

Deep Learning

Medical Image Classification with
ResNets, DenseNets, EfficientNet and ViT

Team: **JAX**

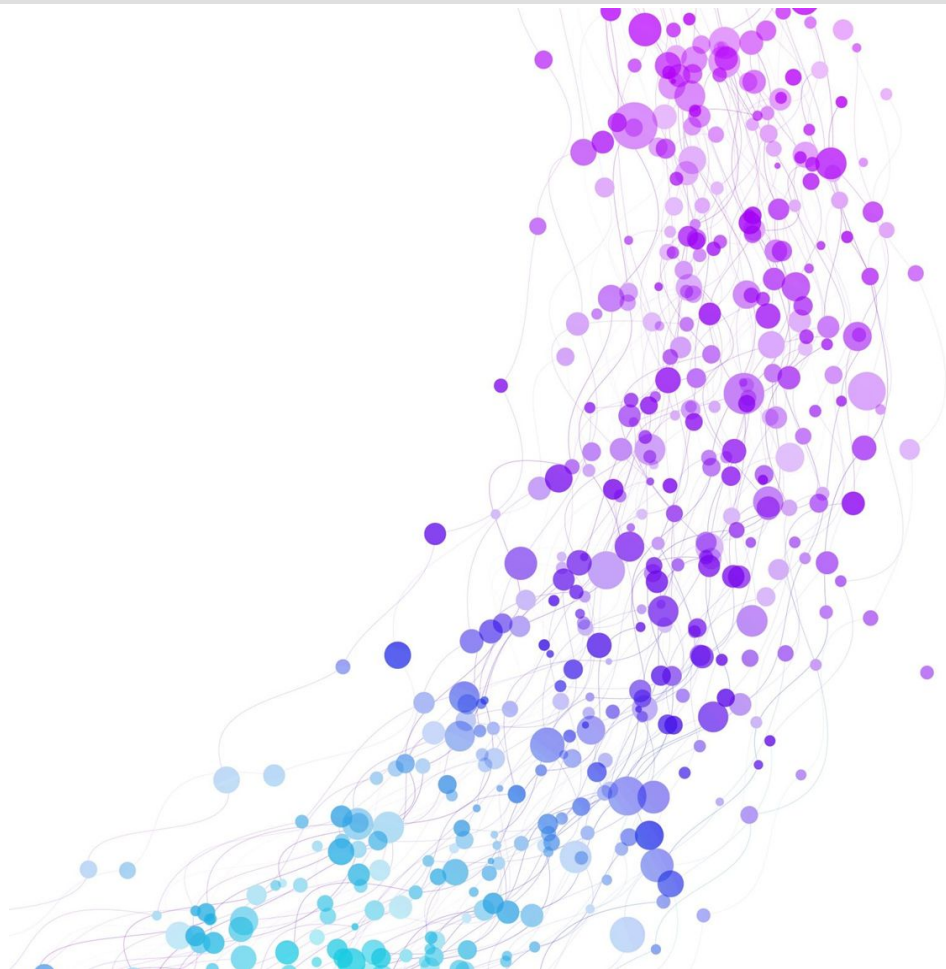
[GIT Repository](#)

Skoltech



Road Map

1. Problem statement
2. Datasets description
3. Method
4. Model
 - 4.1. ResNet
 - 4.2. DenseNet
 - 4.3. EfficientNet
 - 4.4. ViT
5. Evaluation Result
6. Conclusion



Medical image classification plays an increasingly important role in healthcare, especially in diagnosing, treatment planning and disease monitoring.

Problem: The lack of large publicly available datasets with annotations means it is still very difficult, if not impossible, to achieve clinically relevant computer-aided detection and diagnosis (CAD).

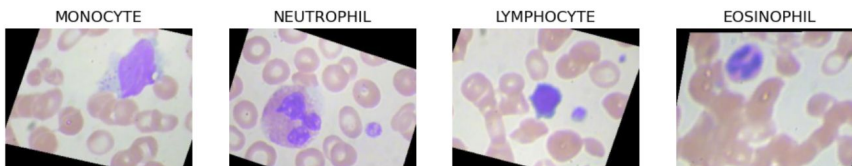
Deep learning models have been shown to be very effective at image classification, and they are increasingly being used to **improve medical tasks**.

⇒ discuss the implementation of four different deep learning models
for medical image classification: ResNet, DenseNet, Efficient Net, and ViT

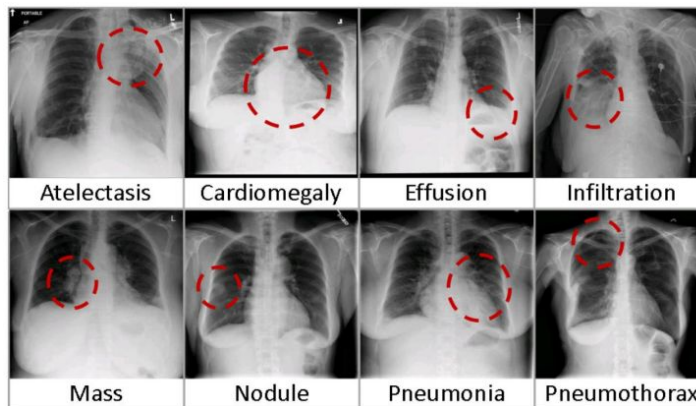
Datasets Description

Name	Blood Cell Images	Random Sample of NIH Chest X-ray Dataset
Description	12,500 augmented images of blood cells	5,606 chest X-rays from random patients
Number labels	4	15

**Blood Cell Images
Dataset**



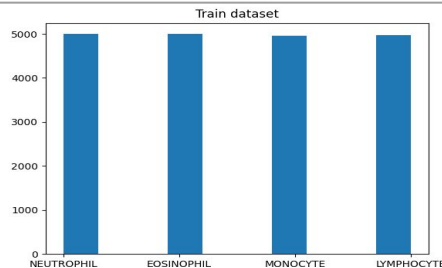
**NIH Chest
X-ray Dataset**



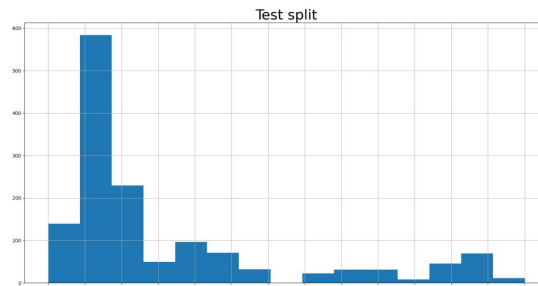
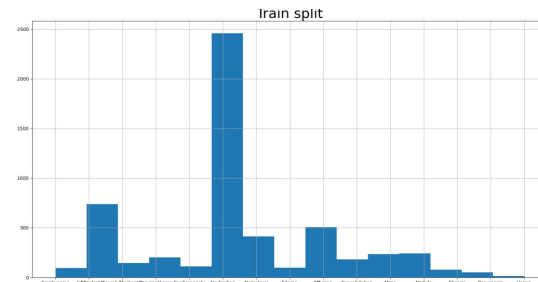
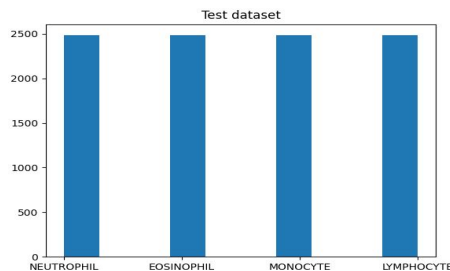
Datasets Description

Name	Blood Cell Images	Random Sample of NIH Chest X-ray Dataset
Description	12,500 augmented images of blood cells	5606 chest X-rays from random patients
Number labels	4	15

Train split:

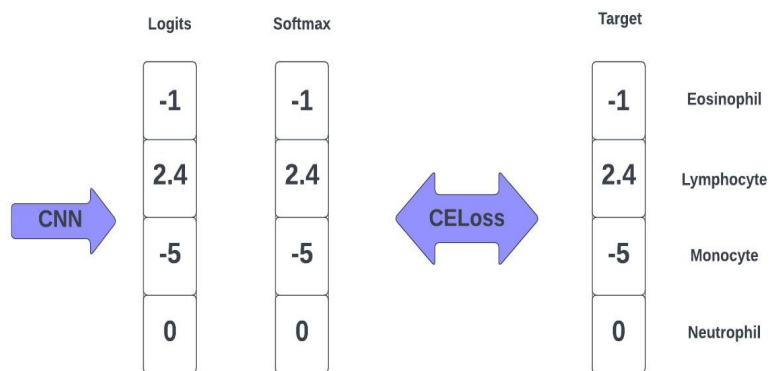


Test split:



Multi-class classification

Blood Cell Images

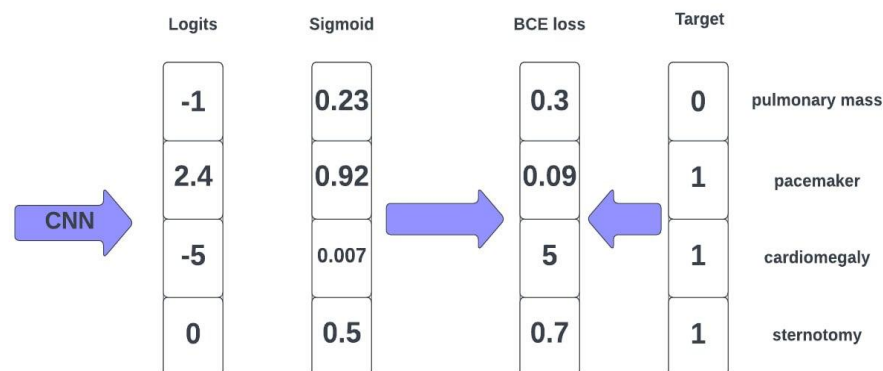


$$f(s)_i = \frac{e^{s_i}}{\sum_j^C e^{s_j}} \quad CE = - \sum_i^C t_i \log(f(s)_i)$$

each instance belongs to only one class

Multi-label classification

Random Sample of NIH Chest X-ray



$$\text{BCELoss} = -\frac{1}{N} \sum_{i=1}^N [y_i \cdot \log(p(y_i)) + (1 - y_i) \cdot \log(1 - p(y_i))]$$

each instance can belong to multiple labels simultaneously

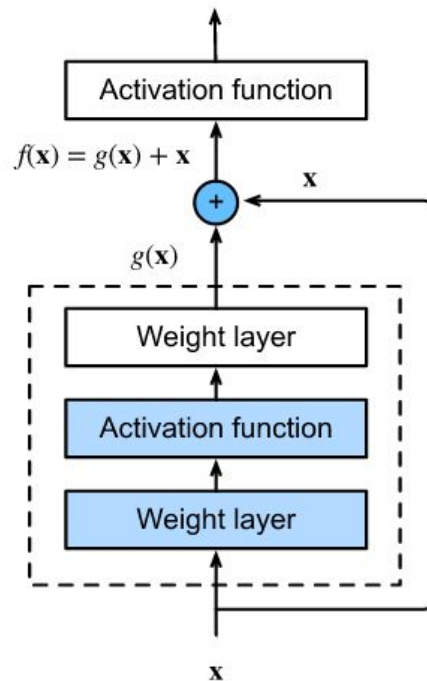
ResNets work by stacking together a series of residual blocks.

Each **residual block** consists of

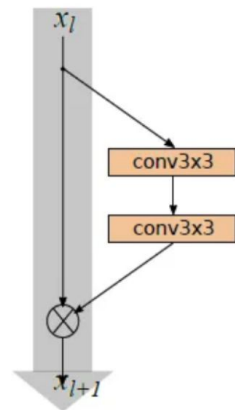
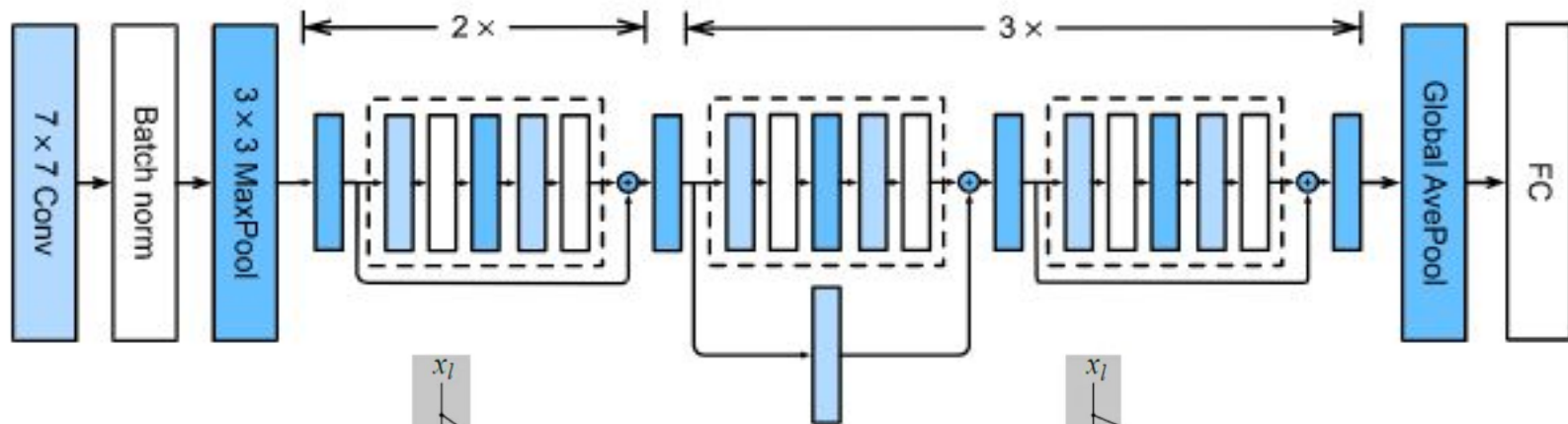
- **convolutional layers**
- followed by a **skip connection**

The skip connection allows the output of the first convolutional layer to be added to the output of the final convolutional layer.

⇒ prevent the *vanishing gradient* problem
of the very deep neural networks

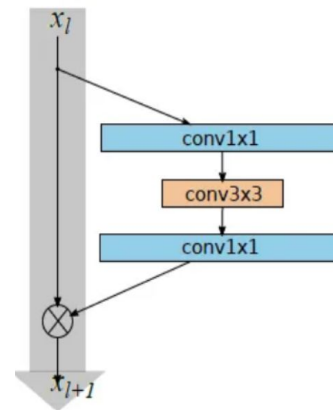


ResNet: Architecture



(a) basic

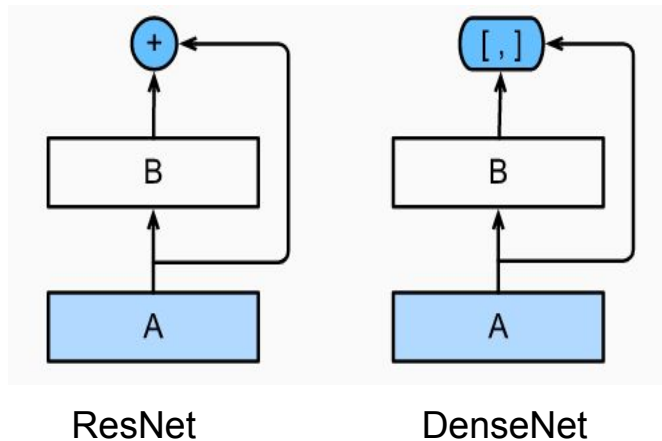
ResNet18, ResNet34



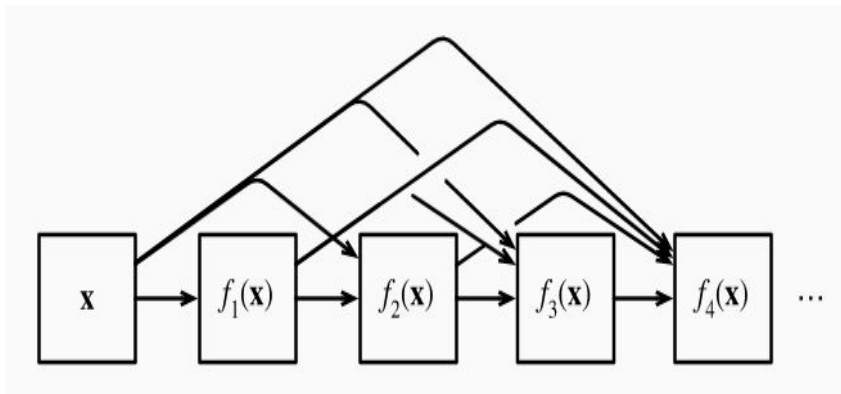
(b) bottleneck

ResNet50, ResNet101, ResNet152

DenseNet - Densely Connected Networks

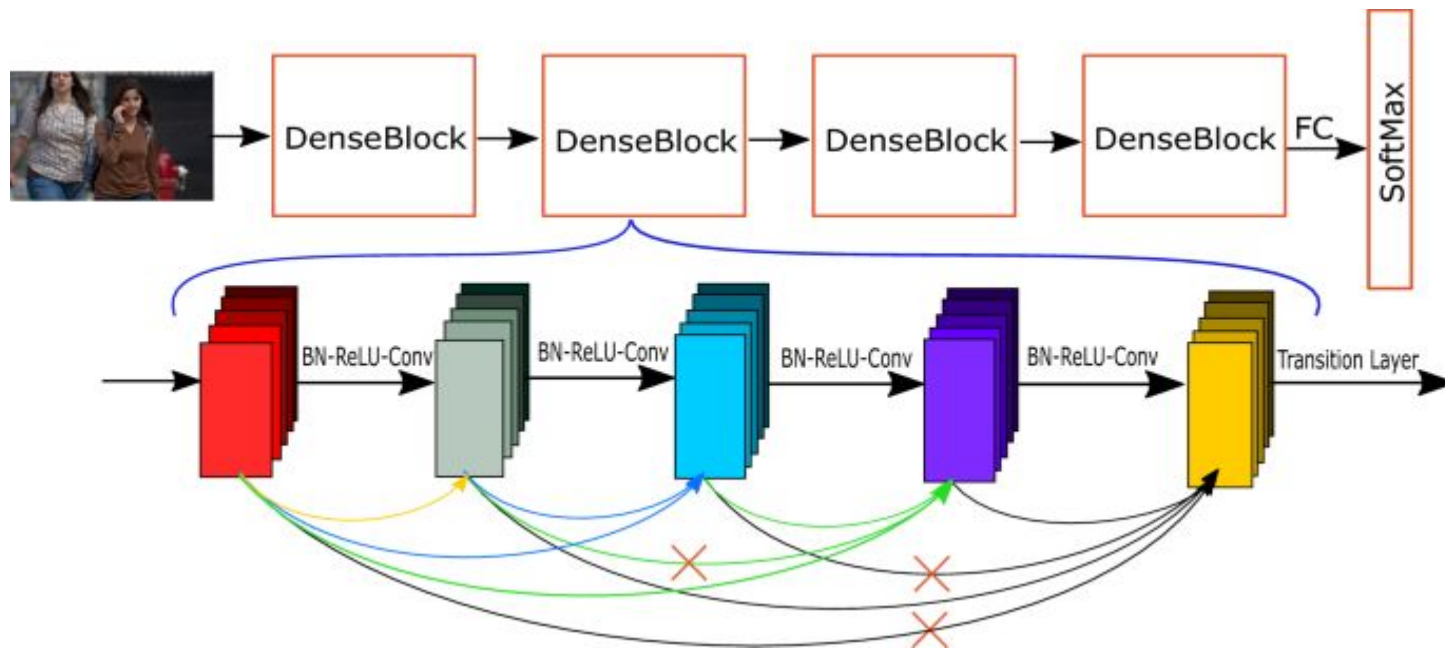


The main difference between **ResNet** and **DenseNet** in **cross-layer connections** is the use of **addition** vs **concatenation**



DenseNets are characterized by their use of dense connections

DenseNet - Densely Connected Networks

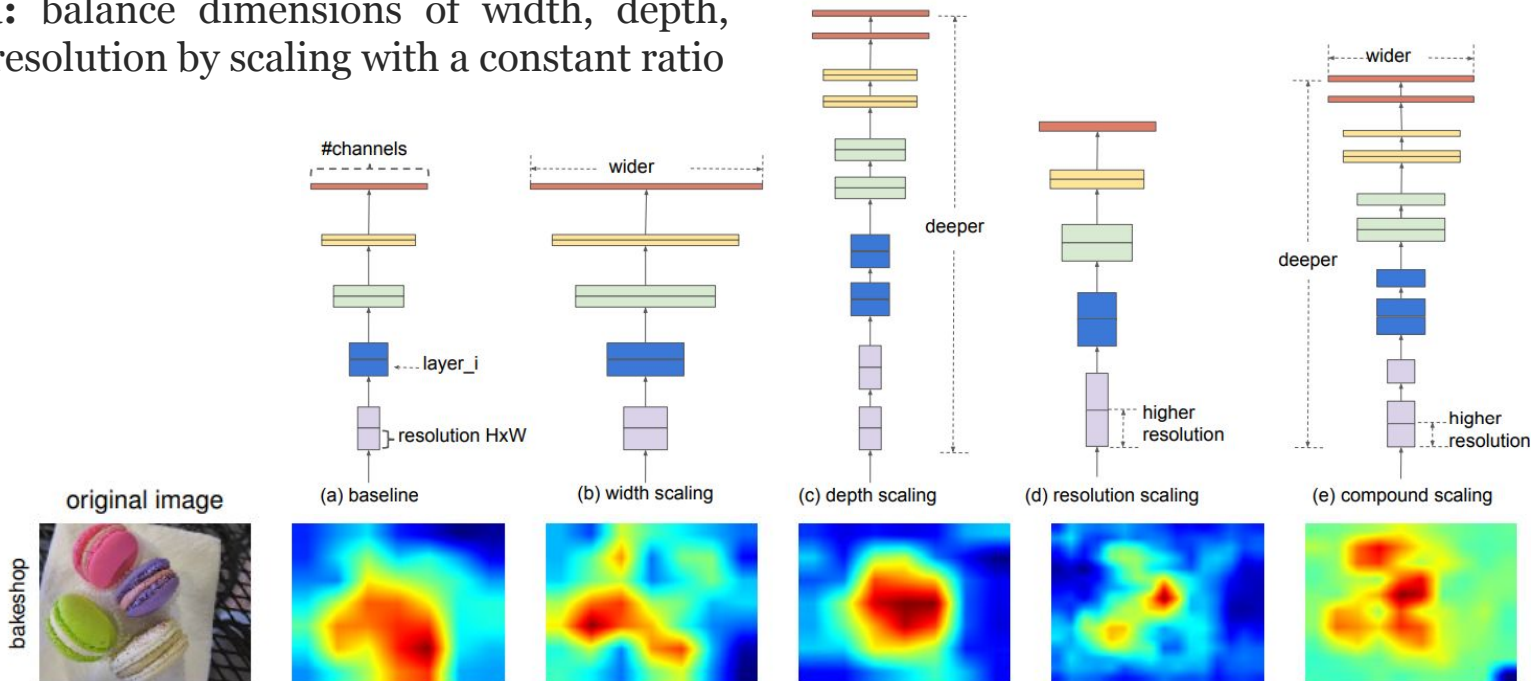


The main components of a DenseNet are

- dense blocks \Rightarrow allow information to flow more easily between
- transition layers different layers of the network

Traditional approaches for a ConvNet scale: width, depth and resolution scaling

Idea: balance dimensions of width, depth, and resolution by scaling with a constant ratio



ConvNet - a set of blocks of layers with input shape as (H, W, C).

$$\max_{d,w,r} \text{Accuracy}(\mathcal{N}(d, w, r))$$

$$s.t. \quad \mathcal{N}(d, w, r) = \bigodot_{i=1 \dots s} \hat{\mathcal{F}}_i^{d \cdot \hat{L}_i} (X_{\langle r \cdot \hat{H}_i, r \cdot \hat{W}_i, w \cdot \hat{C}_i \rangle})$$

$$\text{Memory}(\mathcal{N}) \leq \text{target_memory}$$

$$\text{FLOPS}(\mathcal{N}) \leq \text{target_flops}$$



Compound scaling method

$$\text{depth: } d = \alpha^\phi$$

$$\text{width: } w = \beta^\phi$$

$$\text{resolution: } r = \gamma^\phi$$

$$s.t. \quad \alpha \cdot \beta^2 \cdot \gamma^2 \approx 2$$

$$\alpha \geq 1, \beta \geq 1, \gamma \geq 1$$



Architecture of EfficientNet B0

Stage i	Operator $\hat{\mathcal{F}}_i$	Resolution $\hat{H}_i \times \hat{W}_i$	#Channels \hat{C}_i	#Layers \hat{L}_i
1	Conv3x3	224×224	32	1
2	MBConv1, k3x3	112×112	16	1
3	MBConv6, k3x3	112×112	24	2
4	MBConv6, k5x5	56×56	40	2
5	MBConv6, k3x3	28×28	80	3
6	MBConv6, k5x5	14×14	112	3
7	MBConv6, k5x5	14×14	192	4
8	MBConv6, k3x3	7×7	320	1
9	Conv1x1 & Pooling & FC	7×7	1280	1

Note: use Stochastic Depth with 0,06

EfficientNet family of ConvNet

Transformer hyperparameters:

Patch size: (8, 8),

Model dim: 128,

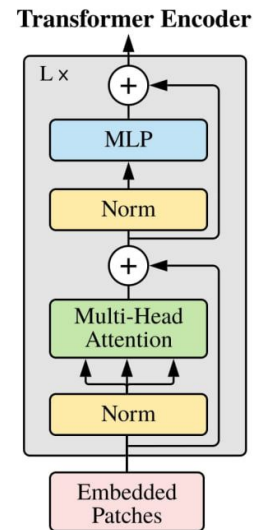
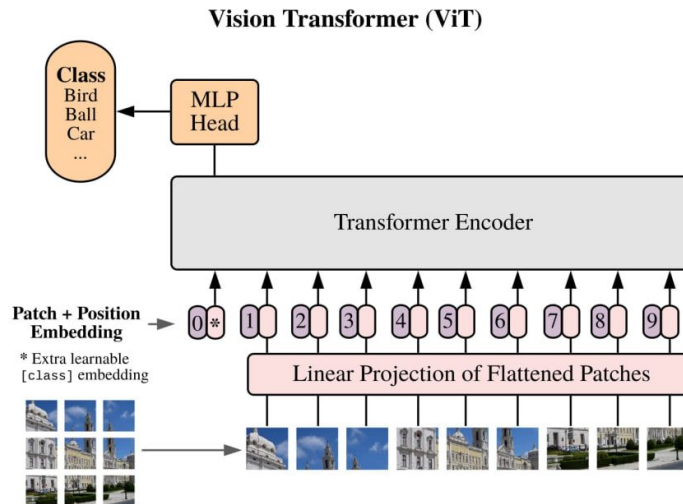
Depth: 6,

N heads: 16,

Mlp ratio: 4,

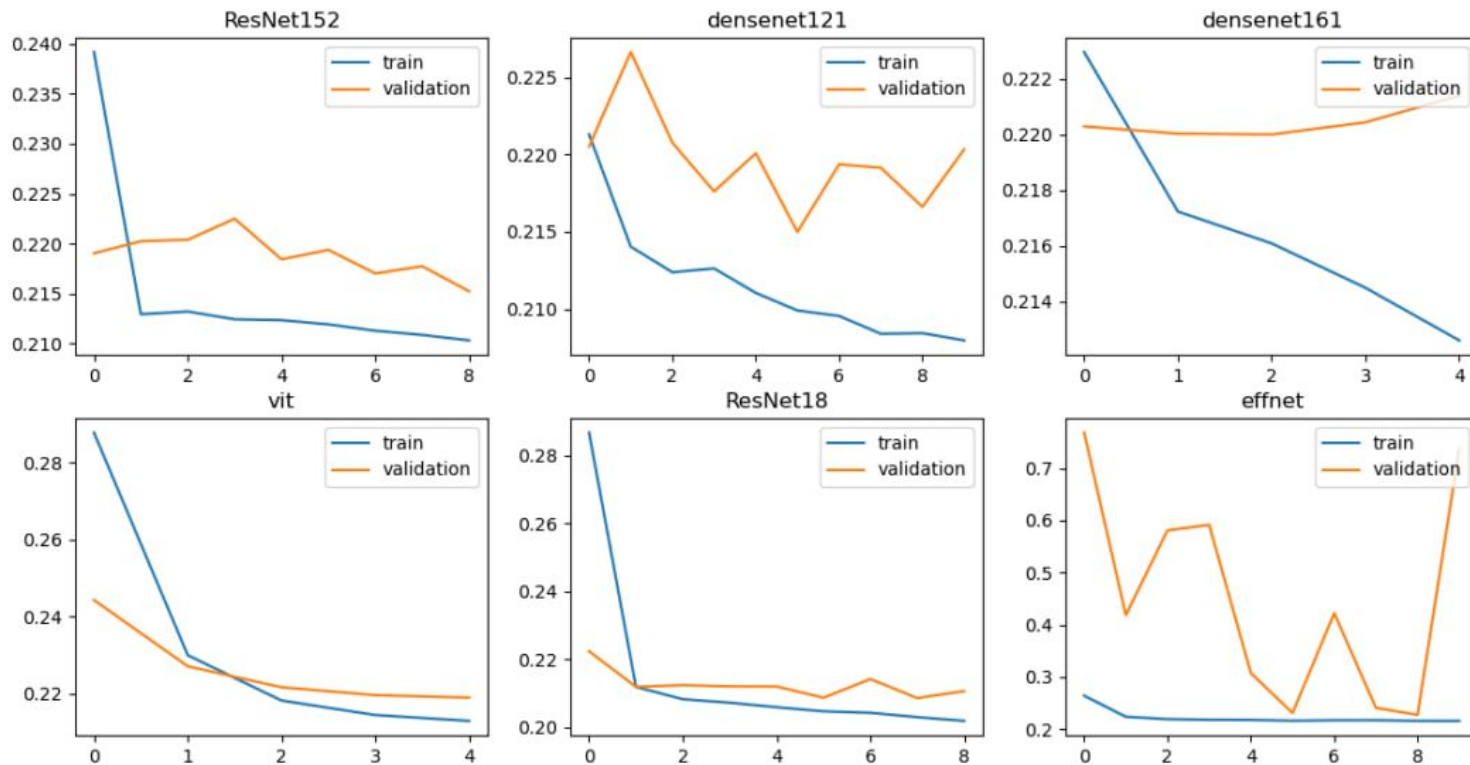
Dropout ratio : 0.3,

Attention dropout ratio: 0.3



Evaluation: Multi-class classification Losses

Multi-label Classification Losses

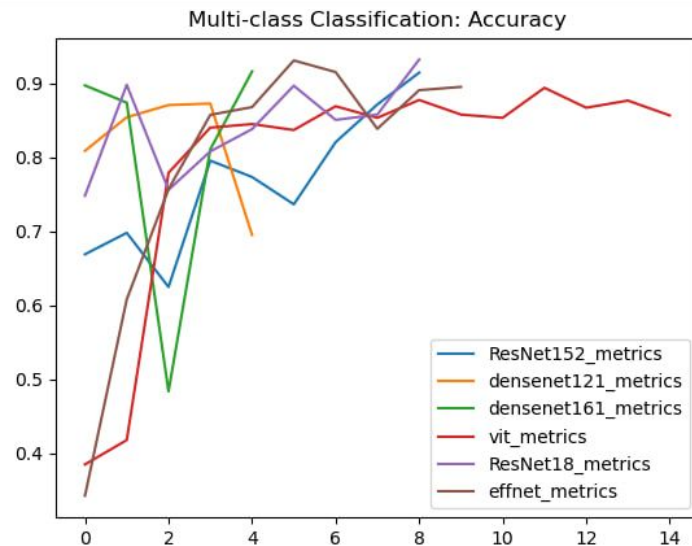


Evaluation: Multi-class classification

Multi-class classification metrics

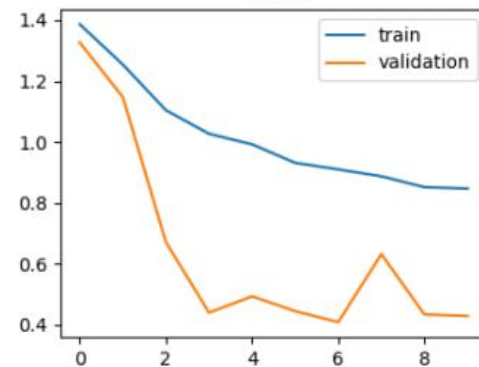
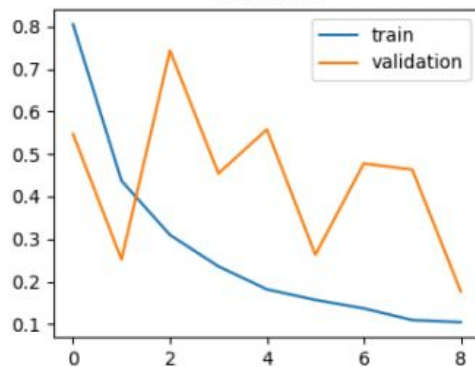
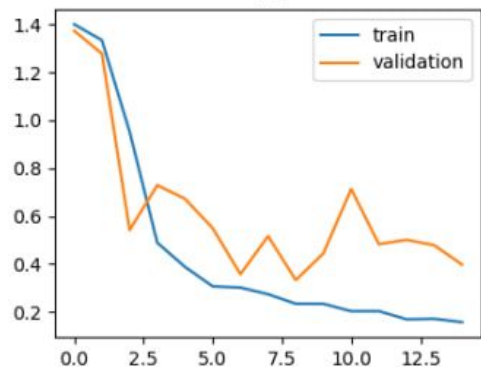
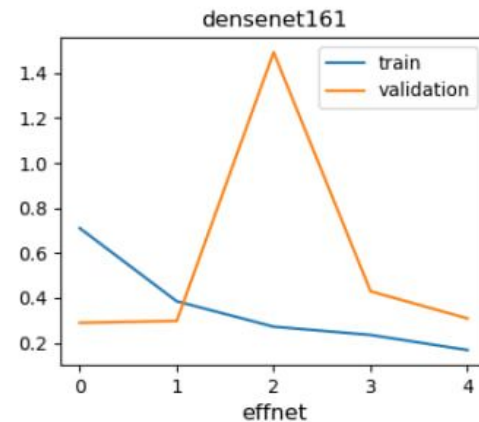
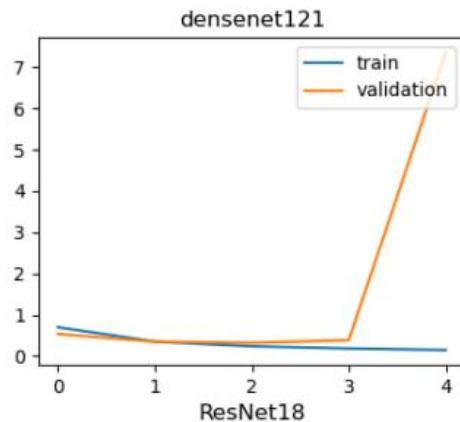
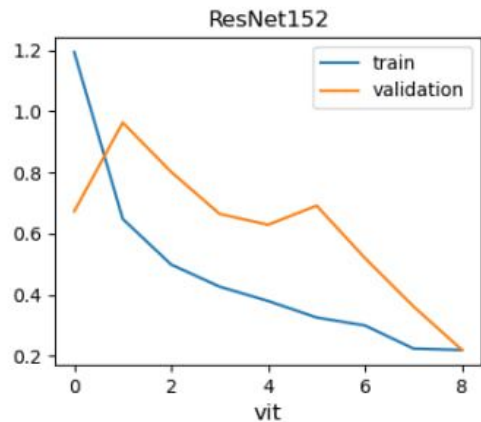
model	accuracy	precision (macro)	precision (micro)	recall (macro)	recall (micro)	F1 (macro)	F1 (micro)
resnet152	0.914757	0.913539	0.914757	0.914980	0.914757	0.913517	0.914757
resnet18	0.932449	0.936440	0.932449	0.932577	0.932449	0.932841	0.932449
densenet121	0.872939	0.887754	0.872939	0.872932	0.872939	0.874157	0.872939
densenet161	0.916365	0.917669	0.916365	0.916586	0.916365	0.915667	0.916365
effnet	0.931242	0.933620	0.931242	0.931398	0.931242	0.931379	0.931242
vit	0.894250	0.912612	0.894250	0.894179	0.894250	0.895673	0.894250

- ResNet18 and Efficient_b0 are the best model with 0.93 accuracy and F1 score
- DenseNet121 and DenseNet161 are unstable with accuracy and F1 are 0.87 and 0.91 for both nets.
- ViT showed SOTA results (accuracy: 0.89, F1: 0.89) with stable learning



Evaluation: Multi-label classification Losses

Multi-class Classification: Losses

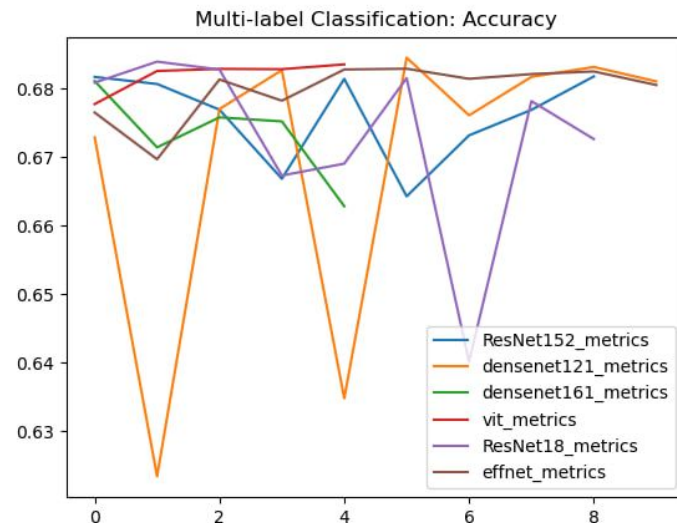


Evaluation: Multi-label classification

Multi-label classification metrics

model	Exact Match Ratio	Hamming loss	Recall	Precision	F1
ResNet152	0.501338	0.501338	0.502230	0.501933	0.502052
ResNet18	0.473684	0.473684	0.476360	0.474874	0.475320
densenet121	0.513827	0.513827	0.513827	0.513827	0.513827
densenet161	0.520963	0.520963	0.520963	0.520963	0.520963
effnet	0.520963	0.520963	0.520963	0.520963	0.520963
vit	0.520963	0.520963	0.520963	0.520963	0.520963

- DenseNets and Efficient_b0 showed satisfied results with 0.52 accuracy and F1 score
- DenseNets and ResNets were unstable, while EfficientNet_b0 showed stable learning
- ViT with metrics, as other algorithms and showed stable learning.



Number of parameters in models

model	total parameters, ths.	trainable parameters, ths.
resnet_18	11178	11178
resnet_152	58152	58152
densenet_121	6957	6957
densenet_161	26480	26480
effnet	5288	5288
transformer	1368	1368

EfficientNet is the most efficient model for both multi-class and multi-label classification problems. It performs well in all metrics and shows stable learning. ViT performs as other algorithms in terms of metrics. However, there are significantly less params to train. That makes this architecture to be alternative for ConvNets in the future.

Team: JAX

Xavier Aramayo Carrasco

33 %

- Find Dataset for Multi-class classification
- Implementation of DenseNet

Julia Kudryavtseva

33 %

- Implementation of EfficientNet, ViT
- Evaluation

Vo Ngoc Bich Uyen

33 %

- Find Dataset for Multi-label classification
- Implementation of ResNet

Thank You!

Skoltech

