



Reducing the Transformer Architecture to a Minimum

Bernhard Bermeitinger^{†1}^a, Tomas Hrycej^{†2}, Massimo Pavone^{‡2}, Julianus Kath^{‡2}, and Siegfried Handschuh^{†2}^b

¹*Institute of Computer Science in Vorarlberg, University of St.Gallen (HSG), Dornbirn, Austria*

²*Institute of Computer Science, University of St.Gallen (HSG), St.Gallen, Switzerland*

[†]*firstname.lastname@unisg.ch*, [‡]*firstname.lastname@student.unisg.ch*

Keywords: attention mechanism, transformers, computer vision, model reduction, deep neural networks

Abstract: Transformers are a widespread and successful model architecture, particularly in Natural Language Processing (NLP) and Computer Vision (CV). The essential innovation of this architecture is the Attention Mechanism, which solves the problem of extracting relevant context information from long sequences in NLP and realistic scenes in CV. A classical neural network component, a Multi-Layer Perceptron (MLP), complements the attention mechanism. Its necessity is frequently justified by its capability of modeling nonlinear relationships. However, the attention mechanism itself is nonlinear through its internal use of similarity measures. A possible hypothesis is that this nonlinearity is sufficient for modeling typical application problems. As the MLPs usually contain the most trainable parameters of the whole model, their omission would substantially reduce the parameter set size. Further components can also be reorganized to reduce the number of parameters. Under some conditions, query and key matrices can be collapsed into a single matrix of the same size. The same is true about value and projection matrices, which can also be omitted without eliminating the substance of the attention mechanism. Initially, the similarity measure was defined asymmetrically, with peculiar properties such as that a token is possibly dissimilar to itself. A possible symmetric definition requires only half of the parameters. All these parameter savings make sense only if the representational performance of the architecture is not significantly reduced. A comprehensive empirical proof for all important domains would be a huge task. We have laid the groundwork by testing widespread CV benchmarks: MNIST, CIFAR-10, and, with restrictions, ImageNet. The tests have shown that simplified transformer architectures (a) without MLP, (b) with collapsed matrices, and (c) symmetric similarity matrices exhibit similar performance as the original architecture, saving up to 90 % of parameters without hurting the classification performance.


1 INTRODUCTION


Recently, *Large Language Models* (LLMs) have shown impressive performance in producing complex text answers to given questions. Their outstanding feature is the massive size of parameter sets (up to billions). The rapidly growing parameter number has limited the possibility of developing such models (as well as objectively investigating their properties) to companies and institutions capable of making considerable investments in computing the model's parameters.

This is why it is of great interest to attempt to find more efficient configurations with fewer parameters without performance loss. A computing model with an excellent success record is based on the trans-

former architecture (Vaswani et al., 2017). Their success is due to an excellent ability to capture contextual information. Initially developed for language processing, transformers have also been successfully used in Computer Vision (CV). The analogy to language processing is the following: the semantics of individual words are determined by other words in the word sequence. Frequently, the basic units are not words but tokens (e.g., n -grams consisting of n consecutive letters). Since the *Vision Transformer* (Dosovitskiy et al., 2021), in an image, the tokens are represented by *patches* — typically square regions of pixels in the image. Other patches can influence or disambiguate a patch's conceptual meaning. For example, the environment in which an individual object is embedded in the image may disambiguate the identification of a specific bird or mushroom species.

The fundamental concept of the transformer is that

^a <https://orcid.org/0000-0002-2524-1850>

^b <https://orcid.org/0000-0002-6195-9034>

of *attention* (Bahdanau et al., 2016). It is based on the insight that a particular token’s semantics are influenced by its close relationships with other tokens. The tokens are encoded as real-valued vectors in a high-dimensional space (frequently around 1,000 dimensions or more). These vectors are called *embeddings*. The algebraic similarity between the embedding vectors measures the semantic proximity between the tokens. This similarity measure is the vector product or the cosine angle between the vectors. The weighting of tokens by such similarity measure is called attention, which, in analogy to human attention, focuses on relevant concepts. From the computational point of view, a transformer is a structure consisting of

- an algorithm for consideration of token context, the *attention mechanism*, and
- a *Multi-Layer Perceptron* (MLP) for nonlinear transformation of intermediary data.

Multi-Head Attention For every transformer in the stack, the following processing is done by the attention mechanism (*multi-head attention* or *MHA*). The input of a training sample in the stack’s s -th Transformer (out of their total number S) is a sequence of input vectors x_{sj} . This sequence is transformed into an equally long sequence of output embeddings z_{si} . Each of them is, for given weights, a formally linear transformation

$$\begin{aligned} z_{si} &= \left(\sum_{j=1}^i a_{sij} x_{sj} W_s^V \right) W_s^O \\ &= \left(\sum_{j=1}^i a_{sij} x_{sj} \right) W_s^V W_s^O \end{aligned} \quad (1)$$

i.e., a weighted average of input embeddings x_{sj} , linearly transformed by matrix $W_s^V W_s^O$. The weight vectors $a_{si} = [a_{si1}, a_{si2}, \dots, a_{sii}]$ are computed as

$$a_{si} = \text{Softmax}(s_{si}) \quad (2)$$

The vector argument of the $\text{Softmax}()$ function measures the similarity between a present token x_Q , “the query” and another token x_K , “the key”.

$$s_{sij} = x_{si} W_s^Q W_s^{KT} x_{sj}^T \quad (3)$$

This form of attention mechanism is referred to as *single-head*. A popular variant consists of an extension to multiple heads indexed by h :

$$z_{si} = \sum_{h=1}^H \left(\sum_{j=1}^i a_{shij} x_{sj} \right) W_{sh}^V W_{sh}^O \quad (4)$$

Each head has its separate matrices W_h^Q , W_h^K , W_h^V , and W_h^O . The weights are also computed separately as

$$a_{shi} = \text{Softmax}(s_{shi}) \quad (5)$$

and

$$s_{shij} = x_{si} W_{sh}^Q W_{sh}^{KT} x_{sj}^T \quad (6)$$

Multi-Layer Perceptron The second component is a standard MLP with a single hidden layer, applied to each intermediary embedding z_{si} :

$$\begin{aligned} h_{si} &= f \left(z_{si} W_s^{(1)} + b_s^{(1)} \right) \\ y_{si} &= h_{si} W_s^{(2)} + b_s^{(2)} \end{aligned} \quad (7)$$

with $f()$ being a nonlinear function, usually the *Gaussian Error Linear Unit* (*GELU*) (Hendrycks and Gimpel, 2023), weight matrices $W_s^{(1)}$ and $W_s^{(2)}$ as well as bias vectors $b_s^{(1)}$ and $b_s^{(2)}$.

(He and Hofmann, 2024) have investigated the possibilities of simplifying the transformer architecture. Their focus has been increasing the signal throughput through the network. The proposed changes primarily consist of modifying or omitting shortcut connections and normalizing layers. In addition, they have addressed the possibility of omitting matrices W^V and W^O . The last idea has also been implemented in our modifications proposed in Section 3.

Our focus is different: we intend to substantially reduce trainable parameters to accelerate the training and improve convergence.

2 TRANSFORMER WITHOUT THE MLP

The MLP requires the majority of the parameters to be fitted. This is justified by the argument that the MLP is the vehicle for implementing nonlinear mappings.

However, it can be argued that the first component, the attention mechanism, can also capture nonlinearities. It is the variable weights that make the mapping nonlinear. The argument of the $\text{Softmax}()$ function is already a quadratic function of input tokens, and the function itself is nonlinear. Even if the $\text{Softmax}()$ were linear, the multiplication of input tokens by the weights a_{sij} (which are quadratic in these tokens) would result in a cubic function of input tokens. The nonlinearity of $\text{Softmax}()$ makes this mapping only more nonlinear.

So, a stack of S transformers is a chain of S at least cubic functions of the input, resulting in a function of polynomial order of at least $3S$. This makes clear that subsequent processing by an MLP is not the only nonlinear element of the processing. The extent of the task’s nonlinearity cannot be assessed in advance. Still, the hypothesis that a reduced transformer without an MLP may cover the nonlinearity needs for some tasks is justified and can be validated by appropriate tests.

Without the MLPs, the transformer architecture

can be described in more explicit terms. This is particularly the case if a single-head option is pursued.

3 SINGLE-HEAD CONFIGURATION

Although the matrices W_s^Q , W_s^K , W_s^V , and W_s^O can theoretically map the embedding vector to an arbitrary vector width, it is common to keep this width constant throughout the model, referring to the *model width* N . Then, in the case of a single head, these matrices are square. With square matrices, it is evident that $W_s^V W_s^O$ can be collapsed to a single matrix W_s^{VO} , and, analogically, $W_s^Q W_s^{KT}$ to W_s^{QK} . This saves 50% of the attention module’s parameters, from $4SN^2$ to $2SN^2$.

Concatenating the transformer-encoder layers without MLP leads to the following recursion:

$$\begin{aligned}
 y_{1i} &= \left(\sum_{j=1}^i a_{1ij} x_{1j} \right) W_1^{VO} \\
 y_{2i} &= \left(\sum_{j=1}^i a_{2ij} y_{1j} \right) W_2^{VO} \\
 &= \left(\sum_{k=1}^i a_{2ik} \left(\sum_{j=1}^k a_{1kj} x_{1j} \right) W_1^{VO} \right) W_2^{VO} \quad (8) \\
 &= W_1^{VO} W_2^{VO} \sum_{k=1}^i a_{2ik} \sum_{j=1}^k a_{1kj} x_{1j} \\
 &\dots
 \end{aligned}$$

When stacking the attention modules, the matrices W_s^{VO} concatenate to their product over $s = 1, \dots, S$. Then, they collapse into a single matrix

$$W^{VO} = \prod_{s=1}^S W_s^{VO} \quad (9)$$

Since every sum $\sum_{j=1}^i a_{sij}$ is equal to unity (as a result of the softmax operation), every successive transformer layer performs a weighted mean of stacked inputs x_{1j} .

The total number of parameters with S matrices W_s^{QK} and a single matrix W^{VO} is $(S+1)N^2$, only slightly more than 25% of the original size without MLP. So far, all this is possible without losing any expressive power of the single-head transformer without MLP — only obsolete parameters are deleted.

In many NLP applications, the output of the last transformer of the stack is expected to produce an embedding of a word or a language token. These output embeddings can be expected to come from the space spanned by the input words or tokens. From

this viewpoint, it may appear questionable to transform the input embeddings by matrices W_s^{VO} and to re-transform them back into the word embeddings. Then, it may be worth attempting to delete the value transformations. This has also been the proposal of (He and Hofmann, 2024), resulting in a simple weighted mean

$$z_{si} = \sum_{j=1}^i a_{sij} x_{sj} \quad (10)$$

The output embedding z_{si} is a convex combination of input embeddings x_{1i} . In other words, it is a member of the convex set spanned by x_{1i} .

This concept has been implemented in the Keras framework by setting the matrices W_s^V and W_s^O to unit matrices. Collapsing $W_s^Q W_s^{KT}$ to W_s^{QK} has been reached by setting the matrix W^K to a unit matrix. The newly defined matrix W_s^{QK} replaces matrix W_s^Q .

4 MULTI-HEAD CONFIGURATION

The relationships of Section 3 are valid wherever the matrices W_{sh}^V , W_{sh}^O , W_{sh}^Q , and W_{sh}^K are square. This may also apply to multiple heads. However, it is usual to commit to a reduced dimension per head. With $H > 1$ heads, it is common to map the embedding vector to a narrower vector of width N/H , assumed to be integer.

In such cases, the matrices W_{sh}^V , W_{sh}^O , W_{sh}^Q , and W_{sh}^K are not square but of dimension $(N, N/H)$. Collapsing $W_{sh}^Q W_{sh}^{KT}$ to W_{sh}^{QK} is then no longer efficient since W_{sh}^{QK} is of dimension (N, N) and has thus N^2 parameters while W_{sh}^Q and W_{sh}^K together have $2N^2/H$, which is a smaller or equal number for $H > 1$.

Moreover, it is impossible to equivalently concatenate the value/projection matrices W_{sh}^{VO} to a unique product because of varying index h along various paths through the heads.

Nevertheless, omitting the W_{sh}^{VO} at all would have the same justification as for single-head configuration: the output embedding z_{si} would become a convex combination of input embeddings x_{1i} , which can be expected to correspond to a meaningful word or token.

5 SYMMETRY OF SIMILARITY

The expression Eq. (3) measures the similarity between queries and keys. The general concept of characterizing similarity between vectors by their product is symmetric: a is equally similar to b as b is to a .

However, the similarity between a key and a query evaluated with the help of $x_{si}W_{sh}^QW_{sh}^{KT}x_{sj}^T$ is asymmetric. This is because the matrices W_{sh}^Q and W_{sh}^K are potentially different.

This asymmetry leads to different similarities between x_{si} and x_{sj} in the roles of key and query: x_{si} is not as similar to x_{sj} as is x_{sj} to x_{si} . The vector x_{si} is also not the most similar to itself. The matrix product $W_{sh}^QW_{sh}^{KT}$ is generally not positive definite, so it is not even guaranteed that the similarity of x_{si} to itself is positive.

The asymmetry can be deliberate and justified from some viewpoints. It is not a matter of course that the roles of queries and keys are symmetric. However, some of the mentioned properties can make its use harmful.

The symmetry can be guaranteed by simply setting $W_{sh}^Q = W_{sh}^K$. Then, half of the parameters dedicated to the query and key matrices can be economized. In the single-head case, the same effect is reached by a symmetric matrix W_s^{QK} , with identical parameters mirrored over the diagonal, i.e., $w_{sij}^{QK} = w_{sji}^{QK}$. Another possibility is to parameterize a lower triangular matrix T_s^{QK} and to multiply it by its transpose, getting

$$W_s^{QK} = T_s^{QK}T_s^{QKT} \quad (11)$$

This amounts to the well-known *Cholesky decomposition* (Cholesky, 1924) of a symmetric matrix.

With both methods, the number of parameters is $\frac{N(N+1)}{2}$ instead of N^2 , or even $2N^2$ of the original version without collapsing W^Q and W^K .

The symmetry is implemented by reusing W_{sh}^Q as W_{sh}^K , omitting the use of W_{sh}^K at all.

6 SETUP OF COMPUTING EXPERIMENTS

The benchmarks for the evaluation have been chosen from the CV domain. They are medium-sized problems that can be run for a sufficient number of experiments. This would not be possible with large models such as those used in language processing.

For the experiments, two well-known image classification datasets MNIST (LeCun et al., 1998) and CIFAR-10 (Krizhevsky, 2009) were used. MNIST contains grayscale images of handwritten digits (0–9) while CIFAR-10 contains color images of exclusively ten different mundane objects like “horse”, “ship”, or “dog”. They contain 60,000 (MNIST) and 50,000

(CIFAR-10) training examples. Their respective pre-configured test split of each 10,000 examples are used as validation sets. While CIFAR-10 is evenly distributed among all classes, MNIST can be considered almost equally distributed.

An important criterion is that the training set size is sufficient for good generalization. The training size (as related to the number of model parameters) must be large enough for the model not to be underdetermined so that we can fairly assess the models’ performances. As a criterion for this, the overdetermination ratio of each benchmark candidate has been evaluated (Hrycej et al., 2023):

$$Q = \frac{KM}{P} \quad (12)$$

with K being the number of training examples, M being the output vector length (usually equal to the number of classes), and P being the number of trainable model parameters.

This formula justifies itself by ensuring that the numerator KM equals the number of constraints to be satisfied (the reference values for all training examples). This number must be larger than the number of trainable parameters for the system to be sufficiently determined. (Otherwise, there is an infinite number of solutions, most of which do not generalize.) This is equivalent to the requirement for the overdetermination ratio Q to be larger than unity.

The losses and accuracies in Table 1 show that the performance with 12 encoders is not superior to that with 6 encoders. The parameter set sizes with 12 encoders have been 563,242 with MLP and 198,100 without MLP. This is substantially more than 287,686 and 101,470, respectively, with 6 encoders. Consequently, the latter variant has been adopted as a baseline.

6.1 RESULTS FOR MNIST

Following the arguments of Sections 2 to 5, the following reduced transformer variants have been tested:

- with and without an MLP in each transformer-encoder,
- with 1 and 4 heads,
- with the original matrix configuration as well matrix pair W^Q and W^K collapsed into one matrix, W^V and W^O omitted (one head variants only), and
- with asymmetric and symmetric similarity measures.

The variants depicted refer to the matrix options:

- *unchanged* corresponds to the original attention module matrix variety;

Table 1: Results of 16 experiments on the two datasets MNIST and CIFAR-10 with 6 or 12 consecutive transformer encoders and 1 or 4 attention heads per encoder layer either with the default MLP inside each encoder layer or skipping it entirely. The loss and accuracy for the training and validation sets are reported after each model is trained for exactly 500 epochs.

Dataset	#Encs-#Heads	MLP?	Q	Train loss	Val. loss	Train. acc. [%]	Val. acc. [%]
MNIST	6-1	yes	2.15	0.0067	0.0747	99.78	98.38
	6-1	no	6.46	0.0277	0.1023	99.07	97.49
	6-4	yes	2.09	0.0018	0.0739	99.95	98.26
	6-4	no	5.91	0.0021	0.0912	99.92	98.29
	12-1	yes	1.08	0.0052	0.0652	99.81	98.71
	12-1	no	3.29	0.0117	0.0970	99.62	97.94
	12-4	yes	1.08	0.0025	0.0656	99.92	98.70
	12-4	no	3.29	0.0026	0.1002	99.93	98.10
CIFAR-10	6-1	yes	1.74	0.1533	2.2418	94.63	60.24
	6-1	no	4.93	0.9341	1.3590	66.16	55.30
	6-4	yes	1.74	0.1109	2.4033	96.01	60.46
	6-4	no	4.92	0.5621	1.6984	80.82	52.37
	12-1	yes	0.89	2.3026	2.3026	9.82	10.00
	12-1	no	2.52	0.5604	1.7219	79.48	54.06
	12-4	yes	0.89	0.0632	2.6379	97.92	58.02
	12-4	no	2.52	0.1787	2.3200	93.59	55.60

Table 2: Loss and accuracy for different variants of transformer-encoder modifications on MNIST: 1 or 4 heads, with or without the MLP, with a single W_{qk} matrix, no value and projection matrices, or a symmetric similarity measurement.

# Heads	MLP?	Modification	# Parameters	Q	Train loss	Val. loss	Train. acc. [%]	Val. acc. [%]
1	yes	unchanged	279,106	2.15	0.0067	0.0747	99.78	98.38
4	yes	unchanged	287,746	2.09	0.0018	0.0739	99.95	98.26
1	yes	Wqk	257,506	2.33	0.0037	0.0794	99.89	98.43
1	yes	Wqk+noWv,Vo	212,866	2.82	0.0063	0.0951	99.78	98.27
1	no	unchanged	92,890	6.46	0.0277	0.1023	99.07	97.49
4	no	unchanged	101,530	5.91	0.0021	0.0912	99.92	98.29
1	no	symmetry	69,910	8.58	0.0331	0.0783	98.85	97.80
4	no	symmetry	69,910	8.58	0.0158	0.0762	99.46	98.24
1	no	Wqk	70,570	8.50	0.0374	0.0996	98.70	97.60
1	no	Wqk+noWv,Vo	26,650	22.51	0.1697	0.1536	94.82	95.32

- *Wqk* variants use a single matrix for the product $W^Q W^{KT}$; these variants are only available for a single attention head, and their similarity measure is asymmetric as in the original version;
- *noWv,Vo* denotes omitting the value matrices W^V as well as the projection matrices W^O ; also, these variants imply a single attention head and asymmetric similarity measurement;
- *symmetric* variants are committed to symmetric similarity measures; W^V and W^O are left untouched.

The performances of the individual variants are given in Table 2. For better comparability, the losses are additionally depicted in Fig. 1.

The following observations can be made:

- The original variants with MLPs perform better than those without MLPs on the training set.
- By contrast, their advance disappears on the validation set, particularly if the symmetric similarity

metrics are used.

- The variant with asymmetric similarity without MLP is inferior to the analogical one with symmetric similarity.
- The minimum variant with query and key matrices W^Q, W^K collapsed to $W^{QK} = W^Q W^{KT}$ and additionally omitted value and projection matrices show a higher loss than other variants. This may be due to its dramatically reduced parameter number, which may lead to an insufficient capacity to capture nonlinearities.

As MNIST is a relatively easy benchmark, the accuracy results are very close to each other. The parameter numbers are substantially different. The symmetric variant without MLP has only about 25 % of the parameter number of the original, full variant with MLP. The variant with collapsed matrices has about 33 % of the original parameters. The parameters include, in addition to the attention modules of all transformer-encoders, the embedding matrix reducing

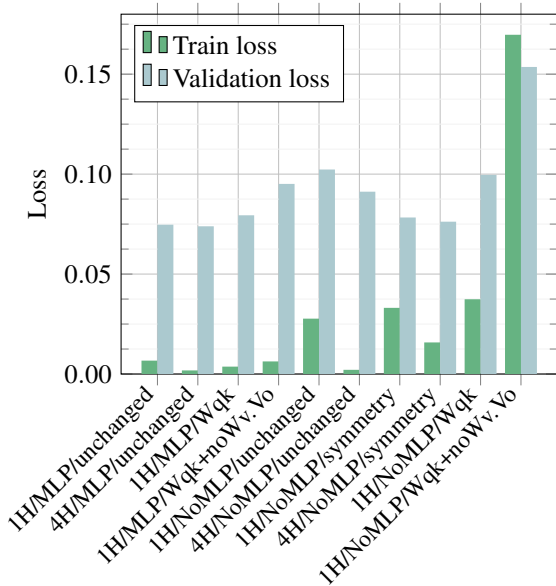


Figure 1: Training and validation losses attained by various reduced transformer-encoders with six encoder layers on MNIST.

the image patch to the embedding vector.

The number of parameters has a strong effect on the generalization capability of the model. This can be quantified with the help of the overdetermination ratio from Eq. (12) in column Q of Table 2. The loss gap between the training and validation sets is the largest for the original version with Q close to unity while it shrinks towards the symmetric version without MLPs.

6.2 RESULTS FOR CIFAR-10

The variants tested are analogical to those for MNIST. The losses and accuracies attained after 500 epochs are given in Table 3, the losses additionally in Fig. 2.

The result characteristics are similar to those for MNIST but more distinct:

- The original variant with MLP reaches the best training set loss but the worst validation set loss.
- Compared to the original variant, the reduced variants without MLP and with symmetric similarity are superior in generalization.
- This also applies to the variant with collapsed key and query matrices.
- Even the minimum variant with all considered matrix reductions (except for symmetry), whose parameter count is only a tenth of the original version with MLP, shows a better validation set performance than the original variant with all matrices and MLP.

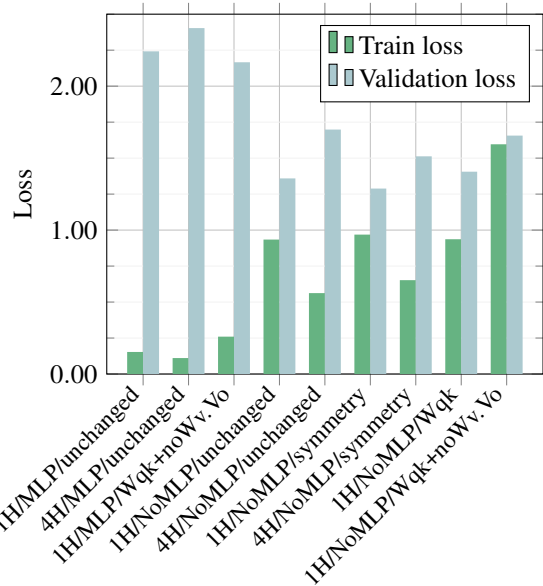


Figure 2: Training and validation losses attained by various reduced transformer-encoders with six encoder layers on CIFAR-10.

The measured accuracies are roughly consistent with the losses on the training set. On the validation set, some of them follow, paradoxically, a different ranking. However, the fact that the loss, not the accuracy, is explicitly trained justifies the arguments via loss rather than accuracy.

6.3 TRIALS WITH IMAGENET

Several trials on the ImageNet dataset (Russakovsky et al., 2015) have been conducted to support the hypotheses with a larger benchmark. Unfortunately, the baseline run with the original transformer architecture, including MLP, has not been successful. In all trials, *Adam* failed to find a substantial improvement in the initial parameter state. By contrast, without MLP, it has been converging at least to a state with a moderate classification performance. This is why we cannot present a serious study on ImageNet. It can only be concluded that discarding MLP is helpful for convergence. The proof that this variant's performance is acceptable is still pending, and further work will be required to provide it.

7 CONCLUSIONS AND LIMITATIONS

The experiments presented have shown limited utility of some parameter-extensive components of the trans-

Table 3: Loss and accuracy for different variants of transformer-encoder modifications on CIFAR-10: 1 or 4 heads, with or without MLP, with a single W_{qk} matrix, no value and projection matrices, or a symmetric similarity measurement.

# Heads	MLP?	Modification	# Parameters	Q	Train loss	Val. loss	Train. acc. [%]	Val. acc. [%]
1	yes	unchanged	287,686	1.74	0.1533	2.2418	94.63	60.24
4	yes	unchanged	287,746	1.74	0.1109	2.4033	96.01	60.46
1	yes	Wqk+noWv,Vo	221,446	2.26	0.2597	2.1659	90.53	54.98
1	no	unchanged	101,470	4.93	0.9341	1.3590	66.16	55.30
4	no	unchanged	101,530	4.92	0.5621	1.6984	80.82	52.37
1	no	symmetry	78,490	6.37	0.9686	1.2885	64.80	55.85
4	no	symmetry	78,490	6.37	0.6521	1.5125	76.10	55.52
1	no	Wqk	79,150	6.32	0.9364	1.4057	66.03	53.70
1	no	Wqk+noWv,Vo	35,230	14.19	1.5961	1.6565	40.52	39.17

former architecture. In particular, the following findings can be formulated:

- The MLP component is frequently presented as necessary for capturing nonlinearities in the modeled relationship. However, the inherent nonlinearity of the similarity measures seems powerful enough in many practical cases.
- While the classification performance without the MLPs is not significantly inferior to that with MLPs, a substantial benefit is saving the parameters. With model size N , the attention mechanism requires $4N^2$ parameters in the form of matrices W^Q , W^KW^V , and W^O . The size of the MLP is usually chosen as an integer multiple of h of the model size. Then, the MLP consists of weights and biases of two layers, with a total of $hN(N+1) + N(hN+1) = 2hN^2 + hN + N \approx 2hN^2$. If the multiple is $h = 4$, MLP has double the number of parameters as the attention mechanism. Consequently, omitting MLP reduces the parameters to 33 % of the original size.
- Symmetric similarity measures tend to perform better than asymmetric ones, with 50 % fewer query and key matrix parameters. This improvement may be reached by excluding undesirable freedoms, such as a token being dissimilar to itself. The parameter reduction can be expected to constrain the search for the optimum fit fruitfully.
- Collapsing the value and the key matrix into one is another possibility of reducing the parameter set of these matrices by 50 %.
- Omitting the value matrix W^V and the projection matrix W^O reduces the parameters of the whole attention module by 50 %. This variant has also been proposed by (He and Hofmann, 2024), with the observation of no significant performance loss in NLP benchmarks.
- Both preceding reductions amount to a reduction to 25 % of the original attention module size.
- In our experiments, the variants with the collapsed query/key matrices, omitted value, and projection matrices are slightly inferior for MNIST but equal for CIFAR-10. These minimum variants have less than 10 % of parameters compared with the classical transformers, including MLP. Compared to the architecture with 12 encoders, it is as little as 5 %.

The savings in computing time have been proportional to the savings in parameter numbers.

Our research has been limited to image processing benchmarks MNIST, CIFAR-10, and ImageNet. The experiments with the last benchmark have partially failed due to computing problems. Empirical evidence with the help of two medium-sized benchmarks and an incomplete test of a larger one is not satisfactory. This requests further research with more robust algorithms. There is considerable potential for second-order optimization methods such as the conjugate gradient algorithm of (Fletcher and Reeves, 1964), thoroughly described in (Press et al., 1992). This algorithm’s convergence is excellent, but implementing the stopping rule in widespread packages seems to improve its ability to prevent early stops before reaching the minimum region.

Limitations to image processing suggest further extension. The proper domain of transformers is NLP. An obstacle to its investigation is the size of benchmark problems, so most published investigations consist of observing the performance of fine-tuning pre-trained models. To use pre-trained parameter sets, these fine-tuned models must be identical or almost identical to the pre-trained models. This makes the testing of different architectures difficult. A possibility is to use a large model used for pre-training as a *teacher* and a medium-sized model as *student*, mimicking its performance. This procedure, referred to as *knowledge distillation*, has been proposed by (Hinton et al., 2015) and used, e.g., by (Sun et al., 2019).

These will be important focuses soon.

REFERENCES

- Bahdanau, D., Cho, K., and Bengio, Y. (2016). Neural Machine Translation by Jointly Learning to Align and Translate. *ICLR*.
- Cholesky, A.-L. (1924). Note Sur Une Méthode de Résolution des équations Normales Provenant de L'Application de la Méthode des Moindres Carrés a un Système D'équations Linéaires en Nombre Inférieur a Celui des Inconnues. — Application de la Méthode a la Résolution D'un Système Défini D'équations Linéaires. *Bulletin Géodésique*, 2(1):67–77.
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., and Houlsby, N. (2021). An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, page 21, Vienna, Austria.
- Fletcher, R. and Reeves, C. M. (1964). Function minimization by conjugate gradients. *The Computer Journal*, 7(2):149–154.
- He, B. and Hofmann, T. (2024). Simplifying Transformer Blocks. arXiv:2311.01906 [cs].
- Hendrycks, D. and Gimpel, K. (2023). Gaussian Error Linear Units (GELUs). arXiv:1606.08415 [cs].
- Hinton, G., Vinyals, O., and Dean, J. (2015). Distilling the Knowledge in a Neural Network. arXiv:1503.02531 [cs, stat].
- Hrycej, T., Bermeitinger, B., Cetto, M., and Handschuh, S. (2023). *Mathematical Foundations of Data Science*. Texts in Computer Science. Springer International Publishing, Cham.
- Krizhevsky, A. (2009). Learning Multiple Layers of Features from Tiny Images. Dataset, University of Toronto.
- LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. (1998). Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324.
- Press, W. H., Teukolsky, S. A., Vetterling, W. T., and Flannery, B. P. (1992). *Numerical recipes in C (2nd ed.): the art of scientific computing*. Cambridge University Press, USA.
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A. C., and Fei-Fei, L. (2015). ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision*, 115(3):211–252.
- Sun, S., Cheng, Y., Gan, Z., and Liu, J. (2019). Patient Knowledge Distillation for BERT Model Compression. arXiv:1908.09355 [cs].
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. (2017). Attention is all you need. In *Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS'17*, pages 6000–6010, Red Hook, NY, USA. Curran Associates Inc.