

INTRODUCTION

Motivation

• Convolutional networks vs. vision transformers: different poolings?







patch token representation

- Vision transformers: [CLS] necessity?
- Supervised vision transformers: low-quality spatial attention?

Focus

- Pooling for both network types improving over default?
- Pooling for high-quality spatial attention?
- Validity in both supervised & self-supervised setting?

Approach

- Framework to inspect & compare pooling methods.
- Introduce SimPool: attention-based, universal pooling for robust vector representation.

Ι	LANDSCAPE OF POOLINGS													
#	Method	CAT	Iter	k	U^0	$\phi_Q(U)$	$\phi_K(X)$	$s(\mathbf{x},\mathbf{y})$	A	$\phi_V(X)$	f(x)	$\phi_X(X)$	$\phi_U(Z)$	ACC (%)
1	GAP [16] max [2] GeM [3] LSE [4] HOW [5]			1 1 1 1 1					$\begin{array}{c} 1_{p}/p\\ 1_{p}\\ 1_{p}/p\\ 1_{p}/p\\ \mathrm{diag}(X^{\top}X) \end{array}$	$egin{array}{c} X \ X \ X \ X \ X \end{array}$ FC $(\operatorname{avg}_3(X))$	$f_{-1}(x)$ $f_{-\infty}(x)$ $f_{\alpha}(x)$ e^{rx} $f_{-1}(x)$		Z Z Z Z Z	55.0 53.9 55.9 55.3 54.8
2	OTK [6] k-means Slot [7]*	✓ ✓	\checkmark	$k \ k \ k$	U random U	$U \\ U \\ W_Q U$	$egin{array}{c} X \ X \ W_K X \end{array}$	$\frac{-\ \mathbf{x} - \mathbf{y}\ ^2}{-\ \mathbf{x} - \mathbf{y}\ ^2}$ $\frac{\mathbf{x}^\top \mathbf{y}}{\mathbf{x}^\top \mathbf{y}}$	SINKHORN $(e^{S/\epsilon})$ $\eta_2(\arg \max_1(S))$ $\sigma_2(S/\sqrt{d})$	$\psi(X) \ X \ W_V X$	$f_{-1}(x)$ $f_{-1}(x)$ $f_{-1}(x)$	$X \ X$	$egin{array}{c c} Z \\ MLP(GRU(Z)) \end{array}$	55.9 55.4 56.7
3	SE [8] CBAM [9]*			1 1	$\pi_A(X) \\ \pi_A(X)$	$\sigma(extsf{mlp}(U)) \ \sigma(extsf{mlp}(U))$	X	$\mathbf{x}^{ op}\mathbf{y}$	$\sigma(\operatorname{conv}_7(S))$	$diag(\mathbf{q})X$ $diag(\mathbf{q})X$		V V diag(\mathbf{a})		55.7 55.6
4	ViT [10]* CaiT [11]*		\checkmark	$1 \\ 1$	$U \ U$	$g_m(W_Q U) \ g_m(W_Q U)$	$\frac{g_m(W_K X)}{g_m(W_K X)}$	$\mathbf{x}^ op\mathbf{y} \mathbf{x}^ op\mathbf{y} \mathbf{x}^ op\mathbf{y}$	$\overline{\boldsymbol{\sigma}_2(S_i/\sqrt{d})_{i=1}^m}$ $\overline{\boldsymbol{\sigma}_2(S_i/\sqrt{d})_{i=1}^m}$	$g_m(W_V X) \\ g_m(W_V X)$	$ \begin{array}{c} f_{-1}(x) \\ f_{-1}(x) \end{array} $	MLP(MSA(X)) X	$\frac{\operatorname{MLP}(g_m^{-1}(Z))}{\operatorname{MLP}(g_m^{-1}(Z))}$	56.1 56.7
5	SimPool	\checkmark		1	$\pi_A(X)$	$W_Q U$	$W_K X$	$\mathbf{x}^{ op}\mathbf{y}$	$\sigma_2(S/\sqrt{d})$	$X - \min X$	$f_{\alpha}(x)$		Z	57.1

l Develop a generic pooling framework.

2 Formulate methods as instantiations.

3 Discuss the properties.

REFERENCES

1] He et al. Deep residual learning for image recognition *CVPR*, 2016.] Tolias et al. Particular object retrieval with integral max-pooling of CNN activations *ICLR*, 2016.] Radenović et al. Fine-Tuning CNN Image Retrieval with No Human Annotation *PAMI*, 2018. [4] Pinheiro et al. From image-level to pixel-level labeling with convolutional networks *CVPR*, 2015.] Tolias et al. Learning and aggregating deep local descriptors for instance-level recognition *ECCV*, 2020. [6] Mialon et al. A trainable optimal transport embedding for feature aggregation and its relationship to attention arXiv. 2020. [7] Locatello et al. Object-centric learning with slot attention *NeurIPS*, 2020.

- [8] Hu et al. Squeeze-and-Excitation Networks CVPR, 2018
- Woo et al. Cham: Convolutional block attention module *ECCV*, 2018 [10] Dosovitskiy et al. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale *ICLR*,

- 4 Derive SimPool.
- 5 Benchmark.
- [11] Touvron et al. Going deeper with Image Transformers *ICCV*, 2021. [13] Hu et al. Gather-excite: Exploiting feature context in convolutional neural networks *NeurIPS*, 2018.
- [15] Caron et al. Emerging properties in self-supervised vision transformers *ICCV*, 2021. [16] He et al. Deep residual learning for image recognition *CVPR*, 2016. [17] Liu et al. A convnet for the 2020s CVPR, 2022. [18] Deng et al. Imagenet: A large-scale hierarchical image database *CVPR*, 2009. [19] Wah et al. The Caltech-UCSD Birds-200-2011 Dataset CIT, 2011. Everingham et al. The pascal visual object classes (voc) challenge *IJCV*, 2009.

Keep It SimPool: Who Said Supervised Transformers Suffer from Attention Deficit?

Bill Psomas^{1,2}, Ioannis Kakogeorgiou¹, Konstantinos Karantzalos¹, Yannis Avrithis² ¹National Technical University of Athens; ²Institute of Advanced Research in Artificial Intelligence (IARAI) psomasbill@mail.ntua.gr

SIMPOOL: A SIMPLE ATTENTION-BASED UNIVERSAL POOLING



- Initial representation: $\mathbf{u}^0 = \pi_A$ by GAP.
- \mathbf{u}^0 (**X**) mapped by W_Q (W_K) to form \mathbf{q} (**K**).
- Attention map: $\mathbf{a} = \boldsymbol{\sigma}_2 \left(K^{\top} \mathbf{q} / \sqrt{d} \right).$

PROPERTY: UNIVERSAL (NETWORKS & SETTINGS)

Method Ep F	RESNET-50	CONVNEXT-S	S VIT-S	Image classification accuracy on ImageNet-1
Baseline 100	77.4	81.1	72.7	Left: Supervised; Below: self-supervised [15].
CaiT [11] 100	77.3	81.2	72.6	line: GAP for convolutional, [CLS] for transform
Slot [7] 100	77.3	80.9	72.9	
GE [13] 100	77.6	81.3	72.6	METHOD EP RESNET-50 CONVINEAT-5 VII
SimPool 100	78.0	81.7	74.3	k-NN Prob k -NN Prob k -NN
Baseline 300	78.1^{\dagger}	83.1	77.9	Baseline 100 61.8 63.0 65.1 68.2 68.9
SimPool 300	78.7 [†]	83.5	78.7	SimPool 100 63.8 64.4 68.8 72.2 69.8

PROPERTY: HIGH-QUALITY ATTENTION MAPS FROM TRANSFORMERS

Attention maps of ViT-S trained on ImageNet-1k. For baseline, we use the mean attention map of the [CLS] token. For SimPool, we use the attention map **a**.



PROPERTY: RESOLVING THE ATTENTION DEFICIT

[CLS] vs. SimPool. Attention maps of ViT-T trained on ImageNet-1k for 100 epochs under supervision. For [CLS], we use the mean attention map of the [CLS] token of each block. For SimPool, we use the attention map a.



global representation



input image

supervised [CLS]

[CLS]

[12] Lee et al. Generalizing pooling functions in convolutional neural networks: Mixed, gated, and tree *AISTATS*, [14] Siméoni et al. Localizing Objects with Self-Supervised Transformers and no Labels *BMVC*, 2021.

[21] Vo et al. Toward unsupervised, multi-object discovery in large-scale image collections *ECCV*, 2020.

• Global representation: $\mathbf{u} = \pi_{SP}(X) := f_{\alpha}^{-1}(f_{\alpha}(V)\mathbf{a}),$ where: $\int \frac{1-\alpha}{\alpha}$ if $\alpha / 1$

$$f_{\alpha}(x) := \begin{cases} x^{\overline{2}}, & \text{if } \alpha \neq 1, \\ \ln x, & \text{if } \alpha = 1. \end{cases}$$

PROPERTY: HIGH-QUALITY ATTENTION MAPS FROM CNNS

Attention maps of ResNet-50 [16] and ConvNeXt-S [17] trained on ImageNet-1k under supervision and selfsupervision [15]. We use the attention map **a**.



PROPERTY: LOCALIZATION

Left: Object localization MaxBoxAccV2 on CUB [19] test and ImageNet-1k validation set with ViT-S. Right: Object discovery CorLoc on VOC [20] trainval set and COCO 20k [21] with self-supervised ViT-S.

Method	SU	PERVISED	Self-Supervised		Method	DING	D-SEG [1	LOST [14]			
	CUB	IMAGENET	CUB	IMAGENET		VOC07	VOC12	2 COCO	VOC07	VOC12	2 (
Baseline	63.1	53.6	82.7	62.0	Baseline	30.8	31.0	36.7	55.5	59.4	
SimPool	77.9	64.4	86.1	66.1	SimPool	53.2	56.2	43.4	59.8	65.0	
Baseline@20	62.4	50.5	65.5	52.5	Baseline@20) 14.9	14.8	19.9	50.7	56.6	
SimPool@20	74.0	62.6	72.5	58.7	SimPool@20) 49.2	54.8	37.9	53.9	58.8	4



EFFICIENCY	&	ABL

Network	POOLING	Depth	Init	ACCURACY	#PARAMS	
BASE	GAP	12	12	73.3	22.1M	
BASE		12	0	72.7	22.1M	
BASE + 1		13	0	73.2	23.8M	
BASE + 2		14	0	73.7	25.6M	
BASE + 3	[CLS]	15	0	73.8	27.4M	
BASE + 4		16	0	73.9	29.2M	
BASE + 5		17	0	74.6	30.9M	
BASE		12	12	74.3	22.3M	
BASE - 1	SimDool	11	11	73.9	20.6M	
BASE - 2	SIIIF 001	10	10	73.6	18.7M	
BASE - 3		9	9	72.5	17.0M	

CONCLUSION

- Outperforms the other pooling methods.
- Presents strong localization properties.

institute of advanced research in artificial intelligence













DINO []

convNeXt-S ConvNeXt-S

DINO [15]

Object localization on ImageNet-1k with supervised ViT-S. Green: ground-truth; red: baseline, blue: SimPool.

LATION

Left: Trade-off between performance and parameters. Supervised ViT-S on ImageNet-1k. BASE+b (BASE-b): *b* blocks added to (removed from) the network. Below: Image classification accuracy on ImageNet-1k vs. exponent $\gamma = (1 - \alpha)/2$ for ResNet-18 and ViT-T.



• Improves performance of convolutional networks and transformers under supervised or self-supervised setting.

• Produces high-quality attention maps that delineate object boundaries.

• Discards the traditional [CLS] from vision transformers.