
ConcatPlexer : Additional Dim1 Batching for Faster ViTs

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Transformers have demonstrated tremendous success not only in the natural lan-
2 guage processing (NLP) domain but also the field of computer vision, igniting
3 various creative approaches and applications. Yet, the superior performance and
4 modeling flexibility of transformers came with a severe increase in computation
5 costs, and hence several works have proposed methods to reduce this burden. In-
6 spired by a cost-cutting method originally proposed for language models, Data
7 Multiplexing (DataMUX), we propose a novel approach for efficient visual recog-
8 nition that employs additional dim1 batching (*i.e.*, concatenation) that greatly
9 improves the throughput with little compromise in the accuracy. We first introduce
10 a naive adaptation of DataMux for vision models, Image Multiplexer, and devise
11 novel components to overcome its weaknesses, rendering our final model, Concat-
12 Plexer, at the sweet spot between inference speed and accuracy. The ConcatPlexer
13 was trained on ImageNet1K and CIFAR100 dataset and it achieved 23.5% less
14 GFLOPs than ViT-B/16 with 69.5% and 83.4% validation accuracy, respectively.

15 1 Introduction

16 Deep learning research community has experienced dazzling advances in model performance across
17 a wide variety of domains and downstream tasks in the last decade [1, 2, 3, 4, 5, 6, 7]. These improve-
18 ments, however, came at the cost of rapidly increasing computational burden, with the introduction
19 of Transformer [2, 3, 8] marking a major milestone in this aspect. With the growing popularity
20 of transformers, methods to reduce their computational costs have become a prominent research
21 topic [9, 10, 11, 12, 13, 14]. However, previous efforts to improve the computational efficiency of
22 transformers have been mostly focused on the NLP domain. Data multiplexing (DataMUX) [15]
23 pioneered this direction of research for language models by projecting multiple input tokens into a
24 single compact representation space and thus enabling the neural network to process them simultane-
25 ously. Although DataMUX [15] has delivered promising preliminary results for the concept of data
26 multiplexing, there is much room for research remaining unexplored especially in the vision domain.
27 For instance, it has mainly trained the transformer on the GLUE benchmark and as a CV task, the
28 authors have only experimented on the MNIST dataset with light Multi-Layer Perceptron (MLP) and
29 Convolutional Neural Network (CNN). This experimental setting is at best a proof-of-concept and
30 thus insufficient to ensure its general applicability in the vision domain.

31 In this paper, we explore the potential of data multiplexing in larger scale general vision applications
32 such as ImageNet1K [16] classification. To that end, we first show the limitations of naive adaptation
33 of DataMUX by constructing a simple baseline named Image Multiplexer that employs DataMUX
34 for visual recognition with minimal modifications. We then progressively transform this architecture
35 to reach a favorable trade-off between the accuracy and inference speed, presenting our final model,
36 ConcatPlexer. ConcatPlexer, in short, is a method for efficiently extracting multiple images' represen-

37 tation at once. ConcatPlexer extracts high-level feature tokens via Transformer encoder layers. Then
 38 instead of projecting multiple inputs to a compact representation space, our ConcatPlexer reduces the
 39 length of input tokens using a learned convolution and concatenates them for simultaneous processing.
 40 Comparison with the naive Image Multiplexer clearly demonstrates that DataMUX in its native form
 41 is ill-suited for vision models but can be made effective with our proposed modifications.

42 The ConcatPlexer and its MultiPlexer baseline are pretrained and compared on ImageNet1K [16],
 43 making them generally applicable vision frameworks. We further finetune them on CIFAR100 [17]
 44 and evaluate the result. For evaluation of this new framework, we suggest a "*multiplexed image*
 45 *classification task*", whose goal is to classify multiple images multiplexed into a single representation.
 46 ConcatPlexer achieves consistent gains over its baselines in both ImageNet and CIFAR100, supporting
 47 its effectiveness in visual recognition tasks.

48 The contribution of this paper is as follows:

- 49 1. This paper defines the "multiplexed image classification task" and deal with the concept of
 50 data multiplexing that projects multiple inputs into a single representation for efficient data
 51 processing in the vision domain.
- 52 2. We propose the ConcatPlexer, a novel framework for multiplexing images, and test its
 53 performance on ImageNet1K and CIFAR100 benchmark. The ConcatPlexer extracts high-
 54 level featured tokens using the transformer encoder patchifier and concatenates multiple
 55 images to process them at once.
- 56 3. We demonstrate that data multiplexing can obtain a favorable trade-off between throughput
 57 and accuracy. Our model can save up to 66.9% of FLOPs compared to ViT-B/16 with mild
 58 drop in accuracy.

59 2 Related Work

60 **Data Multiplexing:** The concept of data multiplexing was first suggested by DataMUX [15]. The
 61 DataMUX processes multiple inputs by projecting multiple texts into a single compact representation
 62 space. This enables models to process much larger batch with a same GPU resource. DataMUX
 63 and its latter version MUX-PLM [18] demonstrate their performance on GLUE Benchmark [19]. In
 64 this work, our proposed ConcatPlexer gains model efficiency by transplanting the data multiplexing
 65 method into the vision domain successfully.

66 **Token Reduction:** Though we are the first to apply the concept of data multiplexing to vision
 67 domain to the best of our knowledge, there are studies that try to cut the computational cost of
 68 transformer-based models by reducing the number of input tokens. ToMe [14] merges similar tokens
 69 to reduce the length of the input sequence. The other approaches [13, 12] prune tokens into a single
 70 token to reduce the length of an input sequence. However, our method processes multiple inputs at
 71 the same time, naturally reducing the computational cost.

72

73 3 Method

74 3.1 Preliminary: DataMUX

75 **Multiplexing:** To project multiple inputs into a single compact representation space, a multiplexing
 76 module is used. Consider (x^1, \dots, x^N) with $x^i \in R^d$ being a tuple of N inputs within a batch.
 77 Multiplexing transforms each input by $\phi^i : R^d \mapsto R^d$ and averages at the end. A backbone takes a
 78 batch (x^1, \dots, x^N) as an input and outputs the multiplexed hidden representation output $h^{1:N}$.

$$h^{1:N} = \Phi(x^1, \dots, x^N) = \frac{1}{N} \sum_{n=1}^N \phi^i(x^i). \quad (1)$$

79 Considering the case for a sequenced token input with length L , the aforementioned multiplexing
 80 process is done in a token-wise order. For an input sequence $x^i = \{w_j^i\}_{j \in [L]}$, each token w_j can be
 81 processed as

$$h_j^{1:N} = \Phi(w_j^1, \dots, w_j^N). \quad (2)$$

82 In DataMUX, either (1) a random fixed orthogonal matrix or (2) a fixed Gaussian random matrix is
 83 used for the linear projector ϕ^i . In our naively implemented Image Multiplexer, a fixed orthogonal
 84 matrix is used.
 85

86 **Demultiplexing:** To disentangle the multiplexed hidden representation output $h^{1:N}$ into N indepen-
 87 dent representation, demuxing module is used. Demultiplexing function Θ^i extracts i^{th} representation
 88 from a multiplexed hidden representation as

$$y^i = \Theta^i(h^{1:N}), \forall i \in [N]. \quad (3)$$

89 For demultiplexing function, there are two choices: (1) MLP Demuxing: using N MLPs to extract
 90 N representation from $h^{1:N}$ and (2) Index Embedding. An MLP Demuxing is used for our naively
 91 implemented Image Multiplexer.

92 **Theoretical claim of DataMUX:** The major factor that Transformer based models can handle the
 93 multiplexed task is Transformer’s multi-head attention mechanism. Each multi-head attention can
 94 extract features of different inputs within the multiplexed representation. Kindly refer to DataMUX
 95 [15] for more detailed theoretical claim.

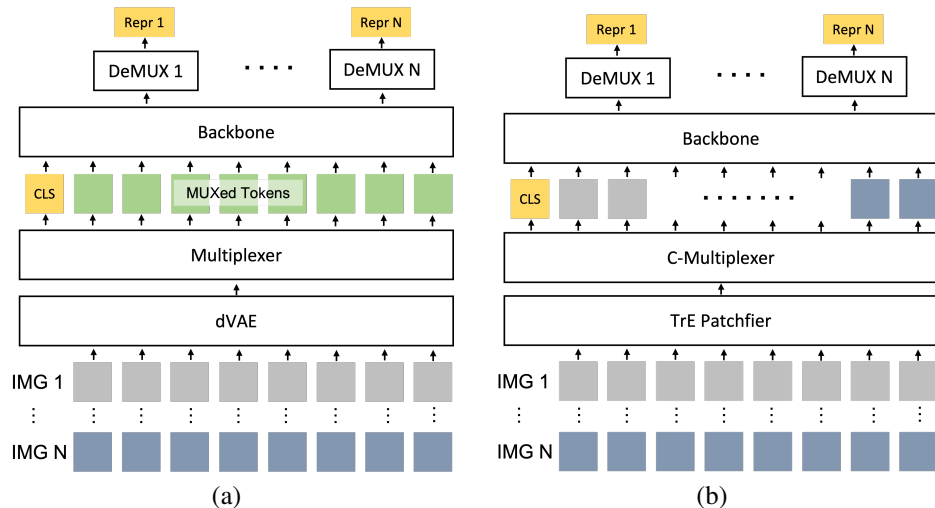


Figure 1: Overall architecture of (a) Image Multiplexer and (b) ConcatPlexer. The Image Multiplexer multiplexes N_{MUX} images using MLP and fixed orthogonal matrices. The ConcatPlexer uses a conv layer to reduce the length of each image token and concatenates them. N_{MUX} is abbreviated as N in this figure.

96 3.2 ConcatPlexer

97 We propose the ConcatPlexer (Figure. 1-b), a model that successfully adapts the concept of DataMUX
 98 to the vision domain by addressing the structural differences between the two modalities. As explained
 99 in [20], the most decisive difference between visual signal and natural language lies in data redundancy.
 100 A pixel that constitutes an image rarely carries significant information by itself while a word token
 101 more likely carries important semantic information. In order to compensate for this difference and
 102 suit DataMUX for vision tasks, we compose our ConcatPlexer with the following architectural
 103 components: Transformer patchifier, ConcatMultiplexer, Demultiplexer, and the Backbone.

104 **Transformer Encoder Patchifier:** To address the redundancy issue of pixel-based data, we extract
 105 high-level features before feeding to the multiplexing backbone. High-level featured tokens will
 106 reduce redundancy. This will make the backbone process and distinguish the multiple inputs parallelly.
 107 Transformer Encoder(TrE) Patchifier is a stacked transformer encoder layer with CNN layer at the
 108 front. Suppose that the dimension of input images is $(bs, 3, W, H)$ where bs is the batch size, 3 is a
 109 color channel, W is the width of an image, and H is the height of an image. First, the CNN layer
 110 patchifies the image into a grid patch turning the input dimension into (bs, L, dim) where bs is the

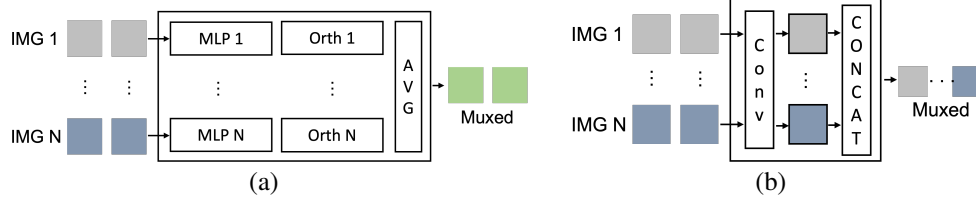


Figure 2: The architecture of (a) Multiplexer and (b) C-Multiplexer. Both inputs N_{MUX} of inputs and combine them into a single input. N_{MUX} is N for this figure.

111 batch size, L is the token length of each image, and dim is the dimension of each token. Then rear
 112 TrE will turn the (bs, L, dim) tokens into a high-level featured token while retaining the size of input
 113 the same.

114 **C-Multiplexer:** While training the Image Multiplexer, we observed that the existing multiplexing
 115 method’s performance degradation outweighs the efficiency gain for the tasks that have low expecta-
 116 tions on random chance. We optimize the trade-off between computational benefit and performance
 117 degradation by preventing severe performance degradation while retaining computational efficiency
 118 to a certain degree. Instead of projecting multiple(N_{MUX}) image tokens into a single compact
 119 representation space, we tried to extract the essence of each image and combine them in a different
 120 manner. In the Figure. 2-b, to extract the essence of each image, conv computation was used on
 121 high-level featured tokens. Using the conv1d layer with the output channel dim , the dimension of $(bs,$
 122 $L, dim)$ tokens from TrE tokenizer becomes $(bs, L/N_{MUX}, dim)$ where N_{MUX} is number of sample
 123 to multiplex. From this, each image gets shorter in length while retaining necessary information
 124 as much as possible. Then we concatenated the tokens of N_{MUX} images to train the backbone to
 125 process N_{MUX} images at the same time and store N_{MUX} representation in a single CLS token.
 126 From this operation, the input of dimension $(bs, L/N_{MUX}, dim)$ becomes $(bs/N_{MUX}, L, dim)$
 127 thereby enables the model to process N_{MUX} times larger batch. As C-Multiplexer is very simple, the
 128 computational overhead is negligible.

129 **Demultiplexer and backbone:** For the Demultiplexer and backbone, the ConcatPlexer uses the same
 130 Demultiplexer and backbone as the Image Multiplexer. The backbone is a ViT-like architecture that
 131 stacks the transformer encoder layers. The backbone takes $(bs/N_{MUX}, L, dim)$ dimension tokens
 132 as an input and outputs $(bs/N_{MUX}, L + 1, dim)$ tokens including the CLS token. N_{MUX} of MLPs
 133 were initialized to separate the representation of each image from a single CLS token of the backbone.
 134 For more detail, refer to Sec. 4.1.

135 3.3 Training

136 To train ConcatPlexer, three loss terms were used. Firstly, classification loss using ground truth class
 137 label was used. To boost the performance of ConcatPlexer, CLIP loss and Label smoothing loss were
 138 used.

139 **CLIP Loss:** The CLIP loss is intended to train the ConcatPlexer’s demultiplexed CLS output to
 140 resemble the representation of CLIP vision encoder. By encouraging the model to learn the general
 141 representation space of the CLIP encoder, CLIP loss can prevent the model from overfitting and
 142 blindly memorizing the ground truth (GT) label. The CLS_x in Eq. 4 is demultiplexed CLS token of
 143 an image x . $CLIP(\cdot)$ is a CLIP vision encoder that outputs feature token of image x . The similarity
 144 between the two features is calculated by the contrastive loss.

$$\mathcal{L}_{CLIP} = Ctrs(CLIP(x), CLS_x). \quad (4)$$

145 **Label Smoothing Loss:** In order to take advantage of the ConcatPlexer’s multiplexed input, we
 146 augmented the image by mixing other high-level image tokens within the multiplexed sample at

147 C-Multiplexer. Tokens of N_{MUX} images are averaged as follows:

$$\begin{aligned}
 M_i &= \sum_{n=1}^{N_{MUX}} f(T_n), \quad f(T_n) = \begin{cases} \alpha * T_n, & \text{if } n = i \\ \frac{(1-\alpha)}{N_{MUX}-1} * T_n, & \text{otherwise} \end{cases} \\
 M_i^{GT} &= \sum_{n=1}^{N_{MUX}} g(Y_n), \quad g(Y_n) = \begin{cases} \alpha * Y_n, & \text{if } n = i \\ \frac{(1-\alpha)}{N_{MUX}-1} * Y_n, & \text{otherwise} \end{cases} \\
 \mathcal{L}_{smooth} &= CE(M_i, M_i^{GT}),
 \end{aligned} \tag{5}$$

148 where T_n and Y_n are a n^{th} high-level featured tokens among N_{MUX} images after TrE tokenizer and
 149 ground truth class label of a $image_n$ in one-hot vector format, respectively. As a result, M_i^{GT} is a
 150 smoothed one-hot vector label, whose i -th component is set to α , and $(1 - \alpha)$ is distributed to other
 151 components corresponding to the labels of other images. This prevents the model from overfitting the
 152 training dataset and shows a slight performance gain.

153 4 Experiment

154 4.1 Baseline: Image Multiplexer

155 Along with ConcatPlexer, we introduce Image Multiplexer (Figure. 1-a), a naively implemented
 156 version of DataMUX in the vision domain, as a baseline. Unlike DataMUX in NLP, the Image
 157 Multiplexer has a long way to go due to several structural differences in vision. For the training,
 158 classification loss and token retrieval loss were used. For the token retrieval loss, the model is
 159 trained to restore the original discrete input tokens. This helped boost the performance of the original
 160 DataMUX. Implementation detail is described in the following section.

161 **Image Multiplexer:** To bring the DataMUX into a vision regime, Image Multiplexer can be broken
 162 down into four parts: (1) Discrete patchifier, (2) Multiplexer, (3) Backbone, and (4) Demultiplexer.

163 **Discrete Patchifier:** NLP inputs are tokenized into discrete tokens. Original DataMUX exploits this
 164 nature with token retrieval task. The discrete patchifier is used to make pixel patched into discrete
 165 tokens. Specifically, DALL-E’s pretrained discrete variational autoencoder (dVAE) [21] was used.
 166 DALL-E’s dVAE patchifies 8x8 pixels into a single discrete 13-bit code (total 8192 codes). This
 167 enables the model to be trained with token retrieval loss.

168 **Multiplexer and Demultiplexer:** The Multiplexer multiplexed (Figure. 2-a) N_{MUX} of discretized
 169 images into a single muxed input. For the Multiplexing module, N_{MUX} of shallow MLPs and random
 170 fixed orthogonal matrices were used. Each discretized image is projected with a shallow MLP and an
 171 orthogonal matrix. Then N_{MUX} of projected representations are averaged to be multiplexed into a
 172 single compact representation space. For the Demultiplexing module, N_{MUX} of shallow MLPs were
 173 used. Each MLP is trained to extract the representation of each image from muxed representation
 174 output of the backbone.

175 The structure of the Multiplexer is the major difference between Image Multiplexer and the Concat-
 176 Plexer. Image Multiplexer projects and combines via linear projection and fixed orthogonal matrices
 177 while C-Multiplexer uses conv computation to reduce tokens of N_{MUX} images and concatenates
 178 them in a single sequence.

179 **Backbone:** The ViT-like architecture was used for the Image Multiplexer backbone. The backbone
 180 shares the same configuration as the ViT-base model [3]. 12 transformer encoder layers were stacked
 181 and the representation dimension was 768.

182 4.2 Experimental Detail

183 For a multiplexed image classification task, Image Multiplexer and ConcatPlexer are trained on
 184 ImageNet1K and CIFAR100 datasets. Both models were trained using an AdamW optimizer with a
 185 learning rate of 1e-4 and weight decay of 0.03 for ImageNet1K dataset. Each model was trained around
 186 50 epochs until it converged on the training set. As shown in table 1 and mentioned before, Image
 187 Multiplexer uses DALL-E [21] tokenizer or CNN patchifier, and ConcatPlexer uses a transformer
 188 encoder (TrE) as a high-level featured tokenizer. The N_{MUX} on table 1 means the number of samples

Name	Tokenizer	Token Length	N_{MUX}	Concat Point	Val Acc
Image Multiplexer	DALL-E	784	2	-	54%
Image Multiplexer	DALL-E	784	4	-	48%
Image Multiplexer	CNN	196	4	-	26%
ConcatPlexer(1)	TrE	196	2	1	62.3%
ConcatPlexer(2)	TrE	196	2	3	65.3%
ConcatPlexer(3)	TrE	196	2	6	69.5%
ConcatPlexer(4)	TrE	196	4	1	56.7%

Table 1: Performance of the Multiplexed Models on ImageNet1K.

Name	Tokenizer	Token Length	N_{MUX}	Concat Point	Val Acc
Image Multiplexer	DALL-E	784	2	-	74%
Image Multiplexer	DALL-E	784	4	-	70%
ConcatPlexer(1)	TrE	196	2	1	76.1%
ConcatPlexer(2)	TrE	196	2	3	78.6%
ConcatPlexer(3)	TrE	196	2	6	83.4%

Table 2: Performance of the Multiplexed Models on CIFAR100.

189 multiplexed in a single sequence. Each model multiplexed from two samples up to four samples at a
190 single sequence. The ‘Concat Point’ means at which layer the TrE tokenizer is concatenated. In other
191 words, before concat point each sample is processed independently and after the concat point N_{MUX}
192 samples are concatenated and processed at once. We call layers before and after the concat point as
193 projection layers and backbone layers, respectively. The total number of layers, including both the
194 projection layers and the backbone layers, was set to 12. The batch size of the Image Multiplexer and
195 the ConcatPlexer is 512-1024 to fit the size of GPU memory. The Image Multiplexer was trained on 8
196 A100 GPUs and ConcatPlexer was trained on 4 A100 GPUs, respectively.

197 4.3 Experiment on ImageNet1K

198 The aforementioned models are pretrained with ImageNet1K and results are reported in Table 1. As
199 shown in Table 1, performance tends to drop as N_{MUX} parameter increases in both Image Multiplexer
200 and ConcatPlexer. This is natural because N_{MUX} being four means that twice more information
201 should be crammed in the same space as N_{MUX} being two. Also, although the total number of
202 projection layers and backbone layers is kept to 12 in total, the thicker projection layer tends to show
203 better performance in the cost of computational efficiency. This is because the projection layer is a
204 layer that comes before muxing and the backbone is a layer that comes after muxing.

205 Referring to Table 1, Image Multiplexer using CNN tokenizer saturated at validation accuracy of 26%.
206 Replacing the tokenizer with DALL-E has boosted validation performance up to 48%, muxing 4 image
207 inputs. The DALL-E tokenizer enables image patches to be discrete, but its token length prohibitively
208 increases the computation cost, which makes the purpose of multiplexed image classification task
209 pointless. On the contrary, the ConcatPlexer shows validation accuracy of 56% with muxing four
210 images at the same time using a single layer TrE tokenizer. The ConcatPlexer with $N_{MUX} = 2$
211 shows validation accuracy of 62% - 69%. Performance gets better as projection Layers increase in
212 the cost of computational efficiency. The ConcatPlexer is compared with conventional ViT[3] most
213 similar backbone but trained with a non-multiplexed image classification task at section 4.5.

214 4.4 Experiment on CIFAR100

215 The pretrained models from Table 1 were finetuned on the CIFAR100 dataset. Similar to Table 1,
216 the ConcatPlexer with smaller Num Muxed and larger projection Layers performs better in the cost
217 of computational cost. Referring to the Table 2, the performance of the ConcatPlexers outperforms
218 the Image Multiplexers with less computational cost. As the model is trained on easier task (smaller
219 number of classes), the degradation gap between ConcatPlexer and ViT reduces. According to Table
220 2 and Table 3, ConcatPlexer(3)’s validation accuracy is 83.4% and ViT-B/16’s validation accuracy is
221 87.13%.

Model	GFLOPs
ViT-B/16	17.58
ConcatPlexer(1)	9.81
ConcatPlexer(2)	11.26
ConcatPlexer(3)	13.45
ConcatPlexer(4)	5.81

Table 3: FLOPs Count Per Image.

Model	Dataset	Val Acc
ViT-B/16	INET1K	77.91%
ConcatPlexer(3)	INET1K	69.5%
ViT-B/16	CIFAR100	87.13%
ConcatPlexer(3)	CIFAR100	83.4%

Table 4: Comparison with ViT.

Model	ImageNet1K Val Acc
W/O MUX	60.7%
MUX	62.3%

Table 5: Thicker Batch vs. Multiplexing.

IR	@1	@5	@10
MMP	32.84%	59.62%	71.42%

Table 6: Flickr zero-shot image retrieval.

222 4.5 Ablation

223 **Comparison with non-multiplexing method:** The ConcatPlexer uses TrE patchifier to get high-level
 224 patch tokens. Then token length is reduced by the conv layer and multiple inputs are concatenated.
 225 Instead of concatenating and just stacking the reduced length inputs after conv computation may look
 226 like a good option. Table 5 indicates that ConcatPlexer performs better than *Without MUX model*.

227 **Comparing with conventional ViT:** Table 4 indicates that ConcatPlexer lacks performance in
 228 ImageNet1K compared to ViT-B/16 . This is because the ConcatPlexer is tackling a harder task:
 229 multiplexed image classification. However, the performance gap narrows if the model is trained on
 230 an easier dataset: CIFAR100. Considering that we are the first to propose the multiplexed image
 231 classification task, there is more room for improvement. Therefore, we believe that the current aspect
 232 seems encouraging.

233 The computational cost of the Image Multiplexer using CNN patchifier was not calculated as its
 234 performance was not comparable with other models. The computational cost of the Image Multiplexer
 235 using DALL-E dVAE was also not calculated as its computational cost was prohibitively large due to
 236 the long sequence length.

237 As the main purpose of the ConcatPlexer is to attain increased computation efficiency and throughput,
 238 we also compared the FLOPs of the ConcatPlexers and ViT-B/16. Table 3 indicates that the Concat-
 239 Plexers require less FLOPs compared to ViT. The FLOPs were counted using fvcore library of Meta
 240 Research.

241 **Comparison with original DataMUX:** Original DataMUX and its descendent [15, 18] demonstrated
 242 its effectiveness on GLUE benchmark [19]. However, an expectation of random chance of CIFAR100
 243 and ImageNet1K is much lower, considering that tasks in the GLUE benchmark usually have two to
 244 three classes to predict. Of course, DataMUX shows an impressive performance in the aspect that it
 245 can multiplex extremely many inputs (up to 10 for MUX-PLM). The ConcatPlexer multiplexes fewer
 246 inputs and the performance gap may be seen as quite large. However, the ConcatPlexer deals with the
 247 trickier task in the sense of expectation of random chance.

248 **Possibility toward multimodal multiplexing:** As a proof of concept, we adapt data multiplexing in
 249 the Vision&Language(VL) domain: Multimodal Multiplexer (MMP). It is a single ViT-architected
 250 model that represents images and texts in a modality-agnostic manner, with no modality-specific
 251 transformer blocks or projection layers. As shown in Table 6, the model achieves 32.84%, 59.62%,
 252 and 71.42% accuracy at recall@1, recall@5, and recall@10 on Flickr30K [22] with zero-shot
 253 image retrieval task. Despite its modality-agnostic information processing mechanism and parameter
 254 efficiency (compared to typical VL models that employ separate transformer towers for each modality),
 255 MMP achieves promising preliminary results on challenging multimodal retrieval task, which we
 256 believe indicates great potential for future follow-up works.

257 5 Limitations and Future works

258 We propose a new vision classification called multiplexed image classification task to multiply the
259 throughput of the neural network. Though this approach can drastically drag up the computational
260 efficiency and throughput by increasing N_{MUX} , the concept of calculating multiple data at the
261 same time makes the previous task much harder, which causes a clear trade-off between N_{MUX} and
262 performance.

263 The performance of transformer-based models in vision tends to be heavily affected by hyperparameter
264 tuning. With more delicate tuning, the ConcatPlexer may have more performance gain. As our
265 Conv-based multiplexing method is somewhat heuristic, we believe that there is more room for
266 improvement if other token length reduction methods are used together. Also, there is room for
267 computation efficiency gain of the ConcatPlexer if we use a method similar to Swin Transformer
268 [23], which separates images from the entire sequence and divides them into partitions, applies local
269 self-attention, and combines information from the CLS token.

270 6 Conclusion

271 From this paper, we would like to bring constructive discussion for the computational efficiency of
272 transformer-based models. Our paper tries to transplant the idea of DataMUX [15] of NLP into the
273 vision field. We propose the multiplexed image classification task and its baseline ConcatPlexer and
274 Image Multiplexer. The proposed model shows the feasibility that idea of DataMUX can be applied
275 in the vision field.

276 References

- 277 [1] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
278 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
279 pages 770–778, 2016.
- 280 [2] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,
281 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. 2017.
- 282 [3] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,
283 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al.
284 An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint*
285 *arXiv:2010.11929*, 2020.
- 286 [4] Chaerin Kong, Jeessoo Kim, Donghoon Han, and Nojun Kwak. Few-shot image generation with
287 mixup-based distance learning. In *Computer Vision—ECCV 2022: 17th European Conference,*
288 *Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XV*, pages 563–580. Springer, 2022.
- 289 [5] Chaerin Kong, DongHyeon Jeon, Ohjoon Kwon, and Nojun Kwak. Leveraging off-the-shelf
290 diffusion model for multi-attribute fashion image manipulation. In *Proceedings of the IEEE/CVF*
291 *Winter Conference on Applications of Computer Vision*, pages 848–857, 2023.
- 292 [6] Seunghyeon Seo, Donghoon Han, Yeonjin Chang, and Nojun Kwak. Mixerf: Modeling
293 a ray with mixture density for novel view synthesis from sparse inputs. *arXiv preprint*
294 *arXiv:2302.08788*, 2023.
- 295 [7] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr:
296 Deformable transformers for end-to-end object detection. *arXiv preprint arXiv:2010.04159*,
297 2020.
- 298 [8] Jiho Jang, Chaerin Kong, Donghyeon Jeon, Seonhoon Kim, and Nojun Kwak. Unifying vision-
299 language representation space with single-tower transformer. *arXiv preprint arXiv:2211.11153*,
300 2022.
- 301 [9] Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santi-
302 ago Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, et al. Big bird: Transformers
303 for longer sequences. *Advances in neural information processing systems*, 33:17283–17297,
304 2020.

- 305 [10] Joshua Ainslie, Tao Lei, Michiel de Jong, Santiago Ontañón, Siddhartha Brahma, Yury Zemlyan-
306 ski, David Uthus, Mandy Guo, James Lee-Thorp, Yi Tay, et al. Colt5: Faster long-range
307 transformers with conditional computation. *arXiv preprint arXiv:2303.09752*, 2023.
- 308 [11] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena,
309 Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified
310 text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020.
- 311 [12] Zhenglun Kong, Peiyan Dong, Xiaolong Ma, Xin Meng, Mengshu Sun, Wei Niu, Xuan Shen,
312 Geng Yuan, Bin Ren, Minghai Qin, et al. Spvit: Enabling faster vision transformers via soft
313 token pruning. *arXiv preprint arXiv:2112.13890*, 2021.
- 314 [13] Youwei Liang, Chongjian Ge, Zhan Tong, Yibing Song, Jue Wang, and Pengtao Xie. Not all
315 patches are what you need: Expediting vision transformers via token reorganizations. *arXiv*
316 *preprint arXiv:2202.07800*, 2022.
- 317 [14] Daniel Bolya, Cheng-Yang Fu, Xiaoliang Dai, Peizhao Zhang, Christoph Feichtenhofer, and
318 Judy Hoffman. Token merging: Your vit but faster. *arXiv preprint arXiv:2210.09461*, 2022.
- 319 [15] Vishvak Murahari, Carlos Jimenez, Runzhe Yang, and Karthik Narasimhan. Datamux: Data mul-
320 tiplexing for neural networks. *Advances in Neural Information Processing Systems*, 35:17515–
321 17527, 2022.
- 322 [16] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-
323 scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern*
324 *recognition*, pages 248–255. Ieee, 2009.
- 325 [17] Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar-10 and cifar-100 datasets. *URL:*
326 *https://www.cs.toronto.edu/kriz/cifar.html*, 6(1):1, 2009.
- 327 [18] Vishvak Murahari, Ameet Deshpande, Carlos E Jimenez, Izhak Shafran, Mingqiu Wang, Yuan
328 Cao, and Karthik Narasimhan. Mux-plms: Pre-training language models with data multiplexing.
329 *arXiv preprint arXiv:2302.12441*, 2023.
- 330 [19] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman.
331 Glue: A multi-task benchmark and analysis platform for natural language understanding. *arXiv*
332 *preprint arXiv:1804.07461*, 2018.
- 333 [20] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked
334 autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on*
335 *Computer Vision and Pattern Recognition*, pages 16000–16009, 2022.
- 336 [21] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark
337 Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on*
338 *Machine Learning*, pages 8821–8831. PMLR, 2021.
- 339 [22] Bryan A Plummer, Liwei Wang, Chris M Cervantes, Juan C Caicedo, Julia Hockenmaier, and
340 Svetlana Lazebnik. Flickr30k entities: Collecting region-to-phrase correspondences for richer
341 image-to-sentence models. In *Proceedings of the IEEE international conference on computer*
342 *vision*, pages 2641–2649, 2015.
- 343 [23] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining
344 Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings*
345 *of the IEEE/CVF international conference on computer vision*, pages 10012–10022, 2021.