

Bilaterally Slimmable Transformer for Elastic and Efficient Visual Question Answering

Zhou Yu, *Member, IEEE*, Zitian Jin, Jun Yu, *Member, IEEE*, Mingliang Xu, Hongbo Wang, Jianping Fan

Abstract—Recent advances in Transformer architectures [1] have brought remarkable improvements to visual question answering (VQA). Nevertheless, Transformer-based VQA models are usually deep and wide to guarantee good performance, so they can only run on powerful GPU servers and cannot run on capacity-restricted platforms such as mobile phones. Therefore, it is desirable to learn an elastic VQA model that supports adaptive pruning at runtime to meet the efficiency constraints of different platforms. To this end, we present the bilaterally slimmable Transformer (BST), a general framework that can be seamlessly integrated into arbitrary Transformer-based VQA models to train a single model once and obtain various slimmed submodels of different widths and depths. To verify the effectiveness and generality of this method, we integrate the proposed BST framework with three typical Transformer-based VQA approaches, namely MCAN [2], UNITER [3], and CLIP-ViL [4], and conduct extensive experiments on two commonly-used benchmark datasets. In particular, one slimmed MCAN_{BST} submodel achieves comparable accuracy on VQA-v2, while being 0.38× smaller in model size and having 0.27× fewer FLOPs than the reference MCAN model. The smallest MCAN_{BST} submodel only has 9M parameters and 0.16G FLOPs during inference, making it possible to deploy it on a mobile device with less than 60 ms latency.

Index Terms—Visual question answering, Slimmable network, Transformer, Multimodal learning, Efficient deep learning.

I. INTRODUCTION

Thanks to recent progress on deep neural networks, machines are able to address complicated multimodal tasks that require a fine-grained understanding of both vision and language cues, such as image-text matching [5][6], visual captioning [7], visual grounding [8][9], and visual question answering (VQA) [2][10]. Among these tasks, VQA is challenging because it requires performing visual reasoning over multimodal data to predict an accurate answer.

Current state-of-the-art VQA approaches can be roughly categorized into two lines of research based on whether they are trained from scratch (*e.g.*, MCAN [2] and MUAN [11])

This work was supported in part by the Zhejiang Provincial Natural Science Foundation of China under Grant LR22F020001, in part by the National Natural Science Foundation of China under Grants 62125201, 62072147, 62020106007 and 61836002, and in part by the Zhejiang Provincial Natural Science Foundation of China under Grant DT23F020007. (Corresponding authors: Jun Yu and Hongbo Wang.)

Z. Yu, J. Yu, and H. Wang are with School of Computer Science and Technology, Hangzhou Dianzi University, Hangzhou, 310018, China (e-mail: yuz@hdu.edu.cn; yujun@hdu.edu.cn; whongbo@hdu.edu.cn).

Z. Jin is with HDU-ITMO Joint Institute, Hangzhou Dianzi University, 310018, China (e-mail: jinzt@hdu.edu.cn).

M. Xu is with School of Computer and Artificial Intelligence Zhengzhou University, 450001, China (e-mail: iexumingliang@zzu.edu.cn).

J. Fan is with AI Lab at Lenovo Research, 100094, China (e-mail: jfan1@Lenovo.com).

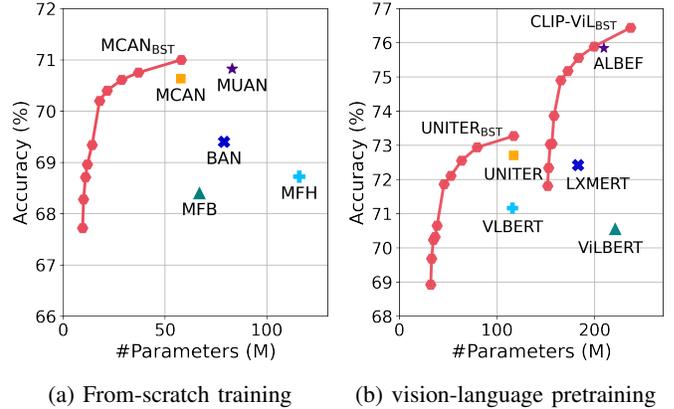


Fig. 1: Accuracy vs. number of model parameters on the VQA-v2 test-dev split to compare with the state-of-the-art methods. The methods are split into two classes depending on whether they are (a) trained from scratch [14][15][11][2] or (b) pretrained with external data [16][17][12][13][3]. By integrating the BST framework with three typical VQA models, namely MCAN [2], UNITER [3], and CLIP-ViL [4], the resulted slimmable MCAN_{BST}, UNITER_{BST} and CLIP-ViL_{BST} models (marked with red lines) either outperform their counterparts with similar model sizes or achieve comparable performance with smaller model sizes, showing the *efficiency* of our framework. Each red hexagon in the line represents a submodel sharing a portion of parameters of its full model (the rightmost one), showing the *elasticity* of our framework.

in Fig. 1a) or pretrained with external multimodal data (*e.g.*, LXMERT [12], UNITER [3], CLIP-ViL [4], and ALBEF [13] in Fig. 1b). Although these two lines of research use different training paradigms, they share the same model architecture, *i.e.*, Transformer [1], which was initially proposed for language modeling and has since become the foundational architecture for the VQA task.

Despite the effectiveness of the Transformer-based VQA methods described above, they typically require large models (*e.g.*, ~200M) to guarantee good performance. This severely limits their applicability in capacity-constrained platforms with specific efficiency constraints (*e.g.*, FLOPs and model size). To address this issue, a series of model compression strategies have been investigated to learn *efficient* Transformer architectures, including low-rank decomposition [18], weight sharing [19][20], model pruning [21][22], and knowledge distillation [23][24][25][26][27]. However, these approaches focus on only one specific scenario and obtain one compact

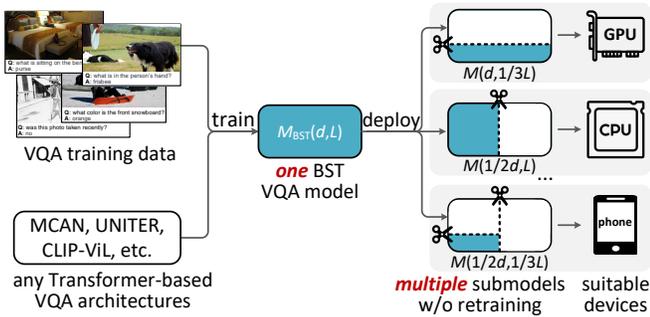


Fig. 2: The schematic diagram of the proposed Bilaterally Slimmable Transformer (BST) framework. The BST framework is general enough to be seamlessly applied to arbitrary Transformer-based VQA models that supports training a *single* BST-based VQA model only once and then pruning it to obtain *multiple* efficient submodels to meet the requirements of different platforms at inference time.

model. In practice, there are a wide variety of hardware platforms, *e.g.*, GPUs, CPUs, and mobile devices. To meet the efficiency requirements of different platforms, compression methods need to redesign the model architectures and then retrain the models, which is both engineer-expensive and computation-expensive. This motivates us to devise an *elastic-and-efficient* framework that supports training a single Transformer-based VQA model, and then adaptively pruning the model to fit different platforms *without retraining*. Compared with the aforementioned model compression approaches, the introduced framework has the following two advantages: 1) it reduces the model design costs as the submodels of different sizes can naturally meet the requirements of different platforms; 2) it also reduces the training cost as the model is trained once and adaptive slimming can be performed at inference time.

Inspired by the slimmable neural network that trains an elastic CNN model with multiple width multipliers to fit different efficiency constraints at runtime [28], we present a bilaterally slimmable Transformer (BST) framework to support model slimming in both the width and depth directions of the Transformer, where the width is the hidden dimensionality of each layer and the depth is the number of layers. As shown in Fig. 2, the BST framework can generally be seamlessly integrated with arbitrary Transformer-based VQA models to support training a single model only once and then pruning it to obtain multiple efficient submodels of various widths and depths at inference time. It is worth noting that each resulting submodel directly inherits a specific portion of the parameters of the full model and does not require further model finetuning. We take three typical Transformer-based VQA approaches, *i.e.*, MCAN [2], UNITER [3], and CLIP-ViL [4], as the reference models to incorporate into the BST framework, resulting in the slimmable MCAN_{BST}, UNITER_{BST}, and CLIP-ViL_{BST} models, respectively. As shown in Fig. 1, the resulting slimmable models either outperform their existing state-of-the-art counterparts at similar model sizes or achieve comparable performance at much smaller model sizes.

To the best of our knowledge, our study is the first attempt to explore efficient and elastic models for VQA. The most closely related studies to our work is the DynaBERT approach [29] and the RWSAN approach [30]. DynaBERT also investigates slimmable Transformer architectures. In contrast to our BST framework, which supports arbitrary Transformer-based architectures for VQA, DynaBERT focuses on pretrained BERT model for NLP tasks. In terms of methodology, our BST is different from DynaBERT in terms of the slimming strategy and training strategy, resulting in better model performance and less training time. RWSAN investigates lightweight VQA models by introducing residual weight-sharing attention (RWSA) layers, resulting in a VQA model with many fewer parameters. In contrast to our BST framework, RWSAN cannot reduce computational costs and fit the efficiency constraints of different platforms adaptively.

Our main contributions are summarized as follows:

- Regarding the motivation, orthogonal to the pursuit of model accuracy, we introduce a new direction to the VQA research to learn an *efficient and elastic* model once and obtain various efficient submodels that can be adaptively fit different platforms.
- Regarding the methodology, we present a general bilaterally slimmable Transformer (BST) framework which can transform any Transformer-based VQA model into a slimmable model to adjust the width and depth at runtime. Compared with DynaBERT [29] which introduces a slimmable BERT model for NLP tasks, our BST framework differs in terms of slimming strategies and training strategy.
- Regarding effectiveness and generality, we integrate the BST framework with three typical VQA models and conduct extensive experiments on two commonly used VQA datasets. The results show that the slimmed submodels either outperform the state-of-the-art approaches with similar model sizes or achieve comparable performance with smaller model sizes.
- Regarding practicability, this study is the first attempt to explore efficient VQA models on different hardware platforms including mobile devices. In particular, our smallest submodel can run on a non-latest mobile device with less than 60 ms latency, showing its potential in a wide range of applications such as robotics, automatic driving, and assistance for visually impaired people.

II. RELATED WORK

In this section, we first briefly review previous studies on VQA, especially the approaches with Transformer architectures. After that, we discuss related work on efficient neural networks and slimmable neural networks.

Visual Question Answering (VQA). The VQA task, which aims to answer a free-form question in natural language with respect to a given image, has attracted increasing interest over the last few years. The core of VQA lies in two lines of research, namely multimodal fusion and attention learning. For multimodal fusion, early methods used linear models with elementwise summation or multiplication to fuse features

from different modalities [31][32]. To better characterize the second-order interactions between multimodal features, Fukui *et al.* [33], Kim *et al.* [34], and Yu *et al.* [35] devised different multimodal bilinear pooling models. For attention learning, question-guided visual attention over image regions became a standard component of many early VQA approaches [36][37]. Yang *et al.* proposed a stacked attention network to iteratively learn visual attention on different levels [36]. More recently, co-attention models that consider both textual and visual attention have been proposed. Lu *et al.* introduced a hierarchical co-attention learning paradigm to learn image attention and question attention iteratively [38]. Yu *et al.* decoupled the co-attention learning into a question self-attention stage and a question-conditioned visual attention stage and optimized the two stages in an end-to-end manner [35]. The aforementioned co-attention models are *coarse-grained* in that they neglect the multimodal interactions at a fine-grained level (*i.e.*, word-region pairs). To this end, Nguyen *et al.* [39], Kim *et al.* [15], and Liu *et al.* [40] introduced dense co-attention models that establish dense interactions among word-region pairs.

Transformer-Based VQA. The Transformer architecture is initially proposed for machine translation in the NLP community [1]. It consists of a sequence of self-attention modules to model complex and dense interactions within a group of input features. This architecture is general enough to be used in various unimodal tasks [41][42][43] and multimodal tasks in image [44][2] and video domains [45][46][47]. For the VQA task, Gao *et al.* [48] and Yu *et al.* [2][5] devised models based on the Transformer architecture and deliver new state-of-the-art performance on VQA benchmark datasets at the time of their publication. Besides these studies on general-purpose VQA, there is also a growing trend towards exploring more granular VQA tasks with specific reasoning skills, *e.g.*, scene-text understanding [49][50], casual reasoning [51][52], knowledge utilization [53][54].

More recently, a BERT model that integrates the Transformer architecture with a self-supervised pretraining paradigm has shown great success in a wide range of NLP tasks. Mirroring the success of BERT, recent studies have naturally extended its framework to the multimodal domain to perform vision-language pretraining (VLP) to learn generic multimodal representations [38][12][3][55]. In particular, they first pretrain Transformer-based models on large image-text corpora to learn task-agnostic representations, and then fine-tune the models on downstream tasks such as VQA. Early VLP approaches designed different pretraining tasks to learn multimodal Transformers on top of preextracted region-based visual features [17][3][55][56][57]. Motivated by the success of pre-trained visual backbones, *e.g.*, ViT [43] and CLIP [58], recent VLP methods tend to exploit these visual backbones to obtain grid-based visual features and perform multimodal pretraining from raw image and language inputs in an end-to-end manner [59][13][60][4][61]. To summarize, Transformer-based approaches dominate the VQA task at present, due to their excellent capability for modeling the complex interactions among multimodal input features. However, Transformer-based VQA models are usually computationally expensive (*i.e.*, have a large number of parameters and FLOPs), hindering

their deployment on mobile devices with limited memory and computation consumption. This motivates us to explore efficient Transformer architectures for VQA.

Efficient Neural Networks. There has been broad interest in building efficient neural networks in the literature. Existing approaches can be generally categorized as either compressing pretrained networks [62][63][64] or training efficient networks directly [65][66][67]. The efficient neural networks above mainly focus on ConvNet architectures. Due to the popularity of Transformer in recent years, efficient Transformer architectures have been investigated in different respects, *e.g.*, low-rank decomposition [18], weight sharing [19][20], model pruning [21][22], and knowledge distillation [23][24][25][26][27]. Low-rank decomposition methods decompose a full-rank parameter matrix into low-dimensional matrices while weight sharing approaches reuse one parameter matrix in different layers of the network. Model pruning methods aim to cut out redundant parameters in the model to obtain a smaller model, and knowledge distillation techniques aim to transfer knowledge from a large teacher model to a small student model. Despite the success of these approaches, their efficient models are dedicated to one specific scenario, and cannot adapt to different efficiency constraints or different hardware platforms at runtime.

Slimmable Neural Networks. Orthogonal to the approaches to efficient neural networks above, slimmable neural networks aim to design *dynamic* models that can adaptively fit different efficiency constraints at runtime. Given a deep neural network, network slimming can be performed on both the depth and width dimensions. For depth slimming, the methods of Wu *et al.* [68], Liu *et al.* [64], and Huang *et al.* [69] learn controllers or gating modules to adaptively drop layers from deep ConvNets. For width slimming, Yu *et al.* introduced a general framework for a family of ConvNets (*e.g.*, ResNet [70] or MobileNet [66]) that supports a predefined set of width multipliers [28]. After that, they further improved the framework to support model slimming with arbitrary widths [71]. To take a further step, Cai *et al.* introduced a once-for-all (OFA) method to support width and depth slimming simultaneously in a unified framework [72]. All of the methods above are only for ConvNet architectures, and their strategies cannot be directly applied to Transformer architectures.

The most closely related study to our work is DynaBERT [29], which also investigates slimmable Transformer architectures. In contrast to our BST framework, which supports arbitrary Transformer-based architectures for VQA, DynaBERT focuses on pretrained BERT model for NLP tasks. Regarding the methodology, our BST is different from DynaBERT in terms of slimming strategy and training strategy, obtaining significant advantages in terms of a higher compression ratio and less training time.

III. BILATERALLY SLIMMABLE TRANSFORMER (BST)

In this section, we describe the bilaterally slimmable Transformer (BST) framework in detail. Before presenting the BST framework, we first revisit the core components of the Transformer architecture [1]. Then, we introduce the BST framework, including the slimming strategies for width and depth.

We take three typical VQA models, MCAN [2], UNITER [3], and CLIP-ViL[4] as examples and integrate them with the proposed BST framework. Without loss of generality, our BST framework can be applied to arbitrary VQA models of Transformer-based architectures. After that, we introduce the BST training strategy which consists of submodel architecture selection strategy and a knowledge distillation training stage. Finally, we provide in-depth comparisons to DynaBERT [29] in terms of methodology.

A. Preliminaries

Transformer is a multi-layer network in which each layer consists of the multi-head attention (MHA) and feed-forward network (FFN) modules [1].

Multi-Head Attention. Denote m query features and n key-value paired features as $Q \in \mathbb{R}^{m \times D}$, $K \in \mathbb{R}^{n \times D}$, and $V \in \mathbb{R}^{n \times D}$, where D is the hidden dimensionality of these features. The multi-head attention module calculates the attended features $F \in \mathbb{R}^{m \times D}$ by using H paralleled attention functions as follows:

$$F = \text{MHA}(Q, K, V) = [\text{head}_1, \text{head}_2, \dots, \text{head}_H]W^o$$

$$\text{head}_j = \text{ATT}(QW_j^Q, KW_j^K, VW_j^V) \quad (1)$$

where $W_j^Q, W_j^K, W_j^V \in \mathbb{R}^{D \times D_H}$ are the projection matrices for the j -th head, and D_H is the dimensionality of the features from each head. $W^o \in \mathbb{R}^{H \cdot D_H \times D}$ is the projection matrix used to aggregate the output features from different heads. We set $D_H = D/H$ so that the model sizes remain constant as H varies. The attention function for each head is defined as the scaled dot-product of the query with all keys:

$$\text{ATT}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{D_H}}\right)V \quad (2)$$

which calculates the scaled dot-product of each query with all keys to obtain the attention weights, and then performs weighted summation over the values.

Feed-Forward Network. The feed-forward network module is a two-layer MLP model applied to the output features of the MHA module to perform a pointwise nonlinear transformation. Given input features $X \in \mathbb{R}^{n \times D}$, the transformed features $F \in \mathbb{R}^{n \times D}$ are obtained as follows:

$$F = \text{FFN}(X) = \text{ReLU}(XW_1 + b_1)W_2^T + b_2 \quad (3)$$

where $W_1, W_2 \in \mathbb{R}^{D \times 4D}$.

Transformer Layer. A typical Transformer layer usually consists of a MHA module and an FFN module as follows:

$$F = \text{Transformer-layer}(X)$$

$$= \text{LN}(\text{FFN}(\hat{F}) + \hat{F}) \quad (4)$$

$$\hat{F} = \text{LN}(\text{MHA}(X, X, X) + X)$$

where residual connection [70] and layer normalization (LN) [73] are applied after the MHA and FFN modules. The LN module takes a D -dimensional feature $x \in X$ as its input and performs normalization as follows to obtain the output feature:

$$y = \text{LN}(x) = \gamma \frac{x - \mu}{\sqrt{\sigma^2 - \epsilon}} + \beta \quad (5)$$

where $\mu, \sigma^2 \in \mathbb{R}$ are the mean and variance calculated on x . $\gamma, \beta \in \mathbb{R}^D$ are the learnable parameters of the scale and shift terms, respectively.

Transformer Architectures. Depending on the composition strategies of the Transformer layers above, existing Transformer architectures can be categorized into three classes, namely encoders [41][74], decoders [75][76], and encoder-decoders [1][77].

Taking a sequence of input tokens, the original Transformer [1] adopts an encoder-decoder architecture. The encoder is composed of a cascade of Transformer layers in depth to obtain the bidirectional representations by jointly conditioning on both the left and right contexts, and the decoder takes the representations from the last encoder layer as input to guide the learning of unidirectional representations by conditioning only on the left context. After that, pure encoder architectures (e.g., BERT [41]) and pure decoder architectures (e.g., GPT [75]) are introduced and integrated with the self-supervised pretraining paradigm, which has been used in a wide range of NLP tasks.

B. The BST Framework

Let an L -layer Transformer with hidden dimensionality D be the reference model, where D and L denote the width and depth of the model, respectively. The goal of BST is to obtain a single *slimmable* Transformer model that can adaptively adjust to a set of submodels of different widths and depths in the inference stage. In the following, we introduce the width slimming and depth slimming strategies. An overview of the BST framework is shown in Fig. 3a.

Width Slimming. With width slimming, we aim to make each Transformer layer adapt to a set of width slimming ratios with respect to the hidden dimensionality d of the reference model. To achieve this goal, we split the parameters of the reference model into different submodels, with each sharing a specific portion of its model parameters. As shown in Eqs.(1), (3), and (5), the learnable parameters in the MHA and FFN modules are all linear projections, which can be simply split into submatrices with respect to the given ratios. Inspired by the settings in [29], we define the candidate width set as $\mathcal{D} = \{1/4D, 1/2D, 3/4D, D\}$.

Given a slimmed width $d \in \mathcal{D}$, we need to prune the model parameters in the MHA, FFN, and LN modules according to d . In the MHA module, the query, key, and value parameters over H heads can be represented as tiled matrices $W^Q, W^K, W^V \in \mathbb{R}^{D \times D_H * H}$. Given the slimmed width d , we keep D_H as a constant and adjust the input dimensionality D and the number of heads H accordingly. This results in the slimmed model parameters $\in \mathbb{R}^{d \times D_H * \hat{H}}$, where $\hat{H} = H * d/D$ refers to the reduced number of heads with the last few heads neglected. The model parameters W^o, W_1, W_2, γ , and β in the FFN and LN modules are adjusted in a similar manner by slimming the dimensionality of the input and output features to d .

In addition to the strategy introduced above, another width slimming strategy was investigated in [29]. In contrast to our *slim-all* strategy that slims the dimensionality of the input, output, and intermediate representations simultaneously, they

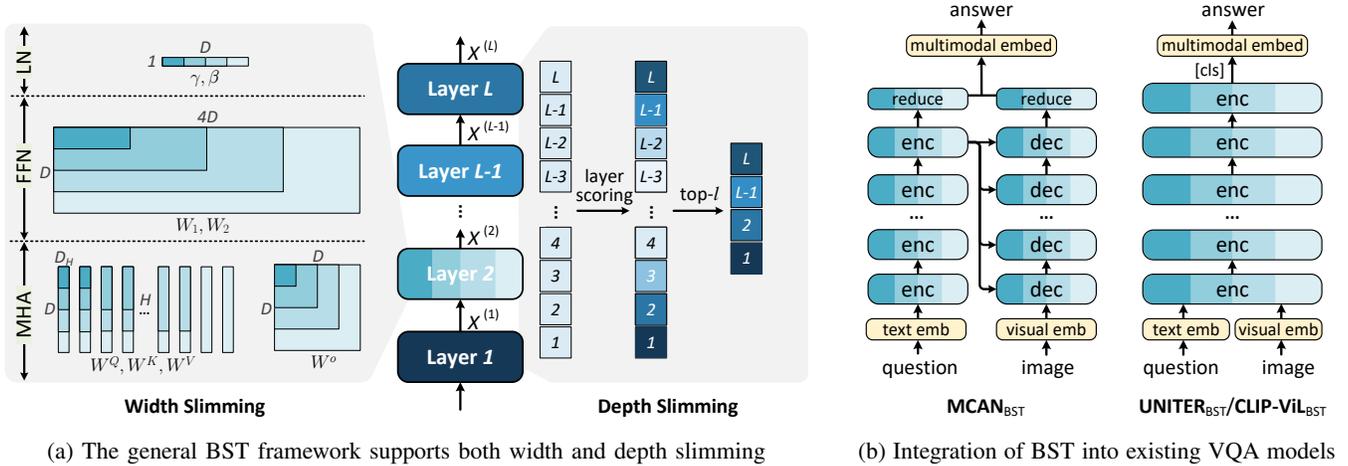


Fig. 3: (a) The overall diagram of the proposed Bilaterally Slimmable Transformer (BST) framework, which consists of width slimming and depth slimming. (b) The integration of the BST framework and three typical Transformer-based VQA models, namely MCAN [2], UNITER [3], and CLIP-ViL [4], respectively. The parameters in the input and output embedding layers (yellow background) are not slimmable as dimensionality of the input and output features remains the same.

introduced an alternative *slim-intermediate* strategy that slims the dimensionality of the intermediate representation while keeping the rest unchanged. The *slim-intermediate* strategy results in a *bottleneck* structure with a low dimensional intermediate representation. This may break the carefully-designed structure of the Transformer (e.g., the $D \rightarrow 4D \rightarrow D$ structure in FFN).

Depth Slimming. With depth slimming, we aim to make the slimmed submodels adapt to a set of depth slimming ratios with respect to the maximum depth of the reference model. As the depth of a typical Transformer model is usually set to a multiple of 6, we define the candidate depth set as $\mathcal{L} = \{1/6L, 1/3L, 2/3L, L\}$. Compared with the settings in [29], we explore the submodels with much shallow depth, resulting in more compact submodels suitable for mobile devices.

To perform depth slimming, we first need to determine which layers are to be slimmed given a specific slimming depth $l \in \mathcal{L}$. As shown on the right side of Fig. 3a, we first assign an importance score to each layer using certain scoring strategies. After that, we select the top- l layers with the largest scores and keep their original order. Here, we introduce three scoring strategies, which result in three types of slimming strategies: 1) the *slim-random* strategy is the most straightforward strategy, as it simply sets the importance scores to random values; 2) the *slim-first* (or *slim-last*) strategy sets the importance scores in ascending (or descending) order; 3) the *slim-middle* strategy sets the smallest scores to the middlemost layer and gradually increases the scores as it move toward the top and bottom layers. This strategy was inspired by the empirical studies in [78], in which the layers closer to the input and output are more important than the middle layers in the Transformer. We use the *slim-middle* strategy as the default option.

C. Integrating BST with Off-the-Shelf VQA Models

BST is a general framework that can be integrated with arbitrary Transformer-based VQA models in theory. In this paper, we choose three typical Transformer-based models, i.e.,

MCAN [2], UNITER [3] and CLIP-ViL[4], shown in Fig. 3b, to integrate with the proposed BST framework. Without loss of generality, the BST framework can also be applied to other Transformer-based models beyond the VQA task. Due to space limitations, we will not expand the description further.

MCAN_{BST}. MCAN was the winning solution in the VQA Challenge 2019. It introduces an encoder-decoder-based Transformer architecture to model complex multimodal interactions and perform accurate visual reasoning. Specifically, the input question is encoded as a sequence of word embeddings using pretrained GloVe embeddings followed by an LSTM network. The input image is encoded as a group of object embeddings using a pretrained object detector [7]. After that, the multimodal embeddings are passed through an L -layer encoder-decoder to obtain the attended output features. In the L -layer question encoder, the word embeddings are transformed with a self-attention mechanism to obtain the attended question features of the same word length. The attended word features and visual embeddings are further fed into an L -layer image decoder to obtain the attended image features with a guided-attention mechanism. On top of the attended question features and image features, two attentional reduction modules are devised to obtain a question feature and an image feature, respectively. Finally, the two feature vectors are simply fused and then fed to a linear classifier to predict the answer. The MCAN model can be trained from scratch in an end-to-end manner on a specific VQA dataset such as VQA-v2 [79].

Since MCAN's core components are the standard Transformer layers, the model can be seamlessly integrated with the BST framework to obtain a slimmable MCAN_{BST} model. The width slimming strategy can be directly applied to each encoder and decoder layer in MCAN, and different depth slimming strategies can also be applied to drop a portion of the encoder and decoder layers simultaneously. Furthermore, the model parameters in the attention reduction module on top of the encoder-decoder are derived from two-layer MLPs, which can also be slimmed in width.

UNITER_{BST}. UNITER is a representative vision-and-language

pretraining (VLP) approach with an L -layer Transformer encoder as its backbone. In contrast to MCAN’s *training-from-scratch* mechanism, UNITER utilizes a *vision-language pre-training* strategy to learn a generalized backbone model from massive image-text pairs, and then finetunes the backbone to adapt to different multimodal tasks. Specifically, for the VQA task, a task-specific head is appended on top of the backbone so that the representation of the predefined [cls] token is fed to a linear classifier to predict the answer. Based on the finetuned UNITER model for VQA, both width slimming and depth slimming are applied to its backbone to transform it into UNITER_{BST}.

CLIP-ViL_{BST}. CLIP-ViL shares the same Transformer architecture with UNITER but introduces a more powerful visual encoder to extract visual representations. Specifically, its visual encoder corresponds to a pretrained ResNet-50×4 model, which is pretrained on 400M image-text pairs by CLIP [58].

Note that the embedding layers in the above models are not slimmable, making the input dimensionality of the first Transformer layer unadjustable. This contradicts our width-slimming strategy. To address this issue, we insert a slimmable linear layer $W_{\text{emb}} \in \mathbb{R}^{D \times D}$ between the embedder and backbone, and make it slimmable in width to adapt to the BST framework¹.

Submodel Complexity Analysis. Given a reference MCAN (UNITER or CLIP-ViL) model of width D and depth L , its model size and FLOPs are both approximately proportional to $O(D^2L)$. This indicates that the computational cost of the smallest submodel $a(1/4D, 1/6L)$ can be up to $96 \times$ smaller than that of the reference model. In practice, the scaling ratio between the submodels and the reference model is not that large. As shown in Fig. 3b, the slimming strategies are only performed in the backbone while the embedders and the classifier are not involved. Their existence introduces an inescapable cost for all the slimmed submodels, limiting the computational overhead of the small submodels. More detailed results are given and analyzed in section IV.

D. Training Strategy for BST Models

The training procedures of BST consists of a *submodel architecture selection* stage and a *knowledge distillation training* stage.

Submodel Architecture Selection. By combining each width in \mathcal{D} with each depth in \mathcal{L} , we obtain a set of submodel architectures \mathcal{A} of different widths and depths as follows:

$$\mathcal{A} = \text{combination}(\mathcal{D}, \mathcal{L}) \quad (6)$$

where $|\mathcal{A}| = |\mathcal{D}| * |\mathcal{L}|$. Each architecture $a(d, l) \in \mathcal{A}$ corresponds to a combination of a specific width $d \in \mathcal{D}$ and depth $l \in \mathcal{L}$. In contrast to previous works that maintain all possible submodel architectures [80][29], we hypothesize that not every submodel architecture is effective. By the effectiveness of a submodel architecture, we mean the extent to which its computational cost (e.g., in terms of FLOPs or

¹For the UNITER_{BST} and CLIP-ViL_{BST} models with pretrained model parameters, W_{emb} is initialized with an identity matrix and updated with the entire model in an end-to-end manner.

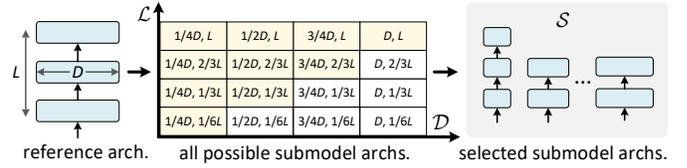


Fig. 4: The diagram of submodel architecture selection by using the triangle selection filtering strategy. The selected submodel architectures (highlighted with yellow background) follow the *deep-and-narrow* principle, which is considered to be more effective than the rest (white background).

model size) matches its delivered performance after model training. Therefore, devising a heuristic strategy to eliminate such ineffective architectures before BST training can reduce the training costs while improving the performance of the remaining submodels.

According to previous studies on designing efficient Transformers [81][29], *deep and narrow* architectures usually deliver better performance than *shallow and wide* architectures under constrained computational costs. This principle can be explained in two ways: 1) The Transformer requires a relatively deep model to guarantee good performance; and 2) the computational cost of a Transformer model is proportional to $O(LD^2)$, suggesting that increasing depth is more economical than increasing width.

To quantize this deep-and-narrow principle, we introduce a simple *triangle selection strategy* to filter out the shallow and wide submodel architectures. As shown in Fig. 4, we introduce a 2-D indicator matrix $I \in \{0, 1\}^{|\mathcal{D}| \times |\mathcal{L}|}$ to track the selection status for all the submodel architectures \mathcal{A} . $I(d, l) = 1$ indicates that the submodel architecture $a(d, l)$ is selected, and it is 0 otherwise. The indicator matrix I is first initialized with all-one values and then converted to an *upper-triangle* matrix. This strategy allows us to obtain *six* shallow and wide submodel architectures, which correspond to the matrix elements above the main diagonal. The submodel architectures \mathcal{S} selected from \mathcal{A} are defined as follows:

$$\mathcal{S} = \text{upper-triangle}(\mathcal{A}) = \{a(d, l) | I(d, l) = 1\} \quad (7)$$

Knowledge Distillation Training. After obtaining the selected submodel architectures \mathcal{S} , we train a slimmable model M_{BST} that can elastically adapt to any submodel architecture $a \in \mathcal{S}$, where M refers to any of the above VQA models. Given the BST model M_{BST} and a submodel architecture a , the obtained submodel is denoted as $M_{\text{BST}}^{(a)}$.

To obtain the BST model, we introduce a knowledge-distillation (KD) training mechanism as follows. In general, we first train an ordinary model M_{teacher} as the teacher model without model slimming and freeze its model parameters. After that, the BST model M_{BST} is initialized with the model parameters from M_{teacher} , and it can be viewed as the student model. Each slimmed submodel $M_{\text{BST}}^{(a)}$ shares a specific portion of the model parameters from M_{BST} and is trained with the supervision of the teacher model using the KD strategy [82]. Specifically, given a training sample, we feed it through the teacher model and each slimmed submodel simultaneously.

The predicted answer distribution from the teacher model is regarded as a soft label, and a proper loss function is imposed on the soft label and each submodel prediction to train the corresponding model².

Unlike the existing approaches [29] that decouple the width slimming and depth slimming into two training stages, we use a simpler one-stage training paradigm for BST with standard mini-batch SGD. In each iteration, we select K submodel architectures from \mathcal{S} and feed the same input samples to both the teacher model and the selected K submodels to obtain $(K + 1)$ predictions in total. After that, we apply the KD loss between the predictions of the teacher and each of the K student models and use the accumulated back-propagated gradients to update the model parameters of K . To stabilize the BST training, the K selected submodels consist of two determined submodels (the smallest one and the largest one) and another $(K-2)$ randomly sampled submodels. We use $K=4$ as the default setting in our experiments.

The whole BST training process is shown in Algorithm 1.

Algorithm 1: Training procedure for BST.

Input: A reference Transformer architecture with width D and depth L . Two predefined sets \mathcal{D} and \mathcal{L} define the slimming widths and depths w.r.t. D and L , resp. The number of sampled submodel architectures K per training iteration. a_s and a_l refer to the smallest and largest submodel architectures, respectively.

Output: An optimized BST model M_{BST} .

Stage I: Submodel Architectures Selection;

$\mathcal{A} = \text{combination}(\mathcal{D}, \mathcal{L});$

$\mathcal{S} = \text{upper-triangle}(\mathcal{A});$

Stage II: Knowledge Distillation Training;

Train M_{teacher} to obtain its optimized model parameters θ ;

Initialize the model parameters θ_{BST} for M_{BST} with θ ;

for $i = 1$ **to** max-iter **do**

 Randomly sample a mini-batch of data x ;

 Initialize the sampled architecture set $\Omega = \emptyset$;

 # add the smallest and largest submodel architectures.

$\Omega \leftarrow \Omega \cup \{a_s, a_l\};$

 # add another $K - 2$ architectures via random sampling.

for $j = 1$ **to** $K - 2$ **do**

 Randomly sample a submodel architecture $a \sim \mathcal{S}_{\Omega}$;

$\Omega \leftarrow \Omega \cup a$

end

 # submodel training using knowledge distillation.

 Feed-forward the teacher model: $y = M_{\text{teacher}}(x)$;

 Freeze M_{teacher} by stopping gradients: $y.\text{detach}()$;

foreach $a \in \Omega$ **do**

 Feed-forward the submodel: $\hat{y} = M_{\text{BST}}^{(a)}(x)$;

 Compute loss: $\text{loss} = \text{KD}(y, \hat{y})$;

 Accumulate backward gradients: $\text{loss}.\text{backward}()$;

end

 Update model parameters θ_{BST} .

end

E. In-Depth Comparison of BST and DynaBERT

As mentioned above, our BST framework has close connections with DynaBERT [29]. We conduct an in-depth comparison and describe their differences in terms of methodology.

²Different loss functions (e.g., BCE and KL-divergence) can be used as the KD loss flexibly, depending on the loss function used in the teacher model.

In terms of the width-slimming strategy, DynaBERT adopts a *slim-intermediate* strategy that only reduces the dimensionality of the intermediate representation while keeping the input and output representations unchanged, which may break the carefully designed bottleneck structure in the original Transformer architecture. In contrast, we use a *slim-all* strategy that reduces all the dimensionalities, keeping the bottleneck structure to achieve better performance.

In terms of the depth-slimming strategy, DynaBERT uses a simple slimming strategy by slimming the layers uniformly. In contrast, BST introduces a *slim-middle* strategy that considers the layer importance and prefers to slim the middlemost layers first, which facilitates model performance.

In terms of the training strategy, DynaBERT uses a two-stage training paradigm in which the width and depth slimming are learned separately, with all submodels being updated in each training step. In contrast, BST uses a simpler one-stage training paradigm to learn width and depth slimming simultaneously and sample a small number of submodels in each step, which significantly improves the training efficiency. Moreover, DynaBERT maintains the submodel architectures of all the width-height combination during training which may include redundant ones. In contrast, BST additionally introduces a submodel selection strategy to remove ineffective submodel architectures before training. This strategy not only reduces the training costs but also improves the performance of the remaining submodels

To summarize, BST has advantages over DynaBERT in terms of training efficiency and model performance. More quantitative comparisons are provided in section IV.

IV. EXPERIMENTAL RESULTS

In this section, we present experiments to evaluate the performance of our BST framework on two benchmark VQA datasets, namely, VQA-v2 [32] and GQA [83]. As mentioned above, we integrate BST with three typical Transformer-based VQA models, namely, MCAN [2], UNITER [3], and CLIP-ViL [4], to demonstrate the effectiveness and generality of BST. Furthermore, we conduct extensive ablation experiments on VQA-v2 to explore the effects of different components.

A. Datasets

VQA-v2 is the most commonly used VQA dataset [79]. It contains human-annotated QA pairs for MS-COCO images [84]. The dataset is split into three sets: `train` (80k images with 444k questions); `val` (40k images with 214k questions); and `test` (80k images with 448k questions). The `test` set is further split into `test-dev` and `test-std` sets. The reported results include three per-type accuracies (yes/no, number, and other), as well as an overall accuracy.

To make a fair comparison among the compared models, we follow the dataset splitting strategy in UNITER that further splits the `val` set into a `minival` subset of 5k images and a `trainval` subset of the remaining 35k images [3]. All the results reported in the experiments come from the models that are trained on the augmented `train+trainval+vg` sets, where `vg` denotes the augmented VQA samples from

TABLE I: Comparison to the state-of-the-art approaches on VQA-v2. For a fair comparison, the compared methods are split into two groups depending on whether they are trained from scratch or pretrained on external data (separated by a double-line). The number of parameters is calculated from an entire model, including the backbone, input and output embedding layers. The number of FLOPs is calculated from one single sample. †: the UNITER_{DYN} model refers to another slimmable model that integrates UNITER and DynaBERT [29], which is reimplemented by ourselves.

#	model	#params	#FLOPs	test-dev				test-std			
				All	Y/N	Num	Other	All	Y/N	Num	Other
<i>From-scratch training with augmented Visual Genome data</i>											
1	UpDn [7]	22M	1.1G	65.32	81.82	44.21	56.05	65.67	82.20	43.90	56.26
2	MFB [35]	68M	2.4G	68.40	84.78	49.05	58.82	-	-	-	-
3	MFH [14]	102M	2.5G	68.76	84.27	49.56	59.89	-	-	-	-
4	BAN [15]	112M	12.3G	69.66	85.46	50.66	60.50	-	-	-	-
5	MUAN [11]	83M	17.3G	70.82	86.77	54.40	60.89	71.10	-	-	-
6	MCAN($D=512, L=6$) [2]	58M	5.5G	70.63	86.82	53.26	60.72	70.90	-	-	-
7	MCAN($1/4D, 1/3L$) [2]	10M	0.2G	67.19	83.38	49.75	57.31	-	-	-	-
8	RWSAN [30]	20M	6.5G	70.19	86.45	52.18	60.38	-	-	-	-
9	MCAN _{BST} (D, L)	58M	5.5G	71.05	87.39	52.96	61.19	71.28	87.36	52.77	61.52
10	MCAN _{BST} ($1/2D, L$)	22M	1.5G	70.45	86.84	52.89	60.43	-	-	-	-
11	MCAN _{BST} ($1/2D, 1/3L$)	14M	0.6G	69.42	85.68	51.96	59.48	-	-	-	-
12	MCAN _{BST} ($1/4D, 1/3L$)	10M	0.2G	68.16	84.84	50.27	57.95	-	-	-	-
<i>Vision-language pretraining with massive external data</i>											
13	ViLBERT [17]	221M	18.7G	70.55	-	-	-	70.92	-	-	-
14	VLBERT [16]	116M	20.7G	71.16	-	-	-	-	-	-	-
15	LXMERT [12]	183M	20.3G	72.42	-	-	-	72.54	88.20	54.20	63.10
16	OSCAR [55]	116M	38.6G	73.16	-	-	-	73.44	-	-	-
17	ALBEF [13]	210M	77.9G	75.84	-	-	-	76.05	91.67	55.43	67.19
18	UNITER($D=768, L=12$) [3]	117M	20.2G	72.70	88.86	55.10	62.87	72.95	89.00	55.37	63.01
19	UNITER($1/4D, 1/3L$) [3]	33M	0.8G	66.89	83.25	49.97	56.74	-	-	-	-
20	CLIP-ViL($D=768, L=12$) [4]	237M	82.1G	76.44	91.38	58.12	67.86	76.74	91.60	58.09	68.07
21	CLIP-ViL($1/4D, 1/3L$) [4]	153M	59.3G	70.32	85.40	52.12	61.59	-	-	-	-
22	†UNITER _{DYN} (D, L)	117M	20.2G	73.19	89.02	56.39	63.46	-	-	-	-
23	†UNITER _{DYN} ($3/4D, 1/2L$)	64M	8.1G	71.99	87.87	54.76	62.34	-	-	-	-
24	†UNITER _{DYN} ($1/4D, 1/2L$)	43M	3.1G	70.48	86.37	52.72	60.93	-	-	-	-
25	UNITER _{BST} (D, L)	117M	20.2G	73.27	89.01	56.73	63.57	73.46	89.17	56.28	63.73
26	UNITER _{BST} ($1/2D, L$)	53M	5.6G	72.11	87.83	55.59	62.42	-	-	-	-
27	UNITER _{BST} ($1/2D, 1/3L$)	39M	2.2G	70.65	86.47	53.47	61.02	-	-	-	-
28	UNITER _{BST} ($1/4D, 1/3L$)	33M	0.8G	69.68	85.49	52.26	60.12	-	-	-	-
29	CLIP-ViL _{BST} (D, L)	237M	82.1G	76.44	91.44	58.27	67.78	76.70	91.54	58.71	67.90
30	CLIP-ViL _{BST} ($1/2D, L$)	172M	64.9G	75.17	90.29	57.31	66.30	-	-	-	-
31	CLIP-ViL _{BST} ($1/2D, 1/3L$)	158M	60.8G	73.86	89.20	55.43	64.95	-	-	-	-
32	CLIP-ViL _{BST} ($1/4D, 1/3L$)	153M	59.3G	72.34	87.41	54.11	63.63	-	-	-	-

Visual Genome [85]. The obtained models are validated on the `minival` set offline, and evaluated on the `test-dev` and `test-std` sets online.

GQA is a challenging VQA dataset that requires more complex reasoning skills [83]. It consists of 113K images and 1.2M balanced question-answer pairs of assorted types and varying compositionality degrees, measuring performance on an array of reasoning skills such as object and attribute recognition, spatial reasoning, logical inference, and comparisons. The dataset is split into the following four sets: `train`

(72k images with 943k questions), `val` (10k images with 132k questions), `test-dev` (398 images with 12k questions), and undisclosed `test-challenge` (1.6k images with 50k questions). Following the suggestions in the official GQA guideline³, all the models are trained on the `train+val` sets and evaluated on the `test-dev` set.

³<https://cs.stanford.edu/people/dorarad/gqa/evaluate.html>

TABLE II: Accuracies of the state-of-the-art methods on GQA. All entries use the officially provided object features for images and are evaluated on the `test-dev` split.

model	#params	#FLOPs	accuracy
UpDn [7]	30M	2.8G	51.62
BAN [15]	120M	14.9G	55.81
MCAN($D=512, L=6$) [2]	59M	6.5G	56.64
MCAN _{BST} (D, L)	59M	6.5G	57.83
MCAN _{BST} ($1/2D, L$)	22M	1.9G	57.67
MCAN _{BST} ($1/2D, 1/3L$)	15M	0.8G	57.09
MCAN _{BST} ($1/4D, 1/3L$)	10M	0.4G	56.38

B. Experimental Setup

The model architectures and training hyperparameters of the MCAN_{BST}, UNITER_{BST}, and CLIP-ViL_{BST} models are the same as those in their original models [2][3][4], respectively.

For MCAN_{BST}, the hidden dimensionality D , number of heads H , and number of layers L are set to 512, 8, and 6, respectively. The MCAN_{BST} models are trained on the VQA-v2 and GQA datasets using slightly different settings. On VQA-v2, binary cross-entropy (BCE) is used as the loss function for both the teacher and the BST model, and both models are trained for up to 15 epochs with a batch size of 64 and a base learning rate of $1e-4$. The learning rate is warmed-up for 3 epochs and decays by $1/5$ every 2 epochs after 10 epochs. On GQA, the learning rate and batch size are the same as those on VQA-v2. KL-divergence is used as the loss function and the teacher and BST models are trained for 11 epochs. The learning rate is warmed-up for 2 epochs and decays by $1/5$ every 2 epochs after 8 epochs.

For UNITER_{BST} and CLIP-ViL_{BST}, we adopt the network architecture from the BERT-base model with $D = 768$, $H = 12$, and $L = 12$. Given a model the finetuned on VQA-v2 as the teacher, UNITER_{BST} and CLIP-ViL_{BST} are initialized from their corresponding teacher models and trained using the AdamW optimizer. UNITER_{BST} is trained for up to 130k iterations with a batch size of 5,120 and a base learning rate of $1.5e-4$. CLIP-ViL_{BST} is trained for up to 15 epochs with a batch size of 32 and a base learning rate of $2e-4$. To reduce the usage of GPU memory, we freeze the model parameters in the visual encoder during BST training.

C. Main Results

In Tables I and II, we compare MCAN_{BST}, UNITER_{BST}, and CLIP-ViL_{BST} to the state-of-the-art VQA methods on VQA-v2 and GQA, respectively. For MCAN_{BST}, the compared methods include UpDn [7], MFB [35], MFH [14], BAN [15], MUAN [11], and MCAN [2], which were the best-performing solutions in the VQA Challenge in recent years. In addition, we introduce the lightweight VQA model RWSAN [30] into the comparison. For UNITER_{BST} and CLIP-ViL_{BST}, the compared methods include ViLBERT [17], VLBERT [16], LXMERT [12], OSCAR [55], ALBEF [13], UNITER [3] and CLIP-ViL [4], which are representative vision-language pre-training methods. Moreover, although DynaBERT [29] is not

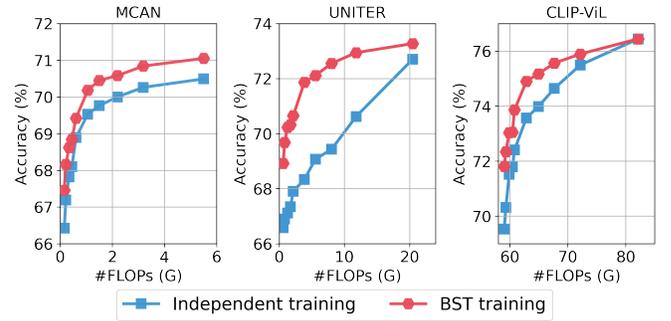


Fig. 5: Accuracy vs. #FLOPs on the VQA-v2 `test-dev` split. For each VQA model (*i.e.*, MCAN [2], UNITER [3], and CLIP-ViL [4]), we report the results of ten submodels obtained from BST training and independent training, respectively.

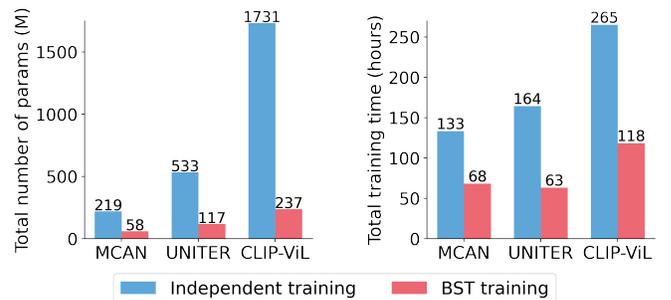


Fig. 6: Total number of parameters (left) and total training time (right) of the 10 submodels obtained by BST training and independent training, respectively. The training time is measured on a single Nvidia TitanV GPU.

specifically designed for VQA, we reimplement it ourselves to integrate it with UNITER, resulting in a slimmable VQA model UNITER_{DYN} for comparison. Due to space limitations, we do not show the results of all the submodels (*i.e.*, the ten submodels selected by the triangle selection strategy in Fig. 4) in these tables. Instead, four typical submodels are selected for comparison with the state-of-the-art approaches.

From the results in the upper part of Table I, we have the following observations: 1) With the same model architecture, the largest submodel MCAN_{BST}(D, L) and the smallest submodel MCAN_{BST}($1/4D, 1/3L$) outperform their independently-trained counterparts, respectively (#9 vs. #6, #12 vs. #7). This improvement is a benefit of the synergistic effect of the weight-sharing submodel architectures and KD training strategy. 2) With only $0.38\times$ the model size and $0.27\times$ the FLOPs, MCAN_{BST}($1/2D, L$) is still competitive with the reference MCAN model (#10 vs. #6), showing the potential of width slimming. In contrast to RWSAN [30], which has a similar model size, MCAN_{BST}($1/2D, L$) achieves higher accuracy and $0.25\times$ the FLOPs. 3) By slimming the depth to $1/3L$ (#11), its corresponding model size and FLOPs are respectively reduced to $0.6\times$ and $0.4\times$ those of its counterpart in #10, respectively, at the expense of a 1-point accuracy drop. Compared with MFB [35], MFH [14], and BAN [15], MCAN_{BST}($1/2D, 1/3L$) achieves superior or comparable performance with up to $0.125\times$ the model size and $0.05\times$ the

TABLE III: Runtime latency (ms) of three typical MCAN_{BST} submodels on three platforms, namely GPU, CPU, and mobile. For each platform, we report the results on three device.

platforms		GPU			CPU			mobile		
device type		RTX 3090	RTX 2080Ti	GTX 1650	Intel E5-2620v4	Intel I3-10100	AMD R7 4800H	Qualcomm Snapdragon 888	MediaTek Dimensity 1100	Qualcomm Snapdragon 660
1	(D, L)	22	29	40	51	101	127	167	361	439
2	$(1/4D, L)$	21	29	35	30	44	50	87	127	197
3	$(1/4D, 1/6L)$	5	7	12	11	14	22	58	93	160

FLOPs. 4) $\text{MCAN}_{\text{BST}}(1/4D, 1/3L)$ still outperforms UpDn [7] by 2.8 points with an extremely small model size of 10M. This model size is close to the lower bound of MCAN, which includes 7.8M uncompressible model parameters in the input and output embedding layers.

From the results in the lower part of Table I, we obtain similar observations to those on MCAN_{BST} . The slimmable $\text{UNITER}_{\text{BST}}$ and $\text{CLIP-ViL}_{\text{BST}}$ models attain superior or comparable performance to their reference independently-trained counterparts (#25 vs. #18, #28 vs. #19, #29 vs. #20, #32 vs. #21). The slimmed submodels in #26-28 and #30-32 attain significant reduction in computational costs at the expense of a drop in accuracy. Compared with the slimmable model $\text{UNITER}_{\text{DYB}}$ of DynaBERT [29] (#22-24), our $\text{UNITER}_{\text{BST}}$ achieves higher performance with similar model sizes (#22 vs. #25, #23 vs. #26, #24 vs. #27). Additionally, the total training time for DynaBERT is $3.6\times$ longer than ours. These results verify that our slimming and training strategies are more effective than those of DynaBERT.

To further examine the generalization of BST, we compare MCAN_{BST} to the state-of-the-art methods on GQA. Table II shows that $\text{MCAN}_{\text{BST}}(D, L)$ is 1.2 points higher than the reference MCAN model without BST. Furthermore, with $0.17\times$ the model size and $0.06\times$ the FLOPs, $\text{MCAN}_{\text{BST}}(1/4D, 1/3L)$ achieves comparable results to the reference MCAN model, surpassing the rest of its counterparts by a distinct margin.

Next, we show the results of all the ten submodels in terms of BST training and standard independent training. By independent training, we mean that each submodel is trained independently *without* sharing its model parameters⁴. From the results in Fig. 5, we can see that all the ten submodels obtained by BST training deliver better performance than their counterparts obtained by independent training. This corroborates the results in Table I.

Finally, the submodels obtained by BST are slimmed from *one single model without retraining*, outperforming the same submodels obtained by independent training in terms of both the total model size and total training time. From the results in Fig. 6, we can see that the total model size of the ten submodels obtained by BST training is $\sim 25\%$ of that obtained by independent training. Furthermore, the total training time for BST training is $\sim 50\%$ of that for independent training, revealing the synergistic effect of different weight-sharing submodels in BST training.

⁴For the independent training for UNITER and CLIP-ViL, each submodel is first initialized with a specific portion of the model parameters from the pretrained model and then finetuned independently.

To summarize, these observations above verify the effectiveness and generality of the proposed BST in terms of different Transformer architectures, training paradigms, and visual encoders.

D. Runtime Latency on Different Hardware Platforms

To precisely measure the efficiency of different submodels, we report the runtime latency on different hardware platforms in Table III. Specifically, we deploy three typical submodels of a MCAN_{BST} model (*i.e.*, the largest submodel, the smallest submodel, and a deep and narrow submodel) on three commonly-used platforms (GPU, CPU, and mobile devices). For each platform, we choose three devices with different computational capabilities.

From the results in Table III, we can see that: 1) For each submodel, different platforms and device types lead to significant discrepancies in terms of runtime latency, which stems from their diverse computational capabilities. 2) When width slimming is performed (#1 vs. #2), the latency on the GPU platform is not reduced significantly, while the latency on the CPU platform decreases prominently. This suggests that the GPU has no significant advantage over the CPU for these *narrow* models. 3) When depth slimming is also further performed (#1 vs. #3), the latencies on the GPU and CPU platforms are not distinct. This suggests that computational capabilities of the GPU and CPU platforms are excessive for such small submodels. 4) The smallest submodel in #3 can be deployed on a non-latest cellphone with a Snapdragon 888 chip. The 58 ms latency can support applications with real-time requirements.

E. Ablation Studies

We run a number of ablations on MCAN_{BST} to analyze the effectiveness of the key component in BST. The results are shown in Table IV and Fig. 7 and discussed in detail below.

Width Slimming. In Table IVa, we show the results of the MCAN_{BST} variants trained with different width slimming strategies, *i.e.*, the *slim-all* strategy introduced in this paper and the *slim-intermediate* strategy introduced in [29]. By comparing two submodels of similar FLOPs, we see that our *slim-all* strategy delivers better model performance than the *slim-intermediate* strategy under different model depths.

Depth Slimming. In Table IVb, we compare the MCAN_{BST} variants with different depth slimming strategies (mentioned in Section III-B) in terms of average accuracy over the ten submodels. From the results, we can see that the *slim-middle* strategy achieves the best performance among the counterparts,

TABLE IV: Ablations of the MCAN_{BST} variants models evaluated on the `minival` split of VQA-v2. The default setting and the best result are bolded.

slimming strategy	submodel	#FLOPs	accuracy (%)
slim-all	$(1/2D, L)$	1.5G	67.98
slim-interm. [29]	$(1/4D, L)$	1.6G	67.86
slim-all	$(1/2D, 1/3L)$	0.6G	66.95
slim-interm. [29]	$(1/4D, 1/3L)$	0.8G	66.83

(a) **Width Slimming.** Under the same model depth and similar number of FLOPs, the obtained submodels trained with the *slim-all* strategy outperform the *slim-intermediate* strategy, showing the significance of keeping the ratio of input-output dimensionality in width slimming.

#sampled submodels	average accuracy (%)	training time (h)
$K=3$	66.88	64
$K=4$	67.14	68
$K=5$	67.15	73
$K=6$	67.16	79

(c) **Number of Sampled Submodels.** $K=4$ is a good trade-off between the average accuracy and training time. Further increasing K does not bring prominent performance improvement but takes more training time.

slimming strategy	average accuracy (%)
slim-random	66.98
slim-first	66.99
slim-last	67.07
slim-middle	67.14

(b) **Depth Slimming.** The *slim-middle* strategy outperforms all the counterparts in terms of average accuracy, verifying that the middle layers in Transformer are less important than the first and last ones. More evidence is shown in Fig. 8.

training strategy	average accuracy (%)
w/ random init, w/o KD	65.89
w/ random init, w/ ID [71]	65.55
w/ random init, w/ KD	66.95
w/ teacher init, w/ KD	67.14

(d) **Training Strategy.** The teacher initialization and KD strategies show advantages over the random initialization and ID strategies in terms of average accuracy, respectively.

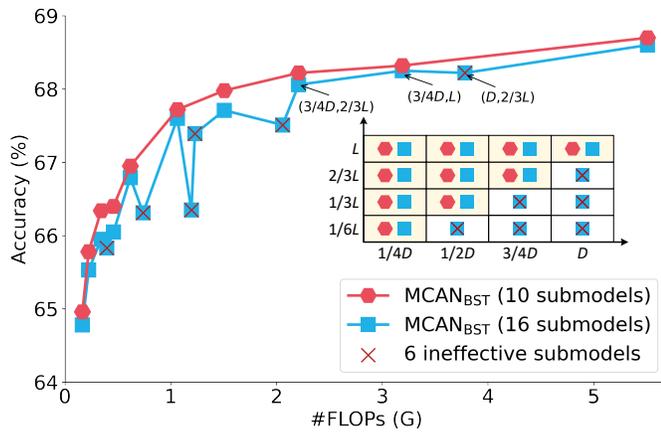


Fig. 7: **Submodel Selection.** The accuracies vs. FLOPs on VQA-v2 `minival` split are reported to compare two MCAN_{BST} variants with 10 and 16 submodels, respectively. The network architectures of the 10 submodels (red squares) correspond to a subset of the 16 submodels (blue squares) after discarding 6 *ineffective* submodels (red crosses) by using our triangle selection strategy.

suggesting that the bottom and top layers of the Transformer are more important than the middle layers. This observation is consistent with the results in [78]. We also provide visualization results in section IV-F by introducing a simple score function based on the magnitudes of the output features of each layer. The calculated layer scores are consistent with those obtained by the slim-middle strategy, verifying its effectiveness from a side view.

Number of Sampled Submodels. In Table IVc, we ablate the effects of different numbers of sampled submodels during the BST training. The results suggest that $K=4$ is a good trade-off between the average accuracy and training time. Further

increasing K does not bring a great performance improvement but takes much more training time. The fast saturation with such a small K is facilitated by the weight-sharing strategy of different submodels.

Training Strategy. Our default training strategy uses the model parameters from a teacher model for initialization and then trains the submodels using a KD training strategy to exploit the implicit knowledge from the teacher model. The results in Table IVd show that both the teacher model initialization and the KD training improve the obtained MCAN_{BST} model, compared to the model variants trained with random initialization or supervised by the ground-truth answer. In contrast to our KD training strategy that uses a fixed teacher model, an alternative strategy introduces a special *inplace distillation* (ID) training strategy that takes the largest submodel as the *dynamic* teacher to perform knowledge distillation [71]. We note that the ID-based strategy results in worse performance than the KD training strategy (65.55% vs. 66.95% in terms of average accuracy), and even underperforms standard training without knowledge distillation (65.55% vs. 65.89%). This suggests that a stable teacher model plays a key role in BST.

Submodel Selection. In Fig. 7, we compare two MCAN_{BST} variants with 10 and 16 submodels. From the results, we have the following observations: 1) all the *ineffective* submodels, which require more FLOPs but obtain lower accuracies than some other submodels, are precisely detected by our simple triangle selection strategy, 2) removing such ineffective submodel architectures before the BST training brings improvements to all the remaining submodels, and 3) taking the submodel $(3/4D, 2/3L)$ as the reference model, the submodel $(3/4D, L)$ outperforms $(D, 2/3L)$ with fewer FLOPs, verifying our hypothesis that increasing depth is more economical than increasing width for the Transformer.

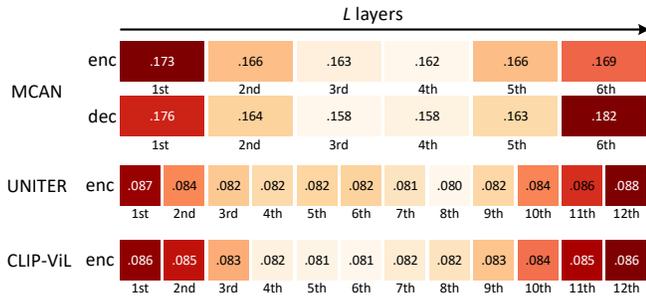


Fig. 8: Visualizations of the calculated layer importance scores given trained MCAN, UNITER, and CLIP-ViL models, respectively. A darker color indicates a larger score of the layer. The layer scores for all the three models exhibit a “heavy ends and light middle” phenomenon, which are consistent with the layer scores achieved by the *slim-middle* strategy.

F. Visualization Analysis

As discussed in section IV-E, the method of measuring the layer importance plays key role in depth slimming. Here, we introduce a data-driven strategy to measure the importance of each layer. Specifically, given a trained MCAN (UNITER or CLIP-ViL) model without model slimming, we first feed all the training samples through the network and memorize the output features for each layer. For the output features obtained from each layer, mean pooling is performed on each flattened feature vector to obtain its unnormalized score. Using this strategy, we show the importance scores for each layer of the MCAN, UNITER, and CLIP-ViL models in Fig. 8. From the visualized results, we see that the layer scores of all models exhibit “heavy tails and a light middle”, which is consistent with the layer scores achieved by the *slim-middle* strategy.

To better understand the behaviors of the weight-sharing submodels learned by BST, we show the attention maps from three $\text{MCAN}_{\text{BST}}(1/4D, L)$ submodels in Fig. 9. Due to space limitations, we only show one example and visualize the attention maps from the first and last layers of the question encoder and image decoder, respectively. To better understand the effect of the attention mechanism, we highlight some representative attention maps with blue bounding boxes. From the results, we have the following observations. In general, the slimmed submodels $\text{MCAN}_{\text{BST}}(1/4D, L)$ and $\text{MCAN}_{\text{BST}}(1/4D, 1/3L)$ have fewer *redundant* heads (*i.e.*, similar attention maps within one layer) than the full-sized model $\text{MCAN}_{\text{BST}}(D, L)$. This verifies the feasibility and necessity of our BST framework. Moreover, the three submodels, which have different widths and depths, all predict the correct answer. This verifies the effectiveness of both the width and depth slimming strategies, as well as the training paradigm. Furthermore, the attention maps from the different submodels have similar properties to the attention maps in the original MCAN paper [2]. Almost all the attention maps from the first layer of the question encoder (*i.e.*, enc-1) attend to the column of words such as ‘*what*’, which act as the question type classifiers. In contrast, some attention maps from the last layer of the question encoder (*i.e.*, enc-6) and image decoder (*i.e.*, dec-1 and dec-6) focus on the columns of keywords such as ‘*head*’.

V. CONCLUSION

In this paper, we present a new direction for the VQA task: learning efficient and elastic models that can adaptively fit different platforms. To this end, we present a general bilaterally slimmable Transformer (BST) framework that can be seamlessly integrated with any Transformer-based VQA model in theory. By integrating the BST framework with three typical Transformer-based VQA approaches, the resulting slimmable models outperform state-of-the-art methods with similar model sizes, or achieve comparable performance to that of much smaller models on both VQA-v2 and GQA datasets. Moreover, the quantitative experiments on diverse hardware platforms and devices show the practicability and necessity of BST.

To the best of our knowledge, the proposed BST framework is the first attempt to explore efficient and elastic models for VQA. We hope our general framework can serve as a baseline to inspire future research on efficient multimodal learning.

REFERENCES

- [1] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” in *Neural Information Processing Systems (NIPS)*, 2017, pp. 6000–6010.
- [2] Z. Yu, J. Yu, Y. Cui, D. Tao, and Q. Tian, “Deep modular co-attention networks for visual question answering,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 6281–6290.
- [3] Y.-C. Chen, L. Li, L. Yu, A. E. Kholly, F. Ahmed, Z. Gan, Y. Cheng, and J. Liu, “Uniter: Universal image-text representation learning,” in *European Conference on Computer Vision (ECCV)*, 2020, pp. 104–120.
- [4] S. Shen, L. H. Li, H. Tan, M. Bansal, A. Rohrbach, K.-W. Chang, Z. Yao, and K. Keutzer, “How much can clip benefit vision-and-language tasks?” in *International Conference on Learning Representations (ICLR)*, 2021.
- [5] Z. Yu, Y. Cui, J. Yu, M. Wang, D. Tao, and Q. Tian, “Deep multimodal neural architecture search,” in *Proceedings of the 28th ACM International Conference on Multimedia*, 2020, pp. 3743–3752.
- [6] J. Dong, X. Li, and D. Xu, “Cross-media similarity evaluation for web image retrieval in the wild,” *IEEE Transactions on Multimedia*, vol. 20, no. 9, pp. 2371–2384, 2018.
- [7] P. Anderson, X. He, C. Buehler, D. Teney, M. Johnson, S. Gould, and L. Zhang, “Bottom-up and top-down attention for image captioning and visual question answering,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 6077–6086.
- [8] L. Yu, Z. Lin, X. Shen, J. Yang, X. Lu, M. Bansal, and T. L. Berg, “MATTNet: Modular attention network for referring expression comprehension,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 1307–1315.
- [9] M. Sun, W. Suo, P. Wang, Y. Zhang, and Q. Wu, “A proposal-free one-stage framework for referring expression comprehension and generation via dense cross-attention,” *IEEE Transactions on Multimedia*, pp. 1–1, 2022.
- [10] N. Ouyang, Q. Huang, P. Li, Y. Cai, B. Liu, H.-f. Leung, and Q. Li, “Suppressing biased samples for robust vqa,” *IEEE Transactions on Multimedia*, vol. 24, pp. 3405–3415, 2021.
- [11] Z. Yu, Y. Cui, J. Yu, D. Tao, and Q. Tian, “Multimodal unified attention networks for vision-and-language interactions,” *arXiv preprint arXiv:1908.04107*, 2019.
- [12] H. Tan and M. Bansal, “Lxmert: Learning cross-modality encoder representations from transformers,” in *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2019, pp. 5100–5111.
- [13] J. Li, R. Selvaraju, A. Gotmare, S. Joty, C. Xiong, and S. C. H. Hoi, “Align before fuse: Vision and language representation learning with momentum distillation,” in *Neural Information Processing Systems (NIPS)*, 2021, pp. 9694–9705.
- [14] Z. Yu, J. Yu, C. Xiang, J. Fan, and D. Tao, “Beyond bilinear: Generalized multimodal factorized high-order pooling for visual question answering,” *IEEE Transactions on Neural Networks and Learning Systems*, vol. 29, no. 12, pp. 5947–5959, 2018.
- [15] J.-H. Kim, J. Jun, and B.-T. Zhang, “Bilinear attention networks,” in *Neural Information Processing Systems (NIPS)*, 2018, pp. 1571–1581.

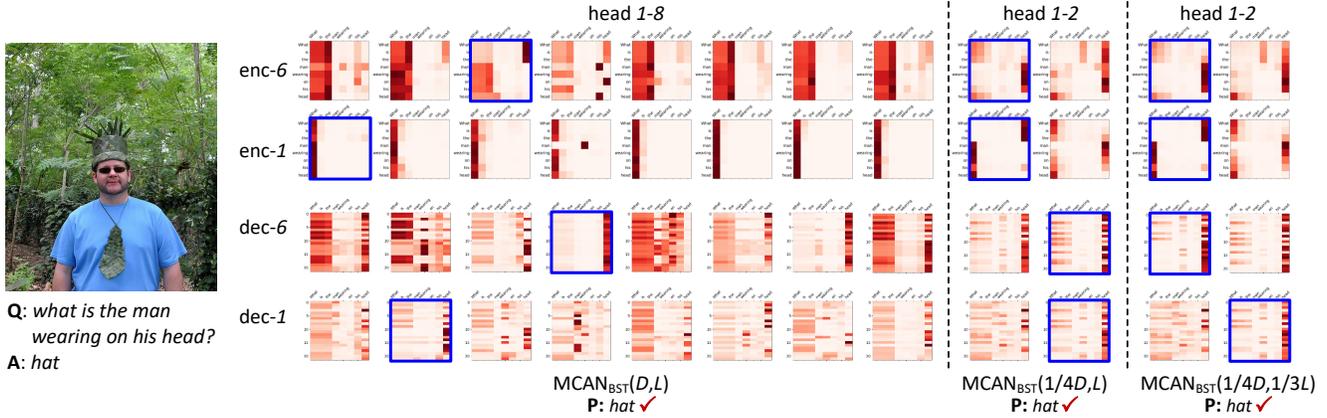


Fig. 9: Visualizations of the learned attention maps ($\text{softmax}(QK^T/\sqrt{D_H})$) from three MCAN_{BST} submodels trained on VQA-v2. For each submodel, we show the attention maps from different heads in the first and last layers of the question encoder and image decoder, respectively. We highlight some representative attention maps with blue bounding boxes to better understand the attention mechanisms behind.

- [16] W. Su, X. Zhu, Y. Cao, B. Li, L. Lu, F. Wei, and J. Dai, "Vi-bert: Pre-training of generic visual-linguistic representations," in *International Conference on Learning Representations (ICLR)*, 2020.
- [17] J. Lu, D. Batra, D. Parikh, and S. Lee, "Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks," in *Neural Information Processing Systems (NIPS)*, 2019, pp. 13–23.
- [18] X. Ma, P. Zhang, S. Zhang, N. Duan, Y. Hou, M. Zhou, and D. Song, "A tensorized transformer for language modeling," in *Neural Information Processing Systems (NIPS)*, 2019, pp. 2229–2239.
- [19] Z. Lan, M. Chen, S. Goodman, K. Gimpel, P. Sharma, and R. Soricut, "Albert: A lite bert for self-supervised learning of language representations," in *International Conference on Learning Representations (ICLR)*, 2019.
- [20] M. Dehghani, S. Gouws, O. Vinyals, J. Uszkoreit, and L. Kaiser, "Universal transformers," in *International Conference on Learning Representations (ICLR)*, 2018.
- [21] B. Cui, Y. Li, M. Chen, and Z. Zhang, "Fine-tune bert with sparse self-attention mechanism," in *Conference on Empirical Methods in Natural Language Processing and the International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, 2019, pp. 3548–3553.
- [22] A. Fan, E. Grave, and A. Joulin, "Reducing transformer depth on demand with structured dropout," in *International Conference on Learning Representations (ICLR)*, 2019.
- [23] R. Tang, Y. Lu, L. Liu, L. Mou, O. Vechtomova, and J. Lin, "Distilling task-specific knowledge from bert into simple neural networks," *arXiv preprint arXiv:1903.12136*, 2019.
- [24] V. Sanh, L. Debut, J. Chaumond, and T. Wolf, "Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter," *arXiv preprint arXiv:1910.01108*, 2019.
- [25] Z. Sun, H. Yu, X. Song, R. Liu, Y. Yang, and D. Zhou, "Mobilebert: a compact task-agnostic bert for resource-limited devices," in *Annual Meeting of the Association for Computational Linguistics (ACL)*, 2020, pp. 2158–2170.
- [26] W. Liu, P. Zhou, Z. Wang, Z. Zhao, H. Deng, and Q. Ju, "Fastbert: a self-distilling bert with adaptive inference time," in *Annual Meeting of the Association for Computational Linguistics (ACL)*, 2020, pp. 6035–6044.
- [27] W. Wang, F. Wei, L. Dong, H. Bao, N. Yang, and M. Zhou, "Minilm: Deep self-attention distillation for task-agnostic compression of pre-trained transformers," in *Neural Information Processing Systems (NIPS)*, 2020, pp. 5776–5788.
- [28] J. Yu, L. Yang, N. Xu, J. Yang, and T. Huang, "Slimmable neural networks," in *International Conference on Learning Representations (ICLR)*, 2018.
- [29] L. Hou, Z. Huang, L. Shang, X. Jiang, X. Chen, and Q. Liu, "Dynabert: Dynamic bert with adaptive width and depth," in *Neural Information Processing Systems (NIPS)*, 2020, pp. 9782–9793.
- [30] B. Qin, H. Hu, and Y. Zhuang, "Deep residual weight-sharing attention network with low-rank attention for visual question answering," *IEEE Transactions on Multimedia*, vol. 24, pp. 1–1, 2022.
- [31] B. Zhou, Y. Tian, S. Sukhbaatar, A. Szlam, and R. Fergus, "Simple baseline for visual question answering," *International Conference on Learning Representations (ICLR)*, 2015.
- [32] S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. Lawrence Zitnick, and D. Parikh, "Vqa: Visual question answering," in *IEEE International Conference on Computer Vision (ICCV)*, 2015, pp. 2425–2433.
- [33] A. Fukui, D. H. Park, D. Yang, A. Rohrbach, T. Darrell, and M. Rohrbach, "Multimodal compact bilinear pooling for visual question answering and visual grounding," *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 457–468, 2016.
- [34] J.-H. Kim, K. W. On, W. Lim, J. Kim, J.-W. Ha, and B.-T. Zhang, "Hadamard product for low-rank bilinear pooling," in *International Conference on Learning Representation (ICLR)*, 2017.
- [35] Z. Yu, J. Yu, J. Fan, and D. Tao, "Multi-modal factorized bilinear pooling with co-attention learning for visual question answering," *IEEE International Conference on Computer Vision (ICCV)*, pp. 1839–1848, 2017.
- [36] Z. Yang, X. He, J. Gao, L. Deng, and A. Smola, "Stacked attention networks for image question answering," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 21–29.
- [37] J. Yu, W. Zhang, Y. Lu, Z. Qin, Y. Hu, J. Tan, and Q. Wu, "Reasoning on the relation: Enhancing visual representation for visual question answering and cross-modal retrieval," *IEEE Transactions on Multimedia*, vol. 22, no. 12, pp. 3196–3209, 2020.
- [38] J. Lu, J. Yang, D. Batra, and D. Parikh, "Hierarchical question-image co-attention for visual question answering," in *Neural Information Processing Systems (NIPS)*, 2016, pp. 289–297.
- [39] D.-K. Nguyen and T. Okatani, "Improved fusion of visual and language representations by dense symmetric co-attention for visual question answering," *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 6087–6096, 2018.
- [40] F. Liu, J. Liu, Z. Fang, R. Hong, and H. Lu, "Visual question answering with dense inter-and intra-modality interactions," *IEEE Transactions on Multimedia*, vol. 23, pp. 3518–3529, 2020.
- [41] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, 2019, pp. 4171–4186.
- [42] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-end object detection with transformers," in *European conference on computer vision (ECCV)*. Springer, 2020, pp. 213–229.
- [43] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, "An image is worth 16x16 words: Transformers for image recognition at scale," 2021.
- [44] J. Yu, J. Li, Z. Yu, and Q. Huang, "Multimodal transformer with multi-view visual representation for image captioning," *IEEE transactions on*

- circuits and systems for video technology*, vol. 30, no. 12, pp. 4467–4480, 2019.
- [45] J. Lei, L. Li, L. Zhou, Z. Gan, T. L. Berg, M. Bansal, and J. Liu, “Less is more: Clipbert for video-and-language learning via sparse sampling,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 7331–7341.
- [46] L. Li, Y.-C. Chen, Y. Cheng, Z. Gan, L. Yu, and J. Liu, “Hero: Hierarchical encoder for video+ language omni-representation pre-training,” in *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2020, pp. 2046–2065.
- [47] L. Zhu and Y. Yang, “Actbert: Learning global-local video-text representations,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 8746–8755.
- [48] P. Gao, Z. Jiang, H. You, P. Lu, S. C. Hoi, X. Wang, and H. Li, “Dynamic fusion with intra-and inter-modality attention flow for visual question answering,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 6639–6648.
- [49] Z. Yu, J. Wei, Yu, J. Zhu, and Z. Kuang, “Knowledge-representation-enhanced multimodal transformer for scene text question answering,” *Journal of Image and Graphics*, vol. 27, no. 9, pp. 2761–2774, 2022.
- [50] R. Hu, A. Singh, T. Darrell, and M. Rohrbach, “Iterative answer prediction with pointer-augmented multimodal transformers for textvqa,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 9992–10 002.
- [51] Y. Liu, Y.-S. Wei, H. Yan, G.-B. Li, and L. Lin, “Causal reasoning meets visual representation learning: A prospective study,” *Machine Intelligence Research*, no. 19, pp. 1–27, 2022.
- [52] L. Chen, X. Yan, J. Xiao, H. Zhang, S. Pu, and Y. Zhuang, “Counterfactual samples synthesizing for robust visual question answering,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 10 800–10 809.
- [53] Z. Yang, Z. Gan, J. Wang, X. Hu, Y. Lu, Z. Liu, and L. Wang, “An empirical study of gpt-3 for few-shot knowledge-based vqa,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, no. 3, 2022, pp. 3081–3089.
- [54] Z. Shao, Z. Yu, M. Wang, and J. Yu, “Prompting large language models with answer heuristics for knowledge-based visual question answering,” *arXiv preprint arXiv:1607.06450*, 2023.
- [55] X. Li, X. Yin, C. Li, P. Zhang, X. Hu, L. Zhang, L. Wang, H. Hu, L. Dong, F. Wei *et al.*, “Oscar: Object-semantics aligned pre-training for vision-language tasks,” in *European Conference on Computer Vision (ECCV)*, 2020, pp. 121–137.
- [56] F. Yu, J. Tang, W. Yin, Y. Sun, H. Tian, H. Wu, and H. Wang, “Ernievil: Knowledge enhanced vision-language representations through scene graphs,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 4, 2021, pp. 3208–3216.
- [57] Y. Cui, Z. Yu, C. Wang, Z. Zhao, J. Zhang, M. Wang, and J. Yu, “Rosita: Enhancing vision-and-language semantic alignments via cross-and intra-modal knowledge integration,” in *Proceedings of the 29th ACM International Conference on Multimedia*, 2021, pp. 797–806.
- [58] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, “Learning transferable visual models from natural language supervision,” in *International Conference on Machine Learning (ICML)*. PMLR, 2021, pp. 8748–8763.
- [59] P. Wang, A. Yang, R. Men, J. Lin, S. Bai, Z. Li, J. Ma, C. Zhou, J. Zhou, and H. Yang, “OFA: unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework,” in *International Conference on Machine Learning (ICML)*, 2022, pp. 23 318–23 340.
- [60] J. Li, D. Li, C. Xiong, and S. Hoi, “Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation,” in *International Conference on Machine Learning (ICML)*, 2022, pp. 12 888–12 900.
- [61] Z.-Y. Dou, Y. Xu, Z. Gan, J. Wang, S. Wang, L. Wang, C. Zhu, P. Zhang, L. Yuan, N. Peng *et al.*, “An empirical study of training end-to-end vision-and-language transformers,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 18 166–18 176.
- [62] S. Han, H. Mao, and W. J. Dally, “Deep compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding,” in *International Conference on Learning Representations (ICLR)*, 2016.
- [63] S. Han, J. Pool, J. Tran, and W. Dally, “Learning both weights and connections for efficient neural network,” in *Neural Information Processing Systems (NIPS)*, 2015, pp. 1135–1143.
- [64] Z. Liu, J. Li, Z. Shen, G. Huang, S. Yan, and C. Zhang, “Learning efficient convolutional networks through network slimming,” in *IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 2736–2744.
- [65] F. N. Iandola, S. Han, M. W. Moskewicz, K. Ashraf, W. J. Dally, and K. Keutzer, “Squeezenet: Alexnet-level accuracy with 50x fewer parameters and 0.5 mb model size,” *arXiv preprint arXiv:1602.07360*, 2016.
- [66] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, “Mobilenets: Efficient convolutional neural networks for mobile vision applications,” *arXiv preprint arXiv:1704.04861*, 2017.
- [67] X. Zhang, X. Zhou, M. Lin, and J. Sun, “Shufflenet: An extremely efficient convolutional neural network for mobile devices,” in *IEEE conference on computer vision and pattern recognition (CVPR)*, 2018, pp. 6848–6856.
- [68] Z. Wu, T. Nagarajan, A. Kumar, S. Rennie, L. S. Davis, K. Grauman, and R. Feris, “Blockdrop: Dynamic inference paths in residual networks,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 8817–8826.
- [69] G. Huang, Y. Sun, Z. Liu, D. Sedra, and K. Q. Weinberger, “Deep networks with stochastic depth,” in *European Conference on Computer Vision (ECCV)*, 2016, pp. 646–661.
- [70] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778, 2016.
- [71] J. Yu and T. S. Huang, “Universally slimmable networks and improved training techniques,” in *IEEE International Conference on Computer Vision (ICCV)*, 2019, pp. 1803–1811.
- [72] H. Cai, C. Gan, T. Wang, Z. Zhang, and S. Han, “Once-for-all: Train one network and specialize it for efficient deployment,” in *International Conference on Learning Representations (ICLR)*, 2019.
- [73] J. L. Ba, J. R. Kiros, and G. E. Hinton, “Layer normalization,” *arXiv preprint arXiv:1607.06450*, 2016.
- [74] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov, “Roberta: A robustly optimized bert pretraining approach,” *arXiv preprint arXiv:1907.11692*, 2019.
- [75] A. Radford, K. Narasimhan, T. Salimans, and I. Sutskever, “Improving language understanding by generative pre-training,” URL https://s3-us-west-2.amazonaws.com/openai-assets/research-covers/languageunsupervised/language_understanding_paper.pdf, 2018.
- [76] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever *et al.*, “Language models are unsupervised multitask learners,” *OpenAI blog*, vol. 1, no. 8, p. 9, 2019.
- [77] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu, “Exploring the limits of transfer learning with a unified text-to-text transformer,” *Journal of Machine Learning Research (JMLR)*, vol. 21, no. 140, pp. 1–67, 2020.
- [78] W. Wang and Z. Tu, “Rethinking the value of transformer components,” in *International Conference on Computational Linguistics (COLING)*, 2020, pp. 6019–6029.
- [79] Y. Goyal, T. Khot, D. Summers-Stay, D. Batra, and D. Parikh, “Making the v in vqa matter: Elevating the role of image understanding in visual question answering,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 6904–6913.
- [80] H. Wang, Z. Wu, Z. Liu, H. Cai, L. Zhu, C. Gan, and S. Han, “Hat: Hardware-aware transformers for efficient natural language processing,” in *Annual Meeting of the Association for Computational Linguistics (ACL)*, 2020, pp. 7675–7688.
- [81] A. Khetan and Z. Karnin, “schubert: Optimizing elements of bert,” in *Annual Meeting of the Association for Computational Linguistics (ACL)*, 2020, pp. 2807–2818.
- [82] G. Hinton, O. Vinyals, and J. Dean, “Distilling the knowledge in a neural network,” *arXiv preprint arXiv:1503.02531*, 2015.
- [83] D. A. Hudson and C. D. Manning, “Gqa: A new dataset for real-world visual reasoning and compositional question answering,” *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 6700–6709, 2019.
- [84] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, “Microsoft coco: Common objects in context,” in *European Conference on Computer Vision (ECCV)*, 2014, pp. 740–755.
- [85] R. Krishna, Y. Zhu, O. Groth, J. Johnson, K. Hata, J. Kravitz, S. Chen, Y. Kalantidis, L.-J. Li, D. A. Shamma *et al.*, “Visual genome: Connecting language and vision using crowdsourced dense image annotations,” *International Journal of Computer Vision (IJCV)*, vol. 123, no. 1, pp. 32–73, 2017.