
BMU-MoCo: Bidirectional Momentum Update for Continual Video-Language Modeling

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Video-language models suffer from forgetting old/learned knowledge when trained
2 with streaming data. In this work, we thus propose a continual video-language mod-
3 eling (CVLM) setting, where models are supposed to be sequentially trained on five
4 widely-used video-text datasets with different data distributions. Although most of
5 existing continual learning methods have achieved great success by exploiting extra
6 information (*e.g.*, memory data of past tasks) or dynamically extended networks,
7 they cause enormous resource consumption when transferred to our CVLM setting.
8 To overcome the challenges (*i.e.*, catastrophic forgetting and heavy resource con-
9 sumption) in CVLM, we propose a novel cross-modal MoCo-based model with
10 bidirectional momentum update (BMU), termed BMU-MoCo. Concretely, our
11 BMU-MoCo has two core designs: (1) Different from the conventional MoCo, we
12 apply the momentum update to not only momentum encoders but also encoders
13 (*i.e.*, bidirectional) at each training step, which enables the model to review the
14 learned knowledge retained in the momentum encoders. (2) To further enhance
15 our BMU-MoCo by utilizing earlier knowledge, we additionally maintain a pair of
16 global momentum encoders (only initialized at the very beginning) with the same
17 BMU strategy. Extensive results show that our BMU-MoCo remarkably outper-
18 forms recent competitors w.r.t. video-text retrieval performance and forgetting rate,
19 even without using any extra data or dynamic networks.

20 1 Introduction

21 Existing video-language modeling (VLM) methods have achieved promising performance for video-
22 text retrieval [56, 36, 58, 27, 22, 50, 25, 4] with non-streaming data. However, in real-world
23 application scenarios, VLM models need to evolve with streaming data (*e.g.*, collected from the
24 Internet [36, 39]) to accommodate more tasks. Under this setting, since it costs too much resource to
25 retrain the model with both old and new data for each task, a common practice is to fine-tune VLM
26 models with only the newly-arrived data. Note that such model fine-tuning leads to severe performance
27 degradation on previous tasks. This is a well-documented phenomenon called catastrophic forgetting
28 [16, 35] under the conventional continual learning setting [45, 42, 31, 14, 21, 57].

29 Therefore, in this work, we propose a continual video-language modeling (CVLM) setting to better
30 simulate the realistic scenario. Under our CVLM setting, models are supposed to be sequentially
31 trained on five widely-used video-text datasets: VATEX [51], ActivityNet [23], MSR-VTT [52],
32 DiDeMo [18], and MSVD [10]. An evaluation protocol is also established for CVLM, which contains
33 three metrics to respectively measure the text-to-video retrieval performance (Recall@1, shortened

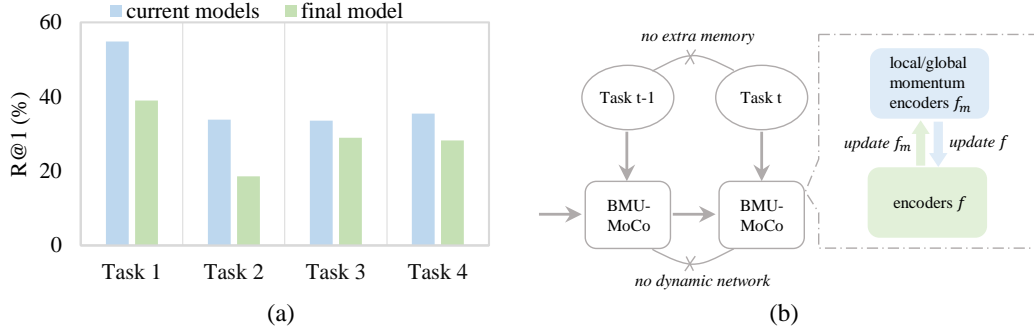


Figure 1: Illustration of the catastrophic forgetting problem in CVLM and the core design for our BMU-MoCo. **(a) The catastrophic forgetting problem in CVLM.** We train a basic cross-modal MoCo model on five tasks and present the comparative results of the final model and current models on learned tasks (Task 1–4). Note that there is no catastrophic forgetting on Task 5 and thus this task is omitted here. **(b) The core design of BMU-MoCo.** Different from the conventional MoCo, we update not only momentum encoders but also encoders through the bidirectional momentum update (BMU) strategy without extra memory or dynamic network across all tasks.

34 as R@1), forgetting rate (FR), and harmonic mean (HM) performance. Moreover, we implement
 35 a basic cross-modal MoCo [17] model (Base-MoCo) as our baseline method since it has shown
 36 superiority on video-language modeling [30, 32]. As illustrated in Figure 1(a), we observe that in
 37 spite of achieving great R@1 results with current models (evaluated right after trained on each task),
 38 the performance of the final Base-MoCo model (trained across all five tasks) drops significantly.

39 To tackle the catastrophic forgetting problem, most recent continual learning works attempt to
 40 preserve the learned knowledge from a variety of perspectives: (1) Maintaining a memory buffer to
 41 save and exploit data from previous tasks [43, 5, 31, 3, 8, 42]; (2) Generating pseudo data of learned
 42 tasks [48, 24, 40, 54]; (3) Extending the network architecture dynamically as each new task arrives
 43 [45, 14, 1, 7, 28]. However, when these methods are transferred to our CVLM setting, the resource
 44 consumption is enlarged rapidly as the number of tasks grows, due to the characteristic of video data.
 45 In addition, another branch of continual learning works focus on imposing a regularization constraint
 46 with quadratic penalty [21, 57, 13, 46] or knowledge distillation [29, 2, 19, 41, 20], which leads to
 47 an unwanted trade-off on the performance of old and new tasks with limited neural resources [37].
 48 Therefore, it is a long-standing and arduous challenge to train a video-language modeling network
 49 under the CVLM setting with both effectiveness and efficiency taken into consideration.

50 To overcome this challenge, we devise BMU-MoCo, a cross-modal MoCo-based model with a
 51 novel bidirectional momentum update (BMU) strategy. As shown in Figure 1(b), our BMU-MoCo
 52 needs neither extra memory data nor dynamically extended neural networks. Concretely, similar to
 53 cross-modal MoCo applied in [30, 32], our BMU-MoCo has a video encoder (*i.e.*, ViT-Base [12])
 54 and a text encoder (*i.e.*, BERT-Base [11]), followed by the momentum video/text encoders. Different
 55 from the original MoCo [17] and its cross-modal versions [30, 32] that utilize momentum update for
 56 only momentum encoders to maintain a consistent queue, our BMU strategy imposes momentum
 57 update on both momentum encoders and encoders. As a result, at each training step, the encoders of
 58 our BMU-MoCo learn new knowledge by end-to-end update with back-propagation whilst reviewing
 59 old knowledge directly from the parameters of momentum encoders by momentum update. In our
 60 opinion, our BMU-MoCo outperforms existing methods for two main reasons: (1) Momentum
 61 encoders are initialized by current encoders at the beginning of each new task and then progress
 62 slowly, which helps our model preserve adequate old knowledge without sacrificing the performance
 63 on new tasks; (2) Since there is no category information under the CVLM setting, learning from
 64 memory data or distilling with a batch of new data only absorbs *part of* previous knowledge while
 65 our BMU-MoCo learns *holistic* knowledge directly from the parameters of momentum encoders. To
 66 further enhance our BMU-MoCo, we also maintain a pair of global (cross-task) momentum encoders
 67 with the same BMU strategy, which are only initialized at the very beginning and thus preserve earlier
 68 knowledge than the normal local (task-specific) momentum encoders.

69 Our main contributions are four-fold: (1) We propose a new continual video-language modeling
70 (CVLM) setting, where models are supposed to be sequentially trained on five widely-used video-text
71 datasets. (2) To effectively and efficiently overcome the catastrophic forgetting problem under the
72 CVLM setting, we devise BMU-MoCo, a cross-modal MoCo-based model with a novel bidirectional
73 momentum update (BMU) strategy. It can review holistic old knowledge directly from the parameters
74 of momentum encoders while learning on new tasks. (3) To further boost our BMU-MoCo, a pair of
75 global momentum encoders are maintained by the same BMU strategy to preserve and review earlier
76 knowledge. (4) Extensive results demonstrate that our BMU-MoCo outperforms recent continual
77 learning methods by large margins w.r.t. both text-to-video retrieval performance and forgetting rate,
78 even without any extra memory data or dynamically extended networks.

79 2 Related Work

80 **Video-Language Modeling.** Most existing methods for video-language modeling follow two
81 paradigms: (1) Single-stream methods [33, 58, 49, 26, 25, 53] typically include a multi-modal trans-
82 former to achieve fine-grained cross-modal interaction between the video and language modalities.
83 Although achieving great performance, they suffer from the huge time complexity caused by the pair-
84 wise inputs during inference, which makes them unsuitable for practical applications. (2) Two-stream
85 methods [15, 38, 4, 30, 32] learn video and text representations independently, and align them after
86 encoding. To ensure the inference efficiency, both the baseline method (*i.e.*, Base-MoCo) and our
87 BMU-MoCo for CVLM are set to be two-stream methods. Importantly, different from Base-MoCo,
88 our BMU-MoCo has a novel BMU strategy to address the catastrophic forgetting problem and two
89 extra global momentum encoders to further boost the model performance.

90 **Continual Learning.** Conventional continual learning methods mainly focus on image classification
91 tasks. They can be roughly categorized into three groups: (1) *Rehearsal-based* methods apply extra
92 memory to store sampled data [42, 31, 43, 9, 5, 3, 8, 6] or generate pseudo data [48, 24, 40, 54]
93 from previous tasks. The memory size and the training complexity tend to be enlarged significantly
94 as the number of tasks grows. (2) *Expansion-based* methods either add extra extended networks
95 for new tasks [45, 14, 1, 7, 28] or select partial model parameters to update for different tasks
96 [55, 44, 47]. They need more computational resources especially for a long sequence of training tasks
97 (e.g., under our CVLM setting). (3) *Regularization-based* methods modify the model parameters
98 with quadratic loss penalty [21, 57, 13, 46] or knowledge distillation constraints [29, 2, 19, 41, 20].
99 Although succeeded in image classification tasks, they still face a large challenge in balancing the
100 model performance between old and new tasks when applied to our CVLM setting. Although our
101 BMU-MoCo can be classified as a regularization-based method, it has a vital difference from existing
102 regularization-based methods: benefiting from the bidirectional momentum updating process, our
103 BMU-MoCo can directly utilize the holistic previous knowledge from the parameters of momentum
104 encoders for model training (*i.e.*, updating the encoders), and *simultaneously* update the momentum
105 encoders at each training step to accommodate new tasks.

106 3 Methodology

107 3.1 Preliminary

108 We propose a new continual video-language modeling (CVLM) setting, where models are supposed
109 to be sequentially trained on n video-text datasets $\mathcal{D} = [\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_n]$. For each task t , it contains
110 a dataset $\mathcal{D}_t = \{V_i, T_i\}_{i=0}^{N_t-1}$ with N_t video-text pairs, where V_i denotes a video with S_i frames
111 and T_i represents an English text. The target of CVLM is to learn a video encoder f_{θ_V} and a
112 text encoder f_{θ_T} , which can respectively project the input video and its related text into a joint
113 embedding space with nearest metric distance. Different from the classical VLM setting which only
114 considers the model performance on the current dataset \mathcal{D}_t , our CVLM setting requires the models
115 to prevent the catastrophic forgetting on previously-used datasets $[\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_{t-1}]$ ($t > 1$) while
116 also performing well on the current dataset \mathcal{D}_t . Note that our proposed BMU-MoCo for CVLM is a
117 memory-free method which only utilizes the current dataset \mathcal{D}_t for each task t . Therefore, without
118 particular statement, we only consider task t with \mathcal{D}_t in the following subsections for simplicity.

119 **3.2 Network Architecture**

120 **Video Encoder.** We follow the most recent video-language modeling works [25, 53, 30, 32] to
 121 learn video representation by fusing the image embeddings of sampled frames per video. Concretely,
 122 given each video V_i with S_i frames, we randomly sample s frames ($s < S_i$) and embed them with an
 123 image encoder f_{img} (i.e., ViT-Base [12]) to obtain the frame embeddings:

$$F_{img}^{i,j} = f_{img}(x_i^j), j = 1, 2, \dots, s, \quad (1)$$

124 where x_i^j denotes the j -th sampled frames of video V_i and $F_{img}^{i,j}$ denotes its image embedding encoded
 125 by f_{img} . Then we project $F_{img}^{i,j}$ by a Linear layer f_{proj} :

$$F_{proj}^{i,j} = f_{proj}(F_{img}^{i,j}), j = 1, 2, \dots, s, \quad (2)$$

126 where $F_{proj}^{i,j} \in \mathbb{R}^d$ denotes the projected d -dimensional image embedding of $F_{img}^{i,j}$. Following COTS
 127 [32] and HiT [30], we obtain the final video embedding of V_i by adopting a fusing layer f_{avg} to
 128 aggregate the image embeddings $\{F_{img}^{i,j}\}$:

$$F_V^i = f_{avg}(F_{proj}^{i,1}, F_{proj}^{i,2}, \dots, F_{proj}^{i,s}), \quad (3)$$

129 where f_{avg} denotes an Average Pooling layer and $F_V^i \in \mathbb{R}^d$ is the video embedding of V_i . In summary,
 130 our video encoder f_{θ_V} encodes the video inputs by adopting f_{img} , f_{proj} and f_{avg} in Eqs. (1)–(3).

131 **Text Encoder.** For the language modality, we adopt BERT-Base [11] as our backbone to encode
 132 each input text T_i . In detail, we first tokenize T_i into a sequence of text tokens $[l_i^1, l_i^2, \dots, l_i^{r_i}]$, where
 133 r_i denotes the length of T_i . Then we obtain the text token embeddings through the backbone f_{bert} :

$$F_{bert}^i = f_{bert}(l_i^1, l_i^2, \dots, l_i^{r_i}), \quad (4)$$

134 where F_{bert}^i denotes the token embeddings of T_i obtained by f_{bert} . We then project them by a Linear
 135 layer \hat{f}_{proj} into the d -dimensional space as:

$$\hat{F}_{proj}^i = [\hat{F}_{proj}^i[1], \hat{F}_{proj}^i[2], \dots, \hat{F}_{proj}^i[r_i]] = \hat{f}_{proj}(F_{bert}^i[1], F_{bert}^i[2], \dots, F_{bert}^i[r_i]), \quad (5)$$

136 where $F_{bert}^i[j]$ denotes the j -th element of F_{bert}^i , and $\hat{F}_{proj}^i[j] \in \mathbb{R}^d$ represents the projected text
 137 embedding of token j in T_i (which has the same dimension d as video embedding F_V^i). To obtain the
 138 final text embedding of T_i , we apply an Average Pooling layer f_{avg} :

$$F_T^i = f_{avg}(\hat{F}_{proj}^i[1], \hat{F}_{proj}^i[2], \dots, \hat{F}_{proj}^i[r_i]), \quad (6)$$

139 where $F_T^i \in \mathbb{R}^d$ denotes the text embedding of T_i . In summary, our text encoder f_{θ_T} encodes the text
 140 inputs by adopting f_{bert} , \hat{f}_{proj} , and f_{avg} in Eqs. (4)–(6).

141 **3.3 BMU-MoCo**

142 **Cross-Modal MoCo.** Similar to the original single-modal MoCo [17], recent state-of-the-art video-
 143 language modeling works COTS [32] and HiT [30] construct a cross-modal MoCo architecture to
 144 maintain video/text momentum encoders by the same momentum update mechanism, which creates
 145 consistent queues for cross-modal contrastive learning objectives. As illustrated in Figure 2, our
 146 BMU-MoCo follows this paradigm and further transfers it to our CVLM setting. Concretely, for a
 147 mini-batch of N_B video-text pairs $\mathcal{B} = \{V_i, T_i\}_{i=1}^{N_B}$, we first obtain the query embeddings q_i^V, q_i^T of
 148 V_i, T_i by video encoder f_{θ_V} and text encoder f_{θ_T} :

$$q_i^V = f_{\theta_V}(V_i), \quad q_i^T = f_{\theta_T}(T_i). \quad (7)$$

149 Then we maintain two momentum encoders $f_{\theta_{V,m}}, f_{\theta_{T,m}}$ (termed local momentum video/text en-
 150 coders in Figure 2) for both video and text modalities, whose parameters $\theta_{V,m}, \theta_{T,m}$ are initialized
 151 by θ_V, θ_T at the beginning of each task t . During the training process, $\theta_{V,m}, \theta_{T,m}$ are continuously
 152 updated by θ_V, θ_T with the momentum update strategy:

$$\theta_{V,m} = m \cdot \theta_{V,m} + (1 - m) \cdot \theta_V, \quad \theta_{T,m} = m \cdot \theta_{T,m} + (1 - m) \cdot \theta_T, \quad (8)$$

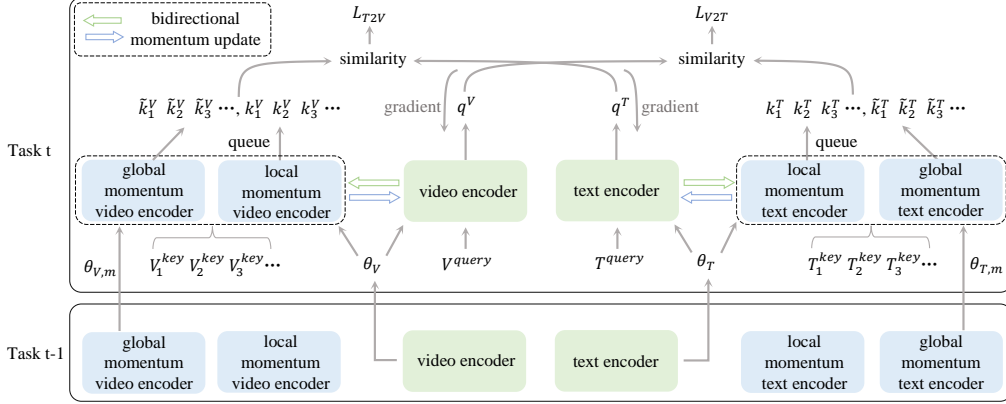


Figure 2: Schematic illustration of our BMU-MoCo. The momentum update strategy is applied to both encoders and momentum encoders (*i.e.*, bidirectional). To exploit earlier knowledge, we further maintain a pair of global momentum encoders with the same BMU strategy, whose parameters are inherited across tasks and only initialized at the very beginning.

153 where m is the coefficient of momentum update. To form the contrastive learning loss of cross-modal
 154 MoCo, we need two consistent queues to preserve the negative video/text samples. In detail, the key
 155 embeddings k_i^V , k_i^T of V_i , T_i are firstly acquired by momentum video and text encoders:

$$k_i^V = f_{\theta_{V,m}}(V_i), \quad k_i^T = f_{\theta_{T,m}}(T_i). \quad (9)$$

156 We then respectively push k_i^V and k_i^T into the negative video queue Q^V and the negative text queue
 157 Q^T (after computing loss), where $Q^V = \{k_1^V, k_2^V, k_3^V, \dots, k_{N_Q}^V\}$ and $Q^T = \{k_1^T, k_2^T, k_3^T, \dots, k_{N_Q}^T\}$
 158 (N_Q is the queue size). The contrastive losses of cross-modal MoCo (Base-MoCo) are:

$$\hat{L}_{V2T} = -\frac{1}{N_B} \sum_{i=1}^{N_B} \log \frac{\exp(\frac{q_i^V \cdot k_i^T}{\tau})}{\exp(\frac{q_i^V \cdot k_i^T}{\tau}) + \sum_{j=1}^{N_Q} \exp(\frac{q_i^V \cdot k_j^T}{\tau})}, \quad (10)$$

$$\hat{L}_{T2V} = -\frac{1}{N_B} \sum_{i=1}^{N_B} \log \frac{\exp(\frac{q_i^T \cdot k_i^V}{\tau})}{\exp(\frac{q_i^T \cdot k_i^V}{\tau}) + \sum_{j=1}^{N_Q} \exp(\frac{q_i^T \cdot k_j^V}{\tau})}, \quad (11)$$

159 where τ is the temperature. Note that the queue size N_Q is decoupled from the batch size N_B .
 160 Therefore, it can take a large value for better representation of the data distribution.

161 **Bidirectional Momentum Update.** Although achieving great success with non-streaming data
 162 (*e.g.*, a single video-text dataset), the original cross-modal MoCo has difficulty in coping with the
 163 catastrophic forgetting problem under our CVLM setting. To overcome this difficulty, we propose
 164 a novel bidirectional momentum update (BMU) strategy for cross-modal MoCo to review the old
 165 knowledge retained in momentum encoders at each training step. Concretely, for video/text encoders
 166 $f_{\theta_V}, f_{\theta_T}$, in addition to the end-to-end update by back-propagation, we further update their parameters
 167 θ_V, θ_T using the parameters $\theta_{V,m}, \theta_{T,m}$ of momentum encoders $f_{\theta_{V,m}}, f_{\theta_{T,m}}$ by momentum update:

$$\theta_V = \hat{m} \cdot \theta_V + (1 - \hat{m}) \cdot \theta_{V,m}, \quad \theta_T = \hat{m} \cdot \theta_T + (1 - \hat{m}) \cdot \theta_{T,m}, \quad (12)$$

168 where \hat{m} is a momentum coefficient, and $\theta_{V,m}, \theta_{T,m}$ are simultaneously updated by Eq. (8). Together,
 169 Eq. (8) and Eq. (12) compose our BMU strategy. Note that the advantages of BMU lie in two aspects:
 170 (1) At the beginning of each new task t , $\theta_{V,m}$ and $\theta_{T,m}$ are respectively initialized by θ_V and θ_T ,
 171 which makes the knowledge of task $t-1$ be preserved. (2) During the training process, $\theta_{V,m}$ and
 172 $\theta_{T,m}$ are constantly and slowly updated by the momentum update strategy, which enables our model
 173 to review the old knowledge but without sacrificing the performance on new tasks.

174 **Global Momentum Encoders.** To further enhance our BMU-MoCo, we propose to maintain a
 175 pair of global momentum encoders which can preserve earlier knowledge. As shown in Figure 2,
 176 they are only initialized at the very beginning of the whole training process under our CVLM setting,

177 and their parameters are transmitted across tasks. Formally, let $f_{\tilde{\theta}_{V,m}}$ and $f_{\tilde{\theta}_{T,m}}$ denote the global
 178 momentum video and text encoders, respectively. Their parameters $\tilde{\theta}_{V,m}$ and $\tilde{\theta}_{T,m}$ are updated by
 179 the BMU strategy along with the parameters θ_V and θ_T of encoders:

$$\theta_V = \hat{m} \cdot \theta_V + (1 - \hat{m}) \cdot \tilde{\theta}_{V,m}, \quad \theta_T = \hat{m} \cdot \theta_T + (1 - \hat{m}) \cdot \tilde{\theta}_{T,m}, \quad (13)$$

$$\tilde{\theta}_{V,m} = m \cdot \tilde{\theta}_{V,m} + (1 - m) \cdot \theta_V, \quad \tilde{\theta}_{T,m} = m \cdot \tilde{\theta}_{T,m} + (1 - m) \cdot \theta_T. \quad (14)$$

180 Note that Eq. (12) and Eq. (13) are implemented subsequently. For each video-text input $\{V_i, T_i\}$,
 181 we obtain a new group of key embeddings $\tilde{k}_i^V, \tilde{k}_i^T$ with the global momentum encoders $f_{\tilde{\theta}_{V,m}}, f_{\tilde{\theta}_{T,m}}$:

$$\tilde{k}_i^V = f_{\tilde{\theta}_{V,m}}(V_i), \quad \tilde{k}_i^T = f_{\tilde{\theta}_{T,m}}(T_i). \quad (15)$$

182 We push \tilde{k}_i^V and \tilde{k}_i^T respectively into two negative queues \tilde{Q}^V and \tilde{Q}^T , where $\tilde{Q}^V =$
 183 $\{\tilde{k}_1^V, \tilde{k}_2^V, \tilde{k}_3^V, \dots, \tilde{k}_{N_Q}^V\}$, $\tilde{Q}^T = \{\tilde{k}_1^T, \tilde{k}_2^T, \tilde{k}_3^T, \dots, \tilde{k}_{N_Q}^T\}$. Note that each query embedding (e.g.,
 184 q_i^T) has two corresponding positive embeddings (k_i^V, \tilde{k}_i^V) and two corresponding negative queues
 185 (Q^V, \tilde{Q}^V). The cross-modal contrastive losses are defined as:

$$L_{V2T} = -\frac{1}{N_B} \sum_{i=1}^{N_B} \log \frac{\exp(\frac{q_i^V \cdot k_i^T}{\tau}) + \exp(\frac{q_i^V \cdot \tilde{k}_i^T}{\tau})}{\exp(\frac{q_i^V \cdot k_i^T}{\tau}) + \exp(\frac{q_i^V \cdot \tilde{k}_i^T}{\tau}) + \sum_{j=1}^{N_Q} [\exp(\frac{q_i^V \cdot k_j^T}{\tau}) + \exp(\frac{q_i^V \cdot \tilde{k}_j^T}{\tau})]}, \quad (16)$$

$$L_{T2V} = -\frac{1}{N_B} \sum_{i=1}^{N_B} \log \frac{\exp(\frac{q_i^T \cdot k_i^V}{\tau}) + \exp(\frac{q_i^T \cdot \tilde{k}_i^V}{\tau})}{\exp(\frac{q_i^T \cdot k_i^V}{\tau}) + \exp(\frac{q_i^T \cdot \tilde{k}_i^V}{\tau}) + \sum_{j=1}^{N_Q} [\exp(\frac{q_i^T \cdot k_j^V}{\tau}) + \exp(\frac{q_i^T \cdot \tilde{k}_j^V}{\tau})]}, \quad (17)$$

186 where τ is the temperature. Now we have the final loss of BMU-MoCo for our CVLM setting:

$$L_{final} = L_{V2T} + L_{T2V}. \quad (18)$$

187 The full (pseudocode) algorithm of our BMU-MoCo is presented in the supplementary material.

188 4 Experiments

189 4.1 Experimental Setup

190 **Datasets.** Our CVLM setting is defined over a sequence of five video-text datasets: (1) VATEX [51]
 191 is a large-scale open-domain dataset, which has 25,991 videos with 250K text descriptions for training,
 192 3,000 videos for validation and 6,000 videos for testing. (2) ActivityNet [23] is an action domain
 193 dataset, which consists of 20K YouTube videos with 100K text descriptions. We follow the standard
 194 setting in [4, 25] to use 10K videos for training and 4.9K for test (the val1 split), where all texts of
 195 each video are concatenated into one query paragraph. (3) MSR-VTT [52] contains 10K videos,
 196 with 20 text descriptions per video. We follow the 1k-A split in recent works [25, 4, 53, 30] with
 197 9K training videos and 1K test videos. (4) DiDeMo [18] consists of 10K Flickr videos with 40K
 198 text annotations. Following [4, 25], we train and evaluate our model on paragraph-to-video retrieval
 199 (the same setting for ActivityNet). (5) MSVD [10] has 1,200 videos with 48K texts for training, 100
 200 videos for validation and 670 ones for testing. Overall, there are around 50K videos with 500K text
 201 descriptions in all five datasets (*i.e.*, each dataset has 100K video-text pairs in average).

202 **Evaluation Metrics.** Similar to the standard video-language modeling setting, we evaluate the
 203 text-to-video retrieval performance of a model on Recall@1 (shortened as R@1). R@1 refers to the
 204 percentage of text queries that correctly retrieve the ground-truth candidate at top-1. For our CVLM
 205 setting, we further define two evaluation metrics: forgetting rate (FR) and harmonic mean (HM).
 206 Formally, let \mathcal{M}_i denotes the model after trained on task i and \mathcal{A}_t^i ($t \leq i$) denotes the R@1 result of
 207 \mathcal{M}_i on task t . The overall R@1 ($\frac{1}{n} \sum_{t=1}^n \mathcal{A}_t^n$) is the average R@1 results of the final model \mathcal{M}_n on
 208 all tasks. Based on these notations, we then define the FR and HM as follows:

209 (1) **Forgetting rate (FR)** of \mathcal{M}_i on task t ($t \leq i$) is the performance decrease between \mathcal{M}_i and \mathcal{M}_t :
 210 $\text{FR} = \mathcal{A}_t^i - \mathcal{A}_t^t$, where a lower FR indicates the model forgets less knowledge. Note that there is no
 211 catastrophic forgetting (FR = 0) when $t = i$. The Overall FR is obtained by just summing the FR
 212 results of the final model \mathcal{M}_n across all tasks: Overall FR = $\sum_{t=1}^n (\mathcal{A}_t^n - \mathcal{A}_t^t)$.

Table 1: Comparative results obtained by the final model \mathcal{M}_n on five video-text datasets/tasks under our CVLM setting: VATEX [51] (*i.e.*, Task1), ActivityNet [23] (*i.e.*, Task2), MSR-VTT [52] (*i.e.*, Task3), DiDeMo [18] (*i.e.*, Task4), MSVD [10] (*i.e.*, Task5). For fair comparison, all baseline models are re-implemented based on the same cross-modal MoCo architecture for our CVLM setting. \dagger denotes applying extra encoders, including encoders from the last task (*e.g.*, LwF [29]) and global momentum encoders (*e.g.*, our BMU-MoCo). ‘Mem.’ denotes applying memory buffer during training. ‘BMU-MoCo (local)’ denotes BMU-MoCo without global momentum encoders.

Method	Mem.	Task1		Task2		Task3		Task4		Task5		Overall	
		R@1 \uparrow	FR \downarrow	R@1 \uparrow	FR \downarrow	R@1 \uparrow	FR \downarrow	R@1 \uparrow	FR \downarrow	R@1 \uparrow	FR \downarrow	R@1 \uparrow	FR \downarrow
Base-MoCo [32]	No	38.99	15.30	18.61	15.23	28.00	5.60	28.22	7.27	40.28	30.82	43.40	34.67
LwF \dagger [29]	No	42.02	12.27	19.95	14.87	28.90	6.00	29.91	6.98	37.91	32.24	40.12	35.81
ER-ring [9]	Yes	41.99	12.30	22.09	11.79	29.80	5.40	30.31	5.08	38.53	32.54	34.57	35.67
DER [5]	Yes	40.15	14.14	21.35	12.65	28.80	5.00	30.71	4.09	39.96	32.19	35.88	35.39
Co2L \dagger [6]	Yes	41.23	13.06	21.74	13.06	27.50	5.30	30.41	5.38	39.29	31.58	34.48	35.06
LUMP \dagger [34]	Yes	40.16	14.13	21.78	12.37	30.50	3.00	29.91	4.99	39.39	32.45	34.49	35.56
BMU-MoCo (local)	No	46.82	7.47	23.27	10.84	30.00	3.40	31.21	4.08	41.94	34.65	25.79	37.05
BMU-MoCo \dagger	No	48.48	5.81	23.45	10.43	30.80	2.90	32.80	3.49	41.83	35.47	22.63	37.59

(2) **Harmonic mean (HM)** calculates the harmonic mean of the overall R@1 (current) and the overall R@1 (final), where the overall R@1 (current) denotes the average of the R@1 values obtained by each current model \mathcal{M}_i ($i = 1, 2, \dots, n$) on each current task i , and the overall R@1 (final) denotes the average R@1 for the final model \mathcal{M}_n on all tasks. Formally, we have:
$$\text{HM} = \frac{2 \cdot \frac{1}{n} \sum_{i=1}^n \mathcal{A}_i^i \cdot \frac{1}{n} \sum_{i=1}^n \mathcal{A}_i^n}{\frac{1}{n} \sum_{i=1}^n \mathcal{A}_i^i + \frac{1}{n} \sum_{i=1}^n \mathcal{A}_i^n}$$
 Note that HM can alleviate the trade-off problem between overall R@1 (current) and overall R@1 (final), which is otherwise an inherent limitation of FR. Specifically, when a model has a lower overall R@1 (current) and a lower overall R@1 (final), it could also have a better/lower FR (which is unsatisfactory) but still lead to a worse/lower HM (see Figure 4(c)).

Implementation Details. Recent works show that video-language models benefit from image-text pre-training [25, 4], which can accelerate the model convergence and is more suitable for the practical scenarios. We thus apply ViT-Base [12]/BERT-Base [11] as our image/text encoder and follow recent state-of-the-art MoCo-based model COTS [32] to pre-train our model with 5.3M image-text pairs. We then sequentially train all the models (BMU-MoCo and all competitors) on five video-text datasets/tasks. For each task, we train a model for 10 epochs and choose the best trained one w.r.t. the validation R@1 results. For those competitors using a memory buffer during training (*e.g.*, ER-ring [9]), we set the memory size to 10% of the average data size 100K (*i.e.*, 10K video-text pairs). Note that the percentage 10% is larger than the buffer size of most recent rehearsal-based continual learning methods [9, 5, 34]. More details are given as follows: (1) In the training phase, all sampled frames of each video are resized to 384×384 and augmented by gray-scaling and color-jitter. (2) For the first epoch of each task under our CVLM setting, we set the learning rate to $5e-5$ and decay it to $5e-6$ afterwards. (3) We select the two momentum coefficients $m = 0.99$, $\hat{m} = 0.99$, and the temperature $\tau = 0.07$. We set the batch size N_B to 48 and the queue size N_Q to 1,440. (4) The total training time on five tasks is around 20 hours with 8 Tesla V100 GPUs for each model.

4.2 Main Results

Table 1 summarizes the comparative results in terms of text-to-video retrieval (R@1), forgetting rate and harmonic mean (HM) obtained by the final model \mathcal{M}_n (per method) on five datasets. We re-implement five recent continual learning methods (fused with cross-modal MoCo) under our CVLM setting, including rehearsal-based methods (ER-ring [9], DER [5]), regularization-based methods (LwF [29]) and their combinations (Co2L [6], LUMP [34]). We can observe that: (1) Our BMU-MoCo outperforms recent methods by large margins without using any extra memory or dynamically extended networks. Concretely, our method achieves the best R@1 and FR results on all tasks, and outperforms the second best by 2.93% for overall R@1, 11.85% for overall FR and 1.78% for overall HM. (2) Without applying global momentum encoders, our BMU-MoCo (local) also beats all competitors, directly showing the effectiveness of our BMU strategy. (3) The improvements over Base-MoCo obtained by utilizing knowledge distillation (*e.g.*, LwF [29]) or extra memory data (*e.g.*, Co2L [6]) are limited due to the lack of category information under our CVLM setting.

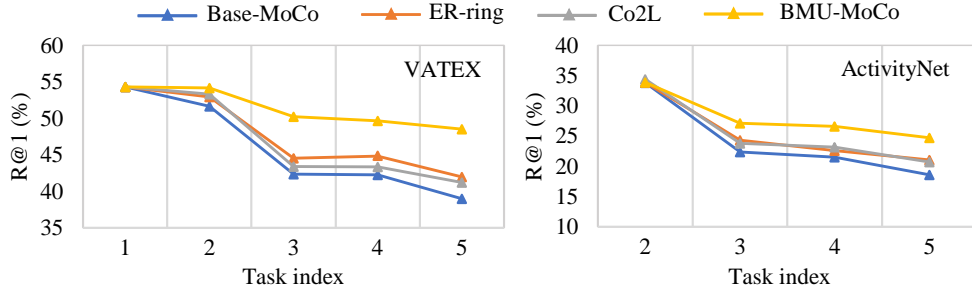


Figure 3: Detailed comparative results for text-to-video retrieval (R@1) obtained by each model \mathcal{M}_i (per method) on the first two tasks: VATEX and ActivityNet. Note that the gap between the *beginning* and the *end* of each line denotes the forgetting rate.

Table 2: Ablation study results for our BMU-MoCo. ‘Local’ denotes applying the local momentum encoders, while ‘Global’ denotes applying the global momentum encoders. Results for text-to-video retrieval (R@1), forgetting rate (FR) and harmonic mean (HM) are reported.

		Task1		Task2		Task3		Task4		Task5		Overall		
BMU	Local	Global	R@1↑	FR↓	R@1↑	FR↓	R@1↑	FR↓	R@1↑	FR↓	R@1↑	R@1↑	FR↓	HM↑
	✓		38.99	15.30	18.61	15.23	28.00	5.60	28.22	7.27	40.28	30.82	43.40	34.67
		✓	38.95	15.97	18.51	15.45	30.10	4.60	28.91	6.78	39.89	31.47	42.75	35.16
	✓	✓	41.44	12.85	21.68	13.22	29.40	4.90	29.31	6.08	39.96	32.35	36.95	35.67
✓	✓		46.82	7.47	23.27	10.84	30.00	3.40	31.21	4.08	41.94	34.65	25.79	37.05
✓		✓	46.35	7.94	23.16	10.99	30.60	3.70	31.41	5.28	41.70	34.64	27.91	37.22
✓	✓	✓	48.48	5.81	23.45	10.43	30.80	2.90	32.80	3.49	41.83	35.47	22.63	37.59

249 Figure 3 shows more detailed comparative results for text-to-video retrieval (R@1) obtained by each
250 model \mathcal{M}_i (per method) on the first two datasets (VATEX [51] and ActivityNet [23]). We compare
251 our BMU-MoCo with three representative competitors, including Base-MoCo [17], ER-ring [9], and
252 Co2L [6]. Concretely, the left sub-figure presents the results of $\mathcal{M}_1 \sim \mathcal{M}_5$ (per method)
253 on task 1 (VATEX), *i.e.*, \mathcal{A}_1^i ($1 \leq i \leq 5$). The right sub-figure presents the results of $\mathcal{M}_2 \sim \mathcal{M}_5$ (per method)
254 on task 2 (ActivityNet), *i.e.*, \mathcal{A}_2^i ($2 \leq i \leq 5$). It can be observed that: (1) For task 1 (VATEX), the
255 performance of our BMU-MoCo drops the most slowly after it is trained on the following tasks (task
256 2 to task 5). (2) For task 2 (ActivityNet), our BMU-MoCo also leads to the slowest performance drop
257 after trained on the following tasks (task 3 to task 5). Overall, our BMU-MoCo indeed significantly
258 alleviates the performance decrease problem during the whole training process.

259 4.3 Ablation Study

260 We first analyze the contributions of the BMU strategy, the local momentum encoders and the global
261 momentum encoders applied in our BMU-MoCo. The ablative results are shown in Table 2. It can be
262 clearly seen that: (1) With our BMU strategy, our model achieves remarkable improvements (4th row
263 vs. 1st row). (2) Simultaneously applying local and global encoders is better than using only one of
264 them (3rd row vs. 1st/2nd row), which indicates that the knowledge preserved in local and global
265 momentum encoders are different. (3) Our BMU strategy helps our model to excavate knowledge
266 preserved in different momentum encoders (6th row vs. 3rd row) and achieve the best performance
267 (6th row vs. 4th/5th row), which further validates the effectiveness of our BMU.

268 Considering the core role of BMU, we thus analyze the impact of the momentum coefficient \hat{m}
269 utilized in our BMU-MoCo. According to COTS [32], the other momentum coefficient m of our
270 model is fixed at 0.99 (only the value of \hat{m} is changed). Figure 4 shows the results for overall R@1
271 (current), overall R@1 (final), and overall HM, respectively. We find that the value of \hat{m} cannot be
272 too big or too small. Concretely, when \hat{m} is too big (*e.g.*, 0.999 and 1), the knowledge preserved
273 in momentum encoders cannot be well-reviewed by our model. When \hat{m} is too small (*e.g.*, 0.9),
274 the end-to-end update by back-propagation is influenced too much, which leads to bad results for
275 overall R@1 (current) and overall HM. It is worth mentioning that the model with smaller \hat{m} (0.9)

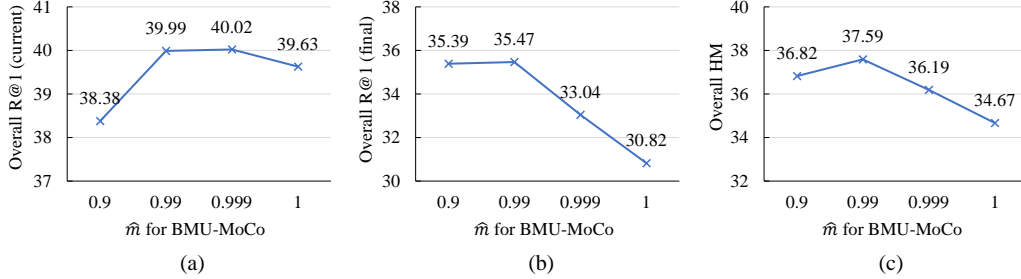


Figure 4: Comparative results of our BMU-MoCo with different momentum coefficient \hat{m} . The results for overall R@1 (current), overall R@1 (final) and overall HM are reported in (a-c), respectively.

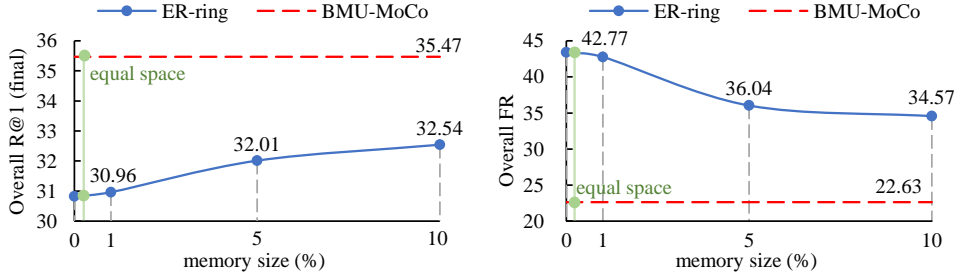


Figure 5: Comparative results obtained by different memory size. The red dotted line denotes our BMU-MoCo (with no memory) while the blue line denotes the representative rehearsal-based method ER-ring. The green line suggests that the storage space consumption of our BMU-MoCo (with global momentum encoders) is equal to ER-ring (with 0.05% memory to store videos).

276 has lower/better FR (14.95) since it sacrifices the model performance on overall R@1 (current). This
 277 phenomenon demonstrates the necessity of utilizing the overall HM to measure the overall (trade-off)
 278 model performance. Therefore, we set the momentum coefficient \hat{m} to 0.99 in all our experiments,
 279 which helps our model to review old knowledge while learning well on new tasks.

280 4.4 Further Evaluation

281 To demonstrate both the efficiency and effectiveness of our BMU-MoCo under our CVLM setting, we
 282 compare our model with a representative rehearsal-based method ER-ring [9] by different memory
 283 size in Figure 5. Note that our BMU-MoCo has two global momentum encoders that need 0.5GB
 284 more storage space than the original cross-modal MoCo (used by all competitors including ER-ring).
 285 As shown in Figure 5, when the memory size of ER-ring becomes 0.05%, it equals to the size of extra
 286 storage space used by our BMU-MoCo (but our model performs significantly better). In real-world
 287 application scenarios, the memory size of rehearsal-based methods like ER-ring enlarges rapidly as
 288 the number of tasks grows, while the fixed extra space size (0.5GB) of our BMU-MoCo is negligible.
 289 More importantly, our BMU-MoCo (with a fixed 0.5GB sapce size) even outperforms ER-ring using
 290 10% memory (about 200GB under our CVLM setting) by large margins for both overall R@1 (final)
 291 and overall FR. This directly indicates the efficiency and effectiveness of our BMU-MoCo.

292 5 Conclusion

293 In this paper, we propose a new continual video-language modeling (CVLM) setting, where models
 294 are supposed to be sequentially trained on five widely-used video-text datasets. To overcome the
 295 catastrophic forgetting and heavy resource consumption challenges, we propose a novel framework
 296 BMU-MoCo, which is a cross-modal MoCo-based model with bidirectional momentum update
 297 (BMU). We maintain both local and global momentum encoders with our BMU strategy to review
 298 broader old knowledge while learning on new tasks. Extensive experimental results show that our
 299 BMU-MoCo outperforms recent competitors by large margins, even without using extra memory data
 300 or dynamically extended networks. The limitation of our work lies in that we have only evaluated
 301 BMU-MoCo under the CVLM setting, and thus we need to transfer it to other continual learning
 302 settings (*e.g.*, continual image-text pre-training) for comprehensive study.

References

- 303
- 304 [1] Ferran Alet, Tomás Lozano-Pérez, and Leslie P Kaelbling. Modular meta-learning. In *CoRL*, pages
305 856–868. PMLR, 2018.
- 306 [2] Rahaf Aljundi, Punarjay Chakravarty, and Tinne Tuytelaars. Expert gate: Lifelong learning with a network
307 of experts. In *CVPR*, pages 3366–3375, 2017.
- 308 [3] Rahaf Aljundi, Min Lin, Baptiste Goujaud, and Yoshua Bengio. Gradient based sample selection for online
309 continual learning. *NeurIPS*, 32, 2019.
- 310 [4] Max Bain, Arsha Nagrani, Gül Varol, and Andrew Zisserman. Frozen in time: A joint video and image
311 encoder for end-to-end retrieval. *arXiv preprint arXiv:2104.00650*, 2021.
- 312 [5] Pietro Buzzega, Matteo Boschini, Angelo Porrello, Davide Abati, and Simone Calderara. Dark experience
313 for general continual learning: a strong, simple baseline. *NeurIPS*, 33:15920–15930, 2020.
- 314 [6] Hyuntak Cha, Jaeho Lee, and Jinwoo Shin. Co2l: Contrastive continual learning. In *ICCV*, pages
315 9516–9525, 2021.
- 316 [7] Michael B Chang, Abhishek Gupta, Sergey Levine, and Thomas L Griffiths. Automatically composing
317 representation transformations as a means for generalization. *arXiv preprint arXiv:1807.04640*, 2018.
- 318 [8] Arslan Chaudhry, Albert Gordo, Puneet K Dokania, Philip Torr, and David Lopez-Paz. Using hindsight to
319 anchor past knowledge in continual learning. *arXiv preprint arXiv:2002.08165*, 3, 2020.
- 320 [9] Arslan Chaudhry, Marcus Rohrbach, Mohamed Elhoseiny, Thalaiyasingam Ajanthan, Puneet K Dokania,
321 Philip HS Torr, and Marc’Aurelio Ranzato. On tiny episodic memories in continual learning. *arXiv
322 preprint arXiv:1902.10486*, 2019.
- 323 [10] David Chen and William B Dolan. Collecting highly parallel data for paraphrase evaluation. In *ACL*, pages
324 190–200, 2011.
- 325 [11] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep
326 bidirectional transformers for language understanding. In *NAACL-HLT*, pages 4171–4186, 2019.
- 327 [12] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
328 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and
329 Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*,
330 2021.
- 331 [13] Mehrdad Farajtabar, Navid Azizan, Alex Mott, and Ang Li. Orthogonal gradient descent for continual
332 learning. In *AISTATS*, pages 3762–3773. PMLR, 2020.
- 333 [14] Chrisantha Fernando, Dylan Banarse, Charles Blundell, Yori Zwols, David Ha, Andrei A Rusu, Alexander
334 Pritzel, and Daan Wierstra. Pathnet: Evolution channels gradient descent in super neural networks. *arXiv
335 preprint arXiv:1701.08734*, 2017.
- 336 [15] Simon Ging, Mohammadreza Zolfaghari, Hamed Pirsiavash, and Thomas Brox. Coot: Cooperative
337 hierarchical transformer for video-text representation learning. *NeurIPS*, 33:22605–22618, 2020.
- 338 [16] Ian J Goodfellow, Mehdi Mirza, Da Xiao, Aaron Courville, and Yoshua Bengio. An empirical investigation
339 of catastrophic forgetting in gradient-based neural networks. *arXiv preprint arXiv:1312.6211*, 2013.
- 340 [17] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross B. Girshick. Momentum contrast for unsuper-
341 vised visual representation learning. In *CVPR*, pages 9726–9735, 2020.
- 342 [18] Anne Lisa Hendricks, Oliver Wang, Eli Shechtman, Josef Sivic, Trevor Darrell, and Bryan Russell.
343 Localizing moments in video with natural language. In *ICCV*, pages 5804–5813, 2017.
- 344 [19] Saihui Hou, Xinyu Pan, Chen Change Loy, Zilei Wang, and Dahua Lin. Lifelong learning via progressive
345 distillation and retrospection. In *ECCV*, pages 437–452, 2018.
- 346 [20] Saihui Hou, Xinyu Pan, Chen Change Loy, Zilei Wang, and Dahua Lin. Learning a unified classifier
347 incrementally via rebalancing. In *CVPR*, pages 831–839, 2019.
- 348 [21] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu,
349 Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic
350 forgetting in neural networks. *PNAS*, 114(13):3521–3526, 2017.

- 351 [22] Bruno Korbar, Fabio Petroni, Rohit Girdhar, and Lorenzo Torresani. Video understanding as machine
352 translation. *arXiv preprint arXiv:2006.07203*, 2020.
- 353 [23] Ranjay Krishna, Kenji Hata, Frederic Ren, Li Fei-Fei, and Juan Carlos Niebles. Dense-captioning events
354 in videos. In *ICCV*, pages 706–715, 2017.
- 355 [24] Frantzeska Lavda, Jason Ramapuram, Magda Gregorova, and Alexandros Kalousis. Continual classification
356 learning using generative models. *arXiv preprint arXiv:1810.10612*, 2018.
- 357 [25] Jie Lei, Linjie Li, Luowei Zhou, Zhe Gan, Tamara L Berg, Mohit Bansal, and Jingjing Liu. Less is more:
358 ClipBERT for video-and-language learning via sparse sampling. *CVPR*, pages 7331–7341, 2021.
- 359 [26] Gen Li, Nan Duan, Yuejian Fang, Ming Gong, and Daxin Jiang. Unicoder-vl: A universal encoder for
360 vision and language by cross-modal pre-training. In *AAAI*, volume 34, pages 11336–11344, 2020.
- 361 [27] Linjie Li, Yen-Chun Chen, Yu Cheng, Zhe Gan, Licheng Yu, and Jingjing Liu. HERO: Hierarchical
362 encoder for video+ language omni-representation pre-training. *EMNLP*, pages 2046–2065, 2020.
- 363 [28] Xilai Li, Yingbo Zhou, Tianfu Wu, Richard Socher, and Caiming Xiong. Learn to grow: A continual
364 structure learning framework for overcoming catastrophic forgetting. In *ICML*, pages 3925–3934. PMLR,
365 2019.
- 366 [29] Zhizhong Li and Derek Hoiem. Learning without forgetting. *TPAMI*, 40(12):2935–2947, 2017.
- 367 [30] Song Liu, Haoqi Fan, Shengsheng Qian, Yiru Chen, Wenkui Ding, and Zhongyuan Wang. HiT: Hierarchical
368 transformer with momentum contrast for video-text retrieval. In *ICCV*, pages 11915–11925, 2021.
- 369 [31] David Lopez-Paz and Marc’Aurelio Ranzato. Gradient episodic memory for continual learning. *NeurIPS*,
370 30, 2017.
- 371 [32] Haoyu Lu, Nanyi Fei, Yuqi Huo, Yizhao Gao, Zhiwu Lu, and Ji-Rong Wen. COTS: Collaborative two-
372 stream vision-language pre-training model for cross-modal retrieval. *arXiv preprint arXiv:2204.07441*,
373 2022.
- 374 [33] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic
375 representations for vision-and-language tasks. *NeurIPS*, 32, 2019.
- 376 [34] Divyam Madaan, Jaehong Yoon, Yuanchun Li, Yunxin Liu, and Sung Ju Hwang. Representational
377 continuity for unsupervised continual learning. In *ICLR*, 2021.
- 378 [35] Michael McCloskey and Neal J Cohen. Catastrophic interference in connectionist networks: The sequential
379 learning problem. In *NLM*, volume 24, pages 109–165. Elsevier, 1989.
- 380 [36] Antoine Miech, Dimitri Zhukov, Jean-Baptiste Alayrac, Makarand Tapaswi, Ivan Laptev, and Josef Sivic.
381 HowTo100M: Learning a text-video embedding by watching hundred million narrated video clips. In
382 *ICCV*, pages 2630–2640, 2019.
- 383 [37] German I Parisi, Ronald Kemker, Jose L Part, Christopher Kanan, and Stefan Wermter. Continual lifelong
384 learning with neural networks: A review. *Neural Networks*, 113:54–71, 2019.
- 385 [38] Mandela Patrick, Po-Yao Huang, Yuki Asano, Florian Metze, Alexander G Hauptmann, Joao F Henriques,
386 and Andrea Vedaldi. Support-set bottlenecks for video-text representation learning. In *ICLR*, 2020.
- 387 [39] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish
388 Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning
389 transferable visual models from natural language supervision. In *ICML*, pages 8748–8763, 2021.
- 390 [40] Jason Ramapuram, Magda Gregorova, and Alexandros Kalousis. Lifelong generative modeling. *Neuro-
391 computing*, 404:381–400, 2020.
- 392 [41] Amal Rannen, Rahaf Aljundi, Matthew B Blaschko, and Tinne Tuytelaars. Encoder based lifelong learning.
393 In *ICCV*, pages 1320–1328, 2017.
- 394 [42] Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl: Incremental
395 classifier and representation learning. In *CVPR*, pages 2001–2010, 2017.
- 396 [43] Matthew Riemer, Ignacio Cases, Robert Ajemian, Miao Liu, Irina Rish, Yuhai Tu, and Gerald Tesauro.
397 Learning to learn without forgetting by maximizing transfer and minimizing interference. *arXiv preprint
398 arXiv:1810.11910*, 2018.

- 399 [44] Clemens Rosenbaum, Tim Klinger, and Matthew Riemer. Routing networks: Adaptive selection of
400 non-linear functions for multi-task learning. *arXiv preprint arXiv:1711.01239*, 2017.
- 401 [45] Andrei A Rusu, Neil C Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Ko-
402 ray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive neural networks. *arXiv preprint*
403 *arXiv:1606.04671*, 2016.
- 404 [46] Gobinda Saha, Isha Garg, and Kaushik Roy. Gradient projection memory for continual learning. *arXiv*
405 *preprint arXiv:2103.09762*, 2021.
- 406 [47] Joan Serra, Didac Suris, Marius Miron, and Alexandros Karatzoglou. Overcoming catastrophic forgetting
407 with hard attention to the task. In *ICML*, pages 4548–4557, 2018.
- 408 [48] Hanul Shin, Jung Kwon Lee, Jaehong Kim, and Jiwon Kim. Continual learning with deep generative
409 replay. *NeurIPS*, 30, 2017.
- 410 [49] Hao Tan and Mohit Bansal. Lxmert: Learning cross-modality encoder representations from transformers.
411 *arXiv preprint arXiv:1908.07490*, 2019.
- 412 [50] Xiaohan Wang, Linchao Zhu, and Yi Yang. T2VLAD: global-local sequence alignment for text-video
413 retrieval. In *CVPR*, pages 5079–5088, 2021.
- 414 [51] Xin Wang, Jiawei Wu, Junkun Chen, Lei Li, Yuan-Fang Wang, and William Yang Wang. VaTeX: A
415 large-scale, high-quality multilingual dataset for video-and-language research. In *ICCV*, pages 4580–4590,
416 2019.
- 417 [52] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. MSR-VTT: A large video description dataset for bridging
418 video and language. In *CVPR*, pages 5288–5296, 2016.
- 419 [53] Jianwei Yang, Yonatan Bisk, and Jianfeng Gao. Taco: Token-aware cascade contrastive learning for
420 video-text alignment. In *ICCV*, pages 11562–11572, 2021.
- 421 [54] Hongxu Yin, Pavlo Molchanov, Jose M Alvarez, Zhizhong Li, Arun Mallya, Derek Hoiem, Niraj K Jha,
422 and Jan Kautz. Dreaming to distill: Data-free knowledge transfer via deepinversion. In *CVPR*, pages
423 8715–8724, 2020.
- 424 [55] Jaehong Yoon, Eunho Yang, Jeongtae Lee, and Sung Ju Hwang. Lifelong learning with dynamically
425 expandable networks. In *ICLR*, 2018.
- 426 [56] Youngjae Yu, Jongseok Kim, and Gunhee Kim. A joint sequence fusion model for video question answering
427 and retrieval. In *ECCV*, pages 487–503, 2018.
- 428 [57] Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual learning through synaptic intelligence. In
429 *ICML*, pages 3987–3995. PMLR, 2017.
- 430 [58] Linchao Zhu and Yi Yang. ActBERT: Learning global-local video-text representations. In *CVPR*, pages
431 8743–8752, 2020.

432 Checklist

- 433 1. For all authors...
- 434 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
435 contributions and scope? [\[Yes\]](#)
- 436 (b) Did you describe the limitations of your work? [\[Yes\]](#) Please see Section 5.
- 437 (c) Did you discuss any potential negative societal impacts of your work? [\[No\]](#)
- 438 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
439 them? [\[Yes\]](#)
- 440 2. If you are including theoretical results...
- 441 (a) Did you state the full set of assumptions of all theoretical results? [\[N/A\]](#)
- 442 (b) Did you include complete proofs of all theoretical results? [\[N/A\]](#)
- 443 3. If you ran experiments...

- 444 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
445 mental results (either in the supplemental material or as a URL)? [No]
- 446 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
447 were chosen)? [Yes] Please see Section 4.1.
- 448 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
449 ments multiple times)? [No]
- 450 (d) Did you include the total amount of compute and the type of resources used (e.g., type
451 of GPUs, internal cluster, or cloud provider)? [Yes] Please see Section 4.1.
- 452 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 453 (a) If your work uses existing assets, did you cite the creators? [Yes] We cited the creators
454 of all five datasets in Section 4.
- 455 (b) Did you mention the license of the assets? [No]
- 456 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 457 (d) Did you discuss whether and how consent was obtained from people whose data you're
458 using/curating? [No]
- 459 (e) Did you discuss whether the data you are using/curating contains personally identifiable
460 information or offensive content? [No]
- 461 5. If you used crowdsourcing or conducted research with human subjects...
- 462 (a) Did you include the full text of instructions given to participants and screenshots, if
463 applicable? [N/A]
- 464 (b) Did you describe any potential participant risks, with links to Institutional Review
465 Board (IRB) approvals, if applicable? [N/A]
- 466 (c) Did you include the estimated hourly wage paid to participants and the total amount
467 spent on participant compensation? [N/A]