

D4.3.2 Speech Event

Due date: **31/03/2024**
Submission Date: **20/02/2025**
Revision Date: **13/03/2025**

Start date of project: **01/07/2023**

Duration: **36 months**

Lead organisation for this deliverable: **Carnegie Mellon University Africa**

Responsible Person: **Clifford Onyonka**

Revision: **1.1**

Project funded by the African Engineering and Technology Network (Afretec) Inclusive Digital Transformation Research Grant Programme		
Dissemination Level		
PU	Public	PU
PP	Restricted to other programme participants (including Afretec Administration)	
RE	Restricted to a group specified by the consortium (including Afretec Administration)	
CO	Confidential, only for members of the consortium (including Afretec Administration)	

Executive Summary

Deliverable D4.3.2 concerns the results of Task 4.3.2, a task whose objective was to train, test, and finally deploy a speech-to-text model on a ROS node that will enable speech utterances in Kinyarwanda and English languages captured by Pepper's microphones to be transcribed into written text.

This report details the output of each phase of the software development process used in the fulfillment of Deliverable D4.3.2. The requirements definition section specifies the functional requirements of the users of Speech Event, the module specification section outlines the functional characteristics of Speech Event, the interface design section outlines the specification of the inputs and outputs of Speech Event, the module design section describes the deep neural networks that perform speech recognition of Kinyarwanda and English utterances, the testing section shows the results and descriptions of unit and end-to-end tests of Speech Event, and the user manual section outlines how to build and run Speech Event.

Contents

1	Introduction	4
2	Requirements Definition	5
3	Module Specification	6
4	Interface Design	7
4.1	Source Code	7
4.2	Configuration Files	8
4.3	Data Files	9
4.4	Topics Subscribing To	10
4.5	Topics Publishing To	10
4.6	Services Supported	11
5	Module Design	12
5.1	Model Architecture	12
5.1.1	Preprocessor	12
5.1.2	Encoder Architecture	12
5.1.3	Decoder/Predictor Architecture	13
5.1.4	Joint Network Architecture	14
5.1.5	Combined Model Architecture	15
6	Testing Report	16
6.1	Unit Testing	16
6.2	End-to-end Testing	18
7	User Manual	19
7.1	Installation	19
7.2	Usage	19
7.3	Graphical User Interface	19
7.4	Driver ROS Node	20
	References	21
	Principal Contributors	22
	Document History	23

1 Introduction

Conversational systems are developed using the following main components: automatic speech recognition (speech-to-text), natural language understanding, dialogue manager, natural language generation, and text-to-speech. Speech-to-text converts an audio signal containing spoken speech utterances to written text transcriptions, natural language understanding analyses the transcribed text to extract meaning, the dialogue manager generates a response based on the inferred meaning of the text, natural language generation formulates text based on the response generated by the dialogue manager, and text-to-speech synthesises spoken speech utterances to be played to the other agent(s) in the conversation cycle [1]. Such a system is displayed in Fig 1.

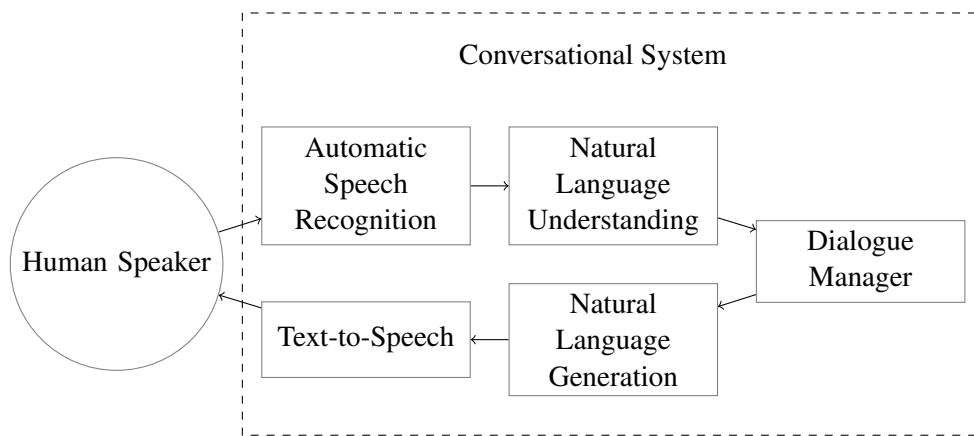


Figure 1: A diagrammatic representation of a conversational system [1]

The Speech Event module only handles one component of the conversational system described in the previous paragraph, that is, speech-to-text (which is referred to as automatic speech recognition in Fig 1). It acquires audio signals published on the `/soundDetection/signal` ROS topic, and then passes these audio signals through a deep neural network that transcribes speech utterances contained within the audio to generate text strings representations, which are then published on the `/speechEvent/text` ROS topic. Kinyarwanda and English are the two languages that are transcribed by Speech Event, with the choice of language being set by the Behaviour Controller which reads it from the Culture Knowledge Base and then passes it to Speech Event via the `/speechEvent/setLanguage` ROS service that is advertised by Speech Event.

2 Requirements Definition

A running Speech Event ROS node performs one main function - Kinyarwanda and English speech recognition. An audio signal is received via the `/soundDetection/signal` ROS topic, and Speech Event transcribes any speech utterances detected in the signal to either Kinyarwanda or English.

In order to successfully perform this function, the following functional requirements need to be fulfilled by Speech Event:

1. Acquire an audio signal from the `/soundDetection/signal` ROS topic
2. Pass the audio signal through an automatic speech recognition (ASR) model to transcribe any speech utterances the audio signal contains to either Kinyarwanda or English text
3. Publish the transcribed text to the `/speechEvent/text` ROS topic

The exact language between Kinyarwanda and English which a running Speech Event ROS node will transcribe is decided based on a configuration option that is set during the initialisation stage when the ROS node is started. The language that is set during the initialisation of the Speech Event ROS node can be overridden by invoking the `/speechEvent/setLanguage` ROS service and passing a different language to it.

To supplement the core requirement that Speech Event performs (Kinyarwanda and English speech recognition), the text transcriptions that Speech Event will transcribe will be displayed on a graphical interface. The graphical interface will be used as a replacement for the terminal interface, especially in presentation settings where displaying text transcriptions on the terminal is impractical due to the small font size used by the terminal and the extra verbosity of the results displayed on the terminal.

A ROS node that emulates Sound Detection will also be developed, and its main purpose will be to showcase the working of Speech Event in isolation, separated from the rest of the CSSR4Africa system and from the Pepper robot. This ROS node will capture audio from the computer on which Speech Event will be running on, and publishing the captured audio to the same ROS topic that Sound Detection publishes to (`/soundDetection/signal`), therefore mimicking Sound Detection's operation of acquiring audio from Pepper and publishing it to the aforementioned ROS topic.

It is also worth noting that Sound Detection may fail, and since the Behaviour Controller does not directly interface with Sound Detection, the Speech Event ROS node bears the burden of signaling failures that may arise in Sound Detection to the Behaviour Controller. This is done by publishing the message 'Error: soundDetection is down' on the `/speechEvent/text` ROS topic.

3 Module Specification

A running Speech Event ROS node performs speech-to-text in real-time, acquiring an audio signal from the `/soundDetection/signal` ROS topic, generating a text description of the acquired audio, and then publishing transcribed text to the `/speechEvent/text` ROS topic as soon as speech utterances are detected in the audio signal. The data that is input to Speech Event, the transformation that this data undergoes, and the data that Speech Event outputs is described below:

1. The language that Speech Event is to be configured to operate in is set when the ROS node is started by setting a language configuration option
2. The language set for Speech Event to operate in when it is first started can be overridden by making a `/speechEvent/setLanguage` ROS service call and passing a different language
3. When the Speech Event ROS node is running, it acquires audio signals published by Sound Detection ROS node to the `/soundDetection/signal` ROS topic
4. The acquired audio signals are then passed through an ASR model that transcribes speech utterances present within the audio signals to text strings
5. The text strings output by the ASR model are then published to the `/speechEvent/text` ROS topic for the Behaviour Controller to use

This process, and the position of Speech Event within the CSSR4Africa system architecture, is captured in Fig 2.

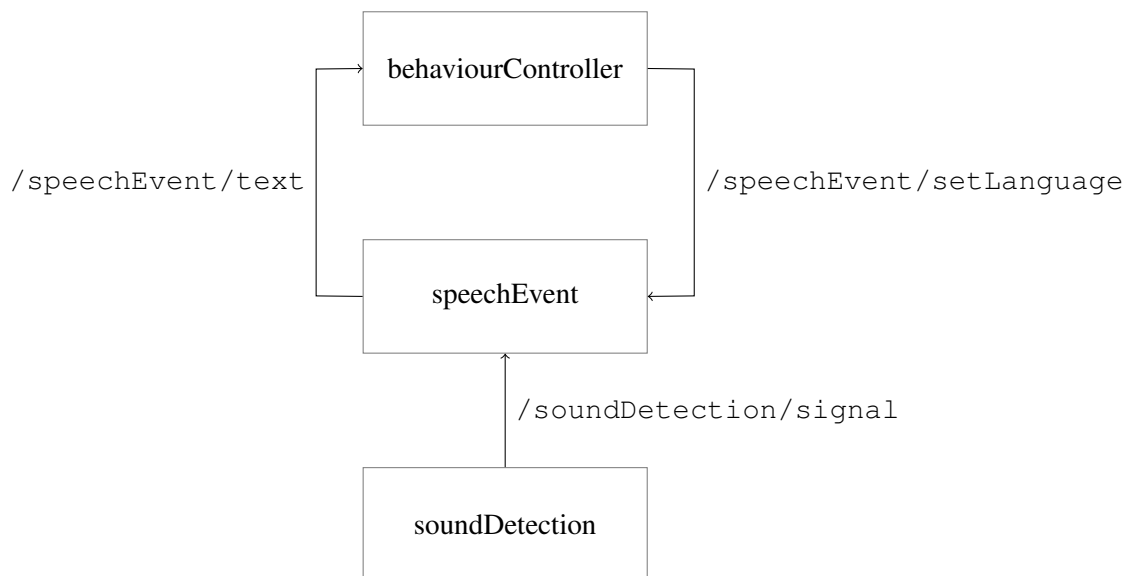


Figure 2: Speech Event within the CSSR4Africa architecture (`/soundDetection/signal` and `/speechEvent/text` are ROS topics, and `/speechEvent/setLanguage` is a ROS service)

4 Interface Design

4.1 Source Code

The file structure of Speech Event is as shown below:

```
speech_event/  
├── config/  
│   ├── speech_event_configuration.ini  
│   └── speech_event_driver.ini  
├── data/  
│   ├── pepper_topics.dat  
│   └── speech_event_input.dat  
├── gui/  
│   └── _main_.py  
├── src/  
│   ├── speech_event/  
│   │   ├── __init__.py  
│   │   ├── speech_event_config_parser.py  
│   │   └── speech_event_implementation.py  
│   ├── speech_event_application.py  
│   ├── speech_event_driver.py  
│   └── speech_event_gui.py  
├── srv/  
│   └── set_language.srv  
├── CMakeLists.txt  
├── README.md  
└── speech_event_requirements.txt
```

The `config/` directory houses configuration files that are used to configure different ROS nodes that are part of Speech Event:

1. `speech_event_configuration.ini`: configurations for the main Speech Event ROS topic
2. `speech_event_driver.ini`: configurations for a Speech Event driver ROS node that mimics Sound Detection, but captures audio from a personal computer's microphones rather than from Pepper's microphones

The `data/` directory contains data that is required by Speech Event when running:

1. `pepper_topics.dat`: list of ROS topics that Speech Event listens to when it is running
2. `speech_event_input.dat`: list of data sources that are required by Speech Event when it is running

The `gui/` directory holds a python script that enables running as a python module the GUI interface that is used to monitor the `/speechEvent/text` ROS topic. When the current working directory is the Speech Event directory, the GUI interface can be run as a python module as follows: `python3 -m gui`.

The `src/` directory holds all the python source code files that implement the Speech Event system together with the GUI application and driver application. Within the `src/` directory is the `speech_event/` subdirectory. This subdirectory is used by the `setup.py` script to configure Speech Event as a python package that can be easily installed by CATKIN when running the `catkin_make` command.

The `srv/` directory contains the `set_language.srv` service file specification for the ROS service `/speechEvent/setLanguage` used to override the language of operation set by the Speech Event ROS node when it first starts.

The file `CMakeLists.txt` is required by CATKIN to denote Speech Event as a ROS node, `README.md` contains documentation of how to use Speech Event, and the python requirements file `speech_event_requirements.txt` contains a list of versioned Python packages that Speech Event depends on.

4.2 Configuration Files

A running Speech Event ROS node requires to be configured prior to starting it, a task that is accomplished via a configuration file that is stored as plain text. The configuration parameters stored in a configuration file as displayed in Table 1.

Key	Value	Description
<code>verbose_mode</code>	<code>true, false</code>	Whether to print informational messages to the terminal
<code>cuda</code>	<code>true, false</code>	Whether to use GPU or CPU for running model inference (True means use GPU, False means use CPU)
<code>sample_rate</code>	<code><integer></code>	The sample rate of audio signals captured from the <code>/soundDetection/signal</code> ROS topic
<code>heartbeat_msg_period</code>	<code><integer></code>	The time period in seconds after which periodic heart beat messages are printed to the terminal

Table 1: Speech Event configuration file options

A running driver ROS node also requires to be configured prior to starting it. Its configuration parameters are stored in a configuration file as displayed in Table 2. (A driver ROS node mimics Sound Detection, capturing audio from a personal computer’s microphone instead of from Pepper, and publishing the audio signal to `/soundDetection/signal` topic in the same way as Sound Detection does).

Key	Value	Description
channels	<integer>	Number of output channels for the audio captured from the personal computer’s microphones
chunk_size	<integer>	Number of audio samples to be read per iteration from the personal computer’s microphones
sample_rate	<integer>	Sample rate to use when capturing audio
speech_amplitude_threshold	<float>	Threshold for segmenting silence and non-silence regions in the captured audio (non-silence regions are assumed to contain speech utterances)
utterance_time_buffer	<integer>	Number of padding samples around non-silence audio regions

Table 2: Driver configuration file options

4.3 Data Files

The topics that Speech Event subscribes to are listed in the `pepper_topics.dat` file, as shown in Table 3.

Key	Value
sound_detection	<code>/soundDetection/signal</code>

Table 3: ROS topics that Speech Event subscribes to

Data and data sources that are required by Speech Event are specified in the `speech_event_input.dat` file, as shown in Table 4.

Key	Value	Description
<code>rw_model_path</code>	<code><string></code>	Path to the Kinyarwanda ASR model
<code>en_model_path</code>	<code><string></code>	Path to the English ASR model
<code>audio_storage_dir</code>	<code><string></code>	Path to directory that will store audio captured from <code>/soundDetection/signal</code> ROS topic

Table 4: Input data and data sources that Speech Event requires

4.4 Topics Subscribing To

A running Speech Event ROS node acquires audio signals from the `/soundDetection/signal` ROS topic. Table 5 shows this fact, while also mentioning the ROS node that publishes these audio signals.

Topic	Node	Platform
<code>/soundDetection/signal</code>	<code>soundDetection</code>	Physical robot

Table 5: Topics subscribing to

4.5 Topics Publishing To

The Speech Event ROS node publishes transcribed text to the `/speechEvent/text` ROS topic, publishing each time a transcription process completes. The published text transcriptions are of the `string` data type. Table 6 shows the ROS topic that Speech Event publishes to.

Topic	Node	Platform
<code>/speechEvent/text</code>	<code>behaviourController</code>	Physical robot

Table 6: Topics publishing to

4.6 Services Supported

A running Speech Event ROS node can have its language of operation (Kinyarwanda or English) updated by invoking the `/speechEvent/setLanguage` ROS service. Table 7 elaborates on this ROS service.

Service	Message Value	Effect
<code>/speechEvent/setLanguage</code>	<code>kinyarwanda, english</code>	Set the language to either Kinyarwanda or English

Table 7: Services supported

5 Module Design

5.1 Model Architecture

Both the Kinyarwanda and English speech recognition processes performed by Speech Event rely on deep learning models. Both of these deep learning models are based on the conformer transducer architecture, and the specific models that are used by Speech Event were acquired from NVIDIA's Nemo catalog of models.

The conformer transducer architecture is made up of a conformer encoder and a transducer decoder. The conformer encoder combines transformers and Convolution Neural Networks (CNNs) in an attempt to utilise the strengths of both, with transformers excelling at capturing global dependencies of an audio signal and CNNs excelling at capturing local dependencies of an audio signal [2]. The transducer decoder, on the other hand, makes use of Recurrent Neural Networks (RNNs) to form an autoregressive model whose next token prediction is conditioned on the previously predicted tokens.

5.1.1 Preprocessor

When an audio signal is received, it is first passed through a preprocessor block that extracts filterbank features from the audio signal. For this process, an audio signal is first converted to a spectrogram by a Short Time Fourier Transform (STFT), and then a series of filters are applied to the resultant spectrogram to obtain the filterbank features that each represent the magnitude of the energy within a distinct frequency band. This preprocessor block is visualised in Fig 3.

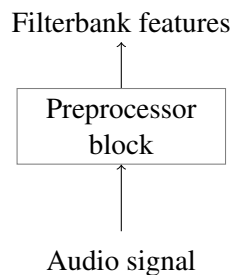


Figure 3: Audio to filterbank features preprocessor

5.1.2 Encoder Architecture

The filterbank features obtained from the preprocessor block are then passed to the encoder part of the model. The encoder transforms acoustic features in the original audio signal captured within the filterbank features into latent features that better represent the speech in the original audio signal.

A conformer encoder is shown in Fig 4. The encoder used in Speech Event has 17 conformer blocks.

The filterbank features are first passed through a SpecAugment block. SpecAugment is a data augmentation method that is suitable for end-to-end speech recognition models such as the Conformer Transducer model used by Speech Event. "The augmentation policy consists of warping the features, masking blocks of frequency channels, and masking blocks of time steps" [3].

The features are then processed by a convolution subsampling block, dropout applied, and then passed through the series of 17 conformer blocks in the encoder.

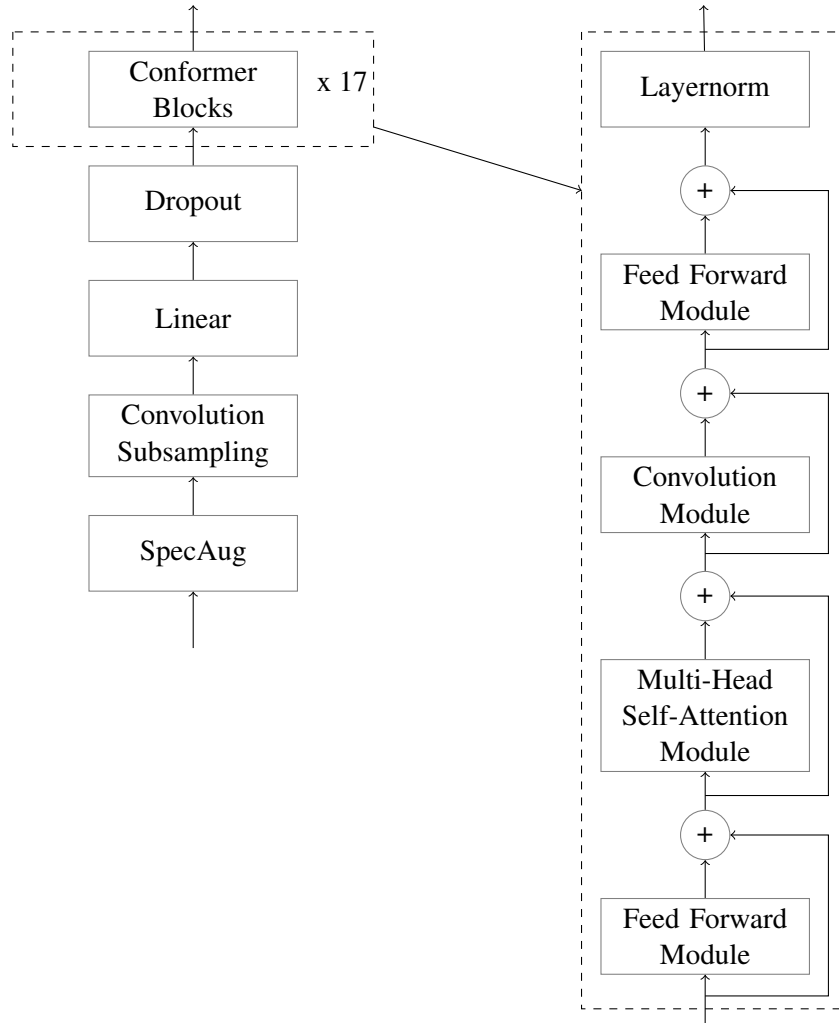


Figure 4: Conformer encoder architecture [2]

Each conformer block consists of a multi-head self-attention module and a convolution module collectively sandwiched between two feedforward modules [2]. The self-attention module is more adept at capturing global feature dependencies while the convolution module is more adept at capturing local feature dependencies, thereby capturing the semantics of speech utterances much better than models that only specialise in capturing only either of the global or local feature dependencies.

5.1.3 Decoder/Predictor Architecture

Keeping in line with Gulati et al., a single LSTM layer is used in the decoder [2]. An embedding layer converts input tokens into vectorised representations that are then passed to the LSTM layer. Additionally, drop out is applied both in the LSTM layer and on the output of the LSTM layer.

The decoder architecture used in Speech Event is shown in Fig 5.

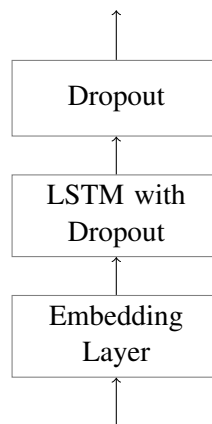


Figure 5: Transducer predictor architecture

5.1.4 Joint Network Architecture

The joint network combines the output of the encoder and the output of the decoder in order to predict the text in the speech utterances contained in the original sound signal passed to the encoder. The combined outputs are then passed through a ReLU activation, dropout applied, and lastly passed through a linear layer.

The joint network architecture used in Speech Event is shown in Fig 6.

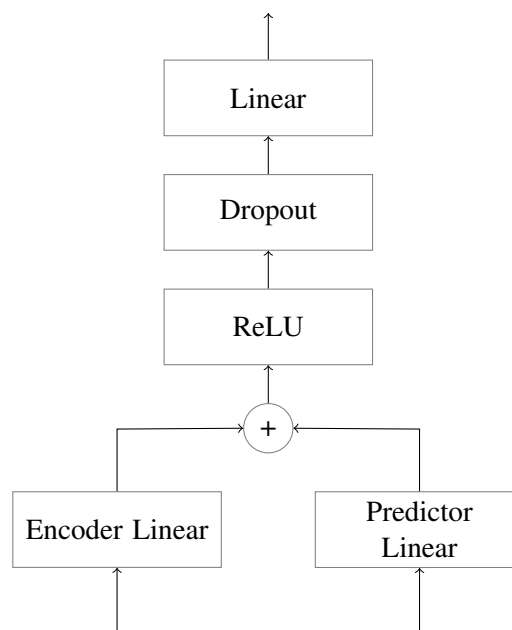


Figure 6: Joint architecture

5.1.5 Combined Model Architecture

The four different sections of the full conformer transducer architecture discussed above are combined together to form the whole network that performs end-to-end automatic speech recognition. This full model's architecture is summarised in Fig 7.

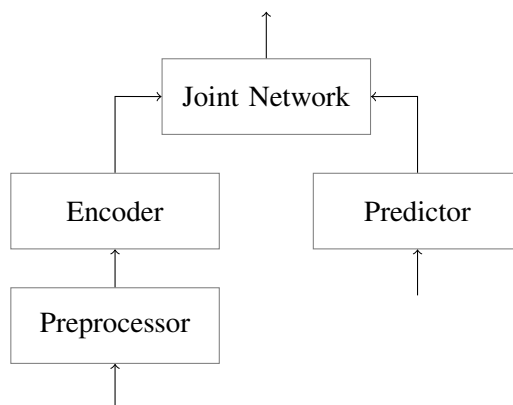


Figure 7: Combined model architecture

6 Testing Report

The `unit_tests` package contains automated unit tests and both automated and manual end-to-end tests for Speech Event (together with unit tests for other ROS nodes that are part of CSSR4Africa). The unit tests test in isolation functions that implement Speech Event, while the end-to-end tests test the whole Speech Event as a whole (from when audio signals are captured from the `/soundDetection/signal` ROS topic to when transcribed text is published on the `/speechEvent/text` ROS topic).

The unittest python testing framework is used for the automated tests. The unittest package is part of the standard python library, and therefore no extra packages need to be installed when running tests.

- To run automated tests: `python3 -m unittest` (the `-v` flag is appended to the command when increased verbosity is required, and a ROS master needs to be running on a separate terminal using `roscore`)

The manual tests are implemented as a ROS node that is run and the test cases checked manually by a human tester. To run these manual tests, the Pepper robot can be used, or they can be run without the Pepper robot by substituting both the Pepper robot and the Sound Detection ROS node with a driver ROS node that emulates their functionality. When running these tests, the most important terminal logs that are linked to the tests being run are printed out in green and blue font colours to make them stand out from the rest of the terminal logs that are not related to the tests being run.

- To run manual tests using the Pepper robot: `roslaunch unit_tests speech_event_launch_robot.launch`
- To run manual tests using a driver ROS node: `roslaunch unit_tests speech_event_launch_test_harness.launch`

6.1 Unit Testing

Four unit tests are provided in Speech Event. The first unit test tests the initialisation function (the function that, at the start of a Speech Event ROS node, reads a configuration file and sets up Speech Event to operate with required settings such as the language to be transcribed).

```
$ python3 -m unittest tests/test_speech_event.py -k test_initialise
-v
test_initialise (tests.test_speech_event.TestSpeechEvent)
Function speech_event.speech_event_implementation.initialise(config:
dict) -> None ... ok
```

Ran 1 test in 4.172s

OK

The second unit test tests the save audio function (this function saves an audio signal to storage as a .wav file, and this wav file is then passed to the ASR model that transcribes speech in the next phase of the speech-to-text process).


```
$ python3 -m unittest tests/test_speech_event.py -k test_save_audio -v
test_save_audio (tests.test_speech_event.TestSpeechEvent)
Function speech_event.speech_event_implementation.save_audio(sample_rate: int, samples: np.ndarray) -> str ... ok
```

```
Ran 1 test in 2.009s
```

OK

The third unit test tests the get audio transcription function (this function uses an ASR deep learning model to transcribe an audio signal, generating a text string representation of any speech utterances contained within the audio signal).

```
$ python3 -m unittest tests/test_speech_event.py -k test_get_audio_transcription -v
test_get_audio_transcription (tests.test_speech_event.TestSpeechEvent)
Function speech_event.speech_event_implementation.get_audio_transcription(filename: str) -> str ... ok
```

```
Ran 1 test in 2.997s
```

OK

The last unit test tests the parse function (this function parses an ini configuration file and generates its python dictionary representation, and this function is re-used everywhere an ini configuration file needs to be read and parsed - the ini specification used by CSSR4Africa does not follow the standardised ini file specification, and therefore cannot be parsed by standard ini libraries provided by the standard python library, which necessitates creation of a custom ini file parser).

```
$ python3 -m unittest tests/test_config_parser.py -k test_parse -v
test_parse (tests.test_config_parser.TestConfigParser)
Function speech_event.speech_event_config_parser.ConfigParser.parse(config_file_path: str) -> dict ... ok
```

```
Ran 1 test in 0.000s
```

OK

6.2 End-to-end Testing

One end-to-end test is provided in Speech Event. A running ROS master node is needed to run this test because this test necessitates creation of a full-fledged ROS node when it is run, and therefore `roscore` needs to be run on a separate terminal before running this test. It tests the whole Speech Event ROS node as a unit, from the capture of audio signals from the `/soundDetection/signal` ROS topic to the final step of publishing transcribed text to the `/speechEvent/text` ROS topic.

```
$ python3 -m unittest tests/test_speech_event.py -k test_e2e -v
test_e2e (tests.test_speech_event.TestSpeechEvent)
Transcribing: 100%|#####| 1/1 [00:01<00:00, 1.16s/it]
Transcribing: 100%|#####| 1/1 [00:00<00:00, 1.05s/it]
Transcribing: 100%|#####| 1/1 [00:00<00:00, 1.19s/it]
Transcribing: 100%|#####| 1/1 [00:00<00:00, 1.94s/it]
ok
```

```
-----
Ran 1 test in 14.281s
```

```
OK
```

7 User Manual

7.1 Installation

The Speech Event package needs to be installed before it can be used. To install Speech Event, clone the Speech Event repository to the CSSR4Africa workspace, and then proceed with the following steps:

1. Install required packages
 - (a) Install required Ubuntu packages: `sudo apt-get install cython3 ffmpeg gfortran libopenblas-dev libopenblas64-dev patchelf pkg-config python3-testresources python3-typing-extensions sox`
 - (b) Install required python packages: `pip3 install -r requirements.txt` (the `requirements.txt` file is located in the Speech Event root directory)
2. Install Speech Event: `roscd && cd ../ && catkin_make`

7.2 Usage

To run a Speech Event ROS node, a Sound Detection ROS node needs to also be running (or a Speech Event driver node, which is described in a subsequent subsection).

1. Run a Speech Event ROS node: `roslaunch speech_event speech_event_application.py -c CONFIG` (where `CONFIG` is the path to a configuration ini file; an editable configuration file is provided in the `config/` directory)
2. View text transcriptions on the terminal (optional): `rostopic echo /speechEvent/text`

7.3 Graphical User Interface

A graphical interface is provided to display the text transcriptions that are being published on the `/speechEvent/text` ROS topic on a more user friendly interface compared to the terminal interface. To view text transcriptions on this graphical interface:

1. Install Tk: `sudo apt-get install python3-tk`
2. Run the GUI application: `roslaunch speech_event speech_event_gui.py` (or `python3 -m gui` if the current working directory is set to the Speech Event root directory)

7.4 Driver ROS Node

A driver ROS node that mimics the Sound Detection ROS node is also provided. It captures audio from a personal computer's microphone instead of from Pepper, and publishes the audio signal on the `/soundDetection/signal` ROS topic in the same way that the Sound Detection ROS node does. To bring up this driver ROS node:

1. Install required Ubuntu packages: `sudo apt-get install libasound-dev portaudio19-dev libportaudio2 libportaudiocpp0 libav-tools`
2. Run a driver ROS node: `roslaunch speech_event speech_event_driver.py -c CONFIG` (where CONFIG is the path to a configuration ini file; an editable configuration file is provided in the `config/` directory)

References

- [1] Cristina Romero-González, Jesus Martínez-Gómez, and Ismael García-Varea. Spoken language understanding for social robotics. In *2020 IEEE International Conference on Autonomous Robot Systems and Competitions (ICARSC)*, pages 152–157, April 2020.
- [2] Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo Wang, Zhengdong Zhang, Yonghui Wu, and Ruoming Pang. Conformer: Convolution-augmented Transformer for Speech Recognition. In *Interspeech 2020*, pages 5036–5040. ISCA, October 2020.
- [3] Daniel S. Park, William Chan, Yu Zhang, Chung-Cheng Chiu, Barret Zoph, Ekin D. Cubuk, and Quoc V. Le. SpecAugment: A Simple Data Augmentation Method for Automatic Speech Recognition. In *Interspeech 2019*, pages 2613–2617. ISCA, September 2019.

Principal Contributors

The main authors of this deliverable are as follows (in alphabetical order).

Clifford Onyonka, Carnegie Mellon University Africa.

David Vernon, Carnegie Mellon University Africa.

Richard Muhirwa, Carnegie Mellon University Africa.

Document History

Version 1.0

First draft.

Clifford Onyonka.

20 February 2025.

Version 1.1

Added a ROS service that will be used by the Behaviour Controller to set the language (Kinyarwanda or English) that Speech Event will operate in

Clifford Onyonka.

13 March 2025.