# IMNTPU at the NTCIR-16 FinNum-3 Task: Data Augmentation for Financial Numclaim Classification

**Yung-Wei Teng [1], Pei-Tz Chiu [1], Ting-Yun Hsiao [1], Mike Tian-Jian Jiang [2] and Min-Yuh Day [1,*]**
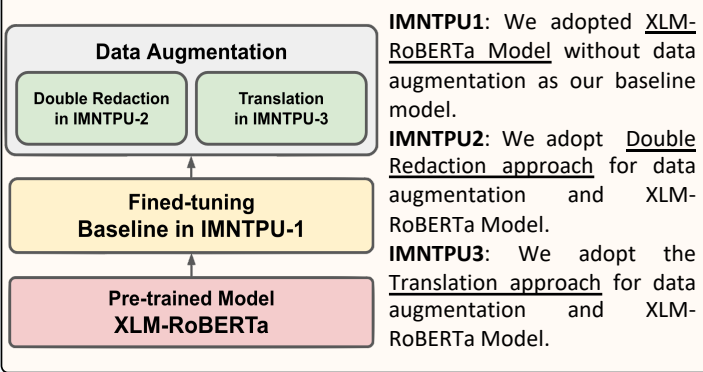
**[1] Information Management, National Taipei University, New Taipei City, Taiwan**
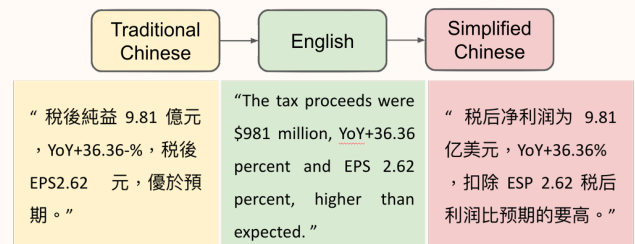**[2] Zeals Co., Ltd. Tokyo, Japan**
*myday@gm.ntpu.edu.tw

This paper provides a detailed description of IMNTPU team at the NTCIR-16 FinNum-3 shared task in formal financial documents. We proposed the use of the XLM-RoBERTa-based model with two different approaches on data augmentation to perform the binary classification task in FinNum-3. The first run (i.e., IMNTPU-1) is our baseline through the fine-tuning of the XLM-RoBERTa without data augmentation. However, we assume that presenting different data augmentations may improve the task performance because of the imbalance in the dataset. Accordingly, we presented double redaction and translation method on data augmentation in the second (IMNTPU-2) and third (IMNTPU-3) runs, respectively. The best macro-F1 scores obtained by our team in the Chinese and English datasets are 93.18% and 89.86%, respectively. The major contribution in this study provide a new understanding toward data augmentation approach for the imbalanced dataset, which may help reduce the imbalanced situation in the Chinese and English datasets.

## Research Architecture and Proposed Method

**Data Augmentation**
- Double Redaction in IMNTPU-2
- Translation in IMNTPU-3

**Fined-tuning Baseline in IMNTPU-1**

**Pre-trained Model XLM-RoBERTa**

**IMNTPU1**: We adopted XLM-RoBERTa Model without data augmentation as our baseline model.

**IMNTPU2**: We adopt Double Redaction approach for data augmentation and XLM-RoBERTa Model.

**IMNTPU3**: We adopt the Translation approach for data augmentation and XLM-RoBERTa Model.

## Tokenization Tricks

**Input:** Good day and welcome to the Apple Inc. Third Quarter Fiscal Year **2018** Earnings Conference Call. Today's call is being recorded.

**XLM-RoBERTa Tokenizer**

**Output:** <s> Good day and welcome to the Apple Inc. Third Quarter Fiscal Year xxnum **2018** Earnings Conference Call. Today's call is being recorded. </s>

**Double Redaction**

**Output:** <s> <mask> Good day and <mask> to the Apple <mask> Third Quarter Fiscal Year xxnum **2018** Earnings Conference Call. Today's call is <mask> recorded. </s>

## Algorithm of Double Redaction

```
1:  Shuffle the tokens in sentence
2:  Delete the duplicated tokens in sentence
3:  Copy the remaining tokens as β
4:  SET the δ and γ
5:  for specific token in β do
6:      if γ less than δ then
7:          Replace original token with <usk> token
8:      else
9:          Cover original token as <mask> token
10:     end if
11: end for
12: while True do
13:     Model predict the original token of <usk> and <mask>
14: end while
```

## Translation Approach

Traditional Chinese → English → Simplified Chinese

" 稅後純益 9.81 億元，YoY+36.36-%，稅後 EPS2.62 元，優於預期。"

"The tax proceeds were $981 million, YoY+36.36 percent and EPS 2.62 percent, higher than expected. "

" 税后净利润为 9.81 亿美元，YoY+36.36%，扣除 ESP 2.62 税后利润比预期的要高。"

## Performance

| Run | Chinese Dataset | | English Dataset | |
|---|---|---|---|---|
| | Dev Set F1-Score (%) | Test Set F1-Score (%) | Dev Set F1-Score (%) | Test Set F1-Score (%) |
| **IMNTPU1** | 90.51 | **93.18** | 87.13 | 88.39 |
| **IMNTPU2** | 88.65 | 91.64 | 88.82 | **89.86** |
| **IMNTPU3** | **92.16** | 91.64 | - | - |

## Conclusions and Contribution

**Conclusions:**

The performance with data augmentation method (Double Redaction) in English dataset is superior than without data augmentation.

**Contribution:**

- The major contribution of the research is that data augmentation approach may help reduce imbalanced situation.

- We have developed a novel method for data augmentation technique, which is double redaction and translation approach, and can decrease the issue of imbalanced dataset.