

---

# **Riva End-to-end User Guide**

**spoonnvidia**

**Oct 11, 2021**



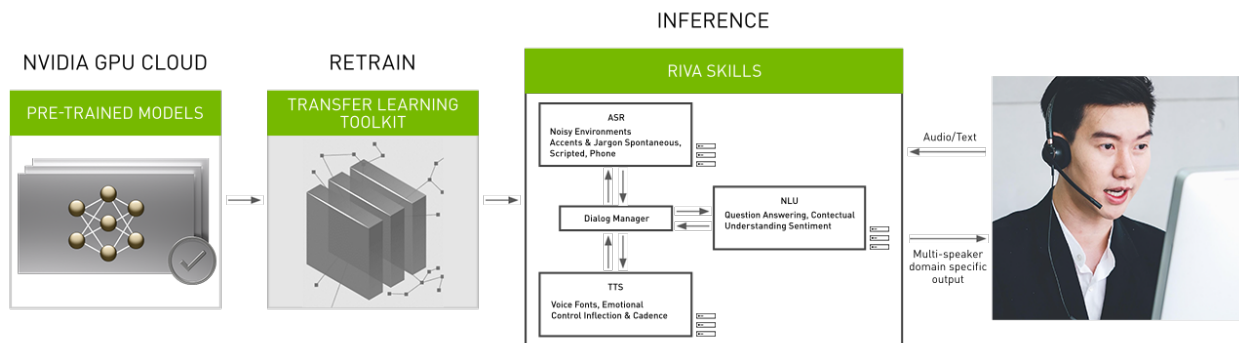
<b>1</b>	<b>Riva - Getting Started</b>	<b>1</b>
1.1	Prerequisites . . . . .	2
1.2	Sample service usage . . . . .	3
<b>2</b>	<b>Riva - ServiceMaker</b>	<b>5</b>
2.1	riva-build . . . . .	5
2.2	riva-deploy . . . . .	6
<b>3</b>	<b>Riva - Server Configurations</b>	<b>7</b>
3.1	Sample config.sh . . . . .	7
<b>4</b>	<b>Riva - Speech Recognition</b>	<b>13</b>
<b>5</b>	<b>Riva - Speech Recognition - Model Architectures</b>	<b>15</b>
5.1	Jasper . . . . .	15
5.2	QuartzNet . . . . .	15
5.3	CitriNet . . . . .	15
<b>6</b>	<b>Riva - Speech Recognition - Deploy ASR Model</b>	<b>17</b>
6.1	riva-build . . . . .	17
6.2	riva-deploy . . . . .	18
<b>7</b>	<b>Riva - Natural Language Processing</b>	<b>19</b>
<b>8</b>	<b>Riva - Rasa</b>	<b>21</b>
8.1	Prerequisites . . . . .	21
8.2	Sample Chatbot . . . . .	22
<b>9</b>	<b>Riva - Rasa - NLU component</b>	<b>23</b>
9.1	Sample nlu.yml . . . . .	23
<b>10</b>	<b>Riva - Rasa - Stories (Dialog Manager)</b>	<b>27</b>
10.1	Sample stories.yml . . . . .	27
<b>11</b>	<b>Riva - Rasa - Actions</b>	<b>29</b>
<b>12</b>	<b>Riva - Rasa - Domain</b>	<b>31</b>
12.1	Sample domain.yml . . . . .	31
<b>13</b>	<b>Riva - Rasa - Configurations</b>	<b>33</b>
<b>14</b>	<b>Riva - Rasa - Model Training</b>	<b>35</b>

<b>15 Riva - Rasa - Credentials</b>	<b>37</b>
<b>16 Riva - Rasa - Endpoints</b>	<b>39</b>
<b>17 Riva - Rasa - Conversation Test</b>	<b>41</b>
<b>18 Riva - Speech Synthesis</b>	<b>43</b>
<b>19 Riva - Speech Synthesis - Model Architectures</b>	<b>45</b>
<b>20 Riva - Custom Model</b>	<b>47</b>
20.1 TAO . . . . .	47
20.2 NeMo . . . . .	47
20.3 Export models . . . . .	48
<b>21 Riva - Custom Model - TAO</b>	<b>49</b>
21.1 TAO Launcher . . . . .	49
<b>22 Riva - Custom Model - TAO - speech-to-text</b>	<b>51</b>
22.1 Usage . . . . .	51
<b>23 Riva - Custom Model - NeMo</b>	<b>53</b>
23.1 Prerequisites . . . . .	54
23.2 Getting started . . . . .	54
<b>24 Riva - Custom Model - NeMo - Data Manifest - ASR</b>	<b>55</b>
<b>25 Riva - Custom Model - NeMo - ASR - Mandarin</b>	<b>57</b>
25.1 0. Launch PyTorch container with NeMo . . . . .	57
25.2 1. Prepare data manifest . . . . .	58
25.3 2. Prepare training configuration . . . . .	59
25.4 3. Train QuartzNet 15x5 zh . . . . .	61
25.5 4. Visualise Training Progress . . . . .	62
<b>26 Riva - Custom Model - NeMo - Data Format - NLP (Intent-slot)</b>	<b>67</b>
26.1 Expected structure . . . . .	67
26.2 Sample dict.intents.csv . . . . .	67
26.3 Sample dict.slots.csv . . . . .	68
26.4 Sample train.tsv . . . . .	68
26.5 Sample train_slots.tsv . . . . .	68
26.6 Sample test.tsv . . . . .	68
26.7 Sample test_slots.tsv . . . . .	69
<b>27 Riva - Custom Model - NeMo - Data Format - NLP (NER)</b>	<b>71</b>
<b>28 Riva - Custom Model - NeMo - Data Format - NLP (Text Classification)</b>	<b>73</b>
<b>29 Riva - Custom Model - NeMo - TTS - Inference</b>	<b>75</b>
29.1 Two-stage Approach . . . . .	75
29.2 End-to-end Approach . . . . .	76
<b>30 Riva - Custom Model - NeMo - Data Manifest - TTS</b>	<b>79</b>
<b>31 Riva - Custom Model - NeMo - TTS - Mandarin</b>	<b>81</b>
<b>32 Riva - gRPC API</b>	<b>83</b>
32.1 riva_asr.proto . . . . .	83

32.2	riva_nlp.proto . . . . .	89
32.3	riva_tts.proto . . . . .	94
32.4	riva_audio.proto . . . . .	95
32.5	health.proto . . . . .	96
<b>33</b>	<b>Riva - Python API - Sample Services</b>	<b>97</b>
33.1	Prerequisites . . . . .	97
33.2	ASR - Recognize . . . . .	97
33.3	NLP - Core Services . . . . .	102
33.4	NLP - AnalyzeIntent . . . . .	107
33.5	TTS - SynthesizeSpeech . . . . .	109
<b>34</b>	<b>Riva - Python Client</b>	<b>111</b>
34.1	asr_client.py . . . . .	111
34.2	Riva ASR service . . . . .	113
<b>35</b>	<b>Riva - Streaming ASR</b>	<b>115</b>
<b>36</b>	<b>Riva - Virtual Assistant</b>	<b>117</b>
36.1	Start Riva services . . . . .	117
36.2	Download sample image . . . . .	118
36.3	Run Riva service . . . . .	118
36.4	Start Virtual Assistant app . . . . .	118
<b>37</b>	<b>Riva - Virtual Assistant - Rasa</b>	<b>119</b>
37.1	Overview of RVA Rasa . . . . .	120
37.2	Prerequisites . . . . .	122
37.3	1. Start Riva Server . . . . .	125
37.4	2. Start Rasa Action Server . . . . .	125
37.5	3. Start Rasa Server . . . . .	125
37.6	4. Start RVA Rasa Server . . . . .	126
37.7	5. Open RVA Rasa . . . . .	127
<b>38</b>	<b>Guide Description</b>	<b>129</b>
<b>39</b>	<b>Author</b>	<b>131</b>



## RIVA - GETTING STARTED



*NVIDIA Riva* is a GPU-accelerated SDKs which we can use to build conversational AI applications using deep learning. Riva platform offers power features including

1. Complete set of tools to build conversational AI combining services of speech recognition (ASR), Natural Language Understanding/Processing (NLU/NLP), and text speech\_synthesis (TTS)
2. Complete workflow to build, train, and deploy models in Riva using NVIDIA TAO, NVIDIA NeMo.
3. Optimised deployed models to deliver greater performance and latency.
4. Client integration using gRPC API

## 1.1 Prerequisites

---

**Note:** Follow below to launch Riva server.

---

1. follow [CLI Install Guide](#) to set up NGC Catalog API
2. Download Riva quick start scripts

```
ngc registry resource download-version "nvidia/riva/riva_quickstart:1.5.0-beta"
```

3. Initialise and start Riva

```
cd riva_quickstart_v1.5.0-beta
bash riva_init.sh
bash riva_start.sh
```

If `bash riva_start.sh` was successfully executed, we will see

```
Waiting for Riva server to load all models...retrying in 10 seconds
Waiting for Riva server to load all models...retrying in 10 seconds
Waiting for Riva server to load all models...retrying in 10 seconds
Riva server is ready...
```

That's it! Riva server is set up.

---

**Hint:** For checking Riva Available services, run `docker logs riva-speech`, and you should see output similar to below outputs:

Model	Version	Status
riva-asr	1	READY
riva-asr-ctc-decoder-cpu-streaming	1	READY
riva-asr-feature-extractor-streaming	1	READY
riva-asr-voice-activity-detector-ctc-streaming	1	READY
riva-trt-riva-asr-am-streaming	1	READY

---

**Note:** Configurations of Riva Server (e.g. models, services, ports) can be configured in `config.sh`. See [Riva - Server Configurations](#) for more details.

---

To stop running Riva server, run

```
bash riva_stop.sh
```



## 1.2 Sample service usage

You might start a container with sample clients for each service

```
bash riva_start_client.sh
```

Try different services using provided notebooks

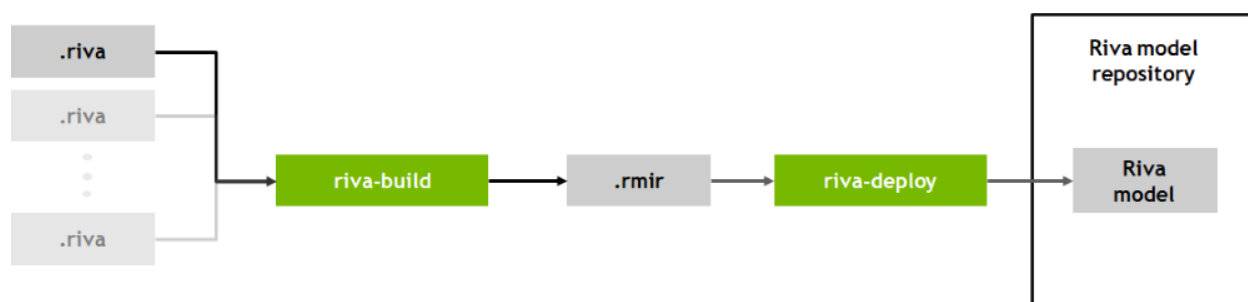
```
jupyter notebook --ip=0.0.0.0 --allow-root --notebook-dir=/work/notebooks
```

**See also:**

- [NVIDIA Riva documentation](#)
- [Stanford CS224N: NLP with Deep Learning](#)



## RIVA - SERVICEMAKER



After exporting the models using TAO or Nemo, *Riva ServiceMaker* could be used to deploy exported models in Riva server.

---

**Note:** See *Riva - Custom Model* for more details of exporting custom AI models in `.riva` formats for deployment.

---

### 2.1 riva-build

`riva-build` tool takes 1 or more `.riva` models as input and outputs a `.rmir` (intermediate format) file.

Syntax of `riva-build`

```
riva-build <task-name> --decoder_type=<decoder> output-dir-for-rmir/model.rmir:key dir-
↪for-riva/model.riva:key
```

Table 1: Currently supported tasks (Riva services) by `riva-build` tool

Tasks	Corresponding services
<code>speech_recognition</code>	ASR
<code>speech_synthesis</code>	TTS
<code>qa</code>	Question answering
<code>token_classification</code>	Token level classification
<code>intent_slot</code>	Joint intent
<code>text_classification</code>	-
<code>punctuation</code>	-

---

**Note:** See *Riva Build* for more details.

---

## 2.2 riva-deploy

riva-deploy tool takes 1 or more .rmir files and a Riva model repo as inputs.

Syntax of riva-deploy

```
riva-deploy /servicemaker-dev/<rmir_filename>:<encryption_key> /data/models
```

---

**Note:** See [Riva Deploy](#) for more details.

---

## RIVA - SERVER CONFIGURATIONS

Configure Riva Server properties by modifying `config.sh` from Riva quick start scripts

---

**Note:** Refer *Riva - Getting Started* to download Riva Quick Start scripts.

---

### 3.1 Sample `config.sh`

```
1 # Copyright (c) 2021, NVIDIA CORPORATION. All rights reserved.
2 #
3 # NVIDIA CORPORATION and its licensors retain all intellectual property
4 # and proprietary rights in and to this software, related documentation
5 # and any modifications thereto. Any use, reproduction, disclosure or
6 # distribution of this software and related documentation without an express
7 # license agreement from NVIDIA CORPORATION is strictly prohibited.
8
9 # Enable or Disable Riva Services
10 service_enabled_asr=true
11 service_enabled_nlp=true
12 service_enabled_tts=true
13
14 # Specify one or more GPUs to use
15 # specifying more than one GPU is currently an experimental feature, and may result in
16 # undefined behaviours.
17 gpus_to_use="device=0"
18
19 # Specify the encryption key to use to deploy models
20 MODEL_DEPLOY_KEY="tlt_encode"
21
22 # Locations to use for storing models artifacts
23 #
24 # If an absolute path is specified, the data will be written to that location
25 # Otherwise, a docker volume will be used (default).
26 #
27 # riva_init.sh will create a `rmir` and `models` directory in the volume or
28 # path specified.
29 #
30 # RMIR ($riva_model_loc/rmir)
31 # Riva uses an intermediate representation (RMIR) for models
```

(continues on next page)

(continued from previous page)

```

31 # that are ready to deploy but not yet fully optimized for deployment. Pretrained
32 # versions can be obtained from NGC (by specifying NGC models below) and will be
33 # downloaded to $riva_model_loc/rmir by `riva_init.sh`
34 #
35 # Custom models produced by NeMo or TLT and prepared using riva-build
36 # may also be copied manually to this location $(riva_model_loc/rmir).
37 #
38 # Models ($riva_model_loc/models)
39 # During the riva_init process, the RMIR files in $riva_model_loc/rmir
40 # are inspected and optimized for deployment. The optimized versions are
41 # stored in $riva_model_loc/models. The riva server exclusively uses these
42 # optimized versions.
43 riva_model_loc="riva-model-repo"
44
45 # The default RMIRs are downloaded from NGC by default in the above $riva_rmir_loc_
46 ↪ directory
47 # If you'd like to skip the download from NGC and use the existing RMIRs in the $riva_
48 ↪ rmir_loc
49 # then set the below $use_existing_rmirs flag to true. You can also deploy your set of_
50 ↪ custom
51 # RMIRs by keeping them in the riva_rmir_loc dir and use this quickstart script with the
52 # below flag to deploy them all together.
53 use_existing_rmirs=false
54
55 # Ports to expose for Riva services
56 riva_speech_api_port="50051"
57 riva_vision_api_port="60051"
58
59 # NGC orgs
60 riva_ngc_org="nvidia"
61 riva_ngc_team="riva"
62 riva_ngc_image_version="1.4.0-beta"
63 riva_ngc_model_version="1.4.0-beta"
64
65 # Pre-built models listed below will be downloaded from NGC. If models already exist in
66 ↪ $riva-rmir
67 # then models can be commented out to skip download from NGC
68
69 ##### ASR MODELS #####
70
71 models_asr=(
72 ### Punctuation model
73 # "${riva_ngc_org}/${riva_ngc_team}/rmir_nlp_punctuation_bert_base:${riva_ngc_model_
74 ↪ version}"
75
76 ### Citrinet-1024 Streaming w/ CPU decoder, best latency configuration
77 # "${riva_ngc_org}/${riva_ngc_team}/rmir_asr_citrinet_1024_asrset1p7_streaming:$
78 ↪ {riva_ngc_model_version}"
79
80 ### Citrinet-1024 Streaming w/ CPU decoder, best throughput configuration
81 # "${riva_ngc_org}/${riva_ngc_team}/rmir_asr_citrinet_1024_asrset1p7_streaming_
82 ↪ throughput:${riva_ngc_model_version}"

```

(continues on next page)

(continued from previous page)

```

76
77 ### Citrinet-1024 Offline w/ CPU decoder,
78 # "${riva_ngc_org}/${riva_ngc_team}/rmir_asr_citrinet_1024_asrset1p7_offline:${riva_
↪ngc_model_version}"
79
80 ### Jasper Streaming w/ CPU decoder, best latency configuration
81 # "${riva_ngc_org}/${riva_ngc_team}/rmir_asr_jasper_english_streaming:${riva_ngc_
↪model_version}"
82
83 ### Jasper Streaming w/ CPU decoder, best throughput configuration
84 # "${riva_ngc_org}/${riva_ngc_team}/rmir_asr_jasper_english_streaming_throughput:$
↪{riva_ngc_model_version}"
85
86 ### Jasper Offline w/ CPU decoder
87 # "${riva_ngc_org}/${riva_ngc_team}/rmir_asr_jasper_english_offline:${riva_ngc_model_
↪version}"
88
89 ### Quartznet Streaming w/ CPU decoder, best latency configuration
90 # "${riva_ngc_org}/${riva_ngc_team}/rmir_asr_quartznet_english_streaming:${riva_ngc_
↪model_version}"
91
92 ### Quartznet Streaming w/ CPU decoder, best throughput configuration
93 # "${riva_ngc_org}/${riva_ngc_team}/rmir_asr_quartznet_english_streaming_throughput:$
↪{riva_ngc_model_version}"
94
95 ### Quartznet Offline w/ CPU decoder
96 # "${riva_ngc_org}/${riva_ngc_team}/rmir_asr_quartznet_english_offline:${riva_ngc_
↪model_version}"
97
98 ### Jasper Streaming w/ GPU decoder, best latency configuration
99 # "${riva_ngc_org}/${riva_ngc_team}/rmir_asr_jasper_english_streaming_gpu_decoder:$
↪{riva_ngc_model_version}"
100
101 ### Jasper Streaming w/ GPU decoder, best throughput configuration
102 # "${riva_ngc_org}/${riva_ngc_team}/rmir_asr_jasper_english_streaming_throughput_gpu_
↪decoder:${riva_ngc_model_version}"
103
104 ### Jasper Offline w/ GPU decoder
105 # "${riva_ngc_org}/${riva_ngc_team}/rmir_asr_jasper_english_offline_gpu_decoder:$
↪{riva_ngc_model_version}"
106 )
107
108 ##### NLP MODELS #####
109
110 models_nlp=(
111 ### Bert base Punctuation model
112 # "${riva_ngc_org}/${riva_ngc_team}/rmir_nlp_punctuation_bert_base:${riva_ngc_model_
↪version}"
113
114 ### BERT base Named Entity Recognition model fine-tuned on GMB dataset with class labels_
↪LOC, PER, ORG etc.
115 # "${riva_ngc_org}/${riva_ngc_team}/rmir_nlp_named_entity_recognition_bert_base:$
↪{riva_ngc_model_version}"

```

(continues on next page)

(continued from previous page)

```

116 ### BERT Base Intent Slot model fine-tuned on weather dataset.
117 # "${riva_ngc_org}/${riva_ngc_team}/rmir_nlp_intent_slot_bert_base:${riva_ngc_model_
118 ↪version}"
119
120 ### BERT Base Question Answering model fine-tuned on Squad v2.
121 # "${riva_ngc_org}/${riva_ngc_team}/rmir_nlp_question_answering_bert_base:${riva_ngc_
122 ↪model_version}"
123
124 ### Megatron345M Question Answering model fine-tuned on Squad v2.
125 # "${riva_ngc_org}/${riva_ngc_team}/rmir_nlp_question_answering_megatron:${riva_ngc_
126 ↪model_version}"
127
128 ### Bert base Text Classification model fine-tuned on 4class (weather, meteorology,
129 ↪personality, nomatch) domain model.
130 # "${riva_ngc_org}/${riva_ngc_team}/rmir_nlp_text_classification_bert_base:${riva_
131 ↪ngc_model_version}"
132 )
133
134 ##### TTS MODELS #####
135
136 models_tts=(
137 # "${riva_ngc_org}/${riva_ngc_team}/rmir_tts_fastpitch_hifigan_ljspeech:${riva_ngc_
138 ↪model_version}"
139 # "${riva_ngc_org}/${riva_ngc_team}/rmir_tts_tacotron_waveglow_ljspeech:${riva_ngc_
140 ↪model_version}"
141 )
142
143 NGC_TARGET=${riva_ngc_org}
144 if [[ ! -z ${riva_ngc_team} ]]; then
145 NGC_TARGET="${NGC_TARGET}/${riva_ngc_team}"
146 else
147 team="\\"
148 fi
149
150 # define docker images required to run Riva
151 image_client="nvcr.io/${NGC_TARGET}/riva-speech-client:${riva_ngc_image_version}"
152 image_speech_api="nvcr.io/${NGC_TARGET}/riva-speech:${riva_ngc_image_version}-server"
153
154 # define docker images required to setup Riva
155 image_init_speech="nvcr.io/${NGC_TARGET}/riva-speech:${riva_ngc_image_version}-
156 ↪servicemaker"
157
158 # daemon names
159 riva_daemon_speech="riva-speech"
160 riva_daemon_client="riva-client"

```



### 3.1.1 enable/ disable Riva services

Specify *line 10-12* to enable or disable Riva services.

```
service_enabled_asr=true  
service_enabled_nlp=true  
service_enabled_tts=true
```

### 3.1.2 encryption key

Specify *line 19* (consistency with the encryption key used to export models)

```
MODEL_DEPLOY_KEY="tlt_encode"
```

### 3.1.3 Riva model location

Specify *line 43* to configure the path of Riva model location

```
riva_model_loc="riva-model-repo"
```

---

**Note:** See [Riva - ServiceMaker](#) for details of how to use `riva-build` and `riva-deploy` tools to deploy custom models in Riva model repository.

---

### 3.1.4 ports config

Specify *line 53-54* to configure the exposed ports for Riva services

```
riva_speech_api_port="50051"  
riva_vision_api_port="60051"
```



## **RIVA - SPEECH RECOGNITION**

Automatic Speech Recognition (ASR) is the first step of building a conversational AI.

- Offline Recognition
- Streaming Recognition



## **RIVA - SPEECH RECOGNITION - MODEL ARCHITECTURES**

### **5.1 Jasper**

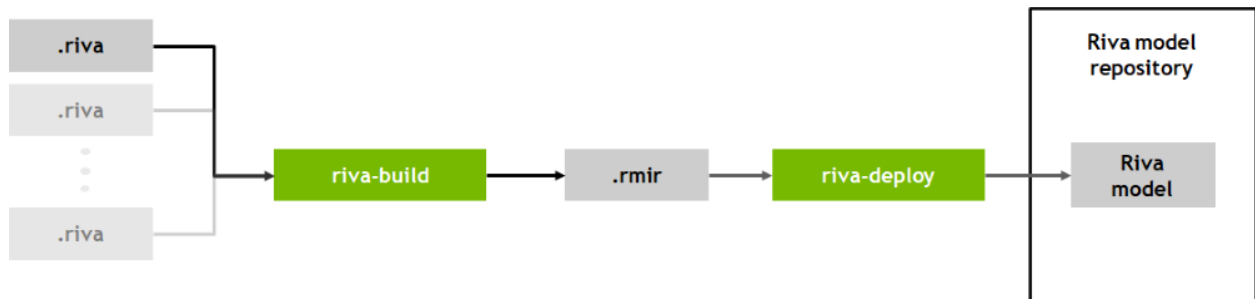
### **5.2 QuartzNet**

### **5.3 Citrinet**



## RIVA - SPEECH RECOGNITION - DEPLOY ASR MODEL

Riva ServiceMaker will be used to deploy ASR .riva models in Riva server.



---

**Note:** Refer to *Riva - ServiceMaker* to see basic usage of riva-build and riva-deploy.

---

### 6.1 riva-build

Simple riva-build pipeline for speech recognition (ASR) models. See [ASR model pipeline configuration](#) for more details.

```
riva-build speech_recognition \
  /servicemaker-dev/<rmir_filename>:<encryption_key> \
  /servicemaker-dev/<riva_filename>:<encryption_key> \
  --name=<pipeline_name> \
  --decoder_type=greedy \
  --acoustic_model_name=<acoustic_model_name>
```

- `speech_recognition`: task of riva-build. See *Riva - ServiceMaker* for more details.

Sample riva-build tool usage for ASR models.

```
docker run --rm --gpus 1 \
  -v $MODEL_LOC:/tao \
  -v $RMIR_LOC:/riva \
  $RIVA_SM_CONTAINER -- \
  riva-build speech_recognition --decoder_type=greedy /riva/asr.rmir:$KEY /tao/
↪ $MODEL_NAME:$KEY
```

where

- MODEL\_LOC: path of directory exported .riva model
- RMIR\_LOC: path of output .rmir file
- RIVA\_SM\_CONTAINER: name of Riva ServiceMaker image.
- MODEL\_NAME: name of exported .riva model
- KEY: encryption\_key used in exporting .riva model.

## 6.2 riva-deploy

Sample riva-deploy tool usage for ASR models.

```
docker run --rm --gpus 1 \  
    -v $RIVA_MODEL_LOC:/data \  
    $RIVA_SM_CONTAINER -- \  
    riva-deploy -f /data/rmir/asr.rmir:$KEY /data/models/
```

where

- RIVA\_MODEL\_LOC: path of Riva model repository.

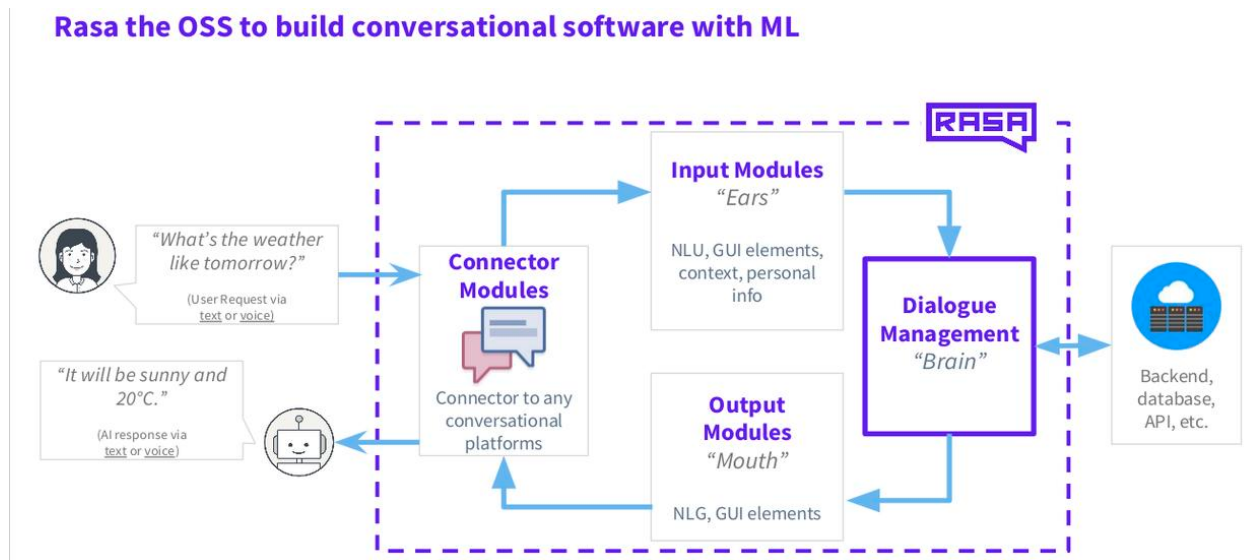


## **RIVA - NATURAL LANGUAGE PROCESSING**

- Text Classification
- Token Classification (Named Entity Recognition)
- Joint Intent and Slots
- Question Answering (Extractive)
- Punctuation and Capitalization



## RIVA - RASA



### 8.1 Prerequisites

Create and activate a virtual environment for Rasa

```
python3 -m venv rasa_envs  
source rasa_envs/bin/activate
```

Download rasa open source

```
pip install rasa
```

## 8.2 Sample Chatbot

Activate the rasa environment, and run

```
rasa init
```

Follow the instructions to create sample chatbot directory, train models, start Rasa server, and start chatting with the chatbot.

Simply type in `/stop` or `Ctrl+C` to exit the Rasa server.

---

**Note:**

- To start Rasa server in shell, run `rasa shell`
  - To interact with Rasa Chatbot for improving its correctness on responses or predictions, run `rasa interactive`. Hit `Ctrl+C` to end the interactive session, and you might export the new data for `nlu.yml`, `domains.yml`, and `stories.yml`
- 

### 8.2.1 Chatbot structure

```
.
├── actions
│   ├── __init__.py
│   ├── __pycache__
│   │   ├── __init__.cpython-38.pyc
│   │   └── actions.cpython-38.pyc
│   └── actions.py
├── config.yml
├── credentials.yml
├── data
│   ├── nlu.yml
│   ├── rules.yml
│   └── stories.yml
├── domain.yml
├── endpoints.yml
├── models
│   └── 20210909-224920.tar.gz
├── tests
│   └── test_stories.yml
```

**See also:**

Useful learning sources:

- [Rasa for beginners](#)
- [Rasa Advanced Workshop - Deployment](#)
- [Rasa Advanced Workshop - Custom Actions, Forms, & Responses](#)
- [Rasa Certification Workshop](#)

## RIVA - RASA - NLU COMPONENT

## 9.1 Sample nlu.yml

```
version: "2.0"

nlu:
- intent: greet
  examples: |
    - hey
    - hello
    - hi
    - hello there
    - good morning
    - good evening
    - moin
    - hey there
    - let's go
    - hey dude
    - goodmorning
    - goodevening
    - good afternoon

- intent: goodbye
  examples: |
    - good afternoon
    - cu
    - good by
    - cee you later
    - good night
    - bye
    - goodbye
    - have a nice day
    - see you around
    - bye bye
    - see you later

- intent: affirm
  examples: |
    - yes
    - y
```

(continues on next page)

(continued from previous page)

```
- indeed
- of course
- that sounds good
- correct

- intent: deny
examples: |
  - no
  - n
  - never
  - I don't think so
  - don't like that
  - no way
  - not really

- intent: mood_great
examples: |
  - perfect
  - great
  - amazing
  - feeling like a king
  - wonderful
  - I am feeling very good
  - I am great
  - I am amazing
  - I am going to save the world
  - super stoked
  - extremely good
  - so so perfect
  - so good
  - so perfect

- intent: mood_unhappy
examples: |
  - my day was horrible
  - I am sad
  - I don't feel very well
  - I am disappointed
  - super sad
  - I'm so sad
  - sad
  - very sad
  - unhappy
  - not good
  - not very good
  - extremely sad
  - so sad
  - so sad

- intent: bot_challenge
examples: |
  - are you a bot?
```

(continues on next page)

(continued from previous page)

- are you a human?
  - am I talking to a bot?
  - am I talking to a human?

---

**Note:**

- after adding/ removing intents, update `intents` in `domain.yml` (See [Riva - Rasa - Domain](#)).
- 

---

**Note:**

- after adding/ removing intents, you may want to add/ remove/ update corresponding responses for such intents. To do that, update `responses` in `domain.yml` (See [Riva - Rasa - Domain](#)).
-





## RIVA - RASA - STORIES (DIALOG MANAGER)

Stories (Dialog Manager) are defined to train Dialog Manager model in order to inform Rasa chatbots what to do next (or next next, next next next, ...)

### 10.1 Sample stories.yml

```
version: '2.0'

stories:
  - story: happy path
    steps:
      - intent: greet
      - action: utter_greet
      - intent: mood_great
      - action: utter_happy

  - story: sad path 1
    steps:
      - intent: greet
      - action: utter_greet
      - intent: mood_unhappy
      - action: utter_cheer_up
      - action: utter_did_that_help
      - intent: affirm
      - action: utter_happy

  - story: sad path 2
    steps:
      - intent: greet
      - action: utter_greet
      - intent: mood_unhappy
      - action: utter_cheer_up
      - action: utter_did_that_help
      - intent: deny
      - action: utter_goodbye
```



**RIVA - RASA - ACTIONS**



## RIVA - RASA - DOMAIN

### 12.1 Sample domain.yml

```
version: "2.0"

intents:
- greet
- goodbye
- affirm
- deny
- mood_great
- mood_unhappy
- bot_challenge

responses:
utter_greet:
- text: "Hey! How are you?"

utter_cheer_up:
- text: "Here is something to cheer you up:"
  image: "https://i.imgur.com/nGF1K8f.jpg"

utter_did_that_help:
- text: "Did that help you?"

utter_happy:
- text: "Great, carry on!"

utter_goodbye:
- text: "Bye"

utter_iamabot:
- text: "I am a bot, powered by Rasa."

session_config:
session_expiration_time: 60
carry_over_slots_to_new_session: true
```



## **RIVA - RASA - CONFIGURATIONS**





## RIVA - RASA - MODEL TRAINING

After updating some components of Rasa, train/ retrain the models before running Rasa.

```
rasa train
```



## **RIVA - RASA - CREDENTIALS**



## **RIVA - RASA - ENDPOINTS**



**RIVA - RASA - CONVERSATION TEST**





## **RIVA - SPEECH SYNTHESIS**

- Synthesize
- SynthesizeOnline



## RIVA - SPEECH SYNTHESIS - MODEL ARCHITECTURES

Table 1: Available models for Speech Synthesis

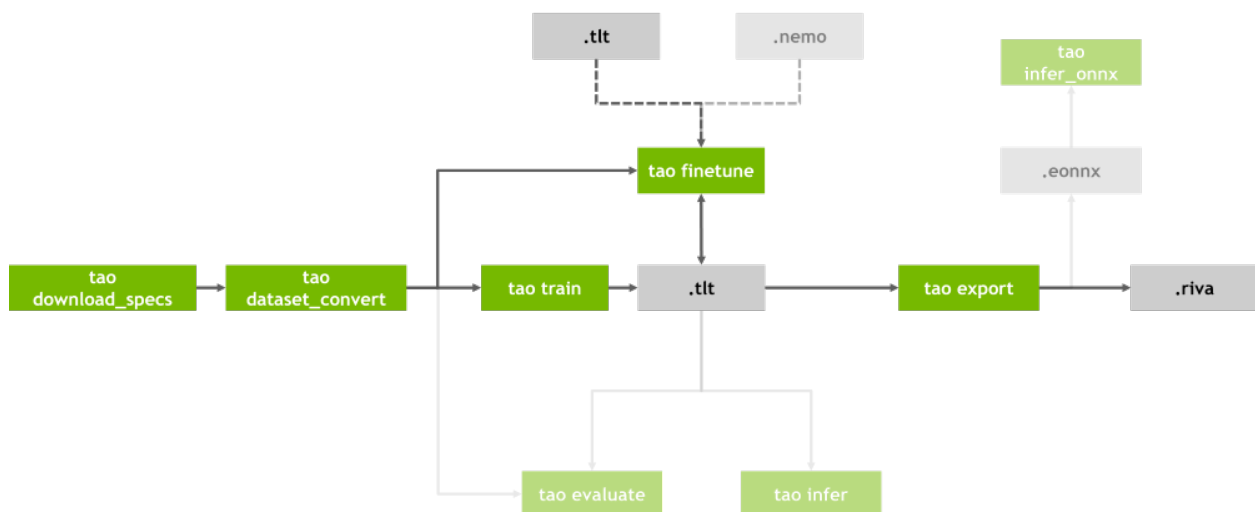
Tasks	Models
Mel Spectrogram Generator	Tacotron 2
Mel Spectrogram Generator	Glow-TTS
Mel Spectrogram Generator	TalkNet
Mel Spectrogram Generator	FastPitch
Mel Spectrogram Generator	FastSpeech2
Audio Generator	WaveGlow
Audio Generator	SqueezeWave
Audio Generator	UniGlow
Audio Generator	MelGAN
Audio Generator	HiFiGAN
Audio Generator	Griffin-Lim
End-to-end model	FastPitch_HifiGan_E2E
End-to-end model	FastSpeech2_HifiGan_E2E



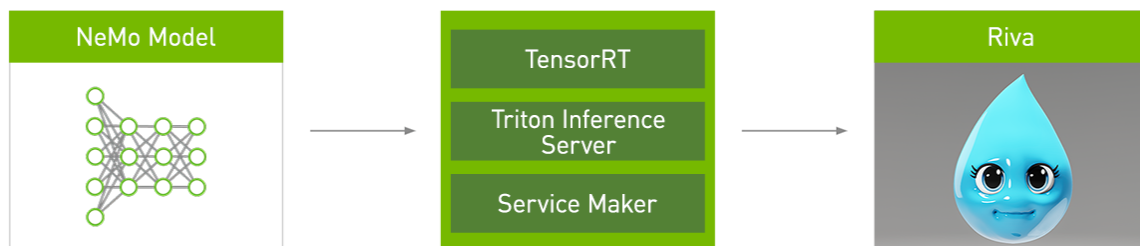
## RIVA - CUSTOM MODEL

To train, export AI models, and then deploy the models in Riva server, we will use NVIDIA TAO and NVIDIA NeMo. Both training frameworks enable the transformation of custom models into `.riva` models that can be deployed in Riva.

### 20.1 TAO



### 20.2 NeMo



## 20.3 Export models

The exported models in `.riva` format can be deployed in Riva with Riva ServiceMaker.

---

**Note:** See *Riva - ServiceMaker* to deploy TAO- or NeMo-trained models.

---

**See also:**

- [NeMo guide](#)
- [NeMo Tutorials](#)

## RIVA - CUSTOM MODEL - TAO

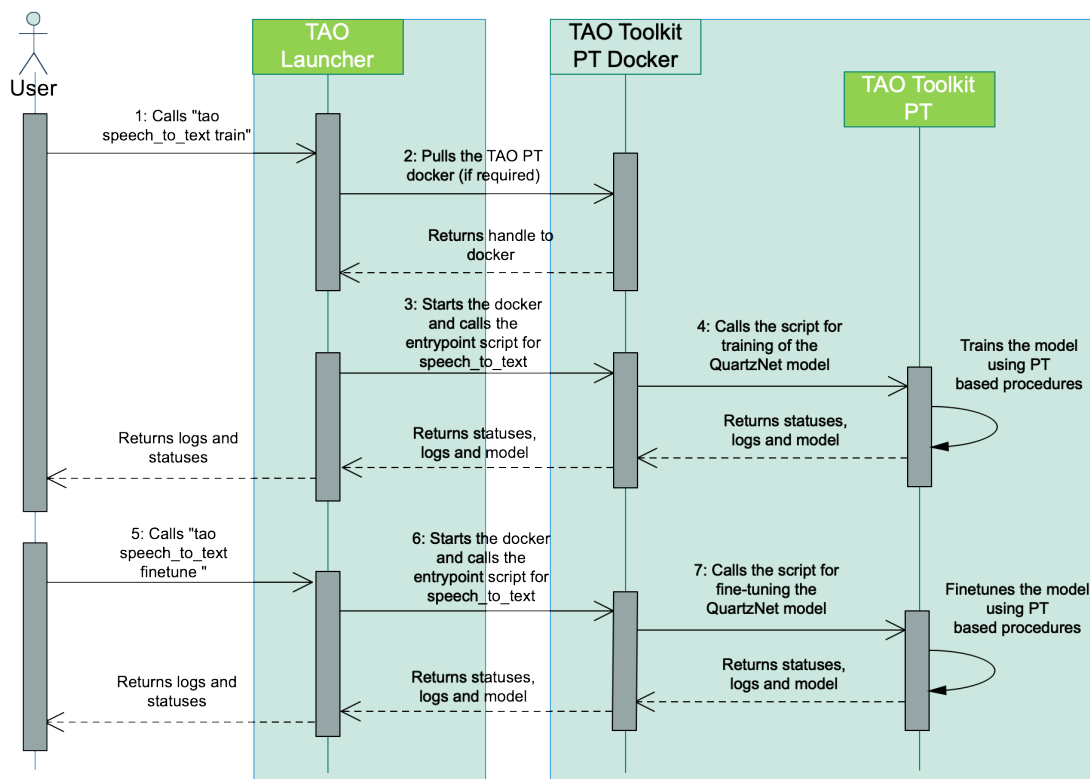
### 21.1 TAO Launcher

```
tao <task> <subtask> <args>
```





## RIVA - CUSTOM MODEL - TAO - SPEECH-TO-TEXT



### 22.1 Usage

The structure used to perform speech\_to\_text task

```
tao speech_to_text <subtask>
```

### 22.1.1 Subtasks

The subtasks for `speech_to_text` include *dataset\_convert*, *evaluate*, *export*, *finetune*, *infer*, *infer\_onnx*, *train*, *download\_specs*

Sample task *download\_specs*

```
!tao speech_to_text download_specs \
  -o $SPECS_DIR/speech_to_text \
  -r $RESULTS_DIR
```

Sample task *export*

```
# SPECS_DIR/speech_to_text/export.yaml
# Name of the .tlt EFF archive to be loaded/model to be exported.
restore_from: trained-model.tlt

# Set export format: ONNX | RIVA
export_format: RIVA

# Output EFF archive containing ONNX.
export_to: exported-model.riva
```

```
!tao speech_to_text export \
  -r $RESULTS_DIR/quartznet/export \
  -k $KEY \
  -e $SPECS_DIR/speech_to_text/export.yaml \
  -g 1 \
  -m $MODELS_DIR/spechtotext_english_quartznet.tlt \
  export_format=RIVA \
  export_to=asr-model.riva
```

### 22.1.2 Arguments

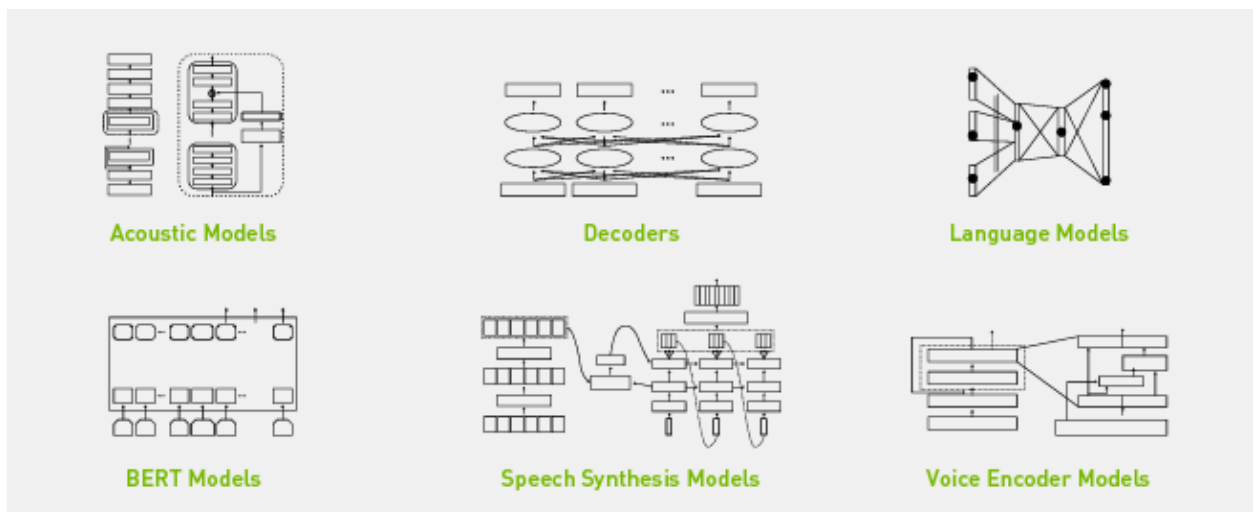
Table 1: TAO speech-to-text task arguments

Arguments	Details
subtask	choose the task from <i>Subtasks</i>
-r	path of results directory
-k	your NGC API key
-e	spec file of a specific subtask
-g	number of gpu to use
-m	path of pretrained ASR model
-o	specify path for downloading spec files
additional arguments	override values in spec file

## RIVA - CUSTOM MODEL - NEMO



NVIDIA NeMo (Neural Module) is built on PyTorch and PyTorch lightning.



## 23.1 Prerequisites

- NGC API Key

## 23.2 Getting started

1. Clone NVIDIA NeMo Github folder

```
git clone https://github.com/NVIDIA/NeMo
```

2. Create PyTorch container

```
docker run --gpus all -it --rm -v <nemo_github_folder>:/NeMo --shm-size=8g \  
-p 8888:8888 -p 6006:6006 --ulimit memlock=-1 --ulimit \  
stack=67108864 --device=/dev/snd nvcr.io/nvidia/pytorch:21.05-py3
```

3. Install NeMo

```
cd /NeMo  
./reinstall.sh
```

### See also:

- [NeMo Tutorials](#)
- [ASR with NeMo](#)

## RIVA - CUSTOM MODEL - NEMO - DATA MANIFEST - ASR

The standardized manifest (a json file storing all metadata of a given dataset) format for NeMo is that

- each line corresponds to one sample of audio
- Number of lines equals to number of samples

A sample of a line in the manifest

```
{"audio_filepath": "path/to/audio.wav", "duration": 3.45, "text": "transcript of audio.  
↪wav"}
```



## RIVA - CUSTOM MODEL - NEMO - ASR - MANDARIN

This sample implements ASR of Mandarin using Nemo.

---

**Note:**

- If you wish to skip the following steps to set up the environment to train mandarin ASR model, a container is prepared which you can launch and start training immediately:

```
docker pull nvcr.io/nvidian/sae/mandarin-asr-spoon:1.0.1
```

- To launch the container, run:

```
docker run --gpus all -it --rm --shm-size=8g \
-p 8888:8888 -p 6006:6006 --ulimit memlock=-1 --ulimit \
stack=67108864 --device=/dev/snd nvcr.io/nvidian/sae/mandarin-asr-spoon:1.0.1
```

- Finally, run the training:

```
python3 main.py --e <train-epoch> --g <no-of-gpus>
```

---

### 25.1 0. Launch PyTorch container with NeMo

Follow *Getting started* to set up container for the sample ASR Mandarin implementation.

Then start a jupyter notebook inside the container

```
jupyter notebook --ip=0.0.0.0 --allow-root
```

Access the jupyter notebook by the URI *http://<your-host-ip>:8888/tree*. Create a notebook, and so you are ready for next step.

## 25.2 1. Prepare data manifest

### 25.2.1 1.1 Download Aishell-1 dataset

Download AISHELL-1 dataset from Openslr. AISHELL-1 dataset is an open-sourced dataset containing 178 hours of Mandarin utterances.

```
!wget https://www.openslr.org/resources/33/data_aishell.tgz
!tar zxvf data_aishell.tgz
```

Retrieve the audio files

```
import os
import tarfile
import shutil
from tqdm import tqdm

# remove manifests if already exists
os.chdir("/NeMo/data_aishell/wav")
if delete_previous:
    shutil.rmtree("train")
    shutil.rmtree("test")
    shutil.rmtree("dev")

# untar the audio files if you have not done this
wav_tars = os.listdir(os.getcwd())
for wav_tar in tqdm(wav_tars):
    current_tar = os.path.join(os.getcwd(), wav_tar)
    tar = tarfile.open(current_tar)
    tar.extractall('.')
    tar.close()
```

### 25.2.2 1.2 Build data manifest

Data manifest are files containing all metadata of audio required by NeMo.

```
import librosa
import json

def build_aishell_manifest(transcript_path, manifest_path, wav_path):
    with open(transcript_path, "r") as fin:
        with open(manifest_path, "w") as fout:

            transcript = fin.readlines()

            for ref in tqdm(transcript):
                s_num = ref[6:11]
                text = ref[11 + 5:].strip()
                audio_file = "BAC009" + ref[6:16] + ".wav"
                audio_path = os.path.join(wav_path, s_num, audio_file)
```

(continues on next page)



(continued from previous page)

```

if os.path.exists(audio_path):

    duration = librosa.core.get_duration(filename=audio_path)

    metadata = {
        "audio_filepath": audio_path,
        "duration": duration,
        "text": text
    }

    json.dump(metadata, fout, ensure_ascii=False)
    fout.write('\n')

# build training manifest
print("building training manifest")
transcript_path = "/NeMo/data_aishell/transcript/aishell_transcript_v0.8.txt"
train_manifest_path = "/NeMo/data_aishell/train_manifest.json"
train_wav_path = "/NeMo/data_aishell/wav/train"
build_aishell_manifest(transcript_path, train_manifest_path, train_wav_path)
print("training manifest prepared at {}".format(train_manifest_path))

# build validation manifest
print("\nbuilding validation manifest")
transcript_path = "/NeMo/data_aishell/transcript/aishell_transcript_v0.8.txt"
valid_manifest_path = "/NeMo/data_aishell/valid_manifest.json"
valid_wav_path = "/NeMo/data_aishell/wav/test"
build_aishell_manifest(transcript_path, valid_manifest_path, valid_wav_path)
print("training manifest prepared at {}".format(valid_manifest_path))

# build testing manifest
print("\nbuilding testing manifest")
transcript_path = "/NeMo/data_aishell/transcript/aishell_transcript_v0.8.txt"
test_manifest_path = "/NeMo/data_aishell/test_manifest.json"
test_wav_path = "/NeMo/data_aishell/wav/dev"
build_aishell_manifest(transcript_path, test_manifest_path, test_wav_path)
print("testing manifest prepared at {}".format(test_manifest_path))

```

**See also:**

see *Riva - Custom Model - NeMo - Data Manifest - ASR* for more about manifest formatting.

## 25.3 2. Prepare training configuration

### 25.3.1 2.1 Load QuartzNet 15x5 zh default config

The model we are using to train with Aishell-1 is QuartzNet 15x5 zh. We could start with a default config file `quartznet_15x5_zh.yaml` at `/NeMo/examples/asr/conf/quartznet/`.

```

from omegaconf import OmegaConf
import copy

```

(continues on next page)

(continued from previous page)

```
try:
    from ruamel.yaml import YAML
except ModuleNotFoundError:
    from ruamel_yaml import YAML

config_path = "/NeMo/examples/asr/conf/quartznet/quartznet_15x5_zh.yaml"

yaml = YAML(typ='safe')
with open(config_path) as f:
    zh_qn_params = yaml.load(f)

print(OmegaConf.to_yaml(zh_qn_params))
```

**Note:** The config file is very lengthy as it contains Chinese lexicons.

### 25.3.2 2.2 Update manifest path

Update the path to data manifests that we just created

```
zh_qn_params["model"]["train_ds"]["manifest_filepath"] = train_manifest_path
zh_qn_params["model"]["validation_ds"]["manifest_filepath"] = valid_manifest_path
```

### 25.3.3 2.3 Update num\_workers

Update num\_workers of data loaders

```
num_workers = 40
zh_qn_params["model"]["train_ds"]["num_workers"] = num_workers
zh_qn_params["model"]["validation_ds"]["num_workers"] = num_workers
```

### 25.3.4 2.4 Update logging directory

We will create a results folder `/NeMo/data_aishell/results` where our training logs and output `.nemo` model will be stored.

```
!mkdir data_aishell/results
```

Update the results path in the configuration.

```
results_path = "/NeMo/data_aishell/results"
zh_qn_params["exp_manager"]["exp_dir"] = results_path
```

## 25.4 3. Train QuartzNet 15x5 zh

After the configuration is ready, we move on to model training session. Define `pytorch_lightning.Trainer` and `nemo_asr.models.EncDecCTCModel` instances:

```
from omegaconf import DictConfig
import pytorch_lightning as pl
import nemo
import nemo.collections.asr as nemo_asr

trainer = pl.Trainer(checkpoint_callback=False, logger=False, gpus=4, max_epochs=50,
    ↪ accelerator="ddp")

zh_qn = nemo_asr.models.EncDecCTCModel(cfg=DictConfig(zh_qn_params['model']),
    ↪ trainer=trainer)
```

Also, we apply `exp_manager` to apply the logging path of training.

```
from nemo.utils.exp_manager import exp_manager
exp_manager(trainer, zh_qn_params.get("exp_manager", None))
```

However, `pl.Trainer` using an accelerator `ddp` cannot be done in interactive environment (notebook).

**Warning:** If you attempt to use `ddp` in the notebook, you will be likely warned:

```
AttributeError: Can't pickle local object 'FilterbankFeatures.__init__.<locals>'.
    ↪ <lambda>'
```

Thus, let's create a python script `main.py` to start the training and save our result `model.nemo`. To do so, go to jupyter notebook menu, open a terminal, create `main.py`, and paste the following script.

```
import librosa
import json
from omegaconf import OmegaConf
import copy
from omegaconf import DictConfig
import pytorch_lightning as pl
import nemo
import nemo.collections.asr as nemo_asr
from nemo.utils.exp_manager import exp_manager
try:
    from ruamel.yaml import YAML
except ModuleNotFoundError:
    from ruamel_yaml import YAML

def main():
    # set manifest path
    train_manifest_path = "/NeMo/data_aishell/train_manifest.json"
    valid_manifest_path = "/NeMo/data_aishell/valid_manifest.json"

    # load QuartzNet 15x5 zh default config
    config_path = "/NeMo/examples/asr/conf/quartznet/quartznet_15x5_zh.yaml"
```

(continues on next page)

(continued from previous page)

```

yaml = YAML(typ='safe')
with open(config_path) as f:
    zh_qn_params = yaml.load(f)

# update manifest path
zh_qn_params["model"]["train_ds"]["manifest_filepath"] = train_manifest_path
zh_qn_params["model"]["validation_ds"]["manifest_filepath"] = valid_manifest_path

# update num_workers
num_workers = 40
zh_qn_params["model"]["train_ds"]["num_workers"] = num_workers
zh_qn_params["model"]["validation_ds"]["num_workers"] = num_workers

# update training logs path
results_path = "/NeMo/data_aishell/results"
zh_qn_params["exp_manager"]["exp_dir"] = results_path

# Initialise trainer
trainer = pl.Trainer(checkpoint_callback=False, logger=False, gpus=4, max_epochs=50,
↪ accelerator="ddp")
zh_qn = nemo_asr.models.EncDecCTCModel(cfg=DictConfig(zh_qn_params['model']),
↪ trainer=trainer)

# training logs export
exp_manager(trainer, zh_qn_params.get("exp_manager", None))

# Start training
zh_qn.summarize()
trainer.fit(zh_qn)

# Save trained model
zh_qn.save_to("/NeMo/data_aishell/results/model.nemo")

if __name__ == '__main__':
    main()

```

Then, in the terminal, execute `main.py` to start the training process

```
python3 main.py
```

## 25.5 4. Visualise Training Progress

The training logs are saved at `/NeMo/data_aishell/results/QuartzNet15x5`, and we use TensorBoard for visualisation using a custom Python environment with TensorBoard installed.

Create a python environment

```
python3 -m venv venv
```

Install tensorboard

```
pip install tensorboard
```

Launch tensorboard

```
cd /NeMo/data_aishell/results/  
tensorboard --bind_all --logdir QuartzNet15x5
```

**Warning:** If you use tensorboard that comes with NeMo package install script `reinstall.sh` (as of Sep 2021), that is:

```
tensorboard --bind_all --logdir QuartzNet15x5/
```

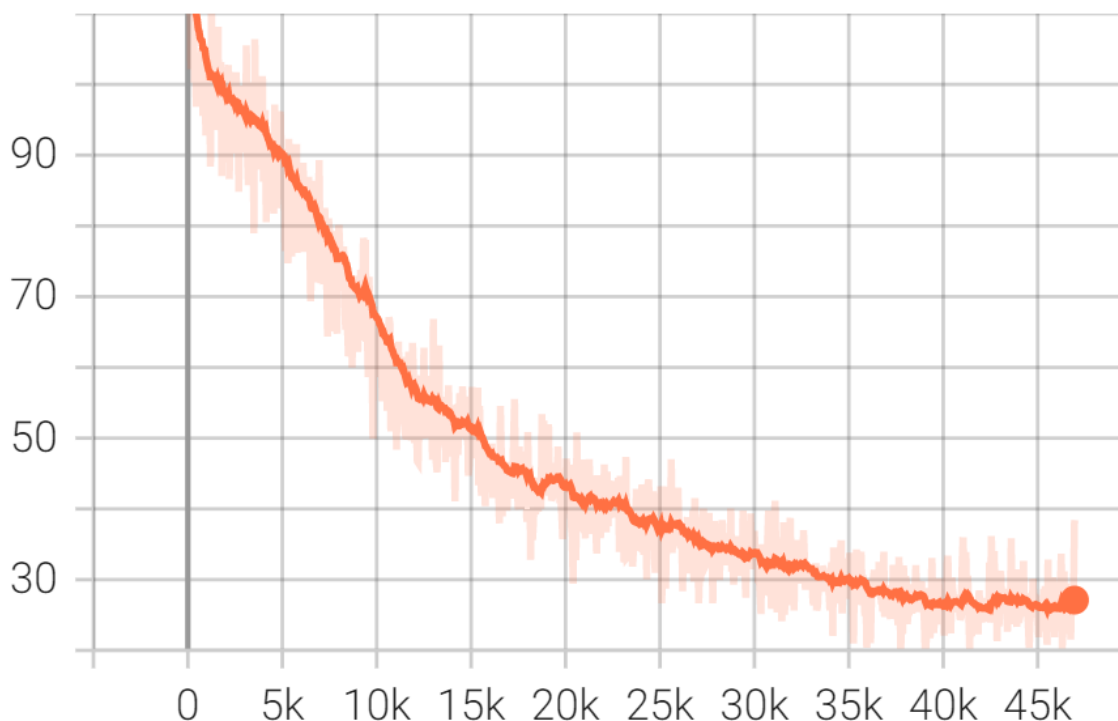
you will meet the following error:

```
ValueError: Duplicate plugins for name projector
```

### 25.5.1 Training loss

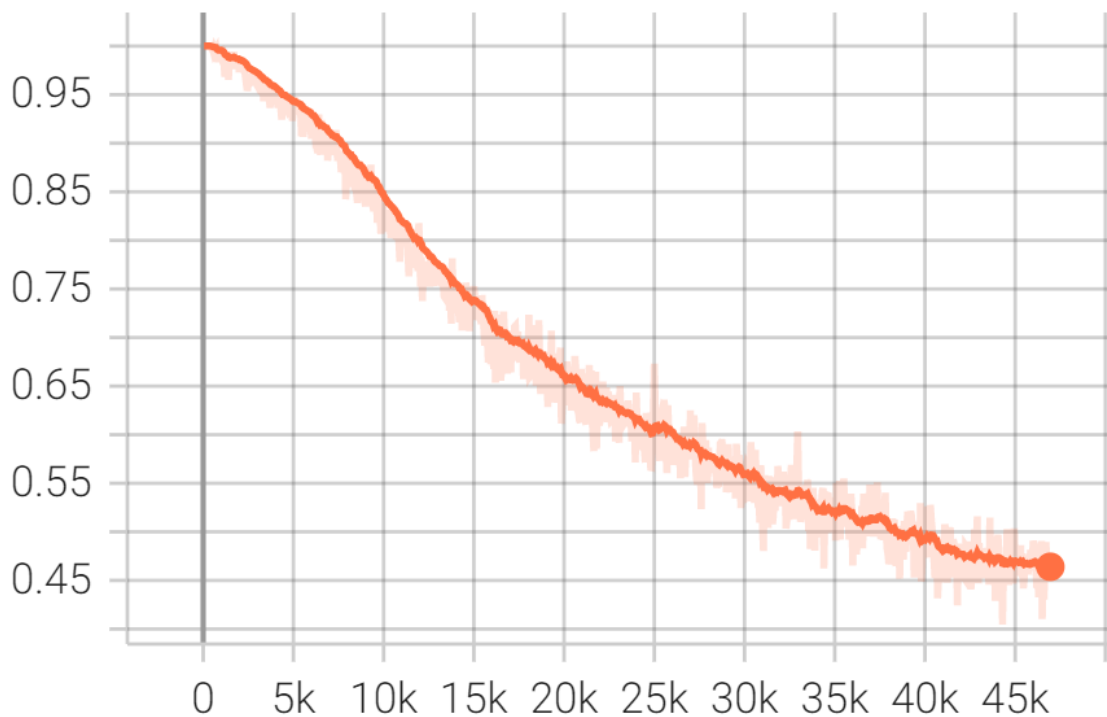
train\_loss

tag: train\_loss



### 25.5.2 Training batch WER

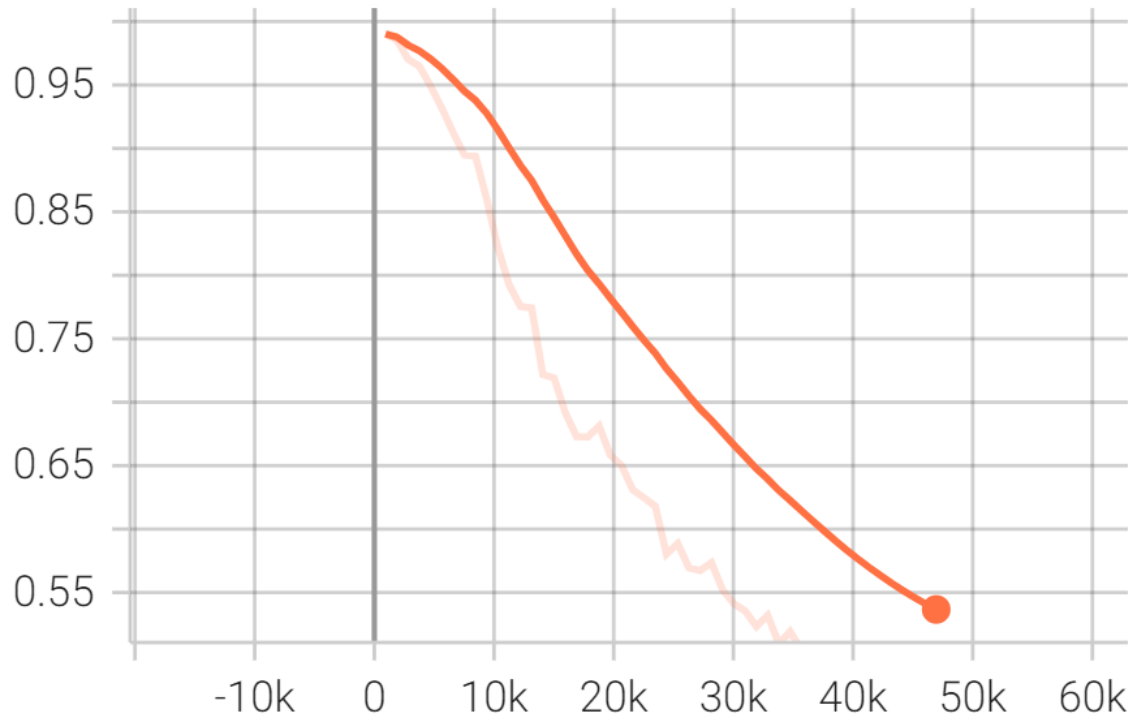
training\_batch\_wer  
tag: training\_batch\_wer



### 25.5.3 Validation WER

val\_wer

tag: val\_wer







## RIVA - CUSTOM MODEL - NEMO - DATA FORMAT - NLP (INTENT-SLOT)

It's required to convert data into specific NeMo format when training NLP models for *Joint Intent and Slot Classification*. The model training requires necessary input files for *intent* training and *slot* training (`dict.intents.csv`, `dict.slots.csv`, `train.tsv`, `train_slots.tsv`, `test.tsv`, `test_slots.tsv`) which are expected to store in a directory. The path of data directory needs to be updated in the training config file.

### 26.1 Expected structure

```
|--<target_data_dir>
    |-- dict.intents.csv
    |-- dict.slots.csv
    |-- train.tsv
    |-- train_slots.tsv
    |-- test.tsv
    |-- test_slots.tsv
```

### 26.2 Sample dict.intents.csv

Suppose we have N intents

```
#intent #label
interest      # 0
transportation # 1
sleep         # 2
...
intentN       # N-1
```

## 26.3 Sample dict.slots.csv

Suppose we have  $N$  slots

```
#slot    #label
interest_organisation    # 0
transportation_method    # 1
sleep_bahaviour          # 2
...
0                          # N-1
```

## 26.4 Sample train.tsv

For every first line of any `train.tsv`, we have a header line

```
sentence <tab> label
```

and, for rest of the lines follows `<sentence> <tab> <label>` where `<label>` refers to label in `dict.intents.csv`

```
I love NVIDIA <tab> <label>
I go to company by bus <tab> <label>
I sleep late everyday <tab> <label>
...
```

## 26.5 Sample train\_slots.tsv

There is no header line for any `train_slots.tsv`. For each line of the file, each token of the line is presented a label to slots (defined in `dict.slots.csv`).

```
# there is space between labels
N-1 N-1 0
N-1 N-1 N-1 N-1 N-1 1
N-1 N-1 2 N-1
...
```

## 26.6 Sample test.tsv

Similar with `train.tsv`

## 26.7 Sample test\_slots.tsv

Similar with `train_slots.tsv`

**See also:**

- [Joint\\_Intent\\_and\\_Slot\\_Classification.ipynb](#)



## **RIVA - CUSTOM MODEL - NEMO - DATA FORMAT - NLP (NER)**

**See also:**

- [Token\\_Classification\\_Named\\_Entity\\_Recognition.ipynb](#)



## RIVA - CUSTOM MODEL - NEMO - DATA FORMAT - NLP (TEXT CLASSIFICATION)

See also:

- [Text\\_Classification\\_Sentiment\\_Analysis.ipynb](#)





## RIVA - CUSTOM MODEL - NEMO - TTS - INFERENCE

### 29.1 Two-stage Approach

#### 29.1.1 Download and load models

```
import soundfile as sf
from nemo.collections.tts.models.base import SpectrogramGenerator, Vocoder

## spectrogram generators
spec_generator = SpectrogramGenerator.from_pretrained(model_name="tts_en_fastpitch").
↳ cuda()
# spec_generator = SpectrogramGenerator.from_pretrained(model_name="tts_en_tacotron2").
↳ cuda()
# spec_generator = SpectrogramGenerator.from_pretrained(model_name="tts_en_glowtts").
↳ cuda()
# spec_generator = SpectrogramGenerator.from_pretrained(model_name="tts_en_tts_en_talknet
↳ ").cuda()
# spec_generator = SpectrogramGenerator.from_pretrained(model_name="tts_en_fastspeech2").
↳ cuda()

## Audio generators
vocoder = Vocoder.from_pretrained(model_name="tts_hifigan").cuda()
# vocoder = Vocoder.from_pretrained(model_name="tts_waveglow").cuda()
# vocoder = Vocoder.from_pretrained(model_name="tts_squeezewave").cuda()
# vocoder = Vocoder.from_pretrained(model_name="tts_uniglow").cuda()
# vocoder = Vocoder.from_pretrained(model_name="tts_melgan").cuda()

# TBA for griffin-lim
```

---

**Note:**

- See *Riva - Speech Synthesis - Model Architectures* for details of models
-

### 29.1.2 Convert to audio

```
parsed = spec_generator.parse("You can type your sentence here to get nemo to produce_↵
↵speech.")
spectrogram = spec_generator.generate_spectrogram(tokens=parsed)
audio = vocoder.convert_spectrogram_to_audio(spec=spectrogram)
```

---

**Note:** Mel Spectrogram generators have two helper functions:

- `parse()`: Accepts raw python strings and returns a `torch.tensor` that represents tokenized text
- `generate_spectrogram()`: Accepts a batch of tokenized text and returns a `torch.tensor` that represents a batch of spectrograms

Vocoder have just one helper function:

- `convert_spectrogram_to_audio()`: Accepts a batch of spectrograms and returns a `torch.tensor` that represents a batch of raw audio
- 

### 29.1.3 Save Audio

```
sf.write("speech.wav", audio.to('cpu').detach().numpy()[0], 22050)
```

## 29.2 End-to-end Approach

### 29.2.1 Download and load models

```
import soundfile as sf
from nemo.collections.tts.models.base import SpectrogramGenerator, Vocoder, ↵
↵TextToWaveform

e2e_model = TextToWaveform.from_pretrained("tts_en_e2e_fastpitchhifigan").cuda()
# e2e_model = TextToWaveform.from_pretrained("tts_en_e2e_fastspeech2hifigan").cuda()
```

### 29.2.2 Convert to audio

```
parsed = e2e_model.parse("You can type your sentence here to get nemo to produce speech.
↵")
audio = e2e_model.convert_text_to_waveform(tokens=parsed)[0]
```

---

**Note:** End-to-end models have two helper functions:

- `parse()`: Accepts raw python strings and returns a `torch.tensor` that represents tokenized text
  - `convert_text_to_waveform()`: Accepts a batch of tokenized text and returns a `torch.tensor` that represents a batch of raw audio
-

### 29.2.3 Save Audio

```
sf.write("speech.wav", audio.to('cpu').detach().numpy()[0], 22050)
```



## RIVA - CUSTOM MODEL - NEMO - DATA MANIFEST - TTS

See *Riva - Custom Model - NeMo - Data Manifest - ASR* to build data manifest for training TTS models.



## **RIVA - CUSTOM MODEL - NEMO - TTS - MANDARIN**

This sample implements the sample training Mel spectrogram generator model Tacotron2 using Mandarin dataset.

See [Original paper](#)





## RIVA - GRPC API

Riva's services are exposed using gRPC to maximise its compatibility with different software infrastructure and easier integrations.

### 32.1 riva\_asr.proto

```
1 // Copyright 2019 Google LLC.
2 // Copyright (c) 2021, NVIDIA CORPORATION. All rights reserved.
3 //
4 // Licensed under the Apache License, Version 2.0 (the "License");
5 // you may not use this file except in compliance with the License.
6 // You may obtain a copy of the License at
7 //
8 //     http://www.apache.org/licenses/LICENSE-2.0
9 //
10 // Unless required by applicable law or agreed to in writing, software
11 // distributed under the License is distributed on an "AS IS" BASIS,
12 // WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
13 // See the License for the specific language governing permissions and
14 // limitations under the License.
15 //
16
17 syntax = "proto3";
18
19 package nvidia.riva.asr;
20
21 option cc_enable_arenas = true;
22 option go_package = "nvidia.com/riva-speech";
23
24 import "riva_audio.proto";
25
26 /*
27 * The RivaSpeechRecognition service provides two mechanisms for converting speech to
28 * ↪ text.
29 */
30 service RivaSpeechRecognition {
31     // Recognize expects a RecognizeRequest and returns a RecognizeResponse. This
32     ↪ request will block
```

(continues on next page)

(continued from previous page)

```

31 // until the audio is uploaded, processed, and a transcript is returned.
32 rpc Recognize(RecognizeRequest) returns (RecognizeResponse) {}
33 // StreamingRecognize is a non-blocking API call that allows audio data to be fed to
↳ the server in
34 // chunks as it becomes available. Depending on the configuration in the
↳ StreamingRecognizeRequest,
35 // intermediate results can be sent back to the client. Recognition ends when the
↳ stream is closed
36 // by the client.
37 rpc StreamingRecognize(stream StreamingRecognizeRequest) returns (stream
↳ StreamingRecognizeResponse) {}
38 }
39
40
41 /*
42 * RecognizeRequest is used for batch processing of a single audio recording.
43 */
44 message RecognizeRequest {
45 // Provides information to recognizer that specifies how to process the request.
46 RecognitionConfig config = 1;
47 // The raw audio data to be processed. The audio bytes must be encoded as specified
↳ in
48 // `RecognitionConfig`.
49 bytes audio = 2;
50 }
51
52
53 /*
54 * A StreamingRecognizeRequest is used to configure and stream audio content to the
55 * Riva ASR Service. The first message sent must include only a
↳ StreamingRecognitionConfig.
56 * Subsequent messages sent in the stream must contain only raw bytes of the audio
57 * to be recognized.
58 */
59 message StreamingRecognizeRequest {
60 // The streaming request, which is either a streaming config or audio content.
61 oneof streaming_request {
62 // Provides information to the recognizer that specifies how to process the
63 // request. The first `StreamingRecognizeRequest` message must contain a
64 // `streaming_config` message.
65 StreamingRecognitionConfig streaming_config = 1;
66 // The audio data to be recognized. Sequential chunks of audio data are sent
67 // in sequential `StreamingRecognizeRequest` messages. The first
68 // `StreamingRecognizeRequest` message must not contain `audio` data
69 // and all subsequent `StreamingRecognizeRequest` messages must contain
70 // `audio` data. The audio bytes must be encoded as specified in
71 // `RecognitionConfig`.
72 bytes audio_content = 2;
73 }
74 }
75
76 // Provides information to the recognizer that specifies how to process the request

```

(continues on next page)

(continued from previous page)

```

77 message RecognitionConfig {
78     // The encoding of the audio data sent in the request.
79     //
80     // All encodings support only 1 channel (mono) audio.
81     AudioEncoding encoding = 1;
82
83     // Sample rate in Hertz of the audio data sent in all
84     // `RecognizeAudio` messages.
85     int32 sample_rate_hertz = 2;
86
87     // Required. The language of the supplied audio as a
88     // [BCP-47](https://www.rfc-editor.org/rfc/bcp/bcp47.txt) language tag.
89     // Example: "en-US".
90     // Currently only en-US is supported
91     string language_code = 3;
92
93     // Maximum number of recognition hypotheses to be returned.
94     // Specifically, the maximum number of `SpeechRecognizeAlternative` messages
95     // within each `SpeechRecognizeResult`.
96     // The server may return fewer than `max_alternatives`.
97     // If omitted, will return a maximum of one.
98     int32 max_alternatives = 4;
99
100     // The number of channels in the input audio data.
101     // ONLY set this for MULTI-CHANNEL recognition.
102     // Valid values for LINEAR16 and FLAC are `1`-`8`.
103     // Valid values for OGG_OPUS are `1`-`254`.
104     // Valid value for MULAW, AMR, AMR_WB and SPEEX_WITH_HEADER_BYTE is only `1`.
105     // If `0` or omitted, defaults to one channel (mono).
106     // Note: We only recognize the first channel by default.
107     // To perform independent recognition on each channel set
108     // `enable_separate_recognition_per_channel` to `true`.
109     int32 audio_channel_count = 7;
110
111     // If `true`, the top result includes a list of words and
112     // the start and end time offsets (timestamps) for those words. If
113     // `false`, no word-level time offset information is returned. The default is
114     // `false`.
115     bool enable_word_time_offsets = 8;
116
117     // If `true`, adds punctuation to recognition result hypotheses.
118     // The default `false` value does not add punctuation to result hypotheses.
119     bool enable_automatic_punctuation = 11;
120
121     // This needs to be set to `true` explicitly and `audio_channel_count` > 1
122     // to get each channel recognized separately. The recognition result will
123     // contain a `channel_tag` field to state which channel that result belongs
124     // to. If this is not true, we will only recognize the first channel. The
125     // request is billed cumulatively for all channels recognized:
126     // `audio_channel_count` multiplied by the length of the audio.
127     bool enable_separate_recognition_per_channel = 12;
128

```

(continues on next page)

(continued from previous page)

```

129         // Which model to select for the given request. Valid choices: Jasper, Quartznet
130         string model = 13;
131
132         // The verbatim_transcripts flag enables or disable inverse text normalization.
133         // 'true' returns exactly what was said, with no denormalization.
134         // 'false' applies inverse text normalization, also this is the default
135         bool verbatim_transcripts = 14;
136
137         // Custom fields for passing request-level
138         // configuration options to plugins used in the
139         // model pipeline.
140         map<string, string> custom_configuration = 24;
141
142     }
143
144     // Provides information to the recognizer that specifies how to process the request
145     message StreamingRecognitionConfig {
146         // Provides information to the recognizer that specifies how to process
147         // the request
148         RecognitionConfig config = 1;
149
150         // If 'true', interim results (tentative hypotheses) may be
151         // returned as they become available (these interim results are indicated with
152         // the 'is_final=false' flag).
153         // If 'false' or omitted, only 'is_final=true' result(s) are returned.
154         bool interim_results = 2;
155     }
156
157     // The only message returned to the client by the 'Recognize' method. It
158     // contains the result as zero or more sequential 'SpeechRecognitionResult'
159     // messages.
160     message RecognizeResponse {
161         // Sequential list of transcription results corresponding to
162         // sequential portions of audio. Currently only returns one transcript.
163         repeated SpeechRecognitionResult results = 1;
164     }
165
166     // A speech recognition result corresponding to the latest transcript
167     message SpeechRecognitionResult {
168
169         // May contain one or more recognition hypotheses (up to the
170         // maximum specified in 'max_alternatives').
171         // These alternatives are ordered in terms of accuracy, with the top (first)
172         // alternative being the most probable, as ranked by the recognizer.
173         repeated SpeechRecognitionAlternative alternatives = 1;
174
175         // For multi-channel audio, this is the channel number corresponding to the
176         // recognized result for the audio from that channel.
177         // For audio_channel_count = N, its output values can range from '1' to 'N'.
178         int32 channel_tag = 2;

```

(continues on next page)

(continued from previous page)

```

179
180 // Length of audio processed so far in seconds
181 float audio_processed = 3;
182 }
183
184 // Alternative hypotheses (a.k.a. n-best list).
185 message SpeechRecognitionAlternative {
186 // Transcript text representing the words that the user spoke.
187 string transcript = 1;
188
189 // The non-normalized confidence estimate. A higher number
190 // indicates an estimated greater likelihood that the recognized words are
191 // correct. This field is set only for a non-streaming
192 // result or, of a streaming result where `is_final=true`.
193 // This field is not guaranteed to be accurate and users should not rely on it
194 // to be always provided.
195 float confidence = 2;
196
197 // A list of word-specific information for each recognized word. Only populated
198 // if is_final=true
199 repeated WordInfo words = 3;
200 }
201
202 // Word-specific information for recognized words.
203 message WordInfo {
204 // Time offset relative to the beginning of the audio in ms
205 // and corresponding to the start of the spoken word.
206 // This field is only set if `enable_word_time_offsets=true` and only
207 // in the top hypothesis.
208 int32 start_time = 1;
209
210 // Time offset relative to the beginning of the audio in ms
211 // and corresponding to the end of the spoken word.
212 // This field is only set if `enable_word_time_offsets=true` and only
213 // in the top hypothesis.
214 int32 end_time = 2;
215
216 // The word corresponding to this set of information.
217 string word = 3;
218 }
219
220
221 // `StreamingRecognizeResponse` is the only message returned to the client by
222 // `StreamingRecognize`. A series of zero or more `StreamingRecognizeResponse`
223 // messages are streamed back to the client.
224 //
225 // Here are few examples of `StreamingRecognizeResponse`s
226 //
227 // 1. results { alternatives { transcript: "tube" } stability: 0.01 }
228 //
229 // 2. results { alternatives { transcript: "to be a" } stability: 0.01 }
230 //

```

(continues on next page)

(continued from previous page)

```
231 // 3. results { alternatives { transcript: "to be or not to be"
232 //                               confidence: 0.92 }
233 //                               alternatives { transcript: "to bee or not to bee" }
234 //                               is_final: true }
235 //
236
237 message StreamingRecognizeResponse {
238
239     // This repeated list contains the latest transcript(s) corresponding to
240     // audio currently being processed.
241     // Currently one result is returned, where each result can have multiple
242     // alternatives
243     repeated StreamingRecognitionResult results = 1;
244 }
245
246 // A streaming speech recognition result corresponding to a portion of the audio
247 // that is currently being processed.
248 message StreamingRecognitionResult {
249     // May contain one or more recognition hypotheses (up to the
250     // maximum specified in `max_alternatives`).
251     // These alternatives are ordered in terms of accuracy, with the top (first)
252     // alternative being the most probable, as ranked by the recognizer.
253     repeated SpeechRecognitionAlternative alternatives = 1;
254
255     // If `false`, this `StreamingRecognitionResult` represents an
256     // interim result that may change. If `true`, this is the final time the
257     // speech service will return this particular `StreamingRecognitionResult`,
258     // the recognizer will not return any further hypotheses for this portion of
259     // the transcript and corresponding audio.
260     bool is_final = 2;
261
262     // An estimate of the likelihood that the recognizer will not
263     // change its guess about this interim result. Values range from 0.0
264     // (completely unstable) to 1.0 (completely stable).
265     // This field is only provided for interim results (`is_final=false`).
266     // The default of 0.0 is a sentinel value indicating `stability` was not set.
267     float stability = 3;
268
269     // For multi-channel audio, this is the channel number corresponding to the
270     // recognized result for the audio from that channel.
271     // For audio_channel_count = N, its output values can range from '1' to 'N'.
272     int32 channel_tag = 5;
273
274     // Length of audio processed so far in seconds
275     float audio_processed = 6;
276 }
```

## 32.2 riva\_nlp.proto

```

1  // Copyright (c) 2021, NVIDIA CORPORATION. All rights reserved.
2  //
3  // NVIDIA CORPORATION and its licensors retain all intellectual property
4  // and proprietary rights in and to this software, related documentation
5  // and any modifications thereto. Any use, reproduction, disclosure or
6  // distribution of this software and related documentation without an express
7  // license agreement from NVIDIA CORPORATION is strictly prohibited.
8
9  syntax = "proto3";
10
11 package nvidia.riva.nlp;
12
13 option cc_enable_arenas = true;
14 option go_package = "nvidia.com/riva_speech";
15
16 /* Riva Natural Language Services implement generic and task-specific APIs.
17 * The generic APIs allows users to design
18 * models for arbitrary use cases that conform simply with input and output types
19 * specified in the service. As an explicit example, the ClassifyText function
20 * could be used for sentiment classification, domain recognition, language
21 * identification, etc.
22 * The task-specific APIs can be used for popular NLP tasks such as
23 * intent recognition (as well as slot filling), and entity extraction.
24 */
25
26 service RivaLanguageUnderstanding {
27
28     // ClassifyText takes as input an input/query string and parameters related
29     // to the requested model to use to evaluate the text. The service evaluates the
30     // text with the requested model, and returns one or more classifications.
31     rpc ClassifyText(TextClassRequest) returns (TextClassResponse) {}
32
33     // ClassifyTokens takes as input either a string or list of tokens and parameters
34     // related to which model to use. The service evaluates the text with the requested
35     // model, performing additional tokenization if necessary, and returns one or more
36     // class labels per token.
37     rpc ClassifyTokens(TokenClassRequest) returns (TokenClassResponse) {}
38
39     // TransformText takes an input/query string and parameters related to the
40     // requested model and returns another string. The behavior of the function
41     // is defined entirely by the underlying model and may be used for
42     // tasks like translation, adding punctuation, augment the input directly, etc.
43     rpc TransformText(TextTransformRequest) returns (TextTransformResponse) {}
44
45     // AnalyzeEntities accepts an input string and returns all named entities within
46     // the text, as well as a category and likelihood.
47     rpc AnalyzeEntities(AnalyzeEntitiesRequest) returns (TokenClassResponse) {}
48
49     // AnalyzeIntent accepts an input string and returns the most likely
50     // intent as well as slots relevant to that intent.

```

(continues on next page)

(continued from previous page)

```

51 //
52 // The model requires that a valid "domain" be passed in, and optionally
53 // supports including a previous intent classification result to provide
54 // context for the model.
55 rpc AnalyzeIntent(AnalyzeIntentRequest) returns (AnalyzeIntentResponse) {}
56
57 // PunctuateText takes text with no- or limited- punctuation and returns
58 // the same text with corrected punctuation and capitalization.
59 rpc PunctuateText(TextTransformRequest) returns (TextTransformResponse) {}
60
61 // NaturalQuery is a search function that enables querying one or more documents
62 // or contexts with a query that is written in natural language.
63 rpc NaturalQuery(NaturalQueryRequest) returns (NaturalQueryResponse) {}
64 }
65
66 // NLPModelParams is a metadata message that is included in every request message
67 // used by the Core NLP Service and is used to specify model characteristics/requirements
68 message NLPModelParams {
69     // Requested model to use. If unavailable, the request will return an error
70     string model_name = 1;
71 }
72
73 // TextTransformRequest is a request type intended for services like TransformText
74 // which take an arbitrary text input
75 message TextTransformRequest {
76     // Each repeated text element is handled independently for handling multiple
77     // input strings with a single request
78     repeated string text = 1;
79     uint32 top_n = 2; //
80     NLPModelParams model = 3;
81 }
82
83 // TextTransformResponse is returned by the TransformText method. Responses
84 // are returned in the same order as they were requested.
85 message TextTransformResponse {
86     repeated string text = 1;
87 }
88
89 // TextClassRequest is the input message to the ClassifyText service.
90 message TextClassRequest {
91     // Each repeated text element is handled independently for handling multiple
92     // input strings with a single request
93     repeated string text = 1;
94
95     // Return the top N classification results for each input. 0 or 1 will return top_
96     ↪class, otherwise N.
97     // Note: Current disabled.
98     uint32 top_n = 2;
99     NLPModelParams model = 3;
100 }
101 // Classification messages return a class name and corresponding score

```

(continues on next page)



(continued from previous page)

```

102 message Classification {
103     string class_name = 1;
104     float score = 2;
105 }
106
107 // Span of a particular result
108 message Span {
109     uint32 start = 1;
110     uint32 end = 2;
111 }
112
113 // ClassificationResults contain zero or more Classification messages
114 // If the number of Classifications is > 1, top_n > 1 must have been
115 // specified.
116 message ClassificationResult {
117     repeated Classification labels = 1;
118 }
119
120 // TextClassResponse is the return message from the ClassifyText service.
121 message TextClassResponse {
122     repeated ClassificationResult results = 1;
123 }
124
125 // TokenClassRequest is the input message to the ClassifyText service.
126 message TokenClassRequest {
127     // Each repeated text element is handled independently for handling multiple
128     // input strings with a single request
129     repeated string text = 1;
130
131     // Return the top N classification results for each input. 0 or 1 will return top_
132     ↪class, otherwise N.
133     // Note: Current disabled.
134     uint32 top_n = 3;
135     NLPModelParams model = 4;
136 }
137
138 // TokenClassValue is used to correlate an input token with its classification results
139 message TokenClassValue {
140     string token = 1;
141     repeated Classification label = 2;
142     repeated Span span = 3;
143 }
144
145 // TokenClassSequence is used for returning a sequence of TokenClassValue objects
146 // in the original order of input tokens
147 message TokenClassSequence {
148     repeated TokenClassValue results = 1;
149 }
150
151 // TokenClassResponse returns a single TokenClassSequence per input request
152 message TokenClassResponse {
153     repeated TokenClassSequence results = 1;

```

(continues on next page)

(continued from previous page)

```

153 }
154
155 // AnalyzeIntentContext is reserved for future use when we may send context back in a
156 // a variety of different formats (including raw neural network hidden states)
157 message AnalyzeIntentContext {
158     // Reserved for future use
159 }
160
161 // AnalyzeIntentOptions is an optional configuration message to be sent as part of
162 // an AnalyzeIntentRequest with query metadata
163 message AnalyzeIntentOptions {
164     // Optionally provide context from previous interactions to bias the model's
165     ↪ prediction
166     oneof context {
167         string previous_intent = 1;
168         AnalyzeIntentContext vectors = 2;
169     }
170     // Optional domain field. Domain must be supported otherwise an error will be
171     ↪ returned.
172     // If left blank, a domain detector will be run first and then the query routed to
173     ↪ the
174     // appropriate intent classifier (if it exists)
175     string domain = 3;
176
177     // Optional language field. Assumed to be "en-US" if not specified.
178     string lang = 4;
179 }
180
181 // AnalyzeIntentRequest is the input message for the AnalyzeIntent service
182 message AnalyzeIntentRequest {
183     // The string to analyze for intent and slots
184     string query = 1;
185     // Optional configuration for the request, including providing context from previous
186     ↪ turns
187     // and hardcoding a domain/language
188     AnalyzeIntentOptions options = 2;
189 }
190
191 // AnalyzeIntentResponse is returned by the AnalyzeIntent service, and includes
192 ↪ information
193 // related to the query's intent, (optionally) slot data, and its domain.
194 message AnalyzeIntentResponse {
195     // Intent classification result, including the label and score
196     Classification intent = 1;
197     // List of tokens explicitly marked as filling a slot relevant to the intent, where
198     ↪ the
199     // tokens may not exactly match the input (based on the recombined values after
200     ↪ tokenization)
201     repeated TokenClassValue slots = 2;
202     // Returns the inferred domain for the query if not hardcoded in the request. In the
203     ↪ case where
204     // the domain was hardcoded in AnalyzeIntentRequest, the returned domain is an exact
205     ↪ match to the

```

(continues on next page)

(continued from previous page)

```

197 // request. In the case where no domain matches the query, intent and slots will be_
↪unset.
198 //
199 // DEPRECATED, use Classification domain field.
200 string domain_str = 3;
201
202 // Returns the inferred domain for the query if not hardcoded in the request. In the_
↪case where
203 // the domain was hardcoded in AnalyzeIntentRequest, the returned domain is an exact_
↪match to the
204 // request. In the case where no domain matches the query, intent and slots will be_
↪unset.
205 Classification domain = 4;
206 }
207
208 // AnalyzeEntitiesOptions is an optional configuration message to be sent as part of
209 // an AnalyzeEntitiesRequest with query metadata
210 message AnalyzeEntitiesOptions {
211 // Optional language field. Assumed to be "en-US" if not specified.
212 string lang = 4;
213 }
214
215 // AnalyzeEntitiesRequest is the input message for the AnalyzeEntities service
216 message AnalyzeEntitiesRequest {
217 // The string to analyze for intent and slots
218 string query = 1;
219 // Optional configuration for the request, including providing context from previous_
↪turns
220 // and hardcoding a domain/language
221 AnalyzeEntitiesOptions options = 2;
222 }
223
224 message NaturalQueryRequest {
225 // The natural language query
226 string query = 1;
227
228 // Maximum number of answers to return for the query. Defaults to 1 if not set.
229 uint32 top_n = 2;
230
231 // Context to search with the above query
232 string context = 3;
233 }
234
235 message NaturalQueryResult {
236 // text which answers the query
237 string answer = 1;
238 // Score representing confidence in result
239 float score = 2;
240 }
241
242 message NaturalQueryResponse {
243 repeated NaturalQueryResult results = 1;

```

(continues on next page)

```

244 }

```

## 32.3 riva\_tts.proto

```

1  // Copyright (c) 2021, NVIDIA CORPORATION. All rights reserved.
2  //
3  // NVIDIA CORPORATION and its licensors retain all intellectual property
4  // and proprietary rights in and to this software, related documentation
5  // and any modifications thereto. Any use, reproduction, disclosure or
6  // distribution of this software and related documentation without an express
7  // license agreement from NVIDIA CORPORATION is strictly prohibited.
8
9
10 syntax = "proto3";
11
12 package nvidia.riva.tts;
13
14 option cc_enable_arenas = true;
15 option go_package = "nvidia.com/riva_speech";
16
17 import "riva_audio.proto";
18
19 service RivaSpeechSynthesis {
20     // Used to request speech-to-text from the service. Submit a request containing the
21     // desired text and configuration, and receive audio bytes in the requested format.
22     rpc Synthesize(SynthesizeSpeechRequest) returns (SynthesizeSpeechResponse) {}
23
24     // Used to request speech-to-text returned via stream as it becomes available.
25     // Submit a SynthesizeSpeechRequest with desired text and configuration,
26     // and receive stream of bytes in the requested format.
27     rpc SynthesizeOnline(SynthesizeSpeechRequest) returns (stream_
↵SynthesizeSpeechResponse) {}
28 }
29
30 message SynthesizeSpeechRequest {
31     string text = 1;
32     string language_code = 2;
33     // audio encoding params
34     AudioEncoding encoding = 3;
35     int32 sample_rate_hz = 4;
36     // voice params
37     string voice_name = 5;
38 }
39
40 message SynthesizeSpeechResponse {
41     bytes audio = 1;
42 }
43
44 /*
45 */

```

(continues on next page)

(continued from previous page)

\*/

## 32.4 riva\_audio.proto

```

1  // Copyright (c) 2021, NVIDIA CORPORATION. All rights reserved.
2  // NVIDIA CORPORATION and its licensors retain all intellectual property
3  // and proprietary rights in and to this software, related documentation
4  // and any modifications thereto. Any use, reproduction, disclosure or
5  // distribution of this software and related documentation without an express
6  // license agreement from NVIDIA CORPORATION is strictly prohibited.
7
8
9  syntax = "proto3";
10
11 package nvidia.riva;
12
13 option cc_enable_arenas = true;
14
15
16 /*
17  * AudioEncoding specifies the encoding of the audio bytes in the encapsulating
18  * message.
19  */
20 enum AudioEncoding {
21     // Not specified.
22     ENCODING_UNSPECIFIED = 0;
23
24     // Uncompressed 16-bit signed little-endian samples (Linear PCM).
25     LINEAR_PCM = 1;
26
27     // `FLAC` (Free Lossless Audio
28     // Codec) is the recommended encoding because it is
29     // lossless--therefore recognition is not compromised--and
30     // requires only about half the bandwidth of `LINEAR16`. `FLAC` stream
31     // encoding supports 16-bit and 24-bit samples, however, not all fields in
32     // `STREAMINFO` are supported.
33     FLAC = 2;
34
35     // 8-bit samples that compand 14-bit audio samples using G.711 PCMU/mu-law.
36     MULAW = 3;
37
38     // 8-bit samples that compand 13-bit audio samples using G.711 PCMU/a-law.
39     ALAW = 20;
40 }

```

## 32.5 health.proto

```
1  // Copyright (c) 2021, NVIDIA CORPORATION. All rights reserved.
2  //
3  // NVIDIA CORPORATION and its licensors retain all intellectual property
4  // and proprietary rights in and to this software, related documentation
5  // and any modifications thereto. Any use, reproduction, disclosure or
6  // distribution of this software and related documentation without an express
7  // license agreement from NVIDIA CORPORATION is strictly prohibited.
8
9
10 //
11 //Based on gRPC health check protocol - more details found here:
12 //https://github.com/grpc/grpc/blob/master/doc/health-checking.md
13 //
14
15 syntax = "proto3";
16 option go_package = "nvidia.com/riva_speech";
17
18 package grpc.health.v1;
19
20
21 option cc_enable_arenas = true;
22
23 service Health {
24   rpc Check(HealthCheckRequest) returns (HealthCheckResponse);
25   rpc Watch(HealthCheckRequest) returns (stream HealthCheckResponse);
26 }
27
28 message HealthCheckRequest {
29   string service = 1;
30 }
31
32 message HealthCheckResponse {
33   enum ServingStatus {
34     UNKNOWN = 0;
35     SERVING = 1;
36     NOT_SERVING = 2;
37   }
38   ServingStatus status = 1;
39 }
```

See also:

- gRPC API
- Python API
- gRPC & Protocol Buffers

## RIVA - PYTHON API - SAMPLE SERVICES

This session illustrates the samples that demonstrate how Riva clients use Riva services through Riva Speech API server. To facilitate the understanding of the usage of Riva Speech API, you might read [Python API](#)

### 33.1 Prerequisites

To use the Riva Speech API from the client side, install `riva-speech` Python wheel from Riva Quick Start scripts (Riva 1.5.0 is used in this demo).

```
pip install riva_api-1.5.0-beta-py3-none-any.whl
```

---

**Note:** Refer *Riva - Getting Started* to download Riva Quick Start scripts.

---

After installing `riva_api-1.5.0-beta-py3-none-any.whl`, run `riva_start_client.sh` to start a client container

```
bash riva_start_client.sh
```

You may then host a jupyter server, and create notebooks to use the sample services.

```
jupyter notebook --allow-root --ip=0.0.0.0
```

### 33.2 ASR - Recognize

1. Import necessary libraries

```
import grpc
import librosa
import io
import IPython.display as ipd
import riva_api.riva_asr_pb2 as rasr
import riva_api.riva_asr_pb2_grpc as rasr_grpc
import riva_api.riva_audio_pb2 as ra
```

2. Establish connections between client and server

```
channel = grpc.insecure_channel("localhost:50051")
server = rasr_srv.RivaSpeechRecognitionStub(channel)
```

### 3. Read audio

```
path = "/work/wav/sample.wav"
audio, sr = librosa.core.load(path, sr=None)
with io.open(path, 'rb') as fh:
    content = fh.read()
```

---

### Note:

- Besides .wav audio files, Riva also accept .alaw, .mulaw, .flac formatted files with single channel.
- 

### 4. Create a RecognizeRequest

```
req = rasr.RecognizeRequest()
req.audio = content
req.config.encoding = ra.AudioEncoding.LINEAR_PCM
req.config.sample_rate_hertz = sr
req.config.language_code = "en-US"
req.config.max_alternatives = 1
req.config.enable_automatic_punctuation = True
req.config.audio_channel_count = 1
```

### See also:

- *riva\_asr.proto*

```
message RecognizeRequest {
    // Provides information to recognizer that specifies how to process the request.
    RecognitionConfig config = 1;
    // The raw audio data to be processed. The audio bytes must be encoded as specified in
    // `RecognitionConfig`.
    bytes audio = 2;
}
```

```
message RecognitionConfig {
    // The encoding of the audio data sent in the request.
    //
    // All encodings support only 1 channel (mono) audio.
    AudioEncoding encoding = 1;

    // Sample rate in Hertz of the audio data sent in all
    // `RecognizeAudio` messages.
    int32 sample_rate_hertz = 2;

    // Required. The language of the supplied audio as a
    // [BCP-47](https://www.rfc-editor.org/rfc/bcp/bcp47.txt) language tag.
    // Example: "en-US".
    // Currently only en-US is supported
    string language_code = 3;
```

(continues on next page)



(continued from previous page)

```

// Maximum number of recognition hypotheses to be returned.
// Specifically, the maximum number of `SpeechRecognizeAlternative` messages
// within each `SpeechRecognizeResult`.
// The server may return fewer than `max_alternatives`.
// If omitted, will return a maximum of one.
int32 max_alternatives = 4;

// The number of channels in the input audio data.
// ONLY set this for MULTI-CHANNEL recognition.
// Valid values for LINEAR16 and FLAC are `1`-`8`.
// Valid values for OGG_OPUS are `1`-`254`.
// Valid value for MULAW, AMR, AMR_WB and SPEEX_WITH_HEADER_BYTE is only `1`.
// If `0` or omitted, defaults to one channel (mono).
// Note: We only recognize the first channel by default.
// To perform independent recognition on each channel set
// `enable_separate_recognition_per_channel` to `true`.
int32 audio_channel_count = 7;

// If `true`, the top result includes a list of words and
// the start and end time offsets (timestamps) for those words. If
// `false`, no word-level time offset information is returned. The default is
// `false`.
bool enable_word_time_offsets = 8;

// If `true`, adds punctuation to recognition result hypotheses.
// The default `false` value does not add punctuation to result hypotheses.
bool enable_automatic_punctuation = 11;

// This needs to be set to `true` explicitly and `audio_channel_count` > 1
// to get each channel recognized separately. The recognition result will
// contain a `channel_tag` field to state which channel that result belongs
// to. If this is not true, we will only recognize the first channel. The
// request is billed cumulatively for all channels recognized:
// `audio_channel_count` multiplied by the length of the audio.
bool enable_separate_recognition_per_channel = 12;

// Which model to select for the given request. Valid choices: Jasper, Quartznet
string model = 13;

// The verbatim_transcripts flag enables or disable inverse text normalization.
// `true` returns exactly what was said, with no denormalization.
// `false` applies inverse text normalization, also this is the default
bool verbatim_transcripts = 14;

// Custom fields for passing request-level
// configuration options to plugins used in the
// model pipeline.
map<string, string> custom_configuration = 24;
}

```

### See also:

- *riva\_audio.proto*

```
/*
 * AudioEncoding specifies the encoding of the audio bytes in the encapsulating message.
 */
enum AudioEncoding {
    // Not specified.
    ENCODING_UNSPECIFIED = 0;

    // Uncompressed 16-bit signed little-endian samples (Linear PCM).
    LINEAR_PCM = 1;

    // `FLAC` (Free Lossless Audio
    // Codec) is the recommended encoding because it is
    // lossless--therefore recognition is not compromised--and
    // requires only about half the bandwidth of `LINEAR16`. `FLAC` stream
    // encoding supports 16-bit and 24-bit samples, however, not all fields in
    // `STREAMINFO` are supported.
    FLAC = 2;

    // 8-bit samples that compand 14-bit audio samples using G.711 PCMU/mu-law.
    MULAW = 3;

    // 8-bit samples that compand 13-bit audio samples using G.711 PCMU/a-law.
    ALAW = 20;
}
```

### 5. Remote procedure call Recognize

```
response = server.Recognize(req)
```

### See also:

- Recognize remote procedure call is defined as:

```
rpc Recognize(RecognizeRequest) returns (RecognizeResponse) {}
```

### See also:

- *riva\_asr.proto*

```
message RecognizeResponse {
    // Sequential list of transcription results corresponding to
    // sequential portions of audio. Currently only returns one transcript.
    repeated SpeechRecognitionResult results = 1;
}
```

```
// A speech recognition result corresponding to the latest transcript
message SpeechRecognitionResult {

    // May contain one or more recognition hypotheses (up to the
    // maximum specified in `max_alternatives`).
    // These alternatives are ordered in terms of accuracy, with the top (first)
```

(continues on next page)

(continued from previous page)

```

// alternative being the most probable, as ranked by the recognizer.
repeated SpeechRecognitionAlternative alternatives = 1;

// For multi-channel audio, this is the channel number corresponding to the
// recognized result for the audio from that channel.
// For audio_channel_count = N, its output values can range from '1' to 'N'.
int32 channel_tag = 2;

// Length of audio processed so far in seconds
float audio_processed = 3;
}

// Alternative hypotheses (a.k.a. n-best list).
message SpeechRecognitionAlternative {
  // Transcript text representing the words that the user spoke.
  string transcript = 1;

  // The non-normalized confidence estimate. A higher number
  // indicates an estimated greater likelihood that the recognized words are
  // correct. This field is set only for a non-streaming
  // result or, of a streaming result where `is_final=true`.
  // This field is not guaranteed to be accurate and users should not rely on it
  // to be always provided.
  float confidence = 2;

  // A list of word-specific information for each recognized word. Only populated
  // if is_final=true
  repeated WordInfo words = 3;
}

```

6. Print the response from Riva server

```
print(response)
```

```

Full Response Message:
results {
  alternatives {
    transcript: "What is natural language processing? "
    confidence: 1.0
  }
  channel_tag: 1
  audio_processed: 4.800000190734863
}

```

## 33.3 NLP - Core Services

1. Import necessary libraries

```
import grpc
import riva_api.riva_nlp_pb2 as rnlp
import riva_api.riva_nlp_pb2_grpc as rnlp_srv
```

2. Establish connections between client and server

```
channel = grpc.insecure_channel('localhost:50051')
server = rnlp_srv.RivaLanguageUnderstandingStub(channel)
```

3. Select one of core NLP services below to continue.

---

**Note:**

- Notice for every core NLP service, `NLPModelParams` model metadata (name of models to be used) is included in every request to specify model characteristics/ requirements.
- `riva_nlp.proto`

```
// NLPModelParams is a metadata message that is included in every request message
// used by the Core NLP Service and is used to specify model characteristics/requirements
message NLPModelParams {
    // Requested model to use. If unavailable, the request will return an error
    string model_name = 1;
}
```

- You can use `docker logs riva-speech` to check which models are available for services.
- 

### 33.3.1 TransformText

4. Create a `TextTransformRequest` (string -> string)

```
req = rnlp.TextTransformRequest()
req.model.model_name = "riva_punctuation"
req.text.append("add punctuation to this sentence")
req.text.append("do you have any NVDA stocks")
req.text.append("i buy one apple four oranges and two watermelons "
                "for my after-dinner it's going to be very cool")
```

**See also:**

- `riva_nlp.proto`

```
message TextTransformRequest {
    // Each repeated text element is handled independently for handling multiple
    // input strings with a single request
    repeated string text = 1;
    uint32 top_n = 2; //
    NLPModelParams model = 3;
}
```

5. remote procedure call TransformText

```
nlp_resp = server.TransformText(req)
```

**Note:**

- TransformText remote procedure call is defined as:

```
rpc TransformText(TextTransformRequest) returns (TextTransformResponse) {}
```

- *riva\_nlp.proto*

```
// TextTransformResponse is returned by the TransformText method. Responses
// are returned in the same order as they were requested.
message TextTransformResponse {
  repeated string text = 1;
}
```

6. Print the response from Riva server

```
print(nlp_resp)
```

```
text: "Add punctuation to this sentence."
text: "Do you have any NVDA stocks?"
text: "I buy one apple, four oranges and two watermelons for my after-dinner It\'s going_
↳to be very cool."
```

### 33.3.2 ClassifyText

4. Create a TextClassRequest (text -> label)

```
request = rnlp.TextClassRequest()
request.model.model_name = "riva_text_classification_domain"
request.text.append("Is it going to snow in Burlington, Vermont tomorrow night?")
request.text.append("What causes rain?")
request.text.append("What is your favorite season?")
```

**See also:**

- *riva\_nlp.proto*

```
// TextClassRequest is the input message to the ClassifyText service.
message TextClassRequest {
  // Each repeated text element is handled independently for handling multiple
  // input strings with a single request
  repeated string text = 1;

  // Return the top N classification results for each input. 0 or 1 will return top_
  ↳class, otherwise N.
  // Note: Current disabled.
  uint32 top_n = 2;
```

(continues on next page)

(continued from previous page)

```
NLPModelParams model = 3;
}
```

### 5. Remote procedure call ClassifyText

```
ct_response = riva_nlp.ClassifyText(request)
```

---

#### Note:

- ClassifyText remote procedure call is defined as:

```
rpc ClassifyText(TextClassRequest) returns (TextClassResponse) {}
```

- *riva\_nlp.proto*

```
message TextClassResponse {
  repeated ClassificationResult results = 1;
}
```

```
// ClassificationResults contain zero or more Classification messages
// If the number of Classifications is > 1, top_n > 1 must have been
// specified.
```

```
message ClassificationResult {
  repeated Classification labels = 1;
}
```

```
// Classification messages return a class name and corresponding score
```

```
message Classification {
  string class_name = 1;
  float score = 2;
}
```

---

### 6. Print the response from the Riva Server

```
print(ct_response.results)
```

```
[labels {
  class_name: "weather"
  score: 0.9975590109825134
}, labels {
  class_name: "meteorology"
  score: 0.984375
}, labels {
  class_name: "personality"
  score: 0.984375
}]
```

### 33.3.3 ClassifyTokens

4. Create a TokenClassRequest (token -> label)

```
req = rnlp.TokenClassRequest()
req.model.model_name = "riva_ner"
req.text.append("Jensen Huang is the CEO of NVIDIA Corporation, "
               "located in Santa Clara, California")
```

See also:

- *riva\_nlp.proto*

```
// TokenClassRequest is the input message to the ClassifyText service.
message TokenClassRequest {
  // Each repeated text element is handled independently for handling multiple
  // input strings with a single request
  repeated string text = 1;

  // Return the top N classification results for each input. 0 or 1 will return top_
  ↪class, otherwise N.
  // Note: Current disabled.
  uint32 top_n = 3;
  NLPModelParams model = 4;
}
```

5. Remote procedure call ClassifyTokens

```
resp = server.ClassifyTokens(req)
```

Note:

- ClassifyTokens remote procedure call is defined as:

```
rpc ClassifyTokens(TokenClassRequest) returns (TokenClassResponse) {}
```

- *riva\_nlp.proto*

```
// TokenClassResponse returns a single TokenClassSequence per input request
message TokenClassResponse {
  repeated TokenClassSequence results = 1;
}
```

```
// TokenClassSequence is used for returning a sequence of TokenClassValue objects
// in the original order of input tokens
message TokenClassSequence {
  repeated TokenClassValue results = 1;
}
```

```
// TokenClassValue is used to correlate an input token with its classification results
message TokenClassValue {
  string token = 1;
  repeated Classification label = 2;
}
```

(continues on next page)

(continued from previous page)

```
    repeated Span span = 3;
}
```

```
// Span of a particular result
message Span {
    uint32 start = 1;
    uint32 end = 2;
}
```

6. Print the response from the Riva server

```
print(resp)
```

```
results {
  results {
    token: "jensen huang"
    label {
      class_name: "PER"
      score: 0.9972559809684753
    }
    span {
      end: 11
    }
  }
  results {
    token: "nvidia corporation"
    label {
      class_name: "ORG"
      score: 0.9613900184631348
    }
    span {
      start: 27
      end: 44
    }
  }
  results {
    token: "santa clara"
    label {
      class_name: "LOC"
      score: 0.9977059960365295
    }
    span {
      start: 58
      end: 68
    }
  }
  results {
    token: "california"
    label {
      class_name: "LOC"
      score: 0.9961509704589844
    }
  }
}
```

(continues on next page)



(continued from previous page)

```

    }
    span {
      start: 71
      end: 80
    }
  }
}
results {
  results {
    token: "hku"
    label {
      class_name: "ORG"
      score: 0.830374002456665
    }
    span {
      start: 26
      end: 28
    }
  }
}
}

```

## 33.4 NLP - AnalyzeIntent

1. Import necessary libraries

```

import grpc
import riva_api.riva_nlp_pb2_grpc as rnlp_srv
import riva_api.riva_nlp_pb2 as rnlp

```

2. Establish connections between client and server

```

channel = grpc.insecure_channel('localhost:50051')
server = rnlp_srv.RivaLanguageUnderstandingStub(channel)

```

3. Create a AnalyzeIntentRequest

```

req = rnlp.AnalyzeIntentRequest(query="How is the weather today in Hong Kong")

```

See also:

- *riva\_nlp.proto*

```

// AnalyzeIntentRequest is the input message for the AnalyzeIntent service
message AnalyzeIntentRequest {
  // The string to analyze for intent and slots
  string query = 1;
  // Optional configuration for the request, including providing context from previous
  // turns
  // and hardcoding a domain/language
  AnalyzeIntentOptions options = 2;
}

```

### 4. Remote procedure call AnalyzeIntent

```
resp = server.AnalyzeIntent(req)
```

#### Note:

- AnalyzeIntent remote procedure call is defined as:

```
rpc AnalyzeIntent(AnalyzeIntentRequest) returns (AnalyzeIntentResponse) {}
```

#### See also:

- *riva\_nlp.proto*

```
// AnalyzeIntentResponse is returned by the AnalyzeIntent service, and includes
↳ information
// related to the query's intent, (optionally) slot data, and its domain.
message AnalyzeIntentResponse {
    // Intent classification result, including the label and score
    Classification intent = 1;
    // List of tokens explicitly marked as filling a slot relevant to the intent, where the
    // tokens may not exactly match the input (based on the recombined values after
    ↳ tokenization)
    repeated TokenClassValue slots = 2;
    // Returns the inferred domain for the query if not hardcoded in the request. In the
    ↳ case where
    // the domain was hardcoded in AnalyzeIntentRequest, the returned domain is an exact
    ↳ match to the
    // request. In the case where no domain matches the query, intent and slots will be
    ↳ unset.
    //
    // DEPRECATED, use Classification domain field.
    string domain_str = 3;

    // Returns the inferred domain for the query if not hardcoded in the request. In the
    ↳ case where
    // the domain was hardcoded in AnalyzeIntentRequest, the returned domain is an exact
    ↳ match to the
    // request. In the case where no domain matches the query, intent and slots will be
    ↳ unset.
    Classification domain = 4;
}
```

### 5. Print the response from Riva server:

```
print(resp)
```

```
intent {
  class_name: "weather.weather"
  score: 0.9965819716453552
}
slots {
  token: "today"
```

(continues on next page)

(continued from previous page)

```

label {
  class_name: "weatherforecastdaily"
  score: 0.8901410102844238
}
}
slots {
  token: "hong"
  label {
    class_name: "weatherplace"
    score: 0.7496770024299622
  }
}
slots {
  token: "kong"
  label {
    class_name: "weatherplace"
    score: 0.5112429857254028
  }
}
domain_str: "weather"
domain {
  class_name: "weather"
  score: 0.9970700144767761
}

```

## 33.5 TTS - SynthesizeSpeech

1. Import necessary libraries

```

import grpc
import numpy as np
import IPython.display as ipd
import riva_api.riva_tts_pb2 as rtts
import riva_api.riva_tts_pb2_grpc as rtts_srv
import riva_api.audio_pb2 as ra

```

2. Establish connections between client and server

```

channel = grpc.insecure_channel("localhost:50051")
server = rtts_srv.RivaSpeechSynthesisStub(channel)

```

3. Create a SynthesizeSpeechRequest

```

req = rtts.SynthesizeSpeechRequest()
req.text = "Buy Nvidia stocks, You will get rich"
req.language_code = "en-US"
req.encoding = ra.AudioEncoding.LINEAR_PCM
req.sample_rate_hz = 22050
req.voice_name = "ljspeech"

```

See also:

- *riva\_tts.proto*

```
message SynthesizeSpeechRequest {  
    string text = 1;  
    string language_code = 2;  
    // audio encoding params  
    AudioEncoding encoding = 3;  
    int32 sample_rate_hz = 4;  
    // voice params  
    string voice_name = 5;  
}
```

#### 4. Remote procedure call Synthesize

```
resp = server.Synthesize(req)
```

#### Note:

- Synthesize remote procedure call is defined as:

```
rpc Synthesize(SynthesizeSpeechRequest) returns (SynthesizeSpeechResponse) {}
```

```
message SynthesizeSpeechResponse {  
    bytes audio = 1;  
}
```

#### 5. Convert buffer to numpy array and listen

```
audio_samples = np.frombuffer(resp.audio, dtype=np.float32)  
ipd.Audio(audio_samples, rate=22050)
```

## RIVA - PYTHON CLIENT

### 34.1 asr\_client.py

```
import argparse
import wave
import sys
import grpc
import time
import riva_api.audio_pb2 as ra
import riva_api.riva_asr_pb2 as rasr
import riva_api.riva_asr_pb2_grpc as rasr_srv

def get_args():
    parser = argparse.ArgumentParser(description="Streaming transcription via Riva AI_
↳Services")
    parser.add_argument("--server", default="localhost:50051", type=str, help="URI to_
↳GRPC server endpoint")
    parser.add_argument("--audio-file", required=True, help="path to local file to stream
↳")
    parser.add_argument(
        "--show-intermediate", action="store_true", help="show intermediate transcripts_
↳as they are available"
    )
    return parser.parse_args()

def listen_print_loop(responses, show_intermediate=False):
    num_chars_printed = 0
    idx = 0
    for response in responses:
        idx += 1
        if not response.results:
            continue

        result = response.results[0]
        if not result.alternatives:
            continue

        transcript = result.alternatives[0].transcript
```

(continues on next page)

```

    if show_intermediate:
        overwrite_chars = ' ' * (num_chars_printed - len(transcript))

        if not result.is_final:
            sys.stdout.write(">> " + transcript + overwrite_chars + '\r')
            sys.stdout.flush()

            num_chars_printed = len(transcript) + 3

        else:
            print("## " + transcript + overwrite_chars + "\n")
            num_chars_printed = 0
    else:
        if result.is_final:
            print(f"## {transcript.encode('utf-8')}\\n")
            sys.stdout.buffer.write(transcript.encode('utf-8'))

CHUNK = 1024
args = get_args()
wf = wave.open(args.audio_file, 'rb')

channel = grpc.insecure_channel(args.server)
client = rasr_srv.RivaSpeechRecognitionStub(channel)

config = rasr.RecognitionConfig(
    encoding=ra.AudioEncoding.LINEAR_PCM,
    sample_rate_hertz=wf.getframerate(),
    language_code="en-US",
    max_alternatives=1,
    enable_automatic_punctuation=True,
)
streaming_config = rasr.StreamingRecognitionConfig(config=config, interim_results=True)

# read data
def generator(w, s):
    yield rasr.StreamingRecognizeRequest(streaming_config=s)
    d = w.readframes(CHUNK)
    while len(d) > 0:
        yield rasr.StreamingRecognizeRequest(audio_content=d)
        d = w.readframes(CHUNK)

responses = client.StreamingRecognize(generator(wf, streaming_config))
listen_print_loop(responses, show_intermediate=args.show_intermediate)

```

## 34.2 Riva ASR service

```
python3 $RIVA_QS/asr_client.py --audio-file $path
```





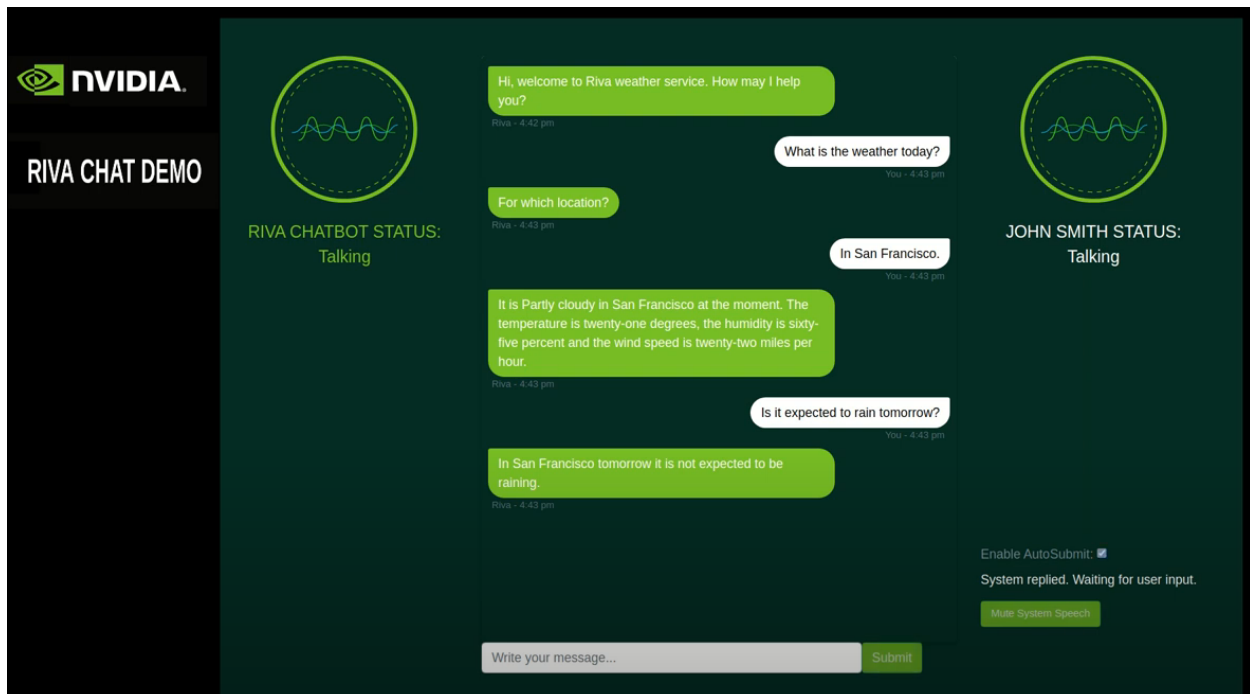
## RIVA - STREAMING ASR

```
python3 $RIVA_QS/examples/riva_streaming_asr_client.py --input-file $path
```



## RIVA - VIRTUAL ASSISTANT

This sample implementation of conversational AI describes how to build Virtual Assistant using NGC sample image.



To see more details, click

- [Riva - Sample Virtual Assistant](#)

### 36.1 Start Riva services

Follow *Riva - Getting Started* and start Riva services

```
bash riva_start.sh
```

## 36.2 Download sample image

Make sure your NGC API is configured before pulling the sample image

```
docker pull nvcr.io/nvidia/riva/riva-speech-client:1.4.0-beta-samples
```

## 36.3 Run Riva service

First we create a container from the sample image we just pulled from NGC.

```
docker run -it --rm -p 8009:8009 nvcr.io/nvidia/riva/riva-speech-client:1.4.0-beta-  
↪samples /bin/bash
```

Edit `riva_config` in `samples/virtual-assistant/config.py`:

```
riva_config = {  
    "RIVA_SPEECH_API_URL": "<riva-host-ip>:50051",  
    "ENABLE_QA": "QA unavailable in this VA version. Coming soon",  
    "WEATHERSTACK_ACCESS_KEY": "<weather-api-key>",  
    "VERBOSE": True # print logs/details for diagnostics  
}
```

- `RIVA_SPEECH_API_URL`: Replace the IP & port with your hosted Riva endpoint
- `WEATHERSTACK_ACCESS_KEY`: Get your access key at [here](#)

## 36.4 Start Virtual Assistant app

After configuring `config.py`, we can now start the web application of Virtual Assistant.

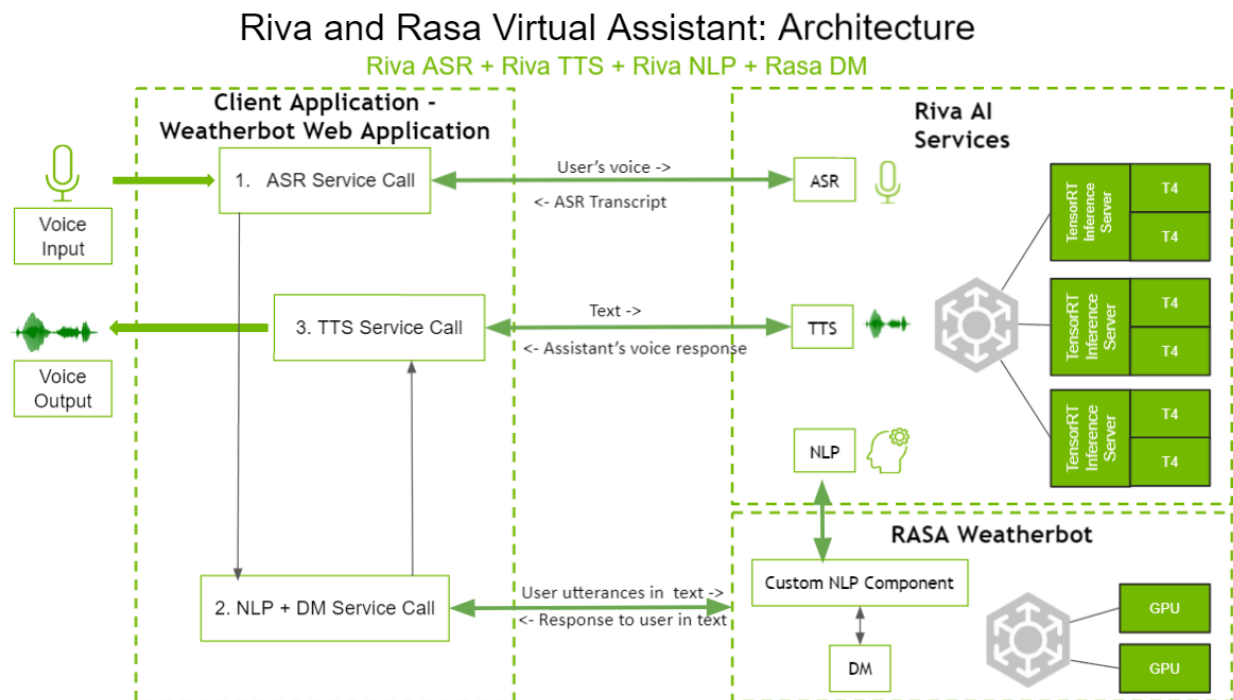
```
python3 main.py
```

Visit `https://<riva-host-ip>:8009/rivaWeather` to use Virtual Assistant.

## RIVA - VIRTUAL ASSISTANT - RASA

This sample implementation of conversational AI describes how to build Virtual Assistant with [Rasa](#) dialogue manager hosted.

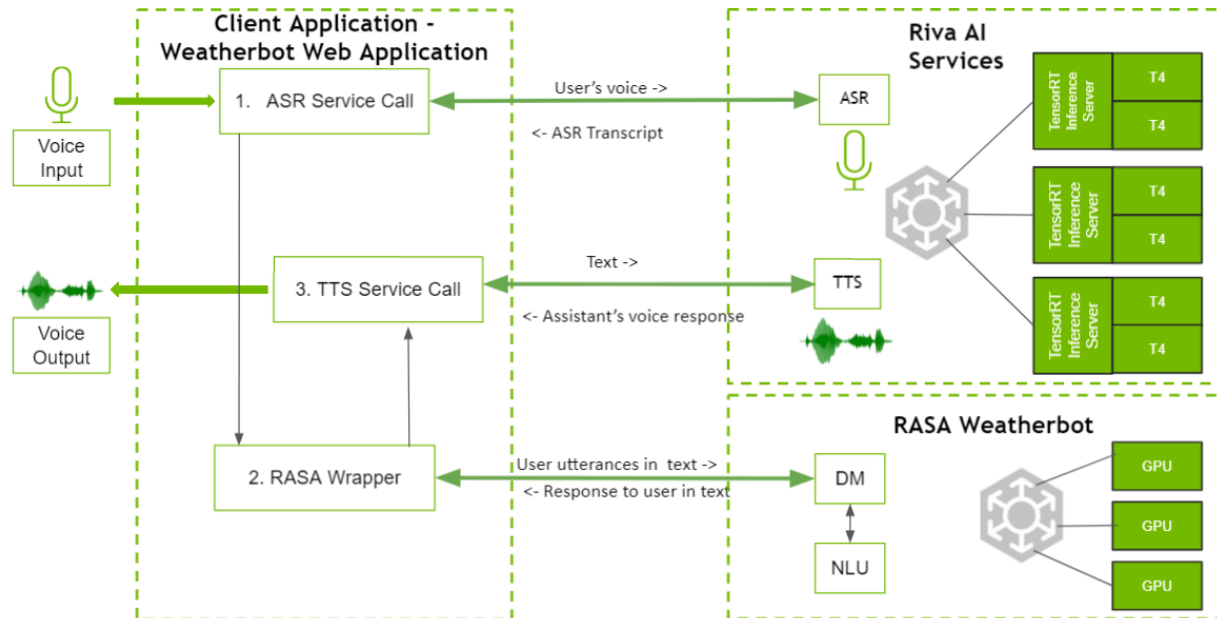
Option 1: Riva ASR + Riva TTS + Riva NLP + Rasa dialog manager



Option 2: Riva ASR + Riva TTS + Rasa NLU + Rasa dialog manager

## Riva and Rasa Virtual Assistant: Architecture

Riva ASR + Riva TTS + Rasa NLU + Rasa DM



### 37.1 Overview of RVA Rasa

```
.
|-- README.md
|-- config.py
|-- rasa-riva-weatherbot-webapp
|   |-- client
|   |   |-- webapplication
|   |   |   |-- cert.pem
|   |   |   |-- key.pem
|   |   |   |-- server
|   |   |   |   |-- __init__.py
|   |   |   |   |-- server.py
|   |   |-- start_web_application.py
|   |-- ui
|   |   |-- README.md
|   |   |-- img
|   |   |   |-- Rivadm.png
|   |   |   |-- User.png
|   |   |-- index.html
|   |   |-- script.js
|   |   |-- static
|   |   |   |-- stylesheets
|   |   |   |   |-- index.css
|   |   |   |-- svg_files
|   |   |       |-- another_sample.svg
|   |   |       |-- circle.svg
|   |   |       |-- logo.svg
```

(continues on next page)

(continued from previous page)

```
|
|                                     |-- logo_sample.svg
|                                     |-- nv_logo.svg
|                                     |-- nv_logo_1.svg
|                                     |-- nvidia_logo.svg
|                                     |-- nvidia_name.svg
|                                     |-- riva_name.png
|                                     |-- sample_.svg
|                                     |-- speech.svg
|                                     |-- speech_logo.svg
|
|                                     |-- style.css
|-- main.py
|-- riva
|   |-- __init__.py
|   |-- asr
|     |-- __init__.py
|     |-- asr.py
|   |-- chatbot
|     |-- __init__.py
|     |-- chatbot.py
|     |-- chatbots_multiconversations_management.py
|   |-- rasa
|     |-- __init__.py
|     |-- rasa.py
|   |-- tts
|     |-- __init__.py
|     |-- tts.py
|     |-- tts_stream.py
|-- rasa-weatherbot
|   |-- __init__.py
|   |-- __pycache__
|     |-- __init__.cpython-37.pyc
|     |-- actions.cpython-37.pyc
|   |-- actions
|     |-- __init__.py
|     |-- actions.py
|     |-- weather.py
|   |-- components.py
|   |-- config
|     |-- config_rasanlp.yml
|     |-- config_rivanlp.yml
|   |-- data
|     |-- nlu_rasanlp.yml
|     |-- nlu_rivanlp.yml
|     |-- rules_rasanlp.yml
|     |-- rules_rivanlp.yml
|     |-- stories_rasanlp.yml
|     |-- stories_rivanlp.yml
|   |-- domain
|     |-- domain_rasanlp.yml
|     |-- domain_rivanlp.yml
|   |-- endpoints.yml
|   |-- riva
```

(continues on next page)

(continued from previous page)

```
| |-- __init__.py
| |-- nlp
| |   |-- __init__.py
| |   |-- nlp.py
|-- tests
|   |-- conversation_tests.md
```

### 37.1.1 RVA-Rasa Pipeline

#### Rasa Action Server & Rasa Server

#### Riva & Web application

rasa-riva-weatherbot-webapp/main.py:

```
import os, sys

root_folder = (os.path.abspath(os.path.join(os.path.dirname(__file__), os.path.pardir)))
sys.path.append(root_folder)

from config import client_config

if __name__ == '__main__':
    if client_config["CLIENT_APPLICATION"] == "WEBAPPLICATION":
        from client.webapplication.start_web_application import start_web_application
        start_web_application()
```

### 37.1.2 Rasa Server & Rasa Action

## 37.2 Prerequisites

1. Pull Riva Sample container

```
docker pull nvcr.io/nvidia/riva/riva-speech-client:1.5.0-beta-samples
```

2. Prepare network configuration files

Create a Riva config `config.py`, and replace `RASA_API_URL` with the local machine ip you use to host Rasa server.

**Warning:** Please use local machine ip (not `0.0.0.0`/ `127.0.0.1`/ `localhost`)

```
# =====
# Copyright (c) 2020, NVIDIA CORPORATION. All rights reserved.
#
# The License information can be found under the "License" section of the
# README.md file.
# =====
```

(continues on next page)



(continued from previous page)

```

client_config = {
    "CLIENT_APPLICATION": "WEBAPPLICATION", # Default and only config value for this_
    ↪version
    "PORT": 5555, # The port your flask app will be hosted at
    "DEBUG": False, # When this flag is set, the UI displays detailed riva data
    "VERBOSE": True # print logs/details for diagnostics
}

riva_config = {
    "RIVA_SPEECH_API_URL": "0.0.0.0:50051", # Replace the IP port with your hosted RIVA_
    ↪endpoint
    "WEATHERSTACK_ACCESS_KEY": "8a82ce50a069bfaaa3db4427ecca723d", # Get your access_
    ↪key at - https://weatherstack.com/
    "VERBOSE": True # print logs/details for diagnostics
}

rasa_config = {
    "VERBOSE": False, # Print logs/details for diagnostics
    "RASA_API_URL": "http://<rasa-host-ip>:5005", # Replace the IP & Port with the rasa-
    ↪weatherbot's IP & Port
}

asr_config = {
    "VERBOSE": False, # Print logs/details for diagnostics
    "SAMPLING_RATE": 16000, # The Sampling Rate for the audio input file. The only value_
    ↪currently supported is 16000
    "LANGUAGE_CODE": "en-US", # The language code as a BCP-47 language tag. The only_
    ↪value currently supported is "en-US"
    "ENABLE_AUTOMATIC_PUNCTUATION": True, # Enable or Disable punctuation in the_
    ↪transcript generated. The only value currently supported by the chatbot is True_
    ↪(Although Riva ASR supports both True & False)
}

tts_config = {
    "VERBOSE": True, # Print logs/details for diagnostics
    "SAMPLE_RATE": 22050, # The speech is generated at this sampling rate. The only_
    ↪value currently supported is 22050
    "LANGUAGE_CODE": "en-US", # The language code as a BCP-47 language tag. The only_
    ↪value currently supported is "en-US"
    "VOICE_NAME": "ljspeech", # The voice name for the speech generated. The only value_
    ↪currently supported is "ljspeech"
}

rivanlp_config = {
    "VERBOSE": False, # Print logs/details for diagnostics
    "NLU_FALLBACK_THRESHOLD": 0.3 # When Intent's confidence/score is less than this_
    ↪value, intent is set to nlu_fallback
}

```

3. Create a Riva config endpoints.yml

```
# =====
# Copyright (c) 2020, NVIDIA CORPORATION. All rights reserved.
#
# The License information can be found under the "License" section of the
# README.md file.
# =====

# This file contains the different endpoints your bot can use.

# Server where the models are pulled from.
# https://rasa.com/docs/rasa/user-guide/configuring-http-api/#fetching-models-from-a-
  ↳server/

#models:
# url: http://my-server.com/models/default_core@latest
# wait_time_between_pulls: 10 # [optional](default: 100)

# Server which runs your custom actions.
# https://rasa.com/docs/rasa/core/actions/#custom-actions/

action_endpoint:
  url: "http://0.0.0.0:5055/webhook"
# Tracker store which is used to store the conversations.
# By default the conversations are stored in memory.
# https://rasa.com/docs/rasa/api/tracker-stores/

#tracker_store:
#   type: redis
#   url: <host of the redis instance, e.g. localhost>
#   port: <port of your redis instance, usually 6379>
#   db: <number of your database within redis, e.g. 0>
#   password: <password used for authentication>
#   use_ssl: <whether or not the communication is encrypted, default false>

#tracker_store:
#   type: mongod
#   url: <url to your mongo instance, e.g. mongodb://localhost:27017>
#   db: <name of the db within your mongo instance, e.g. rasa>
#   username: <username used for authentication>
#   password: <password used for authentication>

# Event broker which all conversation events should be streamed to.
# https://rasa.com/docs/rasa/api/event-brokers/

#event_broker:
# url: localhost
# username: username
# password: password
# queue: queue
```

3. Create a local dir `cfg_dir` and store `config.py` and `endpoints.yml` to `cfg_dir`

## 37.3 1. Start Riva Server

Follow *Riva - Getting Started* to start enable Riva services.

## 37.4 2. Start Rasa Action Server

### 2.1 Create Riva Sample container

```
docker run -it --rm -p 5055:5055 -v $(pwd)/cfg_dir:/workspace/cfg nvcr.io/nvidia/riva/
↪riva-speech-client:1.5.0-beta-samples /bin/bash
```

2.2 Replace `/workspace/samples/virtual-assistant-rasa/config.py` with `/workspace/cfg_dir/config.py`

2.3 Replace `/workspace/samples/virtual-assistant-rasa/rasa-weatherbot/endpoints.yml` with `/workspace/cfg_dir/endpoints.yml`

2.4 Activate the Rasa Python environment

```
. /pythonenvs/rasa/bin/activate
```

2.5 Navigate to the `/workspace/samples/virtual-assistant-rasa/rasa-weatherbot` directory and Start Rasa Action server

```
cd /workspace/samples/virtual-assistant-rasa/rasa-weatherbot
rasa run actions --actions actions
```

## 37.5 3. Start Rasa Server

Open another terminal and create Riva sample container

### 3.1 Create Riva Sample container

```
sudo docker run -it --rm -p 5005:5005 -v $(pwd)/cfg_dir:/workspace/cfg nvcr.io/nvidia/
↪riva/riva-speech-client:1.5.0-beta-samples /bin/bash
```

3.2 Replace `/workspace/samples/virtual-assistant-rasa/config.py` with `/workspace/cfg_dir/config.py`

3.3 Replace `/workspace/samples/virtual-assistant-rasa/rasa-weatherbot/endpoints.yml` with `/workspace/cfg_dir/endpoints.yml`

3.4 Activate the Rasa Python environment.

```
. /pythonenvs/rasa/bin/activate
```

3.5 Navigate to the `/workspace/samples/virtual-assistant-rasa/rasa-weatherbot` directory and Run the Rasa training

```
cd /workspace/samples/virtual-assistant-rasa/rasa-weatherbot
```

If you use Riva NLP:

```
rasa train -c config/config_rivanlp.yml \  
  -d domain/domain_rivanlp.yml \  
  --out models/models_rivanlp/ \  
  --data data/nlu_rivanlp.yml \  
  data/rules_rivanlp.yml \  
  data/stories_rivanlp.yml
```

Or you use Rasa NLU:

```
rasa train -c config/config_rasanlp.yml \  
  -d domain/domain_rasanlp.yml \  
  --out models/models_rasanlp/ \  
  --data data/nlu_rasanlp.yml \  
  data/rules_rasanlp.yml data/stories_rasanlp.yml
```

### 3.6 Start the Rasa Server

- If you use Riva NLP:

```
rasa run -m models/models_rivanlp/ --enable-api \  
  --log-file out.log --endpoints endpoints.yml
```

- If you use Rasa NLU:

```
rasa run -m models/models_rasanlp/ --enable-api \  
  --log-file out.log --endpoints endpoints.yml
```

## 37.6 4. Start RVA Rasa Server

### 4.1 Create Riva Sample container

```
sudo docker run -it --rm -p 5555:5555 -v $(pwd)/cfg_dir:/workspace/cfg nvcr.io/nvidia/  
riva/riva-speech-client:1.5.0-beta-samples /bin/bash
```

4.2 Replace `/workspace/samples/virtual-assistant-rasa/config.py` with `/workspace/cfg_dir/config.py`

4.3 Replace `/workspace/samples/virtual-assistant-rasa/rasa-weatherbot/endpoints.yml` with `/workspace/cfg_dir/endpoints.yml`

### 4.4 Activate Chatbot Python environment

```
. /pythonenvs/client/bin/activate
```

4.5 Navigate to the `/workspace/samples/virtual-assistant-rasa/rasa-riva-weatherbot-webapp` directory and Start Chatbot Client server

```
python3 main.py
```

## 37.7 5. Open RVA Rasa

Browse *[https://\[riva chatbot server host IP\]:5555/rivaWeather](https://[riva chatbot server host IP]:5555/rivaWeather)*.



## GUIDE DESCRIPTION

End-to-end Riva development from custom inference engines to custom Riva deployment.





---

CHAPTER  
**THIRTYNINE**

---

**AUTHOR**

Season POON

Position: machine learning intern @ NVIDIA

Email: [spoon@nvidia.com](mailto:spoon@nvidia.com)