# No Head Left Behind - Multi-Head Alignment Distillation for Transformers

Tianyang Zhao<sup>1,2\*</sup>, Kunwar Yashraj Singh<sup>1†</sup>, Srikar Appalaraju<sup>1‡</sup>, Peng Tang<sup>1</sup>, Vijay Mahadevan<sup>1</sup>, R. Manmatha<sup>1</sup>, Ying Nian Wu<sup>1,2</sup>

<sup>1</sup>AWS AI Labs,

<sup>2</sup>University of California, Los Angeles
tyzhao@ucla.edu, {sinkunwa, srikara, tangpen, vmahad, manmatha, wunyin}@amazon.com

#### Abstract

Knowledge distillation aims at reducing model size without compromising much performance. Recent work has applied it to large vision-language (VL) Transformers, and has shown that attention maps in the multi-head attention modules of vision-language Transformers contain extensive intra-modal and cross-modal co-reference relations to be distilled. The standard approach is to apply a one-to-one attention map distillation loss, i.e. the Teacher's first attention head instructs the Student's first head, the second teaches the second, and so forth, but this only works when the numbers of attention heads in the Teacher and Student are the same. To remove this constraint, we propose a new Attention Map Alignment Distillation (AMAD) method for Transformers with multi-head attention, which works for a Teacher and a Student with different numbers of attention heads. Specifically, we soft-align different heads in Teacher and Student attention maps using a cosine similarity weighting. The Teacher head contributes more to the Student heads for which it has a higher similarity weight. Each Teacher head contributes to all the Student heads by minimizing the divergence between the attention activation distributions for the soft-aligned heads. No head is left behind. This distillation approach operates like crossattention. We experiment on distilling VL-T5 and BLIP, and apply AMAD loss on their T5, BERT, and ViT sub-modules. We show, under vision-language setting, that AMAD outperforms conventional distillation methods on VQA-2.0, COCO captioning, and Multi30K translation datasets. We further show that even without VL pre-training, the distilled VL-T5 models outperform corresponding VL pre-trained VL-T5 models that are further fine-tuned by ground-truth signals, and that fine-tuning distillation can also compensate to some degree for the absence of VL pre-training for BLIP models.

#### Introduction

Recently, large pre-trained Transformers-based (Vaswani et al. 2017) models, such as BERT (Devlin et al. 2019), T5 (Raffel et al. 2020), and GPT (Radford et al. 2018), have shown great capabilities for language modeling. Researchers have further extended these language Transformers to multi-modal Transformers for visual-linguistic tasks,

Copyright © 2024, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

e.g. VL-BERT (Su et al. 2019), VL-T5 (Cho et al. 2021), Oscar (Li et al. 2020b), BLIP (Li et al. 2022b, 2023), OFA (Wang et al. 2022a), Flamingo (Alayrac et al. 2022), Florence (Yuan et al. 2021), PALI (Chen et al. 2022c). These large vision-language (VL) models exhibit promising performances on a variety of visual-linguistic tasks, including visual question answering (VQA), image captioning, visual grounding and image-text matching. Increasing model size (BERT-L (340M), OFA (930M), T5 (11B), GPT-3 (175B) (Brown et al. 2020)) leads to better performance, but also increases memory consumption during deployment and leads to large increases in inference latency.

To alleviate this problem, researchers (Jiao et al. 2020; Sun et al. 2020; Wang et al. 2020b; Fang et al. 2021; Sanh et al. 2019) have applied knowledge distillation (KD) (Hinton, Vinyals, and Dean 2015) approaches to large Transformers, aiming at compressing these large models into smaller ones without compromising much performance. In general, KD involves a large trained and frozen Teacher network, and a small Student network to be trained. The goal is to distill the knowledge from the larger Teacher into the smaller Student to bridge the performance gap between the two caused by difference of model sizes. In the distillation process, the Student learns to mimic the soft response and the latent representation of the Teacher. Specifically, this may involve minimizing the divergence between the Student's and the Teacher's output classification logits, and the divergence between their intermediate representations.

In case of distilling Transformers, their attention maps are often important and contain intermediate representations to be transferred. (Cao et al. 2020) show that certain attention matrices of the pre-trained vision-language Transformers contain extensive intra- and cross-modal co-reference relations. (Fang et al. 2021) further show that minimizing the divergence between these attention maps of Teacher and Student can boost distillation performance.

However, conventional attention map distillation methods for multi-head attention modules, either for language (Jiao et al. 2020; Sun et al. 2020; Wang et al. 2020b; Sanh et al. 2019) or VL (Fang et al. 2021) Transformers, directly minimize the divergence between the attention maps of Teacher and Student for each of their heads in a one-to-one fashion, i.e. the Teacher's first attention head instructs the Student's first head, the second teaches the second, and so forth, as in

<sup>\*</sup>Work conducted during an internship at Amazon.

<sup>†</sup>Corresponding Author

<sup>‡</sup>Corresponding Author

the right side of Figure 1. Hence, these methods can only be applied where the Teacher and Student have an equal number of heads, and do not generalize to the more common case where the large Transformer and the small Transformer have different numbers of attention heads; Otherwise if the Teacher has more heads, its extra heads have to be discarded in the distilling process.

This motivates the design of our approach, Attention Map Alignment Distillation (AMAD), to remove the same-number-of-heads restriction. In brief, AMAD soft aligns different heads in the Teacher and Student attention maps using cosine similarity. Each Teacher head teaches all the Student heads with the contribution being more for the Student heads with which it has higher weights (higher cosine similarities). The Teacher teaches the Student heads by minimizing the divergence between the attention activation distributions for the soft-aligned heads. We may view AMAD as operating like a cross-attention itself.

One intuition behind is that, unlike embeddings for which each vision/language token shares the same order for both Teacher and Student, attention heads do not have a semantic order. (Cao et al. 2020) found that different subsets of attention heads in VL Transformers may encode different co-reference knowledge, e.g. a subset of heads may evolve to pivot on cross-modal interaction between image and text regimes. Therefore, even in the case when Teacher and Student have the same number of attention heads, we still cannot assume that the Teacher's and the Student's heads are aligned semantically without reordering. Conventional attention map distillation methods force Student heads to have exactly the same order as Teacher's, while AMAD allows similarity alignment based distillation free of head order: each Teacher teaches all the Students, the contribution being proportional to the similarity weight between them.

We conduct experiments on distilling VL-T5 (Cho et al. 2021) *base* to *small*, and on distilling BLIP (Li et al. 2022b) *large* to *base*. AMAD loss is applied on all their Transformer sub-modules, including T5-Encoder + Decoder (Raffel et al. 2020) for VL-T5, and Vision Transformer (ViT) (Dosovitskiy et al. 2021) + BERT (Devlin et al. 2019) for BLIP. We evaluate on VQA-2.0 (Goyal et al. 2019), COCO Captioning (Chen et al. 2015), and Multi30K (Elliott et al. 2016) datasets. We show that AMAD boosts performance.

Our contributions include:

- We propose Attention Map Alignment Distillation (AMAD) to distill attention maps from a Teacher Transformer to a Student Transformer with different numbers of attention heads. AMAD uses a soft-alignment approach so that each Teacher head teaches all the Student heads but in proportion to how similar the Student is to the Teacher. We show, under vision-language setting, that AMAD narrows the performance gap between the large Teacher and the small Student. With AMAD, researchers are set free from the same-number-of-heads restriction and have more choices over Transformers for distillation.
- We show that even without VL pre-training, distilled VL-T5 models outperform VL pre-trained VL-T5 models of the same size further fine-tuned with ground-truth data.

We conduct extensive experiments on distilling VL models, which contributes to this relatively under-explored field in the current literature.

## **Related Work**

Vision-Language Pre-training (VLP). Recently, large pretrained Transformers (Liu et al. 2019; Lan et al. 2020; Clark et al. 2020; Yang et al. 2019; Ho et al. 2022; Li et al. 2022a; Appalaraju et al. 2021) have started to show improved capability in a variety of language modeling tasks (Zellers et al. 2018; Wang et al. 2018; Williams, Nangia, and Bowman 2017). Researchers have further extended these models to large image-text and video-text pre-training multi-modal models (Lu et al. 2019; Chen et al. 2020; Zhou et al. 2020; Li et al. 2020a; Cho et al. 2020; Zhang et al. 2021; Sun et al. 2019; Zhu and Yang 2020; Miech et al. 2020; Radford et al. 2021; Appalaraju et al. 2024) for visual-linguistic tasks (Goyal et al. 2019; Hudson and Manning 2019; Lei et al. 2018; Mao et al. 2016; Xu et al. 2016; Zhou, Xu, and Corso 2018). These pre-trained models outperform previous approaches (Yu et al. 2018; Yu, Kim, and Kim 2018; Kim, Jun, and Zhang 2018; Anderson et al. 2018; Liu et al. 2020). As their models sizes grow rapidly, recent works have also explored parameter-efficient learning and model compression methods, including adapters (Sung, Cho, and Bansal 2022; Houlsby et al. 2019; Rebuffi, Bilen, and Vedaldi 2018, 2017), prompt tuning (Gu et al. 2021b; Lester, Al-Rfou, and Constant 2021; Li and Liang 2021).

Knowledge Distillation. Knowledge distillation (KD) (Hinton, Vinyals, and Dean 2015) transfers knowledge from a stronger Teacher network (T) to a Student network (S) by minimizing the divergence of their soft response and intermediate features (Gou et al. 2021). Compared with recent approaches to distillation in vision (Zagoruyko and Komodakis 2017; Peng et al. 2019; Tung and Mori 2019; Yang et al. 2022; Chen et al. 2022b; Wu et al. 2022b; Andonian, Chen, and Hamid 2022; He et al. 2022; Wu et al. 2022a) or language tasks (Wang et al. 2020b; Li et al. 2022c; Wang et al. 2021; Ding et al. 2023), distilling VL models is a relatively under-explored field, as pointed out by the review paper of (Chen et al. 2022a). (Fang et al. 2021) claim to be the first to distill vision-language Transformers, and MAD (Wang et al. 2022b) claims to be the first to use multi-modal distillation for VL models. These two papers both focus on distilling Encoder-only VL Transformers. (Gu et al. 2021a; Ma et al. 2022) also involve distilling knowledge from visual and linguistic domains, but their architectures are based upon Mask RCNN (He et al. 2017) and ResNet (He et al. 2016), and apply to the visual task of object detection.

Specifically, for distilling large Transformers to smaller ones, recent work in language distillation (Jiao et al. 2020; Sun et al. 2020; Wang et al. 2020b; Sanh et al. 2019), vision distillation (Qu et al. 2022), and VL distillation (Fang et al. 2021) all show that applying attention map distillation to transfer the rich co-reference relations to Student can boost performance. These approaches apply attention map distillation either on all attention layers or only on the last self/cross-attention layer and the last cross-attention layer. Specifically, for a given multi-head attention layer, most of

these approaches (Fang et al. 2021; Jiao et al. 2020; Sun et al. 2020; Wang et al. 2020b; Sanh et al. 2019) minimize the sum of divergence between the attention matrices of each head of Teacher and Student. However, this formulation only applies when the Teacher and Student have the same number of attention heads. (Qu et al. 2022) minimizes the divergence between the mean attention matrices of all the heads of Teacher and Student. However, as pointed out by (Cao et al. 2020), different attention heads encode different coreference knowledge, hence applying mean reduction may result in knowledge loss.

Different from most approaches for distilling intermediate features in a one-to-one fashion, (Lin et al. 2022) propose a one-to-all spatial matching strategy for distilling Convnets feature maps, allowing each pixel of the Teacher feature to be distilled to all spatial locations in the Student by similarity mapping; (Ji, Heo, and Park 2021) propose to learn to match Teacher and Student features maps in different ResNet layers. Our design is also inspired by these works.

# **Attention Map Alignment Distillation**

In this section, we introduce our proposed Attention Map Alignment Distillation (AMAD) method. We use plain lower case letters x for scalars, bold lower case letters x for vectors, and bold upper case letters x for matrices.

In a multi-head attention layer of Transformer (Vaswani et al. 2017), each entry of the attention matrix for a head is a dot-product of query and key vectors followed by softmax normalization (Bahdanau, Cho, and Bengio 2014). In matrix form, for each head h, if we denote the number of query vectors as q, the number of key vectors as k, and the attention matrix as  $\mathbf{A}_h \in \mathbb{R}^{q \times k}$ , then we have

$$\mathbf{A}_h = \operatorname{softmax}(\mathbf{Q}_h \mathbf{K}_h^T / \sqrt{d_k}) \tag{1}$$

where  $\mathbf{Q}_h$  and  $\mathbf{K}_h$  are the query and key matrices of head h, and  $d_k$  is the dimension of the key as a scaling factor.

For each training data sample (not batched), the attention maps of all heads for a given H-head multi-head attention layer form a tensor of  $[\mathbf{A}_1,\mathbf{A}_2,...,\mathbf{A}_H] \in \mathbb{R}^{H \times q \times k}$ . AMAD aims at distilling the attention maps of the  $H_t$  heads in the Teacher to those of the  $H_s \leq H_t$  heads in Student.

For representational simplicity, let  $n=q\cdot k$ , and  $\mathbf{t}_i\in\mathbb{R}^n$  denote the column vector representing the flattened attention map  $\mathbf{A}_i\in\mathbb{R}^{q\times k}$  of Teacher head i. Let  $\mathbf{s}_j\in\mathbb{R}^n$  denote the column vector representing the flattened attention map  $\mathbf{A}_j\in\mathbb{R}^{q\times k}$  of Student head j.

For a given Teacher head  $\mathbf{t}_i$ , compute its cosine similarity  $w_{ij}$  with each Student head  $\mathbf{s}_j$  as in Equation 2:

$$w_{ij} = \mathbf{t}_i \cdot \mathbf{s}_j / (\|\mathbf{t}_i\|_2 \cdot \|\mathbf{s}_j\|_2) \tag{2}$$

$$a_{ij} = \frac{\exp(w_{ij})}{\sum_{m=1}^{H_s} \exp(w_{im})}$$
(3)

Then as in Equation 3 above, for the given Teacher head, compute its distilling contribution  $a_{ij}$  to each of the Student heads  $j \in \{1, 2, \cdots, H_s\}$  by applying softmax non-linearity on the similarity weights  $w_{ij}$ :

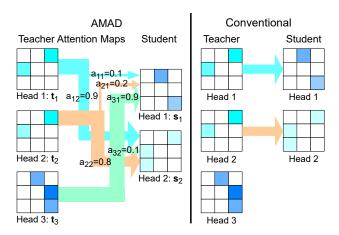


Figure 1: An illustration of AMAD in a toy case, corresponding to Equation 7, where the Teacher has  $H_t = 3$ heads  $(\mathbf{t}_1, \mathbf{t}_2, \mathbf{t}_3)$  and Student has  $H_s = 2$  heads  $(\mathbf{s}_1, \mathbf{s}_2)$ , all with self attention maps of dimension  $n = q \times k = 3 \times 3$ . Different matrix entry colors denote different attention values. On the left, AMAD uses soft-alignment: each Teacher head attention map teaches all the Student heads but in proportion to how close the Student head is to the Teacher head. As in the coloring of the matrices, Teacher heads 1 and 2 are similar to each other, and are relatively similar to Student head 2; While Teacher head 3 is similar to Student head 1. In this case, with AMAD, Teacher heads 1 and 2 instruct Student head 2 more (larger  $a_{12}$ ,  $a_{22}$  and wider arrows above) and instruct Student head 1 less (smaller  $a_{11}, a_{21}$  and narrower arrows above), and Teacher head 3 mainly instructs Student head 1. Also note that the knowledge in the two similar Teacher heads  $t_1$  and  $t_2$  can be compressed mostly to a single Student head  $s_2$ . While on the right, conventional attention map distillation method does not apply when the numbers of heads are different between Teacher and Student: Teacher head 3 has to be discarded in distilling.

Next, we minimize the mean squared error between the given normalized Teacher head attention map  $\mathbf{t}_i$  and the weighted sum of soft-aligned Student head attention maps.

For each Attention head *i* of the Teacher,

$$\mathcal{L}_{\text{AMAD}_i} = \left\| \frac{\mathbf{t}_i}{\|\mathbf{t}_i\|_2} - \sum_{j=1}^{H_s} a_{ij} \cdot \frac{\mathbf{s}_j}{\|\mathbf{s}_j\|_2} \right\|_2^2 \tag{4}$$

The total loss  $\mathcal{L}_{AMAD}$  is the summation of  $\mathcal{L}_{AMAD_i}$  over all the heads of the Teachers,

$$\mathcal{L}_{\text{AMAD}} = \sum_{i=1}^{H_t} \mathcal{L}_{\text{AMAD}_i}$$
 (5)

Now we rewrite the above formulas using a matrix formulation to parallel computations. For each training data sample (not batched), recall that  $n=q\cdot k$ , let matrix  $\mathbf{T}\in\mathbb{R}^{H_t\times n}$  represent the Teacher attention maps of all its heads, whose each row vector is the normalized flattened attention map

 $\mathbf{t}_i^T/\|\mathbf{t}_i\|_2$  of the *i*-th head. Similarly, let matrix  $\mathbf{S} \in \mathbb{R}^{H_s \times n}$  represent the normalized Student attention maps of all its heads. We have (detailed derivation in Appendix),

$$\mathcal{L}_{\text{AMAD}} = \|\mathbf{T} - \text{softmax}_{\text{dim=row}}(\mathbf{T}\mathbf{S}^T)\mathbf{S}\|_2^2$$
 (6)

where the softmax function applies to each row. The pairwise similarity weight matrix  $\mathbf{TS}^T$  before and after softmax are both of dimension  $\mathbb{R}^{H_t \times H_s}$ . Computation with respect to the i-th row of the matrix formulation corresponds to the operations regarding the i-th Teacher head as in Equation 4.

For instance, in the case as in Figure 1, we have,

softmax(
$$\mathbf{TS}^T$$
) =  $(a_{ij})_{H_t \times H_s} = \begin{pmatrix} 0.1 & 0.9 \\ 0.2 & 0.8 \\ 0.9 & 0.1 \end{pmatrix}$  (7)

The above formulas focus on each training data sample (not batched); For implementation, we use batched tensor computations via PyTorch, and we  $L_2$ -normalize each row of the weighted sum softmax( $\mathbf{TS}^T$ )S before calculating the mean squared error. Code is provided in the Appendix.

In contrast to previous attention map distillation methods directly minimizing  $\|\mathbf{T} - \mathbf{S}\|_2^2$  or  $\mathrm{KL}(\mathbf{T}\|\mathbf{S})$ , requiring  $\mathbf{T}$  and  $\mathbf{S}$  of the same shape (Fang et al. 2021; Jiao et al. 2020; Sun et al. 2020; Wang et al. 2020b), AMAD removes the limitation of requiring Teacher and Student to have the same number of attention heads, by letting each Teacher head teaches all the Student heads with the contribution being more for the Student heads with which it has a higher weight (higher cosine similarity), and supports flexible and smooth distillation because of the soft semantic alignment mechanism.

### **Formulation Variants**

We refer to the above formulation as **Variant 1** and the corresponding loss as  $\mathcal{L}_{AMAD-1}$ , and we also explore the following ablative baselines and variants:

**Baseline: One-to-one Distillation**. In this baseline, following (Fang et al. 2021; Jiao et al. 2020; Sun et al. 2020; Wang et al. 2020b), we distill the attention maps in a one-to-one fashion. Note that different from previous work, we have more heads in the Teacher than in Student,  $H_t \geq H_s$ , so we distill the first  $H_s$  heads in Teacher to the  $H_s$  heads in Student, respectively, the extra  $H_t - H_s$  Teacher heads are ignored during distillation, as in the right part of Figure 1:

$$\mathcal{L}_{\text{KD-ATT}} = \|\mathbf{S} - \mathbf{T}[: H_s, :]\|_2^2$$
 (8)

**Variant 2: KL Divergence**. In this variant, we minimize the sum of Kullback–Leibler divergence between the aligned weighted sum of Student multi-head attention map distributions and the Teacher distributions:

$$\mathcal{L}_{\text{AMAD-2}} = \text{KL}(\mathbf{T} || \text{softmax}(\mathbf{T}\mathbf{S}^T)\mathbf{S})$$
 (9)

where the Teacher T and the Student S are all  $L_1$ -normalized by each row in all KL variants, instead of  $L_2$ -normalized. Both input and target contain  $q \cdot H_t$  distributions, and KL( $\cdot$ ) is computed for each of these  $q \cdot H_t$  distributions and then summed up.

**Variants 3: Parameterized Projection**. We borrow the idea from attention mechanisms to apply a learnable linear projection W on each flattened vector of attention map  $\mathbf{s}_i$  of

Student head j before computing similarity and alignment:  $\mathbf{W}\mathbf{s}_{j} + \mathbf{b}$ . In matrix form:

$$\tilde{\mathbf{S}} = \text{ReLU}(\mathbf{S}\mathbf{W}^T + \mathbf{b}) \tag{10}$$

$$\mathcal{L}_{\text{AMAD-3}} = \text{KL}(\mathbf{T} || \text{softmax}(\mathbf{T}\tilde{\mathbf{S}}^T)\tilde{\mathbf{S}})$$
 (11)

where  $\mathbf{W} \in \mathbb{R}^{n \times n}$  is a learnable matrix and  $\mathbf{b}$  is the bias. **Variant 4: Token-level Alignment**. Here, we adopt a finer token-level granularity of aligning correspondence and allow independence for the soft alignment weights  $w_{ij}$  for different query attention vectors in Teacher / Student heads.

ferent query attention vectors in Teacher / Student heads. Formally, denote  $\mathbf{T}_l \in \mathbb{R}^{H_t \times k}$  as the matrix whose h-th row vector is the l-th row vector  $\mathbf{t}_{h,l}^T \in \mathbb{R}^k$  in the Teacher's h-th head attention map  $\mathbf{A}_h$ , and  $\mathbf{S}_l \in \mathbb{R}^{H_s \times k}$  as the matrix whose each h-th row vector is  $\mathbf{s}_{h,l}^T \in \mathbb{R}^k$ . We have:

$$\mathcal{L}_{\text{AMAD-4}} = \sum_{l=1}^{q} \text{KL}(\mathbf{T}_{l} || \text{softmax}(\mathbf{T}_{l} \mathbf{S}_{l}^{T}) \mathbf{S}_{l})$$
 (12)

where  $KL(\cdot)$  is computed for each of the  $H_t$  distributions and then summed up. In this variant, the weight matrices softmax $(\mathbf{T}_l\mathbf{S}_l^T)$  are different for each l-th query attention vector, in contrast to the unified same weight matrix of softmax $(\mathbf{TS}^T)$  for all the queries in previous variants.

We report ablation results for each variant in Table 8. More theoretical analysis is in Appendix.

# **Experimental Setup**

We distill VL-T5 *base* to *small*; and distill BLIP *large* to *base*. Table 1 summarizes their architectural backbones. Visualized architecture and distilling pipeline are in Appendix.

#### **Knowledge Distillation (KD)**

For training efficiency, we only apply distillation in downstream fine-tuning, no distillation involved in pre-training.

As in previous work (Hinton, Vinyals, and Dean 2015; Fang et al. 2021; Jiao et al. 2020; Sanh et al. 2019), we apply classification distillation loss on the classifier output logits. For VQA with a single classifier head,

$$\mathcal{L}_{KD} = CE(\mathbf{z}^S / \tau_d, \mathbf{z}^T / \tau_d)$$
 (13)

where  $\tau_d$  denotes the distillation temperature (Hinton, Vinyals, and Dean 2015), which we simply use 1, as in (Fang et al. 2021).  $\mathbf{z}^S$  and  $\mathbf{z}^T$  refer to the logits from Student and Teacher classifier. CE denotes Cross Entropy, i.e.  $p_i = \frac{\exp(z_i/\tau_d)}{\sum_k \exp(z_k/\tau_d)}$  and  $\mathcal{L}_{\text{KD}} = \sum_i p_i^T \cdot \log(p_i^S)$ .

For auto-regressive captioning and translation tasks, the Teacher and the Student both take ground-truth answer to-ken sequence as input in Teacher Forcing style to maintain consistency for distillation (Beyer et al. 2022),

$$\mathcal{L}_{KD} = \sum_{j=1}^{|y|} CE(\mathbf{z}_j^S / \tau_d, \mathbf{z}_j^T / \tau_d)$$
 (14)

where  $\mathbf{z}_{j}^{S}$  and  $\mathbf{z}_{j}^{T}$  denote logits for the j-th output token from Student and Teacher classifier, and |y| denotes length of seq.

	VL	#Learnable	Vision Stream				Language (Multi-modal) Stream				
	Model	Params	Backbone	#Layers	$d_{ m model}$	#Heads	Backbone	#Layers	$d_{ m model}$	#Heads	
Teacher Student	VL-T5 base VL-T5 small	220M 60M	Faster R-CNN (frozen)	-	-	-	$ $ T5 $ $ w/ $\mathcal{L}_{AMAD}$	12+12 6 + 6	768 512	12 8	
Teacher Student	BLIP <i>large</i> BLIP <i>base</i>	446M 210M	ViT w/ L <sub>amad</sub>	24 12	1024 768	16 12	$egin{array}{c} \operatorname{BERT} \ \mathbf{w} / \ \mathcal{L}_{\mathbf{AMAD}} \end{array}$	12 12	768 768	12 12	

Table 1: Model architectures with details of their Transformer sub-modules. We distill VL-T5 *base* to *small*; and distill BLIP *large* to *base*. AMAD loss is applied on all Transformer modules, including T5-Encoder + Decoder for VL-T5, and ViT + BERT for BLIP. Conventional one-to-one attention map distillation does not apply to some of these modules when Teacher and Student have different numbers of attention heads; AMAD removes the same-number-of-heads constraint and works here.

The overall training objective for the Student  $\mathcal{L}_{TOTAL}$  is a weighted sum of the classification distillation loss  $\mathcal{L}_{KD}$  and the proposed loss  $\mathcal{L}_{AMAD}$ ,

$$\mathcal{L}_{\text{TOTAL}} = \mathcal{L}_{\text{KD}} + \alpha \mathcal{L}_{\text{AMAD}} \tag{15}$$

We apply  $\mathcal{L}_{AMAD}$  to distill the self/cross-attention maps of the last (Wang et al. 2020b; Fang et al. 2021) layers of each stream.  $\alpha$  is tuned so that  $\mathcal{L}_{KD}$  and  $\mathcal{L}_{AMAD}$  scale similarly. We do not add ground-truth loss (Beyer et al. 2022).

#### **Pre-training and Fine-tuning**

As VL-T5 *small* is not released, we pretrain it by ourselves, adopting the same setting as how they pretrain *base*. After uni-modal pretraining its T5 and Faster R-CNN (Ren et al. 2015) sub-modules, it is then VL pretrained on MS COCO (Lin et al. 2014; Chen et al. 2015), Visual Genome (Krishna et al. 2016), VQA-2.0 (Goyal et al. 2019), GQA (Hudson and Manning 2019), and Visual7W (Zhu et al. 2016). For VL-T5 *base* and BLIP, we directly use their released pretrained checkpoints.

After VL pre-training the Teacher and the Student, we finetune the Teacher on downtream tasks, adopting the same settings as in VL-T5 or BLIP. The Teacher model is then frozen and ready to be distilled. We then distill the Teacher to the Student with  $\mathcal{L}_{TOTAL}$  on downstream tasks.

In some of our ablative settings, we do not conduct any VL pre-training for Student: After loading the language-only pre-trained language / multi-modal branch checkpoint and the vision-only pre-trained vision branch checkpoint, we directly finetune the Student on downstream VL tasks with distillation; Teacher is always pre-trained and finetuned.

# **Downstream Fine-tuning Datasets**

We demonstrate visual question-answering performance on VQA-2.0 dataset. We report results on *Karpathy test*, *test-std* and *test-dev* via: https://visualqa.org/challenge.html.

We evaluate image captioning performance on MS COCO dataset (Chen et al. 2015). As in (Cho et al. 2021; Fang et al. 2021; Li et al. 2022b), we use the *Karparthy split* (Karpathy and Fei-Fei 2015), which re-splits train2014 and val2014 images (Lin et al. 2014) into ~11K/5K/5K for train / validation / test. We report BLEU@4 (B) (Papineni et al. 2002), CIDEr (C) (Vedantam, Zitnick, and Parikh 2015), METEOR (M) (Banerjee and Lavie 2005), SPICE (S) (Anderson et al. 2016) evaluation metrics.

We also evaluate multi-modal machine translation performance on Multi30K dataset (Elliott et al. 2016), where models translate English text to German given context images. We report BLEU@4 score using SacreBLEU (Post 2018).

We report our implementation details in Appendix.

# **Results and Analysis**

Table 2 shows results on distilling VL-T5 and BLIP with AMAD, in comparison to recent vision-language models. The higher the better for all metrics. As in previous work (Fang et al. 2021; Cho et al. 2021), captioning performance are shown with their cross entropy optimization variants instead of CIDEr optimization variants. Figure 2 visualizes the effect of model size and number of VL pre-training images to VQA performance for recent models. Overall, Table 2 and Figure 2 support the following arguments:

**1. Reducing VL model size within same family of models causes performance drops:** If all trained with ground-truth supervision without distillation, VL-T5 *small* performs

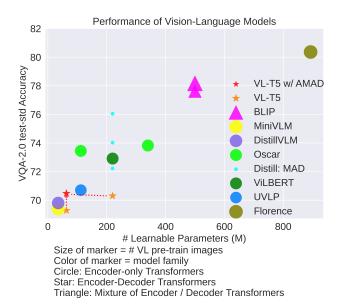


Figure 2: VL model performance with respect to # learnable params (X-axis) and # VL pre-training images (marker size). AMAD pushes VL-T5 towards upper-left (orange to red).

		# Learnable	# VLP	Distilled From	VQ	A-2.0 Acc	: ↑	C	OCO Ca	ptioning	
	Method	Parameters	Images	Which Model	Karpathy	test-std	test-dev	B↑	C↑	M↑	$S \uparrow$
				CNN-LSTM-bas	sed Models						
1	Up-Down *	-	108K	Х	-	70.34		36.2	113.5	27.0	20.3
2	GLIED *	18.3M	-	self distillation	-	-	-	37.9	118.2	28.3	21.2
	Encoder-Only Transformers										
3	ViLBERT *	220M	3M	Х	-	70.92	70.55	-	-	-	-
4	UNITER base *	220M	4M	X	-	72.91	72.70	-		<del>-</del>	-
5	Unified VLP *	112M	3M	X	-	70.7	-	36.5	116.9	28.4	21.2
6	OSCAR base *	112M	4M	X	-	73.44	73.16	36.5	123.7	30.3	23.1
7	MiniVLM *	35M	7M	X	-	-	-	34.3	116.7	28.1	21.3
8	MiniVLM *	35M	14M	X	-	69.4	69.1	35.6	119.8	28.6	21.6
9	DistillVLM *	35M	7M	Oscar base	-	69.8	69.6	35.6	120.8	28.7	22.1
10	MAD-ViLBERT *	220M	-	CLIP-V+T	-	72.22	-	-	-	-	-
11	MAD-UNITER *	220M	-	CLIP-V+T	-	74.02	-	-	-	-	-
			]	Encoder-Decoder	Transformers	3					
12	OFA huge *	930M	15M	×	-	82.0	82.0	43.9	145.3	31.8	24.8
13	VL-T5 base *	220M	180K	Х	67.9	70.30	-	34.5	116.5	28.7	21.9
14	VL-T5 small †	63M	180K	X	66.72	69.28	69.04	32.8	108.2	27.0	20.4
15	Ours VL-T5 small	63M	180K	VL-T5 base †	68.06	70.47	70.41	33.9	114.4	28.3	21.5
16	VL-T5 base *	220M	0	Х	-	-	-	32.6	109.4	28.2	21.0
17	VL-T5 small †	63M	0	X	56.44	-	58.47	30.8	101.4	26.3	19.5
18	Ours VL-T5 small	63M	0	VL-T5 base †	67.79	70.25	70.06	33.3	112.9	28.0	21.3
			Mixtu	e of Encoder / De	coder Transf	ormers					
19	BLIP large *	446M	129M	Х	-	-	-	40.4	136.7	-	_
20	BLIP base *	210M	129M	X	_	78.32	78.25	39.7	133.3	-	-
21	Ours BLIP base	210M	129M	BLIP large *	-	-	-	40.0	134.1	31.0	23.9
22	BLIP base *	210M	14M	Х	<u> </u>	77.62	77.54	38.6	129.7	_	_
23	BLIP base †	210M	0	×	_	-	-	34.7	115.8	28.5	21.5
24	Ours BLIP base	210M	0	BLIP large *	-	-	-	38.7	129.6	30.4	23.3

Table 2: Results on distilling VL-T5 and BLIP with AMAD, with comparisons to recent VL models. Results with \* are reported from their papers; with † are trained by ourselves; also in following tables. AMAD narrows the performance gaps caused by reducing model size or removing VL pre-training (VLP). Furthermore, to our surprise, the finetuning distilled VL-T5 *small* w/o VL pre-training (row 18) even outperforms VL pre-trained and GT fine-tuned VL-T5 *small* (row 14). Also, for BLIP, fine-tuning distillation (row 24) compensates for 14M-scale VL pre-training (row 22), but cannot fully compensate for 129M-scale pre-training (row 20).

worse than *base* (row 13-14, 16-17); and BLIP *base* worse than *large* (row 19-20). This is also visualized in Figure 2.

- **2. For a same model, removing VL pre-training / Pre-training with fewer images degrades performance:** Performance of VL-T5 *base* and *small* both drop significantly if not VL pre-trained (row 13 vs 16; 14 vs 17). Even when VL pretrained, smaller numbers of VL pre-training images cause performance drop in MiniVLM (row 7 vs 8) and in BLIP (row 20, 22, 23; pink triangles in Figure 2).
- **3. Distilling with AMAD narrows the aforementioned gaps caused by shrinking model size or by removing VL pre-training:** Supported by results of distilling VL pre-trained VL-T5 (row 14 vs 15); distilling non-VL-pretrained VL-T5 (row 17 vs 18); and distilling BLIP (row 20 vs 21).

Note that, DistillVLM (row 9) (Fang et al. 2021) distilled Oscar *base* (row 6) to MiniVLM (row 8) (Wang et al. 2020a) and achieved 0.4% VQA accuracy boost and 1.0% Captioning CIDEr score boost. And they claimed to be the first work to apply KD in training VL models. Neither MiniVLM nor

DistillVLM has released their model or code.

**4.** Knowledge distillation compensates to some degree for the absence of VL pre-training: When we do not conduct any VL pre-training for VL-T5 *small*, fine-tuning distilled VL-T5 *small* (row 18) even outperforms the ground-truth supervised VL pre-trained and fine-tuned VL-T5 *small* baseline (row 14). It also outperforms non-VL-pretrained VL-T5 *base* (row 16). The performance is also rather comparable to other recent pre-trained and finetuned VL models. For BLIP, distillation (row 24) compensates for the absence of 14M-scale VL pre-training (row 22), but cannot compensate for 129M-scale pre-training (row 20).

One possible explanation is that the knowledge obtained in the pre-training stage of the Teacher can somehow be distilled to the Student in the downstream fine-tuning process when the Student tries to mimic the Teacher's classification logits and tries to align and mimic the Teacher's attention maps, even if the Student has no access to the pre-training data by itself. The Teacher's attention maps contain valuable

VL-T5 Model	#Params	test-2016	test-2017	test-2018
Teacher † Student	220M 60M	44.00 41.90	39.40 36.85	37.00 34.02
Our Student $\Delta$	60M	43.88 <b>+1.98</b>	38.70 <b>+1.85</b>	36.64 <b>+2.62</b>

Table 3: Multi30K English-German translation BLEU@4 score. The VL-T5 *small* Student distilled with AMAD outperforms ground-truth (GT) fine-tuned VL-T5 *small*.

BLIP Model	#Params	В	С	M	S
Teacher * Student	446M 210M	40.4 22.3	136.7 65.1	20.9	13.5
Student w/ AMAD $\Delta$	210M	29.7 <b>+7.4</b>	95.0 <b>+29.9</b>	25.4 + <b>4.5</b>	18.5 <b>+5.0</b>

Table 4: Results on COCO Captioning karpathy test split. All BLIP Students are w/o any pretraining, i.e. the ViT and BERT modules are randomly initialized for all Students. Results for pretrained and finetuning distilled BLIP models are reported in Appendix. The distilled Student outperforms the GT trained Student by a surprisingly large margin, although the performance still degrades a lot because of neither VL nor uni-modal pre-trained. This indicates that uni-modal pretraining is still necessary even when finetuning distillation is applied.

	#Params	$\mathcal{L}_{ ext{KD}}$	$\mathcal{L}_{ ext{AMAD}}$	В	C	M	S
T †	220M			34.2	115.1	28.3	21.6
S † S ;	60M 60M 60M	<i>y y</i>		32.8 33.1 33.6	108.2 111.8 113.4	27.0 27.9 28.1	20.4 21.2 21.3
<u>S</u>	60M	✓	✓	33.9	114.4	28.3	21.5

Table 5: Ablation results of VL-T5 models on COCO Captioning karpathy test split. 'T' denotes Teacher, and 'S' denotes Student; Also in following tables. All VL-T5 Students are first VL pretrained. Ablation results w/o VL pretraining are in Appendix. Distilling with  $\mathcal{L}_{AMAD}$  outperforms logits distillation with  $\mathcal{L}_{KD}$ , and narrows the CIDEr gap between the *small* Student and the *base* Teacher to only 0.7% with 72% less parameters. The VL-T5 Teacher is reproduced, as their fine-tuned checkpoints are not released, also in following Tables. We also compare AMAD with RKD-D (Park et al. 2019) (the row with ‡ in the Table), a distance-based similarity distillation method: We treat attention maps of all heads as a whole single feature and apply RKD-D on that, since #heads are different between Teacher and Student.

intra-modal and cross-modal coreference relations learned from the pre-training dataset, and  $\mathcal{L}_{AMAD}$  helps the Student to inherit the rich learned representation from the Teacher.

Besides these observations 1-4, we show in Table 3 that AMAD can generalize well to language-heavy translation task. Table 4 unveils one limitation that, although fine-tuning

	#Params	$\mathcal{L}_{ ext{KD}}$	$\mathcal{L}_{ ext{AMAD}}$	Karpathy	std	dev
T†	226M			68.75	71.34	71.23
S † S S	63M 63M 63M	<i>\ \</i>	/	66.72 67.74 <b>68.06</b>	69.28 70.19 <b>70.47</b>	69.04 70.10 <b>70.41</b>

Table 6: Ablation results of VL-T5 models on VQA-2.0 test splits. All Students are first VL pretrained.

	#Params	$\mathcal{L}_{ ext{KD}}$	$\mathcal{L}_{ ext{AMAD}}$	Karpathy	std	dev
T †	226M			68.75	71.34	71.23
S† S	63M	,		56.44	-	58.47
S	63M 63M	1	✓	66.39 <b>67.79</b>	70.25	68.71 <b>70.06</b>

Table 7: Ablation results of VL-T5 models on VQA-2.0 test splits. All Students are w/o VL pretraining, i.e. initialized with language-only pre-trained T5 and vision-only pre-trained Faster R-CNN checkpoints. This shows that vanilla finetuning logits distillation already compensates to some degree for the absence of VL pre-training, and AMAD narrows the gap further.

VL-T5	#Params	Loss Variant	Acc	Δ
Teacher †	226M		68.75	
Student † Student	63M 63M	Ground-Truth $\mathcal{L}_{\text{KD}}$	66.72 67.74	+1.02
Student	63M	$\mathcal{L}_{ ext{KD}} + \mathcal{L}_{ ext{KD-ATT}}$	67.73	+1.01
Student Student Student Student	63M 63M 63M 63M	$\mathcal{L}_{ ext{KD}} + \mathcal{L}_{ ext{AMAD-1}}$ $\mathcal{L}_{ ext{KD}} + \mathcal{L}_{ ext{AMAD-2}}$ $\mathcal{L}_{ ext{KD}} + \mathcal{L}_{ ext{AMAD-3}}$ $\mathcal{L}_{ ext{KD}} + \mathcal{L}_{ ext{AMAD-4}}$	67.96 <b>68.06</b> 68.02 68.05	+1.24 +1.34 +1.30 +1.33
Student	63M	w/ AMAD mean std-err of the mean	$68.02 \pm 0.02$	+1.30 ± 0.02

Table 8: Ablation results on VQA-2.0 Karpathy test split. All VL-T5 Students are first VL pretrained. The baseline of distilling the attention maps from the first 8 heads of Teacher to those of the 8 heads of Student in a one-to-one fashion with  $\mathcal{L}_{\text{KD-ATT}} = \|S - T[:H_s,:]\|_2^2$  does not help improve performance than distilling with  $\mathcal{L}_{\text{KD}}$  only, maybe because the extra  $H_t - H_s$  Teacher heads are discarded during distillation, causing forced knowledge loss. Meanwhile, all  $\mathcal{L}_{\text{AMAD}}$  variants help improve performance consistently. KL variants for  $\mathcal{L}_{\text{AMAD}}$  (2, 3, 4) perform slightly better than MSE  $\mathcal{L}_{\text{AMAD-1}}$ . We have not observed significant performance change brought by learnable projection  $\mathcal{L}_{\text{AMAD-3}}$  or token-level alignment  $\mathcal{L}_{\text{AMAD-4}}$ , compared to  $\mathcal{L}_{\text{AMAD-2}}$ .

distillation might close the performance gap of removing VL pretraining, uni-modal pretraining is still necessary.

### **Ablations**

Effects of logits distillation  $\mathcal{L}_{KD}$  and AMAD  $\mathcal{L}_{AMAD}$ : We ablates on captioning and VQA tasks in w/ and w/o VL pre-

train settings in Table 5-7 and in Appendix. In w/o VL pretrain settings, Students are finetuned directly after loading the vision-only pretrained visual branch and the languageonly pretrained linguistic branch sub-modules.

**AMAD Variants and Baselines** are analyzed in Table 8. We present visualizations, more ablations, and qualitative analysis of  $\mathcal{L}_{AMAD}$  distilled attention maps in Appendix.

## **Appendix**

Please kindly refer to the Appendix via the following link: https://www.amazon.science/publications/no-head-leftbehind-multi-head-alignment-distillation-for-transformers

In the Appendix, we provide a further illustration on the matrix-form loss derivation and formulation of different AMAD variants in Section A; more ablation results and a reverse experiment of distilling a smaller Teacher to a larger Student in Section B; implementation details including training environment, training time, and hyperparameters in Section C; an additional illustration of the distillation workflow for our experimental setup and the architecture of VL-T5 in Section D; visualizations of AMAD distilled cross- and self- attention maps in Section E; and PyTorch code for the proposed AMAD loss in Section F.

#### Conclusion

We have proposed the Attention Map Alignment Distillation (AMAD) method to distill attention maps from a Teacher to a Student Transformer with different numbers of attention heads. AMAD narrows the performance gap between the large Teacher and the small Student in both discriminative VQA and auto-regressive generative captioning / translation tasks. Our ablation further suggests that fine-tuning knowledge distillation can compensate to some degree for the absence of VL pre-training for VL Transformers. However, uni-modal pre-training is still necessary. Exploring theoretically or empirically why fine-tuning distillation can compensate for VL pre-training is potentially intriguing for future work.

#### References

Alayrac, J.-B.; Donahue, J.; Luc, P.; Miech, A.; Barr, I.; Hasson, Y.; Lenc, K.; Mensch, A.; Millican, K.; Reynolds, M.; et al. 2022. Flamingo: a visual language model for few-shot learning. *arXiv preprint arXiv:2204.14198*.

Anderson, P.; Fernando, B.; Johnson, M.; and Gould, S. 2016. SPICE: Semantic Propositional Image Caption Evaluation. In *ECCV*.

Anderson, P.; He, X.; Buehler, C.; Teney, D.; Johnson, M.; Gould, S.; and Zhang, L. 2018. Bottom-Up and Top-Down Attention for Image Captioning and Visual Question Answering. 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, 6077–6086.

Andonian, A.; Chen, S.; and Hamid, R. 2022. Robust Cross-Modal Representation Learning with Progressive Self-Distillation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 16430–16441.

Appalaraju, S.; Jasani, B.; Kota, B. U.; Xie, Y.; and Manmatha, R. 2021. Docformer: End-to-end transformer for document understanding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 993–1003.

Appalaraju, S.; Tang, P.; Dong, Q.; Sankaran, N.; Zhou, Y.; and Manmatha, R. 2024. DocFormerv2: Local Features for Document Understanding - Full Paper. *AAAI*.

Bahdanau, D.; Cho, K.; and Bengio, Y. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*.

Banerjee, S.; and Lavie, A. 2005. METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments. In *ACL Workshop*.

Beyer, L.; Zhai, X.; Royer, A.; Markeeva, L.; Anil, R.; and Kolesnikov, A. 2022. Knowledge distillation: A good teacher is patient and consistent. In *Proceedings of the IEEE/CVF Conference on CVPR*, 10925–10934.

Brown, T. B.; Mann, B.; Ryder, N.; Subbiah, M.; Kaplan, J.; Dhariwal, P.; Neelakantan, A.; Shyam, P.; Sastry, G.; Askell, A.; Agarwal, S.; Herbert-Voss, A.; Krueger, G.; Henighan, T.; Child, R.; Ramesh, A.; Ziegler, D. M.; Wu, J.; Winter, C.; Hesse, C.; Chen, M.; Sigler, E.; Litwin, M.; Gray, S.; Chess, B.; Clark, J.; Berner, C.; McCandlish, S.; Radford, A.; Sutskever, I.; and Amodei, D. 2020. Language Models are Few-Shot Learners. In *NeurIPS*.

Cao, J.; Gan, Z.; Cheng, Y.; Yu, L.; Chen, Y.-C.; and Liu, J. 2020. Behind the scene: Revealing the secrets of pre-trained vision-and-language models. In *ECCV*.

Chen, F.; Zhang, D.; Han, M.; Chen, X.; Shi, J.; Xu, S.; and Xu, B. 2022a. VLP: A Survey on Vision-Language Pretraining. *ArXiv*, abs/2202.09061.

Chen, X.; Cao, Q.; Zhong, Y.; Zhang, J.; Gao, S.; and Tao, D. 2022b. DearKD: Data-Efficient Early Knowledge Distillation for Vision Transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 12052–12062.

Chen, X.; Fang, H.; Lin, T.-Y.; Vedantam, R.; Gupta, S.; Dollar, P.; and Zitnick, C. L. 2015. Microsoft COCO Captions: Data Collection and Evaluation Server. Arxiv.

Chen, X.; Wang, X.; Changpinyo, S.; Piergiovanni, A.; Padlewski, P.; Salz, D.; Goodman, S.; Grycner, A.; Mustafa, B.; Beyer, L.; et al. 2022c. Pali: A jointly-scaled multilingual language-image model. *arXiv preprint arXiv:2209.06794*.

Chen, Y.-c.; Li, L.; Yu, L.; Kholy, A. E.; Ahmed, F.; Gan, Z.; Cheng, Y.; and Liu, J. 2020. UNITER: UNiversal Image-TExt Representation Learning. In *ECCV*.

Cho, J.; Lei, J.; Tan, H.; and Bansal, M. 2021. Unifying Vision-and-Language Tasks via Text Generation. In *International Conference on Machine Learning*.

Cho, J.; Lu, J.; Schwenk, D.; Hajishirzi, H.; and Kembhavi, A. 2020. X-LXMERT: Paint, Caption and Answer Questions with Multi-Modal Transformers. In *EMNLP*.

Clark, K.; Luong, M.-T.; Le, Q. V.; and Manning, C. D. 2020. Electra: Pre-training text encoders as discriminators rather than generators. In *ICLR*.

- Devlin, J.; Chang, M.-W.; Lee, K.; and Toutanova, K. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *NAACL*.
- Ding, Z.; Jiang, G.; Zhang, S.; Guo, L.; and Lin, W. 2023. SKDBERT: Compressing BERT via Stochastic Knowledge Distillation. *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(6): 7414–7422.
- Dosovitskiy, A.; Beyer, L.; Kolesnikov, A.; Weissenborn, D.; Zhai, X.; Unterthiner, T.; Dehghani, M.; Minderer, M.; Heigold, G.; Gelly, S.; Uszkoreit, J.; and Houlsby, N. 2021. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. *ICLR*.
- Elliott, D.; Frank, S.; Sima'an, K.; and Specia, L. 2016. Multi30K: Multilingual English-German Image Descriptions. In *ACL Workshop*, 70–74.
- Fang, Z.; Wang, J.; Hu, X.; Wang, L.; Yang, Y.; and Liu, Z. 2021. Compressing Visual-linguistic Model via Knowledge Distillation. *ICCV*.
- Gou, J.; Yu, B.; Maybank, S. J.; and Tao, D. 2021. Knowledge distillation: A survey. *International Journal of Computer Vision*, 129(6): 1789–1819.
- Goyal, Y.; Khot, T.; Agrawal, A.; Summers-Stay, D.; Batra, D.; and Parikh, D. 2019. Making the V in VQA Matter: Elevating the Role of Image Understanding in Visual Question Answering. *International Journal of Computer Vision*.
- Gu, X.; Lin, T.-Y.; Kuo, W.; and Cui, Y. 2021a. Openvocabulary Object Detection via Vision and Language Knowledge Distillation. In *ICLR*.
- Gu, Y.; Han, X.; Liu, Z.; and Huang, M. 2021b. PPT: Pre-trained Prompt Tuning for Few-shot Learning. *ArXiv*, abs/2109.04332.
- He, K.; Gkioxari, G.; Dollár, P.; and Girshick, R. 2017. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, 2961–2969.
- He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 770–778.
- He, R.; Sun, S.; Yang, J.; Bai, S.; and Qi, X. 2022. Knowledge Distillation as Efficient Pre-training: Faster Convergence, Higher Data-efficiency, and Better Transferability. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 9161–9171.
- Hinton, G.; Vinyals, O.; and Dean, J. 2015. Distilling the Knowledge in a Neural Network. In *NIPS Deep Learning and Representation Learning Workshop*.
- Ho, C.-H.; Appalaraju, S.; Jasani, B.; Manmatha, R.; and Vasconcelos, N. 2022. YORO-Lightweight End to End Visual Grounding. In *European Conference on Computer Vision ECCV CAMP Workshop*.
- Houlsby, N.; Giurgiu, A.; Jastrzebski, S.; Morrone, B.; De Laroussilhe, Q.; Gesmundo, A.; Attariyan, M.; and Gelly, S. 2019. Parameter-efficient transfer learning for NLP. In *International Conference on Machine Learning*, 2790–2799. PMLR.

- Hudson, D. A.; and Manning, C. D. 2019. GQA: A new dataset for real-world visual reasoning and compositional question answering. In *CVPR*. ISBN 9781728132938.
- Ji, M.; Heo, B.; and Park, S. 2021. Show, attend and distill: Knowledge distillation via attention-based feature matching. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, 7945–7952.
- Jiao, X.; Yin, Y.; Shang, L.; Jiang, X.; Chen, X.; Li, L.; Wang, F.; and Liu, Q. 2020. TinyBERT: Distilling BERT for Natural Language Understanding. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, 4163–4174. Online: Association for Computational Linguistics.
- Karpathy, A.; and Fei-Fei, L. 2015. Deep Visual-Semantic Alignments for Generating Image Descriptions. In *CVPR*. ISBN 9781467369640.
- Kim, J.-h.; Jun, J.; and Zhang, B.-t. 2018. Bilinear Attention Networks. In *NeurIPS*, 1–12.
- Krishna, R.; Zhu, Y.; Groth, O.; Johnson, J.; Hata, K.; Kravitz, J.; Chen, S.; Kalantidis, Y.; Jia-Li, L.; Shamma, D. A.; Michael Bernstein; and Fei-Fei, L. 2016. Visual Genome: Connecting Language and Vision Using Crowdsourced Dense Image Annotations. *International Journal of Computer Vision*.
- Lan, Z.; Chen, M.; Goodman, S.; Gimpel, K.; Sharma, P.; and Soricut, R. 2020. Albert: A lite bert for self-supervised learning of language representations. In *ICLR*.
- Lei, J.; Yu, L.; Bansal, M.; and Berg, T. L. 2018. Tvqa: Localized, compositional video question answering. In *EMNLP*.
- Lester, B.; Al-Rfou, R.; and Constant, N. 2021. The Power of Scale for Parameter-Efficient Prompt Tuning. In *EMNLP*.
- Li, C.; Fehérvári, I.; Zhao, X.; Macêdo, I.; and Appalaraju, S. 2022a. SeeTek: Very Large-Scale Open-set Logo Recognition with Text-Aware Metric Learning. 2022 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV), 587–596.
- Li, J.; Li, D.; Savarese, S.; and Hoi, S. 2023. BLIP-2: Bootstrapping Language-Image Pre-training with Frozen Image Encoders and Large Language Models. In *ICML*.
- Li, J.; Li, D.; Xiong, C.; and Hoi, S. 2022b. BLIP: Bootstrapping Language-Image Pre-training for Unified Vision-Language Understanding and Generation. In *ICML*.
- Li, L.; Chen, Y.-C.; Yu Cheng, Z. G.; Yu, L.; and Liu, J. 2020a. HERO: Hierarchical Encoder for Video+Language Omni-representation Pre-training. In *EMNLP*.
- Li, X.; Yin, X.; Li, C.; Zhang, P.; Hu, X.; Zhang, L.; Wang, L.; Hu, H.; Dong, L.; Wei, F.; Choi, Y.; and Gao, J. 2020b. Oscar: Object-Semantics Aligned Pre-training for Vision-Language Tasks. In *ECCV*.
- Li, X. L.; and Liang, P. 2021. Prefix-Tuning: Optimizing Continuous Prompts for Generation. Arxiv.
- Li, Z.; Wang, Z.; Tan, M.; Nallapati, R.; Bhatia, P.; Arnold, A.; Xiang, B.; and Roth, D. 2022c. DQ-BART: Efficient Sequence-to-Sequence Model via Joint Distillation and Quantization. *arXiv preprint arXiv:2203.11239*.

- Lin, S.; Xie, H.; Wang, B.; Yu, K.; Chang, X.; Liang, X.; and Wang, G. 2022. Knowledge Distillation via the Target-aware Transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 10915–10924.
- Lin, T. Y.; Maire, M.; Belongie, S.; Hays, J.; Perona, P.; Ramanan, D.; Dollár, P.; and Zitnick, C. L. 2014. Microsoft COCO: Common Objects in Context. In *ECCV*. ISBN 978-3-319-10601-4.
- Liu, F.; Ren, X.; Liu, Y.; Lei, K.; and Sun, X. 2020. Exploring and distilling cross-modal information for image captioning. *arXiv* preprint arXiv:2002.12585.
- Liu, Y.; Ott, M.; Goyal, N.; Du, J.; Joshi, M.; Chen, D.; Levy, O.; Lewis, M.; Zettlemoyer, L.; and Stoyanov, V. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Lu, J.; Batra, D.; Parikh, D.; and Lee, S. 2019. ViLBERT: Pretraining Task-Agnostic Visiolinguistic Representations for Vision-and-Language Tasks. In *NeurIPS*.
- Ma, Z.; Luo, G.; Gao, J.; Li, L.; Chen, Y.; Wang, S.; Zhang, C.; and Hu, W. 2022. Open-Vocabulary One-Stage Detection with Hierarchical Visual-Language Knowledge Distillation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 14074–14083.
- Mao, J.; Huang, J.; Toshev, A.; Camburu, O.; Yuille, A.; and Murphy, K. 2016. Generation and Comprehension of Unambiguous Object Descriptions. In *CVPR*.
- Miech, A.; Alayrac, J.-B.; Smaira, L.; Laptev, I.; Sivic, J.; and Zisserman, A. 2020. End-to-end learning of visual representations from uncurated instructional videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- Papineni, K.; Roukos, S.; Ward, T.; and Zhu, W. W.-j. 2002. BLEU: a Method for Automatic Evaluation of Machine Translation. In *ACL*. ISBN 1-55860-883-4.
- Park, W.; Kim, D.; Lu, Y.; and Cho, M. 2019. Relational knowledge distillation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 3967–3976.
- Peng, B.; Jin, X.; Liu, J.; Li, D.; Wu, Y.; Liu, Y.; Zhou, S.; and Zhang, Z. 2019. Correlation congruence for knowledge distillation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 5007–5016.
- Post, M. 2018. A Call for Clarity in Reporting BLEU Scores. In *WMT*, volume 1, 186–191.
- Qu, X.; Ding, C.; Li, X.; Zhong, X.; and Tao, D. 2022. Distillation Using Oracle Queries for Transformer-Based Human-Object Interaction Detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 19558–19567.
- Radford, A.; Kim, J. W.; Hallacy, C.; Ramesh, A.; Goh, G.; Agarwal, S.; Sastry, G.; Askell, A.; Mishkin, P.; Clark, J.; et al. 2021. Learning transferable visual models from natural language supervision. *arXiv preprint arXiv:2103.00020*.
- Radford, A.; Narasimhan, K.; Salimans, T.; and Sutskever, I. 2018. Improving Language Understanding by Generative Pre-Training. Arxiv.

- Raffel, C.; Shazeer, N.; Roberts, A.; Lee, K.; Narang, S.; Matena, M.; Zhou, Y.; Li, W.; and Liu, P. J. 2020. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *JMLR*, 21: 1–67.
- Rebuffi, S.-A.; Bilen, H.; and Vedaldi, A. 2017. Learning multiple visual domains with residual adapters. In *NIPS*.
- Rebuffi, S.-A.; Bilen, H.; and Vedaldi, A. 2018. Efficient Parametrization of Multi-domain Deep Neural Networks. 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, 8119–8127.
- Ren, S.; He, K.; Girshick, R.; and Sun, J. 2015. Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. In *NIPS*.
- Sanh, V.; Debut, L.; Chaumond, J.; and Wolf, T. 2019. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*.
- Su, W.; Zhu, X.; Cao, Y.; Li, B.; Lu, L.; Wei, F.; and Dai, J. 2019. VL-BERT: Pre-training of Generic Visual-Linguistic Representations. In *ICLR*.
- Sun, C.; Myers, A.; Vondrick, C.; Murphy, K.; and Schmid, C. 2019. Videobert: A joint model for video and language representation learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 7464–7473.
- Sun, Z.; Yu, H.; Song, X.; Liu, R.; Yang, Y.; and Zhou, D. 2020. MobileBERT: a Compact Task-Agnostic BERT for Resource-Limited Devices. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2158–2170. Online: Association for Computational Linguistics.
- Sung, Y.-L.; Cho, J.; and Bansal, M. 2022. Vl-adapter: Parameter-efficient transfer learning for vision-and-language tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 5227–5237.
- Tung, F.; and Mori, G. 2019. Similarity-Preserving Knowledge Distillation. In 2019 IEEE/CVF International Conference on Computer Vision (ICCV), 1365–1374.
- Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, L.; and Polosukhin, I. 2017. Attention Is All You Need. In *NIPS*.
- Vedantam, R.; Zitnick, C. L.; and Parikh, D. 2015. CIDEr: Consensus-based Image Description Evaluation. In *CVPR*.
- Wang, A.; Singh, A.; Michael, J.; Hill, F.; Levy, O.; and Bowman, S. R. 2018. Glue: A multi-task benchmark and analysis platform for natural language understanding. In *ICLR*.
- Wang, J.; Hu, X.; Zhang, P.; Li, X.; Wang, L.; Zhang, L.; Gao, J.; and Liu, Z. 2020a. MiniVLM: A Smaller and Faster Vision-Language Model. Arxiv.
- Wang, P.; Yang, A.; Men, R.; Lin, J.; Bai, S.; Li, Z.; Ma, J.; Zhou, C.; Zhou, J.; and Yang, H. 2022a. Ofa: Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework. In *International Conference on Machine Learning*, 23318–23340. PMLR.
- Wang, W.; Bao, H.; Huang, S.; Dong, L.; and Wei, F. 2021. MiniLMv2: Multi-Head Self-Attention Relation Distillation for Compressing Pretrained Transformers. In *ACL Findings*.

- Wang, W.; Wei, F.; Dong, L.; Bao, H.; Yang, N.; and Zhou, M. 2020b. Minilm: Deep self-attention distillation for task-agnostic compression of pre-trained transformers. *Advances in Neural Information Processing Systems (NeurIPS)*.
- Wang, Z.; Codella, N.; Chen, Y.-C.; Zhou, L.; Dai, X.; Xiao, B.; Yang, J.; You, H.; Chang, K.-W.; Chang, S.-f.; et al. 2022b. Multimodal Adaptive Distillation for Leveraging Unimodal Encoders for Vision-Language Tasks. *arXiv* preprint arXiv:2204.10496.
- Williams, A.; Nangia, N.; and Bowman, S. R. 2017. A broad-coverage challenge corpus for sentence understanding through inference. In *NAACL*.
- Wu, H.; Gao, Y.; Zhang, Y.; Lin, S.; Xie, Y.; Sun, X.; and Li, K. 2022a. Self-supervised Models are Good Teaching Assistants for Vision Transformers. In *ICML*.
- Wu, K.; Zhang, J.; Peng, H.; Liu, M.; Xiao, B.; Fu, J.; and Yuan, L. 2022b. Tinyvit: Fast pretraining distillation for small vision transformers. In *ECCV*, 68–85. Springer.
- Xu, J.; Mei, T.; Yao, T.; and Rui, Y. 2016. Msr-vtt: A large video description dataset for bridging video and language. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- Yang, Z.; Dai, Z.; Yang, Y.; Carbonell, J.; Salakhutdinov, R. R.; and Le, Q. V. 2019. Xlnet: Generalized autoregressive pretraining for language understanding. In *NeurIPS*.
- Yang, Z.; Li, Z.; Jiang, X.; Gong, Y.; Yuan, Z.; Zhao, D.; and Yuan, C. 2022. Focal and global knowledge distillation for detectors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 4643–4652.
- Yu, L.; Lin, Z.; Shen, X.; Yang, J.; Lu, X.; Bansal, M.; and Berg, T. L. 2018. MAttNet: Modular Attention Network for Referring Expression Comprehension. In *CVPR*.
- Yu, Y.; Kim, J.; and Kim, G. 2018. A joint sequence fusion model for video question answering and retrieval. In *ECCV*.
- Yuan, L.; Chen, D.; Chen, Y.-L.; Codella, N.; Dai, X.; Gao, J.; Hu, H.; Huang, X.; Li, B.; Li, C.; et al. 2021. Florence: A new foundation model for computer vision. *arXiv preprint arXiv:2111.11432*.
- Zagoruyko, S.; and Komodakis, N. 2017. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. *ICLR*.
- Zellers, R.; Bisk, Y.; Schwartz, R.; and Choi, Y. 2018. Swag: A large-scale adversarial dataset for grounded commonsense inference. In *EMNLP*.
- Zhang, P.; Li, X.; Hu, X.; Yang, J.; Zhang, L.; Wang, L.; Choi, Y.; and Gao, I. 2021. VinVL: Making Visual Representations Matter in Vision-Language Models. Arxiv.
- Zhou, L.; Palangi, H.; Zhang, L.; Hu, H.; Corso, J. J.; and Gao, J. 2020. Unified Vision-Language Pre-Training for Image Captioning and VQA. In *AAAI*.
- Zhou, L.; Xu, C.; and Corso, J. J. 2018. Towards automatic learning of procedures from web instructional videos. In *AAAI*.
- Zhu, L.; and Yang, Y. 2020. ActBERT: Learning Global-Local Video-Text Representations. In *Proceedings of the*

- IEEE/CVF Conference on Computer Vision and Pattern Recognition.
- Zhu, Y.; Groth, O.; Bernstein, M.; and Fei-Fei, L. 2016. Visual7W: Grounded Question Answering in Images. In *CVPR*. ISBN 978-1-4673-8851-1.