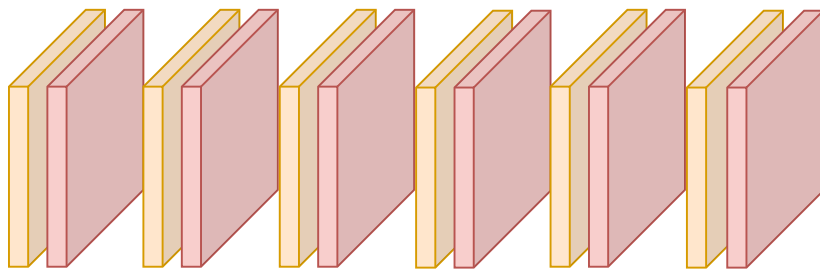


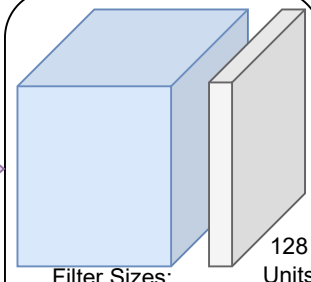
# 1-D CNN

Speech  
Frames



In: 1	In: 2	In: 4	In: 8	In: 16	In: 32
Out: 2	Out: 4	Out: 8	Out: 16	Out: 32	Out: 40
Kernel: 5	Kernel: 5	Kernel: 7	Kernel: 9	Kernel: 11	Kernel: 11
Dilate: 2	Dilate: 2	Dilate: 3	Dilate: 4	Dilate: 5	Dilate: 5

# Reference Encoder

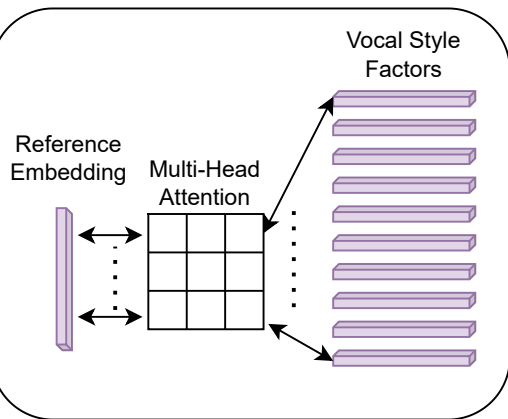


Filter Sizes:  
32 -> 32 -> 64 ->  
64 -> 128 -> 128

128  
Units






Reference  
Embedding

# Vocal Style Factor Layer



Weighted  
Sum  
 $\oplus$

Speaker Identity  
Embedding

-  1-D Dilated Convolution
-  SELU Activation
-  6-Layer 2-D CNN
-  GRU with 128 units
-  Embedding