

CLIPPING: DISTILLING CLIP-BASED MODELS FOR VIDEO-LANGUAGE UNDERSTANDING

Anonymous authors

Paper under double-blind review

ABSTRACT

Pre-training a vision-language model and then fine-tuning it on downstream tasks have become a popular paradigm. However, pre-trained vision-language models with the Transformer architecture usually have a large number of parameters and take long inference time. Knowledge distillation has been an efficient technique to transfer the capability of a large model to a small one while maintaining the accuracy, which has achieved remarkable success in natural language processing. However, the collection of the pre-training data for the pre-training knowledge distillation costs huge manpower in multi-modality applications. In this paper, we propose a novel knowledge distillation method, named CLIPPING¹, where the plentiful knowledge of a large teacher model that has been fine-tuned for video-language tasks with the powerful pre-trained CLIP can be effectively transferred to a small student only at the fine-tuning stage. Especially, a new layer-wise alignment is proposed for knowledge distillation of the intermediate layers from the Transformer to the CNN in CLIPPING, which enables the student model to well absorb the knowledge of the teacher. Besides, we present an effective cross-modality knowledge distillation, which includes both the knowledge of the global video-caption distributions from the teacher model and the knowledge of the local video-caption distributions from the pre-training model (CLIP). Finally, CLIPPING with MobileViT-v2 as the vision encoder without any vision-language pre-training achieves **91.5%–95.3%** of the performance of its teacher on three video-language retrieval benchmarks, with its vision encoder being **19.5x** smaller. CLIPPING also significantly outperforms a state-of-the-art small baseline (ALL-in-one-B) on the MSR-VTT dataset, obtaining relatively **7.4%** performance gain, with **29%** fewer parameters and **86.9%** fewer flops. Moreover, CLIPPING is comparable or even superior to many large pre-training models.

1 INTRODUCTION

Recently, pre-training a vision-language model and then fine-tuning it on downstream tasks are a popular paradigm (Radford et al. (2021); Li et al. (2021); Jia et al. (2021); Li et al. (2022a); Fu et al. (2021)). Pre-trained vision-language models (PVLMs) have achieved great success in many multi-modality tasks (e.g., image-text retrieval, video-text retrieval, image captioning and VQA). However, PVLMs with the Transformer architecture (especially the vision stream) usually have a large number of parameters and a huge amount of computation, which are difficult to be deployed on edge devices such as mobile phones. Recent small models (Mehta & Rastegari (2022); Kumar et al. (2022); Guo et al. (2022); Li et al. (2022b); Yu et al. (2022); Chen et al. (2022b)) show that combining Convolutional Neural Networks (CNNs) and Transformers as a hybrid architecture gets the best of both architectures, but the overall performance of these works is still far away from satisfactory when compared to large pre-training models. Apparently, knowledge distillation (KD) (Hinton et al. (2015)) is an efficient technique to transfer the capability of a large model to a small one while maintaining the accuracy, which has achieved remarkable success in natural language processing (NLP) (Kim & Rush (2016); Jiao et al. (2020)). In NLP, the Transformer distillation methods are usually performed at both the pre-training and the fine-tuning stages. However, the collection and clearing of the pre-training data for the pre-training knowledge distillation cost huge

¹In this paper, CLIPPING means cutting something to make it smaller through distilling.

manpower in multi-modality applications (Jia et al. (2021)). Therefore, it poses a challenge here: *Can the generalization ability of a large pre-training model be transferred to a small model by only performing knowledge distillation at the fine-tuning stage?* To this end, we propose a novel knowledge distillation method, named CLIPPING, where the plentiful knowledge of a large teacher model Clip4clip (Luo et al. (2021)) that has been fine-tuned for video-language tasks with the powerful pre-trained CLIP (Radford et al. (2021)) can be effectively transferred to a small student only at the fine-tuning stage. The contributions of CLIPPING are summarized below:

- 1) We introduce an efficient approach to distill both the vision knowledge and the cross-modality knowledge from teacher to a small model without the vision-language pre-training stage. The resulting model shows strong performance on multi-modality video tasks, e.g., text-to-video retrieval and video-to-text retrieval.
- 2) We propose a layer-wise alignment scheme, called All-Student’s-layers-to-One-Teacher-layer (AS2OT), for knowledge distillation of the intermediate layers from the Transformer to the CNN in CLIPPING, where the student’s layers can be regarded as the bases of the teacher’s feature space, forcing the student model to well absorb the knowledge of the teacher. In our experience, the AS2OT layer-wise alignment significantly surpasses the previous knowledge distillation methods. We believe it will become a popular paradigm for knowledge distillation of intermediate features from the Transformer to the CNN.
- 3) We present an effective cross-modal knowledge distillation, which includes knowledge from both the global and local video-caption distributions. We use the video-caption distributions of the teacher to guide the training of the student. Besides, the student can also benefit from the powerful pre-trained CLIP that obtains certain local frame-word attention ability via learning from massive data. This pre-training knowledge can also be effectively transferred to the student by our method, even though it is just performed at the fine-tuning stage.
- 4) CLIPPING with MobileViT-v2 (Mehta & Rastegari (2022)) as the vision encoder without any vision-language pre-training achieves 91.5%–95.3% of the performance of its teacher on video-language benchmarks, with its vision encoder being 19.5x smaller. CLIPPING also significantly outperforms a state-of-the-art small baseline (ALL-in-one-B) on the MSR-VTT dataset, obtaining relatively 7.4% performance gain, with 29% fewer parameters and 86.9% fewer flops. Moreover, CLIPPING is comparable or even superior to many large pre-training models.

2 RELATIVE WORK

Vision-Language Modeling. Learning from web-collected image-text data, large-scale Vision-Language Pre-training (VLP) models such as CLIP (Radford et al. (2021)) have recently demonstrated great success across various downstream tasks. Nowadays, models such as Clip4clip (Luo et al. (2021)) and MDMMT (Dzabraev et al. (2021)) extended from the pre-trained model CLIP keep appearing. There are also some end-to-end trainable models (Xu et al. (2021); Bain et al. (2021)), which are designed to take advantage of both large-scale image and video captioning datasets. All-in-one (Wang et al. (2022a)) is the first work to consider both efficiency and performance for video-language retrieval tasks. It introduces a unified backbone that enables the representation learning of both video-text multimodal and unimodal inputs.

Knowledge Distillation. Knowledge distillation (KD) (Hinton et al. (2015)) is a simple yet effective technique to improve the performance of a learning model. Earlier works transfer knowledge embedded in the “logits” learned in a large teacher model to a small student model without sacrificing much performance. Recent works (Chen et al. (2021a;b)) use multiple layers of the teacher to supervise each layer in the student, where each layer of the student learns the knowledge from multiple layers of the teacher. (Lin et al. (2022)) proposes a target-aware Transformer and enables the student to mimic each spatial component of the teacher in each distilled layer to boost the student’s performance. For multi-modality KD, (Wang et al. (2021)) designs a fusion-encoder model as the teacher and introduces cross-modal attention knowledge to train the dual-encoder student model. The distillation objective is applied at both the pre-training and the fine-tuning stages and helps the dual-encoder model learn interactions of different modalities. TinyBert (Jiao et al. (2020)) also introduces a two-stage learning framework that performs Transformer distillation at both the pre-training and the task-specific fine-tuning stages. In this paper, we also focus on KD for transferring the knowledge of a pre-training vision-language model to a small one but only at the fine-tuning stage.

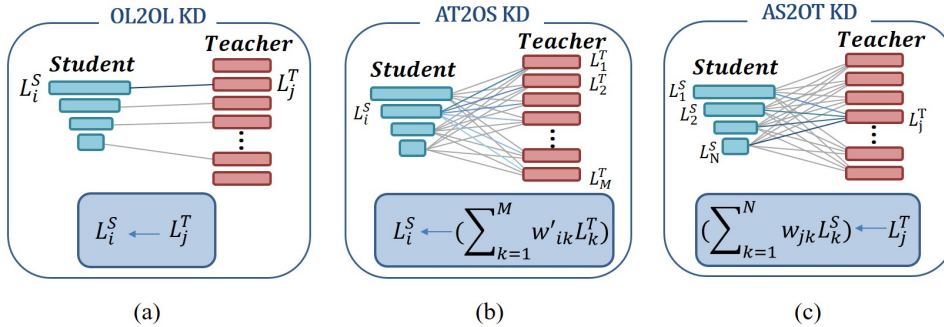


Figure 1: (a) and (b) are previous methods of intermediate features’ knowledge distillation. L_j^T , $j = 1, 2, \dots, M$, and L_i^S , $i = 1, 2, \dots, N$, are the outputs of the j^{th} and i^{th} layers of the teacher and the student, respectively. OL2OL is the most common One-Layer-to-One-Layer KD, where it selects some layers of the teacher and then distills the layers into the student one by one (e.g., from L_j^T to L_i^S). AT2OS (All-Teacher’s-layers-to-One-Student’s-layer) uses all the teacher’s layers to supervise each layer in the student, where each layer of the student (e.g., L_i^S) learns the knowledge from all the selected layers of the teacher ($\sum_{k=1}^M w'_{ik} L_k^T$, where w'_{ik} , $k = 1, 2, \dots, M$, are the knowledge selection weights with $\sum_{k=1}^M w'_{ik} = 1$). (c) Our AS2OT KD enables each of the teacher’s layers to pass its knowledge to all the student’s layers (e.g., $L_j^T \rightarrow \sum_{k=1}^N w_{jk} L_k^S$, $\sum_{k=1}^N w_{jk} = 1$).

Transformers and CNNs. Over the past ten years, CNNs have been the most popular architecture for deep learning on visual tasks. However, in the past two years, an amount of works have shown that Vision Transformers (ViTs) (Dosovitskiy et al. (2021); Liu et al. (2021)) can achieve comparable or even superior performance. Transformers use self-attention, rather than convolution, to aggregate global information across locations. For KD, most existing methods distill knowledge either from a Transformer to another Transformer (T2T) or from a CNN to another CNN (C2C) that computes the loss in the OL2OL or AT2OS style (see Fig. 1). To the best of our knowledge, KD from a Transformer to a CNN (T2C) has not been explored yet. The most related work is (Chen et al. (2022a)) that distills knowledge from a CNN to a Transformer (C2T) in the OL2OL way. However, previous works (e.g., Raghu et al. (2021)) have investigated the internal representation structures of ViTs and CNNs, and found striking differences between the two models, such as ViTs having highly similar representations throughout the model’s layers, while CNNs showing obvious distinction of representations between lower and higher layers (Raghu et al. (2021)). Considering such striking differences between ViTs and CNNs, the previous KD (OL2OL and AT2OS) may not be good distillation ways for our T2C task, which is verified by our experiments.

3 METHOD

We propose to distill both the vision knowledge and the cross-modality knowledge from Clip4clip to a small model with MobileViT-v2 as its vision encoder for multimodal video tasks, which can get the best from both sides: powerful multimodal representations of CLIP and efficient mobile vision Transformer of MobileViT-v2. The overall architecture is illustrated in Fig. 2.

3.1 PRELIMINARIES

Given a batch of videos V and captions T , Clip4clip learns a similarity function $s(v_i, t_j)$ to calculate the similarity between a video $v_i \in V$ and a caption $t_j \in T$. It obtains the frames’ representation and the caption representation in a multi-modal embedding space via CLIP. The frames’ representation is denoted as $Z_i = \{z_{i1}, z_{i2}, \dots, z_{in}\}$, where n is the number of frames in v_i and z_{ik} ($k = 1, 2, \dots, n$) is the class token from the output of the CLIP’s vision encoder corresponding to the k^{th} input frame. The caption representation is denoted as $W_j = \{w_{j1}, w_{j2}, \dots, w_{jm}\}$, where m is the number of words in t_j (including [CLS] and [SEP]) and w_{jm} is used as the representation of t_j in Clip4clip. Clip4clip uses a temporal Transformer encoder to obtain the sequential feature, denoted as $Z'_i = \{z'_{i1}, z'_{i2}, \dots, z'_{in}\} = \text{TemporalTransformer}(Z_i)$. Then, the mean pooling is used to aggregate

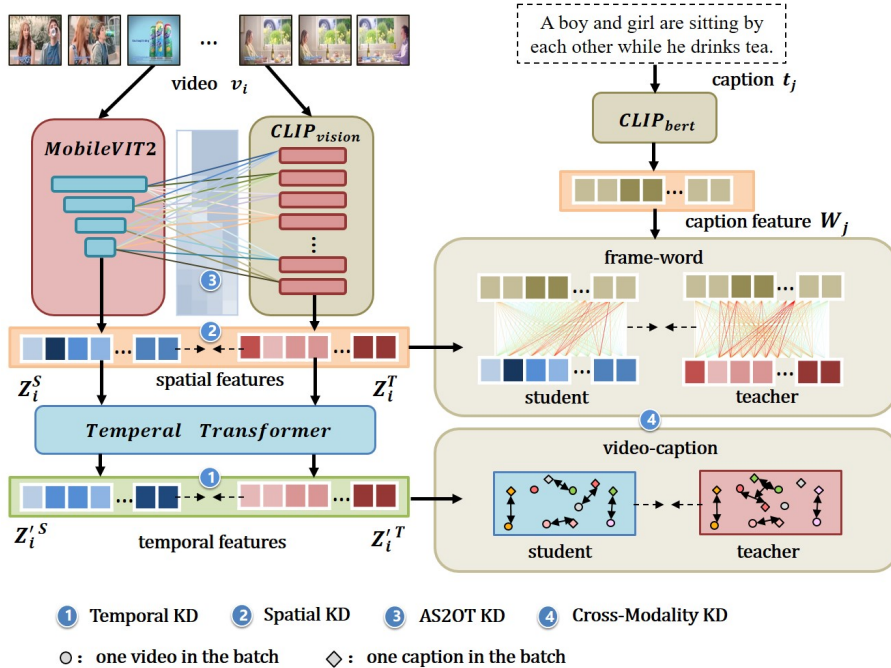


Figure 2: Overview of CLIPPING. There are mainly four knowledge distillation (KD) parts: (1) Temporal KD. (2) Spatial KD. (3) AS2OT KD. (4) Cross-modality KD. The CLIP’s vision encoder is a Transformer.

the features of all frames to obtain the video representation, $z_i = \frac{1}{n} \sum_{k=1}^n z'_{ik}$. Finally, the similarity functions $s(v_i, t_j)$ and $s(t_j, v_i)$ are defined as:

$$s(v_i, t_j) = w_{jm}^\top z_i, \quad s(t_j, v_i) = z_i^\top w_{jm}. \quad (1)$$

Nowadays, models such as Clip4clip based on the pre-trained model CLIP have been applied to many tasks. They pursue better performance with CLIP but have a heavy vision encoder, which makes them difficult to be deployed on edge devices such as mobile phones. Since MobileViT-v2 (Mehta & Rastegari (2022)) is a light-weight and mobile-friendly hybrid network, we employ it as the vision encoder of the student model and maintain the original text encoder and temporal Transformer of Clip4clip.

3.2 CLIPPING

To distill multimodal knowledge from Clip4clip to the MobileViT-v2-based model, we use four kinds of knowledge transfer: 1) temporal KD, 2) spatial KD, 3) AS2OT KD, and 4) cross-modality KD.

3.2.1 TEMPORAL AND SPATIAL KNOWLEDGE DISTILLATION

The temporal KD is motivated by the previous finding that aligning multi-frame dependency from the teacher to the student can enhance the performance of the student (Liu et al. (2020)), which encodes the multi-frame dependency into a latent embedding by using a recurrent unit ConvLSTM. In our architecture, multi-frame dependency is modeled naturally through the temporal Transformer (Temporal Transformer in Fig. 2), so we define the temporal KD loss as:

$$L_{TKD} = \frac{1}{B} \sum_{i=1}^B D_{KL}(Z_i^S, Z_i^T), \quad (2)$$

where B is the batch size, the superscripts S and T denote the student and the teacher, respectively, and $D_{KL}()$ is the KL divergence loss function. In addition to imitating the behavior of CLIP for

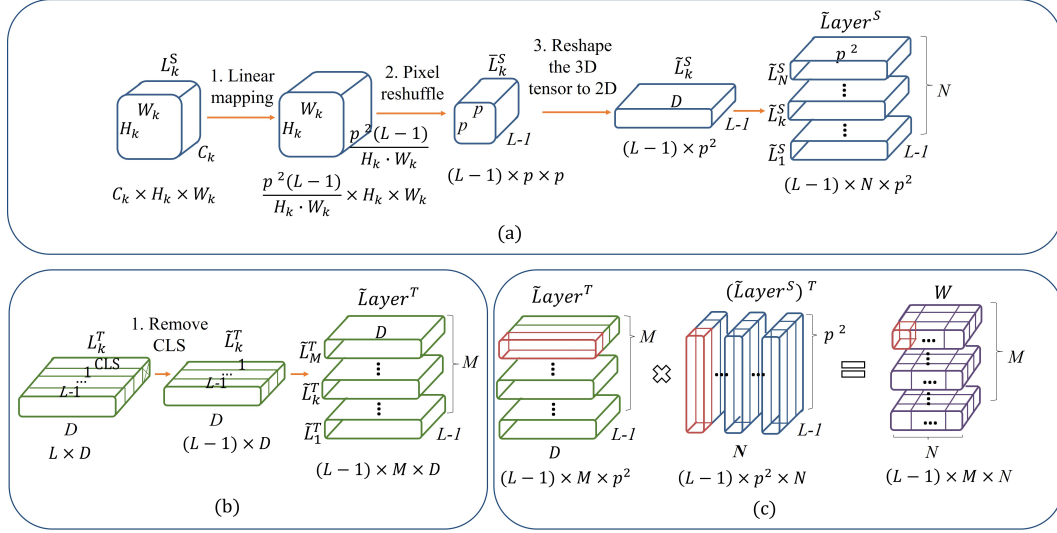


Figure 3: Reshaping the features of (a) the student L_k^S and (b) the teacher L_k^T to the same shape. (c) Similarity Computation.

each frame, we use spatial KD to align the spatial embeddings of the student and the teacher by:

$$L_{SKD} = \frac{1}{B} \sum_{i=1}^B D_{KL}(Z_i^S, Z_i^T). \quad (3)$$

3.2.2 AS2OT KNOWLEDGE DISTILLATION

Recently, it is discovered that distilling intermediate features is more effective (Chen et al. (2021a)). So in addition to the spatial and temporal KD, we also force the vision encoder of the student to mimic the intermediate layers of the teacher. In the traditional OL2OL KD, some layers of the teacher are selected and distilled into the student one by one (Jiao et al. (2020)). However, it is difficult to find the optimal correspondence between the student and the teacher. The recent AT2OS works (Chen et al. (2021a;b)) still cannot help the student learn enough knowledge from the teacher when there is a large architecture gap between them, which is verified in our experiments. For our task, the student (MobileViT-v2) is a hybrid architecture, which has convolution layers in the shallow stages and separable self-attention layers in the later blocks. And a separable self-attention layer differs significantly from the traditional self-attention in two ways: 1) it does not learn an explicit attention map; 2) its input and output are still 3D, which are more similar to CNN features. Except for the separable self-attention in the later blocks, the shallow stages with convolution layers undoubtedly have a huge difference from the Transformer structure in CLIP, which has been investigated in previous works (Yuan et al. (2021); Raghu et al. (2021)). Therefore, we design AS2OT layer-wise alignment for KD of the intermediate layers from the Transformer to the CNN.

AS2OT Property. As shown in Fig. 1(c), the student’s layers in our AS2OT KD can be explained as the bases of the feature space, and each layer of the teacher is a linear combination of the bases. After training, if the student’s layers do ideally form the bases, the knowledge of the teacher’s layers is learned completely by the student. Our experiment in Section 4 shows that the teacher’s features in different layers can be well recovered from the features of the student’s layers. Next, we describe how to implement AS2OT KD.

Feature Reshaping and Similarity Computation. Let the layers of the teacher and the student be $Layer^T = [L_1^T, L_2^T, \dots, L_M^T]$ and $Layer^S = [L_1^S, L_2^S, \dots, L_N^S]$ (usually $M > N$), respectively. $L_k^T \in R^{D \times L}$ is a 2D tensor, where D is the dimensionality of the token feature and L is the number of tokens, while $L_k^S \in R^{C_k \times H_k \times W_k}$ is a 3D tensor, where C_k , H_k and W_k are the channel number, height and width of the k^{th} CNN layer’s feature. We calculate a similarity tensor W between L_k^S and L_k^T , which need to be reshaped first as shown in Figs. 3(a) and (b), respectively. For L_k^S , we first apply a linear operator and then do pixel reshuffle on the result, obtaining \bar{L}_k^S , $k = 1, 2, \dots, N$. Next, \bar{L}_k^S is reshaped to \bar{L}_k^S . As for L_k^T , we remove the class token and obtain a $(L-1) \times D$ tensor, denoted as \bar{L}_k^T , $k = 1, 2, \dots, M$. Finally, all \bar{L}_k^S , $k = 1, 2, \dots, N$, and \bar{L}_k^T , $k = 1, 2, \dots, M$, are represented

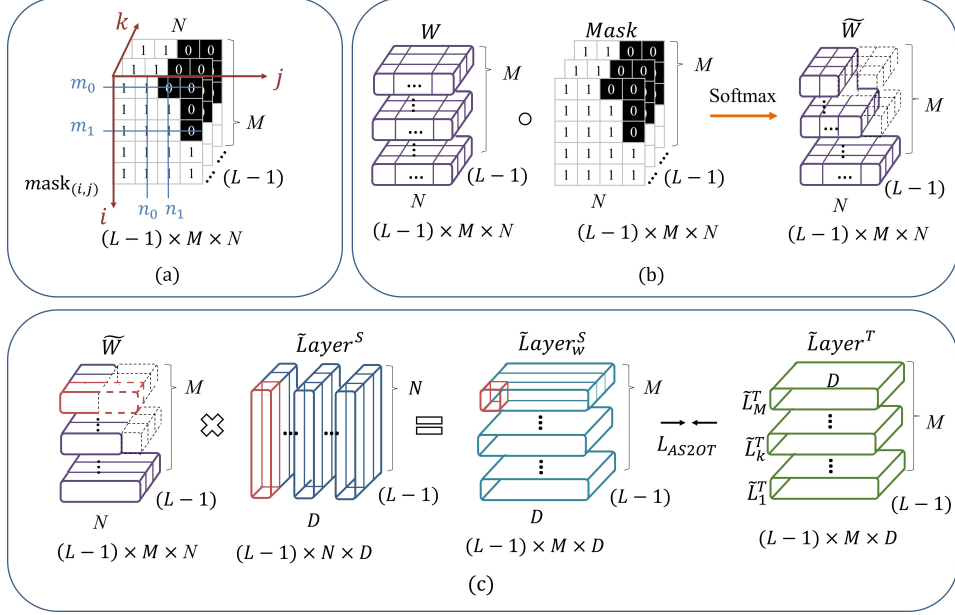


Figure 4: (a) Masks obtained with Eq. 4, where all the $L - 1$ masks along the k dimension are the same. (b) Masking the similarity tensor W . (c) AS2OT alignment. L_{AS2OT} is the AS2OT loss.

as $\tilde{Layer}^S = [\tilde{L}_1^S, \tilde{L}_2^S, \dots, \tilde{L}_N^S]$ and $\tilde{Layer}^T = [\tilde{L}_1^T, \tilde{L}_2^T, \dots, \tilde{L}_M^T]$, respectively. After this reshaping, as shown in Fig. 3(c), we are able to measure the similarities between the intermediate layers of the student and the teacher through $W = \tilde{Layer}^T \times (\tilde{Layer}^S)^\top$, $W \in R^{(L-1) \times M \times N}$.

Masking and AS2OT Alignment. To speed up the training procedure, we design AS2OT KD with sequential masks, which follow this formula:

$$mask_{(i,j)} = \begin{cases} 0, & \text{if } (j > n_0 \text{ and } i \leq m_0) \text{ or } (j > n_1 \text{ and } m_0 < i \leq m_1) \\ 1, & \text{else} \end{cases}, \quad (4)$$

where $i = 1, 2, \dots, M$, $j = 1, 2, \dots, N$, $1 \leq m_0 < m_1 < M$, $1 \leq n_0 < n_1 < N$. And the final attention mask $Mask \in R^{(L-1) \times M \times N}$ is composed of $L - 1$ same masks, Fig. 4(a) shows one example with $m_0 = 1, n_0 = 2, m_1 = 3, n_1 = 3$. $mask_{(i,j)} = 0$ means to mask all the elements along the k dimension (Fig. 4(a)) at (i, j) in W , while $mask_{(i,j)} = 1$ means to maintain their similarities, as shown in Fig. 4(b). This masking is based on our experimental finding: The student’s lower-level layers should learn from the teacher’s lower-layers, while the student’s higher layers should learn from all the teacher’s layers. In Fig. 4(b), we also normalize the masks by performing softmax along the j dimension on each row, obtaining $\tilde{W} = \sigma(W \circ Mask)$, where σ is the softmax function and \circ denotes the element-wise product. Guided by \tilde{W} , as shown in Fig. 4(c), the weighted student layers are calculated as: $\tilde{Layer}_w^S = \tilde{W} \times \tilde{Layer}^S$. Each component of \tilde{Layer}_w^S is the linear combination of N student layers’ features. Finally, the AS2OT KD loss is defined as:

$$L_{AS2OT} = D_{KL}(\tilde{Layer}_w^S, \tilde{Layer}^T). \quad (5)$$

3.2.3 CROSS-MODALITY KNOWLEDGE DISTILLATION

To further adapt the student’s features to the teacher’s multi-modal feature space, we employ the teacher’s knowledge of cross-modality distributions to guide the training of the student. Specifically, we characterize the cross-modal distributions from two perspectives, global video-caption distributions of the teacher and local video-caption distributions of CLIP.

Global Video-Caption Distribution Alignment. We consider both the video-to-caption distribution A_{GVC} and the caption-to-video distribution A_{GCV} , the elements $(s(v_i, t_j))$ and $(s(t_j, v_i))$ of which are defined via the similarities obtained by Eq. 1:

$$A_{GVC} = \sigma \begin{pmatrix} s(v_1, t_1) & \dots & s(v_1, t_B) \\ \dots & \dots & \dots \\ s(v_B, t_1) & \dots & s(v_B, t_j) \end{pmatrix}, \quad A_{GCV} = \sigma \begin{pmatrix} s(t_1, v_1) & \dots & s(t_1, v_B) \\ \dots & \dots & \dots \\ s(t_B, v_1) & \dots & s(t_B, v_B) \end{pmatrix}. \quad (6)$$

Finally, the global video-caption distribution alignment loss is defined as:

$$L_G = D_{KL}(A_{GVC}^S, A_{GVC}^T) + D_{KL}(A_{GCV}^S, A_{GCV}^T). \quad (7)$$

Local Video-Caption Distribution Alignment. CLIP uses text prompts (such as “A picture of a ()”) for zero-shot image classification. CLIP fills them with different words (e.g., “cat” and “dog”) and results in different captions (e.g., “A picture of a cat” and “A picture of a dog”). It can match the captions to the corresponding images, showing some image-word alignment ability. We transfer this pre-training knowledge to the student through a local frame-word alignment as follows. Recall that $Z'_i = \{z'_{i1}, z'_{i2}, \dots, z'_{in}\}$ and $W_j = \{w_{j1}, w_{j2}, \dots, w_{jm}\}$ are the features of the video v_i and the caption t_j , respectively, with n being the number of frames in v_i and m the number of words in t_j (Section 3.1). The similarity between the k^{th} frame and the r^{th} word is defined as $s_{fw}(w_{jr}, z'_{ik}) = (w_{jr})^\top z'_{ik}$. Then we respectively define the local video-to-caption and caption-to-video similarities as:

$$s'(v_i, t_j) = \frac{1}{n} \sum_{k=1}^n \max_{1 \leq r \leq m} \{s_{fw}(w_{jr}, z'_{ik})\}, \quad s'(t_j, v_i) = \frac{1}{m} \sum_{r=1}^m \max_{1 \leq k \leq n} \{s_{fw}(z'_{ik}, w_{jr})\}. \quad (8)$$

Finally, we have the local video-to-caption distribution A_{LVC} and the local caption-to-video distribution A_{LCV} :

$$A_{LVC} = \sigma \begin{pmatrix} s'(v_1, t_1) & \dots & s'(v_1, t_B) \\ \dots & \dots & \dots \\ s'(v_B, t_1) & \dots & s'(v_B, t_B) \end{pmatrix}, \quad A_{LCV} = \sigma \begin{pmatrix} s'(t_1, v_1) & \dots & s'(t_1, v_B) \\ \dots & \dots & \dots \\ s'(t_B, v_1) & \dots & s'(t_B, v_B) \end{pmatrix}, \quad (9)$$

and the local caption-video distribution alignment loss is defined as:

$$L_L = D_{KL}(A_{LCV}^S, A_{LCV}^T) + D_{KL}(A_{LVC}^S, A_{LVC}^T). \quad (10)$$

Combining L_G and L_V , the cross-modality KD loss is:

$$L_{CM} = 0.5 \cdot L_G + 0.5 \cdot L_L. \quad (11)$$

The total loss for training our model is:

$$L = L_{task} + \alpha \cdot L_{TKD} + \beta \cdot L_{SKD} + \gamma \cdot L_{CM} + \delta \cdot L_{AS2OT}, \quad (12)$$

where L_{task} is the task-specific loss, and α , β , γ and δ are loss balance weights.

4 EXPERIMENTS

We conduct comprehensive experiments on three benchmarks for video-text retrieval (video-to-text ($v2t$) and text-to-video ($t2v$)): MSR-VTT (Xu et al. (2016)), MSVD (Chen & Dolan (2011)) and LSMDC (Rohrbach et al. (2015)). The metrics Recall at rank 1 ($R@1$), rank 5 ($R@5$) and rank 10 ($R@10$) are used for evaluation.

4.1 COMPARISON WITH STATE-OF-THE-ARTS

In Table 1, we compare the proposed model CLIPPING with eight state-of-the-art methods on MSR-VTT: TACo (Yang et al. (2021)), VideoClip (Xu et al. (2021)), Frozen (Bain et al. (2021)), VIOLET (Fu et al. (2021)), OA-Trans (Wang et al. (2022b)), BridgeFormer (Ge et al. (2022)), ALL-in-one (Wang et al. (2022a)) and MDMMT (Dzabraev et al. (2021)). It can be seen that CLIPPING significantly surpasses those large-scale video-text/image-text pre-training models for video-text retrieval. For example, our model exceeds VideoClip and Frozen by absolute 9.8% $t2vR@1$ and 8.2% $t2vR@1$, respectively. In addition, CLIPPING also outperforms the small model ALL-in-one-B even though CLIPPING is smaller. Note that the model MDMMT uses CLIP as its backbone, while our CLIPPING uses CLIP as the teacher. The results on the MSVD and LSMDC datasets are given in the supplementary materials.

In Fig. 5, we show the performances and Flops of these models. Among previous models, MDMMT has the best performance and ALL-in-one-S is the fastest. CLIPPING not only obtains relative 4.6% performance gain over MDMMT but also is faster than All-in-one-S.

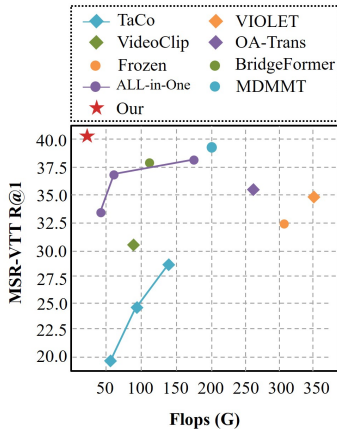


Figure 5: Flops and Performances.

Table 1: Comparison with state-of-the-art models on MSR-VTT (1k split) for text-to-video retrieval. “PT Datasets”: datasets used for pre-training the vision encoder. “HT100M”: HowTo100M dataset (Miech et al. (2019)). “C400M”: CLIP-400M dataset (Radford et al. (2021)). “IN21K”: ImageNet21K dataset (Deng et al. (2009)). “W2M”: WebVid-2M dataset (Bain et al. (2021)). “C3M”: CC3M dataset (Sharma et al. (2018)). “COCO”: COCO dataset (Chen et al. (2015)). “AudioSet”: AudioSet dataset (Kong et al. (2018)).

Model	PT Datasets	Params	R@1	R@5	R@10
TACo	HT100M	212M	28.4	57.8	71.2
VideoClip	HT100M	130M	30.9	55.4	66.8
Frozen	C3M,W2M,COCO	232M	32.5	59.5	70.5
ALL-in-one-S	W2M,HT100M	33M	33.5	-	-
VIOLET	C3M,W2M	198M	34.5	63.0	73.4
OA-Trans	C3M,W2M	232M	35.8	63.4	76.5
BridgeFormer	C3M,W2M	160M	37.6	64.8	75.1
ALL-in-one-B	W2M,HT100M	110M	37.9	68.1	77.1
MDMMT	C400M,AudioSet	226M	38.9	68.3	78.8
CLIPPING (our)	IN21K	78.1M	40.7	68.5	78.9

Table 2: Ablation study of different KD components of CLIPPING on the 1k validation set of MSR-VTT. The first row is the results of the teacher model (Clip4clip). T , S , $AS2OT$, CM_G and CM_L denote temporal KD, spatial KD, AS2OT KD, global cross-modality KD and local cross-modality KD, respectively.

Vision Encoder	PT Dataset	Params	Flops	KD Types	$t2vR@1$	$v2tR@1$
CLIP _{vision}	C400M	87.8M	8.6G	-	44.5	42.2
MobileViTv2	IN21K	4.5M	1.4G	-	25.7	24.5
MobileViTv2	IN21K	4.5M	1.4G	T	28.8	27.3
MobileViTv2	IN21K	4.5M	1.4G	T,S	33.0	32.8
MobileViTv2	IN21K	4.5M	1.4G	$T,S,AS2OT$	37.6	36.2
MobileViTv2	IN21K	4.5M	1.4G	$T,S,AS2OT,CM_G$	39.6	39.1
MobileViTv2	IN21K	4.5M	1.4G	$T,S,AS2OT,CM_G,CM_L$	40.7	40.2

4.2 ABLATION STUDY

Key Components. We provide detailed ablation study to validate each key component of our proposed method, on MSR-VTT. From Table 2, we can see that without any KD, it gets extreme low accuracies. When simply adding the temporal KD and spatial KD, $t2vR@1$ and $v2tR@1$ increase significantly. When we further add AS2OT KD, and then global and local cross-modality KD, CLIPPING’s performance rises gradually. This study shows that all these key components of CLIPPING are effective. Compared with the teacher, the full CLIPPING achieves about 91.5% and 95.3% of the performance of its teacher on $t2vR@1$ and $v2tR@1$, respectively, with the vision encoder of 19.5x smaller.

KD Types. In Table 3, we compare the proposed AS2OT KD with the previous OL2OL KD and AT2OS KD (Fig. 1). For AS2OT and AT2OS, the same layers of the student and the teacher are used for KD, the details of which are give in the supplementary materials. From Table 3, we can see that AT2OS performs better than OL2OL, and our AS2OT outperforms AT2OS. In Section 3.2.2, we use masking to speed up the training. In this study, we also compare “with masking” and “without masking”. The last two rows of Table 3 verify that the masking is beneficial.

Table 3: Ablation study of different KD types on MSR-VTT (1k split). All the models are trained for 36 epochs with the same setting (see the supplementary materials for details).

KD Types	$t2vR@1$	$v2tR@1$
T, S	33.0	32.8
$T, S, OL2OL$	34.6	33.4
$T, S, AT2OS$	35.1	34.4
$T, S, AS2OT_{w/o\ masking}$	37.1	35.7
$T, S, AS2OT_{w/\ masking}$	37.6	36.2

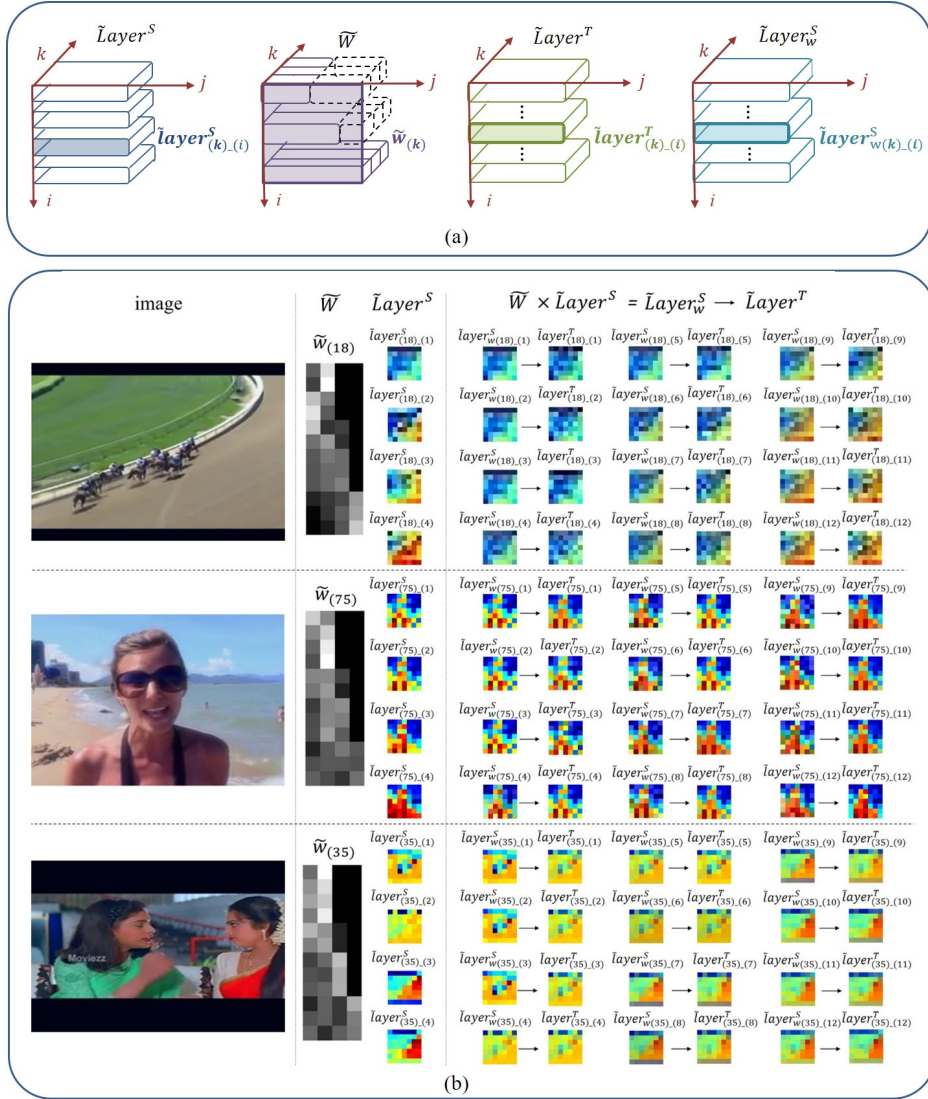


Figure 6: (a) Feature tensors of the student and the teacher. (b) Examples to demonstrate that the teacher’s features are the linear combinations of the student features.

AS2OT Property. The student’s layers in our AS2OT KD can be explained as the bases of the feature space, and each layer of the teacher is a linear combination of the bases (Section 3.2). In Fig. 6, we show that this property holds. In AS2OT KD, the student has 4 layers and the teacher has 12 layers. We randomly pick 3 image examples, and for each example, we select the features of one random token (in the k dimension in Fig. 6(a)). After training, the linear combinations $\tilde{W} \times \tilde{L}ayer^S = \tilde{L}ayer^S_w$ show feature patterns very similar to the teacher’s features ($\tilde{L}ayer^T$). This property verifies that the teacher’s knowledge is fully absorbed by the student.

5 CONCLUSION

In this paper, we propose a novel knowledge distillation method that is specially designed for small vision-language models. It includes temporal KD, spatial KD, AS2OT KD and cross-modality KD. Especially, the AS2OT KD has the property of the student’s layers being the bases of the feature space. After training, the teacher’s features are the linear combinations of the bases, indicating that the student has fully absorbed the knowledge of the teacher. Our method CLIPPING can achieve 91.5%–95.3% of its performance of its teacher on three retrieval benchmarks without any vision-language pre-training KD. CLIPPING significantly outperforms a state-of-the-art small baseline ALL-in-one-B. Moreover, it is comparable or even superior to many large pre-training models. In the future, we will apply CLIPPING to other vision-language models for compression.

REFERENCES

- Max Bain, Arsha Nagrani, Gül Varol, and Andrew Zisserman. Frozen in Time: A Joint Video and Image Encoder for End-to-End Retrieval. In *ICCV*, 2021.
- David Chen and William Dolan. Collecting highly parallel data for paraphrase evaluation. In *ACL*, 2011.
- Defang Chen, JianPing Mei, Yuan Zhang, Can Wang, Zhe Wang, Yan Feng, and Chun Chen. Distillation with semantic calibration. In *AAAI*, 2021a.
- Pengguang Chen, Shu Liu, Hengshuang Zhao, and Jiaya Jia. Distilling Knowledge via Knowledge Review. In *CVPR*, 2021b.
- Xianing Chen, Qiong Cao, Yujie Zhong, Jing Zhang, Shenghua Gao, and Dacheng Tao. DearKD: Data-efficient early knowledge distillation for vision transformers. In *ICLR*, 2022a.
- Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollar, and C. Lawrence Zitnick. Microsoft coco captions: Data collection and evaluation server. In *arXiv:1504.00325*, 2015.
- Yinpeng Chen, Xiyang Dai, Dongdong Chen, Mengchen Liu, Xiaoyi Dong, Lu Yuan, and Zicheng Liu. Mobile-Former: Bridging MobileNet and Transformer. In *CVPR*, 2022b.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *CVPR*, 2009.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021.
- Maksim Dzabaraev, Maksim Kalashnikov, Stepan Komkov, and Aleksandr Petiushko. MDMMT: Multidomain Multimodal Transformer for Video Retrieval. In *CVPR*, 2021.
- Tsu-Jui Fu, Linjie Li, Zhe Gan, Kevin Lin, William Yang Wang, Lijuan Wang, and Zicheng Liu. VIOLET: End-to-End Video-Language Transformers with Masked Visual-token Modeling. In *arXiv:2111.1268*, 2021.
- Yuying Ge, Yixiao Ge, Xihui Liu, Dian Li, Ying Shan, Xiaohu Qie, and Ping Luo. Bridgeformer: Bridging video-text retrieval with multiple choice questions. In *CVPR*, 2022.
- Jianyuan Guo, Kai Han, Han Wu, Yehui Tang, Xinghao Chen, Yunhe Wang, and Chang Xu. CMT: Convolutional Neural Networks Meet Vision Transformers. In *CVPR*, 2022.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the Knowledge in a Neural Network. In *arXiv preprint arXiv:1503.02531*, 2015.
- Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *ICML*, 2021.
- Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu. TinyBERT: Distilling BERT for Natural Language Understanding. In *EMNLP*, 2020.
- Yoon Kim and Alexander M. Rush. Sequence-level knowledge distillation. In *EMNLP*, 2016.
- Qiuqiang Kong, Yong Xu, Wenwu Wang, and Mark D. Plumbley. Audio set classification with attention model: A probabilistic perspective. In *ICASSP*, 2018.
- Pavan Kumar, Anasosalu Vasu, James Gabriel, Jeff Zhu, Oncel Tuzel, and Anurag Ranjan. An Improved One millisecond Mobile Backbone. In *arXiv preprint arXiv:2206.04040*, 2022.
- Junnan Li, Ramprasaath Selvaraju, Akhilesh Gotmare, Shafiq Joty, Caiming Xiong, and Steven Chu Hong Hoi. Align before fuse: Vision and language representation learning with momentum distillation. In *NIPS*, 2021.

- Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. BLIP: Bootstrapping Language-Image Pre-training for Unified Vision-Language Understanding and Generation. In *arXiv preprint arXiv:2201.12086*, 2022a.
- Yanyu Li, Geng Yuan, Yang Wen, Eric Hu, Georgios Evangelidis, Sergey Tulyakov, Yanzhi Wang, and Jian Ren. EfficientFormer: Vision Transformers at MobileNet Speed. In *arXiv:2206.01191*, 2022b.
- Sihao Lin, Hongwei Xie, Bing Wang, Kaicheng Yu, Xiaojun Chang, Xiaodan Liang, and Gang Wang. Knowledge Distillation via the Target-aware Transformer. In *CVPR*, 2022.
- Yifan Liu, Chunhua Shen, Changqian Yu, and Jingdong Wang. Efficient Semantic Video Segmentation with Per-Frame Inference. In *ECCV*, 2020.
- Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin Transformer: Hierarchical Vision Transformer using Shifted Windows. In *ICCV*, 2021.
- Huaishao Luo, Lei Ji, Ming Zhong, Yang Chen, Wen Lei, Nan Duan, and Tianrui Li. Clip4clip: An empirical study of clip for end to end video clip retrieval. *arXiv preprint arXiv:2104.08860*, 2021.
- Sachin Mehta and Mohammad Rastegari. Separable Self-attention for Mobile Vision Transformers. In *ICLR*, 2022.
- Antoine Miech, Dimitri Zhukov, Jean-Baptiste Alayrac, Makarand Tapaswi, Ivan Laptev, and Josef Sivic. HowTo100M: Learning a Text-Video Embedding by Watching Hundred Million Narrated Video Clips. In *ICCV*, 2019.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, and Jack Clark. Learning transferable visual models from natural language supervision. *ICML*, 2021.
- Maithra Raghu, Thomas Unterthiner, Simon Kornblith, Chiyuan Zhang, and Alexey Dosovitskiy. Do vision transformers see like convolutional neural networks? In *NIPS*, 2021.
- Anna Rohrbach, Marcus Rohrbach, Niket Tandon, and Bernt Schiele. A dataset for Movie Description. In *CVPR*, 2015.
- Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypertexted, image alt-text dataset for automatic image captioning. In *ACL*, 2018.
- Alex Jinpeng Wang, Yixiao Ge, Rui Yan, Yuying Ge, Xudong Lin, Guanyu Cai, Jianping Wu, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. All in one: Exploring unified video-language pre-training. In *arXiv preprint arXiv:2203.07303*, 2022a.
- Jinpeng Wang, Yixiao Ge, Guanyu Cai, Rui Yan, Xudong Lin, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. Object-aware video-language pre-training for retrieval. In *CVPR*, 2022b.
- Zekun Wang, Wenhui Wang, Haichao Zhu, Ming Liu, Bing Qin, and Furu Wei. Distilled dual-encoder model for vision-language understanding. In *arXiv preprint arXiv:2112.08723*, 2021.
- Hu Xu, Gargi Ghosh, Po-Yao Huang, Dmytro Okhonko, Armen Aghajanyan, Florian Metze, Luke Zettlemoyer, and Christoph Feichtenhofer. Videoclip: Contrastive pre-training for zero-shot video-text understanding. *arXiv preprint arXiv:2104.08860*, 2021.
- Jun Xu, Tao Mei, Ting Yao, and Yong Rui. MSR-VTT: A Large Video Description Dataset for Bridging Video and Language. 2016.
- Jianwei Yang, Yonatan Bisk, and Jianfeng Gao. Taco: Token-aware cascade contrastive learning for video-text alignment. In *ICCV*, 2021.
- Weihao Yu, Mi Luo, Pan Zhou, Chenyang Si, Yichen Zhou, Xinchao Wang, Jiashi Feng, and Shuicheng Yan. Metaformer is actually what you need for vision. In *CVPR*, 2022.
- Li Yuan, Yunpeng Chen, Tao Wang, Weihao Yu, Yujun Shi, Zi-Hang Jiang, Francis E.H. Tay, Jiashi Feng, and Shuicheng Yan. Tokens-to-Token ViT: Training Vision Transformers From Scratch on ImageNet. In *ICCV*, 2021.