

# Deep Learning Optimization

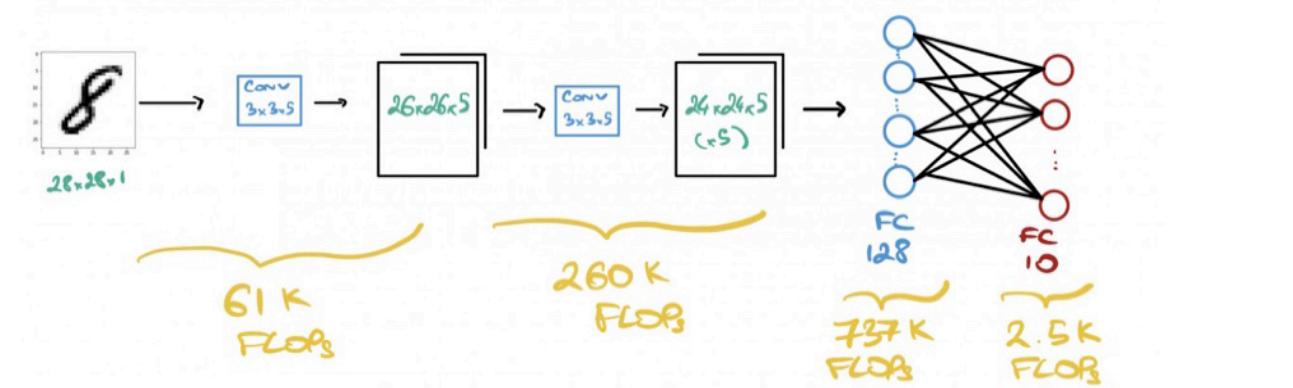
## Optimizing a Deep Learning Model

## Calculating the Inference Time

### Calculating the Inference Time

Time = FLOPs / FLOPS

### Calculating the FLOPs



### Pooling Layers

With Strides:  $FLOPs = (Height/Stride) \times (Width/Stride) \times Depth$

### Fully Connected

MACs = Input Size x Output Size

### Convolutions

MACs = Number of Kernel x Kernel Shape x Output Shape

### MACs: Multiply-Accumulate Computations

An operation that includes an addition and a multiplication

1 MAC = 2 FLOPs

### FLOPs: Floating Point Operations per Second

Tells us how good our hardware is

Tells us how complicated our model is

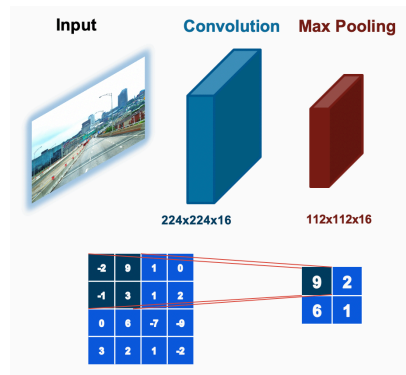
### FLOPs: Floating Point Operations

The number of computations the model will have to perform



Reduce the amount of parameters being passed from one layer to another

### Pooling



### Separable Convolutions

A depthwise followed by a pointwise convolution

Normal Convolution: ~20 M FLOPs

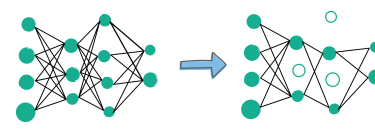
Depthwise Convolution:  $FLOPs = 9 \times 9 \times 9 \times 7 \times 7 \times 2 = 38,402$

Pointwise Convolution:  $FLOPs = 64 \times 19 \times 19 \times 11 \times 2 = 415,872$

Separable Convolution: ~734K FLOPs

Remove the weights that are redundant or that don't matter by setting them to 0

### Model Pruning



### Quantization

Map values from a large set to a smaller one

Example: 2.9089002  $\rightarrow$  2.9

### Weight Sharing

Share weights between layers to reduce the amount of weights we store

Using the K-Means algorithm for example

### Compression

### Reducing the Size of a Model

### Knowledge Distillation

Transfer the knowledge learned by a large and accurate model (teacher) to a smaller and lighter one (student)

