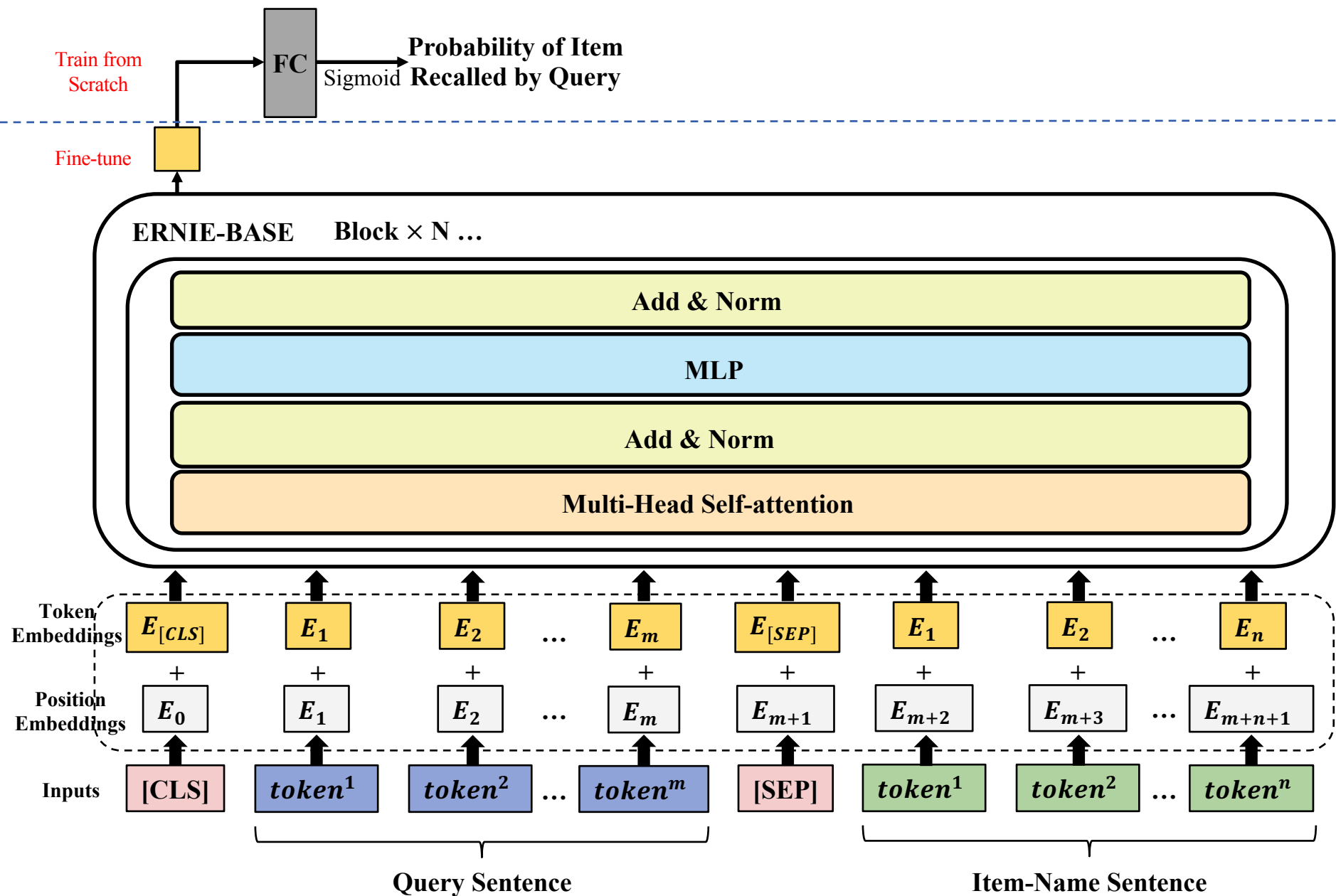# Ernie-Based Recalling

Xiaonan Wang
Search-Ads Algorithms Related Industry Project When
Working as a Machine Learning Scientist
(During Period: 2020.09~2022.07)
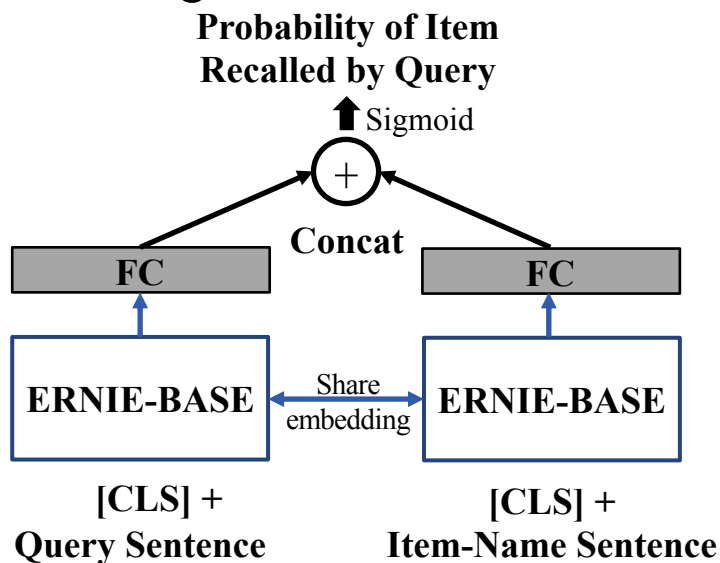
# Version 1: Ernie Single Tower Recalling

Train from Scratch

**FC** Sigmoid → **Probability of Item Recalled by Query**

Fine-tune

**ERNIE-BASE** Block × N ...

**Add & Norm**

**MLP**

**Add & Norm**

**Multi-Head Self-attention**

Token Embeddings: $E_{[CLS]}$ | $E_1$ | $E_2$ | ... | $E_m$ | $E_{[SEP]}$ | $E_1$ | $E_2$ | ... | $E_n$

+ + + + + + + +

Position Embeddings: $E_0$ | $E_1$ | $E_2$ | ... | $E_m$ | $E_{m+1}$ | $E_{m+2}$ | $E_{m+3}$ | ... | $E_{m+n+1}$

Inputs: [CLS] | $token^1$ | $token^2$ | ... | $token^m$ | [SEP] | $token^1$ | $token^2$ | ... | $token^n$

**Query Sentence**              **Item-Name Sentence**

# Version 1: Ernie Single Tower Recalling

## Summarization

- Modeling recalling problem as a <u>classification</u> problem, and select Top-N items under a certain query when recalling. -> Do like ranking does.

- Construct samples using a <u>pointwise</u> approach, and <u>cross-entropy loss</u>.

- Sample optimization: using clicks and relevance as criteria, perform positive and negative sampling on the entire recall space to <u>ensure that the distribution of the training space is consistent with that of the prediction space</u>.
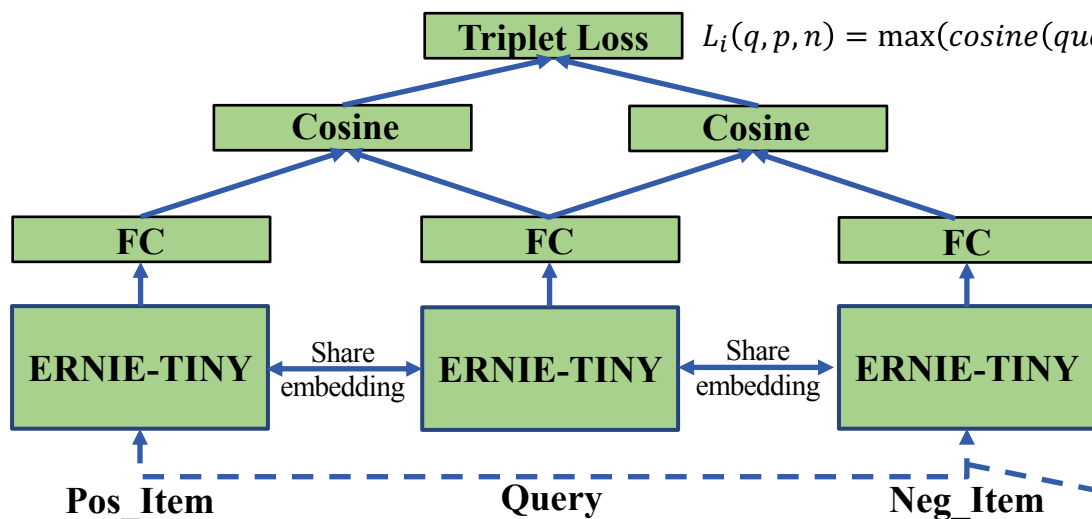
# Version 2: Ernie Two Tower Recalling

## Modeling as Classification Problem:

**Probability of Item Recalled by Query**

↑ Sigmoid

(+)

**Concat**

| FC | FC |

| ERNIE-BASE | Share embedding | ERNIE-BASE |

**[CLS] +**
**Query Sentence**

**[CLS] +**
**Item-Name Sentence**

- Select Top-N items under a certain query when recalling. -> Do like ranking does.

- <u>Pointwise</u> and <u>Cross-Entropy Loss</u>

- Sample Optimization

- Considering there has $M$ queries and $N$ items in candidate set, two-tower approach decreases $M \times N$ samples to $M + N$ samples in inference set, with the sacrifice of accuracy compared with single tower approach.
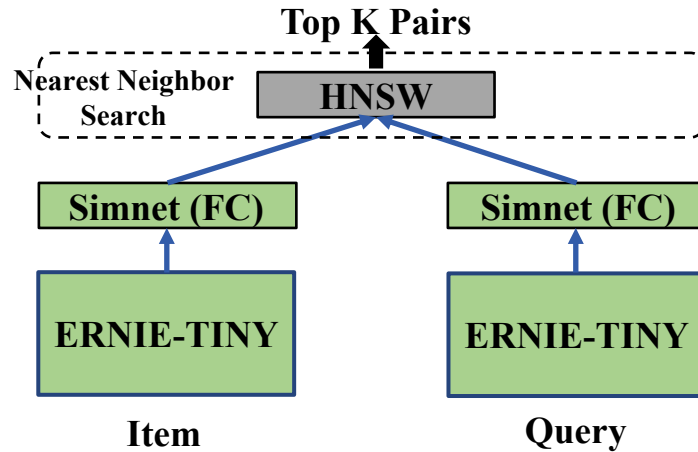
## Modeling as LTR Problem -- Training:

**Triplet Loss**   $L_i(q, p, n) = \max(cosine(query_i, pos_i) - cosine(query_i, neg_i) + \alpha, \ 0)$

| Cosine | Cosine |

| FC | FC | FC |

| ERNIE-TINY | Share embedding | ERNIE-TINY | Share embedding | ERNIE-TINY |

**Pos_Item**       **Query**       **Neg_Item**

actually use the same tower

- Change to Learning-To-Rank approach.

- <u>Pairwise</u>:
  - $ins_i = <query, pos\_item, neg\_item>$

- Considering $M$ queries and $N$ items: two-tower approach decreases $M \times N$ samples to $M + N$ samples in inference set, with the sacrifice of accuracy compared with single tower approach.

# Version 2: Ernie Two Tower Recalling

Modeling as LTR Problem -- Inference:

**Top K Pairs**

**Nearest Neighbor Search**

**HNSW**

**Simnet (FC)**  **Simnet (FC)**

**ERNIE-TINY**  **ERNIE-TINY**

**Item**  **Query**

- Faster Inference: Replace the 12-blocks ernie-base model with the extracted blocks ernie-tiny model

- Use the Learning-To-Rank idea and the pairwise paradigm to distinguish positive and negative samples at a fine-grained level, and describe the ranking relationship。