

# 497 Appendices

## 498 A Proof for Lemma 1.

499 Given a linear transformation  $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{b}$ ,  $\mathbf{x} \in \mathbb{R}^n$ ,  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{b}, \mathbf{y} \in \mathbb{R}^m$ , we have

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{b} \quad (7)$$

$$= (\mathbf{A} - k\mathbf{1}\mathbf{1}^T \mathbf{A})\mathbf{x} + k\mathbf{1}\mathbf{1}^T \mathbf{A}\mathbf{x} + (\mathbf{b} - \mu(\mathbf{b})\mathbf{1}) + \mu(\mathbf{b})\mathbf{1} \quad (8)$$

$$= (\mathbf{A} - k\mathbf{1}\mathbf{1}^T \mathbf{A})\mathbf{x} + (\mathbf{b} - \mu(\mathbf{b})\mathbf{1}) + k(\mathbf{1}^T \mathbf{A}\mathbf{x})\mathbf{1} + \mu(\mathbf{b})\mathbf{1} \quad (9)$$

$$= (\mathbf{A} - k\mathbf{1}\mathbf{1}^T \mathbf{A})\mathbf{x} + (\mathbf{b} - \mu(\mathbf{b})\mathbf{1}) + (k\mathbf{1}^T \mathbf{A}\mathbf{x} + \mu(\mathbf{b}))\mathbf{1} \quad (10)$$

$$= \hat{\mathbf{A}}\mathbf{x} + \hat{\mathbf{b}} + f(\mathbf{x}, k)\mathbf{1} \quad (11)$$

500 where  $\hat{\mathbf{A}} = \mathbf{A} - k\mathbf{1}\mathbf{1}^T \mathbf{A}$ ,  $\hat{\mathbf{b}} = \mathbf{b} - \mu(\mathbf{b})\mathbf{1}$ ,  $f(\mathbf{x}, k) = k\mathbf{1}^T \mathbf{A}\mathbf{x} + \mu(\mathbf{b})$ .

501 If  $k = 1/m$ , then we obtain

$$\mu(\hat{\mathbf{A}}\mathbf{x}) = \frac{1}{m}\mathbf{1}^T (\mathbf{A} - \frac{1}{m}\mathbf{1}\mathbf{1}^T \mathbf{A})\mathbf{x} = \frac{1}{m}(\mathbf{1}^T \mathbf{A} - \mathbf{1}^T \mathbf{A})\mathbf{x} = 0 \quad (12)$$

$$\mu(\mathbf{y}) = \mu(\hat{\mathbf{A}}\mathbf{x}) + \mu(\hat{\mathbf{b}}) + \mu(f(\mathbf{x}, k = 1/m)\mathbf{1}) \quad (13)$$

$$= 0 + 0 + f(\mathbf{x}, k = 1/m) \quad (14)$$

$$= \frac{1}{m}\mathbf{1}^T \mathbf{A}\mathbf{x} + \mu(\mathbf{b}) \quad (15)$$

502 The  $\hat{\mathbf{A}}$  is the recentered matrix of  $\mathbf{A}$ , and all its column vectors have zero-mean. We decompose the  
503 output into two parts.

504 • The first part  $\hat{\mathbf{A}}\mathbf{x} + \hat{\mathbf{b}} = \mathbf{y} - \mu(\mathbf{y})\mathbf{1}$ , with zero mean, is another linear transformation with  
505  $\hat{\mathbf{A}} = \mathbf{A} - \frac{1}{m}\mathbf{1}\mathbf{1}^T \mathbf{A}$ ,  $\hat{\mathbf{b}} = \mathbf{b} - \mu(\mathbf{b})\mathbf{1}$ .

506 • The second part corresponds to the mean information  $\mu(\mathbf{y})\mathbf{1} = (\frac{1}{m}\mathbf{1}^T \mathbf{A}\mathbf{x} + \mu(\mathbf{b}))\mathbf{1}$ .

## 507 B Post-LN Transformers

508 Different from the Pre-LN Transformers, the Post-LN Transformers have the following blocks.

$$\mathbf{x}_{l+1} = \text{LN}(\mathbf{x}_l + \mathcal{F}_l(\mathbf{x}_l)), \quad l = 0, 1, \dots, L - 1, \quad (16)$$

509 Layer normalization is on the main branch instead of the beginning of residual branches. We can  
510 keep a zero-mean branch on the main branch without impact on the functionality.

$$\mathbf{x}_{l+1} = \text{LN}(\mathbf{x}_l + \mathcal{F}_l(\mathbf{x}_l)) \quad (17)$$

$$= \text{LN}((\mathbf{x}_l - \mu(\mathbf{x}_l)\mathbf{1}) + (\mathcal{F}_l(\mathbf{x}_l) - \mu(\mathcal{F}_l(\mathbf{x}_l))\mathbf{1})) \quad (18)$$

$$= \text{LN}(\hat{\mathbf{x}}_l + \hat{\mathcal{F}}_l(\mathbf{x}_l)) \quad (19)$$

$$= \text{RMSNorm}(\hat{\mathbf{x}}_l + \hat{\mathcal{F}}_l(\mathbf{x}_l)) \quad (20)$$

511 For the residual branch  $\mathcal{F}_l$ , we can apply the same method in Pre-LN Transformer. We can modify  
512 the output linear projection to obtain  $\hat{\mathcal{F}}_l$ , which will generate the zero-mean part of the original result.

513 The recentering operation  $\hat{\mathbf{x}}_l = \mathbf{x}_l - \mu(\mathbf{x}_l)\mathbf{1}$  requires extra computation. If elementwise affine  
514 transformation is disabled in LayerNorm,  $\mathbf{x}_l$  is the output of a normalization such that  $\mu(\mathbf{x}_l) = 0$   
515 and  $\hat{\mathbf{x}}_l = \mathbf{x}_l$ . If the transformation is enabled,  $\mathbf{x}_l$  is not guaranteed zero-mean such that the explicit  
516 recentering is necessary.

## 517 C Experiments

### 518 C.1 Implementation of normalization

519 We have provided our implementation with JAX and PyTorch in the supplementary material. The  
520 reported results are based on the following implementations.

521 For JAX, we use the APIs of LayerNorm and RMSNorm in the flax library. For PyTorch, we use the  
 522 implementations of LayerNorm and RMSNorm from NVIDIA’s apex extension <sup>4</sup>. For CRMSNorm,  
 523 we use our own customized implementations. We also provide our customized implementation of  
 524 LayerNorm, RMSNorm.

525 We notice that there are lots of APIs for the standard LayerNorm and RMSNorm. For example,  
 526 PyTorch has provided the official LayerNorm API but lacks of RMSNorm implementation. These  
 527 different implementations are mixed. We do not find one implementation that is dominant over others  
 528 for all the cases. For instance, `torch.nn.LayerNorm` is usually faster than apex’s one when the  
 529 input vectors are small in inference, while it is slower than the apex’s one when the input vectors are  
 530 large. PyTorch’s official implementation is also slower than apex’s for training.

## 531 C.2 Extended Experiments in ViT

Name	Dimension	Depth	Heads	MLP Dimension
Tiny-16	192	12	3	192 × 4
Small-16	384	12	6	384 × 4
Base-16	768	12	12	768 × 4
Large-16	1024	24	16	1024 × 4
Huge-14	1280	32	16	1280 × 4
Giant-14	1664	48	16	8192

Table 2: ViTs with different sizes. The number in the model name is the patch size.

	no norm	Pre-LN	Pre-RMS	Pre-CRMS
<b>PyTorch, single A100, amp</b>	0.8567	1.000	0.9699	0.9783
amp → float32	0.9353	1.000	0.9850	0.9951
single A100 → 16-thread CPU	0.8697	1.000	0.9012	0.8857
PyTorch → JAX	0.9610	1.000	0.9873	1.0005

Table 3: Normalized inference time of ViT.

532 Table 2 list the architecture parameters of Vision Transformer. We first measure the **inference time**.  
 533 We sweep these 6 ViTs with 6 batch sizes (1, 4, 16, 64, 256, 1024) and collect the medians of these  
 534 36 data points. We report the average of these 36 experiments in Table 3. We conduct inference on  
 535 a single A100 with automatic mixed precision (amp) in PyTorch. We further change the precision  
 536 (disabling the amp), computation platforms (16 threads in AMD EPYC 7742 CPUs), and machine  
 537 learning frameworks (JAX).

## 538 C.3 Numerical Issue

539 The theoretical arithmetic equivalence cannot be fully translated into equality in the practical numeri-  
 540 cal computation if we use floating numbers. An intuitive example is that  $\mu(\mathbf{x} + \mathbf{y}) = \mu(\mathbf{x}) + \mu(\mathbf{y})$   
 541 always holds for any vectors  $\mathbf{x}, \mathbf{y}$ . However, if these two vectors are represented as (low precision)  
 542 floating numbers, this equality is not guaranteed in real-world numerical computation. It is possible  
 543 that these small discrepancies may be accumulated and enlarged in the large models, further degrading  
 544 the numerical stability. In our proposed method, we cannot ensure the exactly zero-mean in the main  
 545 branch numerically.

546 The numerical issue is a common problem in machine learning. A typical example is  
 547 operator reordering and layer fusion. PyTorch provides a related API officially, named  
 548 `torch.ao.quantization.fuse_modules`. We can fuse the convolution layer and its following  
 549 batch normalization layer to simplify the computation. These two layers are separate in training and  
 550 can be fused to accelerate the inference. The fusion does not break the arithmetic equivalence but  
 551 changes the numerical results. In spite of the numerical difference, the fusion usually has a neutral  
 552 impact on task-related performance, such as classification accuracy, even in large models. Fine-tuning  
 553 or calibration may be helpful in case there is severe performance degradation.

<sup>4</sup><https://github.com/NVIDIA/apex>

554 Our proposed methods encounter a similar issue as layer fusion since we modify partial parameters.  
555 In our experiments, we can convert the pre-trained Pre-LN ViT-H/14 into Pre-(C)RMS variants  
556 without any accuracy change on the ImageNet validation dataset. Actually, we observe that replacing  
557 PyTorch's official LayerNorm implementation with the apex's one may have a larger impact on the  
558 model performance.