

# Low-Resource Methods for Hawaiian Machine Translation

**Nolan Brophy**

University of Hawai‘i at Hilo  
nolanv@hawaii.edu

**Winston Wu**

University of Hawai‘i at Hilo  
wsu@hawaii.edu

## Abstract

This paper investigates the challenges of low-resource machine translation for ‘Ōlelo Hawai‘i (Hawaiian), a critically endangered Polynesian language. We compile a corpus of publicly available Hawaiian-English bitext and investigate the effectiveness of neural sequence-to-sequence models and large language models for translating Hawaiian. To address data scarcity, we employ various data augmentation techniques, including backtranslation, multilingual training using parallel corpora in related languages, and leveraging dictionary entries. Our experiments demonstrate that multilingual training significantly improves model performance, particularly when incorporating bitext from related Polynesian languages. Fine-tuned large language models were not able to outperform mBART, highlighting that smaller and simpler models are still relevant, especially in low-resource scenarios.

## 1 Introduction

‘Ōlelo Hawai‘i (Hawaiian) is a Polynesian language mainly spoken in Hawaii, USA. It is part of the Austronesian language family and is closely related to other Polynesian languages such as Tahitian, Māori, Samoan, and Tongan. Historically, ‘Ōlelo Hawai‘i was the dominant language in Hawaii, with a rich oral tradition including chants, songs, and traditional knowledge passed down through generations. It is currently classified by UNESCO as a critically endangered language as a result of a sharp decline in native speakers in the late 1800s and early 1900s due to colonization, the suppression of Hawaiian in schools and government, and English becoming the dominant language. Efforts to revitalize Hawaiian starting around 40 years ago have prevented the language from dying through the development of language immersion schools that train the next generation of native speakers. Today, Hawaiian is one of the official languages of the state of Hawaii, and

it is taught in schools, universities, and community programs throughout the state.

Hawaiian has a limited amount of printed materials and text data on the web, making it challenging to develop digital tools and resources for language revitalization. Today, most learning materials are published by academic institutions like universities or organizations such as ‘Aha Pūnana Leo, which runs immersion schools across the state. Our project aims to develop machine translation systems to make the language more accessible and increase the number of speakers of Hawaiian by enabling broader engagement with Hawaiian content.

Developing accurate neural machine translation (NMT) systems requires a large amount on the order of millions of sentence pairs (Koehn and Knowles, 2017), which is not available for Hawaiian as a low-resource language. We first gather existing Hawaiian bitext from the web. Then, to address the low-resource nature of Hawaiian, we employ several methods to counteract data scarcity and improve model performance. We experiment with several data augmentation techniques, including back-translation (Sennrich et al., 2016) to generate new training pairs, incorporating multilingual data from related Polynesian languages like Samoan and Māori to learn patterns across related languages, and using lexical translations from Hawaiian dictionaries. For the models, we experiment with fine-tune popular sequence-to-sequence neural machine translation models such as BART (Lewis et al., 2019), multilingual variants like mBART (Liu et al., 2020a), and LLMs like Qwen3 (Qwen Team, 2025). By combining these approaches, our project aims to create more effective and accessible translation systems, enabling broader access to Hawaiian language content and supporting the language’s revitalization efforts.

## 2 Related Work

It is well-known that machine translation systems struggle when faced with small amounts of data [Haddow et al. \(2022\)](#). Researchers have tackled the low-resource issue from a variety of angles, ranging from improving the data or improving the model.

Back-translation ([Sennrich et al., 2016](#)) is a popular method for augmenting the training data for machine translation systems in low-resource scenarios. In the task of translating from a low-resource source language to a high-resource target language, an initial MT system is trained in the opposite direction (i.e. target to source). Because the target language has much more available monolingual data, e.g. from the web, the MT model is then used to translate a large monolingual corpus from the target language back into the source language, resulting in a larger, but potentially lower-quality, set of translation pairs. This new bitext is added to the original bitext to train a second MT system in the source-to-target direction, which has shown to improve over the first MT system. Back-translation has been successfully applied to other languages such as Spanish-Portuguese, Czech-Polish, and Hindi-Nepali ([Przystupa and Abdul-Mageed, 2019](#)), as well as the low-resource indigenous language Bribri ([Feldman and Coto-Solano, 2020](#)).

Models trained on multiple languages can also translate a single language more accurately ([Zoph et al., 2016](#); [Fan et al., 2021](#)). For example, mBART ([Liu et al., 2020a](#)) is pretrained on 50 additional languages over the English-only BART model and demonstrates enhanced translation capabilities. Multilingual training has been shown to be effective for lower-resource languages such as Indonesian and Malaysian ([Poncelas and Effendi, 2022](#)) and Creole languages ([Robinson et al., 2024](#)). Previous studies have shown that adding multiple reference translations for the same source sentence also improves the model’s generalizability ([Khayrallah et al., 2020](#)).

Large language models have been shown to perform poorly on translation of low-resource languages ([Kocmi et al., 2023](#)). One method to improve performance is to use in-context learning ([Brown et al., 2020](#)), where the model is prompted with some examples. For our work, we use translation pairs as additional examples to guide a general-purpose LLM to generate translations.

## 3 Data

Machine translation systems are typically trained on parallel corpora containing pairs of a sentence in one language and a corresponding sentence in the other language. We collected an initial dataset consisting of roughly 29k Hawaiian-English sentence pairs from a diverse range of sources:

- Ka Baibala Hemolele, the Hawaiian Bible
- Example sentences from the Combined Hawaiian Dictionary ([Trussel, 2020](#))
- Fornander Collection of Hawaiian Antiquities and Folk-lore, Volumes 1 and 2, ([Fornander, 1917](#)), a collection of the Hawaiians’ account of the formation of the Hawaiian Islands and origin of the Hawaiian people
- ‘Ōlelo No‘eau ([Pukui, 1983](#)), a collection of Hawaiian proverbs and sayings
- La‘ieikawai ([Haleole, 1863](#)), a novel, and the first substantial fiction work by a Hawaiian
- Children’s stories written for beginner Hawaiian students

Hawaiian is written with Latin characters and has two orthographic systems. The traditional system was introduced in the 1820s when Hawaiian first gained a writing system, and consists solely of the same characters as the English alphabet. The modern system, introduced in the 1950s, added several characters that aid in disambiguating pronunciation and meaning: the ‘okina, a backwards apostrophe (‘) used to indicate the glottal stop, and the kahakō, a macron above vowels (ā ē ī ō ū) to indicate long vowels. In our data, the Fornander collection (4k sentences) is written using the traditional system, while the rest are written in the modern system. Thus, training a translation system on the combination of Hawaiian written in traditional and modern systems should be more robust to the spelling variety used.

Due to the low-resource nature of Hawaiian, we first examine how models perform when learning to translate Hawaiian to English on this initial dataset. Then, we supplement this data with additional data:

- Back-translated sentences (10k) from English movie subtitles from OpenSubtitles ([Lison and Tiedemann, 2016](#)).
- Translations in Samoan and Māori, two Polynesian languages closely related to Hawaiian, from the NLLB dataset ([Costa-Jussà et al., 2022](#)) provided by OPUS ([Tiedemann, 2012](#)). These translations were automatically

extracted from web corpora. Manual examination indicates that these translations largely consist of sentences from the Bible. Although this data contains several hundred thousand sentence pairs, we randomly select a subset of 10k sentences to not overwhelm the other data.

- Hawaiian words and their English translation from the Combined Hawaiian Dictionary (Trussel, 2020). We use 57k dictionary entries where the translation consists of 3 or fewer words.

## 4 Models and Experiments

We experiment with several popular neural machine translation models for translating from Hawaiian into English. For all experiments, we use the same randomly shuffled dataset, split into 80% train, 10% validation, and 10% test, with additional data added to the training set. The models were trained until performance did not improve on the validation set with a patience of 5 epochs.

### 4.1 Sequence-to-Sequence NMT Models

First, we investigate BART (Lewis et al., 2020), a transformer encoder-decoder model pre-trained on English, specifically the `bart-base` and `bart-large` models. In the same family of models, we also examine mBART-50 (Tang et al., 2020), a BART model pre-trained using a multilingual denoising pretraining objective (Liu et al., 2020b) in 50 languages. We specifically use the `facebook/mbart-large-50-many-to-many-mmt` checkpoint. Note that although these models have been exposed to multiple languages during pre-training, they were not specifically pre-trained on Hawaiian data.

We experiment with fine-tuning these models for translating Hawaiian to English. Because mBART requires the use of a special source language token, we modify the mBART tokenizer by adding three new special tokens `haw_XX`, `mri_XX`, and `smo_XX` to the tokenizer and model vocabulary, which indicate that the source language is Hawaiian, Māori, and Samoan, respectively. The Samoan and Māori language tokens were not used in the initial experiments, but were used in the multilingual data augmentation experiments described below.

### 4.2 Large Language Models

We also experiment with zero-shot, few-shot (in-context learning), and fine-tuning Qwen3 (Qwen

Team, 2025), a general-purpose LLM with support for over 100 languages. Specifically, we use the `Qwen3-4B-Instruct-2507` model, which is instruction fine-tuned and does not support thinking mode. In the zero-shot setting, we simply prompt the model "Translate Hawaiian to English: [Hawaiian sentence]". In the few-shot setting, we provide 10 random translation pairs before prompting the model with the same translation prompt as in the zero-shot setting. As recommended by the Qwen authors, we set `temperature=0.7`, `top_p=0.8`, and `top_k=20`. In order to guide the Qwen3 model to better perform the translation task, we also experiment with parameter efficient fine-tuning using LORA (Hu et al., 2022) with `rank=16` and `alpha=16`, using a batch size of 4 with gradient accumulation of 16. The model was fine-tuned on single-turn conversations, where the input is the same prompt as above, and the output is solely the English translation. We perform LLM fine-tuning using the Unsloth framework (Han et al., 2023). We run our experiments on a local machine with a single NVIDIA A6000 GPU.

### 4.3 Low-Resource Methods for Data Augmentation

Because the existing parallel corpus is relatively small, we experiment with several data augmentation methods described in the previous section, including backtranslation using an mBART model trained in the opposite direction, multilingual Māori-English and Samoan-English bitext, and lexical translations from a Hawaiian dictionary. For these data augmentation experiments, we finetune mBART, which performed the best on the initial dataset.

### 4.4 Evaluation

We evaluate the performance of each model using BLEU (Papineni et al., 2002) and chrF++ (Popović, 2017). BLEU measures the n-gram precision between the translation and the reference, and chrF++ measures character-level and bigram F-score.

## 5 Results and Discussion

A summary of the performance of each model and data scenario is shown in Table 1. For the base BART and mBART models, the larger models show slight improvements over the base model. Surprisingly, mBART performed equally to BART-large, even though it was trained on more languages. The extra multilingual training did not seem to help,

Model	BLEU	chrF++
bart-base	19.17	41.48
bart-large	21.01	43.60
mbart-large-50-mmt	21.27	43.42
mbart + backtranslation	20.95	42.82
mbart + multilingual	27.07	47.46
mbart + dictionary	19.60	41.95
Qwen3 0-shot	7.00	29.03
Qwen3 10-shot	8.32	30.85
Qwen3 Fine-tuned	19.71	39.23

Table 1: Model performance with various models. Data augmentation experiments were performed using the mbart-large-50-mmt model, shortened to mbart in the table.

perhaps because none of the 50 languages it was trained on were closely related to Hawaiian (the closest is Indonesian, which is Austronesian but not in the Polynesian family).

For models trained with additional data, we found a substantial improvement when adding multilingual data in Samoan and Māori. Because these two languages are closely related to Hawaiian, the model was able to effectively learn from the extra non-Hawaiian training data. However, the models trained with additional back-translation and dictionary translations did not improve over the baseline models. For back-translation, the lack of improvement may be due to the corpus being in a different domain: the subtitles come from modern movies, while most of the Hawaiian text is Biblical text or text written around 100 years ago.

For the large language model experiments, the Qwen3 models were not able to generate accurate translations without fine-tuning. This result is understandable given that Hawaiian has a tiny web presence and the model would only have seen a little Hawaiian in its training data. With fine-tuning, Qwen3 approaches the performance of the sequence-to-sequence NMT models. Note that BART (0.1B) and mBART (0.3B) are from a previous era of pre-LLM models with an order of magnitude fewer parameters than Qwen3 (4B), yet still beat a fine-tuned Qwen model.

## 5.1 Discussion

To examine model performance in more detail, we analyze the top model (mBART with initial + extra data) by examining its performance on each

subset of our test data, shown in Table 3. Performance on Ka Baibala Hemolele, the Hawaiian Bible, showed the highest scores for both BLEU and chrF++, which was not a surprise, considering that the Bible comprised the largest portion of the training data. Sentences from the childrens books stories were relatively easier to translate because they generally use simple sentence structures. Performance on the Fornander collection was poor, likely because this data contains a lot of poetry, which is often liberally translated in the training data. Furthermore, this data contains some run-on sentences due to the fact that some of the sentences are verses of songs (mele), which are difficult to split into individual sentences. Laieikawai sentences have a mix of conversational and longer descriptive sentences, characteristic of a fiction novel. ‘Ōlelo No‘eau were also difficult to translate, because these are short phrases that have similarly artistic reference translations.

Appendix ?? and Table 2 contain sample model output and further discussion of the model’s performance. Overall, we find that the models perform better on less poetic text, and tend to more closely follow the literal meaning of the words. These models would be useful for a learner or Hawaiian language enthusiast to get the gist of a Hawaiian text.

## 6 Conclusion

As a critically endangered indigenous language, ‘Ōlelo Hawai‘i suffers from a lack of available bi-text. To address this challenge, we compiled a parallel corpus from publicly available online sources and experimented with fine-tuning multiple neural machine translation models and LLMs to translate Hawaiian. We developed and evaluated several low-resource data augmentation techniques including back-translation, multilingual parallel data from related Polynesian languages (Māori and Samoan), and lexical translations from a Hawaiian dictionary.

From these approaches, multilingual training with Māori and Samoan was most effective, likely because these are closely related languages in the Polynesian language family. Furthermore, our best model, mBART, outperformed Qwen3 (4B), highlighting that sometimes simpler models perform well with less data. Evaluation on a variety of domains showed that the model tends to produce more literal, rather than poetic, translations. Through this work, we aim to support the revitalization of ‘Ōlelo Hawai‘i and make the language and culture

Source	Reference	Model Predictions
Pi'i a'ela 'o Dāvida mai laila aku, a noho i loko o nā wahi pa'a ma 'Enegedi. No ka mea, ke 'ōlelo nei ka Palapala hemolele, 'O ka mea e mana'o'i'o iā ia, 'a'ole ia e hoka.	And David went up from there and lived in the strongholds of En Gedi. As the Scripture says, "Anyone who trusts in him will never be put to shame."	From there David went up and stayed in the strongholds of En Gedi. For the Scriptures say: "He who trusts in God will not be put to shame."
No laila, ke noi aku nei au e pauaho 'ole 'oukou i ku'u pilikia 'ana no 'oukou, 'o kā 'oukou ia e pōmaika'i ai.	I ask you, therefore, not to be discouraged because of my sufferings for you, which are your glory.	Therefore I urge you not to give up my suffering for you, so that you may be blessed.
O oe no ka na e Haunuu, E Haulani, ka mano nui. E Kaalokuloku, e ui e? O kou inoa ia? E o mai.	Are you then, Haunuu, Haulani, the great shark, Kaalokuloku, a question? Is this your name? Make answer.	You are the one, Haunuu. Say, Haulani, the great shark. Say, Kaalokuloku, wail. Is that your name? Answer me.
Huli ae o Pamano a olelo aku: "U! no'u paha ka pii a ola mai au, make olua ia'u."	At this Pamano turned and said: "Yes, here I am going up and if I return alive, I will kill both of you."	Pamano turned around and said: "Yes, it is my responsibility to go up and save my life; I will kill you both."
E ka ohu kolo mai i uka, E ka ohu kolo mai i kai, E kai pupuka, E kai hehena, E kai piliaku.	Ye fog that creeps in the upland, Ye fog that creeps seaward; Ye ugly seas, ye mad seas, Ye kapu-breaking seas.	O sea, O sea of the uplands, O sea of the uplands, O sea of the uplands, O sea of the uplands.
'i akula ua makāula nei, "He wa'a ali'i ho'i kēia e holo mai nei.	Said the seer, "'A chief's canoe comes hither,	Said the seer, "Here comes a chiefly canoe;
Inā he mana'o e ku'i, ku'i mai i ku'u maka."	Strike my face, if you want to!"	if you wish to strike me, strike me the eye."
Ma ke ki'eki'e iki 'ana a'e o ka lā, aia e pi'o ana ke ānuenuē i kai o Kea'au.	A little later in the day the rainbow was at the seacoast of Keaau.	As the sun passes, the rainbow arches over the sea at Keaau.
Komo hou maila ke kai a pio kāna ahi.	The sea came in again and her fire was extinguished.	The sea came in again and his fire was extinguished.
Ua loa'a he 'elua 'o'opu mai ka'u 'ohana keiki mai.	I got two 'o'opu from my nephew.	I got two 'o'opu from my own son.

Table 2: Sample model translations. The sections above separate sentences from the Bible, Fornander collection, Laieikawai, and childrens stories.

Collection	BLEU	chrF++
Baibala	28.13	48.31
Fornander	16.49	41.11
Laieikawai	20.24	43.33
'Ōlelo No'eau	17.61	35.13
Stories	24.55	44.14

Table 3: Translation performance by collection, using the mBART multilingual model.

more accessible to a broader audience. Our models have practical applications for language learners and Hawaiian enthusiasts interested in translating Hawaiian newspapers, literature, and other existing materials. For future work, we plan to work with Hawaiian studies teachers to develop more tools based on their needs and integrate such tools in the classroom.

## Acknowledgments

This work is partially supported by the National Science Foundation (Award No. 2422413). Any opinions, findings, conclusions, or recommendations expressed in this material are those of the

authors and do not necessarily reflect the views of the NSF.

## Ethical Considerations and Limitations

The data compiled in this work were sourced from Project Gutenberg (public domain) or from Ulukau.org (usable for research and educational purposes, prohibited for commercial use). Due to data privacy and sovereignty concerns of the Hawaiian language community, we only experiment with locally-hosted models and were constrained by the VRAM of our GPU. Our results may not be applicable for larger LLMs. This research was also performed before Qwen3.6 was released. The newer Qwen model show substantial improvements in coding performance, and investigating its multilingual capabilities and how it performs on low-resource translation is left for future work.

## References

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, and 1 others. 2020. Language models are

- few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Marta R Costa-Jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, and 1 others. 2022. No language left behind: Scaling human-centered machine translation. *arXiv preprint arXiv:2207.04672*.
- Angela Fan, Shruti Bhosale, Holger Schwenk, Zhiyi Ma, Ahmed El-Kishky, Siddharth Goyal, Mandeep Baines, Onur Celebi, Guillaume Wenzek, Vishrav Chaudhary, and 1 others. 2021. Beyond english-centric multilingual machine translation. *Journal of Machine Learning Research*, 22(107):1–48.
- Isaac Feldman and Rolando Coto-Solano. 2020. [Neural machine translation models with back-translation for the extremely low-resource indigenous language Bribri](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 3965–3976, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Abraham Fornander. 1917. *Fornander Collection of Hawaiian Antiquities and Folk-lore*. Bishop Museum Press.
- Barry Haddow, Rachel Bawden, Antonio Valerio Miceli Barone, Jindřich Helcl, and Alexandra Birch. 2022. [Survey of low-resource machine translation](#). *Computational Linguistics*, 48(3):673–732.
- S N Haleole. 1863. *Ke Kaa o Laieikawai: ka hiwahiwa o Paliuli, kawahineokaliula: Kakauia mailoko mai o na Moolelo Kahiko o Hawaii nei. Kakauia e SN Haleole*. Paia e HM Whitney.
- Daniel Han, Michael Han, and Unsloth team. 2023. [Unsloth](#).
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, and 1 others. 2022. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3.
- Huda Khayrallah, Brian Thompson, Matt Post, and Philipp Koehn. 2020. [Simulated multiple reference training improves low-resource machine translation](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 82–89, Online. Association for Computational Linguistics.
- Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Philipp Koehn, Benjamin Marie, Christof Monz, Makoto Morishita, Kenton Murray, Makoto Nagata, Toshiaki Nakazawa, Martin Popel, and 2 others. 2023. [Findings of the 2023 conference on machine translation \(WMT23\): LLMs are here but not quite there yet](#). In *Proceedings of the Eighth Conference on Machine Translation*, pages 1–42, Singapore. Association for Computational Linguistics.
- Philipp Koehn and Rebecca Knowles. 2017. [Six challenges for neural machine translation](#). In *Proceedings of the First Workshop on Neural Machine Translation*, pages 28–39, Vancouver. Association for Computational Linguistics.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. 2019. [Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension](#). *Preprint*, arXiv:1910.13461.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. [BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.
- Pierre Lison and Jörg Tiedemann. 2016. [OpenSubtitles2016: Extracting large parallel corpora from movie and TV subtitles](#). In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 923–929, Portorož, Slovenia. European Language Resources Association (ELRA).
- Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. 2020a. [Multilingual denoising pre-training for neural machine translation](#). *Preprint*, arXiv:2001.08210.
- Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. 2020b. [Multilingual denoising pre-training for neural machine translation](#). *Transactions of the Association for Computational Linguistics*, 8:726–742.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. [Bleu: a method for automatic evaluation of machine translation](#). In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- Alberto Poncelas and Johanes Effendi. 2022. [Benefiting from language similarity in the multilingual MT training: Case study of Indonesian and Malaysian](#). In *Proceedings of the Fifth Workshop on Technologies for Machine Translation of Low-Resource Languages (LoResMT 2022)*, pages 84–92, Gyeongju, Republic of Korea. Association for Computational Linguistics.
- Maja Popović. 2017. [chrF++: words helping character n-grams](#). In *Proceedings of the Second Conference on Machine Translation*, pages 612–618, Copenhagen, Denmark. Association for Computational Linguistics.

- Michael Przystupa and Muhammad Abdul-Mageed. 2019. [Neural machine translation of low-resource and similar languages with backtranslation](#). In *Proceedings of the Fourth Conference on Machine Translation (Volume 3: Shared Task Papers, Day 2)*, pages 224–235, Florence, Italy. Association for Computational Linguistics.
- Mary Kawena Pukui. 1983. *‘Ōlelo No‘eau: Hawaiian Proverbs and Poetical Sayings*. Bishop Museum Press.
- Qwen Team. 2025. [Qwen3 technical report](#). *Preprint*, arXiv:2505.09388.
- Nathaniel R. Robinson, Raj Dabre, Ammon Shurtz, Rasul Dent, Onenamiyi Onesi, Claire Bizon Monroc, Loïc Grobol, Hasan Muhammad, Ashi Garg, Naome A. Etori, Vijay Murari Tiyyala, Olanrewaju Samuel, Matthew Dean Stutzman, Bismarck Bamfo Odoom, Sanjeev Khudanpur, Stephen D. Richardson, and Kenton Murray. 2024. [Kreyòl-mt: Building mt for latin american, caribbean and colonial african creole languages](#). *Preprint*, arXiv:2405.05376.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. [Improving neural machine translation models with monolingual data](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 86–96, Berlin, Germany. Association for Computational Linguistics.
- Yuqing Tang, Chau Tran, Xian Li, Peng-Jen Chen, Naman Goyal, Vishrav Chaudhary, Jiatao Gu, and Angela Fan. 2020. Multilingual translation with extensible multilingual pretraining and finetuning. *arXiv preprint arXiv:2008.00401*.
- Jörg Tiedemann. 2012. [Parallel data, tools and interfaces in OPUS](#). In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC’12)*, pages 2214–2218, Istanbul, Turkey. European Language Resources Association (ELRA).
- Stephen Kepano Trussel. 2020. [Kepano’s combined Hawaiian dictionary](#).
- Barret Zoph, Deniz Yuret, Jonathan May, and Kevin Knight. 2016. [Transfer learning for low-resource neural machine translation](#). In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1568–1575, Austin, Texas. Association for Computational Linguistics.