# A Semantic Alignment System for Multilingual Query-Product Retrieval

## Team **www**

Qi Zhang, Zijian Yang, Yilun Huang, Zijian Cai, Ze Chen
OPDAI, Interactive Entertainment Group of Netease Inc., Guangzhou, China

https://www.kdd.org/kdd2022/

# Outline

- Task Introduction
- Data Analysis
- Overall Framework
    - Data Processing
    - Model Architecture
- Basic models
- Training Optimization
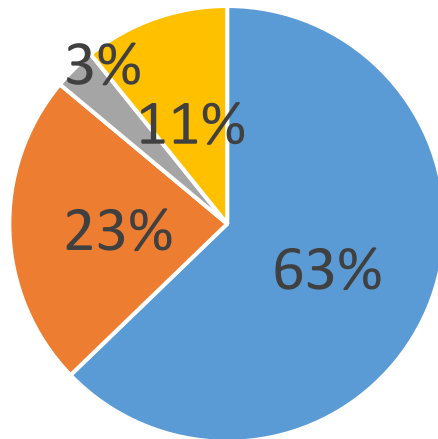- Results
- Summary and Future Work

# Task Introduction

- Task 1 aims at ranking the query-product pairs by **relevance**.
- Query-product dataset consists of **English, Spanish and Japanese**.
- Query-product dataset is classified into **Exact, Substitute, Complement,** or **Irrelevant** (ESCI) categories.
- Evaluation is based on **NDCG**.
- How to better understand the query-product **semantic relevance is the major challenge**.
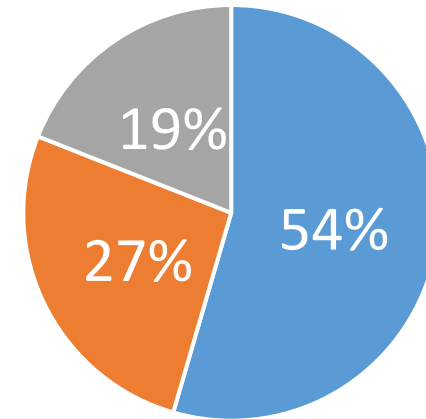
# Data Analysis

ESCI Label



3%
11%
23%
63%

■ Exact  ■ Substitute  ■ Complement  ■ Irrelevabt

**label distribution**

Language



19%
27%
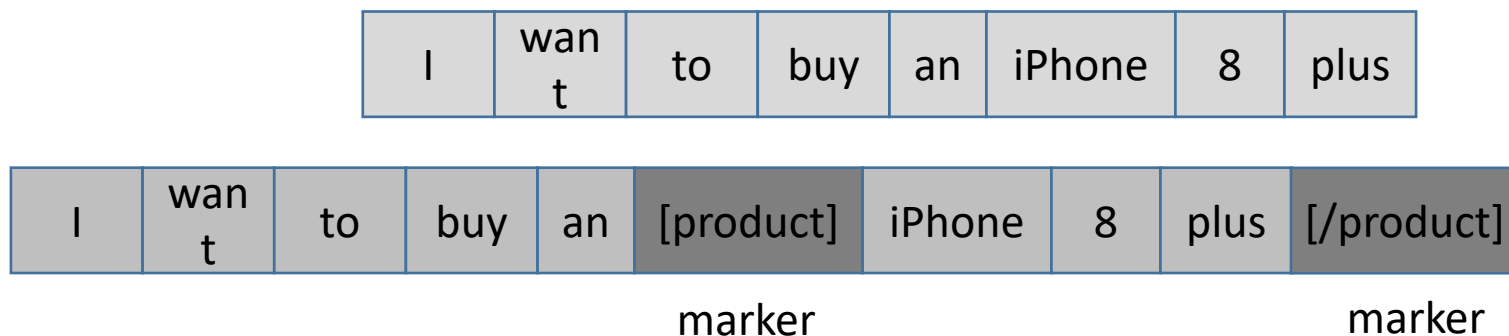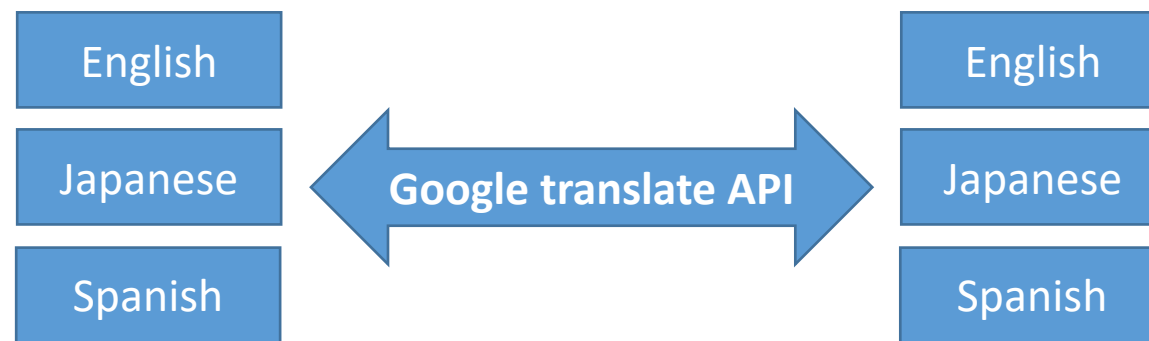54%

■ Enlish  ■ Japanese  ■ Spanish

**language distribution**

➢ The label distribution and the language distribution are both imbalanced.
➢ And according to our statistics, 54% of product brands focus on providing only one product, while only 7.1% of brands provide more than 10 products.
➢ More than 80% of the color names are only customized for a single product.
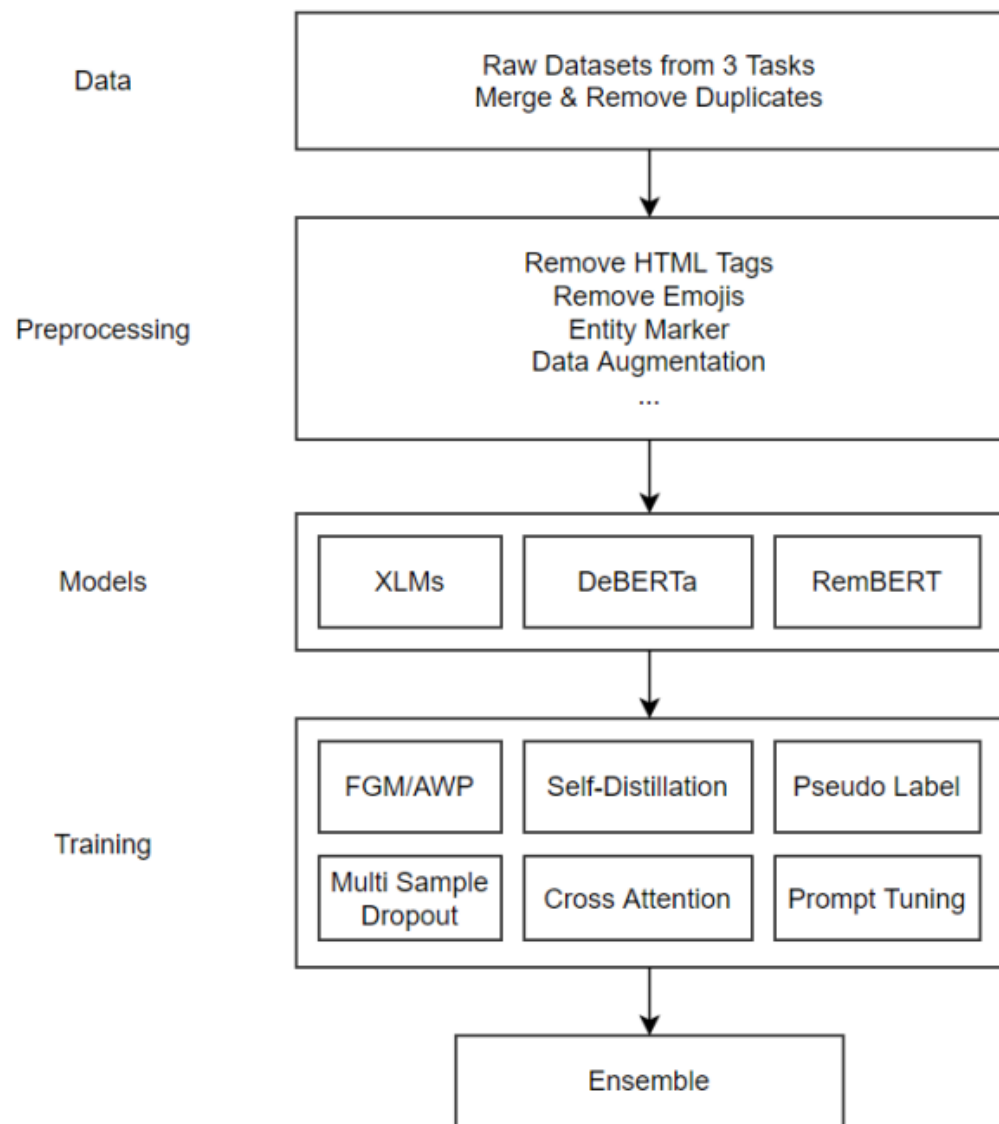
# Data Preprocess

- Remove those
  **HTML marks** and **emoji**.

- Translate all of the data into
  English, Spanish and Japanese
  separately to do **data augmentation**.

- Incorporate the NER information
  using **Entity Marker**.



Geeetech I3 pro W impresora 3D, está diseñado y fabricado por Shenzhen Getech Technology Co., Ltd <br> <br> Con su módulo Wi-Fi y la solución de impresión en nube 3D, puede actualizar I3 pro W para controlar directamente todo el proceso de impresión y compartir su experiencia de impresión a través de la aplicación en cualquier lugar y a cualquier hora. <br> <br> <b>Especificaciones de impresión:</b><br> Tecnología de impresión: FFF / FDM<br> Volumen de construcción: 200 x 200 x 180 mm (7,9 '' * 7,9 '' * 7,1 ''))<br> Resolución de la capa: 0.1-0.3mm<br> Precisión de posicionamiento: 0.1-0.3mm<br> Diámetro del filamento: 1.75mm<br> Diámetro de la boquilla: 0.3mm<br> Tipo de filamento: ABS / PLA / Flexible PLA <br> <br> <b>Software:</b><br> Sistema operativo: Windows / Mac / Linux<br> Aplicación Easy Print 3D<br> Software de control:
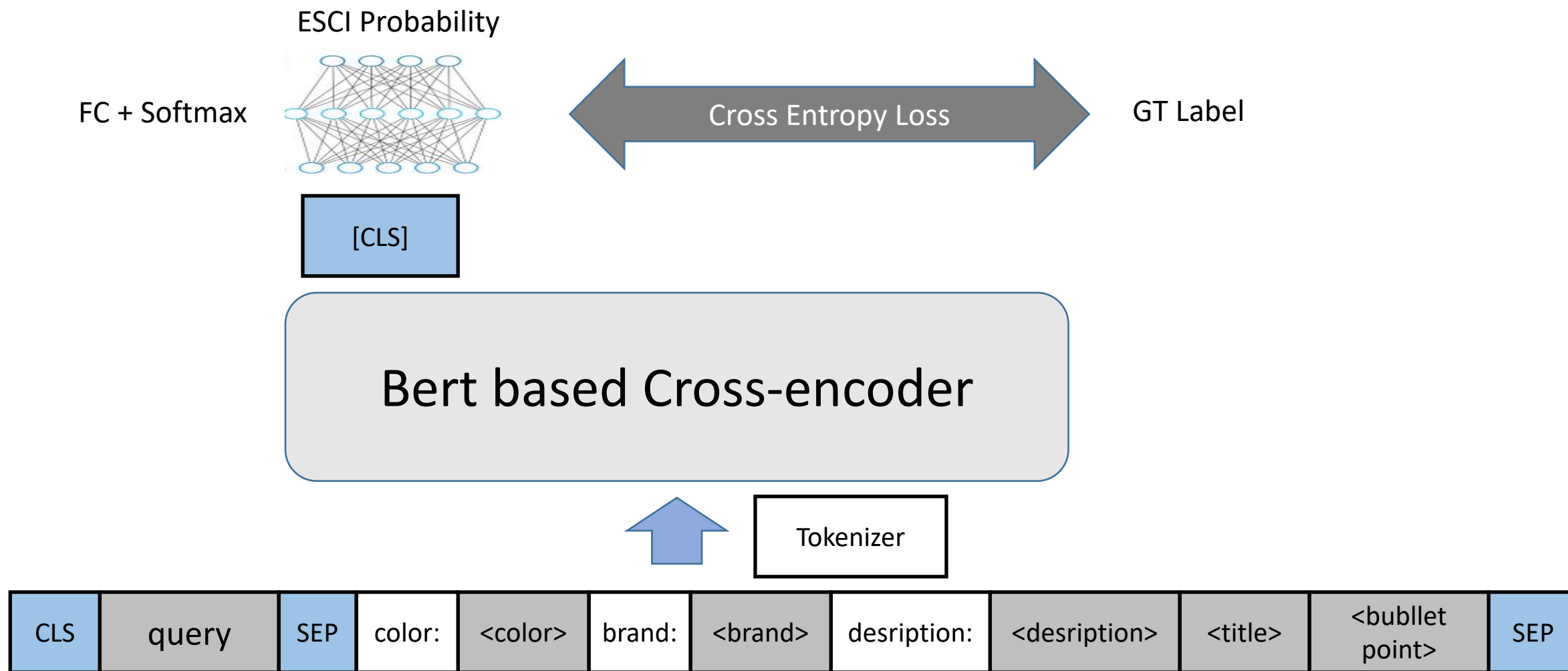
| English | |
| Japanese | **Google translate API** |
| Spanish | |

| English |
| Japanese |
| Spanish |

| I | want | to | buy | an | iPhone | 8 | plus |

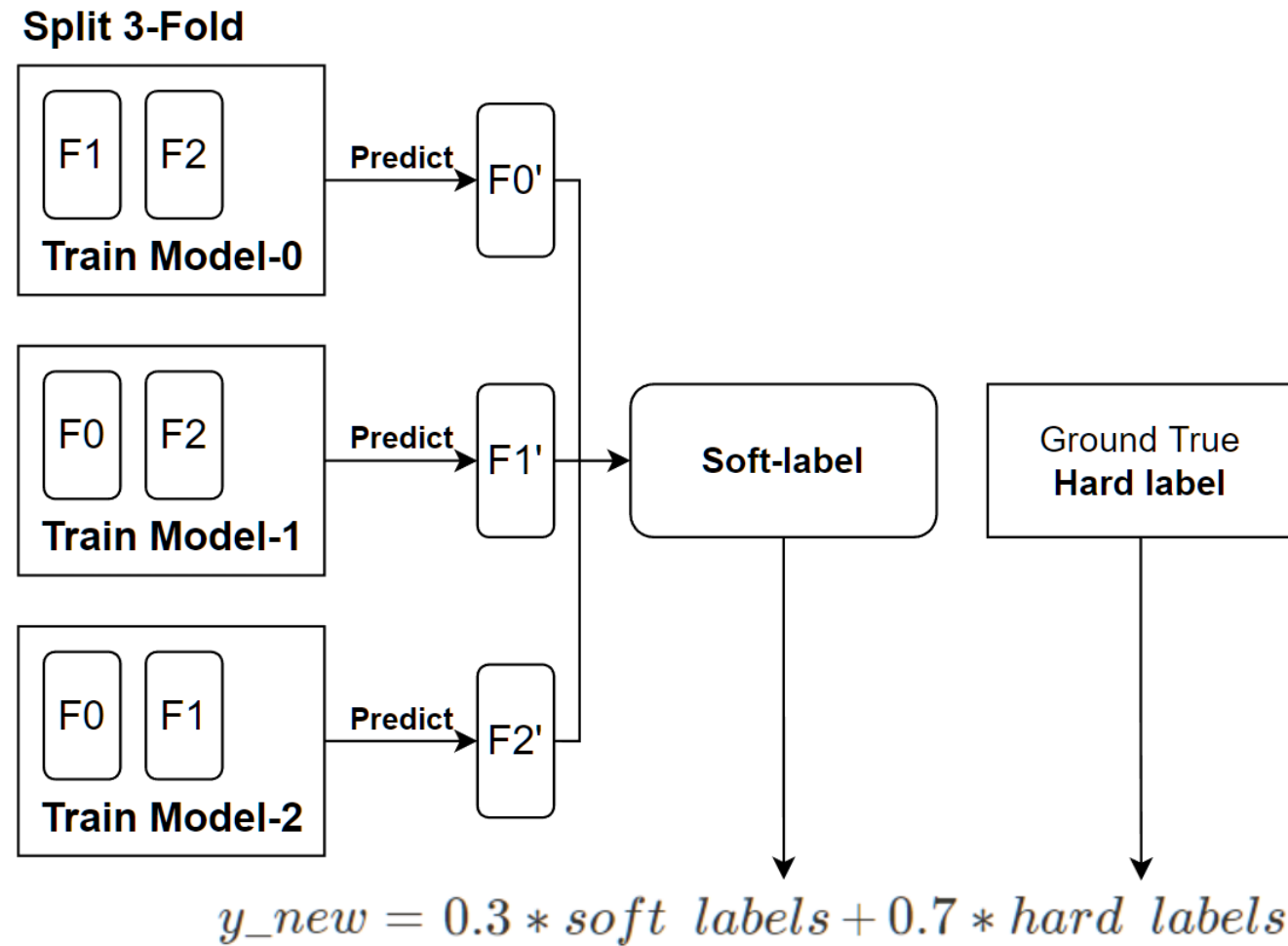| I | want | to | buy | an | [product] | iPhone | 8 | plus | [/product] |

marker                    marker

# Overall Framework



> Step 1: Data merge from 3 tasks.
> Step 2: Data preprocess: cleaning, augmentation, entity marker.
> Step 3: Model fine tune from both multilingual LMs and monolingual LMs.
> Step 4: Different training strategies applied.
> Step 5: Ensemble results from different models and strategies.

# Basic Models

ESCI Probability

FC + Softmax

Cross Entropy Loss

GT Label

[CLS]

Bert based Cross-encoder

Tokenizer

| CLS | query | SEP | color: | <color> | brand: | <brand> | desription: | <desription> | <title> | <bubllet point> | SEP |

# Training Optimization



Split 3-Fold

F1  F2
Train Model-0  →Predict→ F0'

F0  F2
Train Model-1  →Predict→ F1'  → Soft-label

F0  F1
Train Model-2  →Predict→ F2'

Ground True Hard label

$$y\_new = 0.3 * soft\ labels + 0.7 * hard\ labels$$

- **Self Distillation**

To be specific, we use 3-fold bagging training and make prediction on the out-of-fold datasets to generate the **soft labels**.

And then we merge the soft labels with the ground true hard labels with weights 0.3 and 0.7 to get the new training labels.
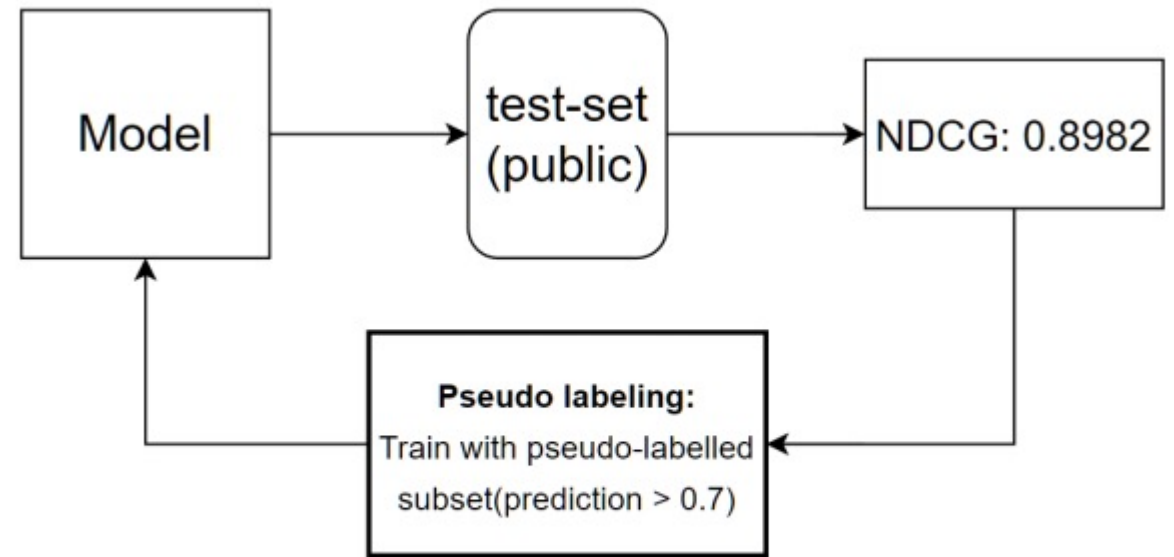
# Training Optimization



Figure 3: Train model with pseudo-labelled subset

- **Pseudo Labeling**

To avoid making the training data more noisy, only samples from the public test set with predicted probabilities **above 0.7** are used as pseudo labels.

And soft labels work better than hard labels during most of our experiments, we guess that hard labels may increase the risk of over-fitting.

# Training Optimization

- **Adversarial Training**

To gain robustness of models, we use Adversarial Weight Perturbation (AWP) in training steps that adversarially perturbs both model weights and the embeddings when the loss is below some threshold (like 0.6).

Besides, we also tried Fast Gradient Method (FGM) which performs slightly worse than AWP does in public leaderboard.

| Methodology | NDCG (Public) |
|:---:|:---:|
| AWP | 0.9022 |
| FGM | 0.9019 |

# Training Optimization

- **Multi-sample dropout & Grouped layer-wise learning rate decay**

There are several effective regularization learning strategies to avoid overfitting of deep neural network, which can not only accelerate training and improve generalization ability, but also achieve lower error rates and losses.
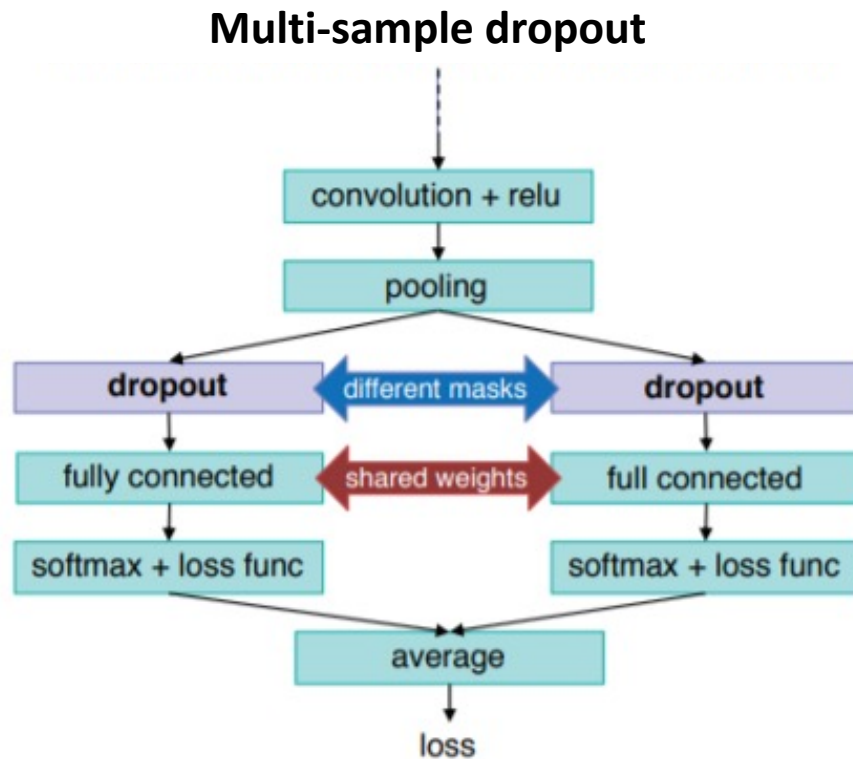
**Multi-sample dropout**



**Table: Grouped layer-wise learning rate decay**

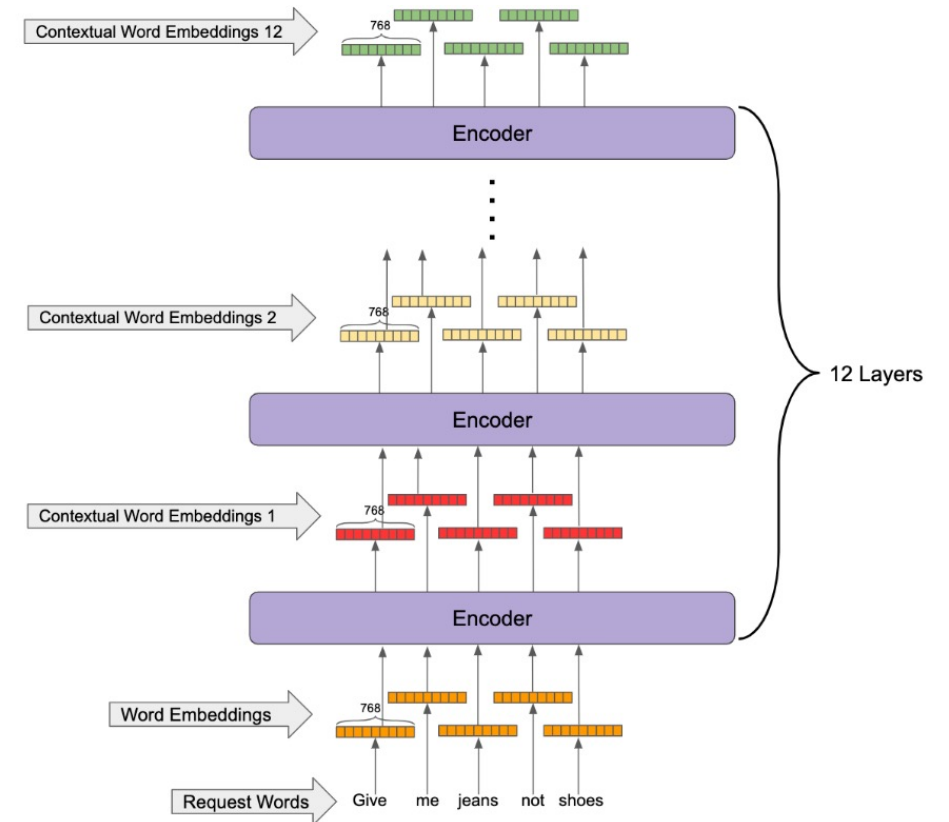| Model Layers | Learning Rate |
|:---:|:---:|
| 0-5 | 5e-6 |
| 6-11 | 1e-5 |
| 12-17 | 1e-5 |
| 18-23 | 2e-5 |

# Training Optimization

- **Weighted multi-layer Pooling**

Utilizing intermediate representations from various layers always provide better performance as it can help in incorporating more information.
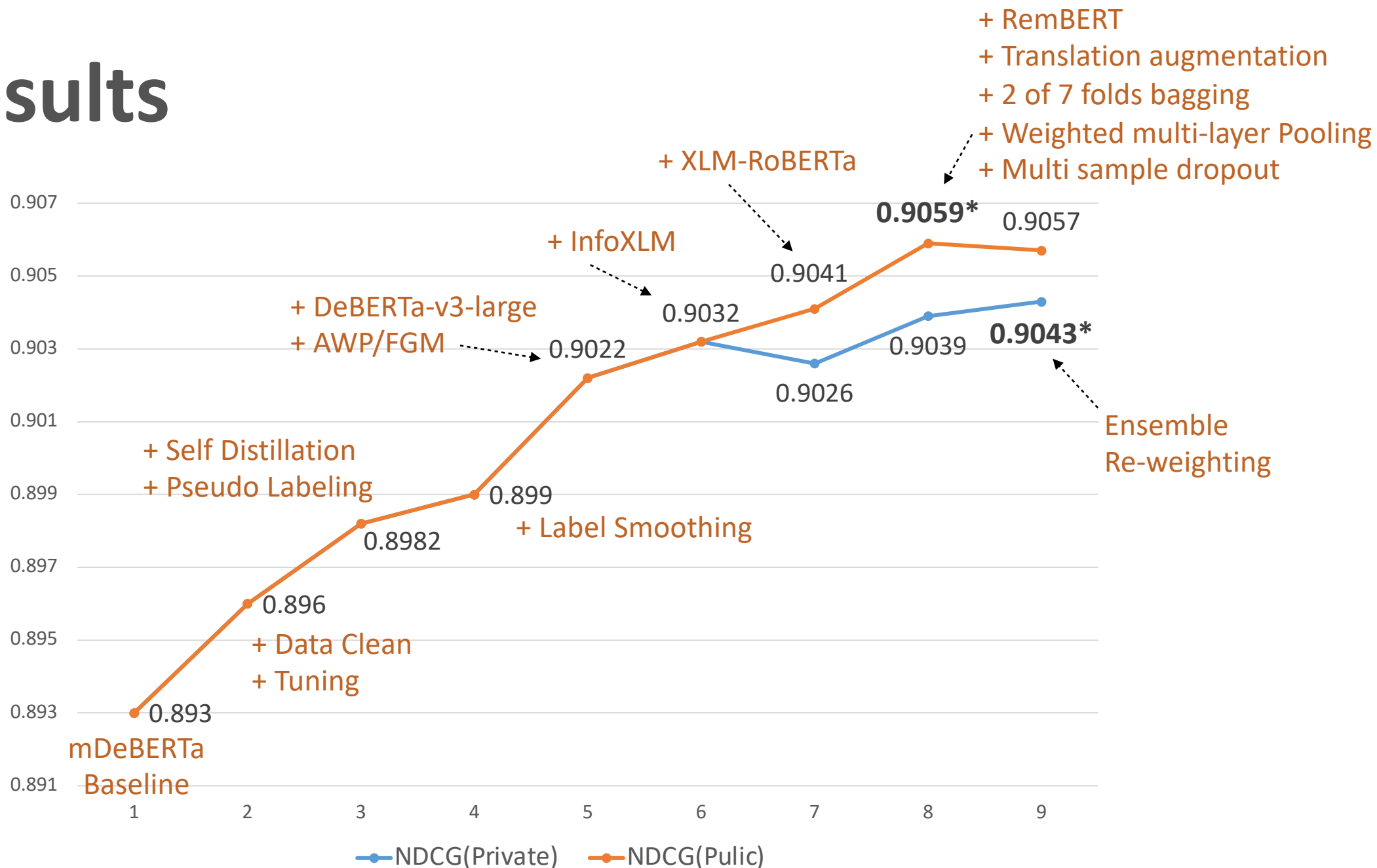
# Ensemble

- Ensemble weights are mainly determined by the public scores and also the local cross-validation scores.

- Lower the weights of the models with high correlation coefficients.

- Our score is improved from **0.9022** to **0.9057** on the public leaderboard, and from **0.9015** to **0.9043** on the private leaderboard after ensemble.

| NDCG (Public) | NDCG (Private) |
|---|---|
| 0.9057 | 0.9043 |

# Results

# Summary and Future Work

- **Summary**
  - We use multilingual and English pre-trained LMs as backbone, with the combination of data processing and sorts of training optimization.
  - For single model, we achieve NDCG score of **0.9022** on the public leaderboard and **0.9015** on the private leaderboard.
  - At last, we do model ensemble to get the final boost from **0.9015** to **0.9043** on the private leaderboard, which ensures us to win the first place.
- **Future Work**
  - End to end multilingual model solution.

# Thank You!