

# Malayalam Text-to Speech system

**Guide: Jiby J P**

**Project Coordinator: Manilal D.L**

Group 16

Kurian Benoy 34

[kurianbenoy@mec.ac.in](mailto:kurianbenoy@mec.ac.in)

# Objective

- To build a Text to Speech(TTS) system in Malayalam
- Obtain the state of art result

- Module1 : EDA, dataset collection
- Module2: Train first TTS system in Malayalam
- Module3: Fine tune TTS system
- Module4: User Interface

**Work Done**

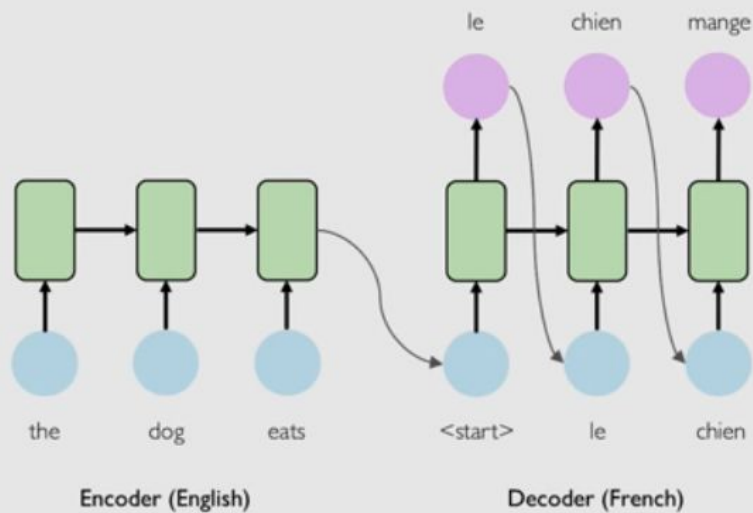
# Tactron2 paper Implementation

Tactron2, a neural network architecture for speech synthesis directly from text. The system consists of recurrent sequence to sequence feature extraction networks that map character embeddings to mel-scale spectrogram, which is currently one of the best Text to speech architectures, followed by modified Wavenet model acting as a synthesis domain. Model achieves 4.53 MOS score. Seq-2-seq models made of RNNs act as backbone networks for this.

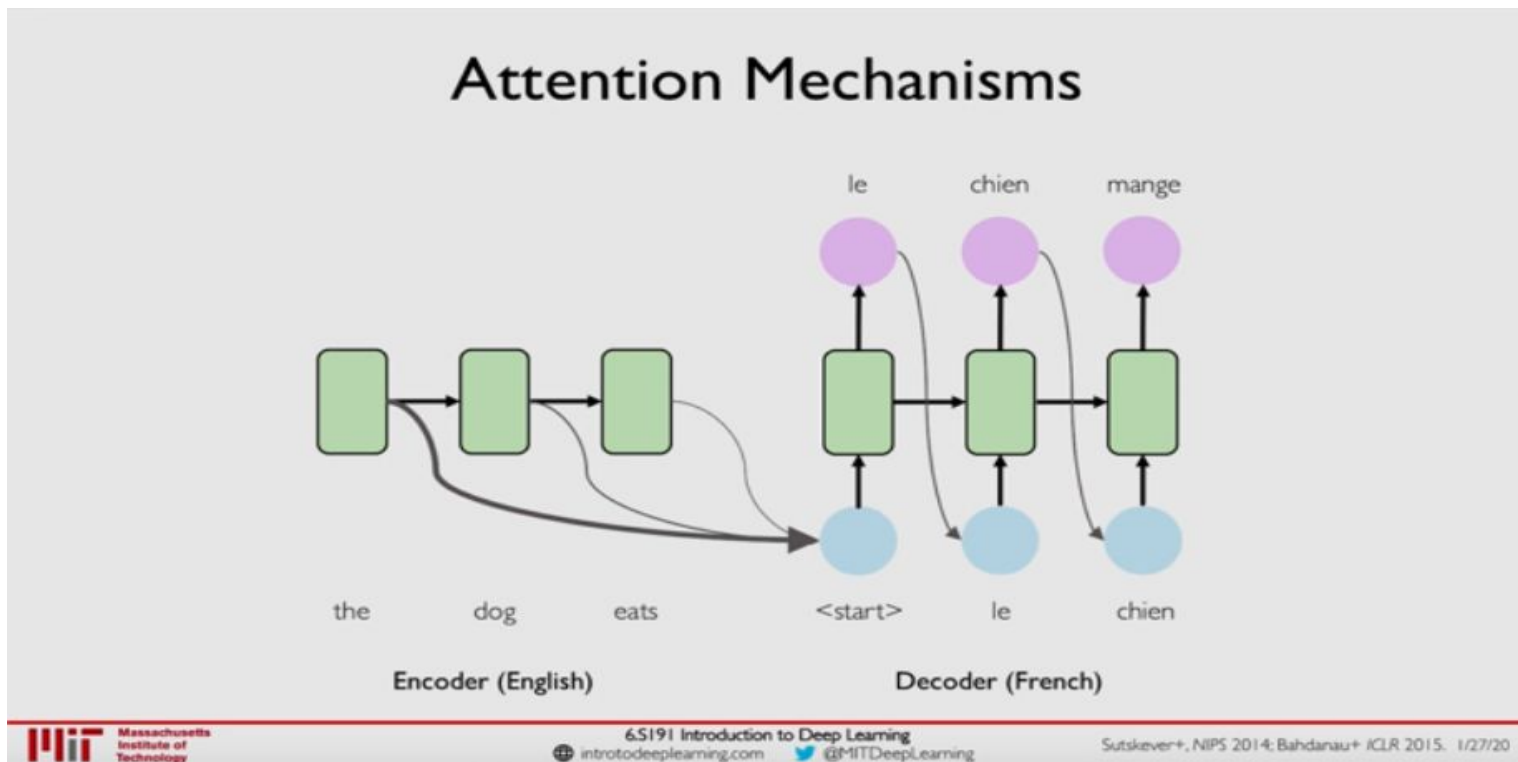
Aim: To understand the paper and implement architecture

# Sequence to Sequence models

## Example Task: Machine Translation



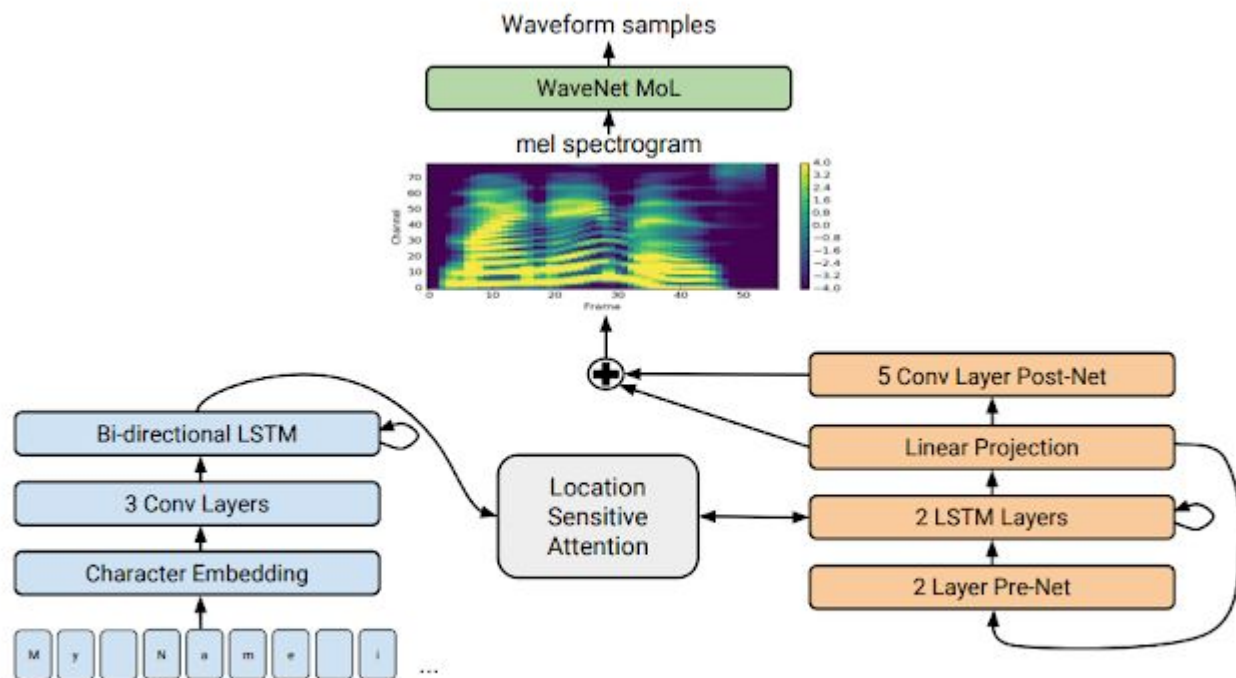
# Sequence to Sequence models



# Components of Tactron2 architecture

1. Location Layer
2. Encoder
3. Decoder
4. Attention Model
5. Posnet
- 6.





A detailed look at Tacotron 2's model architecture. The lower half of the image describes the sequence-to-sequence model that maps a sequence of letters to a spectrogram. For technical details, please refer to [the paper](#).

# Location Layer

```
In [7]: 1 class LocationLayer(nn.Module):
2         def __init__(self, attention_n_filters, attention_kernel_size,
3                     attention_dim):
4             super(LocationLayer, self).__init__()
5             padding = int((attention_kernel_size - 1) / 2)
6             self.location_conv = ConvNorm(2, attention_n_filters,
7                                         kernel_size=attention_kernel_size,
8                                         padding=padding, bias=False, stride=1,
9                                         dilation=1)
10            self.location_dense = LinearNorm(attention_n_filters, attention_dim,
11                                           bias=False, w_init_gain='tanh')
12
13        def forward(self, attention_weights_cat):
14            processed_attention = self.location_conv(attention_weights_cat)
15            processed_attention = processed_attention.transpose(1, 2)
16            processed_attention = self.location_dense(processed_attention)
17            return processed_attention
```

# Attention Module

```
In [8]: 1 class Attention(nn.Module):
2     def __init__(self, attention_rnn_dim, embedding_dim, attention_dim,
3                 attention_location_n_filters, attention_location_kernel_size):
4         super(Attention, self).__init__()
5         self.query_layer = LinearNorm(attention_rnn_dim, attention_dim,
6                                     bias=False, w_init_gain='tanh')
7         self.memory_layer = LinearNorm(embedding_dim, attention_dim, bias=False,
8                                     w_init_gain='tanh')
9         self.v = LinearNorm(attention_dim, 1, bias=False)
10        self.location_layer = LocationLayer(attention_location_n_filters,
11                                          attention_location_kernel_size,
12                                          attention_dim)
13        self.score_mask_value = -float("inf")
14
15    def get_alignment_energies(self, query, processed_memory,
16                             attention_weights_cat):
17
18        processed_query = self.query_layer(query.unsqueeze(1))
19        processed_attention_weights = self.location_layer(attention_weights_cat)
20        energies = self.v(torch.tanh(
21            processed_query + processed_attention_weights))
22
23        energies = energies.squeeze(-1)
24        return energies
25
26    def forward(self, attention_hidden_state, memory, processed_memory,
27               attention_weights_cat, mask):
28        alignment = self.get_alignment_energies(
29            attention_hidden_state, processed_memory, attention_weights_cat)
30
31        if mask is not None:
32            alignment.data.masked_fill_(mask, self.score_mask_value)
33
34        attention_weights = F.softmax(alignment, dim=1)
35        attention_context = torch.bmm(attention_weights.unsqueeze(1), memory)
36        attention_context = attention_context.squeeze(1)
37
38        return attention_context, attention_weights
```

# PostNet Module

```
1 class Postnet(nn.Module):
2     """Postnet
3     - Five 1-d convolution with 512 channels and kernel size 5
4     """
5
6     def __init__(self, hparams):
7         super(Postnet, self).__init__()
8         self.convolutions = nn.ModuleList()
9
10        self.convolutions.append(
11            nn.Sequential(
12
13                (hparams.n_mel_channels, hparams.postnet_embedding_dim,
14                 kernel_size=hparams.postnet_kernel_size, stride=1,
15                 padding=int((hparams.postnet_kernel_size - 1) / 2),
16                 dilation=1, w_init_gain='tanh'),
17                nn.BatchNorm1d(hparams.postnet_embedding_dim))
18            )
19
20        for i in range(1, hparams.postnet_n_convolutions - 1):
21            self.convolutions.append(
22                nn.Sequential(
23                    ConvNorm(hparams.postnet_embedding_dim,
24                             hparams.postnet_embedding_dim,
25                             kernel_size=hparams.postnet_kernel_size, stride=1,
26                             padding=int((hparams.postnet_kernel_size - 1) / 2),
27                             dilation=1, w_init_gain='tanh'),
28                    nn.BatchNorm1d(hparams.postnet_embedding_dim))
29                )
30
31            self.convolutions.append(
32                nn.Sequential(
33                    ConvNorm(hparams.postnet_embedding_dim, hparams.n_mel_channels,
34                             kernel_size=hparams.postnet_kernel_size, stride=1,
35                             padding=int((hparams.postnet_kernel_size - 1) / 2),
36                             dilation=1, w_init_gain='linear'),
37                    nn.BatchNorm1d(hparams.n_mel_channels))
38                )
39
40        def forward(self, x):
41            for i in range(len(self.convolutions) - 1):
42                x = F.dropout(torch.tanh(self.convolutions[i](x)), 0.5, self.training)
43            x = F.dropout(self.convolutions[-1](x), 0.5, self.training)
44
45            return x
46
```

# Encoder Module

```
In [11]: 1 class Encoder(nn.Module):
2         """Encoder module:
3         - Three 1-d convolution banks
4         - Bidirectional LSTM
5         """
6         def __init__(self, hparams):
7             super(Encoder, self).__init__()
8
9             convolutions = []
10            for _ in range(hparams.encoder_n_convolutions):
11                conv_layer = nn.Sequential(
12                    ConvNorm(hparams.encoder_embedding_dim,
13                            hparams.encoder_embedding_dim,
14                            kernel_size=hparams.encoder_kernel_size, stride=1,
15                            padding=int((hparams.encoder_kernel_size - 1) / 2),
16                            dilation=1, w_init_gain='relu'),
17                    nn.BatchNorm1d(hparams.encoder_embedding_dim))
18                convolutions.append(conv_layer)
19            self.convolutions = nn.ModuleList(convolutions)
20
21            self.lstm = nn.LSTM(hparams.encoder_embedding_dim,
22                                int(hparams.encoder_embedding_dim / 2), 1,
23                                batch_first=True, bidirectional=True)
24
25            def forward(self, x, input_lengths):
26                for conv in self.convolutions:
27                    x = F.dropout(F.relu(conv(x)), 0.5, self.training)
28
29                x = x.transpose(1, 2)
30
31                # pytorch tensor are not reversible, hence the conversion
32                input_lengths = input_lengths.cpu().numpy()
33                x = nn.utils.rnn.pack_padded_sequence(
34                    x, input_lengths, batch_first=True)
35
36                self.lstm.flatten_parameters()
37                outputs, _ = self.lstm(x)
38
39                outputs, _ = nn.utils.rnn.pad_packed_sequence(
40                    outputs, batch_first=True)
41
42                return outputs
43
44            def inference(self, x):
45                for conv in self.convolutions:
46                    x = F.dropout(F.relu(conv(x)), 0.5, self.training)
47
48                x = x.transpose(1, 2)
49
50                self.lstm.flatten_parameters()
51                outputs, _ = self.lstm(x)
52
53                return outputs
54
```

# Decoder

```
1 class Decoder(nn.Module):
2     def __init__(self, hparams):
3         super(Decoder, self).__init__()
4         self.n_mel_channels = hparams.n_mel_channels
5         self.n_frames_per_step = hparams.n_frames_per_step
6         self.encoder_embedding_dim = hparams.encoder_embedding_dim
7         self.attention_rnn_dim = hparams.attention_rnn_dim
8         self.decoder_rnn_dim = hparams.decoder_rnn_dim
9         self.prenet_dim = hparams.prenet_dim
10        self.max_decoder_steps = hparams.max_decoder_steps
11        self.gate_threshold = hparams.gate_threshold
12        self.p_attention_dropout = hparams.p_attention_dropout
13        self.p_decoder_dropout = hparams.p_decoder_dropout
14
15        self.prenet = Prenet(
16            hparams.n_mel_channels * hparams.n_frames_per_step,
17            [hparams.prenet_dim, hparams.prenet_dim])
18
19        self.attention_rnn = nn.LSTMCell(
20            hparams.prenet_dim + hparams.encoder_embedding_dim,
21            hparams.attention_rnn_dim)
22
23        self.attention_layer = Attention(
24            hparams.attention_rnn_dim, hparams.encoder_embedding_dim,
25            hparams.attention_dim, hparams.attention_location_n_filters,
26            hparams.attention_location_kernel_size)
27
28        self.decoder_rnn = nn.LSTMCell(
29            hparams.attention_rnn_dim + hparams.encoder_embedding_dim,
30            hparams.decoder_rnn_dim, 1)
31
32        self.linear_projection = LinearNorm(
33            hparams.decoder_rnn_dim + hparams.encoder_embedding_dim,
34            hparams.n_mel_channels * hparams.n_frames_per_step)
35
36        self.gate_layer = LinearNorm(
37            hparams.decoder_rnn_dim + hparams.encoder_embedding_dim, 1,
38            bias=True, w_init_gain='sigmoid')
39
40    def get_go_frame(self, memory):
41        B = memory.size(0)
42        decoder_input = Variable(memory.data.new(
43            B, self.n_mel_channels * self.n_frames_per_step).zero_())
44        return decoder_input
45
46    def initialize_decoder_states(self, memory, mask):
47        B = memory.size(0)
48        MAX_TIME = memory.size(1)
```

# Decoder

```
50 self.attention_hidden = Variable(memory.data.new(
51     B, self.attention_rnn_dim).zero_())
52 self.attention_cell = Variable(memory.data.new(
53     B, self.attention_rnn_dim).zero_())
54
55 self.decoder_hidden = Variable(memory.data.new(
56     B, self.decoder_rnn_dim).zero_())
57 self.decoder_cell = Variable(memory.data.new(
58     B, self.decoder_rnn_dim).zero_())
59
60 self.attention_weights = Variable(memory.data.new(
61     B, MAX_TIME).zero_())
62 self.attention_weights_cum = Variable(memory.data.new(
63     B, MAX_TIME).zero_())
64 self.attention_context = Variable(memory.data.new(
65     B, self.encoder_embedding_dim).zero_())
66
67 self.memory = memory
68 self.processed_memory = self.attention_layer.memory_layer(memory)
69 self.mask = mask
70
71 def parse_decoder_inputs(self, decoder_inputs):
72     # (B, n_mel_channels, T_out) -> (B, T_out, n_mel_channels)
73     decoder_inputs = decoder_inputs.transpose(1, 2)
74     decoder_inputs = decoder_inputs.view(
75         decoder_inputs.size(0),
76         int(decoder_inputs.size(1)/self.n_frames_per_step), -1)
77     # (B, T_out, n_mel_channels) -> (T_out, B, n_mel_channels)
78     decoder_inputs = decoder_inputs.transpose(0, 1)
79     return decoder_inputs
80
81 def parse_decoder_outputs(self, mel_outputs, gate_outputs, alignments):
82     # (T_out, B) -> (B, T_out)
83     alignments = torch.stack(alignments).transpose(0, 1)
84     # (T_out, B) -> (B, T_out)
85     gate_outputs = torch.stack(gate_outputs).transpose(0, 1)
86     gate_outputs = gate_outputs.contiguous()
87     # (T_out, B, n_mel_channels) -> (B, T_out, n_mel_channels)
88     mel_outputs = torch.stack(mel_outputs).transpose(0, 1).contiguous()
89     # decouple frames per step
90     mel_outputs = mel_outputs.view(
91         mel_outputs.size(0), -1, self.n_mel_channels)
92     # (B, T_out, n_mel_channels) -> (B, n_mel_channels, T_out)
93     mel_outputs = mel_outputs.transpose(1, 2)
94
95     return mel_outputs, gate_outputs, alignments
96
```



# Decoder

```
97 def decode(self, decoder_input):
98     cell_input = torch.cat((decoder_input, self.attention_context), -1)
99     self.attention_hidden, self.attention_cell = self.attention_rnn(
100         cell_input, (self.attention_hidden, self.attention_cell))
101     self.attention_hidden = F.dropout(
102         self.attention_hidden, self.p_attention_dropout, self.training)
103
104     attention_weights_cat = torch.cat(
105         (self.attention_weights.unsqueeze(1),
106          self.attention_weights_cum.unsqueeze(1)), dim=1)
107     self.attention_context, self.attention_weights = self.attention_layer(
108         self.attention_hidden, self.memory, self.processed_memory,
109         attention_weights_cat, self.mask)
110
111     self.attention_weights_cum += self.attention_weights
112     decoder_input = torch.cat(
113         (self.attention_hidden, self.attention_context), -1)
114     self.decoder_hidden, self.decoder_cell = self.decoder_rnn(
115         decoder_input, (self.decoder_hidden, self.decoder_cell))
116     self.decoder_hidden = F.dropout(
117         self.decoder_hidden, self.p_decoder_dropout, self.training)
118
119     decoder_hidden_attention_context = torch.cat(
120         (self.decoder_hidden, self.attention_context), dim=1)
121     decoder_output = self.linear_projection(
122         decoder_hidden_attention_context)
123
124     gate_prediction = self.gate_layer(decoder_hidden_attention_context)
125     return decoder_output, gate_prediction, self.attention_weights
126
127 def forward(self, memory, decoder_inputs, memory_lengths):
128
129     decoder_input = self.get_go_frame(memory).unsqueeze(0)
130     decoder_inputs = self.parse_decoder_inputs(decoder_inputs)
131     decoder_inputs = torch.cat((decoder_input, decoder_inputs), dim=0)
132     decoder_inputs = self.prenet(decoder_inputs)
133
134     self.initialize_decoder_states(
135         memory, mask=~get_mask_from_lengths(memory_lengths))
136
137     mel_outputs, gate_outputs, alignments = [], [], []
138     while len(mel_outputs) < decoder_inputs.size(0) - 1:
139         decoder_input = decoder_inputs[len(mel_outputs)]
140         mel_output, gate_output, attention_weights = self.decode(
141             decoder_input)
142         mel_outputs += [mel_output.squeeze(1)]
143         gate_outputs += [gate_output.squeeze(1)]
144         alignments += [attention_weights]
145
```



# Gluing all together(with Tacotron)

```
1 class Tacotron2(nn.Module):
2     def __init__(self,
3                 num_chars,
4                 num_speakers,
5                 r,
6                 postnet_output_dim=80,
7                 decoder_output_dim=80,
8                 attn_type='original',
9                 attn_win=False,
10                attn_norm='softmax',
11                prenet_type='original',
12                prenet_dropout=True,
13                forward_attn=False,
14                trans_agent=False,
15                forward_attn_mask=False,
16                location_attn=True,
17                attn_K=5,
18                separate_stopnet=True,
19                bidirectional_decoder=False):
20         super(Tacotron2, self).__init__()
21         self.postnet_output_dim = postnet_output_dim
22         self.decoder_output_dim = decoder_output_dim
23         self.n_frames_per_step = r
24         self.bidirectional_decoder = bidirectional_decoder
25         decoder_dim = 512 if num_speakers > 1 else 512
26         encoder_dim = 512 if num_speakers > 1 else 512
27         proj_speaker_dim = 80 if num_speakers > 1 else 0
28         # embedding layer
29         self.embedding = nn.Embedding(num_chars, 512, padding_idx=0)
30         std = sqrt(2.0 / (num_chars + 512))
31         val = sqrt(3.0) * std # uniform bounds for std
32         self.embedding.weight.data.uniform_(-val, val)
33         if num_speakers > 1:
34             self.speaker_embedding = nn.Embedding(num_speakers, 512)
35             self.speaker_embedding.weight.data.normal_(0, 0.3)
36             self.speaker_embeddings = None
37             self.speaker_embeddings_projected = None
38         self.encoder = Encoder(encoder_dim)
39         self.decoder = Decoder(decoder_dim, self.decoder_output_dim, r, attn_type, attn_win,
40                               attn_norm, prenet_type, prenet_dropout,
41                               forward_attn, trans_agent, forward_attn_mask,
42                               location_attn, attn_K, separate_stopnet, proj_speaker_dim)
43         if self.bidirectional_decoder:
44             self.decoder_backward = copy.deepcopy(self.decoder)
45         self.postnet = Postnet(self.postnet_output_dim)
46
47     def _init_states(self):
48         self.speaker_embeddings = None
49         self.speaker_embeddings_projected = None
50
51     @staticmethod
52     def shape_outputs(mel_outputs, mel_outputs_postnet, alignments):
53         mel_outputs = mel_outputs.transpose(1, 2)
54         mel_outputs_postnet = mel_outputs_postnet.transpose(1, 2)
55         return mel_outputs, mel_outputs_postnet, alignments
56
```

# Gluing all together(with Tactron)

```
57 def forward(self, text, text_lengths, mel_specs=None, speaker_ids=None):
58     self._init_states()
59     # compute mask for padding
60     mask = sequence_mask(text_lengths).to(text.device)
61     embedded_inputs = self.embedding(text).transpose(1, 2)
62     encoder_outputs = self.encoder(embedded_inputs, text_lengths)
63     encoder_outputs = self._add_speaker_embedding(encoder_outputs,
64                                                  speaker_ids)
65     decoder_outputs, alignments, stop_tokens = self.decoder(
66         encoder_outputs, mel_specs, mask)
67     postnet_outputs = self.postnet(decoder_outputs)
68     postnet_outputs = decoder_outputs + postnet_outputs
69     decoder_outputs, postnet_outputs, alignments = self.shape_outputs(
70         decoder_outputs, postnet_outputs, alignments)
71     if self.bidirectional_decoder:
72         decoder_outputs_backward, alignments_backward = self._backward_inference(mel_specs, encoder_outputs)
73         return decoder_outputs, postnet_outputs, alignments, stop_tokens, decoder_outputs_backward, alignments_backward
74     return decoder_outputs, postnet_outputs, alignments, stop_tokens
75
76 @torch.no_grad()
77 def inference(self, text, speaker_ids=None):
78     embedded_inputs = self.embedding(text).transpose(1, 2)
79     encoder_outputs = self.encoder.inference(embedded_inputs)
80     encoder_outputs = self._add_speaker_embedding(encoder_outputs,
81                                                  speaker_ids)
82     mel_outputs, alignments, stop_tokens = self.decoder.inference(
83         encoder_outputs)
84     mel_outputs_postnet = self.postnet(mel_outputs)
85     mel_outputs_postnet = mel_outputs + mel_outputs_postnet
86     mel_outputs, mel_outputs_postnet, alignments = self.shape_outputs(
87         mel_outputs, mel_outputs_postnet, alignments)
88     return mel_outputs, mel_outputs_postnet, alignments, stop_tokens
89
90 def inference_truncated(self, text, speaker_ids=None):
91     """
92     Preserve model states for continuous inference
93     """
94     embedded_inputs = self.embedding(text).transpose(1, 2)
95     encoder_outputs = self.encoder.inference_truncated(embedded_inputs)
96     encoder_outputs = self._add_speaker_embedding(encoder_outputs,
97                                                  speaker_ids)
98     mel_outputs, alignments, stop_tokens = self.decoder.inference_truncated(
99         encoder_outputs)
100     mel_outputs_postnet = self.postnet(mel_outputs)
101     mel_outputs_postnet = mel_outputs + mel_outputs_postnet
102     mel_outputs, mel_outputs_postnet, alignments = self.shape_outputs(
103         mel_outputs, mel_outputs_postnet, alignments)
104     return mel_outputs, mel_outputs_postnet, alignments, stop_tokens
105
106 def _backward_inference(self, mel_specs, encoder_outputs, mask):
107     decoder_outputs_b, alignments_b = self.decoder_backward(
108         encoder_outputs, torch.flip(mel_specs, dims=(1,)), mask,
109         self.speaker_embeddings_projected)
110     decoder_outputs_b = decoder_outputs_b.transpose(1, 2)
111     return decoder_outputs_b, alignments_b
112
113 def _add_speaker_embedding(self, encoder_outputs, speaker_ids):
114     if hasattr(self, "speaker_embedding") and speaker_ids is None:
115         raise RuntimeError("[!] Model has speaker embedding layer but speaker_id is not provided")
116     if hasattr(self, "speaker_embedding") and speaker_ids is not None:
117         speaker_embeddings = self.speaker_embedding(speaker_ids)
118
119         speaker_embeddings.unsqueeze_(1)
120         speaker_embeddings = speaker_embeddings.expand(encoder_outputs.size(0),
121                                                       encoder_outputs.size(1),
122                                                       -1)
123         encoder_outputs = encoder_outputs + speaker_embeddings
124     return encoder_outputs
```

# Gluing all together(with Tactron)

```
97 def decode(self, decoder_input):
98     cell_input = torch.cat((decoder_input, self.attention_context), -1)
99     self.attention_hidden, self.attention_cell = self.attention_rnn(
100         cell_input, (self.attention_hidden, self.attention_cell))
101     self.attention_hidden = F.dropout(
102         self.attention_hidden, self.p_attention_dropout, self.training)
103
104     attention_weights_cat = torch.cat(
105         (self.attention_weights.unsqueeze(1),
106          self.attention_weights_cum.unsqueeze(1)), dim=1)
107     self.attention_context, self.attention_weights = self.attention_layer(
108         self.attention_hidden, self.memory, self.processed_memory,
109         attention_weights_cat, self.mask)
110
111     self.attention_weights_cum += self.attention_weights
112     decoder_input = torch.cat(
113         (self.attention_hidden, self.attention_context), -1)
114     self.decoder_hidden, self.decoder_cell = self.decoder_rnn(
115         decoder_input, (self.decoder_hidden, self.decoder_cell))
116     self.decoder_hidden = F.dropout(
117         self.decoder_hidden, self.p_decoder_dropout, self.training)
118
119     decoder_hidden_attention_context = torch.cat(
120         (self.decoder_hidden, self.attention_context), dim=1)
121     decoder_output = self.linear_projection(
122         decoder_hidden_attention_context)
123
124     gate_prediction = self.gate_layer(decoder_hidden_attention_context)
125     return decoder_output, gate_prediction, self.attention_weights
126
127 def forward(self, memory, decoder_inputs, memory_lengths):
128
129     decoder_input = self.get_go_frame(memory).unsqueeze(0)
130     decoder_inputs = self.parse_decoder_inputs(decoder_inputs)
131     decoder_inputs = torch.cat((decoder_input, decoder_inputs), dim=0)
132     decoder_inputs = self.prenet(decoder_inputs)
133
134     self.initialize_decoder_states(
135         memory, mask=~get_mask_from_lengths(memory_lengths))
136
137     mel_outputs, gate_outputs, alignments = [], [], []
138     while len(mel_outputs) < decoder_inputs.size(0) - 1:
139         decoder_input = decoder_inputs[len(mel_outputs)]
140         mel_output, gate_output, attention_weights = self.decode(
141             decoder_input)
142         mel_outputs += [mel_output.squeeze(1)]
143         gate_outputs += [gate_output.squeeze(1)]
144         alignments += [attention_weights]
145
```

**Thank you!**