

Google Blaze Face:

- * Lightweight feature extraction network inspired by MobileNet v1/v2.

- * Produces 6 facial Keypoint coordinates (eye centers, ear tragus, mouth center and nose tip).

- * On an $S \times S \times C$ input a $K \times K$ depthwise convolution involves $S^2 c K^2$ multiply-add operations.

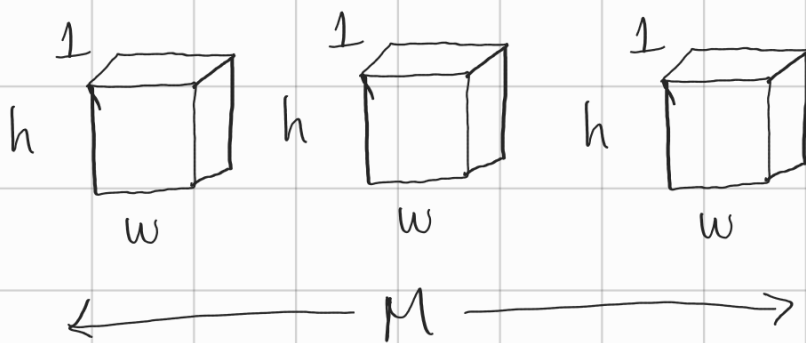
Feature Extractor :

Input: (128, 128, 3)

MobileNet Architecture:

* Depthwise Convolution:

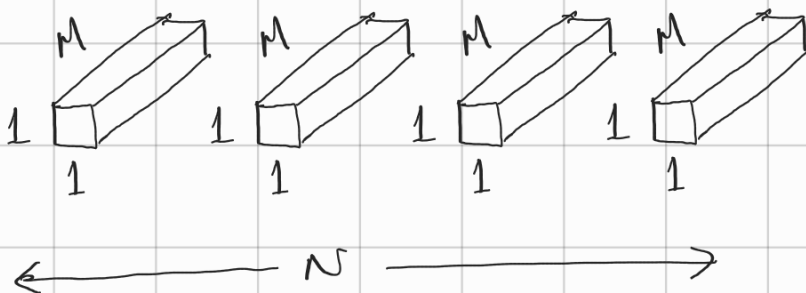
Applies a single filter to each channel.



$\therefore h = \text{height}$
 $w = \text{width}$
 $M = \# \text{channels}$

* Pointwise Convolution:

Applies a 1×1 convolution among the input channels



$\therefore w = h = 1$

$M = \# \text{Input channels}$

$N = \# \text{Output Channels}$

Model:

$$F: X \rightarrow Y$$

$$X \in \mathbb{R}^{h \times w \times ch}$$

$$Y \in \mathbb{Z}^c$$

$\therefore h$: height

w : width

ch : channels

c : classes

Input (224 x 224 x 3)

Architecture:

Layers:

- Depthwise Conv
- Pointwise Conv
- Avg Pooling (7x7)
- Fully Connected
- Softmax

• Depthwise:

$$\hat{G}_{k,l,m} = \sum_{i,j} \hat{K}_{i,j,m} \cdot F_{k+i-1,l+j-1,m}$$

Where $\hat{K} \rightarrow$ Depthwise Conv Kernel ($D_K \times D_K \times M$)

· Softmax:

$$\sigma(\vec{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

$$\sigma: \mathbb{R}^K \rightarrow [0, 1]^K$$

$\therefore \vec{z}$ = Input Vector

$K = \#$ classes

e^{z_i} = Exp Input

e^{z_j} = Exp Output

Note: Google's Blaze Face implements in Depthwise convolution a Kernel Size of 5×5