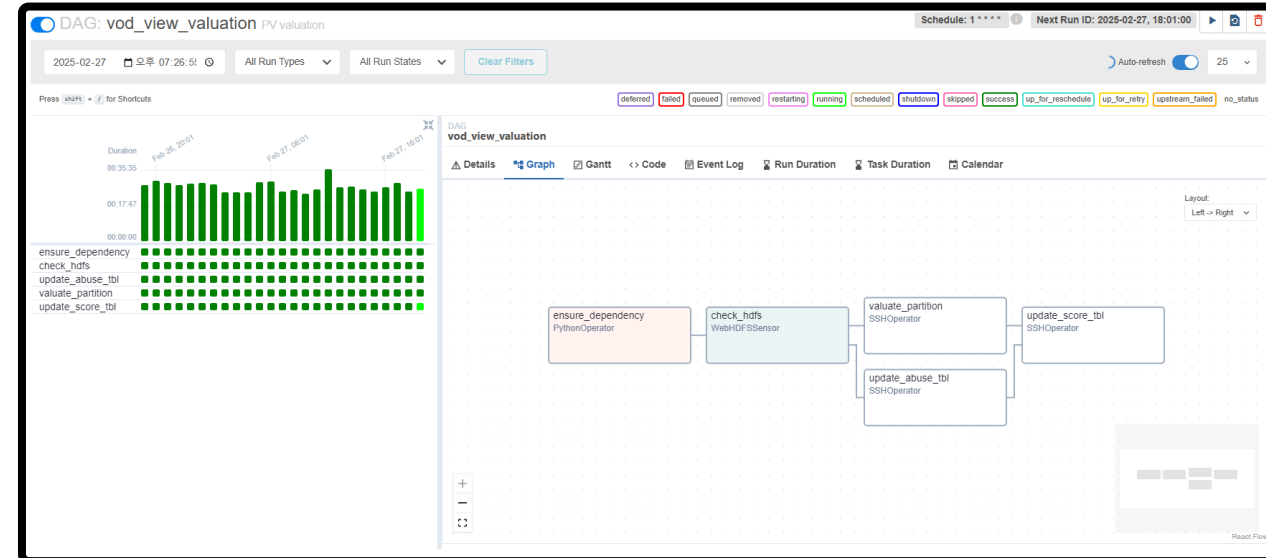
## 스트리머 별 대표 VOD 선정 로직 개발

요약	1. 스트리머와 관련된 모든 VOD들을 점수화 하여 2. 스트리머를 대표할 수 있는 n 개 VOD 선정
근거	스트리머 채널(개인 페이지) 상단에 노출 할 대표 VOD 선정 필요 (기존: 스트리머들의 방치 영역)
방법	VOD 평가를 위한 여러 통계 지표 생성 (최신성, 조회수, 양질성 등)
기술	Airflow, Hive, HBase, RDB

## 조회수 가치 평가

요약	시청 행위(초)의 가치를 절대적 수치가 아닌 상대적 비교를 통해 VOD의 품질을 분석
근거	<ol style="list-style-type: none"> <li>시청 행위는 implicit feedback 이므로, 시청 시간 만으로는 명확한 만족도 측정이 어려움</li> <li>다양한 어뷰징으로 인해 실질 지표 산출이 어려움</li> </ol>
핵심	<ul style="list-style-type: none"> <li>SOOP의 VOD 특징: 쉽게 라이브 방송에서 클릭 몇 번으로 VOD를 생성할 수 있음             <ul style="list-style-type: none"> <li>VOD의 퀄리티 보장이 어려움</li> </ul> </li> <li>다음을 활용하여 VOD의 품질을 예상 및 유저 만족도를 평가             <ul style="list-style-type: none"> <li>유저 별 평균 시청 시간</li> <li>item 별 평균 시청 되는 시간</li> </ul> </li> </ul>
기술	Airflow, PySpark(SQL), Hive

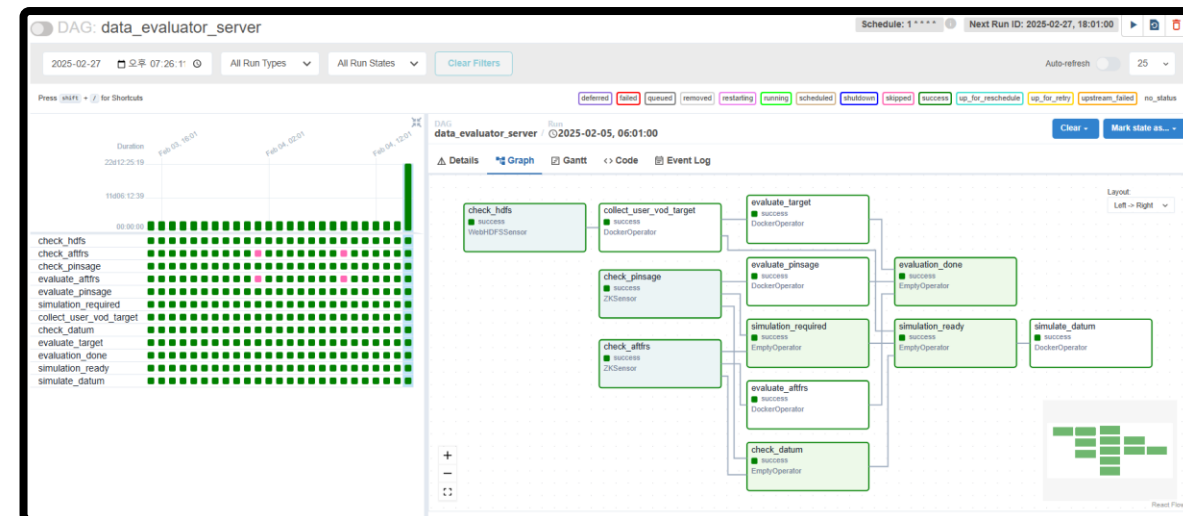
Airflow Dag Graph



## 추천 데이터 오프라인 시뮬레이션 평가

요약	추천 모델들의 데이터를 offline 에서 평가 및 비교
방법	로직: 1. Batch 생성되는 같은 시간 대의 추천 모델 별 데이터 수집 2. 다음 batch의 시청 이력을 토대로 모델 별 데이터 평가  지표: - 정확성 관련 지표 (HitRate/MAP 등) - 데이터 특성 관련 지표 (Freshness/Coverage 등)
기술	Airflow, SQLite, Redash

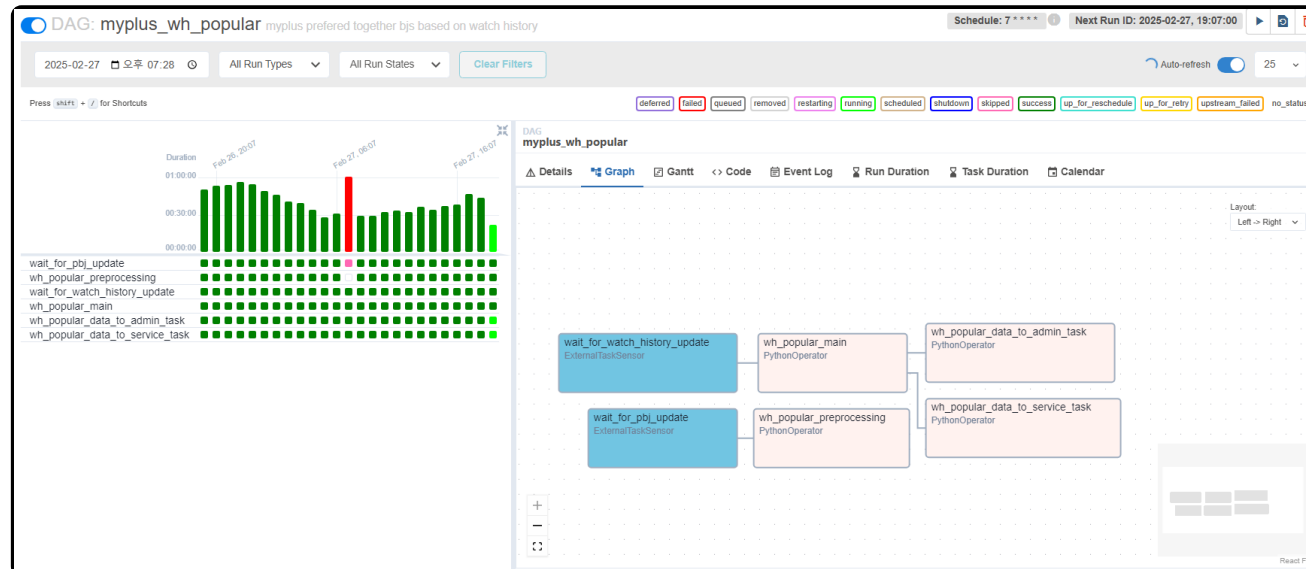
Airflow Dag Graph



## 라이브 방송 추천

요약	선호할 만 한 라이브 방송 추천 모듈 > CTR 3% 상승 및 현재 배포중
내용	<ul style="list-style-type: none"><li>- 타겟 유저: 방송 카테고리 가 아닌, 스트리머 위주의 성향을 띄는 유저</li><li>- 방송 선정: 타겟 유저가 선호하는 스트리머들을 선호하는, 다른 사람들의 선호 스트리머를 통계</li></ul>
기술	Airflow, HDFS

### Airflow Dag Graph





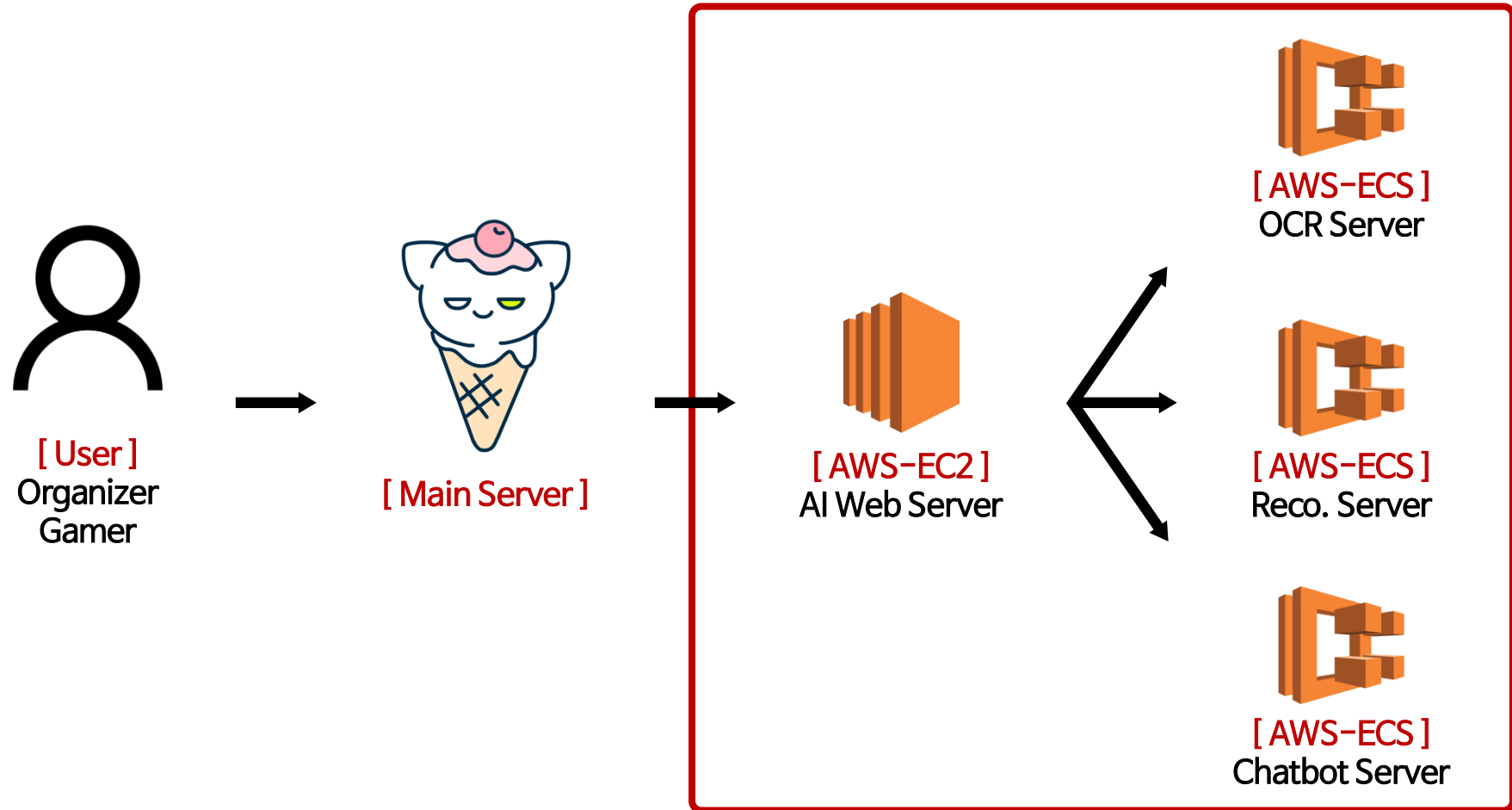
## Undefined

[1] Match Result Recorder

[2] Competition Rule Recommendation

[3] FAQ Chatbot

# AI Serving Pipeline



Match Result Recorder (OCR)	Model	<a href="#">Tesseract</a> (Google, LSTM-based)
	Works	<ul style="list-style-type: none"> <li>• Define Problem</li> <li>• Define Pipeline                             <ul style="list-style-type: none"> <li>• Our Tesseract Model</li> <li>• Cloud API (in case of poor confidence)</li> <li>• Finetuning</li> </ul> </li> <li>• Model Serving</li> </ul>
Competition Rule Recommendation	Model	Matrix Factorization (Alternative Least Squares)
	Works	<ul style="list-style-type: none"> <li>• Define Problem</li> <li>• EDA and Feature Selection (via Correlations)</li> <li>• Model Selection/Tuning</li> <li>• Model Optimization (remove operations)</li> <li>• Model Serving</li> </ul>
Chatbot	Model	Multi-lingual BERT, <a href="#">StarSpace</a> (Facebook)
	Works	<ul style="list-style-type: none"> <li>• Dataset Preprocessing</li> <li>• Model Selection/Tuning</li> <li>• Model Serving</li> </ul>

## ML System Lab.

- [1] Network Embedding Generation
- [2] DNN Model Quantization - 1
- [3] DNN Model Quantization - 2
- [4] Artificial Intelligence Assistant

## Network Embedding Generation

\* Published in 2022 **BIB** (Briefings in Bioinformatics) **Journal**

Paper

[link](#)

Github  
Code

[link](#)

### Project description

[Human Cell lines - Cancer Drugs] Response Prediction

Network (graph) dataset consist of

- Cell line nodes
- Drug nodes
- Protein nodes (connected to Cell lines)

My Task: Train embedding vectors of Cell lines and Drugs

### Problem

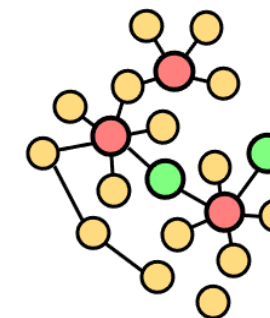
Extremely unbalanced dataset

- About 20,000 Protein nodes
- About 900 Cell line nodes
- About 300 Drug nodes

Fails to reflect the relationships between Cell lines & Drugs

As a result, we got poor response prediction performance

● Cell line ● Drug ● Protein



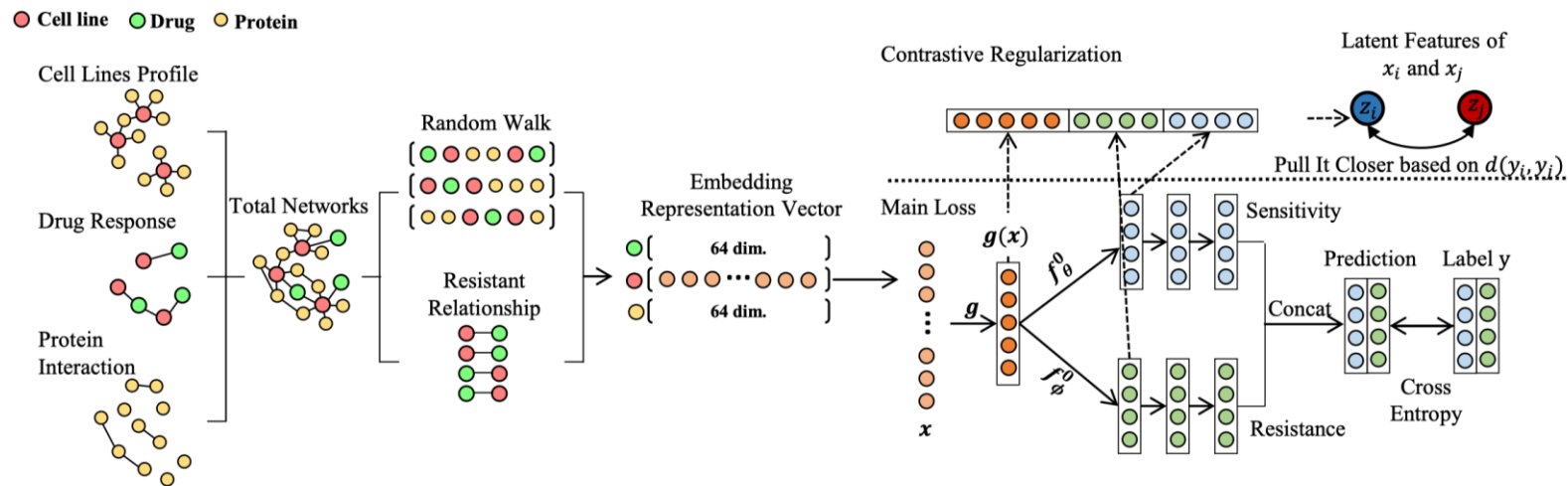
# Solution

Make training process to focus on relationships between Cell lines & Drugs

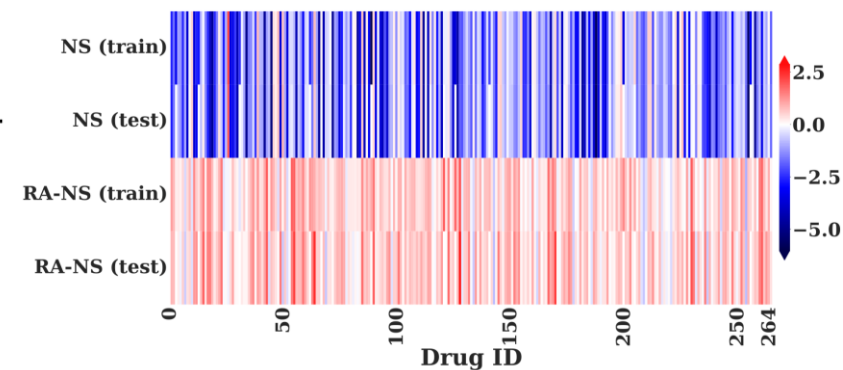
## Response-aware Negative Sampling (RA-NS)

- Cell line & Drug nodes use resistant Drug & Cell line nodes as their negative samples

\* Tested Models: Node2Vec, Graph Convolutional Network, Graph Transformer Network



**Fig. 1.** The framework of RAMP. RAMP consists of two main stages. First, representation vectors are extracted from heterogeneous networks with RA-NS. Second, the multitask architecture of a Bayesian neural network is trained by representation vectors with contrastive regularization.

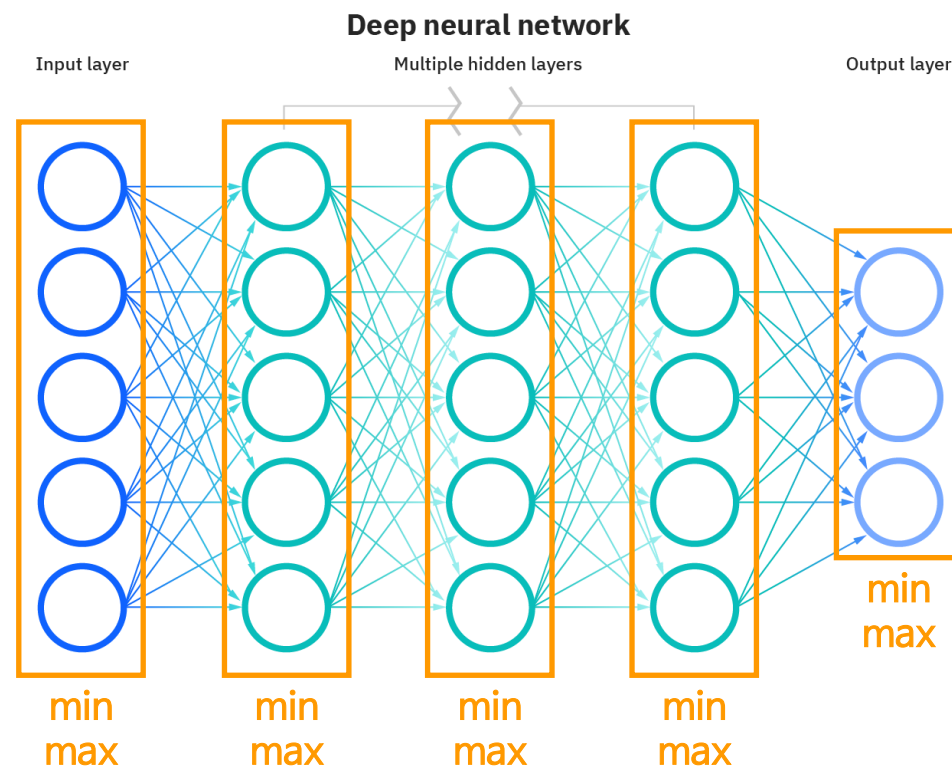


**Fig. 2.** Embedding similarities among drug and cellines. We subtract the similarity of a drug and its resistant cell lines from the similarity of the drug and its responsive cell lines. The results are normalized and plotted in a heatmap format. The higher (or redder) the value is, the better the embedding reflects the network structure.

## DNN Model Quantization - 1

Definition	What is Quantization	General DNN models use Float32 type variables
		Quantized models use low-bit INT types at inference
	What for	<ul style="list-style-type: none"><li>• Model storage</li><li>• In memory load</li><li>• Matrix multiplication</li></ul> with Float32 type cause bottleneck/unusability in low performance H/W

Problem	Poor Performance	Quantized models' performance (e.g., accuracy) drops catastrophically when using sub-8bit INT type
	Why	Too generalized Quantization parameters <ul style="list-style-type: none"><li>Quantization parameters require: <b>Per layer avg-ed min/max range</b> of intermediate outputs across datasets</li></ul>
		Averaged min/max values include outliers



## Solution

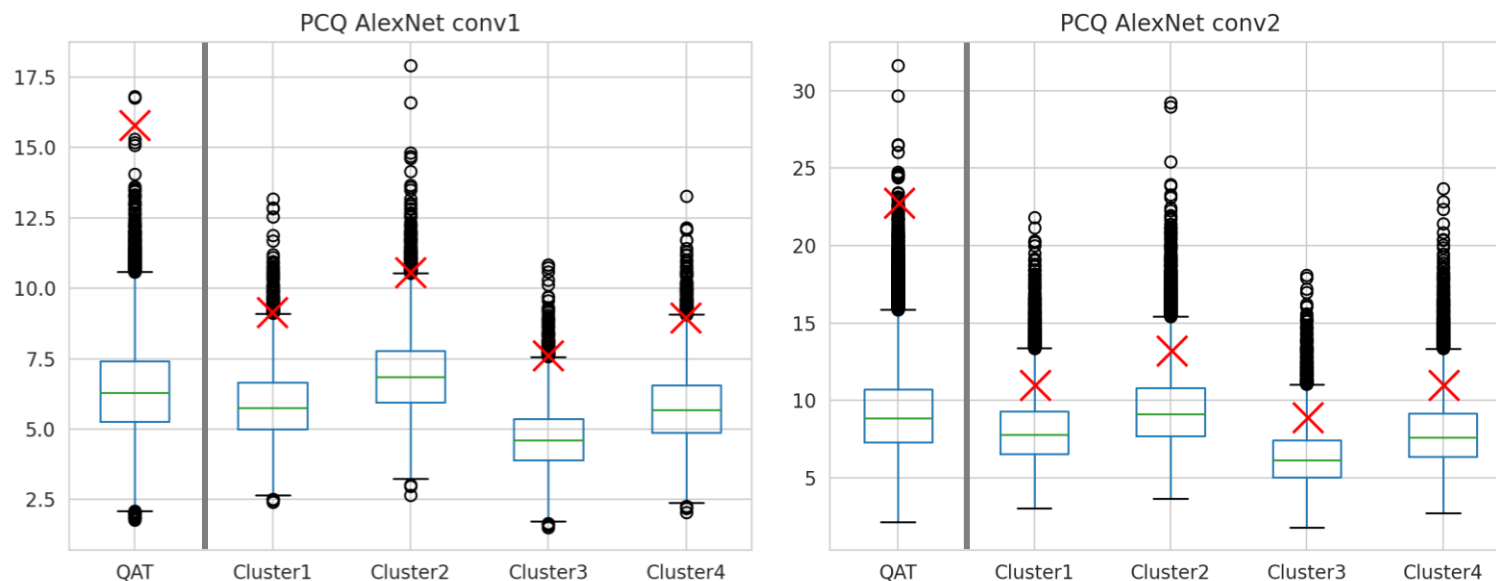
### Granular Exponential Moving Average (Granular EMA)

Train Quantization Parameters while **excluding outliers**

### Neural Network Aware Clustering (NNAC)

Train Quantization Parameters separately **across clusters of input images**

- Some data might need **shorter min/max range**
- Shorter range means **less information loss**



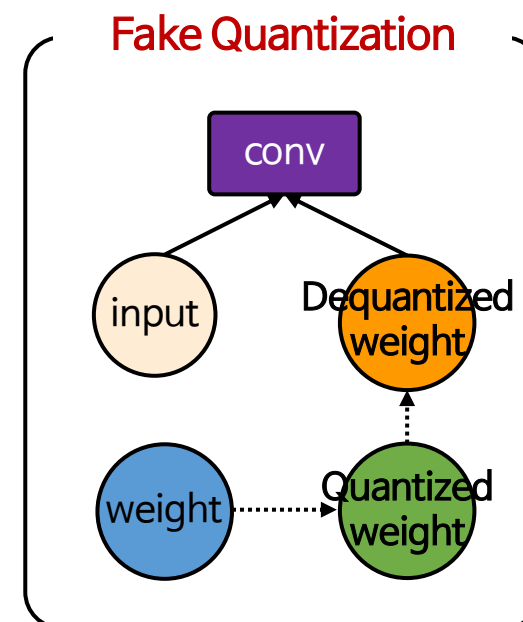
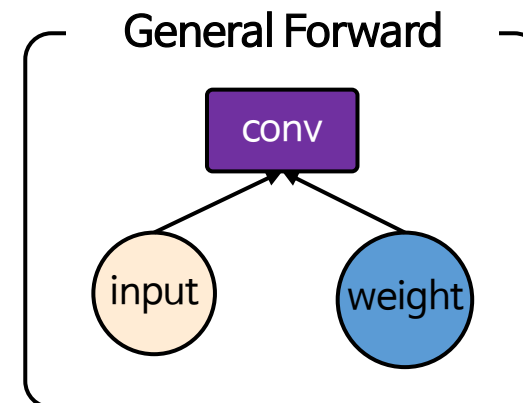
### Figures' Description

- Shows that our method
  - how efficiently exclude outliers
  - how to work with clusters
- QAT : Baseline (Google)
- Cluster\* : Ours
- X : Trained maximum value
- Box-plots : Actual max values per image

## DNN Model Quantization – 2

\* Published in 2022 ICEIC (International Conference on Electronics, Information, and Communication)

Problem	<p><b>Quantization Aware Training (Google)</b></p> <ul style="list-style-type: none"> <li>• Fake-quantize all of the weight matrices with a single low-bit type</li> <li>• Too much <b>quantization errors</b> occur and the trained model gets ruined</li> </ul>
Solution	<p><b>QuantNoise (Facebook)</b></p> <ul style="list-style-type: none"> <li>• Fake-quantize probabilistically selected subsets of matrices (a subset per matrix)</li> <li>• Trained models <b>under-prepared</b> for Quantization</li> </ul>
Solution	<p><b>Fake Single Precision Training (FST)</b></p> <ul style="list-style-type: none"> <li>• Probabilistically select subsets of weight matrices as QuantNoise</li> <li>• Fake-quantize <b>selected subsets</b> with <b>low-bit type</b></li> <li>• Fake-quantize <b>the rests</b> with <b>higher bit type</b> than the selected</li> </ul>



## Artificial Intelligence Assistant

- AI Assistant App, Almond

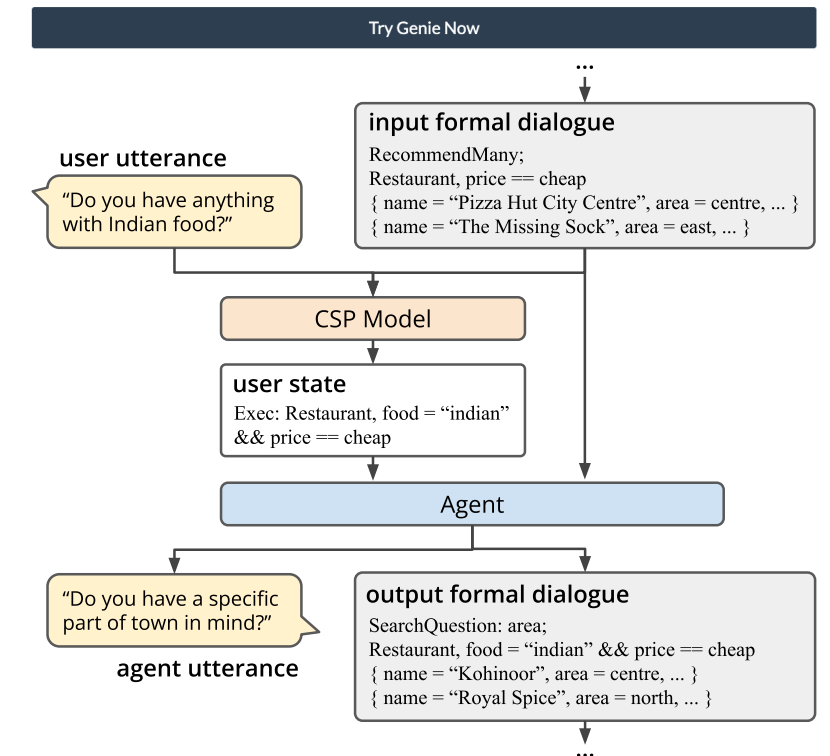
- Currently, the service name has been modified to [Genie](#)
- Developed by Stanford OVAL Lab

- Training Korean Seq2SQL Model

- Dataset preparation
  - Web Crawling
  - Construct templates of sentences (example of sentences)
  - Augment sentences based on templates
- Train & serve model



The Open, Privacy-Preserving Virtual Assistant





## Kakao

[1] Automobile Video Recommendation

[2] Comics Recommendation

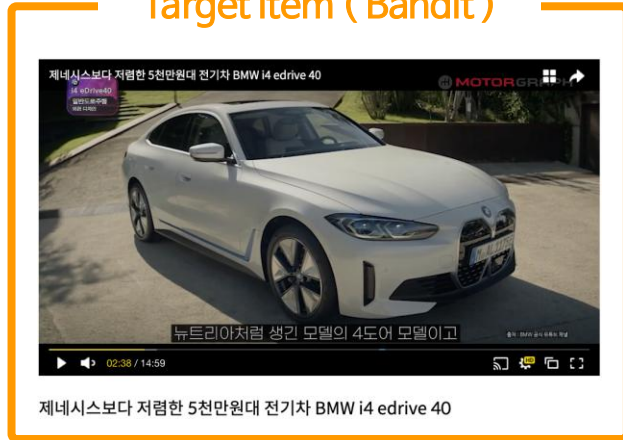
## Automobile Video Recommendation

Exp 1, 2	Thomson Sampling h-params tuning	Purpose	Adjustment of trade-off between exploration & exploitation
		Reason	[Exp-1] High <b>matrix sparsity</b>
			[Exp-2] Considering <b>time bias</b> enhanced by low traffic
Exp 3, 4	Ranking Algorithm (RRF to Weighted- sum)	Purpose	Searching the key model among ensembled models
		Reason	Other well performing services had been used <b>similar model combination</b> <ul style="list-style-type: none"><li>• Therefore, assumed that the composition of used models are good enough</li></ul>
Exp 5	Item2Vec instead of Matrix Factorization	Purpose	Overcome Matrix Factorization model's limitation
		Reason	Needed to generate reco. results <b>within limited item list</b> <ul style="list-style-type: none"><li>• The limited items rated 30~40th on avg., if we force the limitation off</li></ul>
			Needed some models which <b>capture information</b> which MF can't

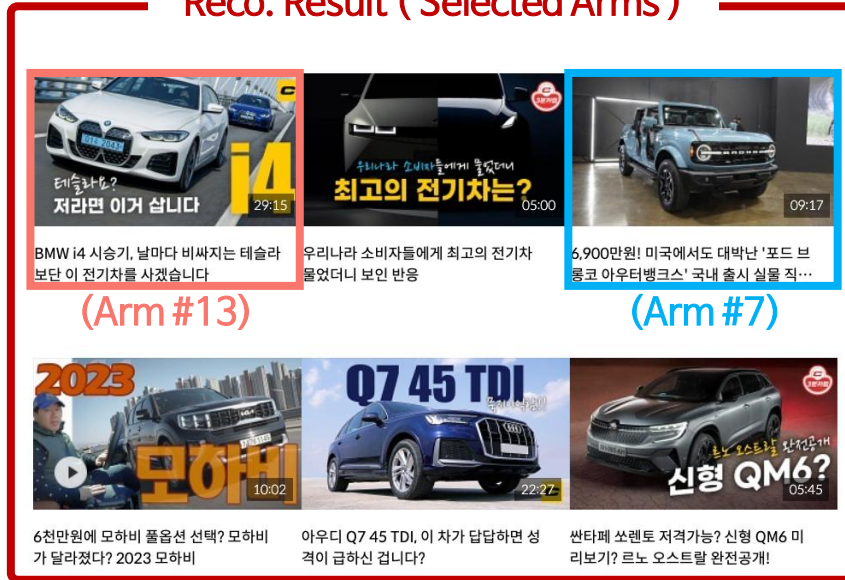
## Automobile Video Recommendation

Exp 1, 2	Thomson Sampling h-params tuning	Purpose	Adjustment of trade-off between exploration & exploitation
		Reason	[Exp-1] High matrix sparsity [Exp-2] Considering time bias enhanced by low traffic

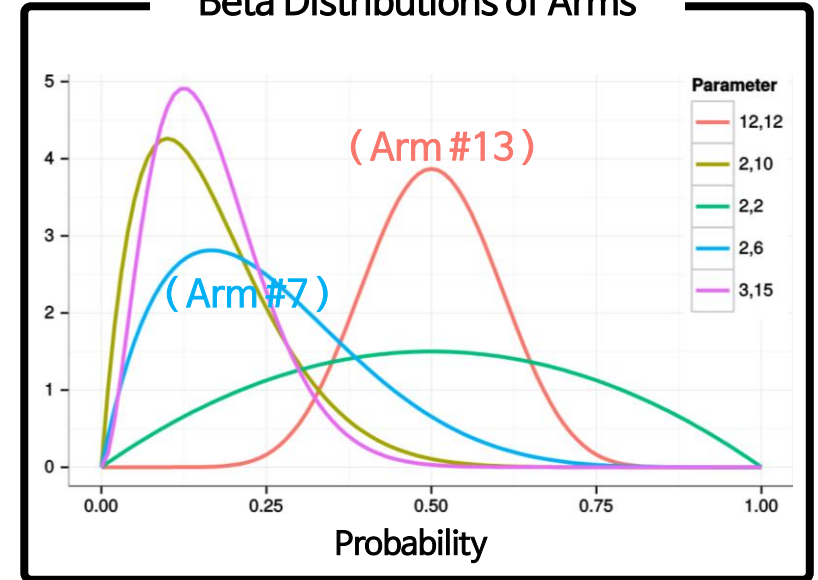
Target Item (Bandit)



Reco. Result (Selected Arms)

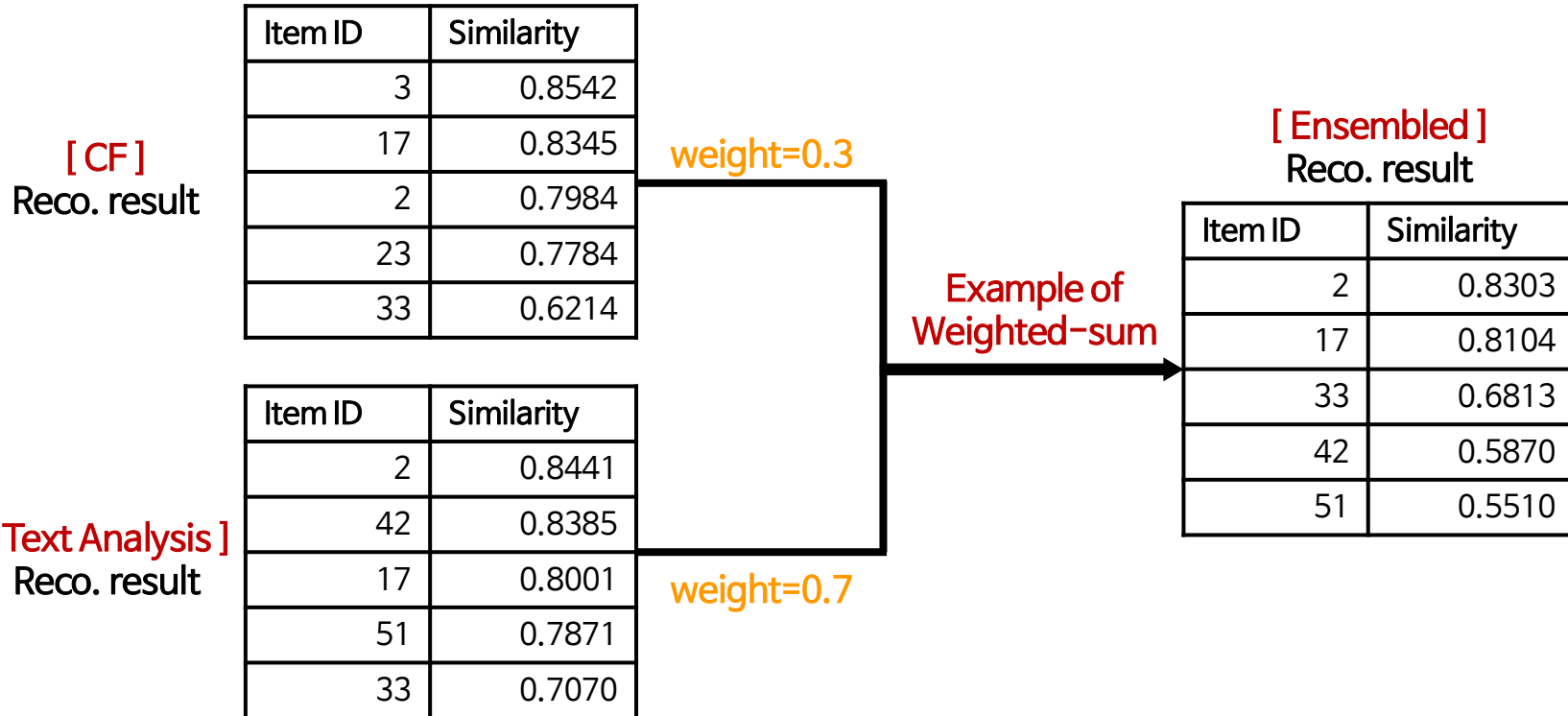


Beta Distributions of Arms



## Automobile Video Recommendation

Exp 3, 4	Ranking Algorithm (RRF to Weighted-sum)	Purpose	Searching the key model among ensembled models
		Reason	Other well performing services had been used <b>similar model combination</b> <ul style="list-style-type: none"> <li>Therefore, assumed that the composition of used models are good enough</li> </ul>



## Automobile Video Recommendation

Exp 5	Item2Vec instead of Matrix Factorization	Purpose	Overcome Matrix Factorization model's limitation
		Reason	Needed to generate reco. results <b>within limited item list</b> <ul style="list-style-type: none"> <li>The limited items rated 30~40th on avg., if we force the limitation off</li> </ul>
			Needed some models which <b>capture information</b> which MF can't

〈 MF Model's Reward Matrix 〉

				
John 	5	1	3	5
Tom 	?	?	?	2
Alice 	4	?	3	?



〈 Item2Vec Model's Input Sequence 〉

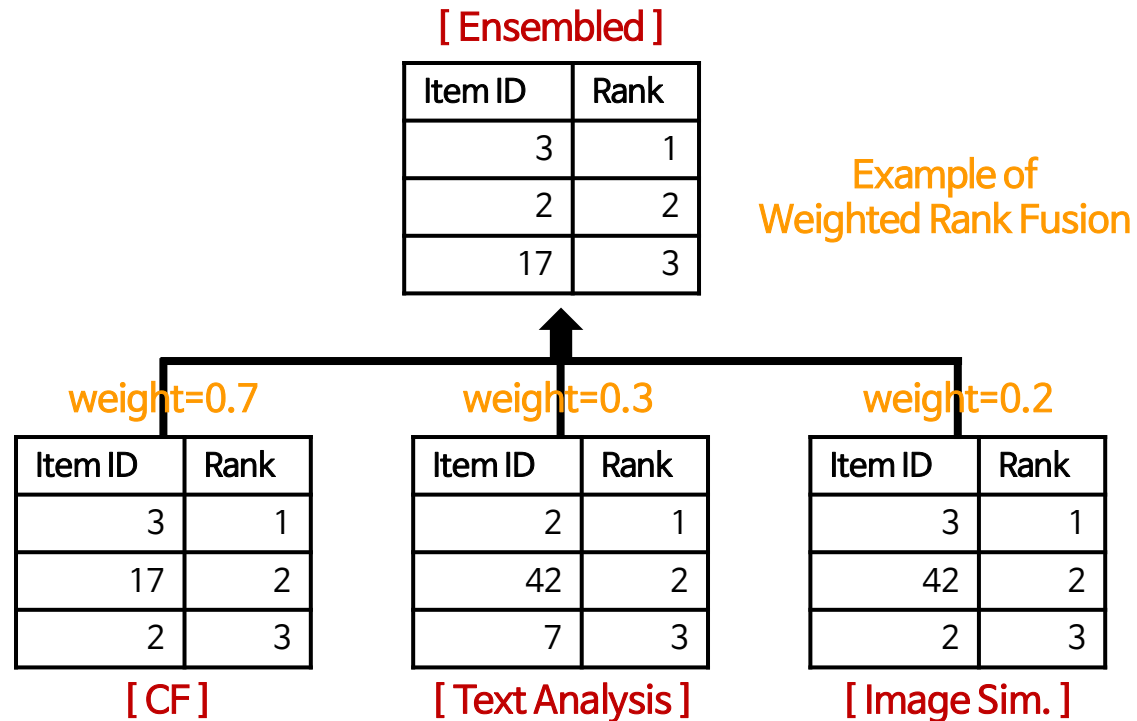


## Comics Recommendation

Exp 6	Word2Vec input dataset reconstruction	Purpose	Better reflection of Japanese characteristics
		Reason	Previously, model used <b>nouns</b> and <b>pronouns</b> only
			According to past researches, <b>verbs</b> and <b>adjectives</b> are also important for JP
Exp 7	Modified ranking algorithm (RRF to WRF)	Purpose	Strengthen the key model
		Reason	By previous experiment logs, the only MF used reco. pipeline without ensemble method outperformed ensembled pipeline
			But the ranking algorithm the system was using weakened MF's power

## Comics Recommendation

Exp 7	Modified ranking algorithm to Weighted Rank Fusion	Purpose	Strengthen the key by giving weight to rank values
		Reason	By previous experiment logs, the only MF used reco. pipeline without ensemble method outperformed ensembled pipeline
			But the the Weighted-sum Ranking Algorithm weakened MF's power





## HanbitSoft

[1] (KR) Multi-speaker Speech Synthesis Model

[2] (EN) Text/Audio Chatbot

## (KR) Multi-speaker Speech Synthesis Model

- Dataset preparation

Web Crawling	Audio files
	Script files
Preprocessing	Cut audio files into files of sentences
	Cut script files into sentences (by comparing STT results)

- H-params optimization

- Demo [https://jarvis08.github.io/pjt\\_hbs\\_multi.html](https://jarvis08.github.io/pjt_hbs_multi.html)

## (EN) Text/Audio Chatbot

