

---

# Attention-Only Transformers and Implementing MLPs with Attention Heads

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

65 The transformer architecture is widely used in machine learning models and consists  
65 of two alternating sublayers: attention heads and MLPs. We prove that an MLP  
65 neuron can be implemented by a masked attention head with internal dimension 1  
65 so long as the MLP’s activation function comes from a restricted class including  
65 SiLU and close approximations of ReLU and GeLU. This allows one to convert an  
65 MLP-and-attention transformer into an attention-only transformer at the cost of  
65 greatly increasing the number of attention heads.

## 68 1 Introduction

72 The transformer architecture was introduced in the landmark 2017 paper *Attention is All You Need*  
72 (Vaswani et al., 2023) and traditionally consists of alternating attention and multilayer-perceptron  
72 (MLP) sublayers. Although initially used for machine translation, transformers have been used across  
72 a wide range of tasks, including language modeling (Radford et al., 2018; Devlin et al., 2019; Liu  
72 et al., 2018), computer vision (Khan et al., 2022; Cornia et al., 2020), and image generation (Parmar  
72 et al., 2018).

74 This work seeks provide a new perspective on the role of MLP layers in transformers, by proving  
74 that they can be implemented by attention layers. In Theorem 2 we show that by including a “bias  
74 token” akin to the persistent memory vectors in Sukhbaatar et al. (2019) and using a slightly unusual  
74 attention-masking pattern, an MLP layer of size  $\ell$  can be written as the sum of  $\ell$  attention heads  
74 with internal dimension 1. We then show in Theorem 4 that one can apply this process throughout  
74 the entire transformer, converting the typical MLP-and-attention transformer into an attention-only  
74 transformer. Finally, the limitations of this method are discussed.

## 77 2 Background

80 **Notation.** Throughout, we will use  $M_{n,k}$  to denote the set of real-valued  $n$ -by- $k$  matrices. We will  
80 write  $\mathbf{0}$  and  $\mathbf{1}$  for matrices where every entry is 0 or 1, respectively, of size specified or implicit in  
80 the text.

83 For matrices  $X \in M_{n_1,k_1}$  and  $Y \in M_{n_2,k_2}$  of any size, we will write  $X \oplus Y$  for the block matrix  
83 in  $M_{n_1+n_2,k_1+k_2}$  with  $X$  and  $Y$  as diagonal blocks and 0 elsewhere. For matrices  $X \in M_{n,k_1}$  and  
83  $Y \in M_{n,k_2}$ , we will write  $[X|Y] \in M_{n,k_1+k_2}$  for the matrix made by appending one to the other.

85 We write ReLU, SiLU and GeLU for the usual activation functions as in Hendrycks & Gimpel (2023).  
85 In particular,  $\text{SiLU}(x) = x\sigma(x)$ , where  $\sigma(x) = 1/(1 + \exp(-x))$ .

87 We will say that a generalized SiLU function is a function of the form  $f(x) = a_1\text{SiLU}(a_2x)$  for some  
87  $a_1, a_2 \in \mathbb{R}$ .

91 The class of generalized SiLU functions includes  $\text{SiLU}(x)$  and approximations of GeLU and ReLU.  
 91 In particular,  $\text{GeLU}(x) \approx \text{SiLU}(1.702x)/1.702$  (Hendrycks & Gimpel, 2023) (reaching a maximum  
 91 absolute error of 0.0203 at  $x = \pm 2.27$ ) and  $\text{ReLU}(x) \approx \text{SiLU}(kx)/k$  for large  $k$  (reaching a  
 91 maximum absolute error of  $\frac{0.2785}{k}$  at  $x = \pm \frac{1.278}{k}$ ).

95 We will now present a slightly abstracted definition of MLPs, attention heads, and transformers,  
 95 which the reader may confirm encompasses the classical transformer framework described in Vaswani  
 95 et al. (2023).

96 **Definition 1.** An MLP with no biases and one hidden layer is a function  $f : M_{n,k} \rightarrow M_{n,k}$  of the  
 96 form  $f(X) = \alpha(XV_1)V_2$  where  $\alpha : \mathbb{R} \rightarrow \mathbb{R}$  is some real-valued function applied entry-wise to  
 96 matrices, and  $V_1 \in M_{k,\ell}$ ,  $V_2 \in M_{\ell,k}$  are fixed parameter matrices. We call  $\ell$  the size of the hidden  
 96 layer, and the function  $\alpha$  is called the activation function.

98 A mask matrix  $\Lambda$  is a matrix with entries in  $\{0, 1\}$  such that every row has at least one nonzero entry.

100 Let  $X, \Lambda \in M_{n,k}$ , and suppose  $\Lambda$  is a mask matrix. Then define the masked softmax function

$$\text{mssoftmax}(X, \Lambda) := \text{rownorm}(\exp(X) \odot \Lambda)$$

104 where  $\text{rownorm}$  denotes row-wise  $\ell^1$  normalization, and  $\odot$  denotes element-wise multiplication.  
 104 That is, the masked softmax function acts like the usual row-wise softmax but applied to only the  
 104 entries of  $X$  where the mask  $\Lambda$  is 1. At the entries where  $\Lambda$  is 0, the output of the masked softmax  
 104 function takes the value 0.

106 A masked attention head is a function  $h : M_{n,k} \rightarrow M_{n,k}$  of the form

$$h(X) = \text{mssoftmax}(XW_{QK}X^T, \Lambda)XW_{OV}$$

113 for some matrices  $W_{OV}, W_{QK} \in M_{k,k}$ , and mask matrix  $\Lambda \in M_{n,n}$ . We call  $W_{OV}$  and  $W_{QK}$  the  
 113 parameter matrices for this attention head.

115 A transformer is a function  $t : M_{N,D} \rightarrow M_{N,D}$  of the form  $X_0 \mapsto X_1 \mapsto \dots \mapsto X_m = t(X_0)$ , where

$$X_{j+1} = \begin{cases} \text{LayerNorm}(X_j + \sum_i h_{j,i}(X_j)) & \text{or} \\ \text{LayerNorm}(X_j + f_j(X_j)) \end{cases}$$

124 for some attention heads  $h_{j,i}$  or MLPs with one hidden layer  $f_j$ . Note the use of Layer Normalization  
 124 (Ba et al., 2016) and skip connections, where one performs some computation  $f$  on  $X_j$  and defines  
 124  $X_{j+1} = \text{LayerNorm}(X_j + f(X_j))$ , as opposed to  $X_{j+1} = f(X_j)$ .

### 127 3 Implementing MLP Layers with Attention Heads

130 In this section we show that MLP layers whose activation functions are generalized SiLU functions  
 130 are in fact a sum of attention heads.

135 **Theorem 2.** Let  $f(X) = \alpha(XV_1)V_2$  be an MLP on  $M_{N,D}$  with no biases and one hidden layer of  
 135 size  $\ell$ , and suppose  $\alpha$  is a generalized SiLU function  $\alpha(x) = a_1\text{SiLU}(a_2x)$ . Then there are  $\ell$  masked  
 135 attention heads  $\{h_i\}_{i=1}^\ell$  on  $M_{N+1,D+1}$  such that

$$f(X) \oplus [0] = \sum_{i=1}^{\ell} h_i(X \oplus [1])$$

139 for all  $X \in M_{N,D}$ .

141 In particular, for the  $i$ th attention head, one uses parameter and mask matrices

$$\begin{aligned}
W_{QK} &= a_2 \left[ \begin{array}{c|c} \mathbf{0} & -V_1^i \\ \hline \mathbf{0} & 0 \end{array} \right] \\
W_{OV} &= a_1 a_2 V_1^i V_2^i \oplus [0] \\
\Lambda &= \left[ \begin{array}{c|c} I_N & \mathbf{1} \\ \hline \mathbf{0} & 1 \end{array} \right]
\end{aligned}$$

158 where the block decompositions are into size  $N$  and  $1$ ,  $V_1^i$  denotes the  $i$ th column of  $V_1$ ,  $V_2^i$  denotes  
158 the  $i$ th row of  $V_2$ , and  $\mathbf{1}$  denotes the column vector of all  $1$ s.

161 We provide a sketch of the proof in the case of  $\ell = a_1 = a_2 = 1$ . For the full proof see Appendix A.

165 *Proof Sketch.* Since  $\ell = 1$ , we will write  $V_1$  and  $V_2$  are single-column matrices, so we will write  $V_1$   
165 and  $V_2$  in place of  $V_1^i$  and  $V_2^i$ . Due to our choice of a particularly constrained masking pattern, our  
165 masked softmax function will only consider two tokens, the former of which has a pre-attention value  
165 from the main diagonal of  $(X \oplus [1])W_{OV}(X \oplus [1])^T = -XV_1 \oplus [0]$  and the latter of which is  $0$ .  
165 Writing  $-x$  for the former entry, we have  $\text{softmax}([-x, 0]) = \text{rownorm}([e^{-x}, 1]) = [\sigma(x), \sigma(-x)]$ .

171 That is, by our choice of  $W_{QK}$  and  $\Lambda$  we have made our head have attention pattern  
171  $\left[ \begin{array}{c|c} \text{diag}(\sigma(XV_1)) & \sigma(-XV_1) \\ \hline \mathbf{0} & 1 \end{array} \right]$ . Then, the complete output of this attention head  $h$  is

$$\begin{aligned}
h(X \oplus [1]) &= \left[ \begin{array}{c|c} \text{diag}(\sigma(XV_1)) & \sigma(-XV_1) \\ \hline \mathbf{0} & 1 \end{array} \right] \left[ \begin{array}{c|c} X & \mathbf{0} \\ \hline \mathbf{0} & 1 \end{array} \right] \left[ \begin{array}{c|c} V_1 V_2 & \mathbf{0} \\ \hline \mathbf{0} & 0 \end{array} \right] \\
&= \left[ \begin{array}{c|c} \text{diag}(\sigma(XV_1))XV_1 V_2 & \mathbf{0} \\ \hline \mathbf{0} & 0 \end{array} \right] \\
&= \text{SiLU}(XV_1)V_2 \oplus [0]
\end{aligned}$$

195 as desired.

196 □

201 **Remark 3.** The additional term  $\oplus[1]$  in Theorem 2 is similar to the persistent vectors of Sukhbaatar  
201 et al. (2019). In that work, the authors propose a new architecture, which they call the all-attention  
201 architecture, in which attention can also be paid to certain static vectors, learned for each attention  
201 head, called the persistent vectors. Our approach could also be implemented in that architecture with  
201 a single persistent vector  $(0, 0, 0, \dots, 0, 1)$  shared across all attention heads in place of the  $\oplus[1]$  terms.

203 Note also that the  $W_{QK}$  and  $W_{OV}$  matrices used in Theorem 2 can be factored into the matrices  $W_Q$ ,  
203  $W_K$ ,  $W_V$ ,  $W_O \in M_{D+1,1}$  from Vaswani et al. (2023) satisfying  $W_{QK} = W_Q W_K^T / \sqrt{D+1}$  and  
203  $W_{OV} = W_V W_O$ . In particular, we can take  $W_Q = W_V = a_2[V_1^i | 0]^T$ ,  $W_K = \sqrt{D+1}[\mathbf{0} | -1]^T$ ,  
203 and  $W_O = a_1[V_2^i | 0]^T$ . Since  $W_K$  is shared across all attention heads, we only need to store two sets  
203 of parameters, the vectors  $W_Q = W_V$  and  $W_O$ .

205 This provides an alternative perspective on MLP neurons: a neuron in an MLP is an attention head  
205 with internal dimension  $1$  and a particularly restrictive masking pattern in which each token attends  
205 only to itself and a static “bias” token.

208 We now turn to have the necessary tools to show that a decoder-only transformer as in Liu et al.  
208 (2018); Radford et al. (2018) can be implemented entirely with attention heads.

212 **Theorem 4.** If a transformer’s MLP layers are activated by a generalized SiLU function, they can be  
212 substituted with attention heads.

215 We again provide just a sketch of the proof and direct the reader to Appendix A for the full proof.

219 *Proof Sketch.* We will create a new transformer on  $M_{N+1, D+1}$  whose residual stream  $X'_j$  on every  
219 sublayer satisfies  $X'_j = X_j \oplus [1]$ . This is sufficient to prove the main claim since the output of

219 this new transformer will be  $X'_{2m} = X_{2m} \oplus [1]$  and therefore contain the output of the original  
219 transformer.

221 For any MLP layer in the original transformer, we use Theorem 2 to replace the MLP layer with  
221 attention heads. For any attention head layer, we can slightly augment the  $W_{QK}$ ,  $W_{OV}$ , and  $\Lambda$   
221 matrices to work on the larger size. Due to skip connections, the resulting matrix retains the  $\oplus[1]$   
221 term, as desired.  $\square$

225 **Remark 5.** *It is instructive to compare this construction to the negative results of Dong et al. (2021),*  
225 *which find that without skip connections or MLPs, a self-attention network converges rapidly to a*  
225 *rank-1 matrix. Since we obviously do away with the MLP layer, our result depends on the use of skip*  
225 *connections. In particular, the “bias term” of  $\oplus[1]$  is zeroed out by the construction in Theorem 2,*  
225 *so applying the construction in Theorem 4 without a skip connection results in  $X'_0 = X_0 \oplus [1]$ , but*  
225  *$X'_1 = X_1 \oplus [0]$ . Then, in the  $j = 2$  sublayer, the construction in Theorem 2 would fail for lack of this*  
225 *bias term, as, without it, the pre-attention matrix  $(X')W_{QK}(X')^T$  is 0.*

228 We additionally show in Appendix B that attention heads can separately implement the components  
228 of an MLP layer, namely activation functions and linear transformations.

## 231 4 Limitations

234 The technique described in Theorem 4 faces several practical limitations. First is the quantity of  
234 attention heads: we use one attention head per dimension of the hidden layer, which can easily  
234 increase the number of attention heads by several orders of magnitude, partially offset by the new  
234 attention heads having smaller internal dimension. For example, in GPT-3 (Brown et al., 2020) the  
234 MLP layer has hidden dimension 49152, so this method would require 49152 additional 1-dimensional  
234 attention heads in each layer. This is an increase from from GPT-3’s normal set of 96 attention heads  
234 per layer, each with internal dimension 128.

236 Second, it may be the case that replacing a feedforward network with attention heads slows down  
236 model inference or training. In particular, this approach replaces matrix multiplication with many  
236 vector-by-vector multiplications. One also computes many terms that are “thrown away” in the  
236 masking step. Combined, these suggest that converting an MLP layer to attention heads might  
236 increase computational costs.

## 239 5 Discussion

242 We have proven that attention heads can implement an MLP layer and that any transformer can be  
242 converted to an attention-only transformer. This approach provides a useful new perspective on the  
242 relative importance of MLP layers and attention heads in language models. MLP layers in a model  
242 like GPT-3 are larger than attention layers by a 2:1 margin if one measures by number of parameters  
242 but by 500:1 if one measures by number of attention heads.

244 One implication of these results is that it is theoretically possible to train an attention-only transformer  
244 that matches the performance of an MLP-plus-attention transformer. It remains unknown whether  
244 such an architecture would be competitive with the more classical transformer architecture in terms  
244 of practical considerations like training or inference speed. Such a test would be a promising future  
244 area of research.

## 1 References

- 11 Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization, 2016.
- 24 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,  
24 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel  
24 Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler,  
24 Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott  
24 Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya  
24 Sutskever, and Dario Amodei. Language models are few-shot learners, 2020.

- 31 Marcella Cornia, Matteo Stefanini, Lorenzo Baraldi, and Rita Cucchiara. Meshed-memory trans-  
31 former for image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision*  
31 *and Pattern Recognition (CVPR)*, June 2020.
- 36 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep  
36 bidirectional transformers for language understanding, 2019.
- 43 Yihe Dong, Jean-Baptiste Cordonnier, and Andreas Loukas. Attention is not all you need: Pure  
43 attention loses rank doubly exponentially with depth. In *International Conference on Machine*  
43 *Learning*, pp. 2793–2803. PMLR, 2021.
- 47 Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus), 2023.
- 57 Salman Khan, Muzammal Naseer, Munawar Hayat, Syed Waqas Zamir, Fahad Shahbaz Khan, and  
57 Mubarak Shah. Transformers in vision: A survey. *ACM Comput. Surv.*, 54(10s), sep 2022. ISSN  
57 0360-0300. doi: 10.1145/3505244. URL <https://doi.org/10.1145/3505244>.
- 63 Peter J. Liu, Mohammad Saleh, Etienne Pot, Ben Goodrich, Ryan Sepassi, Lukasz Kaiser, and Noam  
63 Shazeer. Generating wikipedia by summarizing long sequences, 2018.
- 74 Niki Parmar, Ashish Vaswani, Jakob Uszkoreit, Lukasz Kaiser, Noam Shazeer, Alexander Ku, and  
74 Dustin Tran. Image transformer. In Jennifer Dy and Andreas Krause (eds.), *Proceedings of the 35th*  
74 *International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning*  
74 *Research*, pp. 4055–4064. PMLR, 10–15 Jul 2018. URL <https://proceedings.mlr.press/v80/parmar18a.html>.
- 80 Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language  
80 understanding by generative pre-training. *OpenAI blog*, 2018.
- 86 Sainbayar Sukhbaatar, Edouard Grave, Guillaume Lample, Herve Jegou, and Armand Joulin. Aug-  
86 menting self-attention with persistent memory, 2019.
- 92 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz  
92 Kaiser, and Illia Polosukhin. Attention is all you need, 2023.

## 259 A Proofs of Main Results

263 In this section we present fully detailed proofs of our main results.

266 *Proof of Theorem 2.* We first prove the claim in the case of  $\ell = a_1 = a_2 = 1$ . In this case, since  
266 there is only one column in  $V_1$ , then  $V_1 = V_1^i$ , and similarly  $V_2 = V_2^i$ . Consider the attention matrix  
266  $\text{msoftmax}((X \oplus [1])W_{QK}(X \oplus [1])^T, \Lambda)$ . Multiplying matrices on the level of their blocks, we get  
266 that the first argument of the masked softmax is

$$(X \oplus [1])W_{QK}(X \oplus [1])^T = \begin{bmatrix} X & \mathbf{0} \\ \mathbf{0} & 1 \end{bmatrix} \begin{bmatrix} \mathbf{0} & -V_1^i \\ \mathbf{0} & 0 \end{bmatrix} \begin{bmatrix} X & \mathbf{0} \\ \mathbf{0} & 1 \end{bmatrix}^T = \begin{bmatrix} \mathbf{0} & -XV_1 \\ \mathbf{0} & 0 \end{bmatrix}$$

286 Now consider the masked softmax term in the  $j$ th row for  $j \leq N$ . This row has exactly two unmasked  
286 values, the diagonal entry and the rightmost entry, taking the values 0 and  $-(XV_1)_j$ , respectively.  
286 Applying exp and rownorm results in  $\sigma((XV_1)_j)$  and  $\sigma(-(XV_1)_j)$ , respectively. Thus, the masked  
286 softmax term becomes

$$\begin{aligned} \text{msoftmax}((X \oplus [1])W_{QK}(X \oplus [1])^T, \Lambda) &= \text{msoftmax}\left(\begin{bmatrix} \mathbf{0} & -XV_1 \\ \mathbf{0} & 0 \end{bmatrix}, \begin{bmatrix} I_{n-1} & \mathbf{1} \\ \mathbf{0} & 1 \end{bmatrix}\right) \\ &= \begin{bmatrix} \text{diag}(\sigma(XV_1)) & \sigma(-XV_1) \\ \mathbf{0} & 1 \end{bmatrix} \end{aligned}$$

306 Substituting these values into the expression for  $h(X)$  gives

$$\begin{aligned}
h(X \oplus [1]) &= \text{msoftmax}((X \oplus [1])W_{QK}(X \oplus [1])^T, \Lambda)(X \oplus [1])W_{OV} \\
&= \left[ \begin{array}{c|c} \text{diag}(\sigma(XV_1)) & \sigma(-XV_1) \\ \mathbf{0} & 1 \end{array} \right] (X \oplus [1])W_{OV} \\
&= \left[ \begin{array}{c|c} \text{diag}(\sigma(XV_1)) & \sigma(-XV_1) \\ \mathbf{0} & 1 \end{array} \right] \left[ \begin{array}{c|c} X & \mathbf{0} \\ \mathbf{0} & 1 \end{array} \right] \left[ \begin{array}{c|c} V_1V_2 & \mathbf{0} \\ \mathbf{0} & 0 \end{array} \right] \\
&= \left[ \begin{array}{c|c} \text{diag}(\sigma(XV_1))XV_1V_2 & \mathbf{0} \\ \mathbf{0} & 0 \end{array} \right] \\
&= \left[ \begin{array}{c|c} \text{SiLU}(XV_1)V_2 & \mathbf{0} \\ \mathbf{0} & 0 \end{array} \right] \\
&= \left[ \begin{array}{c|c} f(X) & \mathbf{0} \\ \mathbf{0} & 0 \end{array} \right] \\
&= f(X) \oplus [0]
\end{aligned}$$

345 as desired. This completes the  $\ell = a_1 = a_2 = 1$  case.

347 For a general  $a_1, a_2$ , apply the previous case to an MLP with weight matrices  $a_2V_1$  and  $a_1V_2$ .

349 Finally, for the fully general case with  $\ell > 1$ , for each  $1 \leq i \leq \ell$ , let  $f_i(X) = \alpha(XV_1^i)V_2^i$ , and note

349 that  $f = \sum_{i=1}^{\ell} f_i$ . Let  $h_i$  denote the attention head corresponding to  $f_i$  given by the  $\ell = 1$  case. Then

349 we have that

$$\begin{aligned}
f(X) \oplus [0] &= \sum_{i=1}^{\ell} f_i(X) \oplus [0] \\
&= \sum_{i=1}^{\ell} h_i(X \oplus [1])
\end{aligned}$$

355 as desired. □

359 *Proof of Theorem 4.* We will show that we can create a new transformer  $t'$  on  $M_{N+1, D+1}$  whose  
359 residual stream  $X'_j$  on every sublayer satisfies

$$X'_j = X_j \oplus [1]$$

363 This is sufficient to prove the main claim since the output of this new transformer will be  $X'_{2m} =$   
363  $X_{2m} \oplus [1]$  and therefore contain the output of the original transformer.

365 Without loss of generality, assume that the MLP layers have no bias terms (i.e., that we've already  
365 used the "bias trick" to fold bias terms into the weight matrix).

367 To prove that there is a transformer  $t'$  that satisfies  $X'_j = X_j \oplus [1]$  on every sublayer, we proceed by  
367 induction. For the base case of  $j = 0$ , we tweak the transformer's context window and embedding  
367 weights so that  $X'_0 = X_0 \oplus [1]$ .

369 We split the inductive case depending on whether the original transformer's sublayer used attention  
369 or an MLP. If the original layer was an MLP, then by Theorem 2 there are attention heads  $h'_{j,i}$  such  
369 that  $f_j(X) \oplus [0] = \sum h'_{j,i}(X \oplus [1])$ , so in our transformer  $t'$ , using these attention heads yields

$$\begin{aligned}
X'_{j+1} &= \text{LayerNorm}(X'_j + \sum h'_{j,i}(X'_j)) \\
&= \text{LayerNorm}((X_j \oplus [1]) + \sum h'_{j,i}(X_j \oplus [1])) \\
&= \text{LayerNorm}((X_j \oplus [1]) + (f_j(X) \oplus [0])) \\
&= \text{LayerNorm}(X_j + f_j(X)) \oplus [1] \\
&= X_{j+1} \oplus [1]
\end{aligned}$$

379 as desired.

381 If instead, the transformer used attention heads on the  $j$ th sublayer, we must tweak our original  
381 attention heads to account for the new size. To this end, we will show that for each of the original  
381 attention heads  $h = h_{j,i}$ , we can create an attention head  $h'$  such that

$$h'(X \oplus [1]) = h(X) \oplus [0]$$

385 Let  $W_{QK}$ ,  $W_{OV}$ , and  $\Lambda$  denote the original parameter and masking matrices for  $h$ . Then define

$$\begin{aligned}
W'_{QK} &= W_{QK} \oplus [1] \\
W'_{OV} &= W_{OV} \oplus [0] \\
\Lambda' &= \Lambda \oplus [1]
\end{aligned}$$

393 Then,

$$\begin{aligned}
h'(X \oplus [1]) &= \text{msoftmax}((X \oplus [1])W'_{QK}(X \oplus [1])^T, \Lambda')(X \oplus [1])W'_{OV} \\
&= \text{msoftmax}((X \oplus [1])(W_{QK} \oplus [1])(X \oplus [1])^T, (\Lambda \oplus [1]))(X \oplus [1])(W_{OV} \oplus [0]) \\
&= \text{msoftmax}(XW_{QK}X^T \oplus [1], \Lambda \oplus [1])(XW_{OV} \oplus [0]) \\
&= (\text{msoftmax}(XW_{QK}X^T, \Lambda) \oplus [1])(XW_{OV} \oplus [0]) \\
&= \text{msoftmax}(XW_{QK}X^T, \Lambda)XW_{OV} \oplus [0] \\
&= h(X) \oplus [0]
\end{aligned}$$

404 as desired. Now, creating such  $h'_{j,i}$  for each of the original attention heads  $h_{j,i}$ , we have

$$\begin{aligned}
X'_{j+1} &= \text{LayerNorm}(X'_j + \sum h'_{j,i}(X'_j)) \\
&= \text{LayerNorm}((X_j \oplus [1]) + \sum h'_{j,i}(X_j \oplus [1])) \\
&= \text{LayerNorm}((X_j \oplus [1]) + \sum h_{j,i}(X) \oplus [0])) \\
&= \text{LayerNorm}((X_j + \sum h_{j,i}(X))) \oplus [1] \\
&= X_{j+1} \oplus [1]
\end{aligned}$$

414 as desired. This completes the inductive step and the proof.

415 □

## 418 B Linear Transformations and Activation Functions with Attention Heads

422 Theorem 2 shows that attention heads can implement an MLP layer, but can they separately implement  
422 the components of an MLP, a linear transformation and an activation function? In this section we  
422 show that the answer is yes.

424 We first show that an attention head can perform an arbitrary linear operation row-wise on the matrix.

428 **Theorem 6.** Let  $h : M_{N,D} \rightarrow M_{N,D}$  be an attention head with masking matrix  $\Lambda = I_N$ . Then  
 428  $h(X) = XW_{OV}$ .

432 *Proof.* Because  $\Lambda = I_n$ , after masking, the attention matrix  $\text{msoftmax}(XW_{QK}X^T, \Lambda)$  will have  
 432 nonzero entries only along the diagonal. Since the rows of the attention matrix are normalized to sum  
 432 to 1, it follows that  $\text{msoftmax}(XW_{QK}X^T, \Lambda) = I_n$ . Then,

$$h(X) = \text{msoftmax}(XW_{QK}X^T, \Lambda)XW_{OV} = I_nXW_{OV} = XW_{OV}$$

434 as desired. □

437 Now we will show that one can apply a generalized SiLU function entrywise.

441 **Theorem 7.** Let  $\alpha$  be a generalized SiLU function. Then there are  $D$  attention heads  $h_1, \dots, h_D$  on  
 441  $M_{N+1,D+1}$  such that

$$\alpha(X) \oplus [0] = \sum_{i=1}^D h_i(X \oplus [1])$$

448 *Proof.* This follows immediately from applying Theorem 2 to the MLP  $f(X) = \alpha(XI_N)I_N =$   
 448  $\alpha(X)$ , whose hidden layer is of size  $\ell = D$ . □

451 Note that a transformer usually makes use of skip connections, so that the residual stream experiences  
 451 the transformation  $X \mapsto X + \text{sublayer}(X)$ . Thus, to get the transformation  $X \mapsto \alpha(X)$ , one can  
 451 combine these two theorems, using  $D + 1$  attention heads to produce  $\text{sublayer}(X) = \alpha(X) - X$ , in  
 451 which case  $X \mapsto X + \text{sublayer}(X) = \alpha(X)$ .