


An end-to-end model for cross-lingual transformation of paralinguistic information

Takatomo Kano¹  · Shinnosuke Takamichi¹ ·
Sakriani Sakti¹ · Graham Neubig¹ ·
Tomoki Toda¹ · Satoshi Nakamura¹

Received: 18 July 2016 / Accepted: 5 March 2018
© Springer Science+Business Media B.V., part of Springer Nature 2018

Abstract Speech translation is a technology that helps people communicate across different languages. The most commonly used speech translation model is composed of automatic speech recognition, machine translation and text-to-speech synthesis components, which share information only at the text level. However, spoken communication is different from written communication in that it uses rich acoustic cues such as prosody in order to transmit more information through non-verbal channels. This paper is concerned with speech-to-speech translation that is sensitive to this paralinguistic information. Our long-term goal is to make a system that allows users to speak a foreign language with the same expressiveness as if they were speaking in their own language. Our method works by reconstructing input acoustic features in the target language. From the many different possible paralinguistic features to handle, in this paper we choose duration and power as a first step, proposing a method that can translate these features from input speech to the output speech in continuous space. This is done in a simple and language-independent fashion by training an end-to-end model that maps source-language duration and power information into the target language. Two approaches are investigated: linear regression and neural network models. We evaluate the proposed methods and show that paralinguistic information in the input speech of the source language can be reflected in the output speech of the target language.

Keywords Paralinguistic information · Speech to speech translation · Automatic speech recognition · Machine translation · Text to speech synthesis

✉ Takatomo Kano
kano.takatomo.km0@is.naist.jp

¹ Graduate School of Information Science, Nara Institute of Science and Technology,
Kansai Science City, Japan

1 Introduction

When we speak, we use many different varieties of acoustic and visual cues to convey our thoughts and emotions. Many of those paralinguistic cues transmit additional information that cannot be expressed in words. While these cues may not be a critical factor in written communication, in spoken communication they have great importance; even if the content of the words are the same, if the intonation and facial expression are different an utterance can take an entirely different meaning. As a result, it would be advantageous to take into account these paralinguistic features of speech in any system that is constructed to aid or augment human-to-human communication.

Speech-to-speech translation helps people communicate across different languages, and is thus one prime example of such a system. However, standard speech translation systems only convey linguistic content from source languages to target languages without considering paralinguistic information. Although the input of ASR contains rich prosody information, the words output by ASR are in written form that have no indication of the prosody included in the original speech. As a result, the words output by TTS on the target side will thus be given the canonical prosody for the input text, not reflecting the prosodic traits of the original speech. In other words, because information sharing between the ASR, MT, and TTS modules is limited to only lexical information, after the ASR conversion from speech to text, source-side acoustic details such as rhythm, emphasis, or emotion are lost.

This paper is concerned with speech-to-speech translation that is sensitive to paralinguistic information, with the long-term goal of making a system that allows a user to speak a foreign language with the same expressiveness as if they were speaking in their own language. The proposed method works by recognizing acoustic features (duration and power) in the source language, then reconstructing them in the target language. From the many different possible paralinguistic features to handle, in this paper we choose duration and power as a first step, proposing a method that can translate these features from the input speech to the output speech in continuous space.

First, we extract features at the level of Hidden Markov Model (HMM) states, the use a paralinguistic translation model to predict the duration and power features of HMM states of the output speech. Specifically, we use two approaches: a linear regression model that predicts separately predicts prosody for each word in the vocabulary, and a model that can adapt to more general tasks by training a single model that is applicable to all words in the vocabulary using neural networks.¹

2 Conventional speech-to-speech translation

In conventional speech-to-speech translation systems, the ASR module decodes the text of the utterance from input speech. Acoustic features are represented as $\mathbf{A} = [a_1, a_2 \dots a_{N_a}]$ and the corresponding words are represented as $\mathbf{E} = [e_1, e_2, \dots, e_{N_e}]$.

¹ Part of the content of this article is based on content that has been published in IWSLT and InterSpeech (Kano et al. 2012, 2013). In this paper we describe these methods using a unified formulation, adding a more complete survey, and a discussion of the results in significantly more depth.

N_a and N_e are the lengths of the acoustic feature vectors and spoken words, respectively.

The ASR system finds E that maximizes $P(E|A)$. By Bayes' theorem, we can convert this to (1):

$$P(E|A) \propto P(A|E)P(E), \quad (1)$$

where $P(A|E)$ is the Acoustic Model (AM) and $P(E)$ is the Language Model (LM). The MT module finds the target words sequence J that maximizes probability $P(J|E)$, as in (2):

$$\hat{J} = \operatorname{argmax}_J P(J|E). \quad (2)$$

Similarly to what was done for ASR, we can convert $P(J|E)$ as in (3):

$$\hat{J} = \operatorname{argmax}_J P(E|J)P(J), \quad (3)$$

where $P(E|J)$ is a translation model and $P(E)$ is a language model.

The TTS module generates speech parameters $O = [o_1, o_2, \dots, o_{N_o}]$ given HMM AM states $H_x = [h_1, h_2, \dots, h_{N_h}]$ that represent J . Here N_o is the length of the generated speech parameter sequence, and N_h is the number of states of the HMM AM. The output $O = [o_1, o_2, \dots, o_{N_o}]$ can be represented by (4):

$$\hat{O} = \operatorname{argmax} P(O|H) \quad (4)$$

These three modules only share information through E or J , which are strings of text in the source and target languages, respectively. As a result, all non-verbal information that was original expressed in source speech A is lost the moment it is converted into source text E by ASR.

3 Speech translation considering paralinguistic information

In order to perform speech translation in a way that is also able to consider paralinguistic information, we need to consider how to handle paralinguistic features included in A . Specifically, we need to extract acoustic features during ASR, translate them to another language during MT, and then reflect them in the target speech during TTS.

The first design decision we need to make regards what granularity to use to represent paralinguistic features: phoneme, word, phrase, or sentence level. In the ASR and TTS modules, phonemes are the smallest lexical unit that represent speech, and in the MT module, words are the smallest unit handled by the system. From the point of view of speech processing, phonemes are a good granularity with which to handle paralinguistic features. However, in human speech, paralinguistic features such as emphasis, surprise, and sadness can be more intuitively attributed to the word, phrase and sentence levels (Székely et al. 2014). Thus, as the main focus of our work is on methods for translation of emphasis between languages, for this paper we decide to construct our models purely on the word level. We create word-level AMs for ASR and TTS, extract the paralinguistic features X belonging to each word, and translate these

word-level acoustic features from source to target directly using a regression model in the MT module. Finally we use translated acoustic features to generate output speech in the TTS module.

While the overall framework here is independent of the speech translation task, as the research is ambitious, our experiments described below focus on a limited setting of translating digits. This digit translation task can be motivated by a situation where a customer is contacting a hotel staff member attempting to make a reservation. The customer conveys the reservation number, the hotel staff member confirms it, but the number turns out to be incorrect. In this case, the customer would re-speak the number, using prosody to emphasize the missing information. The problem formulation below will also use this setting as an example, specifically for English–Japanese translation.

3.1 Speech recognition

The first step of the process uses ASR to recognize the lexical and paralinguistic features of the input speech. This can be represented formally as in (5):

$$\hat{E}, \hat{X} = \operatorname{argmax}_{E, X} P(E, X|A), \quad (5)$$

where A indicates the input speech, E indicates the words included in the utterance and X indicates paralinguistic features of the words in E . In order to recognize this information, we construct a word-based HMM AM. The AM is trained with audio recordings of speech and the corresponding transcriptions E using the standard Baum–Welch algorithm. Once we have created our model, we perform simple speech recognition using the HMM AM and a language model that assigns a uniform probability to all digits. Viterbi decoding can be used to find E . Finally we can decide the duration vector x_i of each word e_i based on the time spent in each state of the HMM AM in the path found by the Viterbi algorithm. The power component of the vector is chosen in a similar way, and by taking the mean power value over frames that are aligned to the same state of the AM. We express power as $[power, \Delta power, \Delta \Delta power]$ and join these features together as a super-vector to control power in the translation step. Δ indicates dynamic features. It should be noted that in contrast to other work such as Aguero et al. (2006), for the ASR part, we do not need a manual labeling of the speech prosody, and instead simply segment each word and extract the acoustic features observed.

3.2 Lexical and paralinguistic translation

Lexical translation finds the best translation J of a recognized source sentence E . Generally we can use any variety of statistical machine translation (MT) to obtain this translation in standard translation tasks, but for digit translation we can simply write one-to-one lexical translation rules with no loss in accuracy such as $j_i = e_i$ where i is the word index. Paralinguistic translation converts the source-side acoustic feature vector X into the target-side acoustic feature vector Y according to (6):

$$\hat{Y} = \underset{Y}{\operatorname{argmax}} P(Y|X) \quad (6)$$

There are many types of acoustic features used in ASR and TTS systems, including MFCC, MGC, Filter-bank, F0, power, and duration. In this work we use power and duration to express “emphasis information”. We make this decision due to the fact that MFCC, Filter-bank, and MGC features are more strongly connected to lexical information related to the content of the utterance. F0, power and duration are more correlated with paralinguistic information regarding the method of speech, but because Japanese is a tonal language where F0 has a strong relationship with content distinctions, in this work we focus on duration and power. We control duration and power of each word using a source-side duration and power super-vector $\mathbf{x}_i = [x_1, \dots, x_{N_x}]$, and a target-side duration and power super-vector $\mathbf{y}_i = [y_1, \dots, y_{N_y}]$. Here N_x and N_y represent the length of the paralinguistic feature vector for each source and target word i , respectively.

In these vectors N_x represents the number of HMM states on the source side and N_y represents the number of HMM states on the target side. The sentence duration and power vector consists of the concatenation of the word duration and power vectors such that $\mathbf{Y} = [y_1, \dots, y_n, \dots, y_{N_y}]$. We can assume that the duration and power translation of each word pair are independent of those of other words, allowing us to find the optimal \mathbf{Y} via (7):

$$\hat{Y} = \underset{Y}{\operatorname{argmax}} \prod_n P(y_n|x_n) \quad (7)$$

The word-to-word acoustic translation probability $P(y_n|x_n)$ is calculated according to a linear regression matrix that indicates that y_i is distributed according to a normal distribution, as in (8):

$$P(y_i|x_i) = N(y_i; \mathbf{W}_{e_i,j_i}, x_i', \mathbf{A}) \quad (8)$$

where x' is transposed x and \mathbf{W}_{e_i,j_i} is a regression matrix with bias defining a linear transformation expressing the relationship in duration and power between e_i and j_i . An important point here is how to construct regression matrices for each of the words we want to translate. In order to do so, we optimize each regression matrix in the translation model training data by minimizing root mean squared error (RMSE) with a regularization term, as in (9):

$$\hat{\mathbf{W}}_{e_i,j_i} = \underset{\mathbf{W}_{e_i,j_i}}{\operatorname{argmax}} \sum_{n=1}^N \|y_n^* - y_n\|^2 + \alpha \|\mathbf{W}_{e_i,j_i}\|^2, \quad (9)$$

where N is the number of training samples, n is the id of a training sample, y^* is the target-language reference-word duration and power vector, and α is a hyper-parameter for the regularization term to prevent over-fitting. This maximization can be solved in closed form using simple matrix operations.

3.3 Speech synthesis

In the TTS part of the system we use an HMM-based speech synthesis system (Zen et al. 2009), and reflect the duration and power information of the target word paralinguistic information vector onto the output speech, as in (10):

$$k, asin(10) : \hat{H}_y = \operatorname{argmax} P(H_y|Y) \quad (10)$$

The output speech parameter vector sequence $\mathbf{O} = [o_1, \dots, o_{N_o}]$ is determined by maximizing the target HMM AM \hat{H}_y likelihood function given the target-language sentence $\hat{\mathbf{J}}$, as in (11):

$$\hat{\mathbf{O}} = \operatorname{argmax}_{\mathbf{O}} P(\mathbf{C}|\hat{\mathbf{J}}, \hat{H}_y) \quad (11)$$

$$\text{subject to } \mathbf{C} = \mathbf{M}\mathbf{O} \quad (12)$$

where \mathbf{C} is a joint static and dynamic feature-vector sequence of the target speech parameters, and \mathbf{M} is a transformation matrix from the static feature-vector sequence into the joint static and dynamic feature-vector sequence. When generating speech, the corresponding HMM AM parameters and the length of the target-language state sequence are determined by $\hat{\mathbf{Y}}$ resulting from the paralinguistic translation step. While TTS generally uses phoneme-based HMM models, we instead used a word-based HMM to maintain the consistency of feature extraction and translation. Usually, in TTS phoneme-based HMM AMs, the current HMM AM is heavily influenced by the previous and next phonemes, making it necessary to consider context information from the input sentence. However, in the digit translation task the vocabulary is small, so we construct a word-level independent context HMM AM.

4 End-to-end paralinguistic translation methods

In this section we describe two ways to translate paralinguistic features of the source words to target words. The first is a simple linear regression model that trains a separate model for each word in the vocabulary, and another neural network model that trains a single model for the entire vocabulary but provides the model with information regarding the word identity.

4.1 Linear regression models

Paralinguistic translation converts the source-side paralinguistic features \mathbf{X} into the target-side paralinguistic features \mathbf{Y} , in a manner inspired by previous work on voice conversion (Abe et al. 1988; Toda et al. 2007), as in (13):

$$\hat{\mathbf{Y}} = \operatorname{argmax}_{\mathbf{Y}} P(\mathbf{Y}|\mathbf{X}) \quad (13)$$

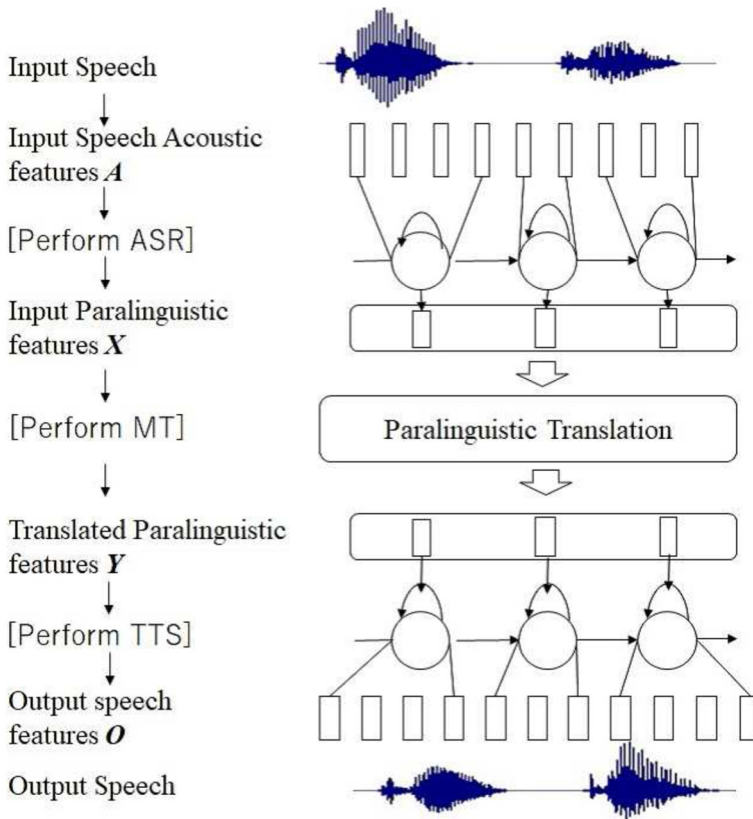


Fig. 1 Overview of the proposed method

In particular, we control duration and power using the source-side word feature vector $x_i = [x_1, \dots, x_{N_h}]$ and target-side word feature vector $y_i = [y_1, \dots, y_{N_h}]$. Here i represents the word id within the vocabulary. In these vectors N_h represents the number of HMM states on the source and target sides. The sentence feature vector consists of the concatenation of the word duration and power vectors \mathbf{Y} where I is the length of the sentence. We assume that the duration and power translation of each word pair is independent, giving (14) (Fig. 1):

$$\hat{\mathbf{Y}} = \underset{\mathbf{Y}}{\operatorname{argmax}} \prod P(y_i | x_i) \quad (14)$$

This can be defined with any function, but we choose to use linear regression, which indicates that y_i is distributed according to a normal distribution, as in (15):

$$P(y_i | x_i) = N(y_i; \mathbf{W}_{e_i, j_i}, x'_i, S) \quad (15)$$

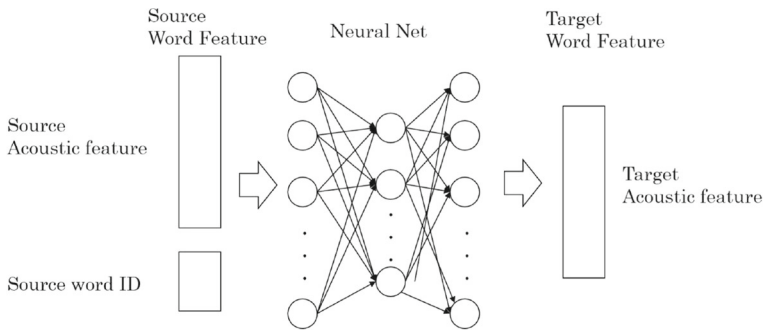


Fig. 2 Neural network for acoustic feature translation

where, x' is transposed x and W_{e_i, j_i} is a regression matrix with bias defining a linear transformation expressing the relationship in duration and power between e_i and j_i .

An important point here is how to construct regression matrices for each of the words we want to translate. In order to do so, we optimize each regression matrix on the translation model training data by minimizing RMSE with a regularization term. This separate training of a model for each word pair allows the model to be expressive enough to learn how each word's acoustics are translated into the target language. However, this has serious problems with generalization, as we will not be able to translate any words that have not been observed in our training data a sufficient number of times to learn the transformation matrix. The simplest way to generalize this model is by not training a separate model for each word, but a global model for all words in the vocabulary. This can be done by changing the word-dependent regression matrix W_{e_i, j_i} into a single global regression matrix W and training the matrix over all samples in the corpus. However, this model can be expected to be insufficiently expressive to properly perform paralinguistic translation. For example, the mapping of duration and power from a one-syllable word to another one-syllable word, and from a one-syllable word to a two-syllable word would vary greatly, but the linear regression model only has the power to perform the same mapping for each word (Fig. 2).

4.2 Global neural network models

As a solution to the problem of the lack of expressiveness in linear regression, we additionally propose a global method for paralinguistic translation using neural networks. Neural networks have higher expressive power due to their ability to handle non-linear mappings, and are thus an ideal candidate for this task. In addition, they allow for the addition of features for many different types of information as is common practice in ASR, MT, and TTS systems, such as word ID vectors, word position, left and right words of input and target words, part of speech, the number of syllables, accent types, etc. This information is known to be useful in TTS (Zen et al. 2009), so we can probably improve the estimation of the output duration and power vector in translation as well. In this research, we use a feed forward neural network that proposes the best output word acoustic feature vector given the input word acoustic

feature vector X . As additional features, we also add a binary vector with the ID of the present word set to 1, and the position of the output word. In this work, because the task is simple we just use this simple feature set, but this could easily be expanded for more complicated tasks. For the sake of simplicity, in this formulation we show an example with the word acoustic feature vector only. First, we set each input unit x_i equal to the input vector value: $l_i = x_i$. The hidden units h_j are calculated according to the input-hidden unit weight matrix W_h , as in (16):

$$\pi_j = \frac{1}{1 + \exp(-\alpha \sum_i w_{i,j}^h l_i)} \quad (16)$$

where α is the gradient of the sigmoid function. The output units ψ_k and final acoustic feature output y_k are set as in (17):

$$\psi_k = \sum_j w_{j,k}^o \pi_j \cdot y_k = \psi_k \quad (17)$$

where W_o is the hidden-output unit weight matrix. As an optimization criterion, we use minimization of RMSE, which is achieved through simple back propagation and weight update, as is standard practice in neural network models.

5 Evaluation

5.1 Experimental setting

We examine the effectiveness of the proposed method through English–Japanese speech-to-speech translation experiments. We use the “AURORA-2” data set, based on a version of the original TIDigits down-sampled at 8 kHz from 55 male and 55 female speakers, with different noise signals artificially added to clean speech data.

As mentioned previously, in these experiments we assume the use of speech-to-speech translation in a situation where the speaker is attempting to reserve a ticket by phone in a different language. When the listener makes a mistake when listening to the ticket digit, the speaker re-speaks, emphasizing the mistaken digit. In this situation, if we can translate the paralinguistic information, particularly emphasis, this will provide useful information to the listener about where the mistake is. In order to simulate this situation, we recorded a bilingual speech corpus where an English–Japanese bilingual speaker emphasizes one word during speech in a string of digits. The content spoken was 500 sentences from the AURORA-2 test set, chosen to be word-balanced by greedy search (Zhang and Nakamura 2003). This was further split into a training set of 445 utterances and a test set of 55 utterances.

To train the ASR model, we use 8440 utterances of clean and noisy speech from the training set of the AURORA-2 dataset and train with the HTK toolkit. In the ASR module we trained an HMM AM, where each word has 16 HMM states, and for silence we allocate 3 states. The lexical translation is performed by Moses (Koehn et al. 2007). We further used the 445 utterances of training data to build an English–Japanese speech

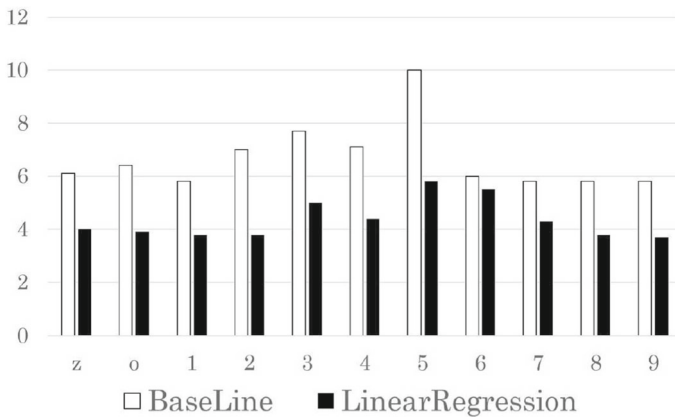


Fig. 3 Root mean squared error rate (RMSE) between the reference target duration and the system output for each digit

translation system that includes our proposed paralinguistic translation model. We set the number of HMM states per word in the ASR AM to 16, the shift length to 5ms, and other various settings to follow (Leonard 1984; Pearce and Hirsch 2000). To simplify the problem, experiments were performed where ASR has no errors. For TTS, we use the same 445 utterances for training an independent context synthesis model. In this case, the speech signals were sampled at 16 kHz. The shift length and HMM states are identical to the settings for ASR.

In the evaluation, we compare the following systems:

- Baseline** No translation of paralinguistic information,
- EachLR** Linear regression with a model for each word,
- AllLR** A single linear regression model trained on all words,
- AllNN** A single neural network model trained on all words,
- AllNN-ID** The AllNN model without additional features.

In addition, we use naturally spoken speech as an oracle output.

5.2 Objective evaluation

We first perform an objective assessment of the translation accuracy of duration and power, the results of which are found in Figs. 3 and 4. For each of the nine digits plus “oh” and “zero,” we compared the difference between the proposed and baseline duration and power and the reference speech duration and power in terms of RMSE. From these results, we can see that the target speech duration and power output by the proposed method is more similar to the reference than the baseline over all eleven categories, indicating that the proposed method is objectively more accurate in translating duration and power. Second, we compare the proposed linear regression against the neural network model in Figs. 5 and 6. We compared the differences between the system duration and power and the reference speech duration and power in terms of RMSE. From these results, we can see that the AllLR model is not effective at mapping

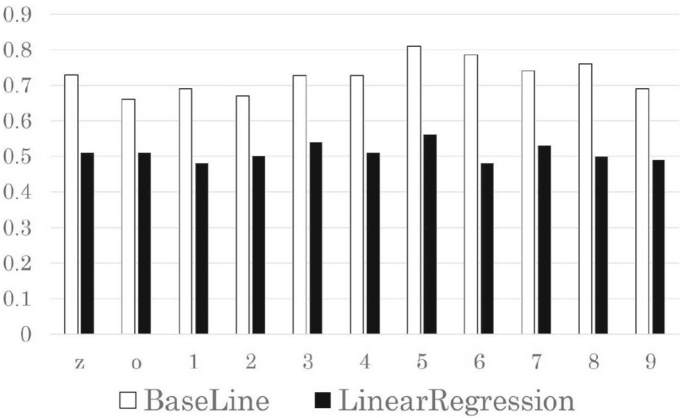


Fig. 4 Root mean squared error rate (RMSE) between the reference target power and the system output for each digit

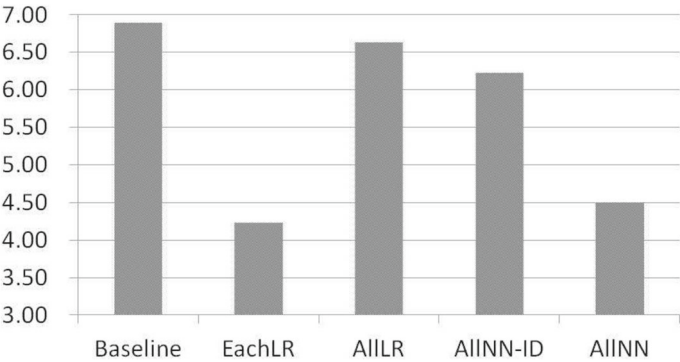


Fig. 5 RMSE between the reference and system duration

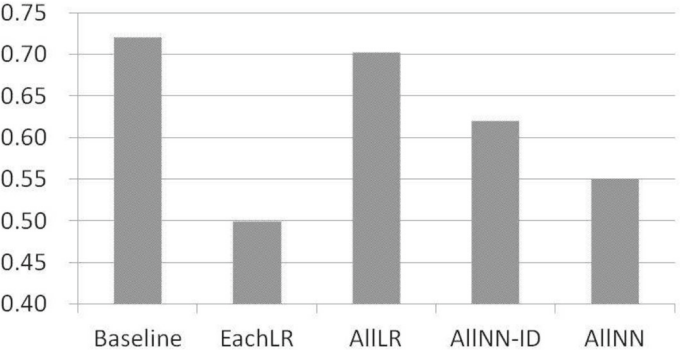


Fig. 6 RMSE between the reference and system power

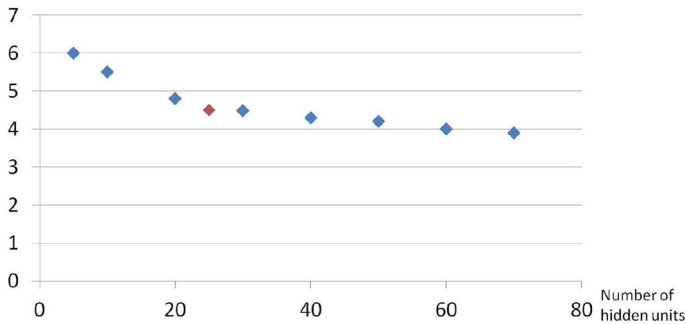


Fig. 7 RMSE of duration for each number of NN hidden units

duration and power information, achieving results largely equal to the baseline. The AIINN model without linguistic information does slightly better but still falls well short of the EachNN baseline. Finally, we can see that our proposed methods outperform the baseline, and AIINN is able to effectively model translation of paralinguistic information, although accuracy of power lags slightly behind that of duration.

We also show the relationship between the number of NN hidden units and RMSE of duration in Fig. 7 (the graph for power was similar). It can be seen that RMSE continues to decrease as we add more units, but with diminishing returns after 25 hidden units. When comparing the number of free parameters in the EachLR model ($17 * 16 * 11 = 2992$) and the AIINN model with 25 hidden units ($28 * 25 + 25 * 16 = 1100$), it can be seen that we were able to significantly decrease the number of parameters as well.

5.3 Subjective evaluation

As a subjective evaluation we asked native speakers of Japanese to evaluate how well emphasis was translated into the target language for the baseline, oracle, and EachLR and AIINN models when translating duration or duration + power. The first experiment asked the evaluators to attempt to recognize the identities and positions of the emphasized words in the output speech. The overview of the result for the word and emphasis recognition rates is shown in Fig. 8. We can see that all of the paralinguistic translation systems show a clear improvement in the emphasis recognition rate over the baseline. There is no significant difference between the linear regression and neural network models, indicating that the neural network learned a paralinguistic information mapping that allows listeners to identify emphasis effectively. The second experiment asked the evaluators to subjectively judge the strength of emphasis with the following three degrees:

- 1 Not emphasized,
- 2 Slightly emphasized,
- 3 Emphasized.

The overview of the experiment regarding the strength of emphasis is shown in Fig. 9, where we can see that all systems show a significant improvement in the subjective perception of strength of emphasis. In this case, there seems to be a slight subjective

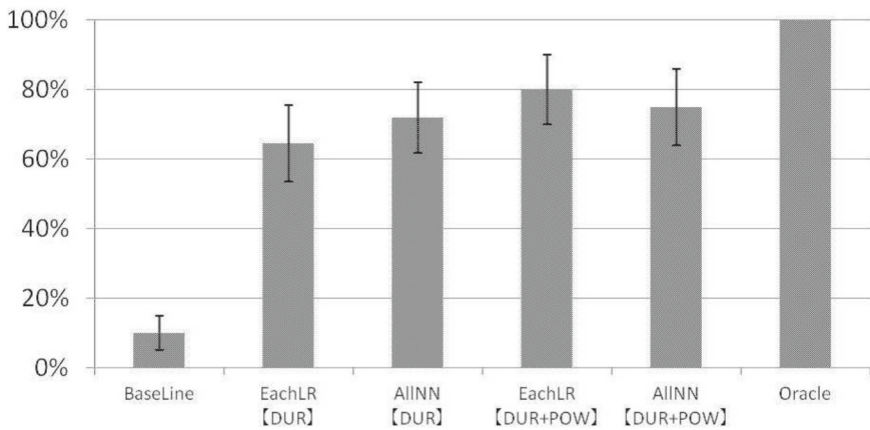


Fig. 8 Prediction rate

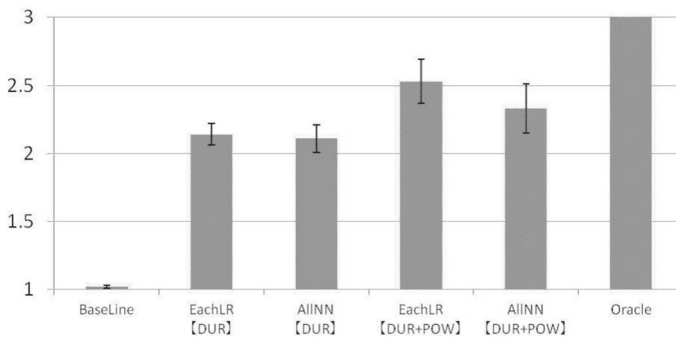


Fig. 9 Prediction strength of emphasis

preference towards EachLR when power is considered, reflecting the slightly smaller RMSE found in the automatic evaluation. We also performed emphasis translation that only used power, but the naturalness of the generated speech was quite low. This resulted in drastic speech volume changes in a short time. Because our proposed method extracts power features for each frame given by duration information, the power extraction has a high dependency on duration. In this method, if we try to handle other acoustic features (e.g. F0) then we also suspect that we will need to model duration together with these features as well.

6 Related work

There have been several studies demonstrating improved speech-translation performance by translating source-side speech non-lexical information into target-side speech non-lexical information. Some previous work (Jiang et al. 2011; Neubig et al. 2012; Dreyer and Dong 2015) has focused on the input speech information (for example, phoneme similarity, number of fillers, and ASR parameters) and tried to explore

a tight coupling of ASR and MT for speech translation, boosting translation quality as measured by BLEU score (Papineni et al. 2002). Other related work focuses on recognizing speech intonation to reduce translation ambiguity on the target side (Takezawa et al. 1998; Wahlster 2001). These methods consider non-lexical information to boost translation accuracy. However as we mentioned before, there is more to speech translation than just accuracy, and we should consider other features such as the speaker's facial and prosodic expressions.

There is some research that considers translating these expressions and improves speech-translation quality in other ways that cannot be measured by BLEU. For example, some work focuses on facial information and tries to translate speaker emotion from source to target (Morishima and Nakamura 2002; Székely et al. 2014). In contrast, Aguero et al. (2006), Anumanchipalli et al. (2012) and Sridhar et al. (2013) focus on the prosody of the input speech, extracting F0 from the source speech at the sentence level and clustering accent groups. These are then translated into target-side accent groups, considering the prosody encoded as factors in a factored translation model (Koehn and Hoang 2007) to convey prosody from source to target.

In our work, we focus on source speech acoustic features, which we extract and translate to target acoustic features directly and continuously. In this framework, we need two translation models: one for word-to-word translation, and another for acoustic translation. We built acoustic translation models with linear regression for each translation pair. This method is simple, and we can translate acoustic features without having an adverse affect on BLEU score. After this work was originally performed, several related works have modeled emphasis by HMM AMs and calculated emphasis levels and translated the emphasis at the word level (Do et al. 2015a, b). These works expand our work to large vocabulary translation tasks. The major difference of this word and our work is the paralinguistic extraction method. In their work, they handle emphasis as a level between 0 and 1 that calculates similarity between an HMM AM for emphasized speech and another HMM AM for normal speech. Each word has one emphasis-level feature and maps these emphasis levels between input and target sequences. In their work, Do et al. need to annotate a paralinguistic label for each type of paralinguistic information they want to handle, so if they were to expand this model to other varieties of paralinguistic information (e.g. emotion or voice quality), they would need further annotated training data to do so. In contrast, in our work we perform regular ASR to obtain alignments and extract observed features, and do not need to specify specific linguistic labels.

State-of-the-art work on speech translation (Do et al. 2017) translates input speech to target words directly with a sequential attentional model. In this work they only focus on linguistic features on the target side, and evaluate quality according to BLEU score. There is also work that focuses on direct speech-to-text translation using sequential attentional models (Duong et al. 2016; Weiss et al. 2017). In this work, any paralinguistic features that exist on the source side may be reflected in the lexical content of the target translations, but paralinguistic information will not be reflected in the target speech.

7 Conclusion

In this paper we proposed a generalized model to translate duration and power information for speech-to-speech translation. Experimental results showed that our proposed method can model input speech emphasis more effectively than baseline methods. In future work we plan to expand beyond the digit translation task in the current paper to a more general translation task using phrase-based or attention-based neural MT. The difficulty here is the procurement of parallel corpora with similar paralinguistic information for large-vocabulary translation tasks. We are currently considering possibilities including simultaneous interpretation corpora and movie dubs. Another avenue for future work is to expand to other acoustic features such as F0, which play an important part in other language pairs.

Acknowledgements The funding was provided by Japan Society for the Promotion of Science (Grand Nos. 24240032 and 26870371).

References

- Abe M, Nakamura S, Shikano K, Kuwabara H (1988) Voice conversion through vector quantization. In: ICASSP-88, international conference on acoustics, speech, and signal processing, New York City, vol 1, pp 655–658
- Agüero PD, Adell J, Bonafonte A (2006) Prosody generation for speech-to-speech translation. In: 2006 IEEE international conference on acoustics speech and signal processing proceedings, Toulouse, France
- Anumanchipalli GK, Oliveira LC, Black AW (2012) Intent transfer in speech-to-speech machine translation. In: 2012 IEEE spoken language technology workshop (SLT), Miami, FL, pp 153–158
- Do QT, Sakti S, Neubig G, Toda T, Nakamura S (2015a) Improving translation of emphasis with pause prediction in speech-to-speech translation systems. In: IWSLT 2015: proceedings of the 12th international workshop on spoken language translation, Da Nang, Vietnam, pp 204–208
- Do QT, Takamichi S, Sakti S, Neubig G, Toda T, Nakamura S (2015b) Preserving word-level emphasis in speech-to-speech translation using linear regression HSMs. In: INTERSPEECH 2015, 16th annual conference of the international speech communication association, Dresden, pp 3665–3669
- Do QT, Sakti S, Nakamura S (2017) Toward expressive speech translation: a unified sequence-to-sequence LSTMs approach for translating words and emphasis. In: Interspeech 2017, 18th annual conference of the international speech communication association, Stockholm, Sweden, pp 2640–2644
- Dreyer M, Dong Y (2015) Apro: all-pairs ranking optimization for mt tuning. In: Proceedings of the 2015 conference of the North American chapter of the association for computational linguistics: human language technologies, Denver, CO, pp 1018–1023
- Duong L, Anastasopoulos A, Chiang D, Bird S, Cohn T (2016) An attentional model for speech translation without transcription. In: Proceedings of the 2016 conference of the North American chapter of the association for computational linguistics: human language technologies, San Diego, CA, pp 949–959
- Jiang J, Ahmed Z, Carson-Berndsen J, Cahill P, Way A (2011) Phonetic representation-based speech translation. In: Proceedings of machine translation summit XIII, Xiamen, China, pp 81–88
- Kano T, Sakti S, Takamichi S, Neubig G, Toda T, Nakamura S (2012) A method for translation of paralinguistic information. In: 2012 International workshop on spoken language translation, Hong Kong, pp 158–163
- Kano T, Takamichi S, Sakti S, Neubig G, Toda T, Nakamura S (2013) Generalizing continuous-space translation of paralinguistic information. In: INTERSPEECH 2013, 14th Annual conference of the international speech communication association, Lyon, France, pp 2614–2618
- Koehn P, Hoang H (2007) Factored translation models. In: EMNLP-CoNLL-2007: proceedings of the 2007 joint conference on empirical methods in natural language processing and computational natural language learning, Prague, Czech Republic, pp 868–876
- Koehn P, Hoang H, Birch A, Callison-Burch C, Federico M, Bertoldi N, Cowan B, Shen W, Moran C, Zens R, Dyer C, Bojar O, Constantin A, Herbst E (2007) Moses: open source toolkit for statistical

- machine translation. In: Proceedings of the 45th annual meeting of the ACL on interactive poster and demonstration sessions, Prague, Czech Republic, pp 177–180
- Leonard R (1984) A database for speaker-independent digit recognition. In: ICASSP '84. IEEE international conference on acoustics, speech, and signal processing, San Diego, CA, pp 328–331
- Morishima S, Nakamura S (2002) Multi-modal translation system and its evaluation. In: Proceedings of the fourth IEEE international conference on multimodal interfaces, Pittsburgh, PA, pp 241–246
- Neubig G, Duh K, Ogushi M, Kano T, Kiso T, Sakti S, Toda T, Nakamura S (2012) The NAIST machine translation system for IWSLT 2012. In: IWSLT-2012: 9th international workshop on spoken language translation, Hong Kong, pp 54–60
- Papineni K, Roukos S, Ward T, Zhu WJ (2002) Bleu: a method for automatic evaluation of machine translation. In: Proceedings of the 40th annual meeting on association for computational linguistics, Philadelphia, PA, pp 311–318
- Pearce D, Hirsch HG (2000) The Aurora experimental framework for the performance evaluation of speech recognition systems under noisy conditions. In: ASR2000—Automatic speech recognition: challenges for the new millenium, Paris, France, pp 181–188
- Sridhar VKR, Bangalore S, Narayanan S (2013) Enriching machine-mediated speech-to-speech translation using contextual information. *Comput Speech Lang* 27(2):492–508
- Székely É, Steiner I, Ahmed Z, Carson-Berndsen J (2014) Facial expression-based affective speech translation. *J Multimodal User Interfaces* 8(1):87–96
- Takezawa T, Morimoto T, Sagisaka Y, Campbell N, Iida H, Sugaya F, Yokoo A, Yamamoto S (1998) A Japanese-to-English speech translation system: ATR-MATRIX. In: 5th international conference on spoken language processing, ICSLP'98 proceedings, Sydney, Australia, pp 2779–2883
- Toda T, Black AW, Tokuda K (2007) Voice conversion based on maximum-likelihood estimation of spectral parameter trajectory. *IEEE Trans Audio Speech Lang Process* 15(8):2222–2235
- Wahlster W (2001) Robust translation of spontaneous speech: a multi-engine approach. In: Proceedings of seventeenth international joint conference on artificial intelligence, invited papers, Seattle, WA, pp 19–28
- Weiss RJ, Chorowski J, Jaitly N, Wu Y, Chen Z (2017) Sequence-to-sequence models can directly transcribe foreign speech. [arXiv:1703.08581](https://arxiv.org/abs/1703.08581)
- Zen H, Tokuda K, Black AW (2009) Statistical parametric speech synthesis. *Speech Commun* 51(11):1039–1064
- Zhang J, Nakamura S (2003) An efficient algorithm to search for a minimum sentence set for collecting speech database. In: 15th international congress of phonetic sciences (ICPhS-15), Barcelona, Spain, pp 3145–3148