



NEW ARCHITECTURES AND TRICKS TO TRAIN THEM

Deep Learning for Computer Vision

Arthur Douillard

Neural Architecture Search

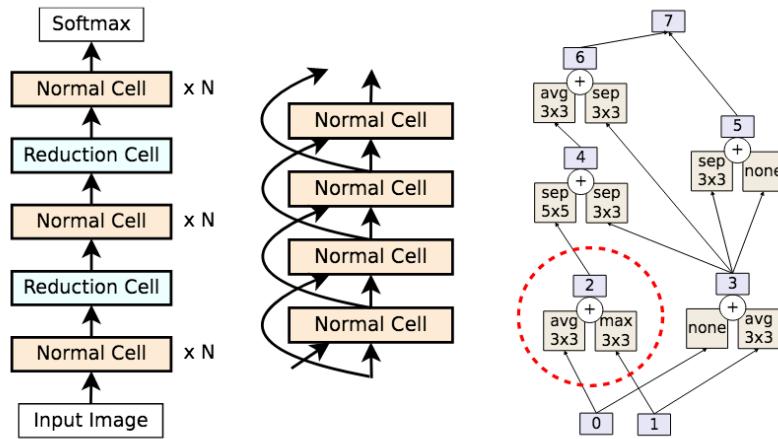
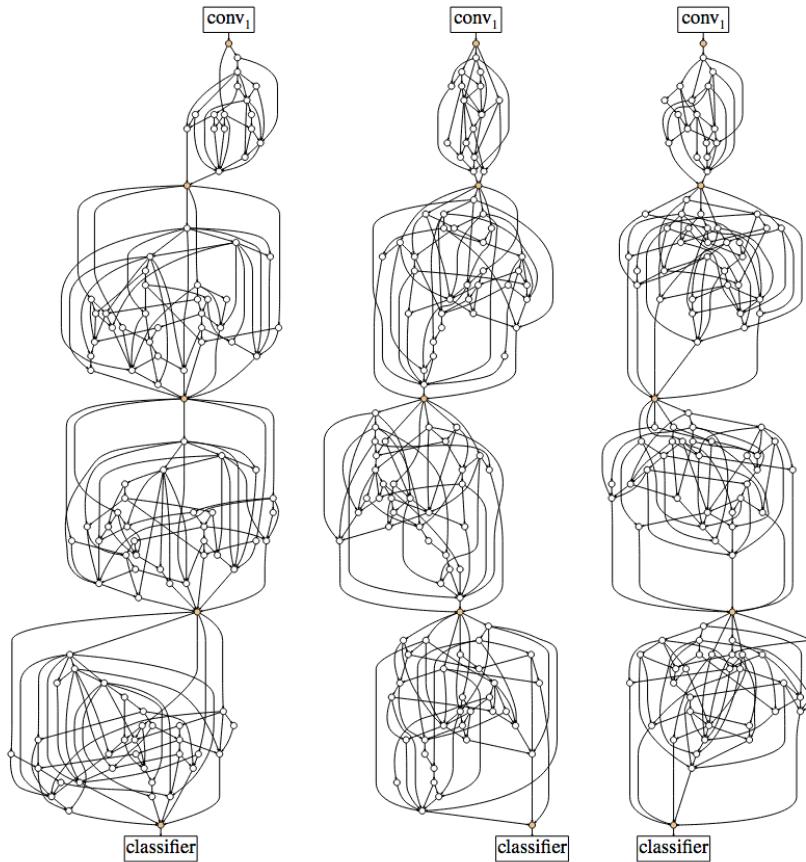


Figure 1: NASNet Search Space [54]. LEFT: the full outer structure (omitting skip inputs for clarity). MIDDLE: detailed view with the skip inputs. RIGHT: cell example. Dotted line demarcates a pairwise combination.

- Given a fixed architecture (left & middle), learn to find the optimal cell (right)
- Learning is done here with an **evolutionary algorithm** that needs to retrain & check model accuracy FOR EACH new mutation!
- NAS are usually very computation intensive, and thus it's mostly big private lab that works on it. With Quoc Le's team at Google Brain the main one.



- Or try randomly wired neural network based on minimal set of rules
- But note that some approaches use reinforcement learning

EfficientNet



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

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

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

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

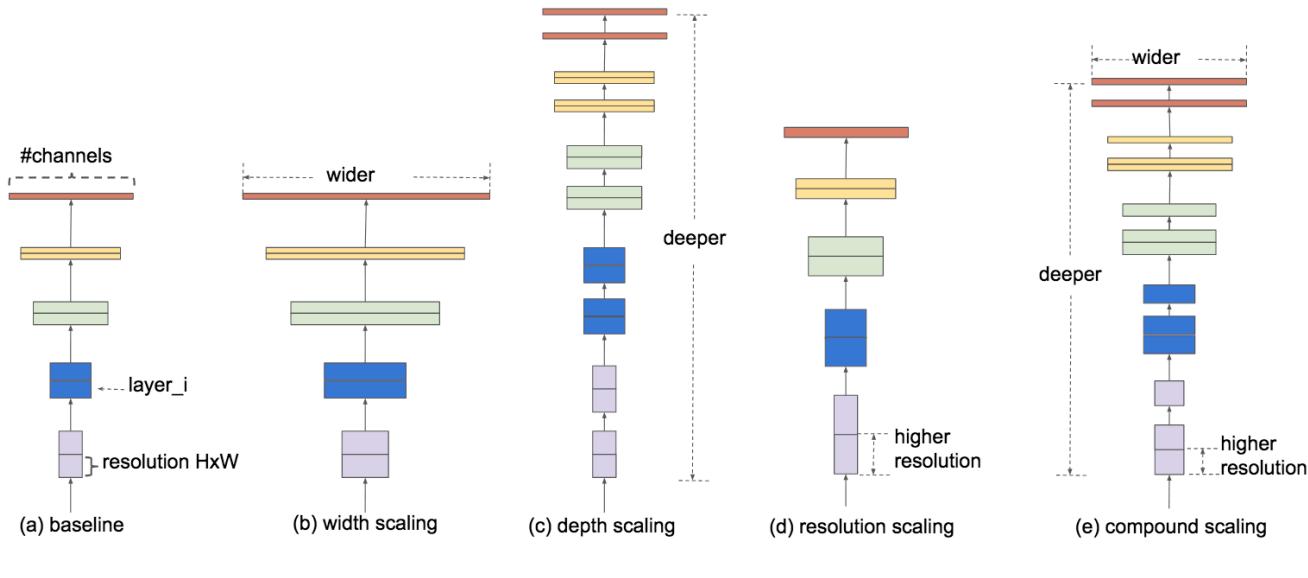
depth: $d = \alpha^\phi$

width: $w = \beta^\phi$

resolution: $r = \gamma^\phi$

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

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



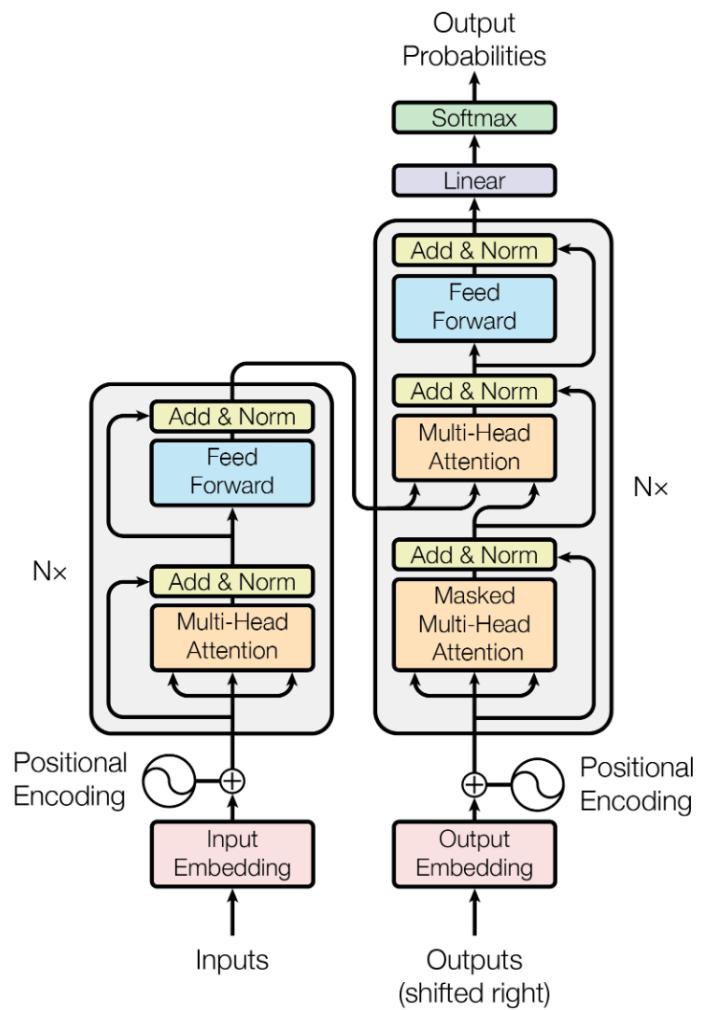
- EfficientNet, one of the best ConvNet as of 2021, was made with NAS
- Based on a **compound scaling** rule they drastically reduce the space to grid search

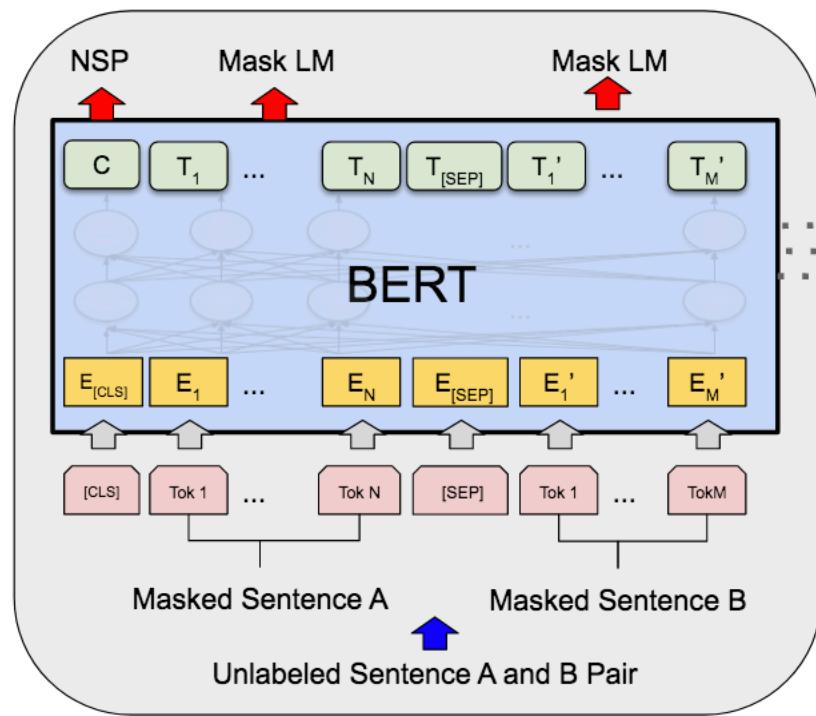
Transformers



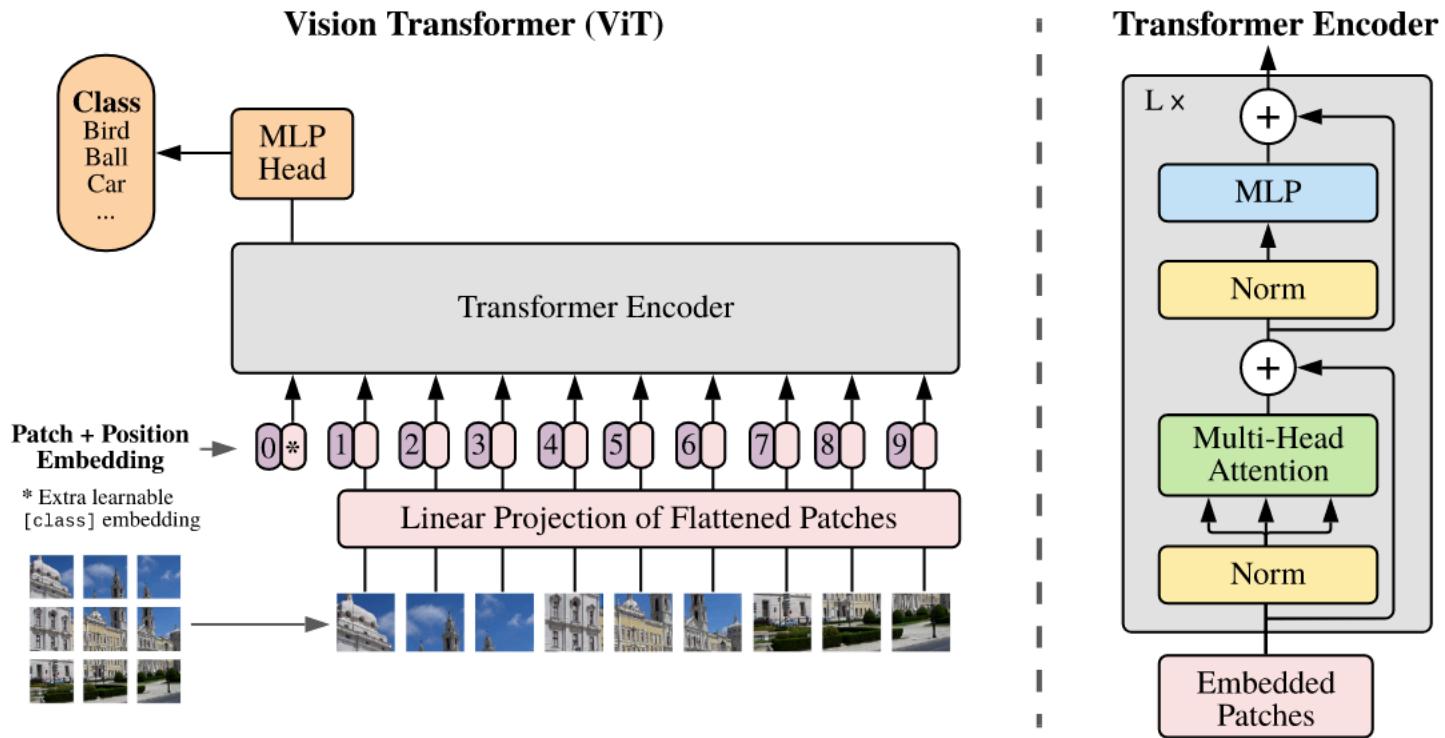
Attention is all you need

- At first designed for machine translation
- Does not rely on Conv1d or RNN
- Main block are FC layers and the famous **Multi-Head Attention**
- Made of a encoder (left) and decoder (right)





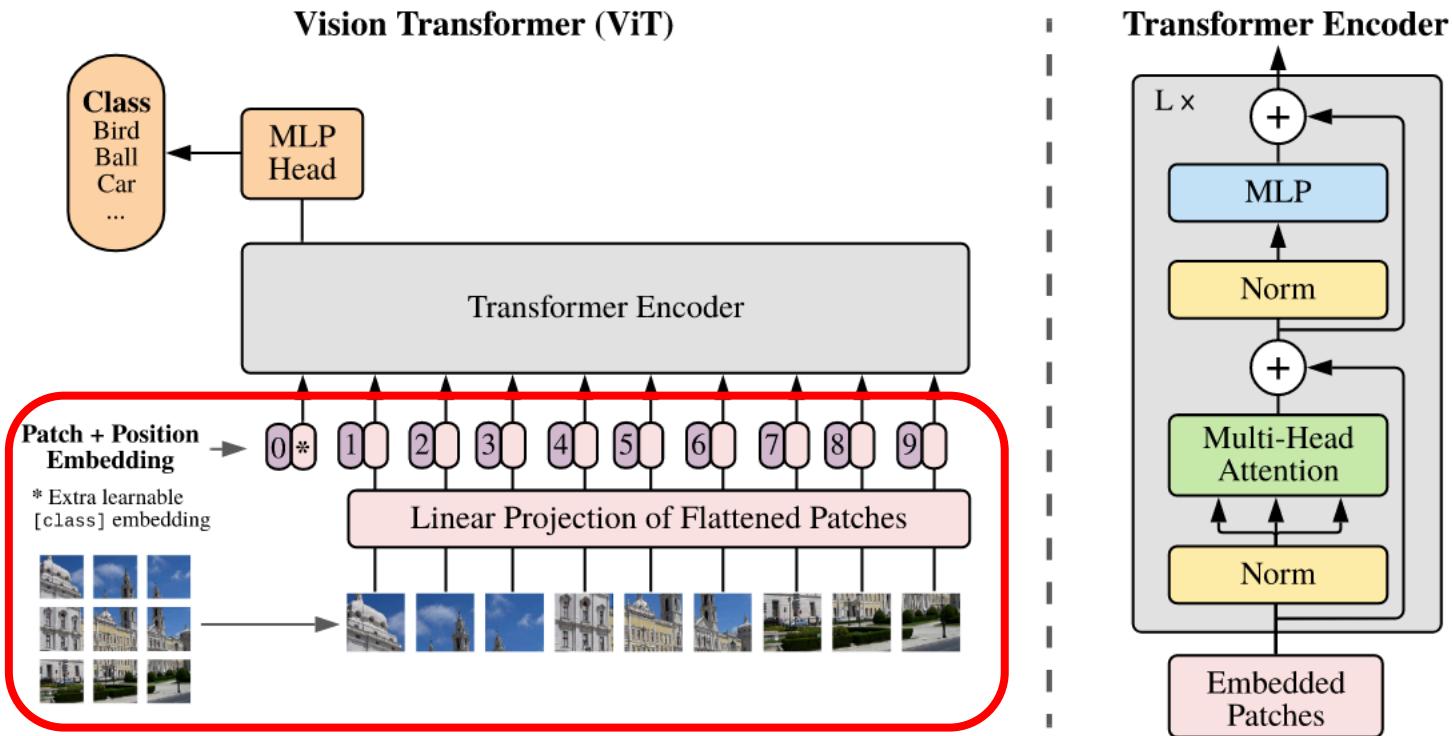
- Modification of the Transformer
- No more encoder/decoder, only many blocks
- Introduction of a special token [CLS] that is learned
- Was a big revolution in the NLP field



- First *full* application of transformers with a BERT-like architecture to vision



Patch, Position, and class token



- Use a 1×1 convolution with a kernel size equal to the patch size to generate the **tokens**
 - Total size is thus $(batch\ size, number\ of\ tokens, embedding\ dimension)$
- Add an **extra token [class]** that is a learned vector of size $(embedding\ dimension)$
- Add to all tokens a learned **positional embeddings**



Self-Attention

1. Apply three different linear transformations, to create the **Query**, the **Key**, and the **Value**

$$Q = XW_q$$

$$K = XW_k$$

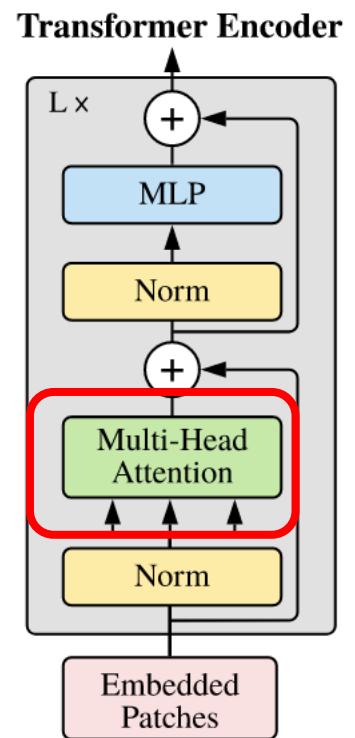
$$V = XW_v$$

2. Compute the attention matrix that measures the inter-tokens similarity

$$A = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right)$$

3. Ponderate the **Value** matrix with the attention matrix

$$Z = AV$$



Multi-heads Self-Attention



$$Q_1 = XW_{q_1} \text{ and } Q_2 = XW_{q_2}$$

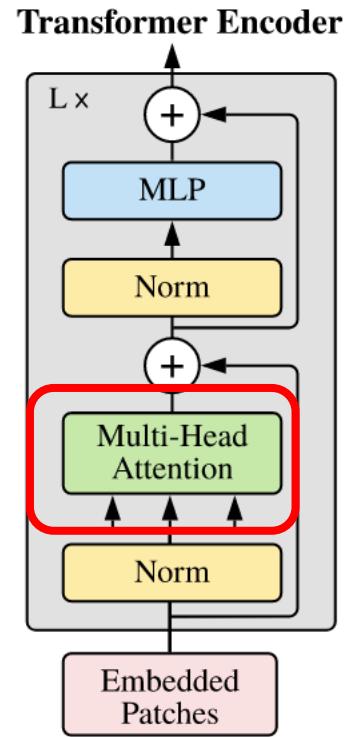
$$K_1 = XW_{k_1} \text{ and } K_2 = XW_{k_2}$$

$$V_1 = XW_{v_1} \text{ and } V_2 = XW_{v_2}$$

$$Z_1 = \text{softmax}\left(\frac{Q_1 K_1^T}{\sqrt{d}}\right) \text{ and } Z_2 = \text{softmax}\left(\frac{Q_2 K_2^T}{\sqrt{d}}\right)$$

$$Z' = [Z_1 Z_2] \in \mathbb{R}^{T \times 2d}$$

$$Z = Z' W_o$$

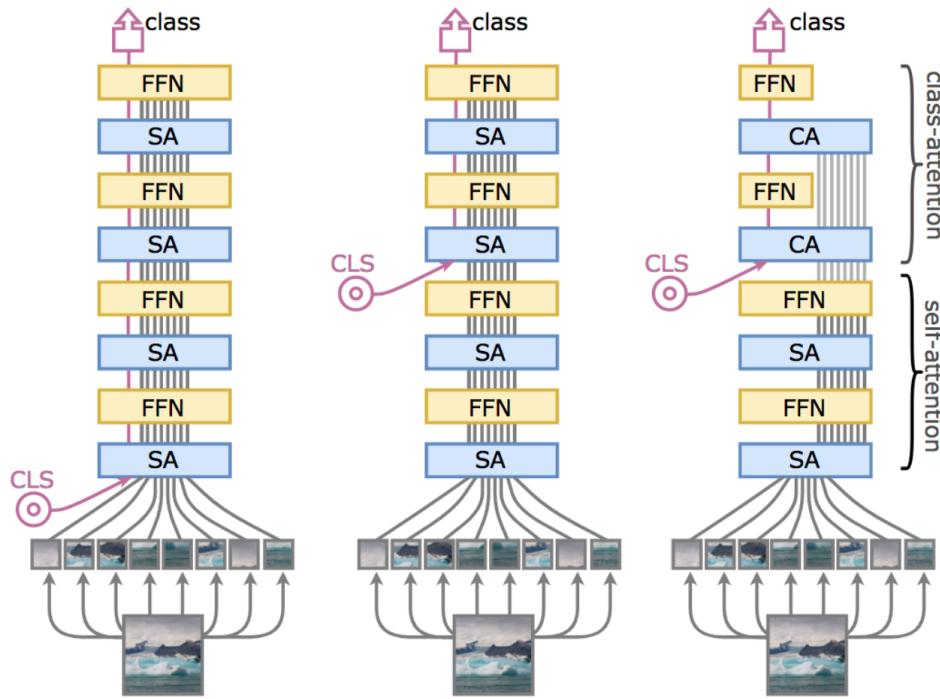


In practice, the self-attention is done multiple times in parallel, with different **heads**.



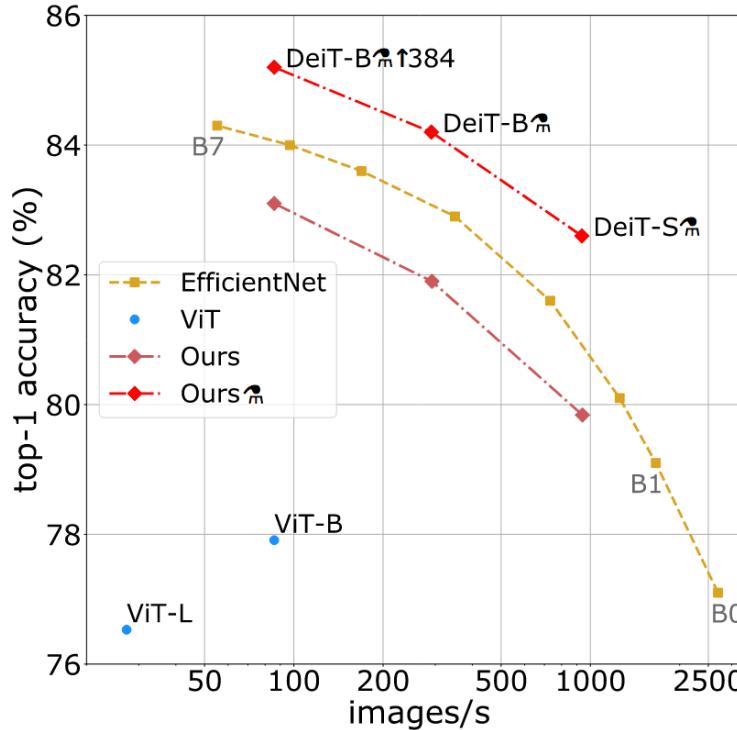
Ablation on ↓			Pre-training		Fine-tuning		Rand-Augment				top-1 accuracy					
							AutoAug	Mixup	CutMix	Erasing	Stoch. Depth	Repeated Aug.	Dropout	Exp. Moving Avg.		
none: DeiT-B	adamw	adamw	✓	✗	✓	✓	✓	✓	✓	✓	✓	✓	✗	✗	81.8 ± 0.2	83.1 ± 0.1
optimizer	SGD	adamw	✓	✗	✓	✓	✓	✓	✓	✓	✓	✓	✗	✗	74.5	77.3
	adamw	SGD	✓	✗	✓	✓	✓	✓	✓	✓	✓	✓	✗	✗	81.8	83.1
data augmentation	adamw	adamw	✗	✗	✓	✓	✓	✓	✓	✓	✓	✓	✗	✗	79.6	80.4
	adamw	adamw	✗	✓	✓	✓	✓	✓	✓	✓	✓	✓	✗	✗	81.2	81.9
	adamw	adamw	✓	✗	✗	✓	✓	✓	✓	✓	✓	✓	✗	✗	78.7	79.8
	adamw	adamw	✓	✗	✓	✓	✗	✓	✓	✓	✓	✓	✗	✗	80.0	80.6
	adamw	adamw	✓	✗	✗	✗	✗	✓	✓	✓	✓	✓	✗	✗	75.8	76.7
regularization	adamw	adamw	✓	✗	✓	✓	✓	✓	✓	✗	✓	✓	✗	✗	4.3*	0.1
	adamw	adamw	✓	✗	✓	✓	✓	✓	✓	✓	✗	✓	✗	✗	3.4*	0.1
	adamw	adamw	✓	✗	✓	✓	✓	✓	✓	✓	✓	✗	✗	✗	76.5	77.4
	adamw	adamw	✓	✗	✓	✓	✓	✓	✓	✓	✓	✓	✓	✗	81.3	83.1
	adamw	adamw	✓	✗	✓	✓	✓	✓	✓	✓	✓	✓	✗	✓	81.9	83.1

- Transformers are hard to train because they lack **inductive bias**, and thus needs way more data than a ConvNet
- DeiT partially close this gap by using tons of **data augmentations** and **regularizations**



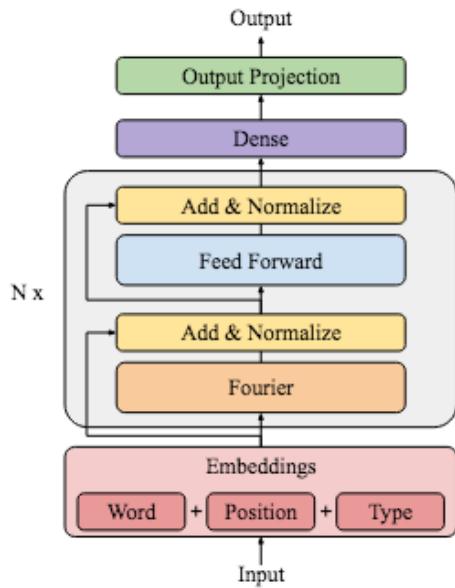
- Inserting the class tokens in later blocks may prove beneficial

Future of Transformers



- More generally transformers for vision seems to be leading the SotA in computer vision
- And they manage to be super fast
- Although, they need more data and regularizations than your common ConvNet
- And using loss curvature regularization (like SAM) seems to help a lot
- Hard to follow literature, hundred of vision transformers in Jan—Aug 2021! Time will tell which ones are really useful

Fourier Transformer



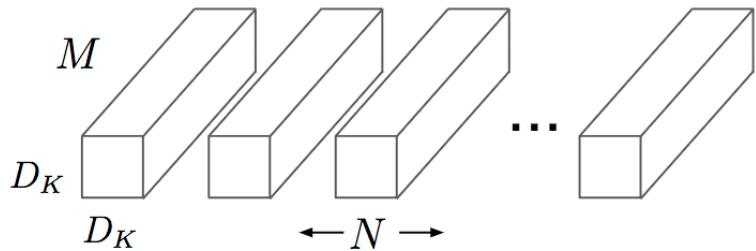
Fourier transform that can be expressed as a matrix multiplication with this constant matrix:

$$W_{nk} = \left(e^{-\frac{2\pi i}{N} nk} / \sqrt{N} \right)$$

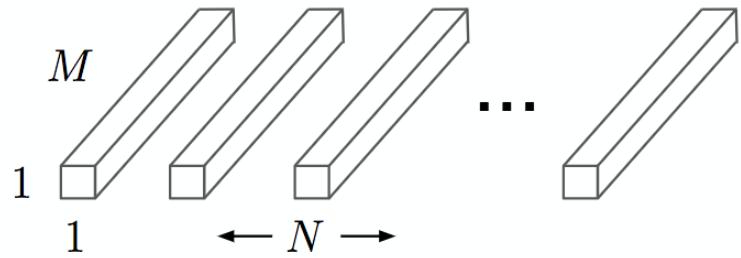
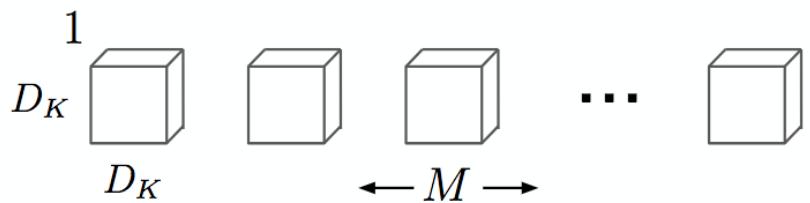
- Not as good as Self-Attention, but still impressive results
- Means that "*Attention is NOT all you need*", but rather a way to combine inputs
 - Likewise convolutions combine pixels through the increasing receptive field

MLP Comeback!

Separable Convolutions



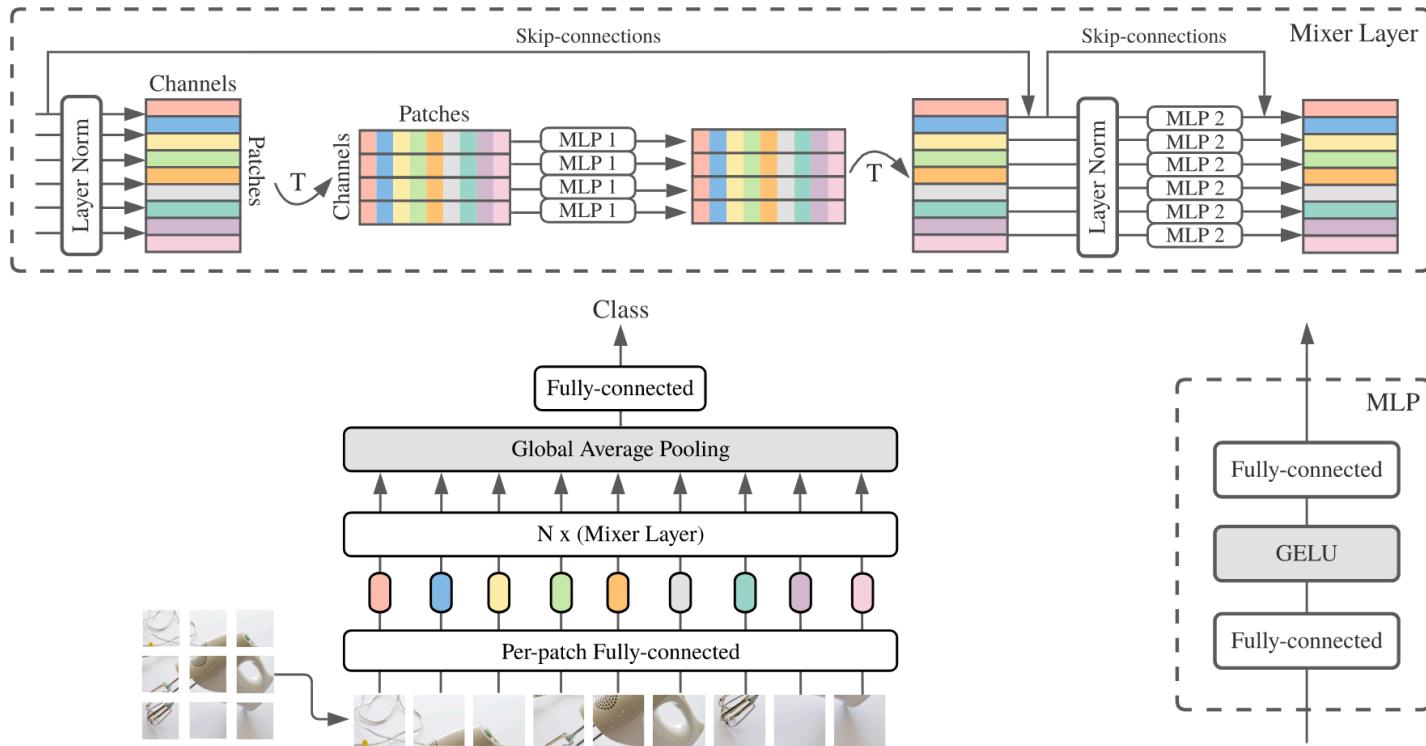
(a) Standard Convolution Filters



(c) 1×1 Convolutional Filters called Pointwise Convolution in the context of Depthwise Separable Convolution

- **Depthwise convolutions:** doesn't mix input channels
- **Pointwise convolutions:** doesn't mix spatial dimensions

MLP-Mixer

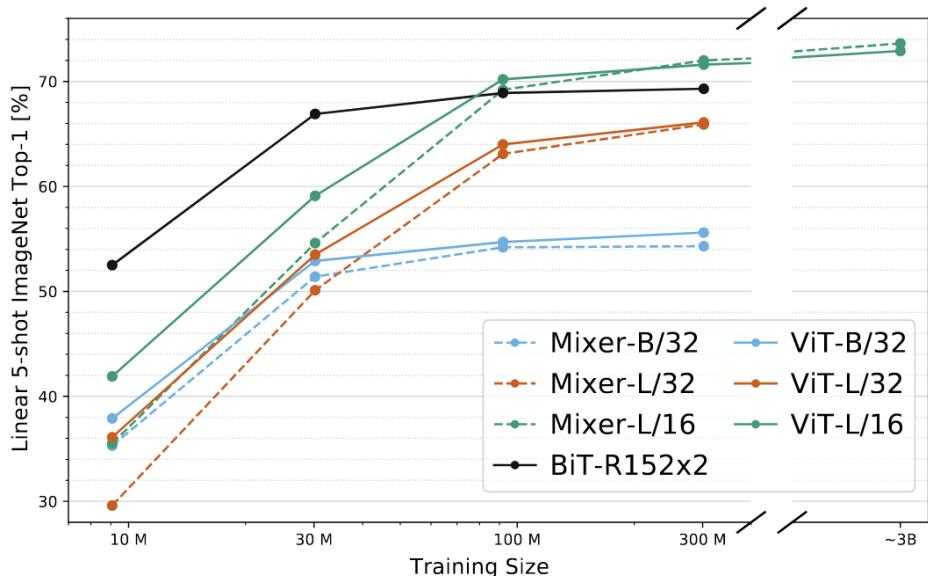


- Similarly to Separable Convolutions, apply a MLP on the channels dimension and a MLP on the patches dimensions
- Multiple other papers had the same idea at the same time (including ResMLP)

MLP-Mixer



	ImNet top-1	ReaL top-1	Avg 5 top-1	VTAB-1k 19 tasks	Throughput img/sec/core	TPUv3 core-days
Pre-trained on ImageNet-21k (public)						
● HaloNet [51]	85.8	—	—	—	120	0.10k
● Mixer-L/16	84.15	87.86	93.91	74.95	105	0.41k
● ViT-L/16 [14]	85.30	88.62	94.39	72.72	32	0.18k
● BiT-R152x4 [22]	85.39	—	94.04	70.64	26	0.94k
Pre-trained on JFT-300M (proprietary)						
● NFFNet-F4+ [7]	89.2	—	—	—	46	1.86k
● Mixer-H/14	87.94	90.18	95.71	75.33	40	1.01k
● BiT-R152x4 [22]	87.54	90.54	95.33	76.29	26	9.90k
● ViT-H/14 [14]	88.55	90.72	95.97	77.63	15	2.30k
Pre-trained on unlabelled or weakly labelled data (proprietary)						
● MPL [34]	90.0	91.12	—	—	—	20.48k
● ALIGN [21]	88.64	—	—	79.99	15	14.82k



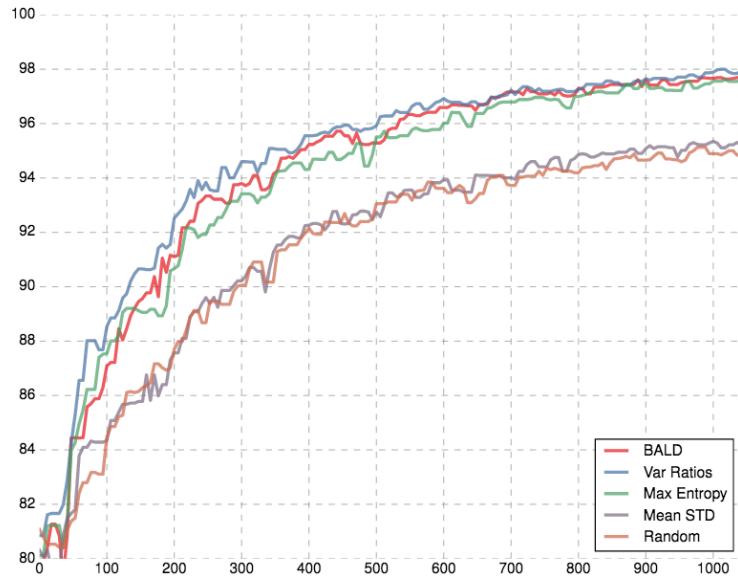
- However needs even more data than transformers based (ViT)
- Convolutions models (here BiT-R152x2) are still much more data efficient when training from scratch
 - However, some advocates that most future applications will be based on those new large (transformer, mlp, etc.) models pretrained on large-scale data
 - See [[Bommasani et al. arXiv 2021](#)]

Tricks that work for most architectures



- It's extremely important to tune the hyperparameters on a val set
- In real-life, you often don't have access to the full test set
 - And this test set may change constantly
- It's also important to ensure that your train and val sets have the same distribution than the test set
 - Beware of the **sampling bias**, e.g. I'm only labeling images of cars oriented towards the front, but in the test / real-life I may see cars in other orientations
 - See [[Torralba and Efros, CVPR 2011](#)]

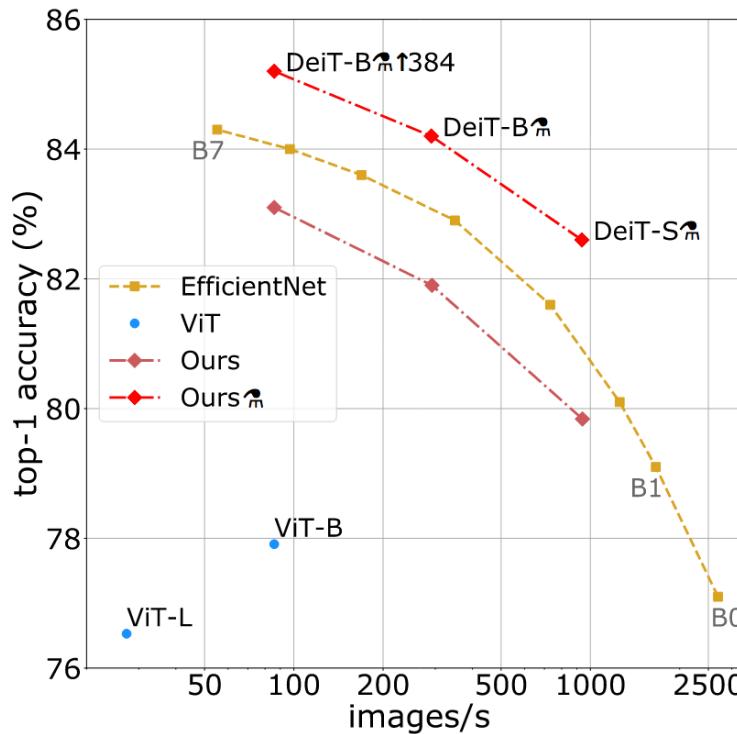
Active Learning



- More data is (almost) always better, see the Sutton's [Bitter Lesson](#)
- But hand-labeling data is very costly, both in \$ and time
- Active learning aims to determine which data to labelize in priority to be added to the training set
 - A lot of the literature is based on Bayesian stats, with Yarin Gal's team
 - But often done on small-scale datasets (MNIST, CIFAR)
 - And a random sampling is often quite competitive despite its simplicity



Thresholding for open-set prediction recall/precision/f1

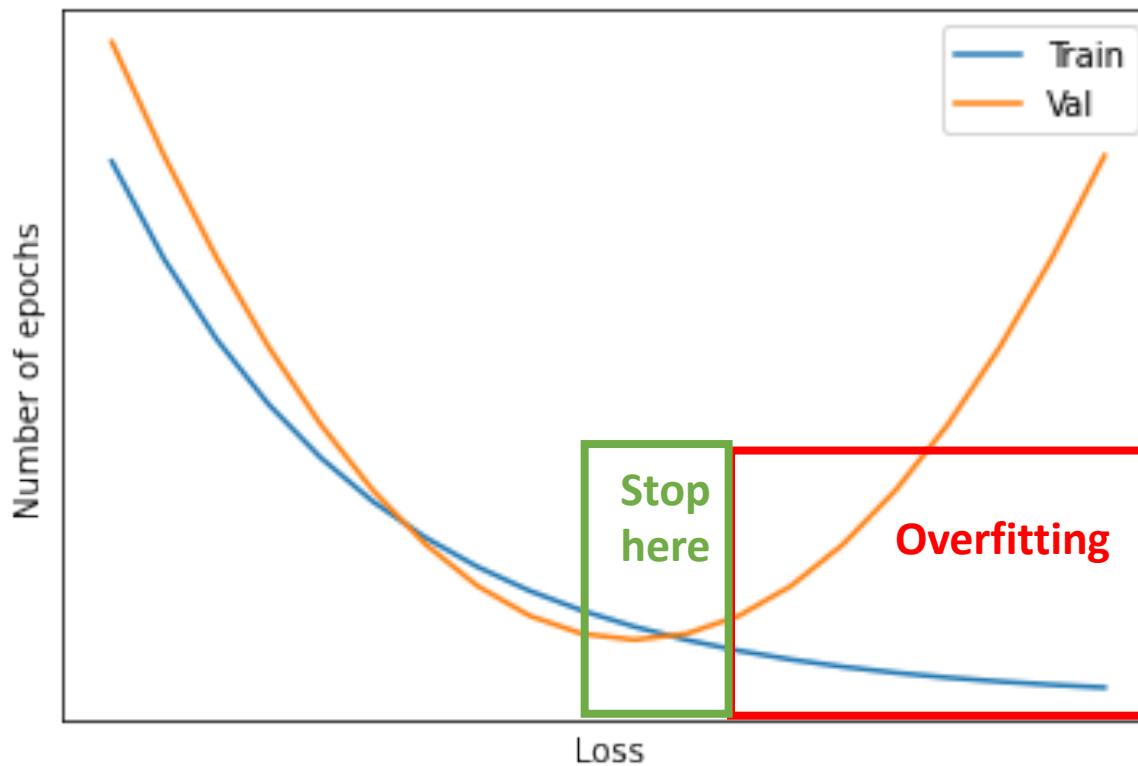


- More generally transformers for vision seems to be leading the SotA in computer vision
- And they manage to be super fast
- Although, they need more data and regularizations than your common ConvNet

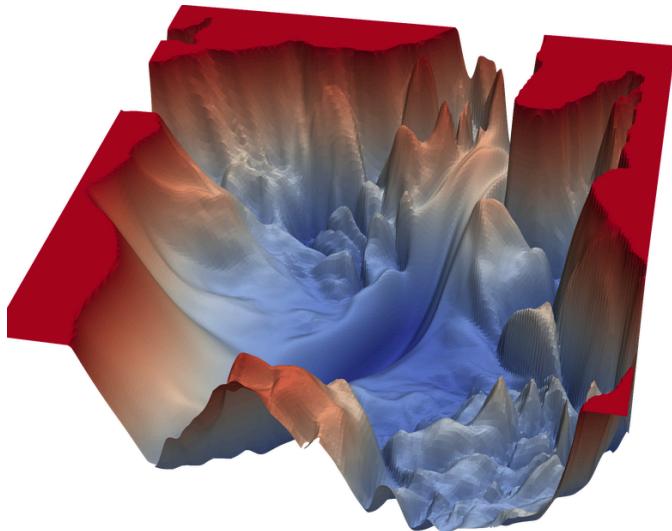
Early Stopping



- Training for too long, and you start to overfit. So when to stop?
- Monitor a metric (accuracy, loss, etc.) on the **VALIDATION SET** (!) and if this metric gets worse for X epochs, stop training
- "*A beautiful free lunch*" according to Turing award's Geoffrey Hinton



Decrease Learning Rate



- A high learning rate during the beginning of the training may help
 - Acts as a regularization by **skipping the local minima that are too sharp** and thus usually generalize less
- Then **decrease learning rate gradually to go deeper in a local minima** towards the training end
 - Either decrease learning rate (usually divided by 10) at particular epochs
 - Or decrease if validation metric (loss, acc, etc.) doesn't improve
- Some scheduling as **Cosine** decreases and increases (a little less) repetitively

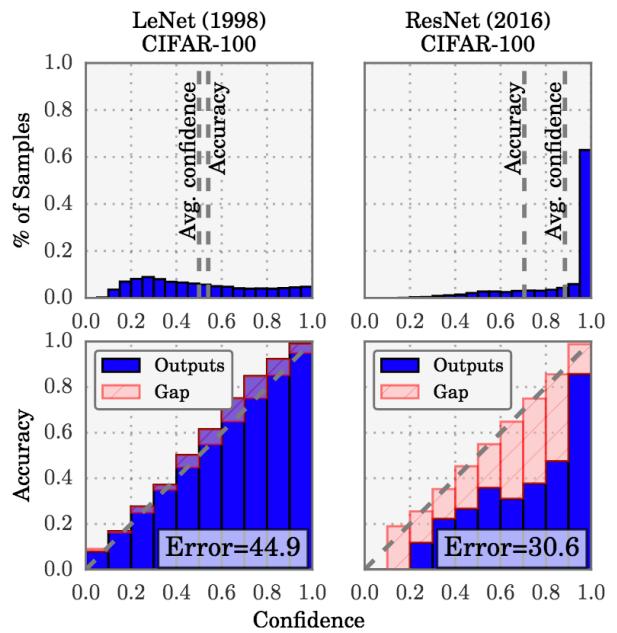
Smooth Labeling



$$y_k^{LS} = y_k(1 - \alpha) + \alpha/K$$

[0, 0, 1, 0] with a $\alpha = 0.1$ produces [0.025, 0.025, 1, 0.025]

- Avoid overconfidence in the model, when confidence is always close to 0.999...
 - And thus reduce overfitting
- Avoid miscalibration where model confidence is not correlated to model accuracy





Image

ResNet-50



Mixup [48]



Cutout [3]



CutMix



Label

Dog 1.0

Dog 0.5
Cat 0.5

Dog 1.0

Dog 0.6
Cat 0.4

- Mix two images and their labels together
 - In practice MixUp mix them with factors like 0.9/0.1 (not 0.5/0.5 as on the image)
- Acts as regularization to reduce overfitting
- A LOT of alternative to MixUp exists (CutOut, CutMix, FixMatch, MixMo, PuzzleMix, etc.)



Knowledge Distillation

- Train one super-mega-large model (called teacher)
- **Distill** the knowledge of the teacher onto a smaller model (called student)
- In practice, train students as usual but add another loss:
 - **KL-divergence** between the probabilities of the teacher and the student
 - The probabilities act as **dark knowledge** with extra information
 - (aka if the teacher says this is a dog with 0.7 confidence, we know it's a dog, but it's probably not the most archetypal dog ever)
 - Often add a **temperature T** on the logits before softmax
 - If $T > 1$, reduces the sharpness of the probabilities leading to more useful info

$$q_i = \frac{\exp(z_i/T)}{\sum_j \exp(z_j/T)}$$

Thresholding & Open-Set

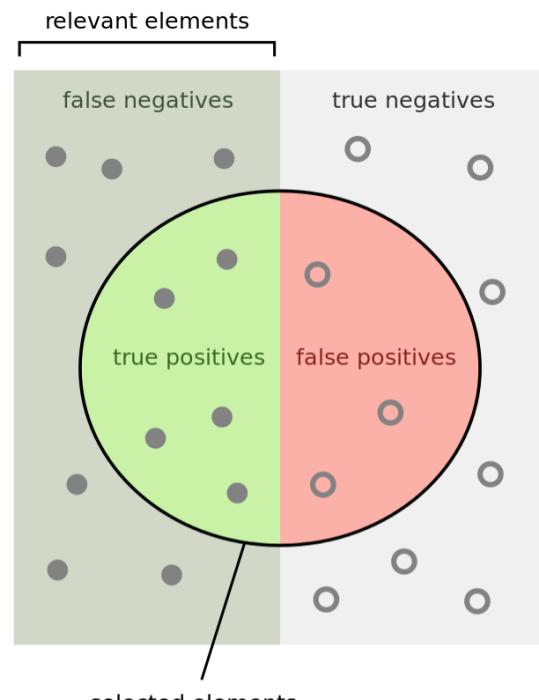


- In the real-life test set, you don't always want to predict on all images
- Example:
 - *at Heuritech we predict fashion trends from images scrapped from Instagram. Given an image of dog, my model should predict any trend!*
- You choose to discard model prediction if its confidence is lower than a threshold
 - One threshold per class is better
 - Compute threshold on the validation set (which needs to have negative images!)



Accuracy vs Precision+Recall=F1

- In real-life classes are never equally balanced
- If there are 90% of *dogs*, and 10% of *cats*, and your model has less than 90% of accuracy it's bad...
 - Answering *dog* every times gives 90% accuracy
- Best to use other metrics like **Precision** and **Recall** or their combination the **F1-Score**
- **Recall** is particularly useful in **open-set** where your model shouldn't predict on all images

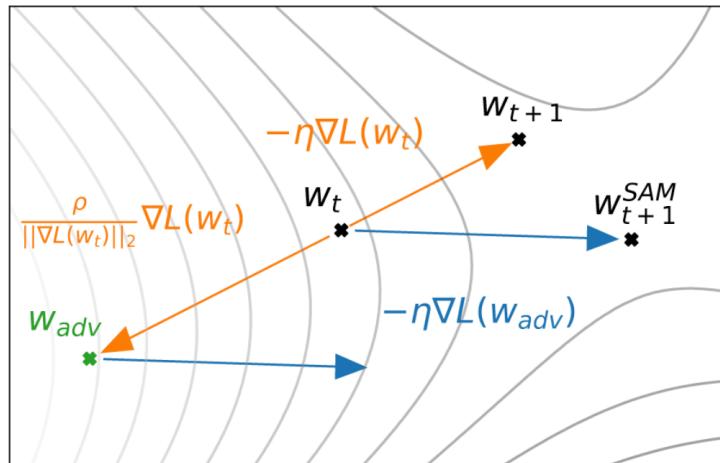


$\text{Precision} = \frac{\text{How many selected items are relevant?}}{\text{How many selected items are selected?}}$	$\text{Recall} = \frac{\text{How many relevant items are selected?}}{\text{How many relevant items are selected?}}$
--	---

$$F_1 = \frac{2}{\text{recall}^{-1} + \text{precision}^{-1}} = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}} = \frac{\text{tp}}{\text{tp} + \frac{1}{2}(\text{fp} + \text{fn})}$$



SAM: Sharpness-Aware Minimization



Input: Training set $\mathcal{S} \triangleq \cup_{i=1}^n \{(x_i, y_i)\}$, Loss function $l : \mathcal{W} \times \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}_+$, Batch size b , Step size $\eta > 0$, Neighborhood size $\rho > 0$.

Output: Model trained with SAM

Initialize weights w_0 , $t = 0$;

while not converged **do**

- Sample batch $\mathcal{B} = \{(x_1, y_1), \dots, (x_b, y_b)\}$;
- Compute gradient $\nabla_w L_{\mathcal{B}}(w)$ of the batch's training loss;
- Compute $\hat{\epsilon}(w)$ per equation 2;
- Compute gradient approximation for the SAM objective (equation 3): $\mathbf{g} = \nabla_w L_{\mathcal{B}}(w)|_{w+\hat{\epsilon}(w)}$;
- Update weights: $w_{t+1} = w_t - \eta \mathbf{g}$;
- $t = t + 1$;

end

return w_t

- Do not optimize network on a particular point of the parameters space but rather a **region**
- All neighbors parameters must also be good, leading to wider optimum and thus better generalization
- Needs twice more forward/backward...

Small break,
then coding session!