



# Data Augmentation for NMT

~Nalin Kumar  
SMS, NISER



# Data Augmentation

- NMT models need huge amount of training data.
- We need a transformation  $f$ , such that for input 'x', the augmented data  $f(x)$ , should preserve the label and x and  $f(x)$  should be diverse enough from each other.



- Generating augmented data in NLP is relatively tougher than data augmentation in image classification tasks.
- Some augmentation techniques:
  - Using thesaurus<sup>[1]</sup> (Replaced words with their synonym)
  - Using word embeddings<sup>[2]</sup> (Replaced words with their neighbouring words)
  - Back-translation<sup>[3]</sup>
  - Data diversification<sup>[4]</sup>
  - Cutoff<sup>[5]</sup>

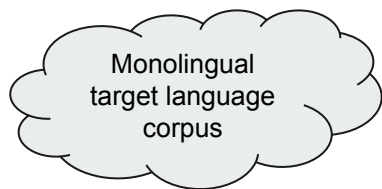
} Used in classification tasks

1. Zhang et al., “Character-level Convolutional Networks for Text Classification” (2015, NeurIPS)
2. Wang et al., “That’s So Annoying!!!: A Lexical and Frame-Semantic Embedding Based Data Augmentation Approach to Automatic Categorization of Annoying Behaviors using #petpeeve Tweets” (2015, EMNLP)
3. Sennrich et al., “Improving Neural Machine Translation Models with Monolingual Data” (2016, ACL)
4. Nguyen et al., “Data Diversification: A Simple Strategy For Neural Machine Translation” (2020, NeurIPS))
5. Shen et al., “A Simple but Tough-to-Beat Data Augmentation Approach for Natural Language Understanding and Generation”

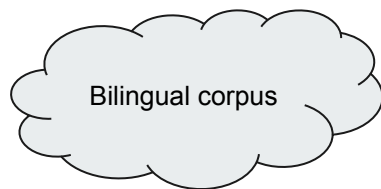


# Back-Translation

- The technique was first proposed in “Improving Neural Machine Translation Models with Monolingual Data” by Rico Sennrich, Barry Haddow and Alexandra Birch (ACL 2016)
- Uses monolingual data (which is considered to be easily available in comparison to bilingual data).



=  $T'$



=  $D = (S, T)$

1. Train  $M_{T \rightarrow S}$  on “Target -> Source” bilingual corpus
2. Train  $M_{S \rightarrow T}$  on  $(M_{T \rightarrow S}(T'), T')$   $\cup D$ .

$M_{T \rightarrow S}(T')$  are referred as synthetic source sentences.



- In this paper, the NMT models are chosen to be RNN enc-dec models.
- They performed two sets of experiments:
  - In the first set, they generated data by pairing the monolingual data with single word <null> token.
  - In the second experiment, they generated data by previous algorithm.



## Results:

| name                                     | training instances                 | BLEU         |             |              |             |
|--|------------------------------------|--------------|-------------|--------------|-------------|
|  |                                    | newstest2014 |             | newstest2015 |             |
|  |                                    | single       | ens-4       | single       | ens-4       |
| syntax-based (Sennrich and Haddow, 2015) |                                    | 22.6         | -           | 24.4         | -           |
| Neural MT (Jean et al., 2015b)           |                                    | -            | -           | 22.4         | -           |
| parallel                                 | 37m (parallel)                     | 19.9         | 20.4        | 22.8         | 23.6        |
| +monolingual                             | 49m (parallel) / 49m (monolingual) | 20.4         | 21.4        | 23.2         | 24.6        |
| +synthetic                               | 44m (parallel) / 36m (synthetic)   | <b>22.7</b>  | <b>23.8</b> | <b>25.7</b>  | <b>26.5</b> |

Table 3: English→German translation performance (BLEU) on WMT training/test sets. Ens-4: ensemble of 4 models. Number of training instances varies due to differences in training time and speed.



# Data Diversification

- This technique was introduced in “Data Diversification: A Simple Strategy For Neural Machine Translation” by Xuan-Phi Nguyen , Shafiq Joty , Wu Kui and Ai Ti Aw (NIPS 2020).
- This technique does not involve any extra monolingual data.





---

**Algorithm 1** Data Diversification: Given a dataset  $\mathcal{D} = (S, T)$ , a diversification factor  $k$ , the number of rounds  $N$ ; return a trained source-target translation model  $\hat{M}_{S \rightarrow T}$ .


---

```
1: procedure TRAIN( $\mathcal{D} = (S, T)$ )
2:   Train randomly initialized  $M$  on  $\mathcal{D} = (S, T)$  until convergence
3:   return  $M$ 

1: procedure DATADIVERSE( $\mathcal{D} = (S, T), k, N$ )
2:    $\mathcal{D}_0 \leftarrow \mathcal{D}$  ▷ Assign original dataset to round-0 dataset.
3:   for  $r \in 1, \dots, N$  do
4:      $\mathcal{D}_r = (S_r, T_r) \leftarrow \mathcal{D}_{r-1}$ 
5:     for  $i \in 1, \dots, k$  do
6:        $M_{S \rightarrow T, r}^i \leftarrow \text{TRAIN}(\mathcal{D}_{r-1} = (S_{r-1}, T_{r-1}))$  ▷ Train forward model
7:        $M_{T \rightarrow S, r}^i \leftarrow \text{TRAIN}(\mathcal{D}'_{r-1} = (T_{r-1}, S_{r-1}))$  ▷ Train backward model
8:        $\mathcal{D}_r \leftarrow \mathcal{D}_r \cup (S, M_{S \rightarrow T, r}^i(S))$  ▷ Add forward data
9:        $\mathcal{D}_r \leftarrow \mathcal{D}_r \cup (M_{T \rightarrow S, r}^i(T), T)$  ▷ Add backward data
10:   $\hat{M}_{S \rightarrow T} \leftarrow \text{Train}(\mathcal{D}_N)$  ▷ Train the final model
11:  return  $\hat{M}_{S \rightarrow T}$ 
```

---

$$\mathcal{D}_1 = (S, T) \bigcup_{i=1}^k (S, M_{S \rightarrow T, 1}^i(S)) \bigcup_{i=1}^k (M_{T \rightarrow S, 1}^i(T), T)$$



$N = 1, K = 3$


- We have  $M^1_{S \rightarrow T, 1}, M^2_{S \rightarrow T, 1}, M^3_{S \rightarrow T, 1}$  forward models and  $M^1_{T \rightarrow S, 1}, M^2_{T \rightarrow S, 1}, M^3_{T \rightarrow S, 1}$  backward models all trained on  $D_0 = (S, T)$ .
- $D_1$  is created by adding the outputs of above model along with the inputs to  $D_0$ .
- The final model is then trained on  $D_1$ .

Table 2: BLEU scores on newstest2014 for WMT’14 English-German (En-De) and English-French (En-Fr) translation tasks. Distill (T>S) (resp. T=S) indicates the teacher model is larger than (resp. equal to) the student model.

| Method                               | WMT’14      |                   |
|--------------------------------------|-------------|-------------------|
|                                      | En-De       | En-Fr             |
| Transformer [28] <sup>†</sup>        | 28.4        | 41.8              |
| Trans+Rel. Pos [23] <sup>†</sup>     | 29.2        | 41.5              |
| Scale Transformer [18]               | 29.3        | 42.7 <sup>6</sup> |
| Dynamic Conv [32] <sup>†</sup>       | 29.7        | 43.2              |
| <b>Transformer with</b>              |             |                   |
| Multi-Agent [31] <sup>†</sup>        | 30.0        | -                 |
| Distill (T>S) [14]                   | 27.6        | 38.6              |
| Distill (T=S) [14]                   | 28.4        | 42.1              |
| Ens-Distill [7]                      | 28.9        | 42.5              |
| <b>Our Data Diversification with</b> |             |                   |
| Scale Transformer [18]               | <b>30.7</b> | <b>43.7</b>       |

Table 3: BLEU scores on IWSLT’14 English-German (En-De), German-English (De-En), and IWSLT’13 English-French (En-Fr) and French-English (Fr-En) translation tasks. Superscript <sup>†</sup> denotes the numbers are reported from the paper, others are based on our runs.

| Method                               | IWSLT’14    |             | IWSLT’13    |             |
|--------------------------------------|-------------|-------------|-------------|-------------|
|                                      | En-De       | De-En       | En-Fr       | Fr-En       |
| <b>Baselines</b>                     |             |             |             |             |
| Transformer                          | 28.6        | 34.7        | 44.0        | 43.3        |
| Dynamic Conv                         | 28.7        | 35.0        | 43.8        | 43.5        |
| <b>Transformer with</b>              |             |             |             |             |
| Multi-Agent <sup>†</sup>             | 28.9        | 34.7        | -           | -           |
| Distill (T>S)                        | 28.0        | 33.6        | 43.4        | 42.9        |
| Distill (T=S)                        | 28.5        | 34.1        | 44.1        | 43.4        |
| Ens-Distill                          | 28.8        | 34.7        | 44.3        | 43.9        |
| <b>Our Data Diversification with</b> |             |             |             |             |
| Transformer                          | <b>30.6</b> | 37.0        | <b>45.5</b> | <b>45.0</b> |
| Dynamic Conv                         | <b>30.6</b> | <b>37.2</b> | 45.2        | 44.9        |

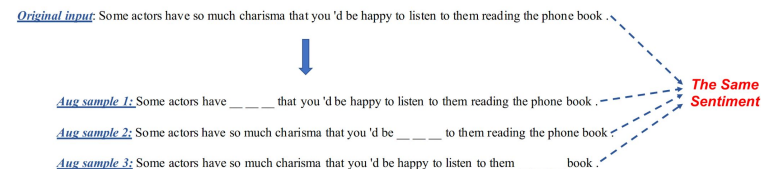
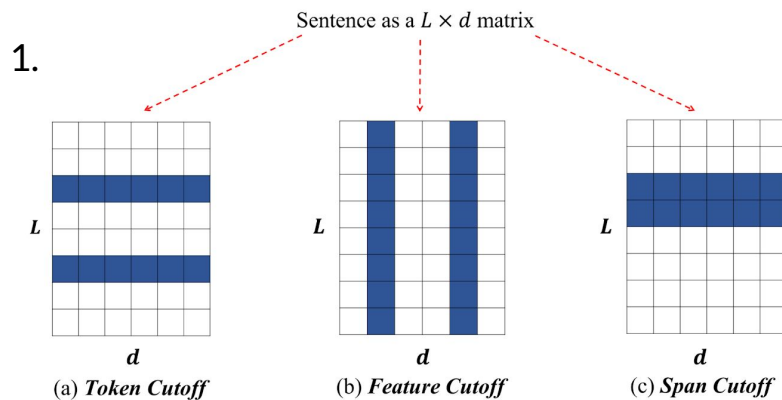
- 
- The technique exhibits strong correlation with ensemble of models. However, the ensemble model would take N times more computation than one individual model while performing inference, whereas the data diversification technique does not have this disadvantage.
  - The method is complementary to back-translation.

| Task            | No back-translation |      | With back-translation |          |             |
|-----------------|---------------------|------|-----------------------|----------|-------------|
|                 | Baseline            | Ours | $ D $                 | Baseline | Ours        |
| IWSLT' 14 En-De | 28.6                | 30.6 | 29×                   | 30.0     | <b>31.8</b> |
| IWSLT' 14 De-En | 34.7                | 37.0 | 29×                   | 37.1     | <b>38.5</b> |
| WMT' 14 En-De   | 29.3                | 30.7 | 2.4×                  | 30.8     | <b>31.8</b> |



## Cutoff

- Introduced in “A Simple but Tough-to-Beat Data Augmentation Approach for Natural Language Understanding and Generation” by Dinghan Shen, Mingzhi Zheng, Yelong Shen, Yanru Qu, Weizhu Chen (Oct 2020).
- They propose a general technique, which can be used in various Natural language Understanding and Natural Language Generation tasks.



- In the experiments, they found that token cutoff performs the best on machine translation tasks.
- This may be attributed to the fact that re-moving spans from both the source and target sentences would result in large information mismatch between the input and output, and thus the resulting pairs may be too challenging.

2.


$$\mathcal{L} = \mathcal{L}_{\text{ce}}(x, y) + \alpha \sum_{i=1}^N \mathcal{L}_{\text{ce}}(x_{\text{cutoff}}^i, y) \\ + \beta \mathcal{L}_{\text{divergence}}(x, x_{\text{cutoff}}^1, x_{\text{cutoff}}^2, \dots, x_{\text{cutoff}}^N, y)$$

$$p_{\text{avg}} = \frac{1}{N+1} \sum_{i=0}^N p(y|x_{\text{cutoff}}^i)$$

$$\mathcal{L}_{\text{divergence}} = \frac{1}{N+1} \sum_{i=0}^N \text{KL}[p(y|x_{\text{cutoff}}^i) || p_{\text{avg}}]$$

,

$$\text{KL}(P||Q) = \sum p_i(x) \log\left(\frac{p_i(x)}{q_i(x)}\right)$$



| Model                                      | BLEU score  |
|--|-------------|
| Actor-critic (Bahdanau et al., 2016)       | 28.5        |
| Transformer Base (Vaswani et al., 2017)    | 34.4        |
| Adversarial training (Wang et al., 2019)   | 35.2        |
| Data Diversification (Nguyen et al., 2019) | 37.2        |
| MAT (Fan et al., 2020)                     | 36.2        |
| Mixed Representations (Wu et al., 2020)    | 36.4        |
| MAT+Knee (Iyer et al., 2020)               | 36.6        |
| Transformer Base & Cutoff (w/o JS loss)    | 36.7        |
| Transformer Base & Cutoff (w/ JS loss)     | <b>37.6</b> |

Table 3: BLEU scores of the proposed cutoff method on the IWSLT2014 German-to-English machine translation task, relative to adversarial-based baseline and other state-of-the-art models.

| Model   | BLEU score  |
|---|-------------|
| Transformer Base (Vaswani et al., 2017)         | 27.3        |
| Admin (Liu et al., 2020a)                       | 27.9        |
| Transformer Base <sup>1</sup> (So et al., 2019) | 28.2        |
| Evolved Transformer (So et al., 2019)           | 28.4        |
| Weighted Transformer (Ahmed et al., 2017)       | 28.4        |
| Adversarial Training (Wang et al., 2019)        | 28.4        |
| Transformer Base & Cutoff (w/o JS loss)         | 28.9        |
| Transformer Base & Cutoff (w/ JS loss)          | <b>29.1</b> |

Table 2: BLEU scores of the proposed cutoff method on the WMT2014 English-to-German machine translation task, compared with adversarial-based baselines. All methods are built on top of 6-layer Transformer Base model (Vaswani et al., 2017).