

Improving Numeracy

by Input Reframing and Quantitative Pre-Finetuning Task

Chung-Chi Chen,¹ Hiroya Takamura,¹ Ichiro Kobayashi,² Yusuke Miyao³

¹AIRC, AIST, Japan

²Ochanomizu University, Japan

³University of Tokyo, Japan

EACL 2023

Input Reframing

| Model | Notation | Tokenized Example |
|---------|---------------------|--|
| BERT | Org. | "147", "##70", "##2" |
| | Digit | "1", "4", "7", "7", "0", "2" |
| | Scientific Notation | "1", ".", "47", "##70", "##200", "##00", "##0", "##e", "+", "05" |
| RoBERTa | Org. | "147", "702" |
| | Digit | "1", "4", "7", "7", "0", "2" |
| | Scientific Notation | "1", ".", "47", "70", "200000", "E", "+", "05" |

Quantitative Pre-Finetuning Task

| Task | Question | Answer |
|--------|--|------------|
| ComNum | [Num 1] is equal to [Num 2]. | TRUE/FALSE |
| | [Num 1] is smaller than [Num 2]. | |
| | [Num 1] is larger than [Num 2]. | |
| QP | FED'S DUDLEY REPEATS EXPECTS GDP GROWTH TO PICK UP IN 2014, FROM [Masked] PCT POST-RECESSION AVERAGE | 1 |
| QNLI | S1: Nifty traded above 7500, Trading Calls Today S2: Nifty above 7400 | Entailment |
| QQA | Elliot weighs 180 pounds whereas Leon weighs 120 pounds. Who has a bigger gravity pull? Option1: Elliot Option2: Leon | Option 1 |

Innumeracy → A Phenomenon with All LMs

| Training: 0 - 199,999 | Testing: | BERT | | RoBERTa | | LinkBERT | | FinBERT | |
|-----------------------|---------------------|--------|----------------|---------|-----------------|----------|----------------|---------|-----------------|
| | | CND-T1 | CND-T2 | CND-T1 | CND-T2 | CND-T1 | CND-T2 | CND-T1 | CND-T2 |
| • CND-T1: 0 - 199,999 | Original | 99.86 | 95.59 (↓ 4.27) | 99.44 | 86.75 (↓ 12.69) | 99.92 | 97.58 (↓ 2.34) | 99.55 | 78.37 (↓ 21.18) |
| • CND-T2: 4M - 5M | Digit-based | 99.96 | 99.03 (↓ 0.93) | 99.92 | 98.46 (↓ 1.46) | 99.99 | 96.54 (↓ 3.45) | 99.96 | 97.03 (↓ 2.93) |
| | Scientific Notation | 99.92 | 99.68 (↓ 0.24) | 99.82 | 99.13 (↓ 0.69) | 99.95 | 99.81 (↓ 0.14) | 99.72 | 98.78 (↓ 0.94) |

RoBERTa → Both Input Reframing and Quantitative PFT Works

| Model | Notation | QP | | QNLI | | | | | QQA | Score |
|------------|---------------------|---------------|---------------|---------------|---------------|---------|---------------|----------------|---------------|--------------|
| | | Comment | Headline | RTE-QUANT | AWP-NLI | NEWSNLI | REDDITNLI | Stress Test | | |
| RoBERTa | Original | 60.46% | 58.03% | 60.15% | 57.64% | 79.58% | 58.77% | 98.93% | 51.96% | 65.69 |
| | Digit-based | 69.25% | 57.65% | 59.40% | 56.69% | 78.90% | 62.38% | 99.91% | 54.34% | 67.31 |
| | Scientific Notation | 64.32% | 55.49% | 60.08% | 57.41% | 78.68% | 60.81% | 100.00% | 53.67% | 66.31 |
| CN-RoBERTa | Original | 86.86% | 77.29% | 62.52% | 56.70% | 78.82% | 64.29% | 99.94% | 50.71% | 72.14 |
| | Digit-based | 64.25% | 55.92% | 68.96% | 58.80% | 77.99% | 60.99% | 99.73% | 50.88% | 67.19 |
| | Scientific Notation | 60.28% | 54.85% | 62.15% | 58.74% | 65.92% | 59.59% | 99.47% | 52.27% | 64.16 |

BERT-Based → Quantitative PFT Works with QNLI

| Model | Notation | QP | | QNLI | | | | | QQA | Score |
|-------------|---------------------|---------|---------------|---------------|---------------|---------------|---------------|----------------|---------------|--------------|
| | | Comment | Headline | RTE-QUANT | AWP-NLI | NEWSNLI | REDDITNLI | Stress Test | | |
| BERT | Original | 70.44% | 57.46% | 64.40% | 59.20% | 72.29% | 60.42% | 99.91% | 53.20% | 67.17 |
| | Digit-based | 65.38% | 54.74% | 57.86% | 56.46% | 71.36% | 60.11% | 99.11% | 53.75% | 64.85 |
| | Scientific Notation | 65.31% | 55.99% | 64.42% | 60.73% | 72.23% | 59.66% | 99.56% | 53.24% | 66.39 |
| CN-BERT | Digit-based | 69.93% | 54.84% | 61.07% | 60.27% | 75.54% | 65.39% | 99.42% | 52.53% | 67.37 |
| | Scientific Notation | 64.87% | 56.40% | 66.39% | 54.70% | 75.41% | 63.94% | 99.42% | 51.90% | 66.63 |
| LinkBERT | Original | 68.81% | 55.70% | 59.94% | 56.85% | 73.43% | 59.01% | 99.91% | 54.14% | 65.97 |
| | Digit-based | 63.76% | 55.41% | 59.54% | 57.42% | 73.63% | 60.17% | 99.73% | 53.44% | 65.39 |
| | Scientific Notation | 65.81% | 56.05% | 57.00% | 56.78% | 75.51% | 58.51% | 99.82% | 54.33% | 65.48 |
| CN-LinkBERT | Digit-based | 68.61% | 54.44% | 63.59% | 55.08% | 71.21% | 58.99% | 100.00% | 50.44% | 65.30 |
| | Scientific Notation | 63.48% | 53.15% | 62.02% | 59.39% | 75.70% | 62.61% | 99.73% | 52.11% | 66.02 |

FinBERT → Both Works with QP

| Model | Reframing | QP-Comment |
|------------|---------------------|---------------|
| FinBERT | Original | 65.26% |
| | Digit-based | 69.89% |
| | Scientific Notation | 70.03% |
| CN-FinBERT | Digit-based | 68.84% |
| | Scientific Notation | 69.76% |

Findings

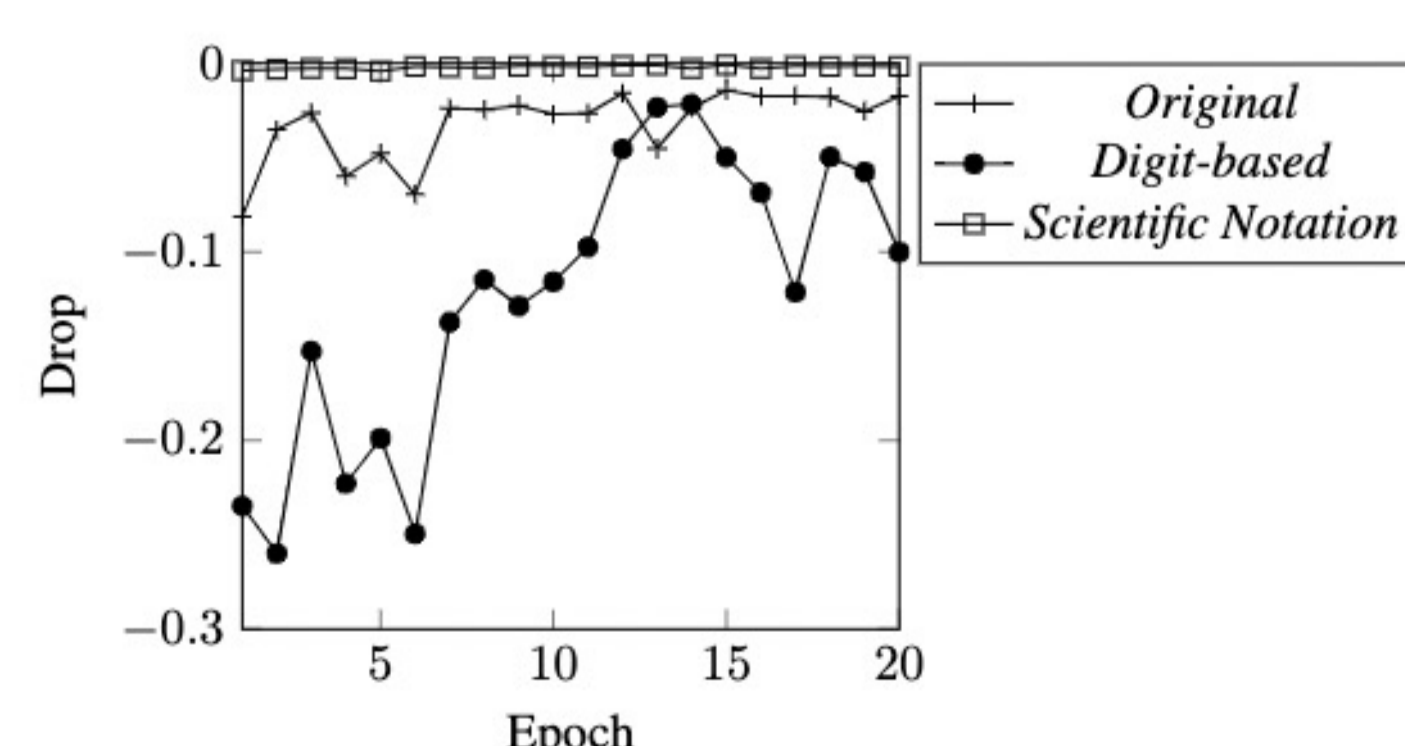
Innumeracy: When the numerals in the test set are **out of the range** of those in the training set, the performance will **drop**

| RoBERTa | QP | QNLI | QQA |
|--------------|----|------|-----|
| Reframing | V | V | V |
| ComNum (PFT) | V | V | X |

| BERT | QP | QNLI | QQA |
|--------------|----|------|-----|
| Reframing | X | V | X |
| ComNum (PFT) | X | V | X |

Future Direction

Convergence



Numeral Factuality:
Numeracy-600K (ACL-2019) – Detecting Exaggerated Information

Quantitative 101 Dataset



<http://q101.nlpfin.com/>

Comparing Number LMs



<http://cn-roberta.nlpfin.com/>