# M3 – Machine Learning for Computer Vision

Project: Deep learning classification - **Final Presentation**
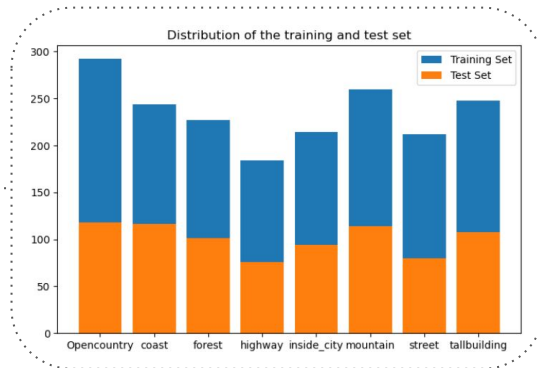
Group 7: Guillem Capellera, Johnny Nuñez and Anna Oliveras

# Outline

- Week 1: Bag of Visual Words framework

- Week 2: Beyond BoVW → SVMs, Spatial Pyramids, Fisher Vectors

Handcrafted methods: Bag of Visual Words

- Week 3: From hand crafted to learnt features

- Week 4: Fine tuning of pre-trained CNNs → Densenet121

- Week 5: Training a CNN from scratch
  - Initial CNN
  - CNN Refinement
  - Residual Connections
  - Residual Network
  - Recap

Data driven methods: Deep Convolutional Networks

- Conclusions

# Datasets


Distribution of the training and test set

- We have 8 classes: coast, forest, highway, inside city, mountain, open country, street, tall buildings

- Big dataset : MIT_split → total of 2288 images

- Small dataset: MIT_small_train_1 → Train with only 50 images for each class !
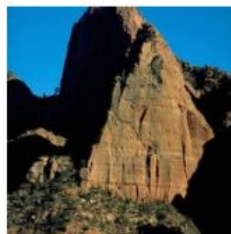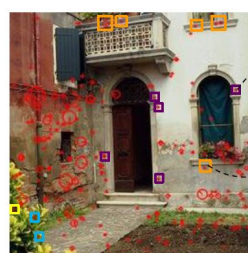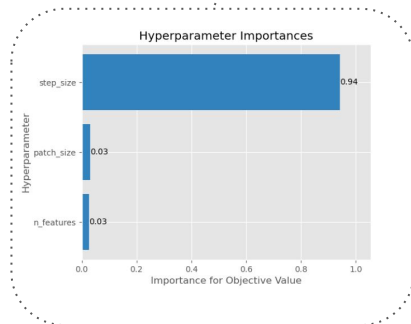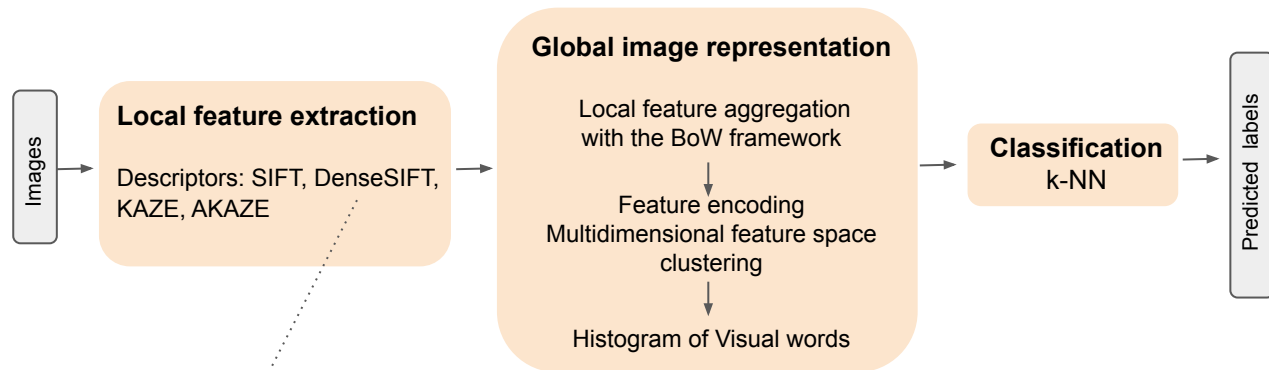
| Opencountry | tallbuilding | highway | street | mountain | coast | inside_city |

# Week 1: Bag of Visual Words framework



Images → **Local feature extraction**

Descriptors: SIFT, DenseSIFT, KAZE, AKAZE

→ **Global image representation**

Local feature aggregation with the BoW framework

↓

Feature encoding
Multidimensional feature space clustering

↓

Histogram of Visual words

→ **Classification**
k-NN

→ Predicted labels

**Grid search best parameters (Optuna)**

Descriptor: DenseSIFT with
- n_features = 251
- patch_size = 3
- step size = 75
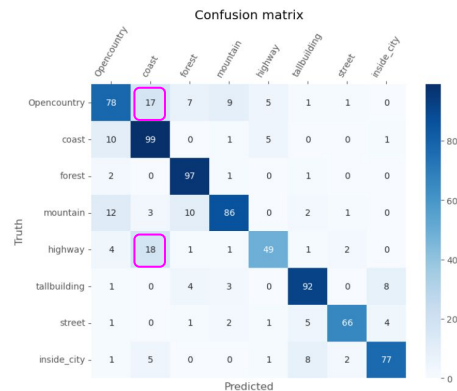
Clustering
- Codebook size k = 1024

Dimensionality reduction
- PCA with n_components = 46

Classifier
- n_neighbors = 18
- distance = euclidean

**Accuracy =** 0.95      **f1-score =** 0.8

Hyperparameter Importances

step_size — 0.94
patch_size — 0.03
n_features — 0.03

Feature encoding

Multidimensional feature space

**Histogram of Visual words**

Confusion matrix
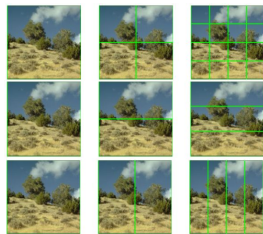
# Week 2: Beyond BoVW

**Spatial Pyramids**

Multiple divisions:
1x1, 2x2, 4x4
vertical 1x2, 1x3, 1x4
horizontal 2x1, 3x1, 4x1



**SVM and Kernels**

- Linear
- RBF
- Histogram intersection

**Feature normalization**

- Power Norm
- L2 Norm

**Fischer Vectors**

Inlude higher order statistics: mean, covariance of local descriptors → GMM clustering
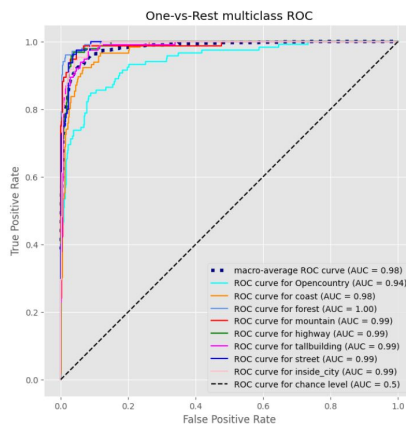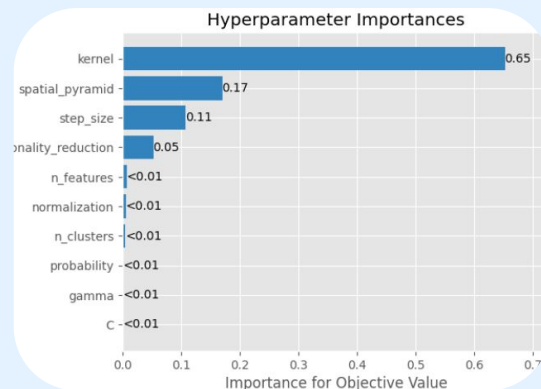
Accuracy = 0.96  (Without hyperparameter optimization)
F1 score = 0.83

**Grid search best parameters (Optuna)→ We use Dense SIFT**

- num features = 178
- **step size = 18**
- num cluster = 798
- num components = 69
- gamma = 0.00445
- C = 4.38
- **dim_reduc = PCA**
- **kernel = RBF**
- Normalization = power
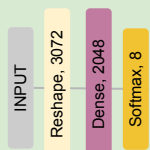- **Spatial_pyramid = vertical 1x4**
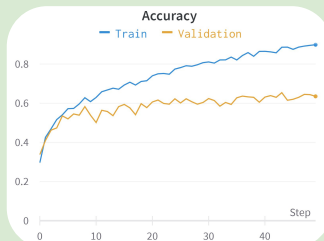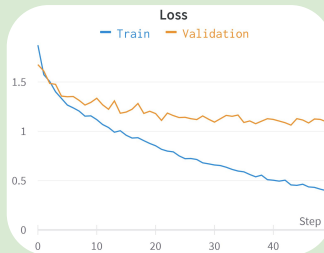
**Accuracy = 0.96**
**F1 score  = 0.86**


Hyperparameter Importances


One-vs-Rest multiclass ROC


Confusion matrix

# Week 3: From hand crafted to learnt features

## Initial Model



INPUT · Reshape, 3072 · Dense, 2048 · Softmax, 8

- 6M parameters
- Accuracy = **0.61**



## Best Model



INPUT · Reshape, 3072 · Dense, 4096 · BNorm, 4096 · Dropout, 4096 · Dense, 2048 · BNorm, 2048 · Dropout, 2048 · Dense, 1024 · BNorm, 1024 · Dropout, 1024 · Dense, 256 · BNorm, 256 · Dropout, 256 · Softmax, 8
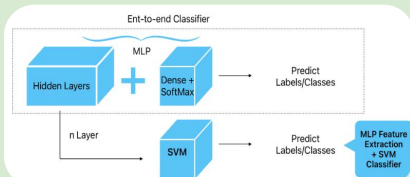
- 23M parameters
- Accuracy = **0.596**



## Deep features + SVM

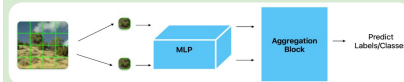Extracting the features after each dense layer output → similar results (best with output dense 4096)

- Accuracy = **0.41**



## Patch based MLP

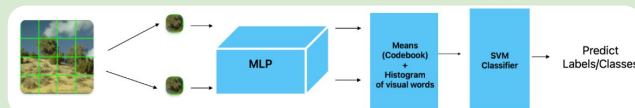→ Best patch size: 32 x 32
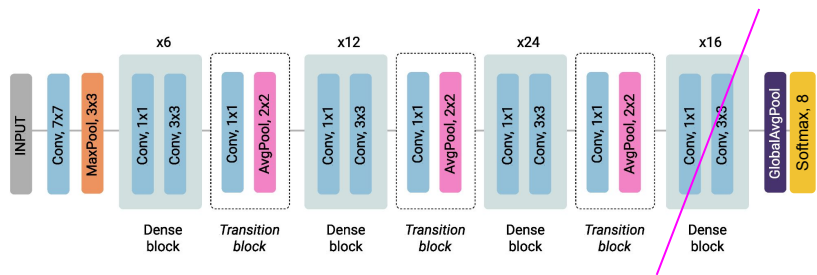→ Best aggregation: mean

- Accuracy = **0.77**



## Patch based deep features + BoVW

→ Num clusters too low —> too much generalization of the features
→ Num clusters too big → too much specificity of the features
→ Best number of clusters (codebook size) = 256

- Accuracy = **0.72**
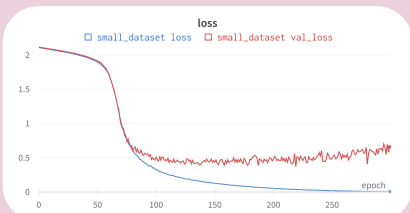
# Week 4: Fine Tuning DenseNet121[1]



x6    x12    x24    x16

Dense block — Transition block — Dense block — Transition block — Dense block — Transition block — Dense block

**29% less parameters** than the original model

## Making the network smaller

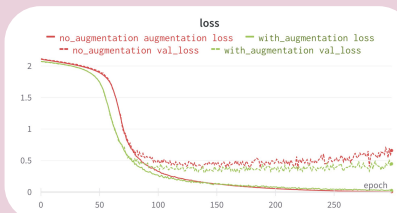| | Model | Epochs | Num parameters | Validation accuracy |
|---|---|---|---|---|
| MIT_split dataset | Original | 300 | 7M | 0.9542 |
| | **Removing 1 DB** | **300** | **5M** | **0.941** |
| | Removing 2 DB | 300 | 1.5M | 0.825 |

### Using MIT_small dataset

- **Overfitting**
- Accuracy drops to **0.845**



### Data augmentation

- Horizontal Flip = True
- Zoom Range = 20%
- Accuracy increases to **0.895**
- Validation loss does not increase



### Improve learning curve

- Early Stopper
- Reduce LR
- BatchNormalization and Dropout
  → Accuracy increases to **0.915**



### Hyperparameter Optimization
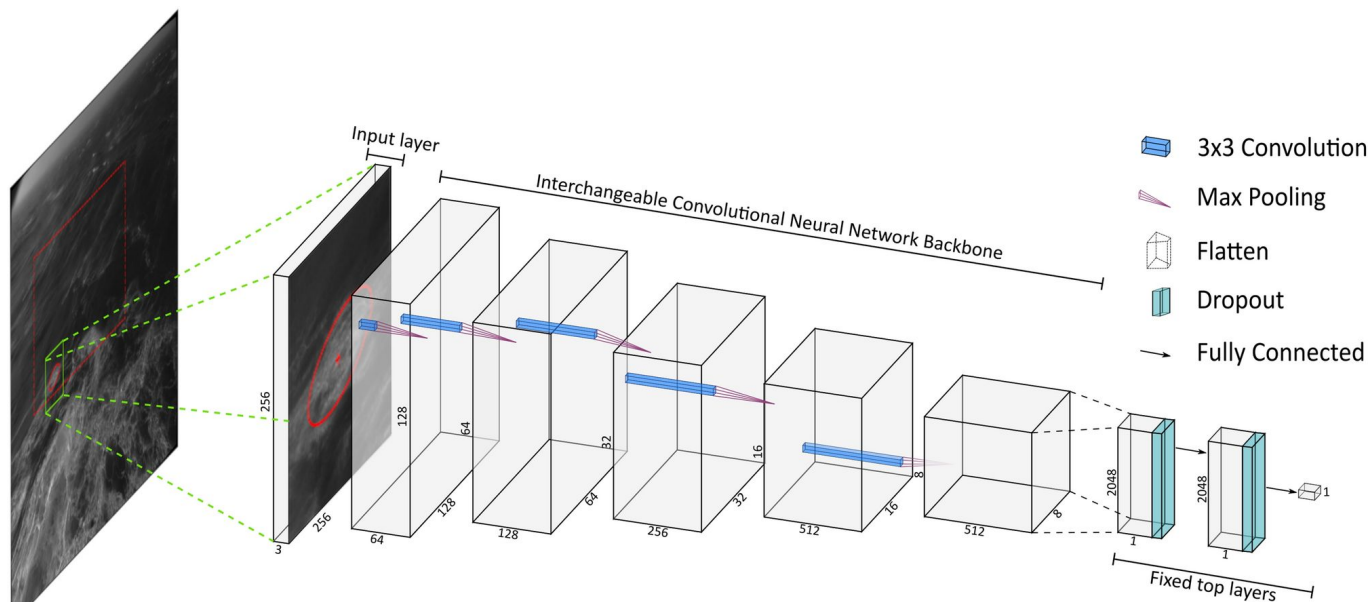
→ Accuracy increases to **0.9518**

Optimizer = Adamax    Wight decay = 0.3
LR = 0.0001    Batch_ size = 10
Dropout = 0.5    BatchNormalization



[1] Huang, G., Liu, Z., Weinberger, K. Q., & van der Maaten, L. (2017). Densely connected convolutional networks. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (pp. 4700-4708). https://arxiv.org/abs/1608.06993

Master in Computer Vision *Barcelona*
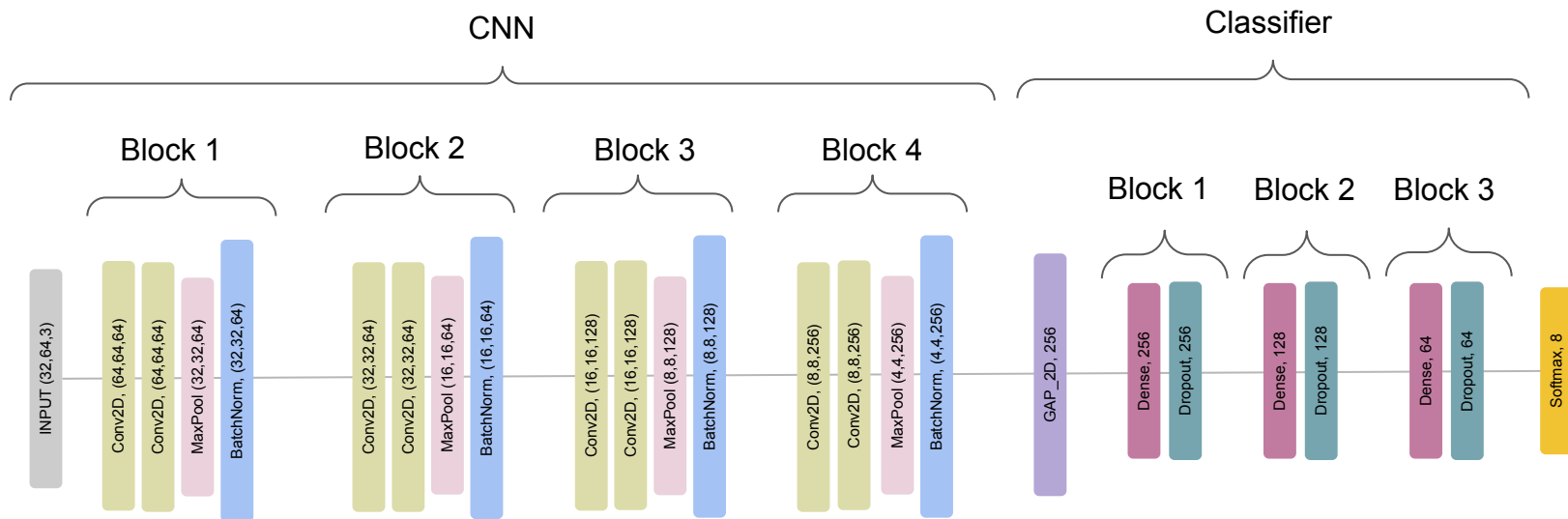
# Week 5

# Building a CNN from scratch



[1] Huang, G., Liu, Z., Weinberger, K. Q., & van der Maaten, L. (2017). Densely connected convolutional networks. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (pp. 4700-4708). https://arxiv.org/abs/1608.06993

# Manual search of the CNN network

- CNN
  - 4 blocks → Conv2D + Conv2D + MaxPool + BatchNorm
- Classifier
  - Global Average Pooling
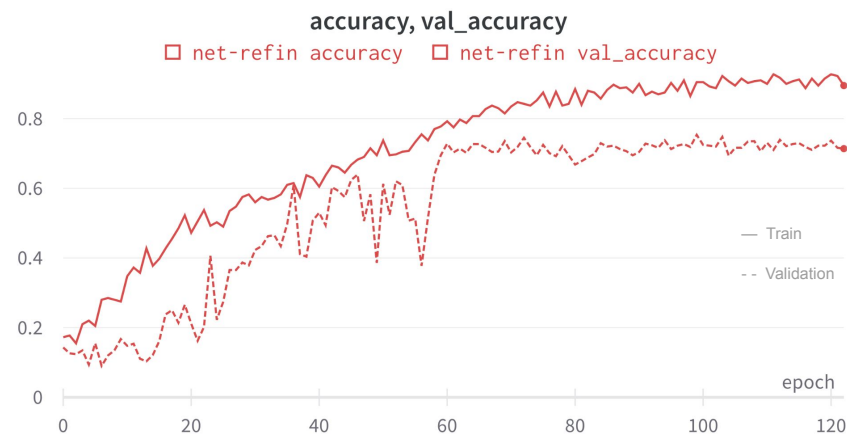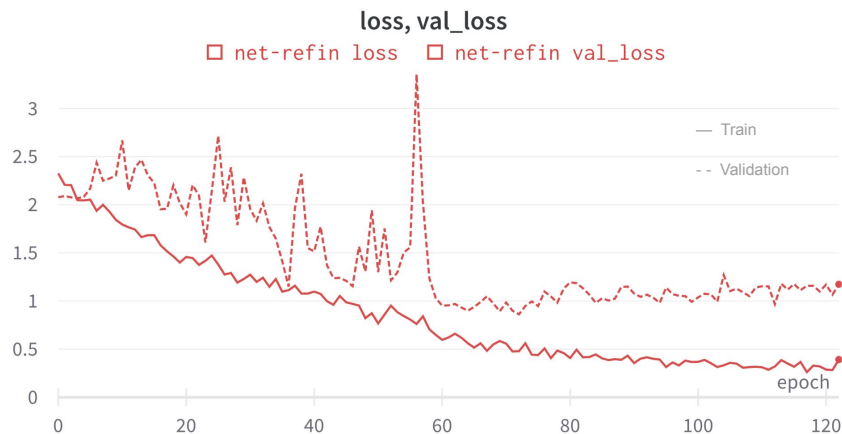  - 3 blocks → Dense + Dropout

# Network refinement

- Number of CNN blocks    → 2, 3, **4**, 5
- For each CNN block      → Add a second layer? **True**, False
- **Dimensionality first filter** → **32**, 64
- Number of dense layers  → 1, 2, **3**, 4

| Config parameter | Importance ⓘ ↓ | Correlation |
|------------------|----------------|-------------|
| num_blocks | | |
| num_denses | | |
| filters1 | | |
| second_layer | | |



loss, val_loss
☐ net-refin loss   ☐ net-refin val_loss



accuracy, val_accuracy
☐ net-refin accuracy   ☐ net-refin val_accuracy

10

# Initializer & activation function

- Initializer → GlorotUniform, GlorotNormal, HeUniform, HeNormal
- Activation → Relu, LeakyRelu, Mish
- Kernel size → 3x3 or 5x5
- Kernel regularizer → True, False

$f(x) = max(0, x)$

$f(x) = 1(x < 0)(\alpha x) + 1(x \geq 0)(x)$ where $\alpha$ = small constant

$f(x) = x \tanh(\ln(1 + e^x))$



loss, val_loss

□ after loss    □ before loss    □ after val_loss    □ before val_loss

— Train
-- Validation

epoch

| Config parameter | Importance ⓘ ↓ | Correlation |
|------------------|----------------|-------------|
| non_linearities.value_leaky_relu | | |
| initializer.value_glorot_normal | | |
| kernel_regularizer | | |

11

# Initializer & activation function

- Initializer → GlorotUniform, GlorotNormal, HeUniform, HeNormal
- Activation → Relu, LeakyRelu, Mish
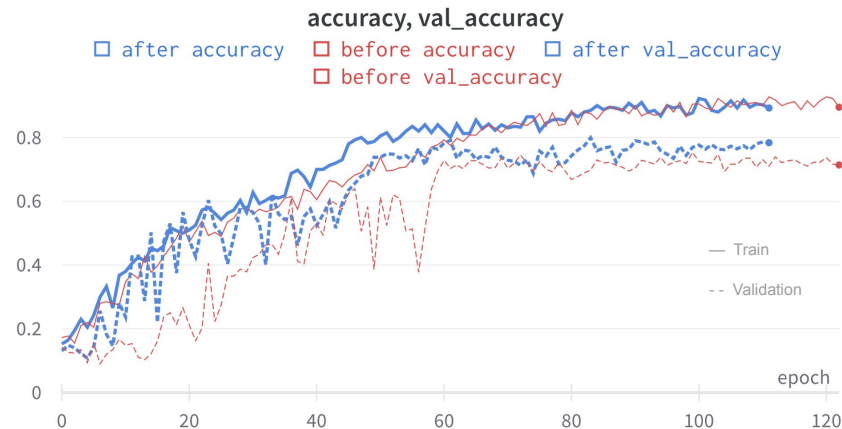- Kernel size → 3x3 or 5x5
- Kernel regularizer → True, False
- 



loss, val_loss

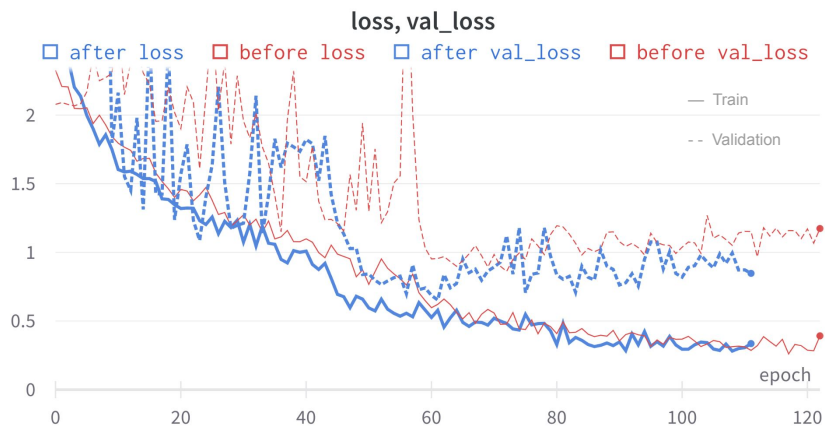☐ after loss   ☐ before loss   ☐ after val_loss   ☐ before val_loss

— Train
-- Validation



accuracy, val_accuracy

☐ after accuracy   ☐ before accuracy   ☐ after val_accuracy
☐ before val_accuracy

— Train
-- Validation

Master in Computer Vision *Barcelona*

# Hyperparameter optimization

| Params | Accuracy |
|--------|----------|
| 1.3M | 0.85 ↑ |

| Sweep | Values | Best | Sensitivity |
|-------|--------|------|-------------|
| LR | [0.01, 0.005, 0.001] | 0.001 | High |
| Batch Size | [8, 16, 32, 64] | 32 | Medium |
| Optimizer | [Adam, SGD] | Adam | Low |
| Dropout values | [0.3 , 0.4, 0.6] | 0.4 | Medium |
| Horizontal Flip | [True, False] | True | High |
| Rotation | [0, 15] | 0 | Medium |
| Width Shift | [0, 0.1] | 0.1 | Medium |
| Height Shift | [0, 0.1] | 0.1 | Medium |
| Shear Range | [0, 0.1] | 0 | Medium |
| Zoom Range | [0, 0.1] | 0 | High |



loss, val_loss
☐ before loss  ☐ after block loss  ☐ before val_loss
☐ after block val_loss
— Train
-- Validation



accuracy, val_accuracy
☐ before accuracy  ☐ after accuracy  ☐ before val_accuracy
☐ after val_accuracy
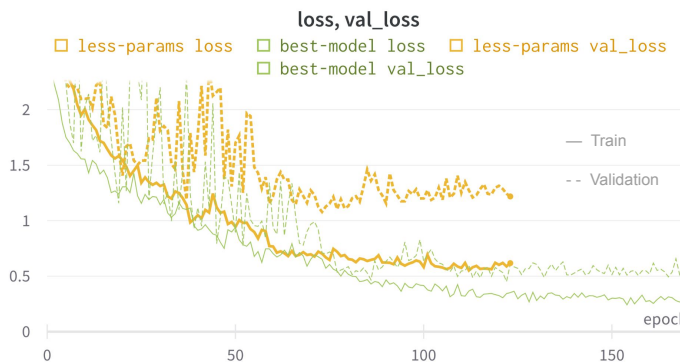— Train
-- Validation

UAB · UOC ⠿UPC upf. Master in Computer Vision Barcelona

# Parameter refinement

- Reduce the number of **filters** at each CNN block

| CNN blocks | Before   | After    |
|------------|----------|----------|
| 1st block  | 32, 32   | 32, 32   |
| 2nd block  | 64, 64   | **32, 32**   |
| 3th block  | 128, 128 | **64, 64**   |
| 4th block  | 256, 256 | **128, 128** |





accuracy, val_accuracy
- less-params accuracy
- less-params val_accuracy
- best-model accuracy
- best-model val_accuracy

loss, val_loss
- less-params loss
- best-model val_loss
- best-model loss
- less-params val_loss

Master in Computer Vision *Barcelona*

# Residual Connections

- Residual Connections between blocks
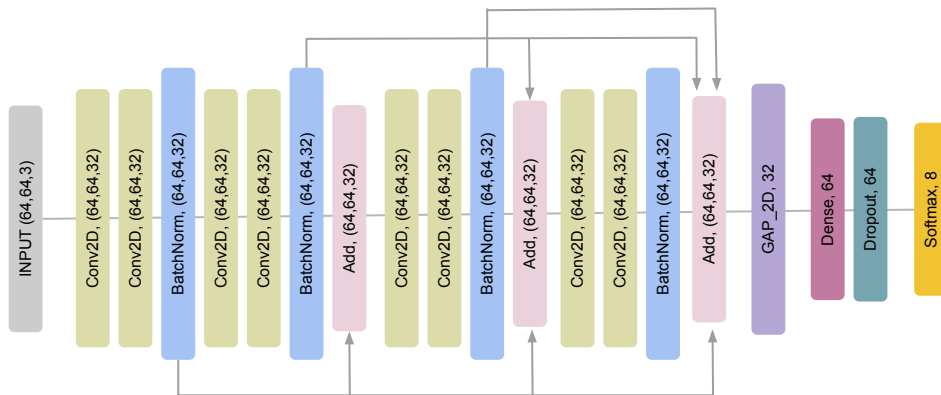- Reduce dimensionality of the filters



**Benefits**

- Converge more easily
- More Stability Training
- Better Generalization
- Easy implementation
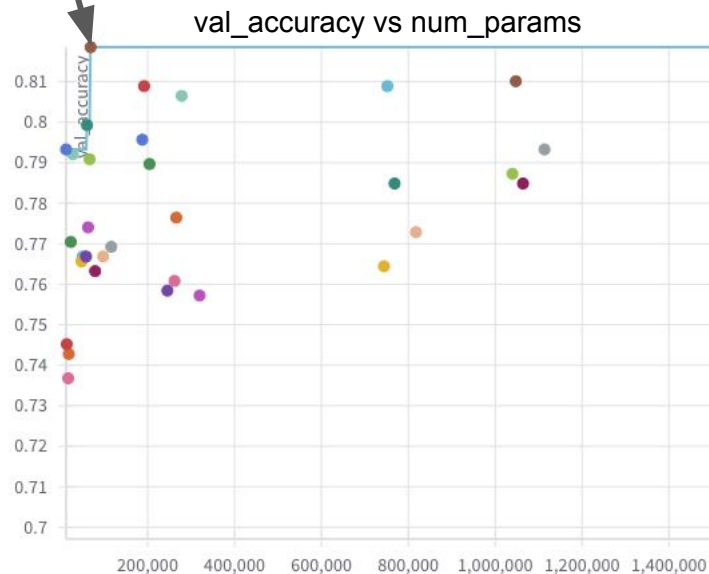- Improve accuracy

# Residual network search

- CNN blocks → 2, 3, **4**, 5
- Num filters → 16, **32**, 64, 128, 256
- Dense layers → 0, **1**, 2, 3

68k!!!



val_accuracy vs num_params

UAB  UOC  UPC  upf.  Master in Computer Vision  *Barcelona*

# Residual network search

- CNN blocks → 2, 3, **4**, 5
- Num filters → 16, **32**, 64, 128, 256
- Dense layers → 0, **1**, 2, 3



loss, val_loss

— Train
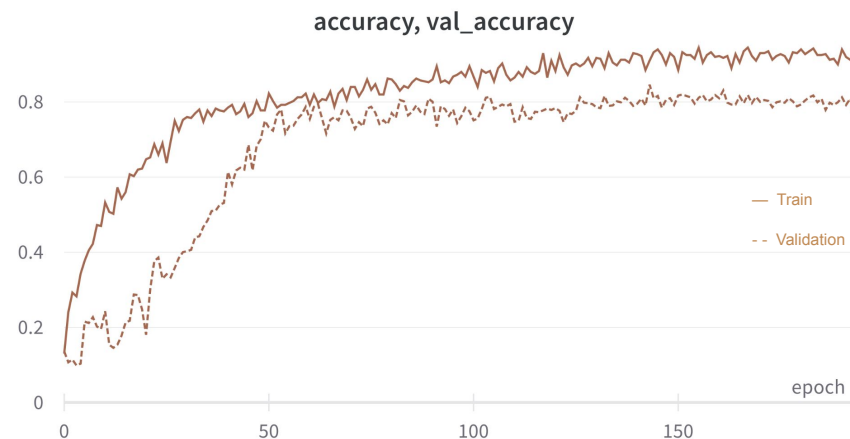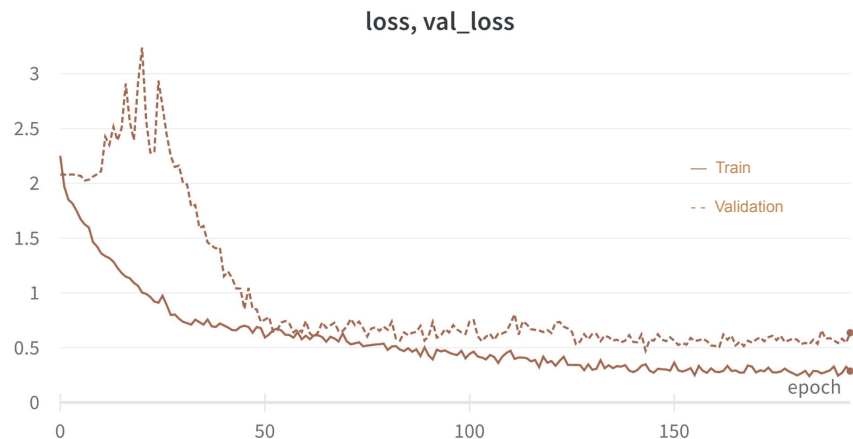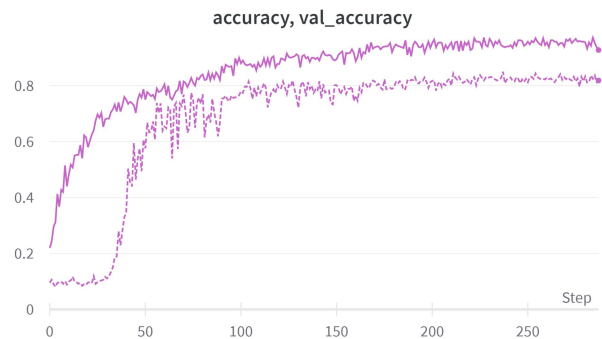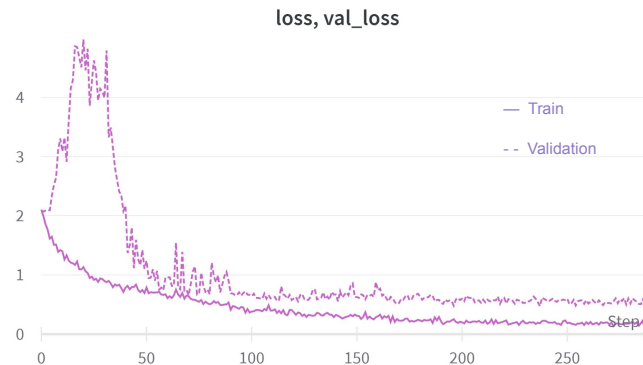-- Validation

epoch



accuracy, val_accuracy

— Train
-- Validation

epoch

Master in Computer Vision *Barcelona*

# Data augmentation refinement

- Horizontal Flip → True
- Rotation [0, 15] → **0**
- Width Shift [0, 0.1] → **0.1**
- Height Shift [0, 0.1] → **0**
- Shear Range [0, 0.1] → 0
- Zoom Range [0, 0.1] → 0

loss, val_loss

— Train

- - Validation

accuracy, val_accuracy

| Config parameter | Importance ⓘ ↓ | Correlation |
|------------------|----------------|-------------|
| height_shift | | |
| shear_range | | |
| rotation | | |
| zoom_range | | |
| width_shift | | |

Master in Computer Vision  *Barcelona*

# Recap week 5



### loss, val_loss

― ResNet loss ― SeqNet-low loss ― SeqNet loss -- ResNet val_loss
-- SeqNet-low val_loss -- SeqNet val_loss

### accuracy, val_accuracy

― ResNet accuracy ― SeqNet-low accuracy ― SeqNet accuracy
-- ResNet val_accuracy -- SeqNet-low val_accuracy
-- SeqNet val_accuracy

| Params | Accuracy |
|--------|----------|
| 1.3M   | 0.85     |

| Params | Accuracy |
|--------|----------|
| 380k   | 0.80     |

| Params | Accuracy |
|--------|----------|
| 68k    | 0.85     |

# Recap week 5

- Good results with BOVW and spatial pyramids (**0.95** accuracy)

- MLP → too simple for our problem in all the cases (**0.77** accuracy)

- Reduced data → challenging and requires the use of **data augmentation** and other techniques to achieve the desired results

- Fine-tuning the DenseNet121 → best results (**0.96** accuracy) and we were able to reduce **30%** of the network parameters to **5M**

- Building a network from scratch is time consuming and difficult to optimize

- Residual connections → help the network to converge easily and generalize.

  → **0.85** accuracy with **68k** parameters

- Our results are limited due to the lack of data

UAB · UOC ·UPC *upf.* Master in Computer Vision *Barcelona*