

Vision And Language Final Report

Yu-Wei Su

yfs5313@psu.edu

1. Task

The task of the paper is to tackle Video Question Answering (VideoQA). The goal of VideoQA is to answer questions about the content of a given video. It involves understanding the visual scenes in the video as well as the linguistic semantics of the question to produce an accurate answer. VideoQA is more complex than ImageQA (image-based question answering) because it requires the model to understand the temporal dynamics and relationships within the video content. [16]

The paper proposes a new long-form Video Question Answering (VideoQA) model called multi-modal iterative special-temporal Transformer (MIST). Previous models tackling VideoQA often focus on only short or video clips (1-2 seconds). However, some challenges arise when it comes to longer videos. First, Long videos contain more events thus requiring the models to have more complex temporal reasoning to answer the question. Traditional methods either rely on dense, computationally expensive video sampling or sparse sampling, which may overlook crucial visual information. Second, the question may contain objects, events, and relations that happen at other times compared to the simultaneous events in short videos.

2. Related Work

2.1. Video Question Answering SOTA (in AGQAv2)

As listed in <https://paperswithcode.com/dataset/agqa>. Two current competing SOTA are in the VideoQA task, especially in the AGQA dataset. The first is MIST, and the second is [3], which can be seen from the table in fig1. Although [3] are winning in half of the sub-categories, including overall AGQA performance. MIST is still better in temporal-related questions such as duration and sequencing tasks.

3. Approach

MIST decomposes the traditional dense spatial-temporal self-attention into cascaded segment and region selection modules that adaptively select frames and image regions relevant to the question. 2) Iterative Selection and Attention

Question Types	Most Likely	PSAC	HME	HCRN [23]	AIO [40]	Temp[ATP] [5]	MIST - AIO
Object-relation	9.39	37.84	37.42	40.33	48.34	50.15	51.43
Relation-action	50.00	49.95	49.90	49.86	48.99	49.76	54.67
Object-action	50.00	50.00	49.97	49.85	49.66	46.25	55.37
Superlative	21.01	33.20	33.21	33.55	37.53	39.78	41.34
Sequencing	49.78	49.78	49.77	49.70	49.61	48.25	53.14
Exists	50.00	49.94	49.96	50.01	50.81	51.79	53.49
Duration comparison	24.27	45.21	47.03	43.84	45.36	49.59	47.48
Activity recognition	5.52	4.14	5.43	5.52	18.97	18.96	20.18
All	10.99	40.18	39.89	42.11	48.59	49.79	50.96

Table 1. QA accuracies of state-of-the-art (SOTA) methods on AGQA v2 test set.

Figure 1. AGQA v2 performance comparison captured from [3]

over Multiple Layers: MIST iteratively conducts selection and attention over multiple layers to support reasoning over multiple events, which can be shown in fig2. Experimental results on four VideoQA datasets demonstrate that MIST achieves state-of-the-art performance while being computationally efficient and interpretable. A more detailed ex-

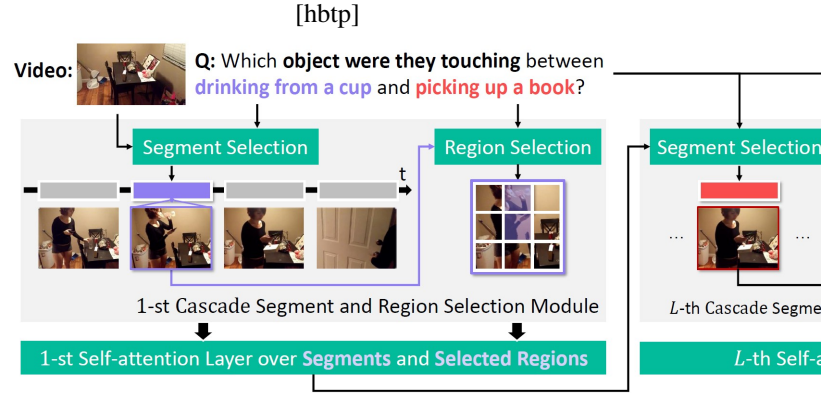


Figure 2. Model diagram as from Original paper

planation of its method can be seen in figure 3. The key part of MIST is that it tries to perform attention to segments of both temporal and regional parts. First, it will dissect a video into k segments, pick one patch by pooling, and then perform attention to each layer's selected patch. After that, the selector will select TOP_k result from k parts to get the temporal features. The model will further perform region selection on the TOP k results. The region section can

also be done using a similar approach by dissecting regions into top j results. With spatial, temporal, and word features obtained, multi-head attention is then performed to capture cross-related relations.

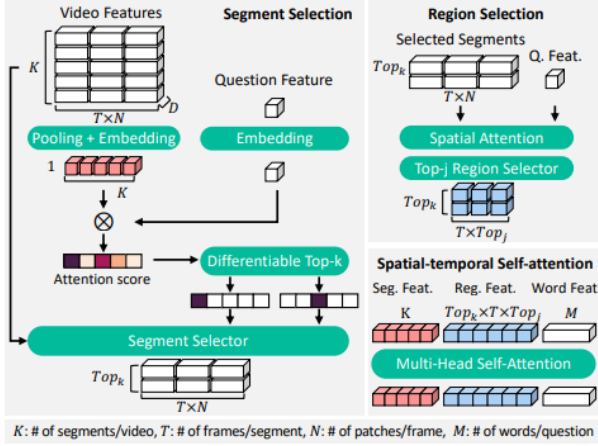


Figure 4. Key components of Iterative Spatial-Temporal Attention Layer. Since region selection follows the same architecture as segment selection, we only show its inputs and outputs.

Figure 3. Figure 4 from [7]

4. Dataset

The datasets [7] has trained on are in figure 4, which is Table 1 captured from [7]. This report focuses on replicating AGQA v2 [8], specifically. As [7] has stated, AGQA is an open-ended VideoQA benchmark for compositional spatio-temporal reasoning. The V2 version has a more balanced distribution. It provides 2.27M QA pairs over 9.7K videos with an average length of 30 seconds.

5. Results

There is no pretrained model provided by the original paper but only with the results from the paper. Thus, I only reported a comparison between my replication result and the paper’s result. A comparison between the paper’s result and my replication result is shown below. As we can see,

Table 1: Comparison of the Multi-Event VideoQA benchmarks.

Datasets	# Avg. Events/Video	Video Length	#Event Classes	#Videos	# QAs
STAR[42]	2.7	12s	157	22K	60K
AGQA v2[14]	6.8	30s	157	9.7K	2.27M
EgoTaskQA[16]	5	25s	793	2K	40K
NExT-QA[43]	8.8	44s	50	5.4K	52K

Figure 4. Enter Caption

Model	original	epoch 5	epoch 10	epoch 14	epoch 16
Accuracy	54.39	51.20	52.88	53.85	53.62

Table 1. replication vs original result.

2023-12-10 12:22:17,922 INFO	Epoch 12, Epoch status: 0.0696, Training VideoQA loss: 0.1352, Training acc: 94.06%
2023-12-10 12:27:49,143 INFO	Epoch 12, Epoch status: 0.1332, Training VideoQA loss: 0.1302, Training acc: 94.21%
2023-12-10 12:33:20,144 INFO	Epoch 12, Epoch status: 0.1998, Training VideoQA loss: 0.1281, Training acc: 94.31%
2023-12-10 12:38:50,865 INFO	Epoch 12, Epoch status: 0.2664, Training VideoQA loss: 0.1245, Training acc: 94.53%
2023-12-10 12:44:22,221 INFO	Epoch 12, Epoch status: 0.3338, Training VideoQA loss: 0.1205, Training acc: 94.37%
2023-12-10 12:49:56,000 INFO	Epoch 12, Epoch status: 0.3996, Training VideoQA loss: 0.1254, Training acc: 94.53%
2023-12-10 12:55:27,203 INFO	Epoch 12, Epoch status: 0.4663, Training VideoQA loss: 0.1265, Training acc: 94.42%
2023-12-10 13:00:50,878 INFO	Epoch 12, Epoch status: 0.5329, Training VideoQA loss: 0.1262, Training acc: 94.44%
2023-12-10 13:06:29,041 INFO	Epoch 12, Epoch status: 0.5995, Training VideoQA loss: 0.1252, Training acc: 94.40%
2023-12-10 13:11:02,137 INFO	Epoch 12, Epoch status: 0.6661, Training VideoQA loss: 0.1259, Training acc: 94.40%
2023-12-10 13:17:33,304 INFO	Epoch 12, Epoch status: 0.7327, Training VideoQA loss: 0.1231, Training acc: 94.60%
2023-12-10 13:23:04,578 INFO	Epoch 12, Epoch status: 0.7993, Training VideoQA loss: 0.1254, Training acc: 94.49%
2023-12-10 13:28:35,651 INFO	Epoch 12, Epoch status: 0.8659, Training VideoQA loss: 0.1247, Training acc: 94.60%
2023-12-10 13:34:06,613 INFO	Epoch 12, Epoch status: 0.9325, Training VideoQA loss: 0.1250, Training acc: 94.55%
2023-12-10 13:39:39,564 INFO	Epoch 12, Epoch status: 0.9991, Training VideoQA loss: 0.1257, Training acc: 94.52%
2023-12-10 13:45:10,996 INFO	val acc: 54.11%
2023-12-10 13:49:31,004 INFO	val acc@8: 95.85%
2023-12-10 13:55:02,500 INFO	Epoch 13, Epoch status: 0.0666, Training VideoQA loss: 0.1768, Training acc: 92.42%
2023-12-10 14:00:24,121 INFO	Epoch 13, Epoch status: 0.1332, Training VideoQA loss: 0.1709, Training acc: 92.61%
2023-12-10 14:05:01,000 INFO	Epoch 13, Epoch status: 0.1998, Training VideoQA loss: 0.1729, Training acc: 92.48%
2023-12-10 14:10:32,921 INFO	Epoch 13, Epoch status: 0.2664, Training VideoQA loss: 0.1725, Training acc: 92.56%
2023-12-10 14:15:04,953 INFO	Epoch 13, Epoch status: 0.3338, Training VideoQA loss: 0.1708, Training acc: 92.63%
2023-12-10 14:20:37,217 INFO	Epoch 13, Epoch status: 0.3996, Training VideoQA loss: 0.1739, Training acc: 92.55%
2023-12-10 14:26:09,304 INFO	Epoch 13, Epoch status: 0.4663, Training VideoQA loss: 0.1716, Training acc: 92.56%
2023-12-10 14:31:43,945 INFO	Epoch 13, Epoch status: 0.5329, Training VideoQA loss: 0.1727, Training acc: 92.57%
2023-12-10 14:37:16,376 INFO	Epoch 13, Epoch status: 0.5995, Training VideoQA loss: 0.1708, Training acc: 92.65%
2023-12-10 14:42:48,570 INFO	Epoch 13, Epoch status: 0.6661, Training VideoQA loss: 0.1725, Training acc: 92.44%
2023-12-10 14:48:19,638 INFO	Epoch 13, Epoch status: 0.7327, Training VideoQA loss: 0.1751, Training acc: 92.43%
2023-12-10 14:53:54,986 INFO	Epoch 13, Epoch status: 0.7993, Training VideoQA loss: 0.1731, Training acc: 92.47%
2023-12-10 14:59:26,424 INFO	Epoch 13, Epoch status: 0.8659, Training VideoQA loss: 0.1763, Training acc: 92.35%
2023-12-10 15:04:58,812 INFO	Epoch 13, Epoch status: 0.9325, Training VideoQA loss: 0.1742, Training acc: 92.48%
2023-12-10 15:10:30,940 INFO	Epoch 13, Epoch status: 0.9991, Training VideoQA loss: 0.1744, Training acc: 92.40%
2023-12-10 15:16:02,980 INFO	val acc: 53.88%
2023-12-10 15:21:34,980 INFO	val acc@8: 96.25%
2023-12-10 15:27:06,980 INFO	Epoch 14, Epoch status: 0.0666, Training VideoQA loss: 0.1539, Training acc: 93.40%
2023-12-10 15:32:38,980 INFO	Epoch 14, Epoch status: 0.1332, Training VideoQA loss: 0.1555, Training acc: 93.34%
2023-12-10 15:38:10,980 INFO	Epoch 14, Epoch status: 0.1998, Training VideoQA loss: 0.1580, Training acc: 93.26%
2023-12-10 15:43:42,980 INFO	Epoch 14, Epoch status: 0.2664, Training VideoQA loss: 0.1508, Training acc: 93.07%

Figure 5. portion of training log-1.

The model has acquired a decent already in the first few epochs.

5.1. Issues

1. The replication is only performed to Epoch 16 due to the limitation of the computing power.
2. The first few training epochs are trained on the same learning rates. In the first few epochs, the connection of collab has been bad, so the training kept stopping. The fact that MIST uses a learning rate scheduler was later discovered in the middle of the training process. So only epochs trained from 11 to 16 are modified and trained on the same scheduled learning rate same to the original result.
3. The replication results are reported with the val dataset due to the reason that the author does not provide the testing set data.

5.2. Screenshots of training log

Some portions of the training log of MIST are provided below (5, 6).

6. Possible Improvements and Results

6.1. Possible Improvements

There are two possible improvements, but I only tried to improve one. First is the fixed segments, and the other is the optimizer-related potential improvement.

2023-12-18 17:45:47,668 INFO	epoch 15, epoch status: 0.9529, Training VideoQA loss: 0.1512, Training acc: 95.08%
2023-12-18 17:51:12,071 INFO	epoch 15, epoch status: 0.9595, Training VideoQA loss: 0.1480, Training acc: 95.71%
2023-12-18 17:56:52,980 INFO	epoch 15, epoch status: 0.9661, Training VideoQA loss: 0.1461, Training acc: 95.65%
2023-12-18 18:01:24,644 INFO	epoch 15, epoch status: 0.7327, Training VideoQA loss: 0.1509, Training acc: 93.40%
2023-12-18 18:07:59,704 INFO	epoch 15, epoch status: 0.7993, Training VideoQA loss: 0.1486, Training acc: 93.74%
2023-12-18 18:13:32,857 INFO	epoch 15, epoch status: 0.8659, Training VideoQA loss: 0.1506, Training acc: 93.40%
2023-12-18 18:19:06,344 INFO	epoch 15, epoch status: 0.9325, Training VideoQA loss: 0.1521, Training acc: 93.51%
2023-12-18 18:24:42,347 INFO	epoch 15, epoch status: 0.9991, Training VideoQA loss: 0.1516, Training acc: 93.53%
2023-12-18 18:34:15,750 INFO	val acc: 93.73%
2023-12-18 18:44:13,860 INFO	val acc@8: 96.28%
2023-12-18 18:47:06,510 INFO	epoch 16, epoch status: 0.9666, Training VideoQA loss: 0.1319, Training acc: 94.54%
2023-12-18 18:52:48,065 INFO	epoch 16, epoch status: 0.1332, Training VideoQA loss: 0.1317, Training acc: 94.49%
2023-12-18 19:01:13,693 INFO	epoch 16, epoch status: 0.1998, Training VideoQA loss: 0.1353, Training acc: 94.36%
2023-12-18 19:08:51,846 INFO	epoch 16, epoch status: 0.2664, Training VideoQA loss: 0.1359, Training acc: 94.30%
2023-12-18 19:09:25,016 INFO	epoch 16, epoch status: 0.3338, Training VideoQA loss: 0.1387, Training acc: 94.17%
2023-12-18 19:14:58,179 INFO	epoch 16, epoch status: 0.3996, Training VideoQA loss: 0.1365, Training acc: 94.31%
2023-12-18 19:20:13,754 INFO	epoch 16, epoch status: 0.4663, Training VideoQA loss: 0.1371, Training acc: 94.20%
2023-12-18 19:26:06,998 INFO	epoch 16, epoch status: 0.5329, Training VideoQA loss: 0.1389, Training acc: 94.12%
2023-12-18 19:31:48,131 INFO	epoch 16, epoch status: 0.5995, Training VideoQA loss: 0.1366, Training acc: 94.25%
2023-12-18 19:37:13,533 INFO	epoch 16, epoch status: 0.6661, Training VideoQA loss: 0.1365, Training acc: 94.20%
2023-12-18 19:42:47,289 INFO	epoch 16, epoch status: 0.7327, Training VideoQA loss: 0.1416, Training acc: 93.80%
2023-12-18 19:48:10,613 INFO	epoch 16, epoch status: 0.7993, Training VideoQA loss: 0.1364, Training acc: 94.22%
2023-12-18 19:53:56,518 INFO	epoch 16, epoch status: 0.8659, Training VideoQA loss: 0.1410, Training acc: 94.12%
2023-12-18 19:59:18,851 INFO	epoch 16, epoch status: 0.9325, Training VideoQA loss: 0.1396, Training acc: 94.16%
2023-12-18 20:00:00,912 INFO	epoch 16, epoch status: 0.9991, Training VideoQA loss: 0.1411, Training acc: 94.15%
2023-12-18 20:10:54,900 INFO	val acc: 93.62%
2023-12-18 20:12:12,382 INFO	val acc@8: 96.11%
2023-12-18 20:17:45,369 INFO	epoch 17, epoch status: 0.9666, Training VideoQA loss: 0.1283, Training acc: 95.80%
2023-12-18 20:23:10,780 INFO	epoch 17, epoch status: 0.1332, Training VideoQA loss: 0.1233, Training acc: 94.88%
2023-12-18 20:30:51,680 INFO	epoch 17, epoch status: 0.1998, Training VideoQA loss: 0.1237, Training acc: 94.85%
2023-12-18 20:38:14,996 INFO	epoch 17, epoch status: 0.2664, Training VideoQA loss: 0.1289, Training acc: 94.71%
2023-12-18 20:44:14,996 INFO	epoch 17, epoch status: 0.3338, Training VideoQA loss: 0.1275, Training acc: 94.75%
2023-12-18 20:50:00,420 INFO	epoch 17, epoch status: 0.3996, Training VideoQA loss: 0.1243, Training acc: 94.82%
2023-12-18 20:55:13,640 INFO	epoch 17, epoch status: 0.4663, Training VideoQA loss: 0.1286, Training acc: 94.60%
2023-12-18 21:01:00,512 INFO	epoch 17, epoch status: 0.5329, Training VideoQA loss: 0.1276, Training acc: 94.74%
2023-12-18 21:06:42,166 INFO	epoch 17, epoch status: 0.5995, Training VideoQA loss: 0.1261, Training acc: 94.70%
2023-12-18 21:12:15,793 INFO	epoch 17, epoch status: 0.6661, Training VideoQA loss: 0.1313, Training acc: 94.50%
2023-12-18 21:17:49,002 INFO	epoch 17, epoch status: 0.7327, Training VideoQA loss: 0.1294, Training acc: 94.40%
2023-12-18 21:23:12,002 INFO	epoch 17, epoch status: 0.7993, Training VideoQA loss: 0.1306, Training acc: 94.50%
2023-12-18 21:28:19,117 INFO	epoch 17, epoch status: 0.8659, Training VideoQA loss: 0.1306, Training acc: 94.64%
2023-12-18 21:34:18,330 INFO	epoch 17, epoch status: 0.9325, Training VideoQA loss: 0.1312, Training acc: 94.63%
	epoch 17, epoch status: 0.9991, Training VideoQA loss: 0.1331, Training acc: 94.47%

Figure 6. Enter Caption

6.1.1 fixed segments

it is unknown whether the model will perform poorly if the events are cut into different pieces. It may miss some vital event features due to the uniform segments without understanding segments of interest. Also, fixing K will lead the model to have a longer length per segment as the video lengths become longer. It could lead to a worse performance as irrelevant information is increased per segment.

Potential improvements could be tweaking the models to adaptively choose the hyperparameter K or cutting video segments into non-uniform segments to retain event information. Or even different K segments could also be used to capture long-term events and short-term events. Some challenges like computational expansiveness and how to feed it with different sizes would also arise. Due to the time limitation, I did not try to solve fixed segment issues.

6.1.2 Sharpness-Aware Minimization

Many works suggest that sharpness can be positively correlated with test errors of DNN ([15], [17], [14], [4], [12]). However, the following works showed that this correlation with generalization is not strong enough for Standard Sharpness. Some work [5] also shows that SAM could be helpful in vision and language tasks. Even though there are pieces of empirical evidence that don't agree with it [9], [10], [1]. I still would like to test whether it would help the generalizations.

The sharpness can be formulated as follows([2]):

$$s(\omega, S) := \max_{\|\delta\|_2 \leq \rho} \frac{1}{|S|} \sum_{i: (x_i, y_i) \in S} \ell_i(\omega + \delta) - \ell_i(\omega)$$

where $S = \{x_i, y_i\}_{i=1}^n$ be the training set, and $\ell_i(\omega)$ be the loss function corresponding to weights $\omega \in R^d$. The

Model	original	SAM epoch 2	SAM epoch 6
Accuracy	54.39	51.76	51.17
Model	MIST epoch 5	MIST epoch 6	MIST epoch 14
Accuracy	51.20	51.28	53.85

Table 2. replication vs original result.

key is that if one solution is sharper than the other solution, the maximum of its neighborhoods will have larger values since it goes up more drastically. It is seen that a solution with higher sharpness has lower generalization compared to a flatness solution empirically [6] [13]. Thus a Sharpness-Aware Minimization is proposed to minimize the object function with the care of sharpness.

With $\ell(\omega)$ be the original Empirical Risk Minimization(ERM) loss:

$$\ell(\omega) = \frac{1}{n} \sum_{i=1}^n \ell_i(\omega)$$

The SAM objective for perturbation would be:

$$\ell^{SAM}(\omega) = \max_{\|\delta\|_p \leq \rho} \frac{1}{|n|} \sum_{i=1}^n \ell_i(\omega + \delta)$$

An adaptive sharpness measure, which is what I tested, is also proposed to tackle rescaling issues [11]. The SAM package I used can be found at <https://github.com/davda54/sam>.

6.2. Modification Results

The SAM-modified version results are compared with the Original and non-modified versions. A comparison between the paper's, non-modified replication results, and SAM-MIST is shown below in table 2. The SAM-MIST actually has a better validation accuracy within a few epochs. However, since SAM requires two back propagation computing, the training time of the SAM version is twice for an epoch. Thus, only six epochs are trained on SAM, and whether it could perform better in the following epochs is unknown.

6.3. Issues

1. The replication is only performed to Epoch 6 due to the limitation of the computing power.
2. Similar to the MIST training. The first few training epochs are trained on the same learning rates. In the first few epochs, the connection of collab has been bad, so the training kept stopping. The fact that MIST uses a learning rate scheduler was later discovered during the training process. So Epoch 6 may actually be trained on the learning rate from Epoch 4 of the original settings.

```

2023-12-10 01:31:22,361 INFO val acc@0: 96.65%
2023-12-10 01:43:40,010 INFO Epoch 3, epoch status: 0.0666, Training VideoQA loss: 0.4884, Training acc: 74.11%
2023-12-10 01:55:44,202 INFO Epoch 3, epoch status: 0.1332, Training VideoQA loss: 0.4784, Training acc: 74.18%
2023-12-10 02:07:52,915 INFO Epoch 3, epoch status: 0.1998, Training VideoQA loss: 0.4684, Training acc: 74.39%
2023-12-10 02:19:55,951 INFO Epoch 3, epoch status: 0.2664, Training VideoQA loss: 0.4797, Training acc: 74.05%
2023-12-10 02:32:00,873 INFO Epoch 3, epoch status: 0.3330, Training VideoQA loss: 0.4735, Training acc: 74.52%
2023-12-10 02:44:03,010 INFO Epoch 3, epoch status: 0.3996, Training VideoQA loss: 0.4717, Training acc: 74.51%
2023-12-10 02:56:08,436 INFO Epoch 3, epoch status: 0.4663, Training VideoQA loss: 0.4686, Training acc: 74.62%
2023-12-10 03:08:10,458 INFO Epoch 3, epoch status: 0.5329, Training VideoQA loss: 0.4672, Training acc: 74.78%
2023-12-10 03:20:15,137 INFO Epoch 3, epoch status: 0.5995, Training VideoQA loss: 0.4632, Training acc: 74.89%
2023-12-10 03:32:17,430 INFO Epoch 3, epoch status: 0.6661, Training VideoQA loss: 0.4614, Training acc: 74.83%
2023-12-10 03:44:22,010 INFO Epoch 3, epoch status: 0.7327, Training VideoQA loss: 0.4552, Training acc: 75.17%
2023-12-10 03:56:24,943 INFO Epoch 3, epoch status: 0.7993, Training VideoQA loss: 0.4563, Training acc: 75.08%
2023-12-10 04:08:29,721 INFO Epoch 3, epoch status: 0.8659, Training VideoQA loss: 0.4517, Training acc: 75.27%
2023-12-10 04:20:32,141 INFO Epoch 3, epoch status: 0.9325, Training VideoQA loss: 0.4498, Training acc: 75.13%
2023-12-10 04:32:36,721 INFO Epoch 3, epoch status: 0.9991, Training VideoQA loss: 0.4449, Training acc: 75.61%
2023-12-10 04:47:38,369 INFO val acc: 51.28%
2023-12-10 04:47:38,370 INFO val acc@0: 96.41%
2023-12-10 04:59:49,831 INFO Epoch 4, epoch status: 0.0666, Training VideoQA loss: 0.4236, Training acc: 76.25%
2023-12-10 05:11:56,794 INFO Epoch 4, epoch status: 0.1332, Training VideoQA loss: 0.4383, Training acc: 75.93%
2023-12-10 05:23:58,116 INFO Epoch 4, epoch status: 0.1998, Training VideoQA loss: 0.4237, Training acc: 76.19%
2023-12-10 05:36:01,567 INFO Epoch 4, epoch status: 0.2664, Training VideoQA loss: 0.4258, Training acc: 76.29%
2023-12-10 05:48:02,473 INFO Epoch 4, epoch status: 0.3330, Training VideoQA loss: 0.4253, Training acc: 76.12%
2023-12-10 06:00:05,069 INFO Epoch 4, epoch status: 0.3996, Training VideoQA loss: 0.4218, Training acc: 76.36%
2023-12-10 06:12:07,261 INFO Epoch 4, epoch status: 0.4663, Training VideoQA loss: 0.4209, Training acc: 76.37%
2023-12-10 06:24:09,984 INFO Epoch 4, epoch status: 0.5329, Training VideoQA loss: 0.4285, Training acc: 76.31%
2023-12-10 06:36:10,854 INFO Epoch 4, epoch status: 0.5995, Training VideoQA loss: 0.4167, Training acc: 76.56%
2023-12-10 06:48:13,783 INFO Epoch 4, epoch status: 0.6661, Training VideoQA loss: 0.4183, Training acc: 76.57%
2023-12-10 07:00:14,729 INFO Epoch 4, epoch status: 0.7327, Training VideoQA loss: 0.4162, Training acc: 76.68%
2023-12-10 07:12:18,919 INFO Epoch 4, epoch status: 0.7993, Training VideoQA loss: 0.4095, Training acc: 76.87%
2023-12-10 07:24:19,696 INFO Epoch 4, epoch status: 0.8659, Training VideoQA loss: 0.4111, Training acc: 76.74%
2023-12-10 07:36:21,998 INFO Epoch 4, epoch status: 0.9325, Training VideoQA loss: 0.4114, Training acc: 76.63%
2023-12-10 07:48:25,271 INFO Epoch 4, epoch status: 0.9991, Training VideoQA loss: 0.4081, Training acc: 76.97%
2023-12-10 08:03:24,264 INFO val acc: 51.17%
2023-12-10 08:03:24,264 INFO val acc@0: 96.20%
2023-12-10 08:15:33,416 INFO Epoch 5, epoch status: 0.0666, Training VideoQA loss: 0.3896, Training acc: 77.58%
2023-12-10 08:27:38,506 INFO Epoch 5, epoch status: 0.1332, Training VideoQA loss: 0.3902, Training acc: 77.59%

```

Figure 7. log of SAM-MIST training

3. The replication results are reported with the val dataset because the author does not provide the testing set data.

6.4. Screenshots of training log

Some portions of the training log of SAM-MIST are provided below (7,).

7. Code Repository

The link to my training code can be found in. The original MIST github can be found in <https://github.com/thejackys/SAM-MIST>. The SAM optimizer can be found in <https://github.com/davda54/sam>. I only included the modified version for SAM as the replication version is basically follows the MIST’s own Github repo <https://github.com/showlab/mist>.

References

- [1] Maksym Andriushchenko, Francesco Croce, Maximilian Müller, Matthias Hein, and Nicolas Flammarion. A modern look at the relationship between sharpness and generalization. *arXiv preprint arXiv:2302.07011*, 2023. 3
- [2] Maksym Andriushchenko and Nicolas Flammarion. Towards understanding sharpness-aware minimization, 2022. 3
- [3] Ziyi Bai, Ruiping Wang, and Xilin CHEN. Glance and focus: Memory prompting for multi-event video question answering. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. 1
- [4] Junbum Cha, Sanghyuk Chun, Kyungjae Lee, Han-Cheol Cho, Seunghyun Park, Yunsung Lee, and Sungrae Park. Swad: Domain generalization by seeking flat minima. *Advances in Neural Information Processing Systems*, 34:22405–22418, 2021. 3
- [5] Xiangning Chen, Cho-Jui Hsieh, and Boqing Gong. When vision transformers outperform resnets without pre-training or strong data augmentations. *arXiv preprint arXiv:2106.01548*, 2021. 3

- [6] Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware minimization for efficiently improving generalization, Apr 2021. 3
- [7] Difei Gao, Luwei Zhou, Lei Ji, Linchao Zhu, Yi Yang, and Mike Zheng Shou. Mist: Multi-modal iterative spatial-temporal transformer for long-form video question answering, 2022. 2
- [8] Madeleine Grunde-McLaughlin, Ranjay Krishna, and Ma-neesh Agrawala. Agqa 2.0: An updated benchmark for compositional spatio-temporal reasoning, 2022. 2
- [9] Yiding Jiang, Behnam Neyshabur, Hossein Mobahi, Dilip Krishnan, and Samy Bengio. Fantastic generalization measures and where to find them. *arXiv preprint arXiv:1912.02178*, 2019. 3
- [10] Simran Kaur, Jeremy Cohen, and Zachary Chase Lipton. On the maximum hessian eigenvalue and generalization. In *Proceedings on*, pages 51–65. PMLR, 2023. 3
- [11] Jungmin Kwon, Jeongseop Kim, Hyunseo Park, and In Kwon Choi. Asam: Adaptive sharpness-aware minimization for scale-invariant learning of deep neural networks. In *International Conference on Machine Learning*, pages 5905–5914. PMLR, 2021. 3
- [12] Kaifeng Lyu, Zhiyuan Li, and Sanjeev Arora. Understanding the generalization benefit of normalization layers: Sharpness reduction. *Advances in Neural Information Processing Systems*, 35:34689–34708, 2022. 3
- [13] Behnam Neyshabur, Srinadh Bhojanapalli, David McAllester, and Nathan Srebro. Exploring generalization in deep learning, 2017. 3
- [14] Namuk Park and Songkuk Kim. How do vision transformers work? *arXiv preprint arXiv:2202.06709*, 2022. 3
- [15] Chen Xing, Devansh Arpit, Christos Tsirigotis, and Yoshua Bengio. A walk with sgd. *arXiv preprint arXiv:1802.08770*, 2018. 3
- [16] Yaoyao Zhong, Junbin Xiao, Wei Ji, Yicong Li, Weihong Deng, and Tat-Seng Chua. Video question answering: Datasets, algorithms and challenges, 2022. 1
- [17] Pan Zhou, Jiashi Feng, Chao Ma, Caiming Xiong, Steven Chu Hong Hoi, et al. Towards theoretically understanding why sgd generalizes better than adam in deep learning. *Advances in Neural Information Processing Systems*, 33:21285–21296, 2020. 3