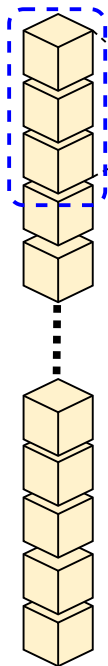
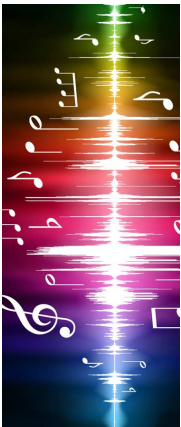


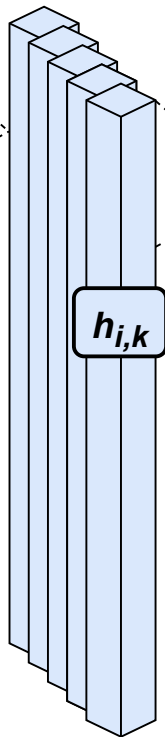
Audio signal  
(1D)

Input:  $1 \times N$



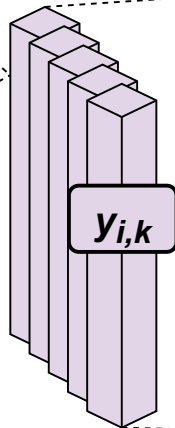
Convolutional layer  
(K filters)

$K \times (N-M+1)$



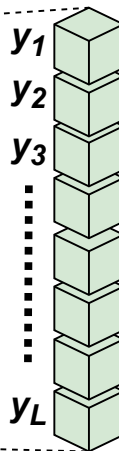
Max-pooling layer

$K \times (N-M+1) / 2$



Flattened output

$1 \times (L = K \cdot (N-M+1) / 2)$



Fully-Connected layer +  
Non-linear activation +  
Normalization +  
Dropout

Output:  $1 \times O$

