# Sampling and Filtering of Neural Machine Translation Distillation Data
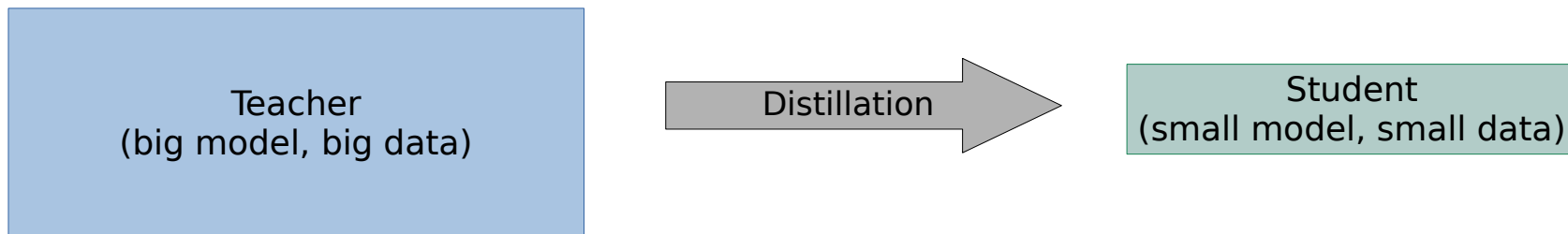
Vilém Zouhar

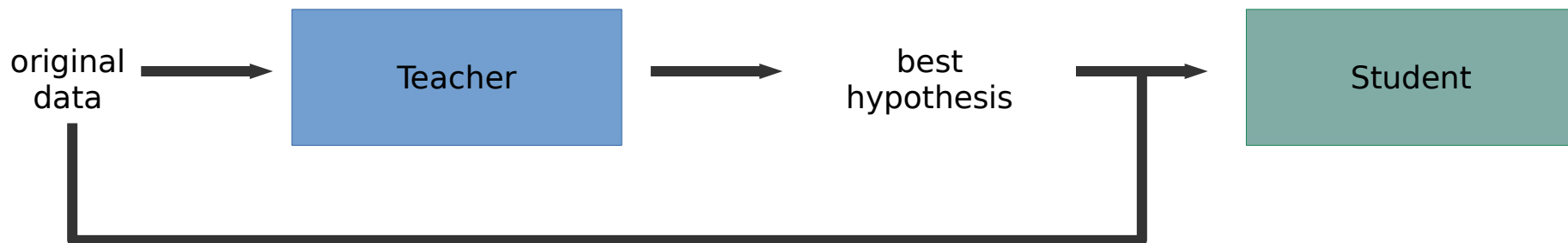## Model {stealing,distillation}

- Transferring knowledge from large model (teacher) to smaller one (student)
- Filtering improves performance (German et. al.)
  - Model inferences have no additional features that imply the quality

Teacher
(big model, big data)

Distillation

Student
(small model, small data)

## Most prevalent

- – Take all the best decoder hypotheses
- – Combine them with original (not teacher train) data
- – Output: two translations per one sentence → Twice as much data

## More sophisticated

- Filter sentences which are suspicious
    - Not recognized by language recognizer (German et. al.)
    - More than / less than k*length of the reference /source
- Filter based on reference quality, e.g. TER (Freitag et. al.)
    - Not main topic of their work

→ **TER is very arbitary, how do other metrics behave?**

→ **What about other sampling methods?**

→ **What about the combination of them?**

## Distillation

- Query teacher on source sentences

- Get 12 translations for every sentence

- Treat them w.r.t. their quality

  - Filter worse ones

  - Oversample better ones

- Quality measured by

  - Decoder score
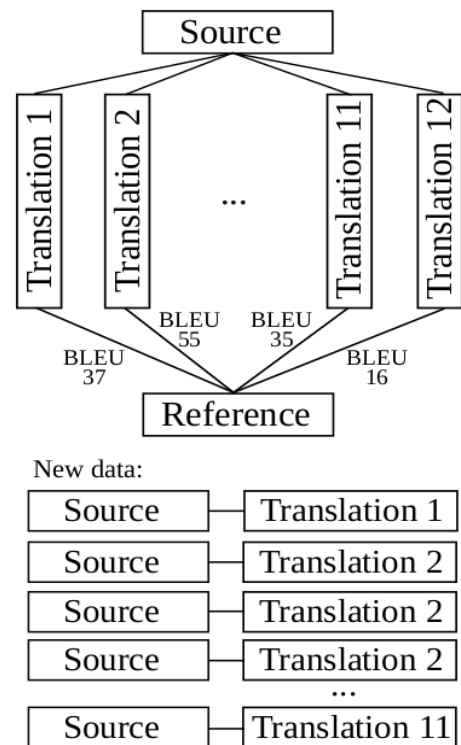
  - Reference based metrics (BLEU, ChrF, TER, SP)



Figure 1. Scheme of an example of hypothesis sampling with BLEU metrics.
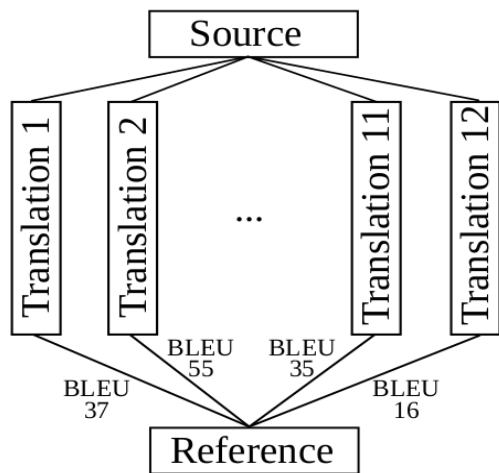
**Process**

- Create new data set (data sampling)

  - Use different sampling methods

- Train student on the new data

  - CS→EN, EN→CS, EN→DE

- Evaluate and see the effect of the chosen method

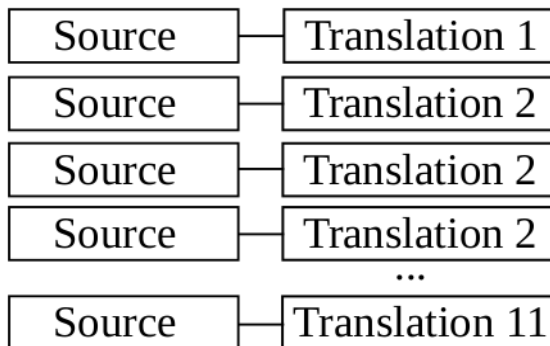# Data Sampling

## Notation

- $F_{metrics}^{k_1, k_2, k_3, ...}$

    - take top translation $k_1$ times, second one $k_2$ times, etc. ...



New data:

| Source | Translation 1 | second best $k_2 = 1$ |
| Source | Translation 2 | top best $k_1 = 3$ |
| Source | Translation 2 | |
| Source | Translation 2 | |
| Source | Translation 11 | |

## Notation

- $T^n_{metrics}$
  - Take **n** top translation hypotheses according to **metrics** (equals $F^{1,1,1,\dots,1(n)}_{metrics}$ )

- $G^m_{metrics}$
  - Take all sentence translations with **metrics**
  - At least **m**

- $Dedup[X]$
  - Deduplicate sentence pairs of **X**

## Notation

- Concatenation of sampling methods

- $T_{BLEU}^{2} + G_{score}^{-10}$

  - Join the top **2** sentences measured by *BLEU*

  - Add them to hypotheses with **decoder score** at least **-10**

| Dataset | CS→EN | EN→CS | EN→DE |
|---|---|---|---|
| $T_{\text{score}}^1$ + Original | 44.4 | 36.4 | 28.3 |
| $Dedup[T_{\text{BLEU}}^4 + T_{\text{score}}^4]$ + Original | 43.7 | 35.3 | 29.1 |
| $S_{\text{score}}^{4,3,2,1} + 2 \times$ Original | 43.9 | 36.1 | 28.3 |
| $S_{\text{BLEU}}^{4,3,2,1} + 2 \times$ Original | **45.5** | **37.3** | 28.8 |
| $S_{\text{BLEU}}^{4,3,2,1} + 4 \times$ Original | **45.5** ⋆ | **37.4** ⋆ | 28.9 |
| $T_{\text{score}}^4 + T_{-}^{12}$ | 41.6 | 33.2 | 28.3 |
| $T_{\text{BLEU}}^4 + T_{-}^{12}$ | 42.6 | 33.9 | 28.7 |
| $T_{\text{BLEU}}^4 + T_{\text{score}}^4$ | 43.3 | 33.2 | 28.9 |
| $Dedup[\sum T_{\text{metric}}^2]$ | 43.6 | 34.7 | 29.1 |
| $Dedup[\sum T_{\text{metrics}}^2] + T_{-}^{12}$ | 40.8 | 32.0 | 27.2 |
| $Dedup[T_{\text{BLEU}}^4 + T_{\text{score}}^4] + T_{\text{BLEU}}^1 + T_{\text{score}}^1$ | 43.5 | 34.7 | 29.2 |
| $Dedup[T_{\text{BLEU}}^4 + T_{\text{score}}^4] + Dedup[T_{\text{BLEU}}^1 + T_{\text{score}}^1]$ | 42.6 | 34.9 | **29.6** ⋆ |
| $Dedup[T_{\text{BLEU}}^4 + T_{\text{score}}^4]$ | 43.5 | 35.0 | **29.3** |

Table 6: BLEU scores for students trained on datasets made of combination of sampling methods. $\sum_{\text{metric}}$ sums over all used metrics (BLEU, ChrF, TER, SP, score).

17

# Contributions

## Meaning of results

- Smaller students
    - Improve results by 0.9 BLEU score
    - Not much, but orthogonal to other tricks

- Students with the same architecture like teacher
    - Improve results by 2.0 BLEU score

19

# Conclusion

**Although widely used:**

- Taking only highest-scoring sentence & original → not the best results
- Combination of oversampled good hypotheses & original
- Choice of reference metric does not significantly influence results

**Caveats:**

- Evaluation on custom dataset, not standardized ones
  - Computation power limitations
- Oversampling vectors chosen at random
  - They are an independent variable as well
- Models very small

# Future Work

**Bigger models:**

- Behavior on larger data & models
- Computationally expensive, but focus only on the best methods

**Experiment with other ML domains:**

- Does distillation oversampling work in other ML domains as well?

**More oversampling schemes**

- Sampling parameters chosen arbitrarily