

Input
4s Mono Audio



Backbone CNN

Linear Layer
 512×512

1D Batchnorm

ReLU

Dropout $p=0.5$

Linear Layer
 512×31

Regression

v_1

v_2

...

v_N

Loudness

Timbre
Descriptors