

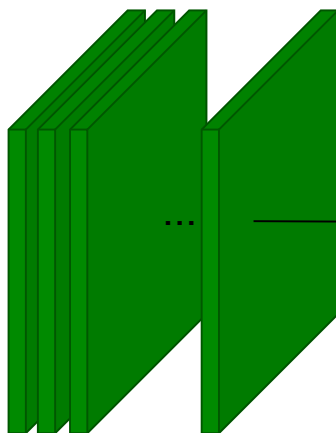
Input Patches

Feature Extractors

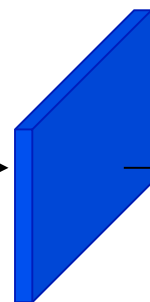
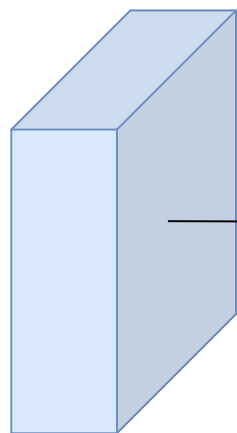
Patch Embeddings

Classifiers

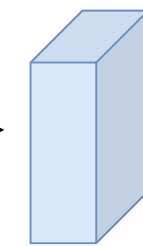
$s = 0$



$b \times 3 \times 76 \times 76$



$1 \times b \times 128$



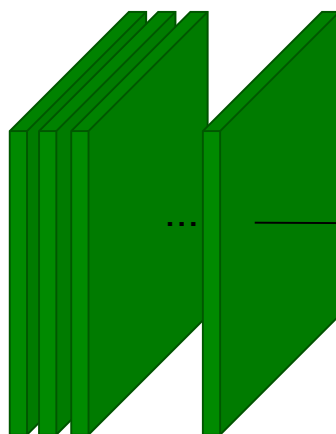
$1 \times b \times 7$ ($s = 0$ Instance Predictions)

$s = 0$ Bag Prediction

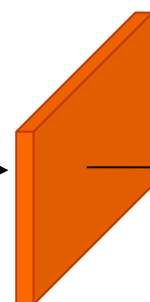
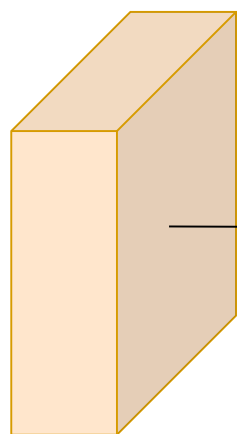


$1 \times 1 \times 7$

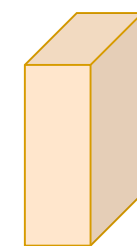
$s = 1$



$4b \times 3 \times 76 \times 76$



$1 \times 4b \times 128$



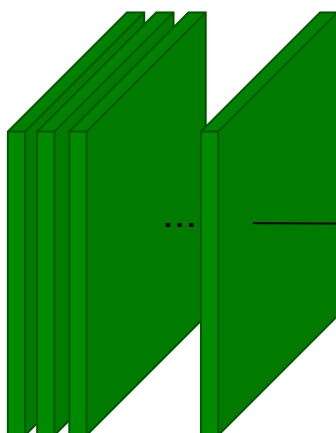
$1 \times 4b \times 7$ ($s = 1$ Instance Predictions)

$s = 1$ Bag Prediction

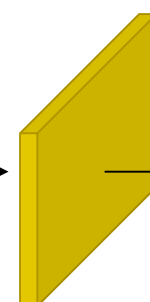
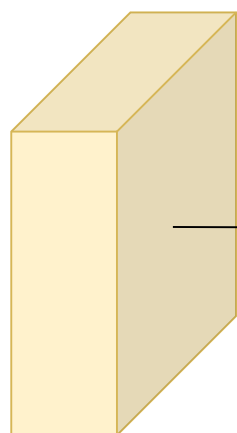


$1 \times 1 \times 7$

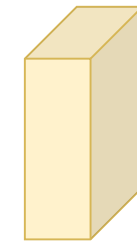
$s = 2$



$16b \times 3 \times 76 \times 76$



$1 \times 16b \times 128$



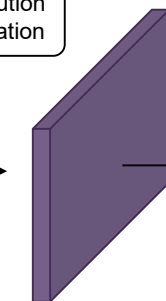
$1 \times 16b \times 7$ ($s = 2$ Instance Predictions)

$s = 2$ Bag Prediction

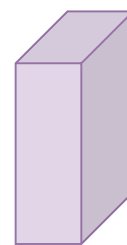


$1 \times 1 \times 7$

Multi-resolution
Concatenation



$1 \times 16b \times 384$



$1 \times 16b \times 7$ ($s = m$ Instance Predictions)

$s = m$ Bag Prediction



$1 \times 1 \times 7$

Multi-Res Multi-Out Model Only