

Input
4s Mono Audio



Backbone CNN

Linear Layer
 512×512

1D Batchnorm

ReLU

Dropout $p=0.5$

Linear Layer
 $512 \times \text{\#Classes}$

Sigmoid

Classification

p_1

p_2

...

p_N

