

# Neural Machine Translation Methods for Translating Text to Sign Language Glosses

Dele Zhu<sup>1</sup>, Vera Czehmann<sup>1, 2</sup> and Eleftherios Avramidis<sup>2</sup>

1. Technical University of Berlin, Berlin, Germany
2. German Research Center for Artificial Intelligence (DFKI), Berlin, Germany



# Introduction

## Glosses:

- written representation of Sign Languages (SL)
- limited representation ability, but useful for:
  - interpreters and educational uses
  - intermediate step for spoken to SL translation
  - investigating methods that may apply to more accurate representations in the future

## Objective:

- Use low-resource MT techniques successful for spoken languages

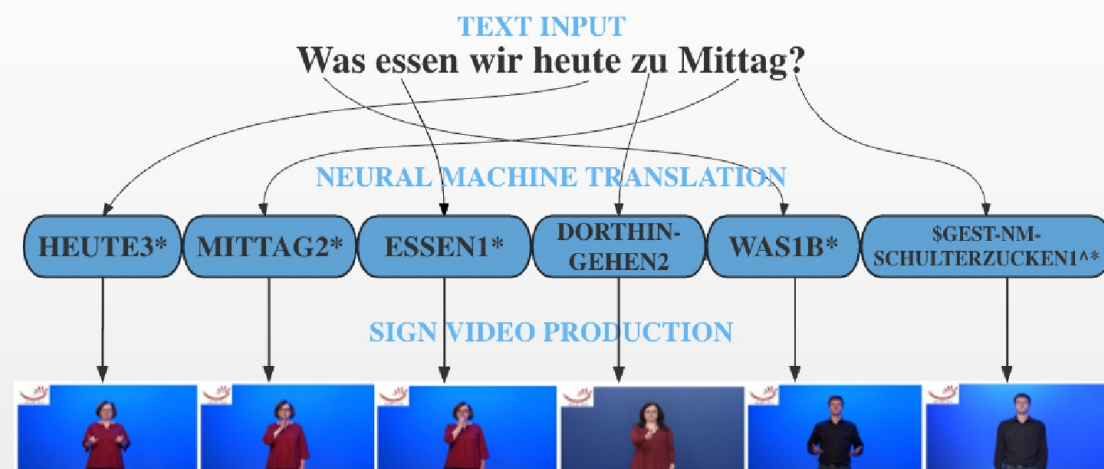


Fig.1 Text-to-video SLT using glosses as an intermediate step (Source of images: Müller et al., 2020).

# Approach Overview

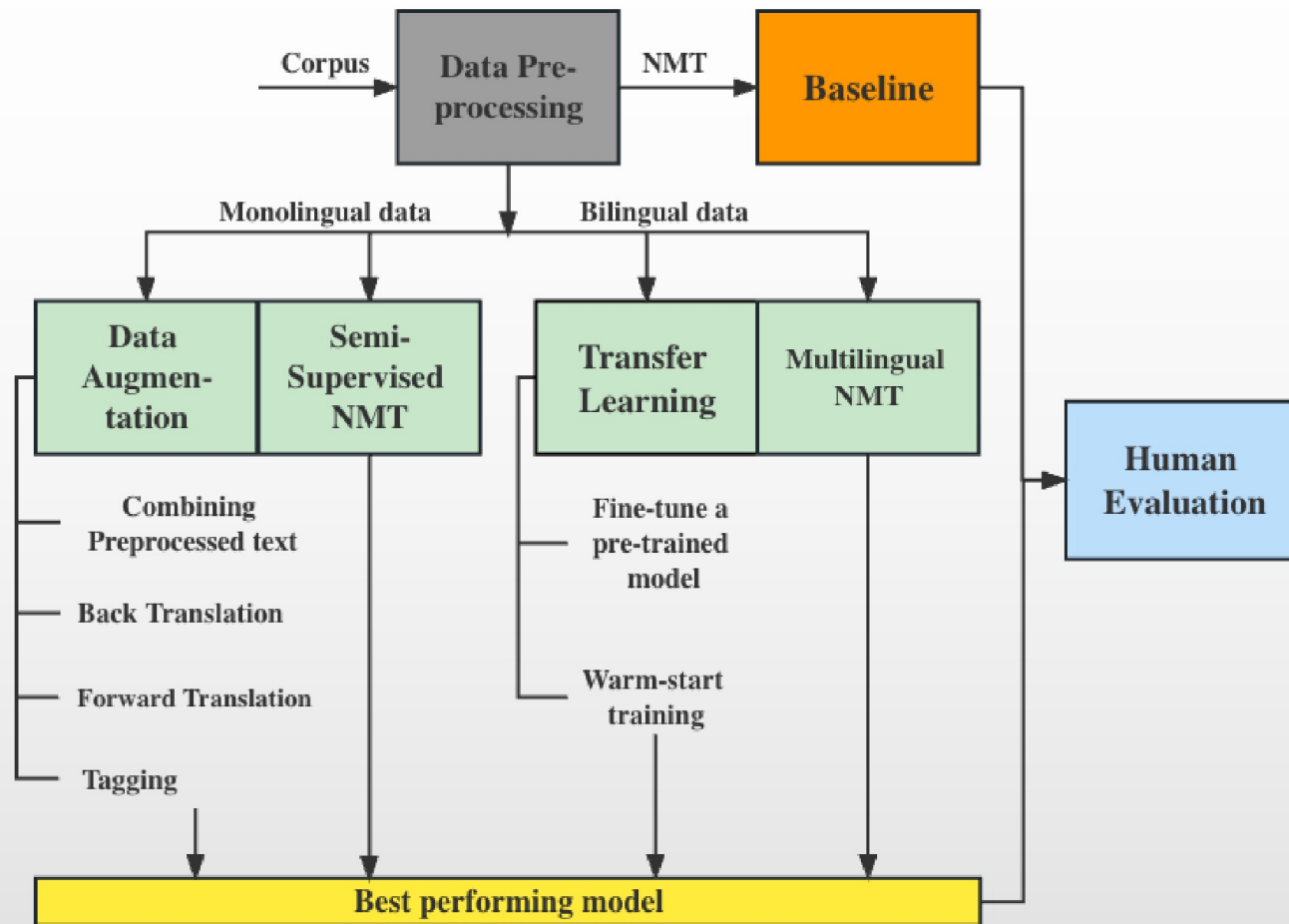


Fig.2 Scheme of experiments

# Datasets



## Main experiments:

- German sign language - Deutsche Gebärdensprache (DGS)
  1. RWTH-PHOENIX-Weather 2014T (Camgoz et al., 2018)
  2. The Public DGS Corpus (Hanke et al., 2020)

## Generalization experiments:

- American sign language (ASL)  
NCSLGR (Vogler and Neidle 2012)

# Experiments - Baseline & Data Augmentation & Semi-supervised NMT



## Baseline:

- Transformer adapted to LR settings (1 encoder, 2 decoders; Gu et al 2018)

## Data Augmentation:

- **Combined preprocessed text:** Collection of source text with different preprocessing techniques
- **Back translation** (Sennrich et al., 2016a): Additional source data translating target data (glosses) from target-to-source model
- **Forward translation** (Zhang and Zong, 2016): Additional synthetic target data translating a source-language monolingual dataset with the baseline system
- **Tagging** (Caswell et al., 2019): Add special prefix before synthetic sentences

## Semi-supervised NMT:

- Copy a monolingual dataset to both source & target side (Currey et al. 2017) and combine with the SL parallel dataset

# Experiments - Transfer Learning & Multilingual NMT



## Transfer Learning:

- **Fine-tuning** of a pre-trained model: Opus-MT de-en model
- **Warm-start training** (Nguyen and Chiang, 2017):
  - Train the parent model with large-scale and SL train set
  - Fine-tune the pre-trained model with SL dataset

## Multilingual NMT (Johnson et al., 2017):

- Train NMT system with both SL dataset and large-scale de-en dataset
- Add **target-language-indicator** before each source sentence

German-to-English:       <**2en**> Wie heißt du? →What is your name?

German-to-DGSGlosses:   <**2gloss**> im süden freundliches wetter→sued region besser

# Results



| Corpus  | System         | BPE Vocab | Dev          |              |              | Test         |              |              |
|---------|----------------|-----------|--------------|--------------|--------------|--------------|--------------|--------------|
|         |                |           | BLEU         | ChrF         | TER          | BLEU         | ChrF         | TER          |
| PHOENIX | Baseline       | 2k        | 22.78        | 51.87        | 55.84        | 20.14        | 52.04        | 56.12        |
|         | Combine        | 2k        | <u>24.01</u> | 52.32        | <u>53.20</u> | <u>21.88</u> | 51.51        | <u>54.53</u> |
|         | Semi+Tag       | 32k       | <u>26.55</u> | <u>55.76</u> | <u>50.83</u> | <u>24.15</u> | <u>55.13</u> | <u>51.17</u> |
|         | Fine-tune      | 65k       | <u>26.39</u> | <b>56.84</b> | 50.88        | <u>24.67</u> | <b>55.97</b> | <u>52.86</u> |
|         | Warm           | 32k       | <b>27.62</b> | <b>56.92</b> | <b>49.25</b> | <u>24.89</u> | <u>55.46</u> | <b>50.40</b> |
|         | Multi+back     | 32k       | <b>28.41</b> | <b>57.54</b> | <u>49.39</u> | <b>26.32</b> | <b>56.70</b> | <u>51.15</u> |
|         | Multi-big+back | 32k       | <b>28.53</b> | <b>57.64</b> | <b>48.93</b> | <b>25.98</b> | <b>56.67</b> | <u>50.94</u> |
| DGS     | Baseline       | 5k        | 4.04         | 31.20        | 79.34        | 3.13         | 30.38        | 78.64        |
|         | Combine        | 5k        | 3.71         | 29.97        | 80.21        | 2.75         | 29.31        | 80.01        |
|         | Semi+Tag       | 32k       | <u>5.00</u>  | <u>32.69</u> | <u>79.47</u> | <u>4.10</u>  | <u>31.30</u> | <u>78.67</u> |
|         | Fine-tune      | 65k       | <u>5.82</u>  | <u>35.05</u> | <u>79.92</u> | <u>4.53</u>  | <b>34.14</b> | <u>78.98</u> |
|         | Warm           | 32k       | <u>5.87</u>  | <u>33.42</u> | <b>74.07</b> | <u>4.55</u>  | <u>31.90</u> | <u>74.54</u> |
|         | Multi+back     | 32k       | <u>5.35</u>  | <u>33.43</u> | <u>78.30</u> | <u>4.85</u>  | <u>32.16</u> | <u>76.76</u> |
|         | Multi-big+back | 32k       | <b>6.82</b>  | <b>35.57</b> | <u>76.37</u> | <b>5.78</b>  | <b>33.87</b> | <u>76.12</u> |

Table 1: Selected automatic metric scores of extensive experimentation search of the two DGS corpora. We **boldface** all the values that are not statistically significantly different from the best value of each evaluation metric and underline the results that are statistically significantly higher than baseline at the 95% confidence level.



## Comparison with previous work – generalization test on ASL

| Approach                           | Dev<br>BLEU $\uparrow$ | Test<br>BLEU $\uparrow$ |
|------------------------------------|------------------------|-------------------------|
| Amin et al. (2021)                 | -                      | 10.42                   |
| Egea Gómez et al. (2021) $\dagger$ | -                      | 13.13                   |
| Stoll et al. (2020)                | 16.34                  | 15.26                   |
| Zhang and Duh (2021)               | -                      | 16.43                   |
| Li et al. (2021)                   | -                      | 18.89                   |
| Saunders et al. (2020b)            | 20.23                  | 19.10                   |
| Saunders et al. (2022)             | 21.93                  | 20.08                   |
| Egea Gómez et al. (2022)           | -                      | 20.57                   |
| Walsh et al. (2022)                | 25.09                  | 23.19                   |
| <b>Our PHOENIX Multi+back</b>      | <b>28.41</b>           | <b>26.32</b>            |

Table 4: Results comparison with recent work. ( $\dagger$ ) We compute the BLEU by ourselves, as the authors of paper only present the BLEU score in character level.

| System     | BLEU         | Test<br>ChrF | TER          |
|------------|--------------|--------------|--------------|
| Baseline   | 10.50        | 30.65        | <b>78.95</b> |
| Back       | 9.37         | 28.25        | <b>78.67</b> |
| Warm-start | <u>12.11</u> | <u>33.53</u> | 83.44        |
| Multi      | <u>12.35</u> | <u>38.33</u> | <b>78.26</b> |

Table 2: Automatic metric scores for NCSLGR corpus.



# Human Evaluation



| System                           | Size | automatic       |                 |                | human           |                |
|----------------------------------|------|-----------------|-----------------|----------------|-----------------|----------------|
|                                  |      | BLEU $\uparrow$ | ChrF $\uparrow$ | TER $\uparrow$ | Mean $\uparrow$ | Std $\uparrow$ |
| PHOENIX Egea Gómez et al. (2021) | 642  | 13.13           | 46.86           | 73.33          | 2.74            | 1.64           |
| PHOENIX baseline                 |      | 20.14           | 52.04           | 56.12          | 3.85            | 1.58           |
| PHOENIX Multi+back               |      | <b>26.32</b>    | <b>56.70</b>    | <b>51.15</b>   | <b>4.44</b>     | <b>1.35</b>    |
| DGS Baseline (sampled 10%)       | 511  | 3.44            | 29.56           | 78.55          | 2.49            | 1.81           |
| DGS Multi-big (sampled 10%)      |      | <b>6.97</b>     | <b>33.16</b>    | <b>73.45</b>   | <b>3.28</b>     | <b>1.60</b>    |

Table 3: System comparison based on the human evaluation. The **bold-faced** systems are significantly better than the respective baselines.

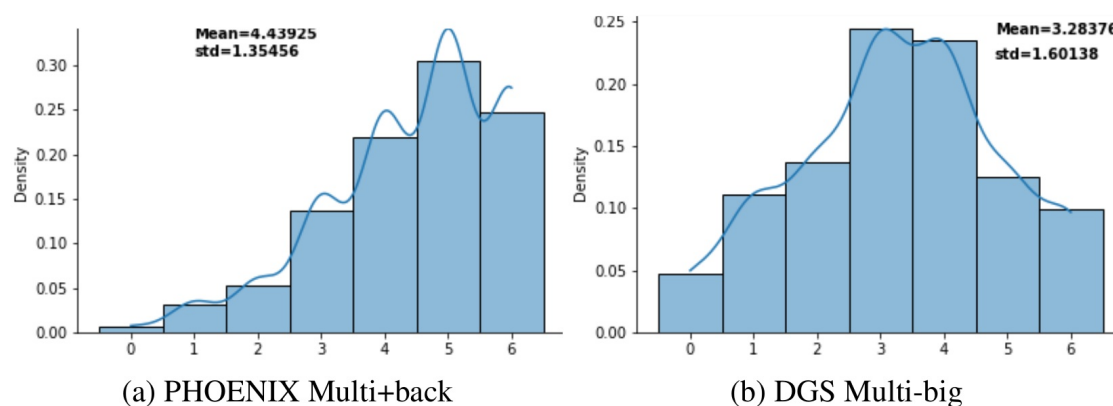


Figure 3: Density of human evaluation scores for the two best-scoring systems.

## Conclusion & Contributions



this paper is the first work on text-to-gloss machine translation

- to achieve significant improvements, as compared to the baseline methods and related work, on the two known natural German SL datasets annotated with glosses,
- to perform extensive experimentation with most known LRL-related MT methods and their combinations and in particular:
  - to apply **semi-supervised NMT** by copying the monolingual data to both the source and target side
  - to use **transfer learning** via the warm-start strategy, and
  - to use a **multilingual NMT** setting with the focus on improving the text-to-gloss direction.



Association for  
Computational  
Linguistics

**ACL 2023**

**Thanks for your attention!**