



AFRL-AFOSR-UK-TR-2023-0034

Referential Grounding in Multimodal Machine Translation

Specia, Lucia
IMPERIAL COLLEGE OF SCIENCE TECHNOLOGY & MEDICINE
EXHIBITION RD
LONDON, ,
GB

12/22/2022
Final Technical Report

DISTRIBUTION A: Distribution approved for public release.

Air Force Research Laboratory
Air Force Office of Scientific Research
European Office of Aerospace Research and Development
Unit 4515 Box 14, APO AE 09421

REPORT DOCUMENTATION PAGE

PLEASE DO NOT RETURN YOUR FORM TO THE ABOVE ORGANIZATION.

1. REPORT DATE 20221222	2. REPORT TYPE Final	3. DATES COVERED	
		START DATE 20191115	END DATE 20220514
4. TITLE AND SUBTITLE Referential Grounding in Multimodal Machine Translation			
5a. CONTRACT NUMBER	5b. GRANT NUMBER FA8655-20-1-7006	5c. PROGRAM ELEMENT NUMBER 61102F	
5d. PROJECT NUMBER	5e. TASK NUMBER	5f. WORK UNIT NUMBER	
6. AUTHOR(S) Lucia Specia			
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) IMPERIAL COLLEGE OF SCIENCE TECHNOLOGY & MEDICINE EXHIBITION RD LONDON GB			8. PERFORMING ORGANIZATION REPORT NUMBER
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES) EOARD UNIT 4515 APO AE 09421-4515		10. SPONSOR/MONITOR'S ACRONYM(S) AFRL/AFOSR IOE	11. SPONSOR/MONITOR'S REPORT NUMBER(S) AFRL-AFOSR-UK-TR-2023-0034
12. DISTRIBUTION/AVAILABILITY STATEMENT A Distribution Unlimited: PB Public Release			
13. SUPPLEMENTARY NOTES			
14. ABSTRACT <p>This project aimed to advance the state of the art in multimodal machine translation (MMT). MMT is an area where a text in the source language is supplemented by visual information (images or video) to be used as additional context to better understand and translate the text into a target language. The core of the advances proposed are on referential grounding, i.e., on guiding the alignment between image regions and source (and/or target) words such that the visual context can be more useful for translation.</p> <p>Work done during the project in covered the following directions:</p> <ol style="list-style-type: none"> 1. Improving supervised attention mechanisms to map source or target words to image regions, addressing both attention at encoding time (i.e. learning alignments between source words and objects in the image) and at decoding time (i.e. learning alignments between target words and objects in the image), as well as improving the underlying multimodal neural machine translation architectures and fusion strategies to use such information and exploring more recent and better types of visual features. 2. Leveraging information from multiple vision-and-language tasks and datasets to improve multilingual grounding. 3. Creating resources to facilitate work on referential grounding. 			
15. SUBJECT TERMS			
16. SECURITY CLASSIFICATION OF:		17. LIMITATION OF ABSTRACT	18. NUMBER OF PAGES
a. REPORT U	b. ABSTRACT U	c. THIS PAGE U	SAR 2
19a. NAME OF RESPONSIBLE PERSON NANDINI IYER			19b. PHONE NUMBER (Include area code) 314-235-6161

Referential grounding in multimodal machine translation

FA8655-20-1-7006

Prof Lucia Specia, Imperial College London

Project start 15th Nov 2019

Current project ending 14 May 2022

Reporting period: Jan 2022 – May 2022

Progress report

This project aimed to advance the state of the art in multimodal machine translation (MMT). MMT is an area where a text in the source language is supplemented by visual information (images or video) to be used as additional context to better understand and translate the text into a target language. The core of the advances proposed are on **referential grounding**, i.e., on guiding the alignment between image regions and source (and/or target) words such that the visual context can be more useful for translation.

Work done during the project in covered the following directions:

1. Improving supervised attention mechanisms to map source or target words to image regions, addressing both attention at encoding time (i.e. learning alignments between source words and objects in the image) and at decoding time (i.e. learning alignments between target words and objects in the image), as well as improving the underlying multimodal neural machine translation architectures and fusion strategies to use such information and exploring more recent and better types of visual features.
2. Leveraging information from multiple vision-and-language tasks and datasets to improve multilingual grounding.
3. Creating resources to facilitate work on referential grounding.

This report focuses in the last 4 months of the project, and covers further work towards **direction 1**, namely we proposed the first approach to **Simultaneous Video Translation**, i.e. real-time translation or interpreting where a translation needs to be generated for incomplete source sentences and where a video is available as additional context. An example of application is the translation of the audio stream in live broadcasting such as news. Different from our previous work, where a single image was available as static visual information for each text segment to translate, in our more recent work a video containing multiple pieces of visual information is available for each text segment. This brings many challenges to MMT, including decisions around how to process the video (frame sampling approach, video encoding approach) as well as how to combine multiple pieces of visual (frames or even image regions in frames) and textual (source and/or target words) information. The latter can be seen as a form of **referential grounding** between frames and text sub-segments.

Using videos as visual information for MMT is appealing as it offers richer visual context, especially for longer text segments. It also opens new avenues for research on referential grounding: for a correctly grounded translation, the model needs to identify correspondences between specific video frames or parts of frames that are relevant for the words seen so far in the incomplete textual input that is incrementally made available. A draft summary of the work done is attached to this report (paper to be submitted). In this paper, we use a dataset of videos where people describe their apartments for rental to train and evaluate our simultaneous video translation models.

Publications relevant to reporting period

Wang, Y., Miao, Y., Specia, Lucia. (2022). Translate in the Scenes: Multi-modal Simultaneous Machine Translation (draft of paper to be submitted, attached).

Other cited publications (previous report)

Haralampieva, V., Caglayan, O., Specia, L. (2022). Supervised Visual Attention for Simultaneous Multimodal Machine Translation. *Journal of Artificial Intelligence Research (JAIR)*, 74:1059-1089.

Wang, J., P.S., Figueiredo, J.M., Specia, L. (2022). MultiSubs: A Large-scale Multimodal and Multilingual Dataset. *Proceedings of the 13th Conference on Language Resources and Evaluation (LREC 2022)*, p. 6776–6785

Mitzalis, F., Caglayan, O., Madhyastha, P.S., Specia, L. (2021). BERTGEN: Multi-task Generation through BERT. *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics (ACL 2021)*, p. 6440–6455.

Caglayan, O., Ive, J., Haralampieva, V., Madhyastha, P.S., Barrault, L., & Specia, L. (2020). Simultaneous Machine Translation with Visual Context. *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP 2020)*, p. 2350-2361.

Specia, L., Wang, J., Lee, S.J., Ostapenko, A., Madhyastha, P. (2021). Read, Spot and Translate. *Machine Translation Journal*, 35, 145–165, Springer.

Translate in the Scenes: Multi-modal Simultaneous Machine Translation

Anonymous ACL submission

Abstract

We propose the task of simultaneous video translation, which aims to explore multimodal information to help produce machine translations for videos in real-time. We posit that visual information can help the language channel achieve a good trade-off between translation quality and latency in this very challenging task. To enable research in this area, we created the *ApartmentTour Simultaneous Video Translation Dataset*, which contains 33495 parallel subtitles for 82 hours of videos describing apartments for rental. We then build on previous approaches that address simultaneous translation with multimodal information from static images using RNN and Transformer architectures and show that these lead to considerable improvements compared to text-only translation models. Finally, we propose a frame selection approach to efficiently deal with the problem of processing multiple frames for one sentence.

1 Introduction

Given the recent breakthroughs in machine translation (MT) with neural approaches (Bahdanau et al., 2014; Vaswani et al., 2017) and the impressive improvements in translation quality observed in the last 6-7 years, the community has turned to a number of practical applications, including *simultaneous* MT. This task simulates the job of a human interpreter, who translates simultaneously as source words are produced, in real-time. This has a number of applications such as international summits, conferences and lectures. Compared to the traditional *consecutive* MT, where the entire sentence is available as the input, simultaneous MT (SiMT) is much more challenging, given that the incomplete source language content usually reduces the quality of the translation significantly. In addition, the model is required to consider not only the quality of the translation but also the latency of

the reader, i.e. the time the reader has to wait to receive the translation before it is produced. A good trade-off must be obtained between the quality and the latency.

Cho and Esipova (2016) first explored the possibility of SiMT with a heuristic decoding algorithm. After that, several methods were proposed that focus on finding the best timing to perform the translation action through supervised or unsupervised neural networks (Gu et al., 2016; Zheng et al., 2019), or fixing the latency to optimise for translation quality (Ma et al., 2018).

Hoping to utilise additional modalities to compensate for the incomplete semantics during the simultaneous translation, previous work has also proposed simultaneous image translation (Caglayan et al., 2020; Imankulova et al., 2020). Previous work has shown that the visual information from an image can help the model anticipate the missing source context as well as disambiguate source terms, overall helping predict the next source word in the pipeline of simultaneous translation. However, existing work has been limited to static visual input (one image per source sentence) and to datasets where the text is a description of the image (i.e. the Multi30K dataset (Elliott et al., 2016)), where the visual information is clearly well aligned to the source sentence, and can be seen as an alternative representation of the information given by the complete source sentence. When made available to the model at the first timestep of encoding/decoding, the complete visual representation information is therefore bound to provide helpful information to complement the incomplete source context. No previous work has been done in the more realistic, more practical but also more challenging setting addressed in this paper: that of SiMT from videos.

In the simultaneous video translation task, we aim to simultaneously generate the translation using as input the incomplete source text and the

082 corresponding *partial* visual information, i.e. only
083 the frames visible thus far. The long sequences of
084 the videos bring several complexities to the task.
085 First, there are multiple frames for each sentence
086 and the frames will appear gradually as a stream,
087 i.e. both video and text modalities are simultane-
088 ous, requiring different architectures. Second, the
089 text in the subtitles of the videos are not always
090 aligned to the visual content, so there is a weak cor-
091 respondence between the text and the video which
092 can introduce noise in the process and reduce the
093 translation quality.

094 To address these challenges, we start by col-
095 lecting a new dataset – *ApartmentTour Simulta-*
096 *neous Video Translation Dataset*, which includes
097 534 apartment tour videos, accompanied by En-
098 glish subtitles with timestamps and parallel Chi-
099 nese translations. We then build the Multi-modal
100 Simultaneous Machine Translation model which
101 can process multiple frames. We exploit different
102 visual features to represent the videos, as well as
103 different fusion methods to fuse two modalities.
104 To tackle the weak semantic equivalence between
105 the content in the text and the video, we propose
106 two masking strategies to mask frames which we
107 expect to contain information that is not related
108 to the text. The experimental results show that
109 our approaches to SiMT from videos perform bet-
110 ter than their monomodal counterparts without the
111 video. The masking strategies have proved to filter
112 out some irrelevant frames and improve translation
113 quality.

114 In Section 2, we provide the background of the
115 monomodal and multimodal simultaneous machine
116 translation. Section 3 outlines the processing of
117 our ApartmentTour Simultaneous Video Translation
118 Dataset. Our proposed Multi-modal Simultaneous
119 Machine Translation model is described in Section
120 4. Section 5 provides the results of our models
121 and the qualitative analysis. Section 6 is for the
122 conclusion.

123 2 Related Work

124 2.1 Simultaneous machine translation

125 Simultaneous machine translation is challenging
126 as machines have to translate the sentences simul-
127 taneously as the words appear. Partial sentences
128 will inevitably impair the quality of the translation
129 compared to typical machine translation task which
130 translates the entire sentences. In addition, both the
131 quality and the latency of the translation have to be

considered.

132
133 In the earlier years, segmentation-based meth-
134 ods have been proposed. These methods (Fügen
135 et al., 2007; Bangalore et al., 2012; Sridhar et al.,
136 2013; Oda et al., 2014) aim to segment a whole
137 source sentence, translate each part independently,
138 and concatenate all parts to form a whole target
139 sentence to simulate the simultaneous translation
140 scenario.

141 Cho and Esipova (2016) start to introduce the
142 neural network into simultaneous machine transla-
143 tion to exploit the dependency between different
144 parts of the translation. They use an rnn-based ma-
145 chine translation model trained on full sentences
146 and connected it with a greedy decoding part to
147 do the simultaneous machine translation. They de-
148 fined two criteria, wait-if-worse and wait-if-diff
149 to heuristically decide the timing to write a target
150 word. However, this method is based on predefined
151 rules, which cannot greatly exploit the information
152 of hidden states of the recurrent neural networks.

153 Gu et al. (2016) propose a reinforcement
154 learning-based method to train an agent to gen-
155 erate the READ/WRITE action sequences. The
156 agent is trained using the reward function which
157 considers both the quality and the latency.

158 As reinforcement learning-based methods are
159 not stable and difficult to train, Ma et al. (2018)
160 propose the wait- k policy, which forces the trans-
161 lation (the first WRITE action) to start after
162 reading k source words, and then performs the
163 READ/WRITE actions alternatively. The policy
164 enables the model to have the ability to anticipate
165 the future words.

166 However, as the latency for the wait- k method
167 is predefined, it cannot learn the best timing for
168 READ/WRITE actions for different data. In addi-
169 tion, the anticipation of the future words can often
170 be incorrect. Zheng et al. (2019) propose a method
171 that combines the adaptivity of the reinforcement
172 learning-based method and the easy training prop-
173 erties of the predefined latency methods. They design
174 a supervised action sequence generation network,
175 which takes in the ground truth action sequences
176 generated using a heuristic algorithm.

177 2.2 Multimodal Simultaneous Translation

178 Multimodal simultaneous translation is a newly pro-
179 posed task that seeks to introduce another modality
180 to compensate for the miss semantic content.

181 [Imankulova et al. \(2020\)](#) propose the RNN-based
 182 Multimodal Simultaneous Neural Machine Trans-
 183 lation model (MSNMT) by combining the wait- k
 184 simultaneous translation model ([Ma et al., 2018](#))
 185 and a multimodal model ([Libovický and Helcl,](#)
 186 [2017](#)). A hierarchical attention mechanism ([Li-](#)
 187 [bovický and Helcl, 2017](#)) is adopted to fuse the text
 188 modality and image modality in the decoder side.
 189 Based on the two-GRU-layer architecture described
 190 in ([Caglayan et al., 2016](#)), [Caglayan et al. \(2020\)](#)
 191 introduce object classification features extracted
 192 from ResNet50 ([He et al., 2016](#)) and object detec-
 193 tion features extracted from the “bottom-up-top-
 194 down(BUTD)” ([Anderson et al., 2018](#)) extractor.
 195 Also, they exploit encoder attention and decoder
 196 attention mechanisms which can only attend to the
 197 available hidden states to simulate the simultane-
 198 ous translation setting. [Ive et al. \(2021\)](#) propose
 199 the first RL-based multimodal simultaneous trans-
 200 lation model to introduce the visual cues to the
 201 agent and environment. Both the global image clas-
 202 sification features ([He et al., 2016](#)) and the local
 203 visual concept features are explored. The research
 204 on multimodal simultaneous translation is mostly
 205 based on static images but not videos.

206 3 ApartmentTour Dataset

207 We pick up 534 apartment tour videos from the
 208 ApartmentTour dataset ([Zhong et al., 2020](#)) and
 209 extracted WAV files using FFmpeg¹. Next, the
 210 English transcription (splitted by paragraphs) with
 211 punctuation and the corresponding timestamps for
 212 each video are extracted using Microsoft Azure
 213 Speech Translation service² by uploading the WAV
 214 files (as shown in Figure 1). The timestamp for
 215 each paragraph includes the offset and the duration
 216 in the unit of 100 nanoseconds (1 nanosecond =
 217 1×10^{-9} second). We translated the English tran-
 218 scription to Chinese using the Youdao Translation
 219 API³⁴. As most paragraphs of transcription gener-
 220 ated by the Microsoft Azure Speech Translation
 221 service are very long and contain many sentences,
 222 we split each paragraph by the period mark to get
 223 several sentences. In addition, since the subtitles

¹<http://ffmpeg.org/>

²<https://azure.microsoft.com/en-us/services/cognitive-services/speech-translation/>

³<https://ai.youdao.com/product-fanyi-text.s>

⁴We also created German translation using Azure, which we did not use in our experiments but will release with the data.

	Train	Val	Test
Sentence / Video clip	20974	6521	6000
English words	426521	129649	120872
English vocabulary	16691 in total		
Chinese words	387089	117270	109500
Chinese vocabulary	19928 in total		

Table 1: Statistics about the ApartmentTour Simultaneous Video Translation Dataset

224 are converted from audio, they contain many hesita-
 225 tions and other symbols that are not very meaning-
 226 ful and sometimes correspond to a whole sentence,
 227 such as “So.” and “Yeah.”. To focus on more mean-
 228 ingful and hard to translate sentences, we remove
 229 those sentences whose word lengths are less than
 230 eight words. As a result, our final dataset includes
 231 20,974 sentences for training, 6,521 sentences for
 232 validation, and 6,000 sentences for testing. The
 233 average length of the sentences is around 20 words.
 234 Some statistics of the dataset are shown in Table 1.

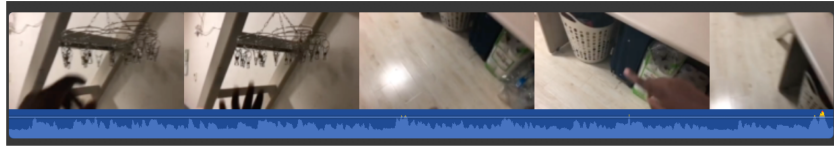
235 Since using all frames from the video would
 236 make the process very inefficient, we resort to
 237 frame sampling, a common practice in video pro-
 238 cessing ([Lei et al., 2021](#)). For convenience, we sam-
 239 ple one frame every five words (as shown in Figure
 240 2). The frame at timestep 0 is always sampled
 241 as the first frame for the corresponding sentence.
 242 As we only get sentence-level timestamps using
 243 Microsoft Azure Speech Translation service, we
 244 roughly compute timestamp for each frame using
 245 offset + $(\frac{i \times 5}{\text{word length}}) \times \text{duration}$ for the i -th frame
 246 (starting from 0).

247 4 Models

248 As shown in Figure 3, our model consists of three
 249 parts: language channel, video channel and the
 250 translation networks.

251 4.1 Language channel

252 The architecture is the same as that in [Vaswani](#)
 253 [et al. \(2017\)](#), but using a pre-norm variant. In the
 254 simultaneous translation task, we can only use the
 255 unidirectional information, so we will generate a
 256 lookahead mask to make each word in the sen-
 257 tence not attend to the words behind itself in the
 258 self-attention layers. For simultaneous translation,
 259 we adopt the wait- k policy proposed by [Ma et al.](#)
 260 [\(2018\)](#). Though the authors report better results
 261 using training time wait- k , that is, training and test-
 262 ing on partial sentences, we find that testing time
 263 wait- k works better on our dataset (this was also re-



Timestamp	English	Chinese
[2526800000, 226800000]	Only a washer, so I hang up my clothes right here to dry on this little clothespin thing. And then I also kind of lay things out on here as well. And then I also use it for some storage like I keep my water and my paper towels and toilet paper. I've got like a carry on bag. I've got my dirty clothes hamper. I've got a bunch of stuff guys.	只有洗衣机，所以我把衣服挂在这里，用这个小晾衣架晾干。然后我也会把东西放在这里。我也用它来储存水、纸巾和厕纸。我有个随身包。我有脏衣服篮。我有一堆东西，伙计们。
offset duration		

Figure 1: An illustration of the generation of the source and translation text. The WAV file of the video was extracted and uploaded to the Microsoft Azure speech translation service to generate the English subtitles, Chinese translation and timestamps. The timestamp for each paragraph includes the offset and the duration in the unit of 100 nanoseconds ($1 \text{ nanosecond} = 1 \times 10^{-9} \text{ second}$)



Figure 2: An illustration of the frame sampling. One frame is sampled for every five words. The frame at timestep 0 is also sampled.

then padded to get a feature matrix of shape [176, 768] for each sentence (as the maximum number of frames for one sentence is 22, the maximum number of text features is $22 * 8 = 176$).

ported to be the case in (Caglayan et al., 2020)). For the language channel, we train using the full sentences and test using partial sentences. To be more specific, only the source text feature of $t + k$ source words at test time can be seen when generating the t -th target word. We also apply the GRU-based simultaneous translation method (Caglayan et al., 2020) in the experiment section.

4.2 Video channel

In what follows we describe what variants of visual features we extract from the sampled frames, as well as the strategy we propose to mask out irrelevant frames.

Visual features We adopt the Vision-and-Language Transformer (ViLT) (Kim et al., 2021) to create a masked language model for each frame sampled from the video and get the text feature for 8 words (nouns and adjectives describing the image) from one frame, and the shape is [8, 768] for each frame. The features for the same sentence but different frames are stacked together and

Masking strategies As with any video dataset, our data will contain many frames that are irrelevant to what the Youtubers are talking about, and thus would represent noise rather than helpful information to improve SiMT. Here we propose two strategies to mask the irrelevant frames and reduce the impact of noise:

In the apartment tour videos, there are quite a lot of clips where the Youtubers face the camera and talk about some background knowledge of filming the videos. The frames sampled from these clips are always irrelevant to the corresponding caption and contain very little information. By checking the dataset, we find that there are a large number of frames in which the Youtuber’s face takes up most of the place, making up more than 5% of the frame. We here roughly define these frames as “irrelevant frames”. As shown in Figure 4, the left frame is our irrelevant frame and the middle frame is considered as a relevant frame as there is no face appearing in the frame. The Youtuber’s face also appears in the right frame, but as it does not take up more than 5% of the frame, we still consider it a relevant frame. Haar Cascades in OpenCV⁵ is used to do face detection in every frame and mask the frame

⁵https://docs.opencv.org/3.4/d2/d99/tutorial_js_face_detection.html

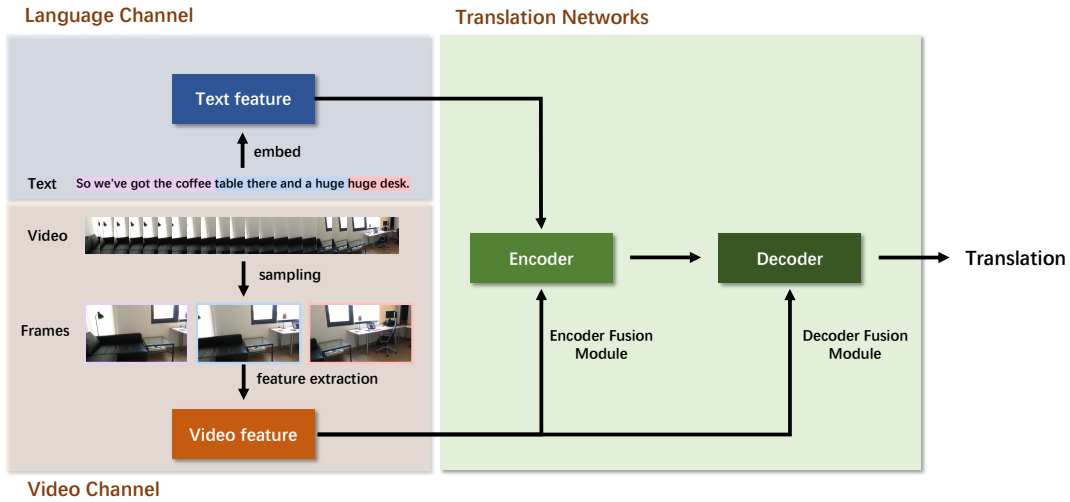


Figure 3: Multi-modal Simultaneous Machine Translation model which consists of language channel, video channel and the translation networks. Language channel is for text processing and video channel is for visual feature extraction. In the translation networks, two modalities are fused in the decoder side or the encoder side.



Figure 4: An example of “irrelevant” and “relevant” frames.

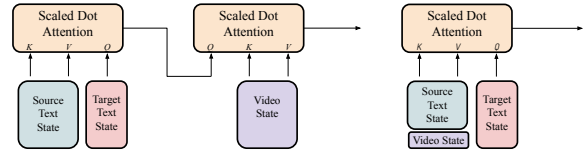


Figure 5: An illustration of the serial attention layer (left) and the concatenation attention layer (right)

314 if the face region takes up more than 5%, so the
 315 model will not utilise the feature of this frame.

316 We also use the Universal Sentence Encoder (Cer
 317 et al., 2018) to embed the subtitle (the source text)
 318 and the caption of the corresponding frame, then
 319 get the semantic textual similarity by calculating
 320 the inner product of the encodings. If the similar-
 321 ity between the subtitle and the frame’s caption is
 322 lower than a threshold, this frame will be masked
 323 as it is not so related to the source text.

324 4.3 Fusion in the translation networks

325 In what follows, we describe our fusion strategies
 326 to build a multimodal model from the language and
 327 video channel representations.

328 4.3.1 Decoder Fusion

329 Two modalities will be fused in the cross attention
 330 layers in the decoder. We try two methods for the
 331 cross-attention layer: serial attention and concate-
 332 nation attention (shown in Figure 5).

333 The serial attention layer includes two scaled
 334 dot attention. The first scaled dot

335 layer will attend the query (the target state) to the
 336 source text state (the key and value are both the
 337 text state). Then the second scaled dot attention
 338 layer will attend the output from the first scaled
 339 dot attention layer to the visual state (the query is
 340 the output from the last attention layer, the key and
 341 value are both the visual state). The output from the
 342 second attention layer will be returned. As for the
 343 concatenation attention layer, the source text state
 344 and the visual state are concatenated together (set
 345 as key and value). The attention layer attends the
 346 target state (set as the query) to the concatenated
 347 state to get the context vector and return it.

348 Also, we do some modifications to do the simul-
 349 taneous machine translation for the Transformer
 350 model. The number of the frames that are available
 351 at timestep t is computed by $N = \text{math.ceil}((t +$
 352 $\delta)/5 + 0.0001)$, and the index for the latest frame
 353 is computed by $i = N - 1$ ($\delta \leq k$ to guarantee
 354 that the model will not access to the future image
 355 information). The frames[$i - 1$] and frames[i]
 356 are the visible frames when translating the t -th
 357 word (when i is 0, only the frames[i] is visible).
 358 During the training, all the source words together with the

		BLEU	METEOR	AL
w1	Baseline	20.982	25.534	1.544
	MSMT	21.541	25.997	1.552
w2	Baseline	22.849	26.924	2.224
	MSMT	23.123	27.047	2.180
w3	Baseline	24.708	27.853	2.904
	MSMT	25.069	28.142	2.896

Table 2: The results of the baseline Monomodal Transformer model and our proposed Multi-modal Simultaneous Machine Translation model (MSMT) from wait-1 to wait-3.

frames[$i - 1$] and frames[i] will be visible. Only the $t + k$ words together with the frames[$i - 1$] and frames[i] will be visible when generating the t -th word at test time.

4.3.2 Encoder Fusion

As it is an encoder side fusion method, the visual feature will be introduced earlier. Merged attention (Chen et al., 2020; Hendricks et al., 2021) is adopted to concatenate the source text feature and the visual feature together and pass the merged feature through the stacked self-attention layers and feed-forward layers. The scaled dot attention layer is used for the decoder to attend the target state to the hidden state got from the encoder (set the target state as query, the hidden state as key and value).

5 Experiments

5.1 Multimodal Simultaneous Machine Translation

In table 2, we compare our Multi-modal Simultaneous Machine Translation model (decoder side fusion + serial attention layer + ViLT feature + encoder mask) to the baseline. When doing simultaneous translation from wait-1 to wait-3, introducing the ViLT feature with encoder mask can improve the translation quality without sacrificing the simultaneity of the translation (the Average Lagging of the multimodal model is even a little bit smaller than the one of the monomodal model for wait-2 and wait-3 tasks).

5.2 Ablation study

5.2.1 Architecture

As visual information can play a more important role when k is smaller (Caglayan et al., 2020), here we do an early feature and architecture search on wait-1 to decide the best practice. The feature is fixed as the ResNet50 feature to search for the best architecture. The results are shown in Table 3.

		BLEU	
Baseline	RNN	16.073	
	Transformer	20.982	
MSMT		No mask	OpenCV mask
	RNN Concat (decoder fusion)	18.294	-
	Transformer Merged (encoder fusion)	20.378	21.086
	Transformer Concat (decoder fusion)	20.939	20.678
	Transformer Serial (decoder fusion)	21.276	21.239

Table 3: Architecture search. We fix the feature and experiment across different architectures on wait-1 to find the best architecture.

The Transformer baseline model is a lot better than the RNN baseline, improving the BLEU score by around 5 points. The transformer can greatly learn the contextual information using self-attention layers.

From Table 3, when using decoder side fusion with serial attention layer, the Transformer-based model can get the best BLEU scores (the bold numbers in Table 3). It shows that an extra scaled dot layer attending the text hidden state to the visual hidden state is more effective than directly concatenating these two states together. Though better improvement in the image-guided translation has been shown using the encoder side fusion, it is not the case for our task here. There is only one image for one sentence for the image-guided translation task, and the image is manually picked to guarantee that it is very relevant to the text. But in our task, there are multiple images for one sentence, and the images are extracted without the manually further process. So encoder side fusion will introduce a lot of noise when generating the hidden state. When the masking strategy is applied to mask the irrelevant frames, there is an improvement in the encoder side fusion model compared to the model using the raw features, which can somehow confirm that irrelevant frames will deteriorate the translation quality.

5.2.2 Feature & Masking strategy

Table 4 shows that, for wait-1 and wait-2 simultaneous translation tasks, introducing different kinds of no mask image features can all bring improvement to the translation quality. However, using OpenCV mask will lead to worse results for the ResNet feature and the image captioning feature than those using raw image features. Also, the encoder mask does not improve the translation qual-

		BLEU		
Baseline		20.964		
		No mask	OpenCV mask	Encoder mask
w1	MSMT w/ ResNet	21.276	21.239	21.302
	MSMT w/ Image cap.	21.404	21.035	21.097
	MSMT w/ ViLT	21.408	21.306	21.541
Baseline		22.849		
		No mask	OpenCV mask	Encoder mask
w2	MSMT w/ ResNet	23.122	22.992	23.134
	MSMT w/ Image cap.	23.285	22.760	23.010
	MSMT w/ ViLT	23.085	23.197	23.123
Baseline		24.708		
		No mask	OpenCV mask	Encoder mask
w3	MSMT w/ ResNet	24.439	24.479	24.520
	MSMT w/ Image cap.	24.376	24.465	24.866
	MSMT w/ ViLT	24.698	24.960	25.069

Table 4: The results of the baseline Transformer model and the multimodal Transformer model with the serial attention layer and different features from wait-1 to wait-3.

ity too much for wait-1 and wait-2. But when it comes to wait-3, the masking strategy is important. From Table 4, using features without masking on the wait-3 task makes the translation quality even worse than the baseline model’s translation quality. For example, we get 24.376 using the image captioning feature, which is worse than the baseline 24.708. With masking strategy (especially the encoder masking strategy), we can get the translations of higher quality. That is because, as k increases, the language channel can get more source words which can directly improve the translation quality while most kinds of image features are comparatively noisy. When the monomodal baseline can get quite a great result when k is large (the language channel is more powerful), the introduction of the noisy video modality will affect more and impair the translation quality. However, when k is smaller, the importance of exposing to more video information to help anticipate overweighs the importance of reducing noise.

The feature and the masking strategy selected will both affect the translation quality. Encoder masking strategy is better than the OpenCV mask strategy as it can rule out more frames irrelevant to the text. Taking all three tasks (wait-1 to wait-3) into consideration, using ViLT feature with the encoder mask is the best practice.

5.3 Qualitative Analysis

5.3.1 Disambiguation

As we all know, polysemy is a difficulty in machine translation tasks. It will be challenging for the translation model to determine the correct meaning in different sentences. Using the image or the video

can help disambiguate the polysemous words by providing additional cues.



Source: I just have glasses that I literally never use.
 Target: 我(我)只是(just)有(have)我(我)从不(literally never)使用(use)的杯子(drinking glasses)。
 Baseline: 我(我)只(just)有(have)一个眼镜(eyeglasses)。我(我)从来不(literally never)用(use)。
 MSMT: 我(我)只(just)有(have)杯子(drinking glasses)。我(我)从来不(literally never)用(use)。

Figure 6: Multimodality can help disambiguate the polysemous word. In this example, the multimodal model can utilise the visual information to correctly translate “glasses” to “drinking glasses” in Chinese while the monomodal one translates it to “eyeglasses”.

The word “glasses” has two common definitions(GLA): (1) a device for correcting vision, which is also known as “eyeglasses”. (2) the plural form of the container for drinks made of glass. From Figure 6, when using the monomodal translation model, the word “glasses” in the sentence “I just have glasses that I literally never use.” was incorrectly translated as “眼镜” (the first meaning we specified), but when using the MSMT, the corresponding frame provides the information about the “glasses” which is a container, and the model successfully disambiguate the word and generate correct translation “杯子” (the second meaning we specified).

5.3.2 Extra visual information to anticipate

This aspect is more related to the simultaneous translation task. As we adopt the wait- k policy to control the latency, some WRITE actions will be executed before getting enough context from the text, leading to the inaccurate translation. However, an extra modality can provide the context required to translate some specific words.

Compound word As this task is a simultaneous task, every step is simultaneous, including the speech recognition. The words in the source language are generated one after another, so there is always the case that the translation model will only get one part of a compound word and translate this part directly. As shown in Figure 7, when using the wait-1 policy, the monomodal model fails to translate the compound word as a whole but only generate the translation for “gas” (“气体”), while MSMT successfully “anticipate” the following word using the extra information from the frame to generate



Source: I used to have like a **gas stove (燃气灶)**.
 Baseline (w1): 我 以前 以前 已经 有 一个 气体 (gas).
 MSMT (w1): 我 以前 以前 有 过 一个 **燃气灶 (gas stove)**.
 Target: (consecutive): 我 以前 有 一个 **燃气灶 (gas stove)**.

Figure 7: Multimodality can help predict the following word during the simultaneous translation. In this example, MSMT can correctly predict that the word “gas” would be followed by the word “stove” according to the visual information, and generate the correct translation for the compound word “gas stove”, while the monomodal model directly translates the “gas”.



Source: By the way , check out this crazy lamp .
 Baseline (w2): 顺便 说 一下 , 看看 这个 疯狂 的 灯 。
 MSMT (w2): 顺便 说 一下 , 看看 这 **盏** 疯狂 的 灯 。
 Baseline (w3): 顺便 说 一下 , 看看 这个 疯狂 的 灯 。
 MSMT (w3): 顺便 说 一下 , 看看 这 **盏** 疯狂 的 灯 。
 Target: (consecutive): 顺便 说 一 句 , 看看 这 **盏** 疯狂 的 灯 。

Figure 8: Multimodality can help generate the accurate measure word. In this example, using the extra information from the video frame, MSMT successfully generates the correct measure word/classifier for the word “lamp”.

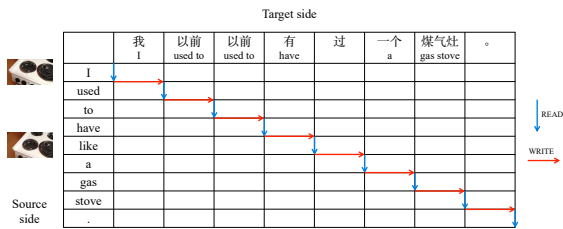


Figure 9: The process of the wait-1 multimodal simultaneous translation.

the correct translation “燃气灶 (gas stove)”. From Figure 9, MSMT successfully anticipate and generate the correct translation “燃气灶 (gas stove)” before reading the word “stove”.

Measure word This point is a more language-specific one. In Chinese, a measure word, also known as a classifier, must be used together with the numerals or the demonstrative before the nouns (Li and Thompson, 1989). Measure words should be chosen according to the noun. For example, for the noun “pencil”, the commonly used measure word is “支” and for the noun “table” it should be “张”. In a consecutive translation task, it will be quite easy to figure out the correct measure word using the context. However, in the simultaneous translation task, it will be more difficult due to the incomplete semantics. For example, as shown in 8, when using wait-2, when translating something like “this crazy lamp”, the model can only see the English word “this” without knowing the noun after it. Model without visual information might directly translate the English word “this” to “这个”, where “个” is not the appropriate mea-

sure word for “lamp”. However, with the visual information, MSMT can get the accurate measure word “盏”. For wait-3 model, we can also see that the visual information help accurately generate the measure word.

6 Conclusion

In this paper, we have explored the newly proposed Simultaneous Video Translation task. We have first created the *ApartmentTour Simultaneous Video Translation Dataset* specifically for this task, which consists of descriptive apartment tour videos and the corresponding multi-language captions with timestamps. Secondly, we have designed the Multimodal Simultaneous Machine Translation model to process the videos. We have experimented on different fusion timings (encoder side fusion and decode side fusion), cross attention layers (concatenated, serial and merged), features (ResNet50 feature, image captioning feature and ViLT feature). Furthermore, we have proposed the preliminary idea of the masking strategy which can help exclude some noises for the visual space features. Experiment results have shown that our Multi-modal Simultaneous Machine Translation model can gain improvement over the monomodal one and at the same time keep the same latency level. Finally, we have found that when using the Transformer, decoder side fusion with serial cross attention layer and the ViLT feature with encoder mask, we can obtain the greatest translation quality improvement.

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

References

Glasses | meaning in the cambridge english dictionary. <https://dictionary.cambridge.org/dictionary/english/glasses>. (Accessed on 08/25/2021).

Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. 2018. Bottom-up and top-down attention for image captioning and visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6077–6086.

Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*.

Srinivas Bangalore, Vivek Kumar Rangarajan Sridhar, Prakash Kolan, Ladan Golipour, and Aura Jimenez. 2012. Real-time incremental speech-to-speech translation of dialogs. In *Proceedings of the 2012 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 437–445.

Ozan Caglayan, Walid Aransa, Yaxing Wang, Marc Masana, Mercedes García-Martínez, Fethi Bougares, Loïc Barrault, and Joost Van de Weijer. 2016. Does multimodality help human and machine for translation and image captioning? *arXiv preprint arXiv:1605.09186*.

Ozan Caglayan, Julia Ive, Veneta Haralampieva, Pranava Madhyastha, Loïc Barrault, and Lucia Specia. 2020. Simultaneous machine translation with visual context. *arXiv preprint arXiv:2009.07310*.

Daniel Cer, Yinfei Yang, Sheng-yi Kong, Nan Hua, Nicole Limtiaco, Rhomni St John, Noah Constant, Mario Guajardo-Céspedes, Steve Yuan, Chris Tar, et al. 2018. Universal sentence encoder. *arXiv preprint arXiv:1803.11175*.

Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. 2020. Uniter: Universal image-text representation learning. In *European conference on computer vision*, pages 104–120. Springer.

Kyunghyun Cho and Masha Esipova. 2016. Can neural machine translation do simultaneous translation? *arXiv preprint arXiv:1606.02012*.

Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee.

Desmond Elliott, Stella Frank, Khalil Sima'an, and Lucia Specia. 2016. Multi30k: Multilingual english-german image descriptions. *arXiv preprint arXiv:1605.00459*.

Christian Fügen, Alex Waibel, and Muntsin Kolss. 2007. Simultaneous translation of lectures and speeches. *Machine translation*, 21(4):209–252.

Jiatao Gu, Graham Neubig, Kyunghyun Cho, and Victor OK Li. 2016. Learning to translate in real-time with neural machine translation. *arXiv preprint arXiv:1610.00388*.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778.

Lisa Anne Hendricks, John Mellor, Rosalia Schneider, Jean-Baptiste Alayrac, and Aida Nematzadeh. 2021. Decoupling the role of data, attention, and losses in multimodal transformers. *arXiv preprint arXiv:2102.00529*.

Aizhan Imankulova, Masahiro Kaneko, Toshio Hirasawa, and Mamoru Komachi. 2020. Towards multimodal simultaneous neural machine translation. *arXiv preprint arXiv:2004.03180*.

Julia Ive, Andy Mingren Li, Yishu Miao, Ozan Caglayan, Pranava Madhyastha, and Lucia Specia. 2021. Exploiting multimodal reinforcement learning for simultaneous machine translation. In *EACL*.

Wonjae Kim, Bokyung Son, and Ildoo Kim. 2021. Vilt: Vision-and-language transformer without convolution or region supervision. *arXiv preprint arXiv:2102.03334*.

Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.

Jie Lei, Linjie Li, Luowei Zhou, Zhe Gan, Tamara L Berg, Mohit Bansal, and Jingjing Liu. 2021. Less is more: Clipbert for video-and-language learning via sparse sampling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7331–7341.

Charles N Li and Sandra A Thompson. 1989. *Mandarin Chinese: A functional reference grammar*, volume 3. Univ of California Press.

Jindřich Libovický and Jindřich Helcl. 2017. Attention strategies for multi-source sequence-to-sequence learning. *arXiv preprint arXiv:1704.06567*.

Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer.

Mingbo Ma, Liang Huang, Hao Xiong, Renjie Zheng, Kaibo Liu, Baigong Zheng, Chuanqiang Zhang, Zhongjun He, Hairong Liu, Xing Li, et al. 2018. Stacl: Simultaneous translation with implicit anticipation and controllable latency using prefix-to-prefix framework. *arXiv preprint arXiv:1810.08398*.

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664 Yusuke Oda, Graham Neubig, Sakriani Sakti, Tomoki
665 Toda, and Satoshi Nakamura. 2014. Optimizing seg-
666 mentation strategies for simultaneous speech transla-
667 tion. In *Proceedings of the 52nd Annual Meeting of*
668 *the Association for Computational Linguistics (Vol-*
669 *ume 2: Short Papers)*, pages 551–556.

670 Kishore Papineni, Salim Roukos, Todd Ward, and Wei-
671 Jing Zhu. 2002. Bleu: a method for automatic eval-
672 uation of machine translation. In *Proceedings of the*
673 *40th annual meeting of the Association for Compu-*
674 *tational Linguistics*, pages 311–318.

675 Vivek Kumar Rangarajan Sridhar, John Chen, Srinivas
676 Bangalore, Andrej Ljolje, and Rathinavelu Chengal-
677 varayan. 2013. Segmentation strategies for stream-
678 ing speech translation. In *Proceedings of the 2013*
679 *Conference of the North American Chapter of the*
680 *Association for Computational Linguistics: Human*
681 *Language Technologies*, pages 230–238.

682 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob
683 Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
684 Kaiser, and Illia Polosukhin. 2017. Attention is all
685 you need. In *Advances in neural information pro-*
686 *cessing systems*, pages 5998–6008.

687 Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho,
688 Aaron Courville, Ruslan Salakhudinov, Rich Zemel,
689 and Yoshua Bengio. 2015. Show, attend and tell:
690 Neural image caption generation with visual atten-
691 tion. In *International conference on machine learn-*
692 *ing*, pages 2048–2057. PMLR.

693 Baigong Zheng, Renjie Zheng, Mingbo Ma, and Liang
694 Huang. 2019. Simpler and faster learning of adap-
695 tive policies for simultaneous translation. *arXiv*
696 *preprint arXiv:1909.01559*.

697 Yujie Zhong, Linhai Xie, Sen Wang, Lucia Specia,
698 and Yishu Miao. 2020. Watch and learn: Map-
699 ping language and noisy real-world videos with self-
700 supervision. *arXiv preprint arXiv:2011.09634*.

A Visual Feature extraction 701

ResNet50 object classification features We re-
702 size each image to 224x224 and input the image to
703 the ResNet50 (He et al., 2016) network pretrained
704 on the ImageNet database (Deng et al., 2009) to
705 get the [1, 1000] top layer feature, where each ele-
706 ment is the score for one of 1000 object categories.
707 As there are several images for one sentence, we
708 concatenate the features of multiple images for the
709 same sentence by the first axis and pad to a feature
710 matrix of shape [22, 1000] as the maximum length
711 in the corpus will not exceed 105 (the correspond-
712 ing to 22 frames being sampled). 713

Image captioning features We use the image
714 captioning network proposed by Xu et al. (2015)
715 pretrained on Microsoft COCO dataset (Lin et al.,
716 2014) to extract features. The flattened feature
717 from the last layer of the encoder of the image
718 captioning network for each sampled frame is col-
719 lected. We expect that the image captioning model
720 can anticipate the corresponding sentence for one
721 image and generate feature which is more suitable
722 in the simultaneous video translation task. 723

B Experimental Setting 724

We use the pysimt framework⁶ for our experiments.
725 As for the evaluation metrics, BLEU score (Pap-
726 ineni et al., 2002) is used for the quality evaluation
727 and Average Lagging (Ma et al., 2018) for latency
728 measuring. For the RNN model, we run each exper-
729 iment with the max epoch set to 100 and an early
730 stop criterion: when there is no improvement for
731 10 epochs, the experiment will stop. The initial
732 learning rate is 0.0004, and if the BLEU score on
733 the validation set does not improve for two epochs,
734 the learning rate will decay by multiplying 0.5, and
735 the minimum learning rate is 1e-6. Adam (Kingma
736 and Ba, 2014) optimiser is adopted, and the batch
737 size is set to 64. We use the ResNet50 feature as
738 the visual feature for the RNN model. As we pay
739 more attention to the Transformer architecture, we
740 do not try too many different kinds of features for
741 the RNN model, just to see whether the introduc-
742 tion of the video can improve the quality of the
743 simultaneous translation. 744

For the Transformer model, we run each exper-
745 iment with 350 epochs. As specified in (Vaswani
746 et al., 2017), we use the Noam learning rate sched-
747 uler, setting the initial learning rate to 0.2, warmup
748

⁶<https://github.com/ImperialNLP/pysimt>

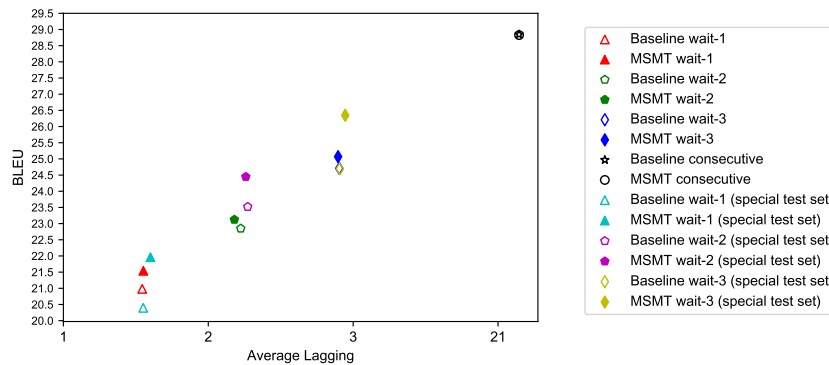


Figure 10: BLEU scores versus Average Lagging for the baseline monomodal model and the multimodal model on two different test sets.

step to 4000 and model dimension to 512. Each experiment is trained on one GeForce GTX 1080 Ti GPU and takes around 7 hours to finish. The number of learnable parameters is around 70M.

We train the models on the training set and record the checkpoints which can get the highest BLEU scores on the validation set, then evaluate these models on the test set. We run experiments for different settings three times and calculate the average.

C Special test set

We create a special test set in which the sentences are all highly related to the frames' content, and all have the pattern 'a/this + noun' (measure words are needed in the Chinese translation) to show that the visual information can better improve the translation quality of these sentences. The results are shown in Figure 10.

As we can see from Figure 10, from wait-1 to wait-3, the improvements of the Multi-modal Simultaneous Machine Translation model (MSMT) on the special test set are more significant than the improvements on the original test set. It shows that when the frames are highly related to the text, the visual cues can better help anticipate to improve the translation quality of the sentences with the pattern 'a/this + noun' (measure words are needed in the Chinese translation). For wait-3, though the language channel will get more source words and be more powerful, the visual information still helps increase the BLEU score by around 1.6 on the special test set.

D Ethical Consideration

One consideration is the collecting and using of the data. The videos in the *ApartmentTour Simul-*

taneous Video Translation Dataset are all public on Youtube. When uploading a public video to Youtube, the Youtuber grants each other user to use the content. All the videos collected are under fair use. Microsoft Azure Service used for generating transcripts follows the GDPR. Only the transcripts, translations, corresponding timestamp, the packed image features and the public Youtube video id of the videos in the dataset will be released. The pysimt framework used in this paper is under MIT License.