

Bharatiya Vidya Bhavan's
SARDAR PATEL INSTITUTE OF TECHNOLOGY

NLP Experiment 9 & 10

Submitted To

Prof. Reeta Koshy

Submitted By

Pranay Singhvi (2021300126)

Tathagat Sengupta (2021300110)

GitHub Link: <https://github.com/Code-forlife/WorldSpeak-Connect>

Introduction

Imagine a world where language barriers vanish. SeamlessM4T, a cutting-edge machine translation model from Meta AI, aspires to do just that. It's a single, powerful tool capable of tackling various translation tasks. Need to translate text from French to Japanese? No problem. Want to have a real-time conversation in Spanish with someone who only speaks Mandarin? SeamlessM4T can handle that too.

This ambitious model goes beyond simple text translation. It can translate spoken language (speech-to-text) and even generate speech from written text (text-to-speech). It aims to support nearly 100 languages for input and 35 for output, making communication across borders a breeze. By leveraging a powerful transformer architecture, SeamlessM4T strives to deliver high-quality translations, overcoming the limitations of previous models. This technology has the potential to revolutionize communication, fostering global collaboration and understanding.

Objective and Scope

The objective of SeamlessM4T:

- **Break Down Language Barriers:** SeamlessM4T aims to eliminate the need for multiple translation tools by offering a one-stop solution for various translation tasks.
- **Seamless Communication Across Languages:** It tackles text-to-text translation but goes further by handling real-time speech translation (speech-to-text and text-to-speech).
- **Massive Language Coverage:** This ambitious project strives to support nearly 100 languages for input and 35 for output, fostering communication across a vast spectrum.
- **High-Quality Translations:** By utilizing a powerful architecture, SeamlessM4T seeks to deliver accurate and natural-sounding translations, surpassing the limitations of existing models.

1. Scope:

SeamlessM4T breaks new ground in machine translation by aiming to be a one-stop shop for all your language needs. Here's a glimpse of its impressive scope:

- **Multimodal Magic:** It tackles not just text translation (think emails or documents) but also spoken language. You can translate speech directly into text (think interviews) or vice versa (think real-time conversations).
- **Language Buffet:** The goal? To support a massive range of languages. We're talking nearly 100 for understanding and translating what's spoken, and 35 for generating spoken languages.
- **Beyond Borders:** SeamlessM4T aspires to bridge the communication gap across cultures by enabling smooth interaction between people speaking

different languages.

Dataset

While the specific details of the datasets used to train SeamlessM4T might not be explicitly mentioned in the introduction of the research paper, we can be certain they are massive and multilingual. SeamlessM4T's ability to handle translation across nearly 100 languages necessitates a vast amount of training data in various languages, encompassing both text and speech formats. Imagine a library containing text documents, audio recordings, and their corresponding translations in dozens of languages – that's the kind of data that likely fueled SeamlessM4T's learning process. The specifics of this data (sources, sizes, etc.) are typically not public knowledge to protect training strategies, but the overall concept is clear: a mountain of multilingual text and speech data is what empowers SeamlessM4T's impressive translation capabilities.

Model Description

SeamlessM4T, developed by Meta AI, stands out as a groundbreaking machine translation model. Unlike its predecessors that handle specific translation tasks (text-to-text, speech-to-text, etc.), SeamlessM4T boasts a versatile architecture tackling them all – text-to-text (T2TT), speech-to-text (S2TT), text-to-speech (T2ST), and speech-to-speech (S2ST) translation. This "one-stop shop" approach streamlines communication by eliminating the need for separate models for each task.

Under the Hood: A Transformer Symphony

At its core, SeamlessM4T leverages a powerful transformer-based architecture. Transformers are a popular choice in machine translation due to their ability to capture long-range dependencies within sequences – crucial for understanding the nuances of language. SeamlessM4T actually utilizes two sequence-to-sequence (seq2seq) transformer models working in tandem.

Step 1: Understanding the Input - The First Transformer

The first transformer acts as an interpreter, handling the incoming data regardless of its form (text or speech). This initial stage involves dedicated encoders for each modality (text and speech) to process the input effectively. The text encoder breaks down written text into a sequence of tokens, while the speech encoder analyzes audio information to extract relevant features.

Step 2: Bridging the Language Gap - Translation into Text

Once the input is understood, the first transformer generates a textual representation of the meaning, essentially translating the content into a common

language format. This internal representation acts as a bridge between the source and target languages.

Step 3: Crafting the Output - The Second Transformer and Vocoder

Now that the meaning is captured in text form, the second transformer takes over. Its role is to generate the desired output, be it written text (T2TT) or spoken language (S2ST and T2ST). For text output, the second transformer directly generates the translated text sequence.

For speech output, however, an additional stage is necessary. Here, a vocoder comes into play. Inspired by HiFi-GAN architecture, the vocoder transforms the generated speech tokens from the second transformer into high-fidelity audio waveforms, replicating the natural flow of human speech.

SeamlessM4T v2: Taking it a Step Further

SeamlessM4T boasts two versions – v1 and v2. While the core principles remain the same, v2 introduces the innovative UnitY2 architecture specifically designed for speech generation. UnitY2 allows the model to produce speech more efficiently by working with smaller speech units compared to v1. This translates to faster speech generation without sacrificing quality.

Screenshots:

Model Block Diagram

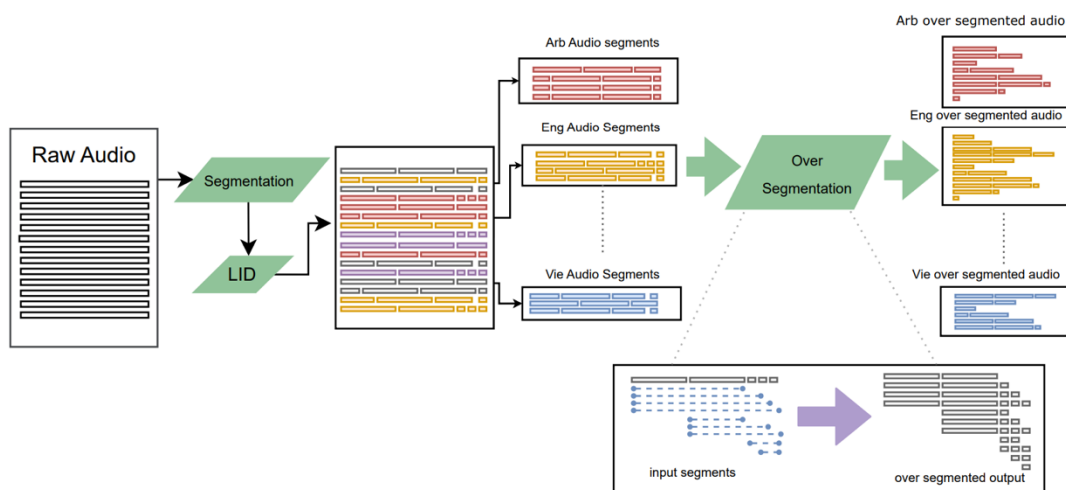


Figure 1: Workflow of speech processing.

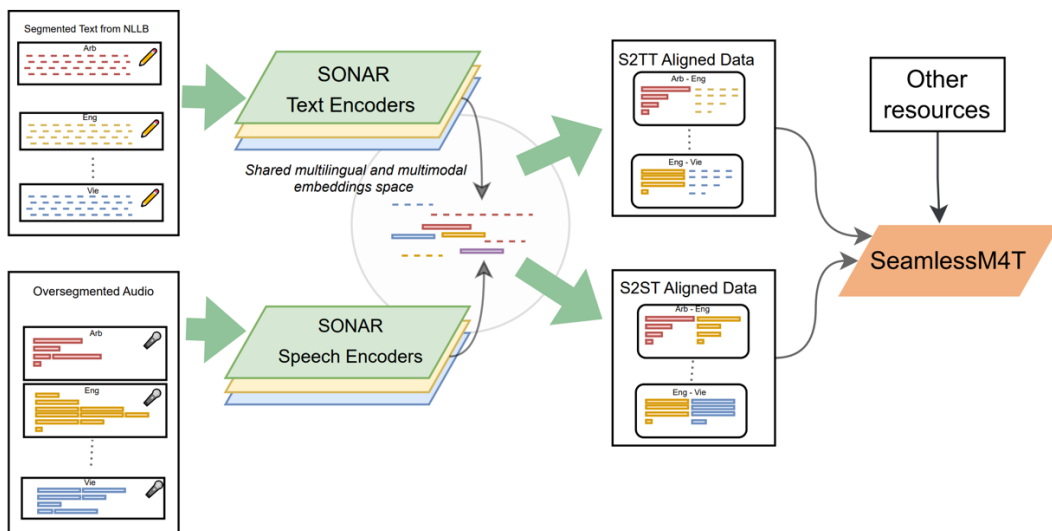
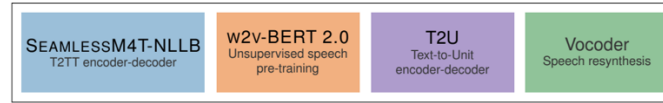


Figure 2: Workflow of the SONAR encoding and mining processes.

(1) Pre-trained models



(2) Multitasking UNITY

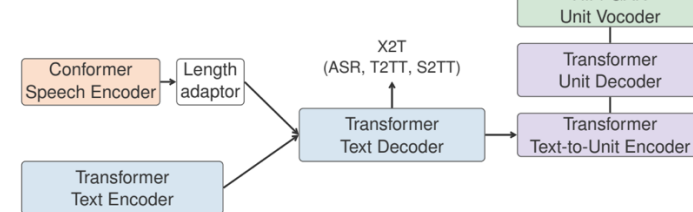
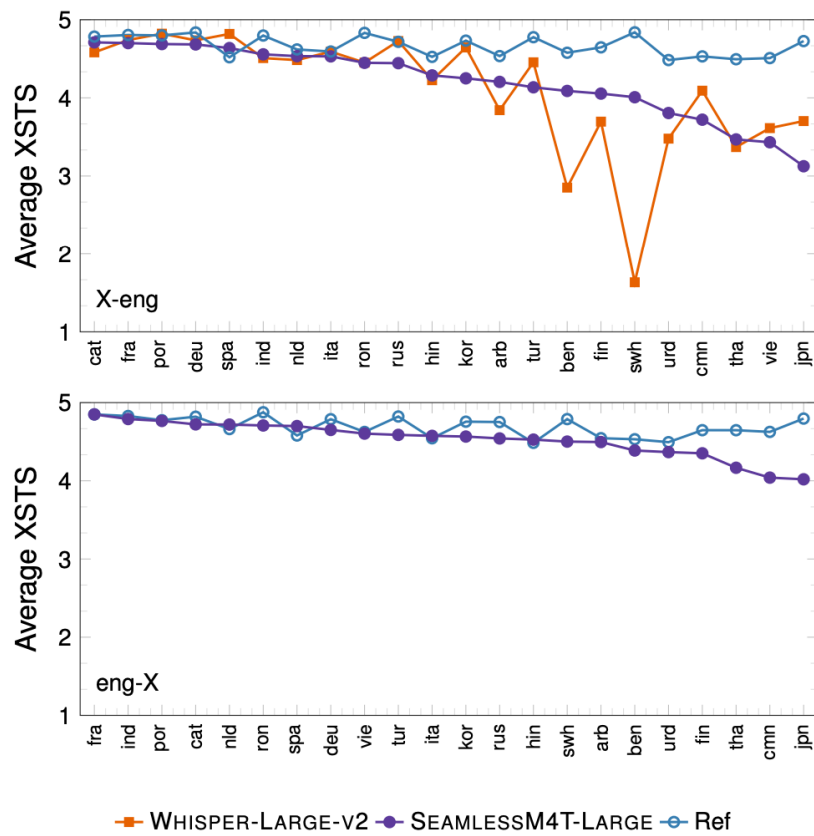
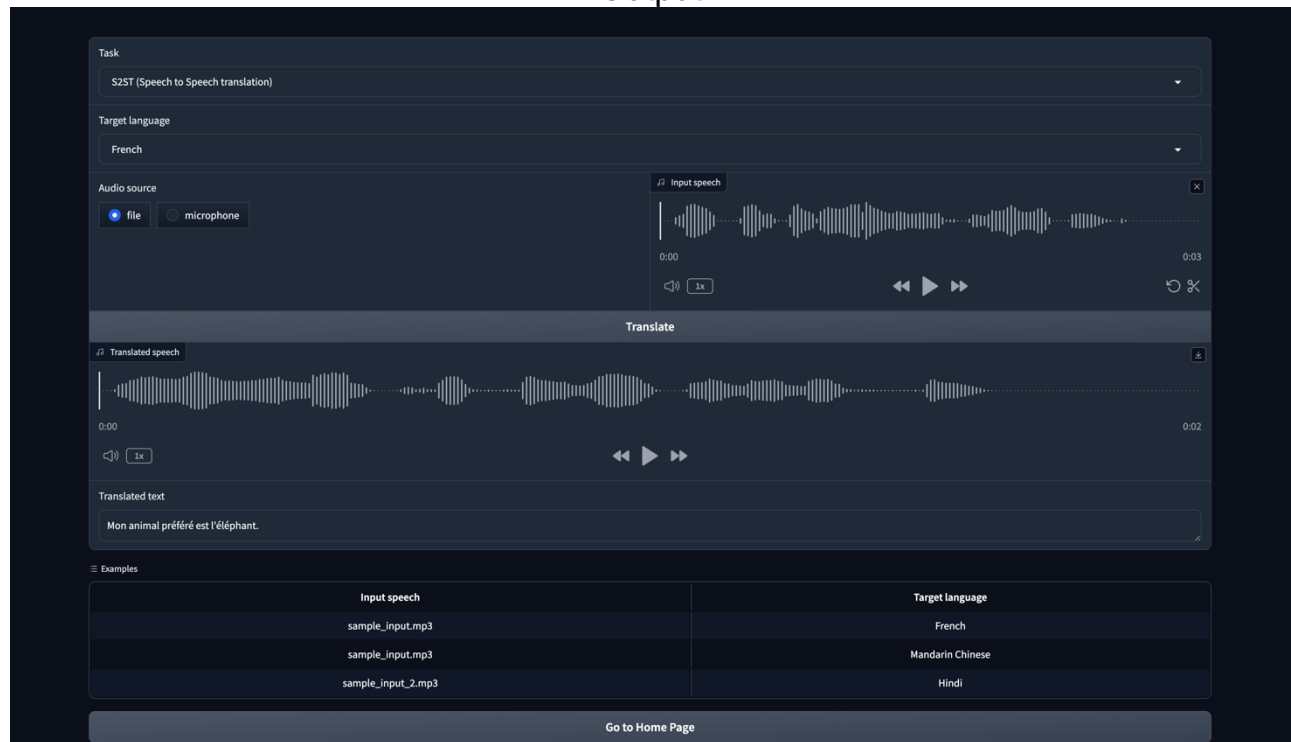


Figure 4: Overview of SEAMLESSM4T. (1) shows the pre-trained models used when finetuning multitasking UNITY. (2) outlines multitasking UNITY with its two encoders, text decoder, T2U encoder-decoder, and the supporting vocoders for synthesizing output speech in S2ST.

Comparison



Output



Conclusion:

In conclusion, SeamlessM4T stands as a testament to the power of AI in breaking down language barriers. This versatile model tackles various translation tasks, from written text to spoken conversation, across nearly 100 languages. Its core lies in a powerful transformer architecture, with two models working together. The first deciphers the input, regardless of format, while the second crafts the desired translated output – be it text or high-fidelity speech thanks to a vocoder. The v2 version further streamlines speech generation with its UnitY2 architecture. SeamlessM4T's ability to translate seamlessly positions it as a game-changer, fostering global communication and understanding across cultures. As research progresses, models like SeamlessM4T hold the key to a future where language is no longer an obstacle to human connection.