

# Identical and Fraternal Twins: Fine-Grained Semantic Contrastive Learning of Sentence Representations

Qingfa Xiao<sup>1,3</sup>, Shuangyin Li<sup>1,✉</sup> and Lei Chen<sup>2,3</sup>

<sup>1</sup>South China Normal University

<sup>2</sup>The Hong Kong University of Science and Technology

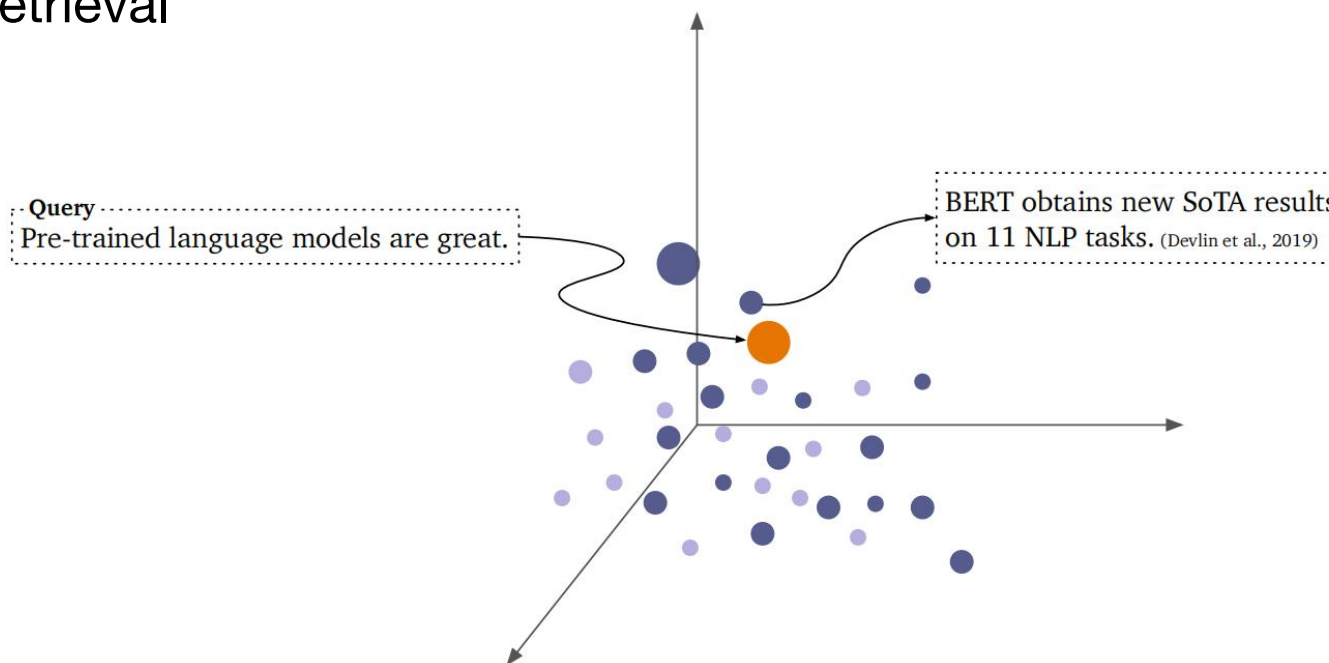
<sup>3</sup>The Hong Kong University of Science and Technology (Guangzhou)  
qingfaxiao@m.scnu.edu.cn, shuangyinli@scnu.edu.cn, leichen@cse.ust.hk



# Background: Sentence Embeddings

Learning universal representations of sentences has wide applications in NLP

- Semantic matching
- Sentence clustering
- Information retrieval
- ...

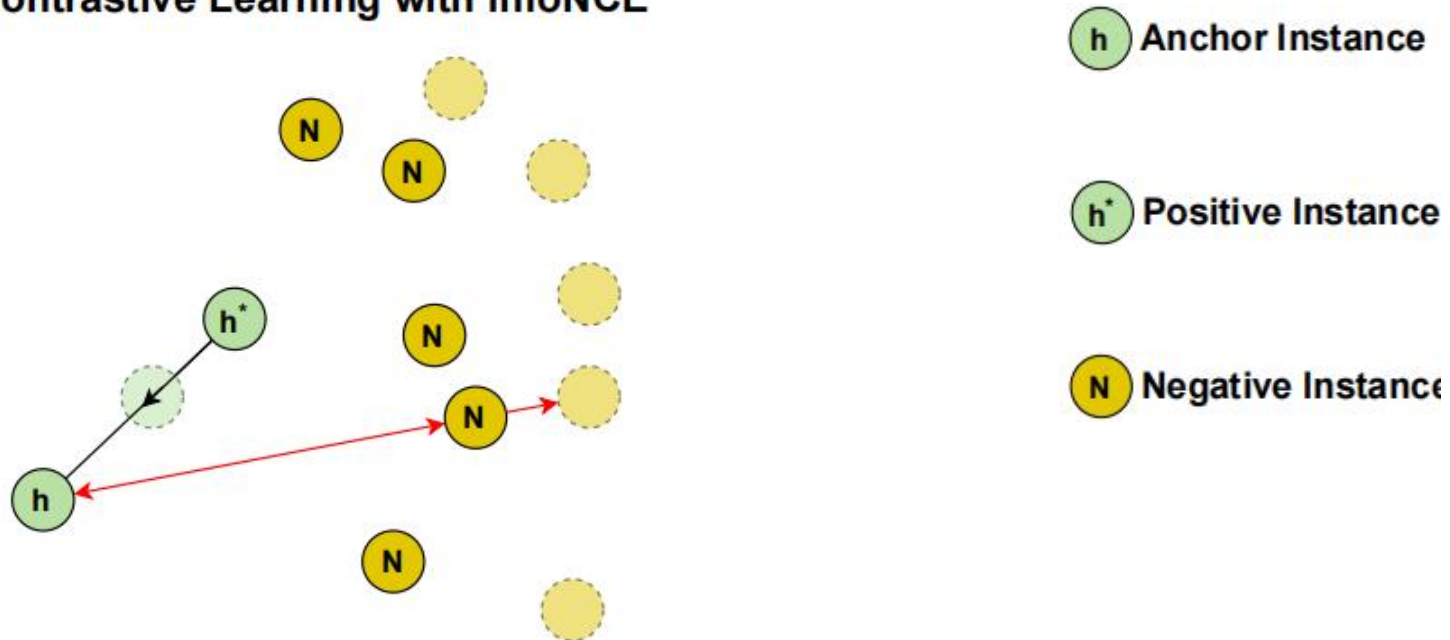


# Background: Contrastive Learning

Contrastive learning is an intuitive and effective training objective that aims to create desired semantic representations by

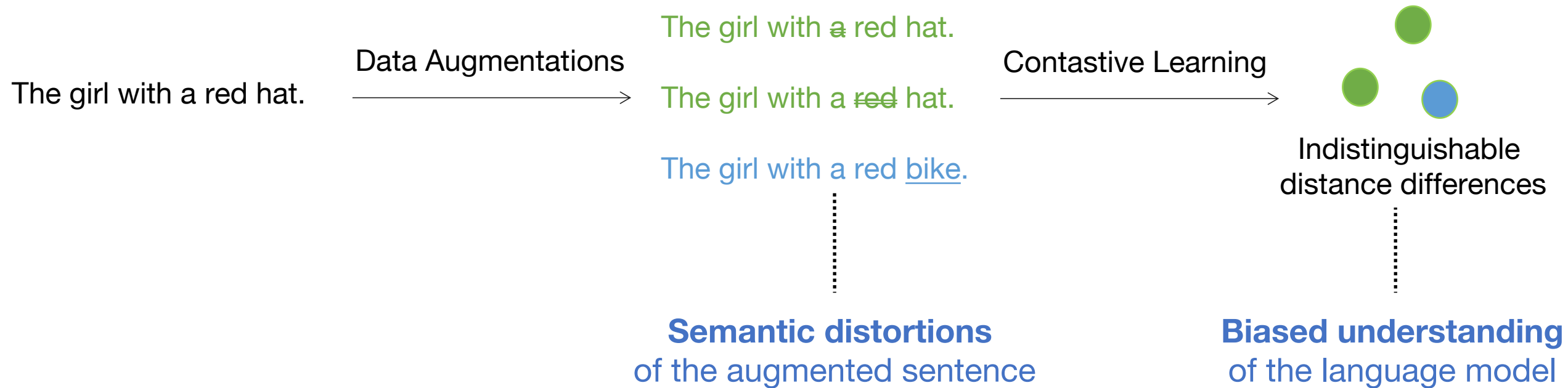
- **bringing** semantically positive instances **closer** together
- **pushing away** those that are not semantically negative.

(a) Contrastive Learning with InfoNCE



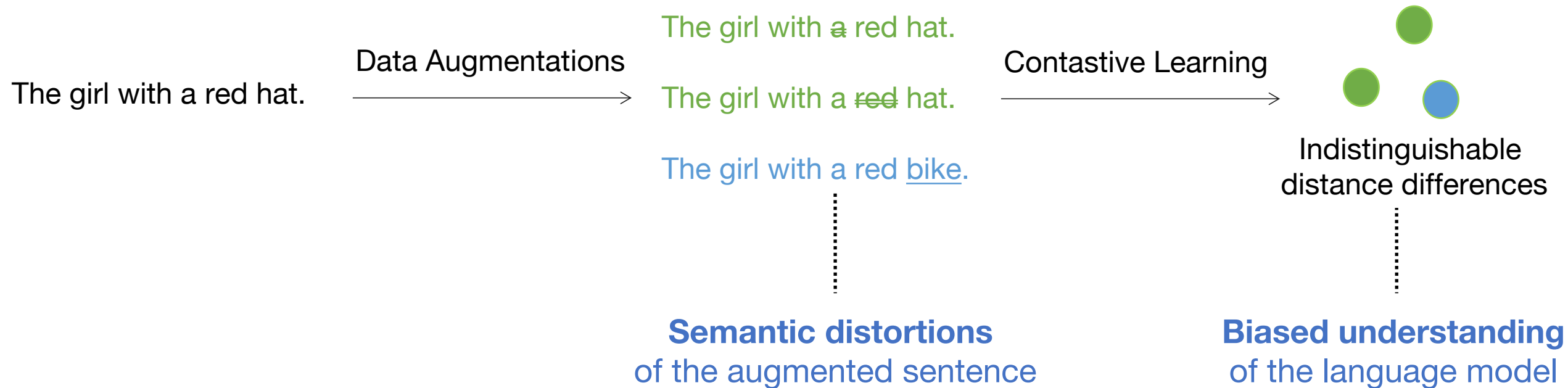
# Motivation:

**The lack of fine-grained semantic discrimination ability via contrastive learning.**



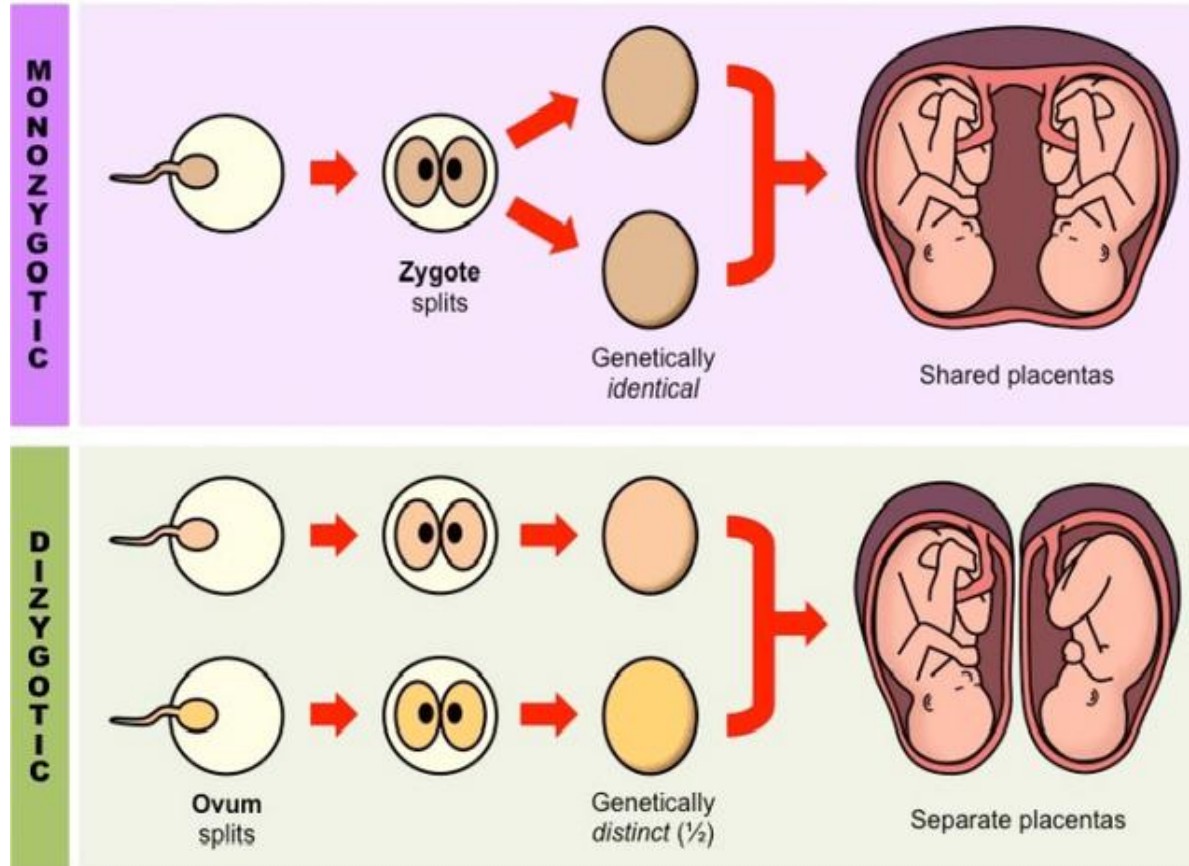
# Motivation:

**The lack of fine-grained semantic discrimination ability via contrastive learning.**



**When the different types of positive pairs come to contrastive learning, they should be treated under the different standards.**

# Motivation: Which twins are more similar?



Identical Twins



Fraternal Twins

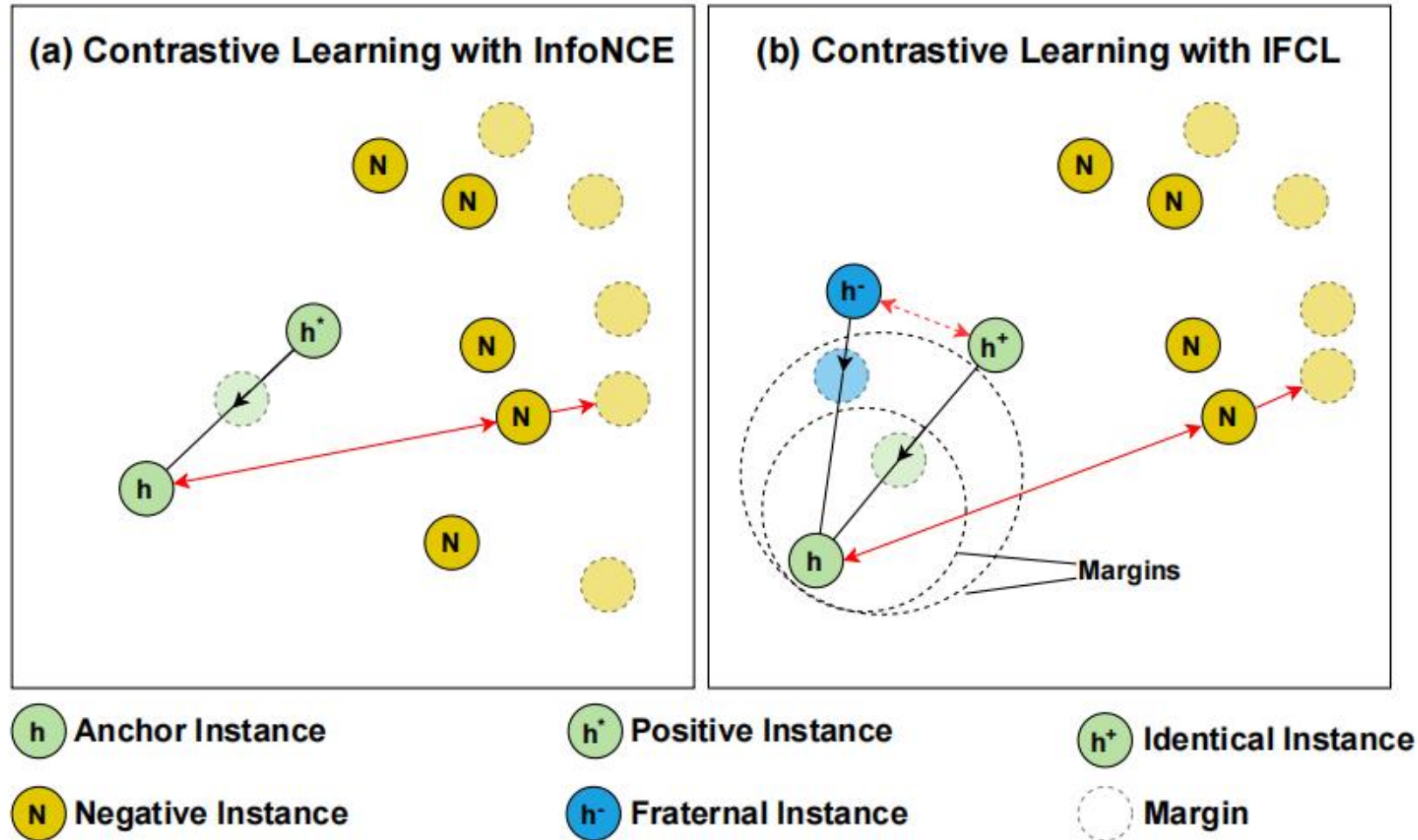


**Everyone can easily distinguish between these two pairs of twins.**

**But, can a language model do the same?**

# Objective:

Keep the margins between the two pairs of twins to help model distinguish the subtle differences.



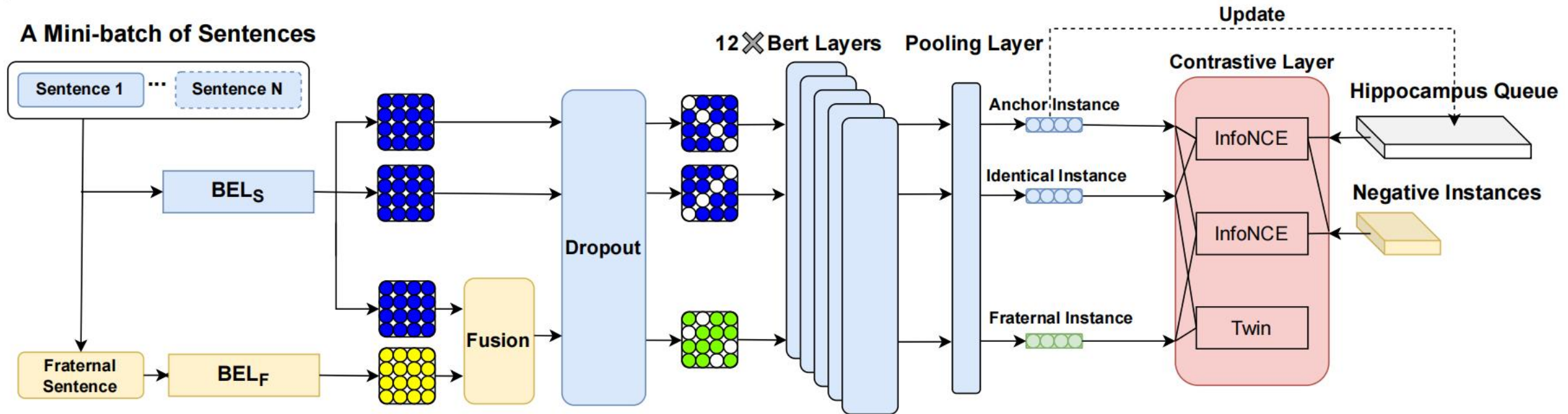
# Challenges:

## Contrastive learning with data augmentation

- Data augmentation for gaining positive pairs with less **semantic distortions**
- Adaptive contrastive learning for different types of positive samples to **address sub-optimal issues**



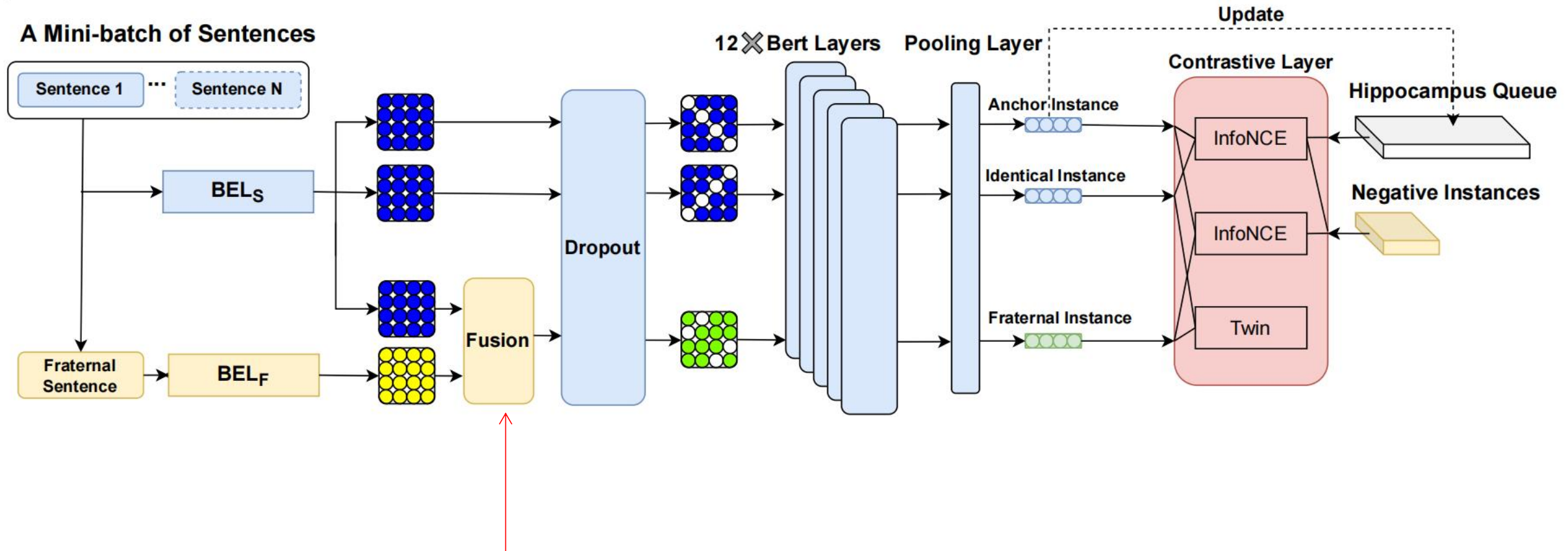
# IFCL Framework:



Three main components in this IFCL framework:

- A fusion data augmentation technique
- A training loss function named Twin Loss
- A hippocampus queue mechanism

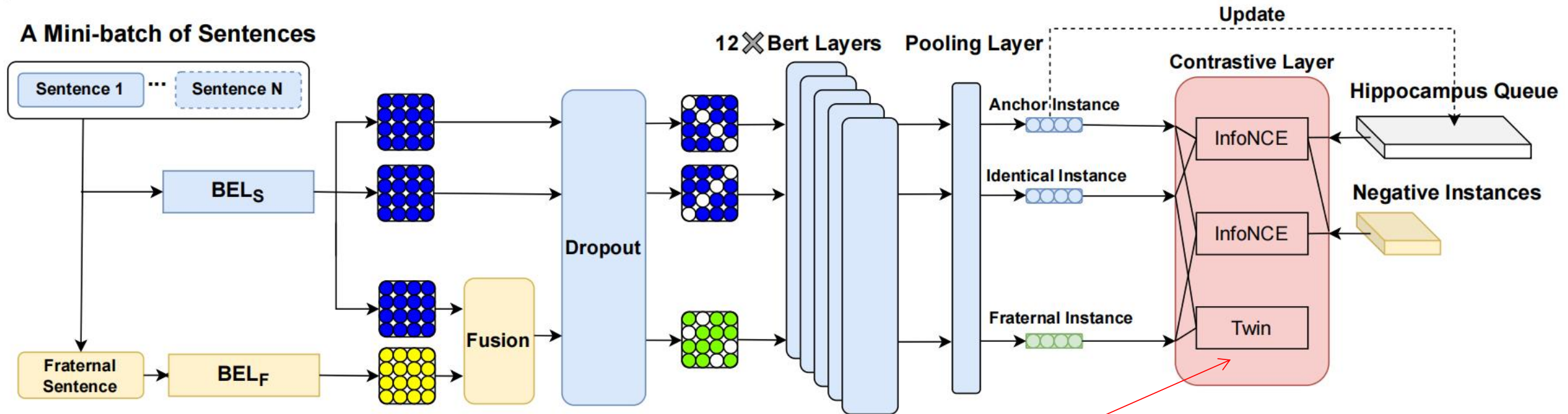
# IFCL Framework:



- A fusion data augmentation technique

minimizing semantic distortions and increasing diversity of expressions.

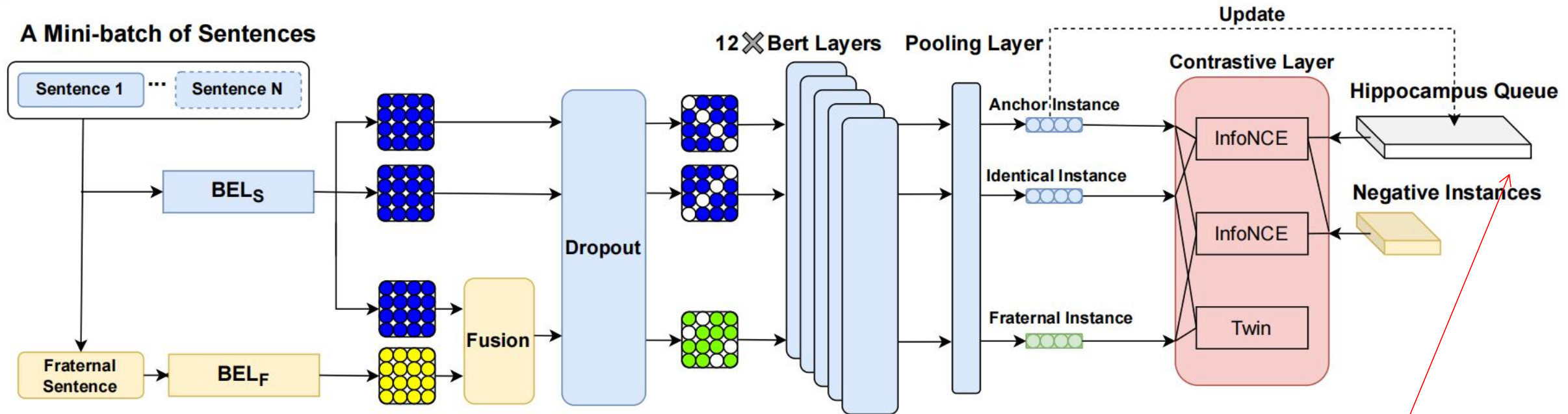
# IFCL Framework:



- A training loss function named Twin Loss

capturing **fine-grained semantics** and alleviating the sub-optimal issues according to their margins

# IFCL Framework:

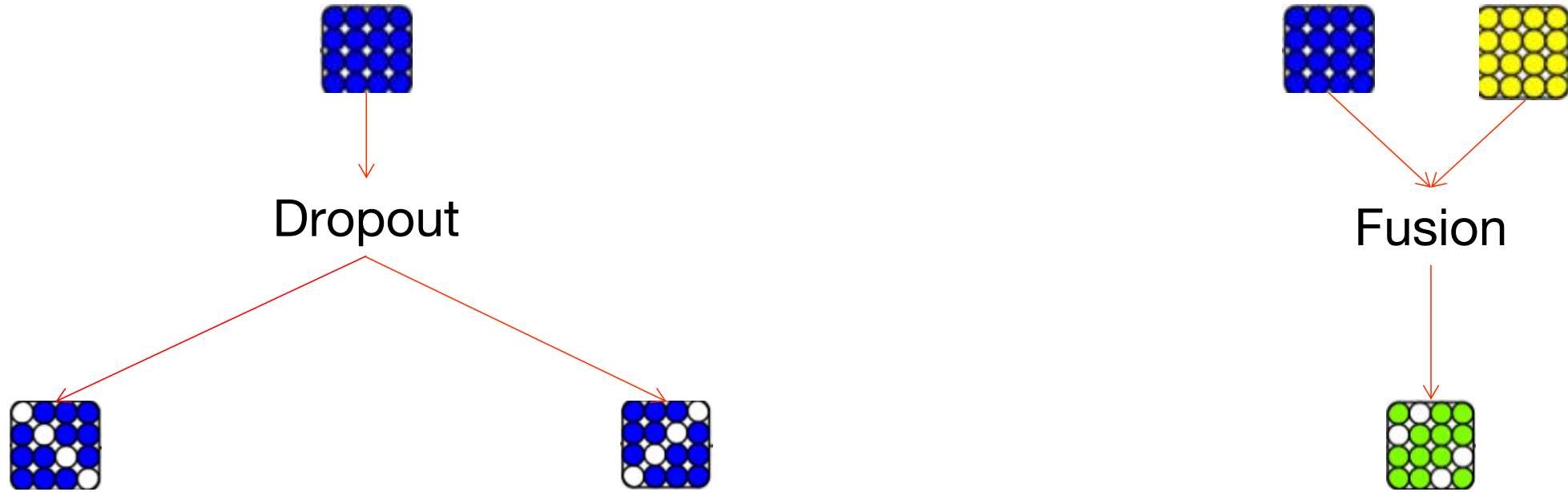


- A hippocampus queue mechanism

**storing** the previous mini-batches into a short-term memory and **reusing** the negative effectively

# Method:

## How to generate the Identical and Fraternal Twins?



**Identical twins:** the most similar positive pair

**Fraternal twins:** the diverse pair with less semantic distortions

## Method:

### InfoNCE Loss with positive and negative instances

For the set of identical twins  $\{h_i, h_i^+\}_{i=1}^N$  or fraternal twins  $\{h_i, h_i^-\}_{i=1}^N$ , we define the function by using negative instances  $\{\mathbf{H}_m\}_{m=1}^{k*N}$  stored in the hippocampus queue.

$$\ell_i^I = -\log \frac{e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_i^+)/\tau}}{\sum_{j=1}^N e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_j^+)/\tau} + \varphi},$$
$$\varphi = \sum_{m=1}^{k*N} p_m * e^{\text{sim}(\mathbf{h}_i, \mathbf{H}_m)/\tau},$$

# Method:

## Twins Loss for fine-grained semantic understanding

- This loss function aim to **keeping the margins** between two types of positive pairs
- M represents the **innate margins** between identical and fraternal twins
- Each M depends on the previous step to **prevent sub-optimal optimization problems**

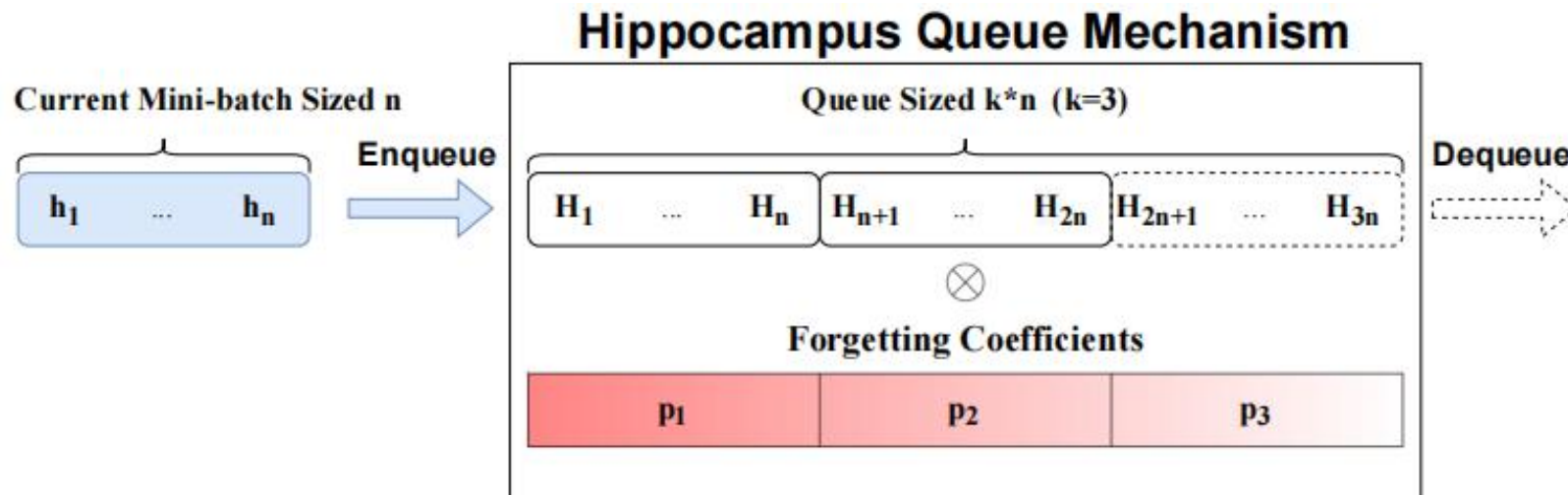
$$\ell_i^T = \left| e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_i^+)} - e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_i^-)} - \mathbf{M}_i \right|,$$

$$\mathbf{M}_i = e^{\text{sim}(\text{emb}_i, \text{emb}_i^+)} - e^{\text{sim}(\text{emb}_i, \text{emb}_i^-)}$$

# Method:

## Hippocampus Queue Mechanism for reusing instances

- The queue storing the negative is **continuously updated**
- The sample is gradient-free to **save GPU memory**
- The forgetting coefficient focus more on the **latest instance**





# Results:

## Experiments on semantic textual similarity tasks

- Evaluate using Spearman's correlation metric
- Performe well in both Chinese and English tasks.

Results of Chinese tasks		
Method	Chinese STS-B	SimCLUE
BERT	55.52	29.89
BERT-whitening <sup>•</sup>	68.27	-
SimCSE-BERT <sup>•</sup>	68.91	40.74
SimCSE-BERT <sup>◇</sup>	60.41	40.54
<b>IFCL-BERT<sup>◇</sup></b>	<b>71.41</b>	<b>44.42</b>

Results of English tasks								
Method	STS12	STS13	STS14	STS15	STS16	STS-B	SICK-R	Avg.
BERT <sub>base</sub>	39.70	59.38	49.67	66.03	66.19	53.87	62.06	56.70
BERT-flow <sub>base</sub>	58.40	67.10	60.85	75.16	71.22	68.66	64.47	66.55
BERT-whitening <sub>base</sub>	57.83	66.90	60.90	75.08	71.31	68.24	63.73	66.28
ConSERT <sub>base</sub>	64.64	78.49	69.07	79.72	75.95	73.97	67.31	72.74
SimCSE-BERT <sub>base</sub> <sup>◇</sup>	68.40	82.41	74.38	80.91	78.56	76.85	72.23	76.25
VaSCL-BERT <sub>base</sub> <sup>◇</sup>	69.08	81.95	74.64	82.64	80.57	80.23	71.23	77.19
DCLR-BERT <sub>base</sub> <sup>◇</sup>	70.81	83.73	75.11	82.56	78.44	78.31	71.59	77.22
MoCoSE-BERT <sub>base</sub> <sup>◇</sup>	71.48	81.40	74.47	<b>83.45</b>	78.99	78.68	<b>72.44</b>	77.27
PT-BERT <sub>base</sub> <sup>◇</sup>	71.20	<b>83.76</b>	<b>76.34</b>	82.63	78.90	79.42	71.94	77.74
<b>IFCL-BERT<sub>base</sub><sup>◇</sup></b>	<b>71.57</b>	82.35	75.08	83.03	<b>80.17</b>	<b>80.27</b>	72.16	<b>77.80</b>
BERT <sub>large</sub>	57.73	61.17	61.18	68.07	70.25	59.59	60.34	62.62
ConSERT <sub>large</sub>	70.69	82.96	74.13	82.78	76.66	77.53	70.37	76.45
SimCSE <sub>large</sub> <sup>◇</sup>	70.88	84.16	76.43	84.50	79.76	79.26	73.88	78.41
DCLR-BERT <sub>large</sub> <sup>◇</sup>	71.87	<b>84.83</b>	<b>77.37</b>	<b>84.70</b>	<b>79.81</b>	79.55	74.19	78.90
MoCoSE-BERT <sub>large</sub> <sup>◇</sup>	<b>74.50</b>	84.54	77.32	84.11	79.67	80.53	73.26	79.13
<b>IFCL-BERT<sub>large</sub><sup>◇</sup></b>	73.88	84.31	76.64	84.01	79.56	<b>81.37</b>	<b>76.30</b>	<b>79.44</b>

# Analyse:

## What makes the Twins Loss effective?

Reducing the mutual information between positive pairs while preserving task-relevant information is optimal for the task

- More diverse semantics are preserved

$\text{MII}(h, h^-)$  is higher than  $\text{MII}(h, h^+)$

- Mutual information of positive pairs contains more task-relevant information.

$\text{MII}(h, h^+) \approx \text{MII}(h, h^-) \approx \text{MII}_{task}$

**Table 4.** Mutual information and task-relevant information. The IFCL-BERT w/o TL means training IFCL-BERT without Twins Loss. The experiments are conducted with EnData and STS-B datasets on Bert-base.

Method	$\text{MII}(h, h^+)$	$\text{MII}(h, h^-)$	$\text{MII}_{task}$
IFCL-BERT	4.15	4.17	4.31
IFCL-BERT w/o TL	4.23	4.20	4.58
SimCSE	4.24	-	4.52

**Thank YOU**

**Q & A**