

# Assessing Large Pre-trained Models for Sign Language Processing: Is Text-Only Superior to Multimodal?

Anonymous ACL submission

## Abstract

Motivated by the recent success of text-only modeling in certain vision-language tasks, this paper proposes that sign language processing can also use (large) text-only language models for inference, freeing sign language models from the necessity of low-resource multimodal learning from scratch. To compare the performance of pre-trained text-only models against multimodal ones, we introduce the first text-only and multimodal large (7B) language models to be pre-trained and then fine-tuned on a sign language recognition task. We propose new prompting strategies and fine-tuning strategies for text-only signed language processing, incorporating both linguistics of signed languages and theoretically motivated strategies to mitigate catastrophic forgetting (of spoken language). We test the generalization of these models to other sign language recognition and generation tasks, showing text-only models are capable sign language models that are still adept at spoken language tasks and, by changing the prompt, can even generalize to new prosodic and iconic sign recognition tasks. Finally, we analyze trade-offs between our text-only and multimodal models. Our code and model checkpoints will be open-source.

## 1 Introduction

Traditionally, multimodal models are the de facto choice for sign language processing tasks, given the continuous and visual nature of signs. However, recent work by Wang et al. (2022) and Cheng et al. (2023) calls into question the presumption that multimodal models are best for all visual tasks. In particular, they use text-based descriptors and text-only language models to achieve strong performance on few-shot video-to-text tasks *without visual pretraining or fine-tuning*. Motivated by this, we suggest that sign language processing can also be achieved through text-only language modeling.

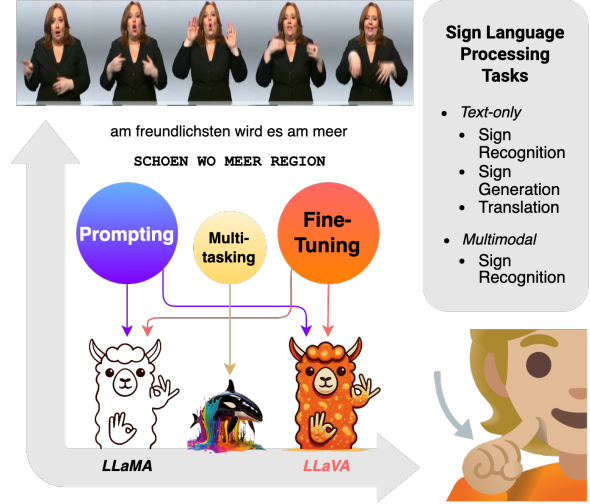


Figure 1: We show that text-only open foundation models when prompted or fine-tuned, can learn to perform sign language processing tasks compared to their multimodal alternatives. Further, multitask fine-tuning alleviates forgetting of spoken language capabilities (e.g., QA tasks in English).

Rather than a visual medium, we focus on sign language as a *language*, using language learning theories and linguistic theories of signed languages to create one-shot text prompts and fine-tune text-only models to replace the visual modality of signs in multimodal models. As a primary benefit, sign language processing may be freed of traditional, low-resource, multimodal learning from scratch while also gaining the benefits of the recent successes in language modeling, such as capability at in-context learning.

Our central hypothesis stems from the fact that text-only pre-trained models are already able to represent rich semantic and prosodic information based on the context in spoken languages, like English (Garí Soler and Apidianaki, 2021; Saba, 2023). Sometimes text-only models can do this even better than multimodal pre-trained models (Wang et al., 2022). If we can use this semanti-

cally rich text-based representation space to capture the important visual concepts of signed-languages (e.g., intensity), then we defeat the purpose of asking multimodal models to learn the visual semantics of signed languages from scratch. So, we ask the question, *why can't text-only pre-trained models use the textual modality to learn effective visual sign language representations?*

In answering this question, we demonstrate the benefits of applying large pre-trained language models to tasks in signed languages. Moreover, our results point to a future where language models can also be pre-trained on signed languages *without significant degradation of their spoken language capabilities*, marking an essential step for the wider adoption of signed languages. In more detail, our contributions are described as follows.

1. We use linguistic rules to prompt and fine-tune large (7B) text-only pre-trained models and compare them to multimodal pre-trained models on sign recognition for the first time.
2. We theoretically and empirically study the problem of catastrophic forgetting during fine-tuning on sign language data, providing solutions to resolve this issue.
3. We use annotator costs, carbon emission, and performance differences to analyze trade-offs between the use of a multimodal model and our linguistically-backed text-only models.

Our results show fine-tuning large, pre-trained, text-only models offers new generalization capabilities compared to previous sign recognition training strategies; e.g., via in-context learning. We also do a case study on emergent iconicity by pre-trained models for signed languages. All code, data, and model checkpoints will be publicly available.

## 2 Related Work

Besides text-only models like LLaMA (Touvron et al., 2023a), Mixtral (Jiang et al., 2024), QWEN (Bai et al., 2023), Orca (Mukherjee et al., 2023), Phi (Gunasekar et al., 2023), multimodal models have been gaining popularity, especially in computer vision communities. Large Vision-Language models such as LLaVA (Liu et al., 2023b), Video-LLaMA (Zhang et al., 2023), Video-LLaVA (Lin et al., 2023), LanguageBind (Zhu et al., 2024), MultiModal-GPT (Gong et al., 2023), Mirasol3B (Piergiovanni et al., 2023), LAVIS (Li et al., 2023), LaViLa (Zhao et al., 2023), and UniVL (Luo et al., 2020) propose to align representations of combina-

tions of images, videos, text, and/or speech signals with human judgments. Further details of these and similar models have been discussed in a survey paper by Yin et al. (2023). However, none of these models claim to include sign language processing tasks in their pre-training or fine-tuning data. Through our theoretical and empirical studies, this paper aims to address this gap.

The absence of literature using large models for sign processing is mainly due to the low-resource nature of signed languages (Yin et al., 2021). However, there have been several lines of research applying transformer-based language models to sign language translation (Camgoz et al., 2018; Yin and Read, 2020; Chen et al., 2023b), sign language understanding (Hu et al., 2023; Moryossef et al., 2021), sign generation (Stoll et al., 2020), Sign-Writing translation (Jiang et al., 2023), incorporating facial expressions (Viegas et al., 2023), modeling prosody (Inan et al., 2022), and sign language segmentation (Moryossef et al., 2023). To the best of our knowledge, Lee et al. provides the only other work that leverages (smaller, but still large) language models with shared vocabularies for sign language processing. They focus on older models (without RLHF, Ouyang et al., 2022). However, none of these works involve modern large language models (text-only nor multimodal), which we introduce in this paper for the first time.

## 3 Method

In this section, we introduce the details of both the text-only and multimodal foundation models used in experiments (see Figure 2), along with the studied prompting and fine-tuning strategies. We also provide a theoretical basis for choosing appropriate training data to prevent foundation models from forgetting the traditional language capabilities on which they were pre-trained.

### 3.1 Sign Data, Tasks, and Models

**DGS Data** Due to widespread adoption as a benchmark in the sign language processing community, we use the RWTH-PHOENIX-14T<sup>1</sup> corpus of weather forecast signs in German Sign Language (DGS). This dataset contains around 7000 training samples, 500 validation samples, and 600 test samples. Each sample has a video, a text in spoken German, and a gloss – which is an intermediary

<sup>1</sup><https://www-i6.informatik.rwth-aachen.de/~koller/RWTH-PHOENIX-2014-T/>

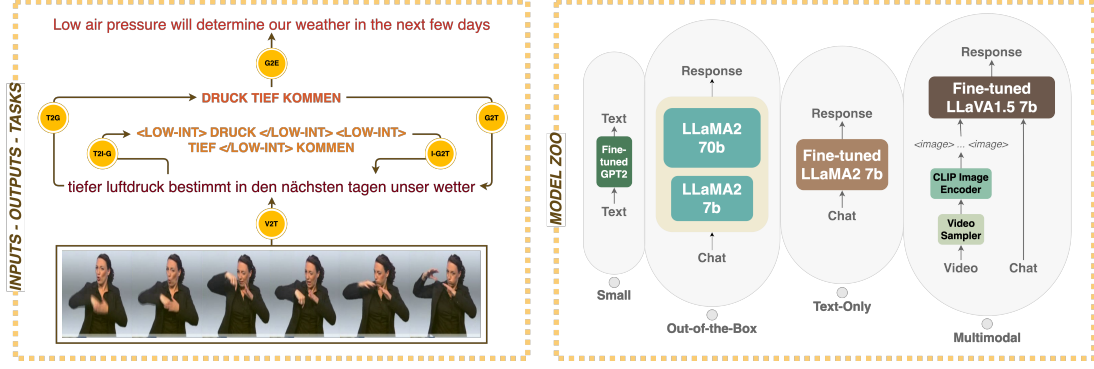


Figure 2: This figure presents a summary of all the inputs, outputs, tasks, and models we are using and introducing in this paper. The box on the left contains a sample from the RWTH-PHOENIX-14T dataset. From top to bottom, the sentences are English text, DGS glosses, intensified DGS glosses, and German text. Yellow knobs represent tasks, in which the acronyms of the tasks are inlaid (please refer to Section §3.1 for detailed task names).

textual representation of signs – in German Sign Language. Video samples consist of frames of multiple signers sampled at 25 fps, with a size of 210 by 260 pixels. We also include an enhanced version of this dataset, which contains *intensifier* information in its gloss representations as introduced by (Inan et al., 2022). Intensifiers in signed languages are depicted through non-manual markers and can change the meaning of a sign, and this dataset contains additional tokens to capture intensifier information. We also translate the German text to English text to provide data for a cross-lingual task (discussed next). We use Google Translate.<sup>2</sup>

**Tasks** As RWTH-PHOENIX-14T is a parallel corpus between spoken German and DGS, most previous research has focused on translation tasks between these languages. In this paper, we focus on translating DGS to German (broadly considered as a sign understanding or recognition task) and German to DGS (broadly considered as sign generation). In addition to these, we introduce additional tasks to test generalization. Specifically, we consider:

- **(G2T) DGS Gloss to German Text:** a text-only translation task from textual intermediary representations of DGS (glosses) to German text.
- **(T2G) German Text to DGS Gloss:** the inverse problem of the above and is text-only.
- **(V2T) DGS Videos to German Text:** a multimodal task where the input is a video of a signer signing in DGS, and the output is German text.
- **(I-G2T) Intensified DGS Gloss to German Text:** a text-only task with augmented DGS tokens. Additional symbols `<HIGH-INT>` and `<LOW-INT>`

are wrapped around glosses to depict intensity in the video that is not depicted in traditional gloss representations (Inan et al., 2022).

- **(T2I-G) German Text to Intensified DGS Gloss:** the inverse problem of (I-G2T), still text-only.
- **(G2E) DGS Gloss to English Text:** a novel task of cross-modal translation, where DGS glosses from the German Sign Language family are translated to English text from the spoken Indo-European language family. Without any pre-training, this is a difficult test of generalization and composition of contextualized meanings across traditional and signed languages.

To test generalizability and in-context learning, G2T is the only DGS task we use for any fine-tuning (see § 3.3). All the other tasks are used to evaluate the models’ performance.

**Models** In this paper, we use two main foundation models: LLaMA-2 7B Chat (Touvron et al., 2023b) for text-only inputs and LLaVA 1.5 7B (Liu et al., 2023a,c) for multimodal inputs. To compare with traditional sign language processing approaches, which use smaller language models *sans any foundational pre-training*, we also use a randomly initialized GPT-2 model (Radford et al., 2019) trained on the G2T task of the RWTH-PHOENIX-14T dataset. This controlled difference allows us to quantify the utility of concepts learned during foundational training (e.g., in LLaMA and LLaVA) on sign language processing. Lastly, for G2T task, we use LLaMA-2 70B with 4-bit quantization<sup>3</sup> to show how the number of parameters affects the results.

<sup>2</sup><https://cloud.google.com/translate/>

<sup>3</sup><https://ollama.com/library/llama2:70b>

Prompt Strategy	BLEU <sub>1</sub>	ROUGE <sub>1</sub>	BS-F1
zero-shot prompt	24.5	0.277	0.841
rule-based prompt	22.8	0.255	0.836
notation prompt	24.3	0.277	0.840
one-shot prompt	27.1	0.309	0.851

Table 1: Preliminary evaluation of prompting strategies on the validation set of RWTH-PHOENIX-14T using LLaMA-2 7B. The prompts are given in Appendix § B. BS-F1 refers to BERTScore-F1.

## 3.2 Prompting and Initial Results

To replace the visual modality of signed languages, we propose to prompt text-only foundation models using linguistic and cognitive science rules of glossing and signing. We first prompt these foundation models for the tasks described in § 3.1. We incorporate the following linguistic rules of signed languages into the design of the prompts that we provide to the models:

- **zero-shot prompt:** The prompt is structured as, "This is a sentence in German Sign Language glosses: <glosses>. You MUST translate these to spoken German. You MUST give the answer directly without any other text." Does not contain any linguistic rules.
- **rule-based prompt:** The prompt is structured as five rules of glossing semantics. These rules are described in (Hanke et al., 2020).
- **notation prompt:** This is structured as a set of rules about gloss morphologies. These rules are borrowed from Stein et al. (2010).
- **one-shot prompt:** This prompt gives a single example of a DGS gloss and a corresponding German text. This example is formatted following the semantic and morphological rules above.

Initially, we experiment with four different prompt strategies, then we pre-select two (the top-performing prompting strategy and the basic one) among these. All prompts are given in Appendix B.

For the multimodal foundation model, we provide a single chat template. We use a mixed prompting strategy, where the video of signers is sampled at 50 frame intervals, fed into a CLIP-based Image Encoder (Radford et al., 2019), and then incorporated into the prompt tokenization by the use of <image> for each frame. Then, the image portion of the prompt is succeeded by the text-based prompt "This video is in German Sign Language. What is the sentence being signed in German?"

## 3.3 Supervised Fine-Tuning with LoRA

Besides in-context learning via few-shot prompts, we also consider fine-tuning LLaMA2 and LLaVA1.5 models using Supervised Fine-Tuning<sup>4</sup>, which is a supervised training method in addition to the RLHF algorithm (Ouyang et al., 2022) for chat-based model training, which aligns the models' representations with human judgments. In this case, the human annotations are either glosses or text. For fast model training and reduced memory consumption, we use Low-Rank Adaptation of Language Models (LoRA) as introduced by Hu et al. (2022). We give details of model hyperparameters and training details in Appendix A.

**Sign-Only Fine-Tuning** As noted, for text-only models we fine-tune on the G2T task from § 3.1, and for multimodal we fine-tune on the V2T task. This provides the model a simple introduction to the meaning of signed glosses by grounding them to their parallel German language context.

**Multitasking Fine-Tuning** As we discuss in the next section, we hypothesize that the former (sign-only) tuning strategy can lead to catastrophic forgetting. Due to the shared token vocabulary, the model may overwrite existing knowledge and semantics in the contextualized representations of traditional language tokens. Intuitively, we expect that forcing the model to "replay" traditional language tasks from pre-training will prevent forgetting. To accomplish this, we also train on an additional (traditional task) dataset (OpenOrca<sup>5</sup>) randomly mixing the sign and traditional data during tuning. This dataset consists of system prompts, questions, and responses, augmented from the FLAN collection (Longpre et al., 2023). It is commonly used to fine-tune smaller open models such as LLaMA for better task success, surpassing proprietary models such as GPT-3.5. The dataset is mainly in English and consists of multiple tasks: entailment and semantic understanding, temporal and spatial reasoning, causal judgment, multilingual understanding, world knowledge, logical and geometric reasoning, and similar other tasks (Mukherjee et al., 2023). While the original dataset contains around 3 million samples, we use the same split sizes as RWTH-PHOENIX-14T to ensure balance in sign/traditional task prioritization.

<sup>4</sup>[https://huggingface.co/docs/trl/main/en/sft\\_trainer](https://huggingface.co/docs/trl/main/en/sft_trainer)

<sup>5</sup><https://huggingface.co/datasets/Open-Orca/OpenOrca>



### 3.4 Theory: Multi-Tasking Mitigates Forgetting

Motivated by neuroscience, *experience replay* has been suggested as a strategy to reduce forgetting in machine learning, with positive results (Rolnick et al., 2019). Moreover, replay has been studied in mathematical theories of how language models learn with similar success (Sicilia and Alikhani, 2022). Our multi-tasking strategy (discussed above) can be viewed as a type of experience replay since many tasks from OpenOrca are presumed to be similar to prior experience during pre-training.<sup>6</sup> In this section, we re-frame our learning environment using the theoretical tools provided by Sicilia and Alikhani (2022) to motivate our hypothesis. Namely, we show that multi-task fine-tuning (i.e., replay) can help mitigate forgetting in shared-vocabulary sign processing with foundation models.

**Sign Language Processing Algorithm** Our current task setup is of a translation algorithm, where the model learns how to translate from a signed language to a spoken language and vice versa. Specifically, in the case of foundation models learning this, the algorithm contains two specific steps:

1. **Pre-Training:** Foundation models are trained on multiple tasks that do not include (many or any) sign-language-specific tasks. Using the terminology of Sicilia and Alikhani (2022), this process picks the weights to minimize the *test divergence* or “error”  $\mathbf{TD}_{PT}$  where  $PT$  is the pre-training data distribution:

$$\begin{aligned} \mathbf{TD}_{PT}(\theta) &= \mathbf{E}[\ell(D, \hat{D})] \\ D &\sim \text{LM}(X; \theta), \hat{D} \sim \text{ANOT}(X) \end{aligned} \quad (1)$$

where  $\text{LM}$  is the foundation model (e.g., a language model),  $\text{ANOT}$  is a human completion/annotation provided the same context  $X$  (e.g., a prompt), and  $X$  ranges over the dataset  $PT$ . The test  $\ell$  compares any measure of the quality or other properties of the generated text between foundation model and human; e.g., it can represent automatic metrics like BLEU, ROUGE, or error at next-word prediction as well as more abstract tests (like human preference).

2. **Fine-Tuning:** In this stage, the foundational model is fine-tuned on sign language processing tasks such as gloss-to-text translation. For the

*sign-only fine-tuning*, we call this data distribution  $DGS$ . So, abstractly, our sign-only fine-tuning process described previously attempts to minimize  $\mathbf{TD}_{DGS}(\theta)$ .

**Problem** When we write out the pre-training and fine-tuning objectives clearly in the terminology of Sicilia and Alikhani (2022), it is clear that the two processes optimize *different* objectives (e.g., over different datasets). There is no way to ensure that picking  $\theta$  to minimize  $\mathbf{TD}_{DGS}$  will not have a negative impact (i.e., increase)  $\mathbf{TD}_{PT}$ . This potential for increase in error on the pre-training tasks characterizes the behavior we call “forgetting.”

**Solution** As mentioned, we also consider a *multi-tasking fine-tuning* strategy where  $DGS$  data and tasks similar to the pre-training data are mixed. This multi-tasking data can be represented by the mixture distribution:

$$\text{MIX} = \alpha \text{PT} + (1 - \alpha) \text{FT} \quad (2)$$

where  $\alpha \in (0, 1)$  is a weighing factor between the probabilities assigned by two datasets. Instead of sampling  $X$  from only  $PT$  or only  $FT$ , we flip an  $\alpha$ -weighted coin to pick from which we sample. Holding all else constant, this implies the equality:

$$\mathbf{TD}_{\text{MIX}} = \alpha \mathbf{TD}_{PT} + (1 - \alpha) \mathbf{TD}_{FT}. \quad (3)$$

By this choice, we can see:

$$|\mathbf{TD}_{\text{MIX}} - \mathbf{TD}_{PT}| \quad (4)$$

$$= (1 - \alpha) |\mathbf{TD}_{FT} - \mathbf{TD}_{PT}| \quad (5)$$

$$< |\mathbf{TD}_{FT} - \mathbf{TD}_{PT}|. \quad (6)$$

Since  $\mathbf{TD}_{\text{MIX}}$  is always closer in magnitude to  $\mathbf{TD}_{PT}$  than  $\mathbf{TD}_{FT}$ , we can see that minimizing  $\mathbf{TD}_{\text{MIX}}$  can better prevent large increases  $\mathbf{TD}_{PT}$ , or “forgetting.” This simple inequality provides a theoretical motivation for our multi-tasking suggestion in § 3.3. Our empirical results in § 4 also confirm our theoretical hypotheses.

## 4 Findings

In this section, we conduct experiments to answer six research questions. We outline all of these questions in the following sections and give answers to them with our findings.

### 4.1 Automatic Metrics

For all the tasks, to compare the generated text with the ground truth, we make use of automatic

<sup>6</sup>Most open-source models do not share training data.

Performance of All Models on All Tasks										
Task	Prompt Strategy	Finetuned GPT2			Not Finetuned LLaMA2 7b			Multitasking LLaMA2 7b		
		B <sub>1</sub>	R <sub>LSum</sub>	BS <sub>F1</sub>	B <sub>1</sub>	R <sub>LSum</sub>	BS <sub>F1</sub>	B <sub>1</sub>	R <sub>LSum</sub>	BS <sub>F1</sub>
T2G	one-shot	1.419	0.027	0.798	8.556	0.127	<b>0.818</b>	<b>10.921</b>	<b>0.165</b>	0.794
T2G	zero-shot	1.879	0.030	0.810	8.335	0.122	<b>0.802</b>	<b>10.485</b>	<b>0.161</b>	0.794
G2E	one-shot	3.604	0.066	0.822	<b>9.226</b>	0.084	0.807	3.104	0.034	0.828
G2E	zero-shot	3.931	0.056	0.808	<b>12.369</b>	0.103	0.816	5.442	0.064	0.83
I-G2T	one-shot	2.242	0.048	0.791	9.573	0.111	0.691	<b>17.637</b>	0.155	0.524
I-G2T	zero-shot	1.642	0.043	0.768	11.589	0.143	0.769	<b>21.157</b>	0.279	0.845
T2I-G	one-shot	1.305	0.054	0.815	42.277	0.576	0.897	<b>43.636</b>	0.156	0.778
T2I-G	zero-shot	0.050	0.062	0.802	<b>56.128</b>	0.704	0.910	43.229	0.155	0.778

Table 2: This table shows the performance of all the models for all the tasks that we introduce in Section §3.1 for the test set. The one-shot strategy contains an example for the task. B<sub>1</sub> corresponds to BLEU-1, R<sub>LSum</sub> corresponds to ROUGE, and BS<sub>F1</sub> corresponds to BERTScore.

metrics. We use both traditional n-gram metrics of BLEU (Papineni et al., 2002), ROUGE (Lin, 2004), and also use learned generation metrics such as BERTScore (Zhang\* et al., 2020). For the implementation of all of these, we use the Huggingface evaluate library<sup>7</sup>. We do not include classification-based metrics, as our language models generate full-textual responses rather than classes.

**RQ1: How do different prompting strategies affect the performance of the pre-trained (not fine-tuned) text-only model?** Using these automatic metrics, we first evaluate the performance of the prompting strategies for the non-finetuned LLaMA 1.5-7b model. We present these results for all the tasks in Table 1. These show that rule-based prompts and notation-based prompts perform similar to or less than zero-shot prompts. One-shot prompts are the best prompting strategy where an example translation is provided; this reinforces assumptions of few-shot prompts performing better than zero-shot.

TEST SET				
Models	B <sub>1</sub> ↑	B <sub>2</sub> ↑	R <sub>LSum</sub> ↑	BS <sub>F1</sub> ↑
LLaMA2 7b	<b>12.057</b>	1.968	0.144	0.764
LLaMA2 70b	11.281	<b>2.054</b>	<b>0.175</b>	<b>0.798</b>

Table 3: This table shows the performance differences between LLaMA2 7b, and LLaMA2 70b variants. The bigger model generates more intelligible sentences, yet fails to carry out the translation task.

**RQ2: How does the number of parameters affect the performance of the model in text-only**

<sup>7</sup><https://huggingface.co/docs/evaluate/>

**sign language processing tasks?** We show the effects of the number of parameters of the text-only model for the G2T task in Table 3. A higher number of parameters does not always correlate with better automatic metric results. A higher number of parameters also increases the fine-tuning duration.

**RQ3: How does supervised fine-tuning the text-only model on the G2T affect the performance?**

To answer this questions, we fine-tune several foundational models. These results compare the baseline of a small GPT-2 model which is fine-tuned on the G2T task, with our larger models LLaMA 2 7b, and Multitasking LLaMA 2 7b. We first show the results for the fine-tuned task of G2T in Figure 3.

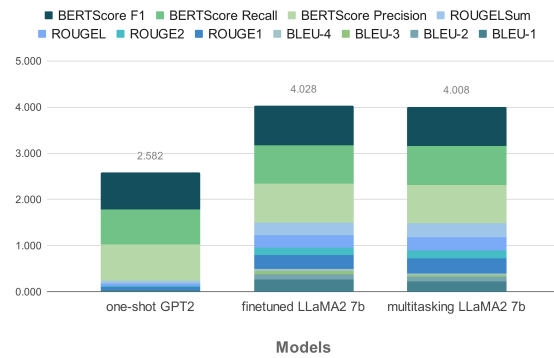


Figure 3: This figure shows the bar plot of ablations on the G2T task. It can be seen that the performance of the larger LLaMA-based models is higher overall compared to a smaller model (GPT2). Also, multitasking to prevent forgetting does not affect model performance.

**RQ4: Can the performance in G2T generalize to other sign language processing tasks? Does it perform better than smaller transformer-based**

## language models, which are not pre-trained?

To answer, we show the results for all the sign language tasks in Table 2. It can be seen that the multitasking model outperforms the smaller model across all tasks. There is variability across tasks on whether the original LLaMA model performs better than the multitasking version. This can be caused by the differences in the task setup and input outputs being more easily with semantic information only from the pre-trained representation.

**RQ5: How does the fine-tuned multimodal model perform in comparison to the text-only model? What are the implications of videos as inputs rather than glosses?** To answer this, we fine-tune LLaVA 7b on the RWTH-PHOENIX-14T videos. The performance differences are shown in Table 4. Here, it can be seen that the fine-tuned model is performing better than the non-finetuned model across all metrics. The implications of using videos rather than glosses mean that in the absence of signer annotations on the glosses, videos can be used as input as well, with a decrease in the overall performance (compare Figure 3 and Table 4), but text-only models outperform video models. We give a more detailed analysis of this in our trade-offs section §6.

Multimodal Sign Understanding (SignVideo2Text)				
Models	TEST SET			
	B <sub>1</sub> ↑	B <sub>2</sub> ↑	R <sub>LSum</sub> ↑	BS <sub>F1</sub> ↑
LLaVA1.5 7b	2.140	0.006	0.022	0.658
ft-LLaVA1.5 7b	<b>12.776</b>	<b>2.404</b>	<b>0.103</b>	<b>0.779</b>

Table 4: This table shows the automatic metric results for the translation task of German Sign Language video to German Text. ft-LLaVA1.5 7b is the fine-tuned model.

**RQ6: Given the theoretical background of forgetting, how does including multiple tasks during fine-tuning affect performance?** To answer this question, we use the generic open language model Benchmarks by EleutherAI Evaluation Harness (Gao et al., 2023) and test the performance difference between the multitasking, finetuned, and non-finetuned models. We show the results in the bar plot in Figure 4. We can empirically observe that there is a drop in performance between non-finetuned and fine-tuned LLaMA2 models. This shows the data shift that we have outlined in Section §3.4 due to the differences in data distribu-

tion between the pretrained LLaMA2 and the sign-finetuned LLaMA2. This strongly suggests that there is forgetting of the original capabilities of the pretrained model.

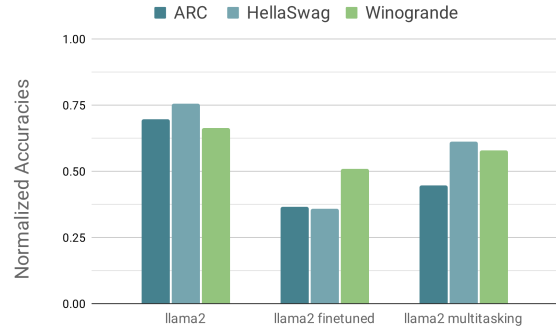


Figure 4: This is the bar plot showing the ablation study on the multitasking/mixing model on the Open Language Model Benchmarks of ARC (Clark et al., 2018), HellaSwag (Zellers et al., 2019), and Winogrande (Sakaguchi et al., 2019), all degrade (forgetting) when LLaMA2 is fine-tuned on the sign language tasks, and when trained on multiple tasks, it performs better.

## 5 Towards Prosodic, Iconic and Semantically-Rich Sign Language Representations via LLMs

Signed languages and the current machine learning setups for sign language processing systems have been constrained to multimodal translation systems mostly, as can be seen from our tasks as well. However, sign interpretation and production by humans are not translation-based processes between modalities. Cognitive science, neuroscience, and linguistics research into the signed languages by Kubicek and Quandt (2019, 2021) show that prosody during signing affects interpretation and action recognition, and Karadöller et al. (2023); Chen et al. (2023a); Campisi et al. (2023) show that different signed languages use different levels of iconicity and iconic signs can facilitate interpretation. In this section, we present a case study on the current iconicity characteristics that are developed during the fine-tuning of the LLaMA2 model by using emojis as placeholders for intensifiers.

### 5.1 Iconicity Case Study: Emojis as Intensifiers

During the fine-tuning of the LLaMA2 7b+ model, it has been observed in the generated outputs for the intensified tasks there are emojis, even though the model is not instructed to include emojis, and the

Ground Truth Text	Intensified Gloss	Generated Text
woher soll die wärme derzeit auch kommen	WOHER <HIGH-INT>WARM </HIGH-INT>	wo es auch am wochenende warmer wird 🌞🌞🌞
where is the heat supposed to come from?	WHERE WARM	where it gets warmer on the weekend too 🌞🌞🌞
morgen gibt es mal sonne mal wolken in der nordhälfte zeigt sich die sonne häufiger	MORGEN SONNE WOLKE WECHSELHAFT NORD <HIGH-INT>SONNE </HIGH-INT>OFT	morgen wird es auch wieder sehr sonnig und teils auch wolkig und auch <b>hochintensiv</b> 🌞
Tomorrow there will be clouds, in the north half the sun will be more common	TOMORROW SUN CLOUD VARIABLE NORTH <HIGH-INT>SUN </HIGH-INT>OFTEN	Tomorrow it will be very sunny and sometimes also cloudy and also <b>highly intensive</b> 🌞

Table 5: This table shows three samples of German Text, DGS Gloss, and the generated text by the LLaMA2 7b+ model. Each sample includes a translation in English as well. LLaMA learns to depict intensifier tokens as emojis without any instructions or training data examples.

training set does not contain emoji tokens for the RWTH-PHOENIX-14-T. Some samples are shown in Table 5. Here, it is observed that the model is mapping the intensifier tokens that exist in the intensified dataset to emojis. However, this is not a one-to-one mapping, and it is more so using the iconicity of the emoji to depict semantics that does not exist in the textual glosses.

It can be claimed that iconicity, which is normally depicted in the spatial modality during the signing, is now depicted with a different modality in a semantically rich textual form. Also, in the last sample, the generation directly includes "highly intensive," which shows that sometimes the model does not map the intensifier tokens directly to emojis. Overall, it can be qualitatively claimed that this mapping of semantics to icons via emojis is a property of LLMs fine-tuned on multiple tasks. This provides a paradigm shift in sign language processing, where including prosodically-rich tasks of signed languages can be accomplished with the help of large foundation models instead of seeing them as translation problems. Yet, new task definitions and datasets specific to signed languages should be made available for further investigations of these capabilities.

## 6 The Glossing Trade-Off

This section presents a trade-off between using textual representations of signs such as glosses or Sign-Writing that are linguistically-backed or directly using video of signers. This trade-off may not be an option most of the time, as having access to intermediary textual representations such as glosses as part of the sign corpora is not prevalent across all datasets available online. To decide whether to use glosses or videos, we can use insights from the linguistics literature and data collection experience from the RWTH-PHOENIX-14-T dataset.

In the original data collection effort as described

by Forster et al. (2012) and Stein et al. (2010), the annotations of glosses are done by a congenitally deaf person with no previous annotation experience. On average, they report that it took the annotator 24 hours to annotate 15 minutes of footage. When we compare these statistics to the fine-tuning statistics of the text-only and multimodal models, we can observe the trade-offs better. This is presented in Table 6. It can be seen that the text-only model has nearly double the performance of the multimodal, and it needs less storage space and leads to less carbon emissions, even though it takes longer to annotate.

Trade-off Statistics						
	$T_A$ (h)	$T_{FT}$ (h)	$T_I$ (s/tok)	S (GB)	Carbon Emissions (kg)	Perf. (B <sub>1</sub> )
Annotator + Text-Only	2400	8	4	0.1	0.211	22.85
Multimodal	0	8	8	50	0.240	13.62

Table 6: This table shows different statistics comparing the human annotation with the text-only model and video-based multimodal model. Carbon emissions are calculated using the US EPA’s greenhouse gas equivalencies calculator.  $T_A$ : average time for annotation,  $T_{FT}$ : average time for fine-tuning,  $T_I$ : average time for inference, S: storage space needed for data.

## 7 Conclusion

In this paper, we have prompted, fine-tuned, and compared text-only and multimodal language models for sign language processing tasks. We have provided theoretical grounding and analyzed our results from cognitive science and theoretical perspectives. From our findings, it can be claimed that text-only language models perform better than multimodal models. Moving forward, training bigger models with larger multilingual corpora is a promising next step for a broader set of novel sign language processing tasks.



## 8 Limitations

The major limitation of our work has been the computing power required to fine-tune, test, and carry out inference. Even with the smallest large language models, it becomes quickly infeasible to test multiple independent variables. Hence, our techniques have been tested on the smaller end of the large language family of models. Larger models can have higher performance gains. An additional limitation of our models is the context length. With long linguistic rules added to the prompt, certain samples of glosses made the inference lengthy. The maximum number of generated tokens has been a limiting factor of the output of models as well, which resulted in poor performance metrics. These can be alleviated with higher computing powers. Another major limitation is the dataset size and number of available tasks in sign language processing. The sign language processing community has focused on translation tasks so far, and not many other task definitions and datasets exist that can be useful for signers. This affects our benchmarking, as the only tasks we can test the generalization on are either other translation tasks or traditional NLP tasks that are non-specific to signed languages. Having diverse tasks and accompanying datasets is needed for the future of sign language processing.

## 9 Ethical Statement

We are using LLaMA2-based models for both our text-only and multimodal setups, which are trained on data acquired by Meta and are not made publicly available; even though the model itself is open-source, the pretraining dataset is not open. This leads to unaccountable biases that have been collected during the dataset formation and in the pretraining, our models may have inherent biases passed down from these pretraining setups. Our RWTH-PHOENIX-14-T dataset contains the faces of the signers, which is a piece of private information. This private information is used in accordance with the original dataset creator’s directions and privacy concerns. Furthermore, sign language processing can be a sensitive topic, especially when the community-centric approach is not taken for the design of systems. For this, we collaborate with the deaf and hard-of-hearing communities or signers in general while developing such systems as this one.

## References

- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. 2023. [Qwen technical report](#).
- Necati Cihan Camgoz, Simon Hadfield, Oscar Koller, Hermann Ney, and Richard Bowden. 2018. Neural sign language translation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Emanuela Campisi, Anita Slonimska, and Asli Özyürek. 2023. Cross-linguistic differences in the use of iconicity as a communicative strategy. In *the 8th Gesture and Speech in Interaction (GESPIN 2023)*.
- Xuanyi Chen, Junfei Hu, Falk Huettig, and Asli Özyürek. 2023a. [The effect of iconic gestures on linguistic prediction in Mandarin Chinese: a](#). [Online; accessed 14. Feb. 2024].
- Yutong Chen, Ronglai Zuo, Fangyun Wei, Yu Wu, Shujie Liu, and Brian Mak. 2023b. [Two-stream network for sign language recognition and translation](#).
- Feng Cheng, Xizi Wang, Jie Lei, David Crandall, Mohit Bansal, and Gedas Bertasius. 2023. [VindLU: A Recipe for Effective Video-and-Language Pretraining](#). [Online; accessed 15. Feb. 2024].
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. [Think you have solved question answering? try arc, the ai2 reasoning challenge](#).
- Jens Forster, Christoph Schmidt, Thomas Hoyoux, Oscar Koller, Uwe Zelle, Justus Piater, and Hermann Ney. 2012. [RWTH-PHOENIX-weather: A large vocabulary sign language recognition and translation corpus](#). In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC’12)*, pages 3785–3789, Istanbul, Turkey. European Language Resources Association (ELRA).
- Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. 2023. [A framework for few-shot language model evaluation](#).

679	Aina Garí Soler and Marianna Apidianaki. 2021. <a href="#">Let's play mono-poly: BERT can reveal words' polysemy level and partitionability into senses</a> . <i>Transactions of the Association for Computational Linguistics</i> , 9:825–844.	736
680		737
681		738
682		739
683		
684	Tao Gong, Chengqi Lyu, Shilong Zhang, Yudong Wang, Miao Zheng, Qian Zhao, Kuikun Liu, Wenwei Zhang, Ping Luo, and Kai Chen. 2023. <a href="#">Multimodal-gpt: A vision and language model for dialogue with humans</a> .	740
685		741
686		742
687		743
688	Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio César Teodoro Mendes, Allie Del Giorno, Sivakanth Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo de Rosa, Olli Saarikivi, Adil Salim, Shital Shah, Harkirat Singh Behl, Xin Wang, Sébastien Bubeck, Ronen Eldan, Adam Tauman Kalai, Yin Tat Lee, and Yuanzhi Li. 2023. <a href="#">Textbooks Are All You Need</a> . <i>arXiv</i> .	744
689		745
690		746
691		747
692		
693		748
694		749
695		750
696	Thomas Hanke, Marc Schulder, Reiner Konrad, and Elena Jahn. 2020. <a href="#">Extending the Public DGS Corpus in size and depth</a> . In <i>Proceedings of the LREC2020 9th Workshop on the Representation and Processing of Sign Languages: Sign Language Resources in the Service of the Language Community, Technological Challenges and Application Perspectives</i> , pages 75–82, Marseille, France. European Language Resources Association (ELRA).	751
697		752
698		753
699		754
700		755
701		756
702		757
703		
704		758
705	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. <a href="#">LoRA: Low-rank adaptation of large language models</a> . In <i>International Conference on Learning Representations</i> .	759
706		760
707		761
708		762
709		763
710	Hezhen Hu, Weichao Zhao, Wengang Zhou, and Houqiang Li. 2023. <a href="#">SignBERT+: Hand-Model-Aware Self-Supervised Pre-Training for Sign Language Understanding</a> . <i>IEEE Trans. Pattern Anal. Mach. Intell.</i> , 45(9):11221–11239.	764
711		765
712		766
713		767
714		768
715	Mert Inan, Yang Zhong, Sabit Hassan, Lorna Quandt, and Malihe Alikhani. 2022. <a href="#">Modeling intensification for sign language generation: A computational approach</a> . In <i>Findings of the Association for Computational Linguistics: ACL 2022</i> , pages 2897–2911, Dublin, Ireland. Association for Computational Linguistics.	769
716		770
717		771
718		772
719		
720		773
721		774
722	Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, Léo Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2024. <a href="#">Mixtral of experts</a> .	775
723		776
724		777
725		
726		778
727		779
728		
729		780
730		781
731		782
732		783
733	Zifan Jiang, Amit Moryossef, Mathias Müller, and Sarah Ebling. 2023. <a href="#">Machine translation between spoken languages and signed languages represented in SignWriting</a> . In <i>Findings of the Association for Computational Linguistics: EACL 2023</i> , pages 1706–1724, Dubrovnik, Croatia. Association for Computational Linguistics.	784
734		785
735		786
		787
		788
		789
	Dilay Z. Karadöller, David Peeters, Francie Manhardt, Asli Özyürek, and Gerardo Ortega. 2023. <a href="#">Iconicity and gesture jointly facilitate learning of L2 signs at first exposure</a> . <i>Language Learning</i> .	
	Emily Kubicek and Lorna C. Quandt. 2019. <a href="#">Sensorimotor system engagement during ASL sign perception: An EEG study in deaf signers and hearing non-signers</a> . <i>Cortex</i> , 119:457–469.	
	Emily Kubicek and Lorna C. Quandt. 2021. <a href="#">A Positive Relationship Between Sign Language Comprehension and Mental Rotation Abilities</a> . <i>J. Deaf Stud. Deaf Educ.</i> , 26(1):1–12.	
	Huije Lee, Jung-Ho Kim, Eui Jun Hwang, Jaewoo Kim, and Jong C. Park. <a href="#">Leveraging Large Language Models With Vocabulary Sharing For Sign Language Translation</a> . In <i>2023 IEEE International Conference on Acoustics, Speech, and Signal Processing Workshops (ICASSPW)</i> , pages 04–10. IEEE.	
	Dongxu Li, Junnan Li, Hung Le, Guangsen Wang, Silvio Savarese, and Steven C.H. Hoi. 2023. <a href="#">LAVIS: A one-stop library for language-vision intelligence</a> . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)</i> , pages 31–41, Toronto, Canada. Association for Computational Linguistics.	
	Bin Lin, Yang Ye, Bin Zhu, Jiaxi Cui, Munan Ning, Peng Jin, and Li Yuan. 2023. <a href="#">Video-llava: Learning united visual representation by alignment before projection</a> .	
	Chin-Yew Lin. 2004. <a href="#">ROUGE: A package for automatic evaluation of summaries</a> . In <i>Text Summarization Branches Out</i> , pages 74–81, Barcelona, Spain. Association for Computational Linguistics.	
	Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2023a. Improved baselines with visual instruction tuning.	
	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023b. <a href="#">Visual instruction tuning</a> .	
	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023c. Visual instruction tuning. In <i>NeurIPS</i> .	
	Shayne Longpre, Le Hou, Tu Vu, Albert Webson, Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V. Le, Barret Zoph, Jason Wei, and Adam Roberts. 2023. <a href="#">The flan collection: Designing data and methods for effective instruction tuning</a> .	
	Huaishao Luo, Lei Ji, Botian Shi, Haoyang Huang, Nan Duan, Tianrui Li, Jason Li, Taroon Bharti, and Ming Zhou. 2020. Univl: A unified video and language pre-training model for multimodal understanding and generation. <i>arXiv preprint arXiv:2002.06353</i> .	

790	Amit Moryossef, Zifan Jiang, Mathias Müller, Sarah Ebling, and Yoav Goldberg. 2023. <a href="#">Linguistically motivated sign language segmentation</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 12703–12724, Singapore. Association for Computational Linguistics.	846
791		847
792		848
793		849
794		850
795		
796	Amit Moryossef, Ioannis Tsochantaridis, Roei Aharoni, Sarah Ebling, and Srini Narayanan. 2021. <a href="#">Real-Time Sign Language Detection Using Human Pose Estimation</a> . In <i>Computer Vision – ECCV 2020 Workshops</i> , pages 237–248. Springer, Cham, Switzerland.	851
797		852
798		853
799		854
800		855
801	Subhabrata Mukherjee, Arindam Mitra, Ganesh Jawahar, Sahaj Agarwal, Hamid Palangi, and Ahmed Awadallah. 2023. <a href="#">Orca: Progressive learning from complex explanation traces of gpt-4</a> .	856
802		857
803		858
804		859
805	Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. <a href="#">Training language models to follow instructions with human feedback</a> .	860
806		861
807		862
808		
809		
810		
811		
812		
813	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. <a href="#">Bleu: a method for automatic evaluation of machine translation</a> . In <i>Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics</i> , pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.	863
814		864
815		865
816		866
817		867
818		868
819		869
820	AJ Piergiovanni, Isaac Noble, Dahun Kim, Michael S. Ryoo, Victor Gomes, and Anelia Angelova. 2023. <a href="#">Mirasol3b: A multimodal autoregressive model for time-aligned and contextual modalities</a> .	870
821		871
822		872
823		873
824	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. <i>OpenAI blog</i> , 1(8):9.	874
825		875
826		876
827		877
828	David Rolnick, Arun Ahuja, Jonathan Schwarz, Timothy Lillicrap, and Gregory Wayne. 2019. Experience replay for continual learning. <i>Advances in Neural Information Processing Systems</i> , 32.	878
829		879
830		880
831		881
832	Walid Saba. 2023. <a href="#">Towards ontologically grounded and language-agnostic knowledge graphs</a> . In <i>Proceedings of the 15th International Conference on Computational Semantics</i> , pages 94–98, Nancy, France. Association for Computational Linguistics.	882
833		883
834		884
835		885
836		
837	Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2019. <a href="#">Winogrande: An adversarial winograd schema challenge at scale</a> .	886
838		887
839		888
840	Anthony Sicilia and Malihe Alikhani. 2022. <a href="#">LEATHER: A framework for learning to generate human-like text in dialogue</a> . In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2022</i> , pages 30–53, Online only. Association for Computational Linguistics.	889
841		890
842		891
843		892
844		893
845		894
		895
		896
		897
		898
		899
		900
		901
		902
		903
		904



*International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 7347–7360, Online. Association for Computational Linguistics.

Kayo Yin and Jesse Read. 2020. [Better sign language translation with STMC-transformer](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 5975–5989, Barcelona, Spain (Online). International Committee on Computational Linguistics.

Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. 2023. [A survey on multimodal large language models](#).

Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. [HellaSwag: Can a machine really finish your sentence?](#) In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4791–4800, Florence, Italy. Association for Computational Linguistics.

Hang Zhang, Xin Li, and Lidong Bing. 2023. [Video-llama: An instruction-tuned audio-visual language model for video understanding](#).

Tianyi Zhang\*, Varsha Kishore\*, Felix Wu\*, Kilian Q. Weinberger, and Yoav Artzi. 2020. [Bertscore: Evaluating text generation with bert](#). In *International Conference on Learning Representations*.

Yue Zhao, Ishan Misra, Philipp Krähenbühl, and Rohit Girdhar. 2023. Learning video representations from large language models. In *CVPR*.

Bin Zhu, Bin Lin, Munan Ning, Yang Yan, Jiaxi Cui, HongFa Wang, Yatian Pang, Wenhao Jiang, Junwu Zhang, Zongwei Li, Wancai Zhang, Zhifeng Li, Wei Liu, and Li Yuan. 2024. [Languagebind: Extending video-language pretraining to n-modality by language-based semantic alignment](#).

## A Hyperparameters & Training Implementation Details

We trained all of the models on an Apple MacBook Pro with an M3 Max chip. Libraries used were PyTorch, Huggingface TRL, Transformers, Datasets, Evaluate, and W&B. The hyperparameters for the LLaMA models are: learning rate of 1e-3, lr scheduler type: "reduce lr on the plateau", per device training batch size of 2, number of epochs of 5, and weight decay of 0.01, and maximum sequence length of 300 tokens. LoRA configuration for the LLaMA model is: rank of 8, LoRA alpha of 32, and LoRA dropout of 0.1. For the LLaVA model: mm projector learning rate of 2e-5, one epoch, batch size of 2, learning rate of 5e-5, linear lr scheduler type, maximum sequence length of 2048. LoRA configuration for LLaVA model: LoRA rank: 128, and LoRA alpha: 256.

## B All Prompt Types

Here we present all the prompt types that have been used in the experiments:

- **zero-shot prompt:** This is a sentence in German Sign Language glosses: <glosses>. You MUST translate these to spoken German. You MUST give the answer directly without any other text.
- **rule-based prompt:** "Instructions Here are some basic rules of German GLOSSES: 1) German signs correspond to meanings not to words. 2) Some GLOSSES are formed from more than one German word. In this case the words are joined by a hyphen. The hyphen indicates one single sign that is labeled with two or more German words. 3) Glosses combined with a plus sign are two separate signs that are joined together to make what appears to be a single sign 4) In DGS, some signs are repeated for specific meaning. for instance LEARN + LEARN changes the sign from the VERB "To Learn" to the NOUN "Learning." 5) Words that are to be Fingerspelled are indicated in one of two ways: - Separated by hyphens between each Fingerspelled letter: G-L-A-D-Y-S - Preceded by the initials FS in parenthesis: (fs) GLADYS. Task You MUST translate <glosses> of DGS to German without using any special characters, according to these rules."
- **notation-based prompt:** "Instruction Below is a list of common symbols used in the writing of DGS Glosses: - The Crosshatch: This symbol indicates a loan sign, a sign originating from the fingerspelling of an English word. - Parentheses: ( ) Additional information about the production of a sign is can added to the written gloss between a set of parentheses. Such information can be abbreviated as in (2h)DO++, or it may appear as German instructions to add information to a sign: GIVE (left), or to a Classifier CL:1 (man hurries past). - CL: The abbreviation CL: indicates a classifier. The information following the colon indicates the hand shape and number of hands. - The Umlaut (two dots above a given hand shape) ( indicate the bending of the fingers of that hand. The 3 (called the "bent three") is the hand shape used in the



1007 sign “INSECT”. This technique is only used  
1008 in reference to a specific handshape such as a  
1009 classifier.

1010 Task You MUST translate <glosses> to Ger-  
1011 man according to these symbols."

1012 • **one-shot prompt:** "Example ""Here’s a sam-  
1013 ple DGS gloss: “ORT REGEN DURCH  
1014 REGEN KOENNEN UEBERSCHWEM-  
1015 MUNG KOENNEN” which translates to  
1016 ""mancherorts regnet es auch länger und  
1017 ergiebig auch lokale überschwemmungen sind  
1018 wieder möglich"" in German

1019 Task You MUST translate <glosses> to Ger-  
1020 man according to this example. "