

Clearly images! :)

Convs

```

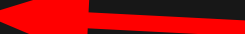
173
174 class AudioEncoder(nn.Module):
175     def __init__(self, n_mels: int, n_ctx: int, n_state: int, n_head: int, n_layer: int):
176         super().__init__()
177         self.conv1 = Conv1d(n_mels, n_state, kernel_size=3, padding=1)
178         self.conv2 = Conv1d(n_state, n_state, kernel_size=3, stride=2, padding=1)
179         self.register_buffer("positional_embedding", sinusoids(n_ctx, n_state))
180         self.blocks = nn.ModuleList([ResidualAttentionBlock(n_state, n_head) for _ in range(n_layer)])
181         self.ln_post = LayerNorm(n_state)
182
183     def forward(self, x: Tensor):
184         """
185         x: torch.Tensor, shape = (batch_size, n_mels, n_ctx)
186         the mel spectrogram of the audio
187         """
188         print('00:', x.shape) # torch.Size([1, 80, 3000])
189         x = F.gelu(self.conv1(x))
190         print('01:', x.shape) # torch.Size([1, 384, 3000])
191         x = F.gelu(self.conv2(x))
192         print('02:', x.shape) # torch.Size([1, 384, 1500])
193         x = x.permute(0, 2, 1)
194         print('03:', x.shape) # torch.Size([1, 1500, 384])
195         assert x.shape[1:] == self.positional_embedding.shape, "incorrect audio shape"
196         x = (x + self.positional_embedding).to(x.dtype)
197         for block in self.blocks: x = block(x)
198         x = self.ln_post(x)
199         return x
200

```

<https://github.com/openai/whisper/blob/main/whisper/model.py#L174>



```
(xld) root@Ubuntu-2404-noble-amd64-base ~/whisper # python workbench.py  
100%|██████████████████████████████████████████| 72.1M/72.1M [00:12<00:00, 5.91MiB/s]  
00: torch.Size([1, 80, 3000])  
01: torch.Size([1, 384, 3000])  
02: torch.Size([1, 384, 1500])  
03: torch.Size([1, 1500, 384])  
00: torch.Size([1, 80, 3000])  
01: torch.Size([1, 384, 3000])  
02: torch.Size([1, 384, 1500])  
03: torch.Size([1, 1500, 384])  
He is quite content to die.
```

 twice