

Dynamic Neural Networks: A Survey

Yizeng Han*, Gao Huang*, *Member, IEEE*, Shiji Song, *Senior Member, IEEE*, Le Yang, Honghui Wang, and Yulin Wang

Abstract—Dynamic neural network is an emerging research topic in deep learning. Compared to static models which have fixed computational graphs and parameters at the inference stage, dynamic networks can adapt their structures or parameters to different inputs, leading to notable advantages in terms of accuracy, computational efficiency, adaptiveness, etc. In this survey, we comprehensively review this rapidly developing area by dividing dynamic networks into three main categories: 1) *sample-wise* dynamic models that process each sample with data-dependent architectures or parameters; 2) *spatial-wise* dynamic networks that conduct adaptive computation with respect to different spatial locations of image data; and 3) *temporal-wise* dynamic models that perform adaptive inference along the temporal dimension for sequential data such as videos and texts. The important research problems of dynamic networks, e.g., architecture design, decision making scheme, optimization technique and applications, are reviewed systematically. Finally, we discuss the open problems in this field together with interesting future research directions.

Index Terms—Dynamic networks, Adaptive inference, Efficient inference, Convolutional neural networks.

1 INTRODUCTION

DEEP neural networks (DNNs) are playing an important role in various areas including computer vision (CV) [1], [2], [3], [4], [5] and natural language processing (NLP) [6], [7], [8]. Recent years have witnessed many successful deep models such as AlexNet [1], VGG [2], GoogleNet [3], ResNet [4], DenseNet [5] and Transformer [6]. These architectural innovations have enabled the training of deeper, more accurate and more efficient models. The recent research on neural architecture search (NAS) [9], [10] further speeds up the process of designing more powerful structures. However, most of the prevalent deep learning models perform inference in a static manner, i.e., both the computational graph and the network parameters are fixed once trained, which may limit their representation power, efficiency and interpretability [11], [12], [13], [14].

Dynamic networks, as opposed to static ones, can adapt their structures or parameters to the input during inference, and therefore enjoy favorable properties that are absent in static models. In general, dynamic computation in the context of deep learning has the following advantages:

1) Efficiency. One of the most notable advantages of dynamic networks is that they are able to allocate computations on demand at test time, by selectively activating model components (e.g. layers [12], channels [15] or sub-networks [16]) *conditioned* on the input. Consequently, less computation is spent on canonical samples that are relatively easy to recognize, or on less informative spatial/temporal locations of an input. In addition to *computational efficiency*, dynamic models have also shown promising results for improving *data efficiency* in the scenario of few-shot learning [17], [18].

2) Representation power. Due to the data-dependent network architecture/parameters, dynamic networks have significantly enlarged parameter space and improved rep-

resentation power. For example, with a minor increase of computation, model capacity can be boosted by applying feature-conditioned attention weights on an ensemble of convolutional kernels [13], [19]. It is worth noting that the popular soft attention mechanism could also be unified in the framework of dynamic networks, as different channels [20], spatial areas [21] or temporal locations [22] of features are dynamically re-weighted at test time.

3) Adaptiveness. Dynamic models are able to achieve a desired trade-off between accuracy and efficiency for dealing with varying computational budgets on the fly. Therefore, they are more adaptable to different hardware platforms and changing environments, compared to static models with a fixed computational cost.

4) Compatibility. Dynamic networks are compatible with most advanced techniques in deep learning, including architecture design [4], [5], optimization algorithms [23], [24] and data preprocessing [25], [26], which ensures that they can benefit from the most recent advances in the field to achieve state-of-the-art performance. For example, dynamic networks can inherit architectural innovations in lightweight models [27], or be designed via NAS approaches [9], [10]. Their efficiency could also be further improved by acceleration methods developed for static models, such as network pruning [28], weight quantization [29], knowledge distillation [30] and low-rank approximation [31].

5) Generality. As a substitute for static deep learning techniques, many dynamic models are general approaches that can be applied seamlessly to a wide range of applications, such as image classification [12], [32], object detection [33] and semantic segmentation [34]. Moreover, the techniques developed in CV tasks are proven to transfer well to language models in NLP tasks [35], [36], and vice versa.

6) Interpretability. We finally note that the research on dynamic networks may bridge the gap between the underlying mechanism of deep models and brains, as it is believed that the brains process information in a dynamic way [37], [38]. It is possible to analyze which components of a dynamic model are activated [32] when processing an input sample, and to observe which parts of the input are ac-

• Yizeng Han, Gao Huang, Shiji Song, Le Yang, Honghui Wang and Yulin Wang are with the Department of Automation, Tsinghua University, Beijing 100084, China. E-mail: {hanyz18, yangl15, wanghh20, wangyl19}@mails.tsinghua.edu.cn; {gaohuang, shijis}@tsinghua.edu.cn. Corresponding author: Gao Huang.

*. Equal contribution.

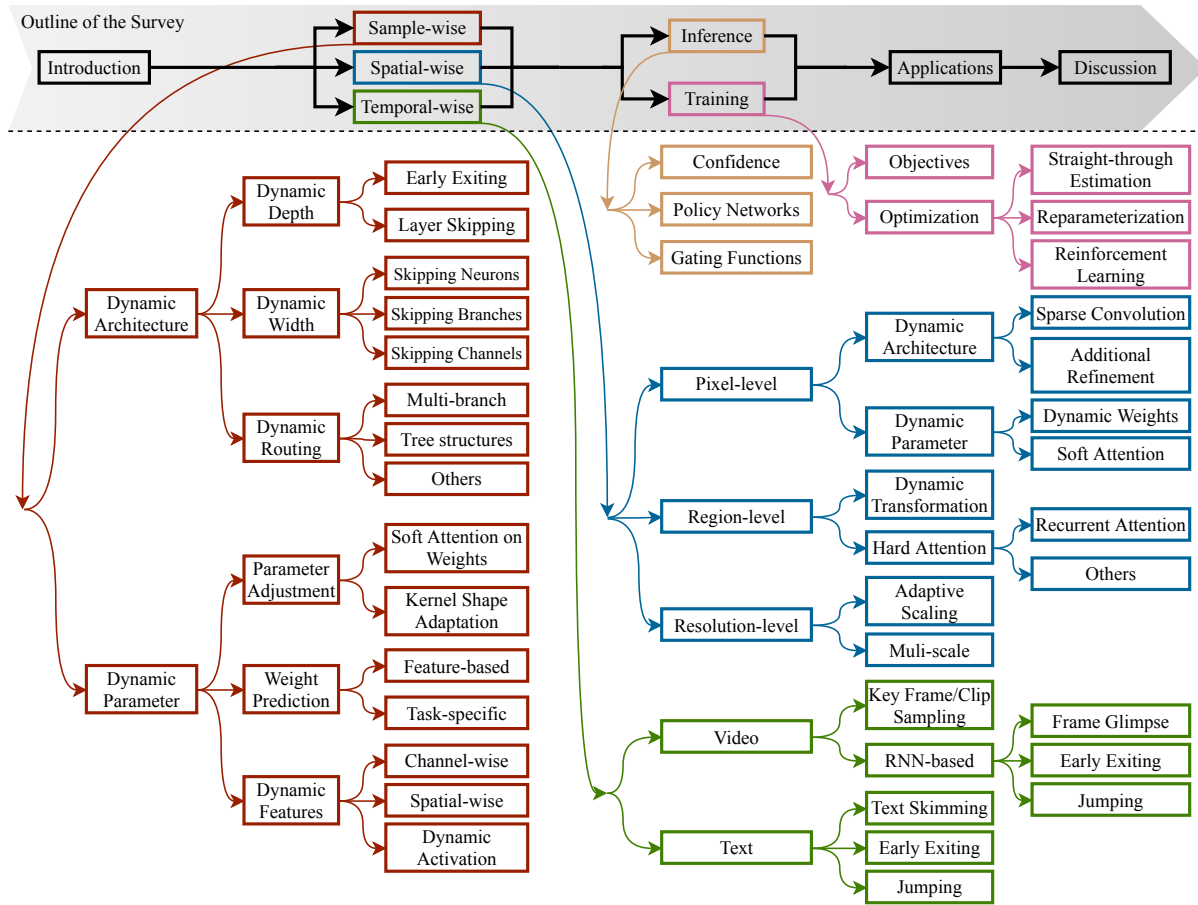


Fig. 1. Overview of the survey. We first review the dynamic networks that perform adaptive computation at three different granularities (i.e. sample-wise, spatial-wise and temporal-wise). Then we summarize the decision making strategy, training technique and applications of dynamic models. Existing open problems in this field together with some future research directions are finally discussed. Best viewed in color.

TABLE 1
Notations used in this paper.

Notations	Descriptions
\mathbb{R}^m	m -dimensional real number domain
a, \mathbf{a}	Scalar, vector/matrix/tensor
\mathbf{x}, \mathbf{y}	Input, output feature
\mathbf{x}^ℓ	Feature at layer ℓ
\mathbf{h}_t	Hidden state at time step t
$\mathbf{x}(\mathbf{p})$	Feature at spatial location \mathbf{p} on \mathbf{x}
Θ	Learnable parameter
$\Theta \mathbf{x}$	Dynamic parameter conditioned on \mathbf{x}
$\mathbf{x} \star \mathbf{W}$	Convolution of feature \mathbf{x} and weight \mathbf{W}
\otimes	Channel-wise or element-wise multiplication
$\mathcal{F}(\cdot, \Theta)$	Functional Operation parameterized by Θ
$\mathcal{F} \circ \mathcal{G}$	Composition of function \mathcal{F} and \mathcal{G}

countable for certain predictions [39]. These properties may shed light on interpreting the decision process of DNNs.

In fact, adaptive inference, the key idea underlying dynamic networks, has been studied before the popularity of modern DNNs. The most classical approach is building a model ensemble through a cascaded [40] or parallel [41] structure, and selectively activating the models conditioned on the input. Spiking neural networks (SNNs) [42], [43] also perform data-dependent inference by propagating pulse signals. However, the training strategy for SNNs is quite different from that of popular convolutional neural networks (CNNs), and they are less used in vision tasks. Therefore, we leave out the work related to SNNs in this survey.

In the context of deep learning, dynamic inference with modern deep architectures, has raised many new research questions and has attracted great research interests in the

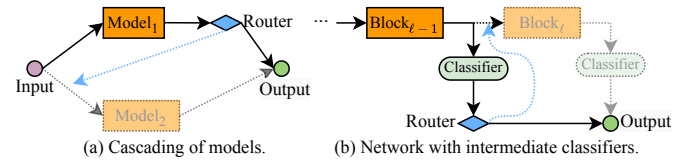


Fig. 2. Two early-exiting schemes. The dashed lines and shaded modules are not executed, conditioned on the decisions made by the routers.

past three years. Despite the extensive work on designing various types of dynamic networks, a systematic and comprehensive review on this topic is still lacking. This motivates us to write this survey, to review the recent advances in this rapidly developing area, with the purposes of 1) providing an overview as well as new perspectives for researchers who are interested in this topic; 2) pointing out the close relations of different subareas and reducing the risk of reinventing the wheel; and 3) summarizing the key challenges and possible future research directions.

This survey is organized as follows (see Fig. 1 for an overview). In Sec. 2, we introduce the most common *sample-wise* dynamic networks which adapt their architectures or parameters conditioned on each input sample. Dynamic models working on a finer granularity, i.e., *spatially* adaptive and *temporally* adaptive models, are reviewed in Sec. 3 and Sec.4, respectively¹. Then we investigate the decision making strategies and the training techniques of dynamic

1. These two categories can also be viewed as sample-wise dynamic networks as they perform adaptive computation within each sample at a finer granularity, and we adopt such a split for narrative convenience.

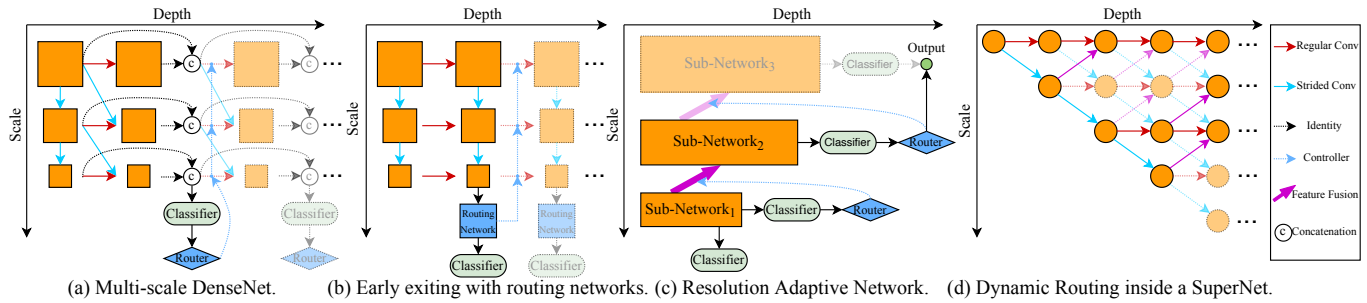


Fig. 3. Multi-scale architectures with dynamic inference graphs. The first three models (a, b, c) perform adaptive early exiting with specific architecture designs and exiting policies. Dynamic routing is achieved inside a SuperNet (d) to activate data-dependent inference paths.

networks in Sec. 5. The applications of dynamic models are further summarized in Sec. 6. Finally, we conclude this survey with a discussion on a number of open problems and future research directions in Sec. 7. For better readability, we list the notations that will be used in this survey in Table 1.

2 SAMPLE-WISE DYNAMIC NETWORKS

Aiming at processing different inputs in data-dependent manners, sample-wise dynamic networks are typically designed from two perspectives: 1) adjusting model architectures to allocate appropriate computation based on each sample, and therefore reducing redundant computation for increased efficiency (Sec. 2.1); 2) adapting network *parameters* to every input sample with fixed computational graphs, with the goal of boosting the representation power with minimal increase of computational cost (Sec. 2.2).

2.1 Dynamic Architectures

Considering different inputs may have diverse computational demands, it is natural to perform inference with dynamic architectures conditioned on each sample. Specifically, one can adjust the network depth (Sec. 2.1.1), width (Sec. 2.1.2), or perform dynamic routing within a super network (SuperNet) that includes multiple possible paths (Sec. 2.1.3). Networks with dynamic architectures not only save redundant computation for canonical (“easy”) samples, but also preserve their representation power when recognizing non-canonical (“hard”) samples. Such a property leads to remarkable advantages in efficiency compared to the acceleration techniques for static models [28], [29], [44], which handle “easy” and “hard” inputs with identical computation, and fail to reduce intrinsic computational redundancy.

2.1.1 Dynamic Depth

As modern DNNs are getting increasingly deep for recognizing more “hard” samples, a straightforward solution to reducing redundant computation is performing inference with dynamic depth, which can be realized by 1) *early exiting*, i.e. allowing “easy” samples to be output at shallow exits without executing deeper layers [12], [45], [46]; or 2) *layer skipping*, i.e. selectively skipping intermediate layers conditioned on each sample [11], [47], [48].

1) Early exiting. The complexity (or “difficulty”) of input samples varies in most real-world scenarios, and shallow networks are capable of correctly identifying many canonical inputs. Ideally, these samples should be output at certain early exits without executing deeper layers.

For an input sample \mathbf{x} , the forward propagation of an L -layer deep network \mathcal{F} could be represented by

$$\mathbf{y} = \mathcal{F}^L \circ \mathcal{F}^{L-1} \circ \dots \circ \mathcal{F}^1(\mathbf{x}), \quad (1)$$

where \mathcal{F}^ℓ denotes the operational function at layer ℓ , $1 \leq \ell \leq L$. In contrast, early exiting allows to terminate the inference procedure at an intermediate layer. For the i -th input sample \mathbf{x}_i , the forward propagation can be written as

$$\mathbf{y}_i = \mathcal{F}^{\ell_i} \circ \mathcal{F}^{\ell_i-1} \circ \dots \circ \mathcal{F}^1(\mathbf{x}_i), 1 \leq \ell_i \leq L. \quad (2)$$

Note that ℓ_i is adaptively determined based on \mathbf{x}_i . Extensive architectures have been studied to endow DNNs with such early exiting behaviors, as discussed in the following.

a) *Cascading of DNNs.* The most intuitive approach to enabling early exiting is cascading multiple models (see Fig. 2 (a)), and adaptively retrieving the prediction of an early network without activating latter ones. For example, Big/little-Net [49] cascades two CNNs with different depths. After obtaining the *SoftMax* output of the first model, early exiting is conducted when the score margin between the two largest elements exceeds a threshold. Moreover, a number of classic CNNs [1], [3], [4] are cascaded in [46] and [50]. After each model, a decision function is trained to determine whether the obtained feature should be fed to a linear classifier for immediate prediction, or be sent to subsequent models.

b) *Intermediate classifiers.* The models in the aforementioned cascading structures are mutually independent. Consequently, once a “difficult” sample is decided to be fed to a latter network, a whole inference procedure needs to be executed from scratch without reusing the already learned features. A more compact design is involving intermediate classifiers within one backbone network (see Fig. 2 (b)), so that early features can be propagated to deep layers if needed. Based on such a multi-exit architecture, adaptive early exiting is typically achieved according to confidence-based criteria [45], [51] or learned functions [46], [52], [53].

c) *Multi-scale architecture with early exits.* Researchers [12] have observed that in chain-structured networks, the multiple classifiers may interfere with each other, which degrades the overall performance. A reasonable interpretation could be that in regular CNNs, the high-resolution features lack the coarse-level information that is essential for classification, leading to unsatisfying results for early exits. Moreover, early classifiers would force the shallow layers to generate *task-specialized* features, while a part of *general* information is lost, resulting in degraded performance for deep exits. To tackle this issue, multi-scale dense network (MSDNet) [12] adopts 1) a *multi-scale* architecture, which consists of multiple sub-networks for processing feature maps with different resolutions (scales), to quickly generate coarse-level features that are suitable for classification; 2) *dense connections*, to reuse early features and improve the performance of deep classifiers (see Fig. 3 (a)). Such a specially-designed architecture effectively enhances the

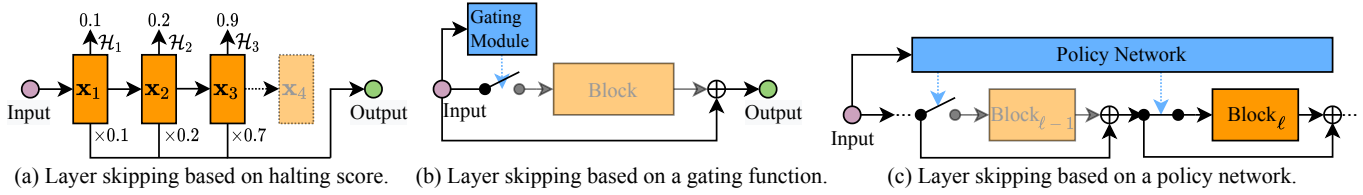


Fig. 4. Dynamic layer skipping. Feature x_4 in (a) are not calculated conditioned on the halting score, and the gating module in (b) decides whether to execute the block based on the intermediate feature. The policy network in (c) generates the skipping decisions for all layers in the main network.

overall accuracy of all the classifiers in the network.

Based on the multi-scale architecture design, researchers have also studied the exiting policies [54], [55] (see Fig. 3 (b)) and training schemes [56] of early-exiting dynamic models. More discussion about the inference and training schemes for dynamic models will be presented in Sec. 5.

Previous methods typically achieve the adaptation of network depths. From the perspective of exploiting spatial redundancy in features, resolution adaptive network (RANet, see Fig. 3 (c)) [32] first processes each sample with low-resolution features, while high-resolution representations are conditionally utilized based on early predictions.

Adaptive early exiting is also extended to language models (e.g. BERT [7]) for improving their efficiency on NLP tasks [57], [58], [59], [60]. In addition, it can be implemented in recurrent neural networks (RNNs) for *temporally* dynamic inference when processing sequential data such as videos [61], [62] and texts [63], [64], [65] (see Sec. 4).

2) Layer skipping. The general idea of the aforementioned early-exiting paradigm is skipping the execution of all the deep layers after a certain classifier. More flexibly, the network depth can also be adapted on the fly by strategically skipping the calculation of *intermediate layers* without placing extra classifiers. Given the i -th input sample x_i , dynamic layer skipping could be generally written as

$$y_i = (\mathbb{1}^L \circ \mathcal{F}^L) \circ (\mathbb{1}^{L-1} \circ \mathcal{F}^{L-1}) \circ \dots \circ (\mathbb{1}^1 \circ \mathcal{F}^1)(x_i), \quad (3)$$

where $\mathbb{1}^\ell$ denotes the indicator function determining the execution of layer \mathcal{F}^ℓ , $1 \leq \ell \leq L$. This scheme is typically implemented on structures with skip connections (e.g. ResNet [4]) to guarantee the continuity of forward propagation, and here we summarize three common approaches.

a) *Halting score* is first proposed in [11], where an accumulated scalar named halting score adaptively decides whether the hidden state of an RNN will be directly fed to the next time step. The halting scheme is extended to vision tasks by viewing residual blocks within a ResNet stage² as linear layers within a step of RNN [33] (see Fig. 4 (a)). Rather than skipping the execution of layers with independent parameters, multiple blocks in each ResNet stage could be replaced by one weight-sharing block, leading to a significant reduction of parameters [66]. In every stage, the block is executed for an adaptive number of steps according to the halting score.

In addition to RNNs and CNNs, the halting scheme is further implemented on Transformers [6] by [35] and [36] to achieve dynamic network depth on NLP tasks.

b) *Gating function* is also a prevalent option for dynamic layer skipping due to its plug-and-play property. Take ResNet as an example (see Fig. 4 (b)), let x^ℓ denote the input feature of the ℓ -th residual block, gating function \mathcal{G}^ℓ

generates a binary value to decide the execution of residual block \mathcal{F}^ℓ . This procedure could be represented by³

$$x^{\ell+1} = \mathcal{G}^\ell(x^\ell)\mathcal{F}^\ell(x^\ell) + x^\ell, \mathcal{G}^\ell(x^\ell) \in \{0, 1\}. \quad (4)$$

SkipNet [47] and convolutional network with adaptive inference graph (Conv-AIG) [48] are two typical approaches to enabling dynamic layer skipping. Both methods induce lightweight computational overheads to efficiently produce the binary decisions on whether skipping the calculation of a residual block. Specifically, Conv-AIG utilizes two FC layers in each residual block, while the gating function in SkipNet is implemented as an RNN for parameter sharing.

Rather than skipping layers in classic ResNets, dynamic recursive network [67] iteratively executes one block with shared parameters in each stage. Although the weight-sharing scheme seems similar to the aforementioned IamNN [66], the skipping policy of [67] differs significantly: gating modules are exploited to decide the recursion depth.

Instead of either skipping a layer, or executing it thoroughly with a full numerical precision, a line of work [68], [69] studies adaptive *bit-width* for different layers conditioned on the *resource budget*. Furthermore, fractional skipping [70] adaptively selects a bit-width for each residual block by a gating function based on *input features*.

c) *Policy network* can be built to take in an input sample, and directly produces the skipping decisions for all the layers in a backbone network [71] (see Fig. 4 (c)).

2.1.2 Dynamic Width

In addition to dynamic network *depth* (Sec. 2.1.1), a finer-grained form of conditional computation is performing inference with dynamic *width*: although every layer is executed, its multiple units (e.g. neurons, channels or branches) are selectively activated conditioned on the input.

1) Skipping neurons in fully-connected (FC) layers. The computational cost of an FC layer is determined by its input and output dimensions. It is commonly believed that different neuron units are responsible for representing different features, and therefore not all of them need to be activated for every sample. Early studies learn to adaptively control the neuron activations by auxiliary branches [72], [73], [74] or other techniques such as low-rank approximation [75].

2) Skipping branches in mixture-of-experts (MoE). In Sec. 2.1.1, adaptive model ensemble is achieved via a *cascading* way, and later networks are conditionally executed based on early predictions. An alternative approach to improving the model capacity is the MoE [41], [76] structure, which means that multiple network branches are built as experts in *parallel*. These experts could be selectively executed, and their outputs are fused with data-dependent weights.

Conventional *soft* MoEs [41], [76], [77] adopt real-valued weights to dynamically rescale the representations obtained

2. Here we refer to a stage as a stack of multiple residual blocks with the same feature resolution.

3. For simplicity and without generality, the subscript for sample index will be omitted in the following.

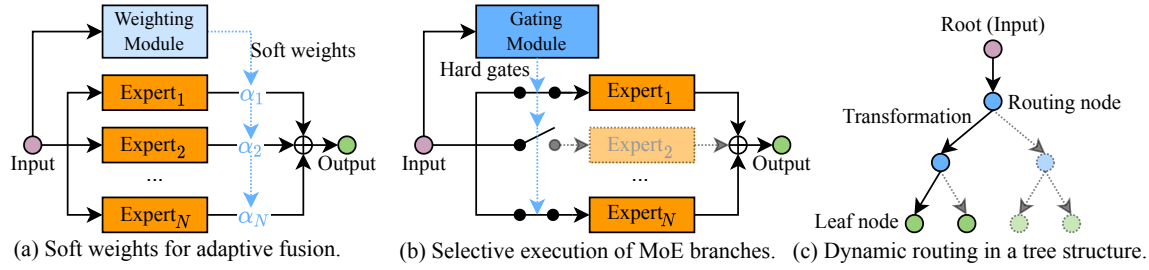


Fig. 5. MoE with soft weighting (a) and hard gating (b) schemes both adopt an auxiliary module to generate the weights or gates. In the tree structure (c), features (nodes) and transformations (paths) are represented as circles and lines with arrows respectively. Only the full lines are activated.

from different experts (Fig. 5 (a)). In this way, all the branches still need to be executed, and thus the computation cannot be reduced at test time. *Hard* gates with only a fraction of non-zero elements are developed to increase the inference efficiency of the MoE structure (see Fig. 5 (b)) [78], [79], [80]: let \mathcal{G} denote a gating module whose output is a N -dimensional vector α controlling the execution of N experts $\mathcal{F}_1, \mathcal{F}_2, \dots, \mathcal{F}_N$, the final output can be written as

$$\mathbf{y} = \sum_{n=1}^N [\mathcal{G}(\mathbf{x})]_n \mathcal{F}_n(\mathbf{x}) = \sum_{n=1}^N \alpha_n \mathcal{F}_n(\mathbf{x}), \quad (5)$$

and the n -th expert will not be executed if $\alpha_n = 0$.

Hard MoE has been implemented in diverse network structures. For example, HydraNet [78] replaces the convolutional blocks in a CNN by multiple branches, and selectively execute them conditioned on the input. For another example, dynamic routing network (DRNet) [80] performs a branch selection in each cell structure which is commonly used in NAS [10]. On NLP tasks, sparsely gated MoE [16] and switch Transformer [81] embeds hard MoE in a long short-term memory (LSTM) [82] network and a Transformer [6], respectively. Instead of making choice with *binary* gates as in [80], only the branches corresponding to the *top-K* elements of the real-valued gates are activated in [16], [78], [81].

3) Skipping channels in CNNs. Modern CNNs usually have considerable channel redundancy. Based on the common belief that the same feature channel can be of disparate importance for different samples, adaptive width of CNNs could be realized by dynamically activating convolutional channels. Compared to the *static* pruning methods [28], [44] which remove “unimportant” filters permanently, such a *data-dependent* pruning approach improves the inference efficiency without degrading the model capacity.

a) *Multi-stage architectures along the channel dimension.* Recall that the early-exiting networks [12], [32] discussed in Sec. 2.1.1 can be viewed as multi-stage models along the *depth* dimension, where late stages are conditionally executed based on early predictions. One can also build multi-stage architectures along the *width* (channel) dimension, and progressively execute these stages on demand.

Along this direction, an optimal architecture is searched among multiple structures with different widths, and any sample can be output at an early stage when a confident prediction is obtained [83]. Channel gating network (CGNet) [84] first executes a subset of convolutional filters in every layer, and the remaining filters are only activated on strategically selected areas.

b) *Dynamic pruning with gating functions.* In the aforementioned progressive activation paradigm, the execution of a later stage is decided based on previous output. As a result, a complete forward propagation is required for every stage, which might be suboptimal for reducing the practical

inference latency. Another prevalent solution is to decide the execution of channels at every layer by gating functions. For example, runtime neural pruning (RNP) [15] models the layer-wise pruning as a Markov decision process, and an RNN is used to select specific channel groups at every layer. Moreover, pooling operations followed by FC layers are utilized to generate *channel-wise hard attention* (i.e. making discrete decisions on whether to activate each channel) for each sample [85], [86], [87], [88]. The recent work [89] uses a gate module to decide the width for a whole stage of a ResNet. Different reparameterization and optimizing techniques are required for training these gating functions, which will be reviewed in Sec. 5.2.

Rather than placing plug-in gating modules *inside* a CNN, GaterNet [90] builds an *extra* network, which takes in the input sample and directly generates all the channel selection decisions for the backbone CNN. This implementation is similar to BlockDrop [71] that exploits an additional policy network for dynamic layer skipping (Sec. 2.1.1).

c) *Dynamic pruning based on feature activations directly* has also been realized without auxiliary branches and computational overheads [91], where a regularization item is induced in training to encourage the sparsity of features.

On basis of the existing literature on dynamically skipping either network *layers* [47], [48] or convolutional *filters* [15], [85], [86], [87], recent work [92], [93], [94] has realized dynamic inference with respect to network *depth* and *width* simultaneously: only if a layer is determined to be executed, its channels will be selectively activated, leading to a more flexible adaptation of computational graphs.

2.1.3 Dynamic Routing

The aforementioned methods mostly adjust the depth (Sec. 2.1.1) or width (Sec. 2.1.2) of *classic architectures* by activating their computational units (e.g. layers [47], [48] or channels [15], [87]) conditioned on the input. Another line of work develops different forms of SuperNets with various possible inference paths, and performs dynamic routing inside the SuperNets to adapt the computational graph to each sample.

To achieve dynamic routing, there are typically routing nodes that are responsible for allocating features to different paths. For node s at layer ℓ , let $\alpha_{s \rightarrow j}^\ell$ denote the probability of assigning the reached feature \mathbf{x}_s^ℓ to node j at layer $\ell + 1$, the path $\mathcal{F}_{s \rightarrow j}^\ell$ will be activated only when $\alpha_{s \rightarrow j}^\ell > 0$. The resulting feature reaching node j is represented by

$$\mathbf{x}_j^{\ell+1} = \sum_{s \in \{s: \alpha_{s \rightarrow j}^\ell > 0\}} \alpha_{s \rightarrow j}^\ell \mathcal{F}_{s \rightarrow j}^\ell(\mathbf{x}_s^\ell). \quad (6)$$

The probability $\alpha_{s \rightarrow j}^\ell$ can be obtained in different manners. Note that the dynamic early-exiting networks [12], [32] are a special form of SuperNets, where the routing decisions are only made at intermediate classifiers. The CapsuleNets

[14], [95] also perform dynamic routing between capsules, i.e. groups of neurons, to character the relations between (parts of) objects. Here we mainly focus on specific architecture designs of the SuperNets and their routing policies.

1) Path selection in multi-branch structures. The simplest dynamic routing can be implemented by selectively executing *one* of multiple candidate modules at each layer [96], [97], which is equivalent to producing a one-hot probability distribution $\alpha_{s \rightarrow \cdot}^{\ell}$ in Eq. 6. The main difference of this approach to hard MoE (Fig. 5 (b)) is that only one branch is activated without any fusion operations.

2) Neural trees and tree-structured networks. As decision trees always perform inference along one forward path that is dependent on input properties, combining tree structure with neural networks can naturally enjoy the adaptive inference paradigm and the representation power of DNNs simultaneously. Note that in a tree structure, the outputs of different nodes are routed to *independent* paths rather than being *fused* as in MoE structures (compare Fig. 5 (b), (c)).

a) *Soft decision tree* (SDT) [98], [99], [100] adopts neural units as routing nodes (blue nodes in Fig. 5 (c)), which decides the portion that the inputs are assigned to their left/right sub-tree. Each leaf node generates a probability distribution over the output space, and the final prediction is the expectation of the results from all leaf nodes. Although the probability for a sample reaching each leaf node in an SDT is data-dependent, all the paths are still executed, which limits the inference efficiency.

b) *Neural trees with deterministic routing policies* [101], [102] are designed to make hard routing decisions during inference, avoiding computation on those unselected paths.

c) *Tree-structured DNNs*. Instead of developing decision trees containing neural units, a line of work builds special network architectures to endow them with the routing behavior of decision trees. For instance, a small CNN is first executed to classify each sample into coarse categories, and specific sub-networks are conditionally activated based on the coarse predictions [103]. A subsequent work [104] not only partitions samples to different sub-networks, but also divides and routes the feature channels.

Different to those networks using neural units only in routing nodes [101], [102], or routing each sample to pre-designed sub-networks [103], [104], adaptive neural tree (ANT) [105] adopts CNN modules as feature transformers in a hard neural tree (see lines with arrows in Fig. 5 (c)), and learns the tree structure together with the network parameters simultaneously in the training stage.

3) Others. Performing dynamic routing within more general SuperNet architectures is also a recent research trend. Representatively, an architecture distribution with partly shared parameters is *searched* from a SuperNet containing $\sim 10^{25}$ sub-networks [106]. During inference, every sample is allocated by a controller network to one sub-network with appropriate computational cost. Instead of training a standalone controller network, gating modules are plugged inside the *hand-designed* SuperNet (see Fig. 3 (d)) to decide the routing path based on intermediate features [107].

2.2 Dynamic Parameters

Although the networks with dynamic *architectures* in Sec. 2.1 can adapt their inference graphs to each sample and

achieve an efficient allocation of computation, they usually have special architecture designs, requiring specific training strategies or careful hyper-parameters tuning (Sec. 7).

Another line of work adapts network *parameters* to different inputs while keeping the architectures fixed, which has been shown effective in improving the representation power of networks with a minor increase of computational cost. Given an input sample \mathbf{x} , the output of a conventional network (module) with static parameters can be written as $\mathbf{y} = \mathcal{F}(\mathbf{x}, \Theta)$. In contrast, the output of a model with dynamic parameters could be represented by

$$\mathbf{y} = \mathcal{F}(\mathbf{x}, \hat{\Theta}|\mathbf{x}) = \mathcal{F}(\mathbf{x}, \mathcal{W}(\mathbf{x}, \Theta)), \quad (7)$$

where $\mathcal{W}(\cdot, \Theta)$ is the operation producing input-dependent parameters, and its design has been extensively explored.

In general, the parameter adaptation can be achieved from three aspects (see Fig. 6): 1) adjusting the trained parameters based on the input (Sec. 2.2.1); 2) directly generating the network parameters from the input (Sec. 2.2.2); and 3) rescaling the features with soft attention (Sec. 2.2.3).

2.2.1 Parameter Adjustment

A typical approach to parameter adaptation is adjusting the weights based on their input during inference as presented in Fig. 6 (a). This implementation usually evokes little computation to obtain the adjustments, e.g., attention weights [13], [19], [108], [109] or sampling offsets [110], [111], [112].

1) Attention on weights. To improve the representation power without noticeably increasing the computation, soft attention can be performed on multiple convolutional kernels, producing an adaptive ensemble of parameters [13], [19]. Assuming that there are N kernels $\mathbf{W}_n, n = 1, 2, \dots, N$, such a dynamic convolution can be formulated as

$$\mathbf{y} = \mathbf{x} * \tilde{\mathbf{W}} = \mathbf{x} * \left(\sum_{n=1}^N \alpha_n \mathbf{W}_n \right). \quad (8)$$

This procedure increases the model capacity yet remains high efficiency, as the result obtained through fusing the outputs of N convolutional branches (as in MoE structures, see Fig. 5 (a)) is equivalent to that produced by performing once convolution with $\tilde{\mathbf{W}}$. However, only $\sim 1/N$ times of computation is consumed in the latter approach.

Weight adjustment could also be achieved by performing soft attention over the *spatial locations* of convolutional weights [108], [109]. For example, segmentation-aware convolutional network [108] applies locally masked convolution to aggregate information with larger weights from similar pixels, which are more likely to belong to the same object. Unlike [108] that requires a sub-network for feature embedding, pixel-adaptive convolution (PAC) [109] adapts the convolutional weights based on the attention mask generated from the input feature at each layer.

Instead of adjusting weights conditioned on every sample itself, meta-neighborhoods [113] adapt the network parameters to each input sample based on its similarity to the neighbors stored in a dictionary.

2) Kernel shape adaptation. Apart from adaptively scaling the weight *values*, parameter adjustment can also be realized to reshape the convolutional kernels and achieve *dynamic reception of fields*. Towards this direction, deformable convolutions [110], [111] sample feature pixels from adaptive locations when performing convolution on each pixel. Deformable kernels [112] samples weights in the kernel space

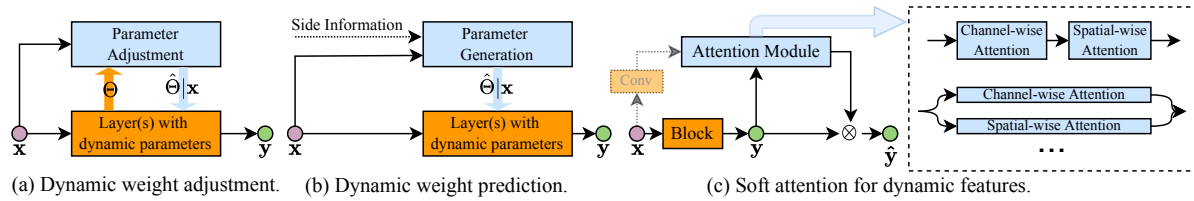


Fig. 6. Three implementations of dynamic parameters: adjusting (a) or generating (b) the backbone parameters based on the input, and (c) dynamically rescaling the features with the attention mechanism.

to adapt the *effective* reception field (ERF) while leaving the reception field unchanged. Table 2 summarizes the formulations of the above three methods. Due to their irregular memory access and computation pattern, these kernel shape adaptation approaches typically require customized CUDA kernels for the implementation on GPUs. However, recent literature has shown that the practical efficiency of deformable convolution could be effectively improved by co-designing algorithm and hardware based on embedded devices such as FPGAs [114].

2.2.2 Weight Prediction

Compared to making modifications on model parameters on the fly (Sec. 2.2.1), weight prediction [115] is more straightforward: it directly generates (a subset of) input-adaptive parameters with an independent model at test time (see Fig. 6 (b)). This idea was first suggested in [116], where both the weight prediction model and the backbone model were feedforward networks. Recent work has further extended the paradigm to modern network architectures and tasks.

1) General architectures. Dynamic filter networks (DFN) [117] and HyperNetworks [118] are two classic approaches realizing runtime weight prediction for CNNs and RNNs, respectively. Specifically, a filter generation network is built in DFN [117] to produce the filters for a convolutional layer. As for processing sequential data (e.g. a sentence), the weight matrices of the main RNN are predicted by a smaller one at each time step conditioned on the input (e.g. a word) [118]. WeightNet [119] unifies the dynamic schemes of [13] and [20] by predicting the convolutional weights via simple grouped FC layers, achieving competitive results in terms of the accuracy-FLOPs⁴ and accuracy-parameters trade-offs.

Rather than generating standard *convolutional* weights, LambdaNetworks [120] learns to predict the weights of *linear* projections based on the contexts of each pixel together with the relative position embeddings, showing advantages in terms of computational cost and memory footprint.

2) Task-specific information has also been exploited to predict model parameters on the fly, enabling dynamic networks to generate task-aware feature embeddings. For example, edge attributes are utilized in [121] to generate filters for graph convolution, and camera perspective is incorporated in [122] to generate weights for image convolution. Such task-aware weight prediction has been shown effective in improving the data efficiency on many tasks, including visual question answering [123], [124] and few-shot learning [17], [18].

2.2.3 Dynamic Features

The main goal of either *adjusting* (Sec. 2.2.1) or *predicting* (Sec. 2.2.2) model parameters is producing more dynamic

and informative features, and therefore enhancing the representation power of deep networks. A more straightforward solution is rescaling the features with input-dependent soft attention (see Fig. 6 (c)), which requires minor modifications on computational graphs. Note that for a linear transformation \mathcal{F} , applying attention α on the output is equivalent to performing computation with re-weighted parameters, i.e.

$$\mathcal{F}(\mathbf{x}, \Theta) \otimes \alpha = \mathcal{F}(\mathbf{x}, \Theta \otimes \alpha). \quad (9)$$

1) Channel-wise attention is one of the most common soft attention mechanisms. Existing work typically follows the form in squeeze-and-excitation network (SENet) [20]:

$$\tilde{\mathbf{y}} = \mathbf{y} \otimes \alpha = \mathbf{y} \otimes \mathcal{A}(\mathbf{y}), \alpha \in [0, 1]^C. \quad (10)$$

In Eq. 10, $\mathbf{y} = \mathbf{x} \star \mathbf{W}$ is the output feature of a convolutional layer with C channels, and $\mathcal{A}(\cdot)$ is a lightweight function composed of pooling and linear layers for producing α . Taking the convolution into account, the procedure can also be written as $\tilde{\mathbf{y}} = (\mathbf{x} \star \mathbf{W}) \otimes \alpha = \mathbf{x} \star (\mathbf{W} \otimes \alpha)$, from which we can observe that applying attention on features is equivalent to performing convolution with dynamic weights.

Other implementations for attention modules have also been developed, including using standard deviation to provide more statistics [125], or replacing FC layers with efficient 1D convolutions [126]. The empirical performance of three computational graphs for soft attention is studied in [127]: 1) $\tilde{\mathbf{y}} = \mathbf{y} \otimes \mathcal{A}(\mathbf{y})$, 2) $\tilde{\mathbf{y}} = \mathbf{y} \otimes \mathcal{A}(\mathbf{x})$ and 3) $\tilde{\mathbf{y}} = \mathbf{y} \otimes \mathcal{A}(\text{Conv}(\mathbf{x}))$. It is found that the three forms yield different performance in different backbone networks.

2) Spatial-wise attention. Spatial locations in features could also be dynamically rescaled with attention to improve the representation power of deep models [128]. Instead of using pooling operations to efficiently gather global information as in channel-wise attention, convolutions are often adopted in spatial-wise attention to encode local information. Moreover, these two types of attention modules can be integrated in one framework [21], [129], [130], [131] (see Fig. 6 (c)).

3) Dynamic activation functions. The aforementioned approaches to generating dynamic features usually apply soft attention before static activation functions. A recent line of work has sought to increase the representation power of models with dynamic activation functions [132], [133]. For instance, DY-ReLU [132] replaces ReLU ($\mathbf{y}_c = \max(\mathbf{x}_c, 0)$) with the max value among N linear transformations $\mathbf{y}_c = \max_n \{a_c^n \mathbf{x}_c + b_c^n\}$, where c is the channel index, and a_c^n, b_c^n are linear coefficients calculated from \mathbf{x} . On many vision tasks, these dynamic activation functions can effectively improve the performance of different network architectures with negligible computational overhead.

To summarize, soft attention has been exploited in many fields due to its simplicity and effectiveness. Moreover, it can be incorporated with other methods conveniently. E.g., by replacing the weighting scalar α_n in Eq. 5 with channel-wise [134] or spatial-wise [135] attention, the output of

4. Floating point operations, which is widely used as a measure of inference efficiency of deep networks.

TABLE 2
Kernel shape adaptation by dynamically sampling feature pixels [110], [111] or convolutional weights [112].

Method	Formulation	Sampled Target	Dynamic Mask
Regular Convolution	$y(p) = \sum_{k=1}^K W(p_k) x(p + p_k)$	-	-
Deformable ConvNet-v1 [110]	$y(p) = \sum_{k=1}^K W(p_k) x(p + p_k + \Delta p_k)$	Feature map	No
Deformable ConvNet-v2 [111]	$y(p) = \sum_{k=1}^K W(p_k) x(p + p_k + \Delta p_k) \Delta m_k$	Feature map	Yes
Deformable Kernels [112]	$y(p) = \sum_{k=1}^K W(p_k + \Delta p_k) x(p + p_k)$	Conv kernel	No

multiple branches with independent kernel sizes [134] or feature resolutions [135] are adaptively fused.

Note that we leave out the detailed discussion on the self attention mechanism, which is widely studied in both NLP [6], [7] and CV fields [136], [137], [138] to re-weight features based on the similarity between queries and keys at different locations (temporal or spatial). Readers who are interested in this topic may refer to review studies [139], [140], [141]. In this survey, we mainly focus on the feature re-weighting scheme in the framework of dynamic inference.

3 SPATIAL-WISE DYNAMIC NETWORKS

In visual learning, it has been found that not all locations contribute equally to the final prediction of CNNs [142], which suggests that *spatially* dynamic computation has great potential for reducing computational redundancy. In other words, making a correct prediction may only require processing a fraction of pixels or regions with an adaptive amount of computation. Moreover, based on the observations that low-resolution representations are sufficient to yield decent performance for most inputs [27], the static CNNs that take in all the input with the same resolution may also induce considerable redundancy.

To this end, spatial-wise dynamic networks are built to perform adaptive inference with respect to different spatial locations of images. According to the granularity of dynamic computation, we further categorize the relevant approaches into three levels: *pixel level* (Sec. 3.1), *region level* (Sec. 3.2) and *resolution level* (Sec. 3.3).

3.1 Pixel-level Dynamic Networks

Commonly seen spatial-wise dynamic networks perform adaptive computation at the pixel level. Similar to the categorization in Sec. 2, pixel-level dynamic networks are grouped into two types: models with pixel-specific *dynamic architectures* (Sec. 3.1.1) and *dynamic parameters* (Sec. 3.1.2).

3.1.1 Pixel-wise Dynamic Architectures

Based on the common belief that foreground pixels are more informative and computational demanding than those in the background, some dynamic networks learn to adjust their architectures for each pixel. Existing literature generally achieves this by 1) *dynamic sparse convolution*, which only performs convolutions on a subset of sampled pixels; 2) *additional refinement*, which strategically allocates extra computation (e.g. layers or channels) on certain spatial positions.

1) Dynamic sparse convolution. To reduce the unnecessary computation on less informative locations, convolution can be performed only on strategically sampled pixels. Existing sampling strategies include 1) making use of the intrinsic sparsity of the input [143]; 2) predicting the positions of zero elements on the output [144], [145]; and 3) estimating the saliency of pixels [146], [147], [148]. A typical approach is using an extra branch to generate a spatial mask, determining the execution of convolution on each pixel (see Fig. 7).

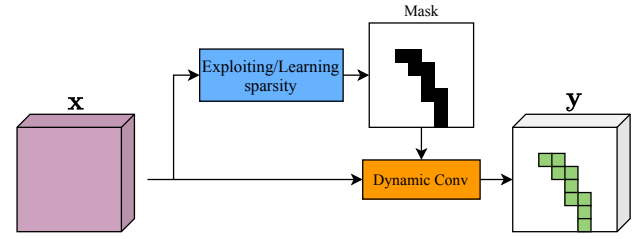


Fig. 7. Dynamic convolution on selected spatial locations. The 1 elements (black) in the spatial mask determine the pixels (green) that require computation in the output feature map.

Pixel-wise dynamic depth could also be achieved based on a halting scheme [33] (see Sec. 2.1.1). These dynamic convolutions usually neglect the unselected positions, which might degrade the network performance. Interpolation is utilized in [148] to efficiently fill those locations, therefore alleviating the aforementioned disadvantage.

2) Dynamic additional refinement. Instead of only sampling certain pixels to perform convolutions, another line of work first conducts relatively cheap computation on the whole feature map, and adaptively activate extra modules on selected pixels for further *refinement*. Representatively, dynamic capacity network [149] generates coarse features with a shallow model, and salient pixels are sampled based on the gradient information. For these salient pixels, extra layers are applied to extract finer features. Similarly, specific positions are additionally processed by a fraction of convolutional filters in [84]. These methods adapt their network architectures in terms of *depth* or *width* at the pixel level, achieving a spatially adaptive allocation of computation.

The aforementioned dynamic additional refinement approaches [84], [149] are mainly developed for image classification. On the semantic segmentation task, pixel-wise *early exiting* (see also Sec. 2.1.1) is proposed in [34], where the pixels with high prediction confidence are output without being processed by deeper layers. PointRend [150] shares a similar idea, and applies additional FC layers on selected pixels with low prediction confidence, which are more likely to be on borders of objects. All these researches demonstrate that by exploiting the spatial redundancy in image data, dynamic computation at the pixel level beyond sample level significantly increases the model efficiency.

3.1.2 Pixel-wise Dynamic Parameters

In contrast to entirely skipping the convolution operation on a subset of pixels, dynamic networks can also apply data-dependent parameters on different pixels for improved representation power or adaptive reception fields.

1) Dynamic weights. Similar to the sample-wise dynamic parameter methods (Sec. 2.2), pixel-level dynamic weights are achieved by test-time *adjustment* [108], [109], *prediction* [151], [152], [153], [154] or *dynamic features* [21], [129], [130], [135]. Take weight prediction as an example, typical approaches generate an $H \times W \times k^2$ kernel map to produce spatially dynamic weights (H, W are the spatial size of the

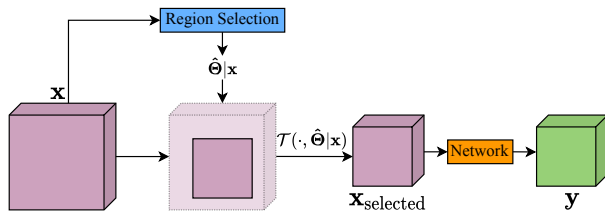


Fig. 8. Region-level dynamic inference. The region selection module generates the transformation/localization parameters, and the subsequent network performs inference on the transformed/cropped region.

output feature and k is the kernel size). Considering the pixels belonging to the same object may share identical weights, dynamic region-aware convolution (DRConv) [155] generates a segmentation mask for an input image, dividing it into m regions, for each of which a weight generation network is responsible for producing a data-dependent kernel.

2) Dynamic reception fields. Traditional convolution operations usually have a fixed shape and size of kernels (e.g. the commonly used 3×3 2D convolution). The resulting uniform reception field across all the layers may have limitations for recognizing objects with varying shapes and sizes. To tackle this, a line of work learns to adapt the reception field for different feature pixels [110], [111], [112], as discussed in Sec. 2.2.1. Instead of adapting the sampling location of features or kernels, adaptive connected network [156] realizes a dynamic trade-off among self transformation (e.g. 1×1 convolution), local inference (e.g. 3×3 convolution) and global inference (e.g. FC layer). The three branches of outputs are fused with data-dependent weighted summation. Besides images, the local and global information in non-Euclidean data, such as graphs, could also be adaptively aggregated.

3.2 Region-level Dynamic Networks

Pixel-level dynamic networks mentioned in Sec. 3.1 often require specific implementations for sparse computation, and consequently may face challenges in terms of achieving real acceleration on hardware [148]. An alternative approach is performing adaptive inference on *regions/patches* of input images. There mainly exists two lines of work along this direction (see Fig. 8): one performs parameterized *transformations* on a region of feature maps for more accurate prediction (Sec. 3.2.1), and the other learns patch-level *hard attention*, with the goal of improving the effectiveness and/or efficiency of models (Sec. 3.2.2).

3.2.1 Dynamic Transformations

Dynamic transformations (e.g. affine/projective/thin plate spline transformation) can be performed on images to undo certain variations [157] for better generalization ability, or to exaggerate the salient regions [158] for discriminative feature representation. For example, spatial transformer [157] adopts a localization network to generate the transformation parameters, and then applies the parameterized transformation to recover the input from the corresponding variations. Moreover, transformations are learned to adaptively zoom-in the salient regions on some tasks where the model performance is sensitive to a small portion of regions.

3.2.2 Hard Attention on Selected Patches

Inspired by the fact that informative features may only be contained in certain regions of an image, dynamic networks with hard spatial attention are explored to strategically select patches from the input for improved efficiency.

1) Hard attention with RNNs. The most typical approach is formulating a classification task as a sequential decision process, and adopting RNNs to make iterative predictions based on selected patches [159], [160]. For example, images are classified within a fixed number of steps, and at each step, the classifier RNN only sees a cropped patch, deciding the next attentional location until the last step is reached [159]. An adaptive step number is further achieved by including early stopping in the action space [160]. Glance-and-focus network (GFNet) [39] builds a general framework of region-level adaptive inference by sequentially focusing on a series of selected patches, and is compatible with most existing CNN architectures. The recurrent attention mechanism together with the early exiting paradigm enables both *spatially* and *temporally* adaptive inference [39], [160].

2) Hard attention with other implementations. Rather than using an RNN to predict the region position that the model should pay attention to, class activation mapping (CAM) [142] is leveraged in [161] to iteratively focus on salient patches. At each iteration, the selection is performed on the previously cropped input, leading to a progressive refinement procedure. A multi-scale CNN is built in [162], where the sub-network in each scale takes in the cropped patch from the previous scale, and is responsible for simultaneously producing 1) the feature representations for classification and 2) the attention map for the next scale. Without an iterative manner, the recent differentiable patch selection [163] adopts a differentiable top-K module to select a fixed number of patches in one step.

3.3 Resolution-level Dynamic Networks

The researches discussed above typically divide feature maps into different areas (pixel-level or region-level) for adaptive inference. On a coarser granularity, some dynamic networks could treat each image as a whole by processing feature representations with adaptive resolutions. Although it has been observed that a low resolution might be sufficient for recognizing most “easy” samples [27], conventional CNNs mostly process all the inputs with the same resolution, inducing considerable redundancy. Therefore, resolution-level dynamic networks exploit spatial redundancy from the perspective of feature resolution rather than the saliency of different locations. Existing approaches mainly include 1) scaling the inputs with adaptive ratios (Sec. 3.3.1); 2) selectively activating the sub-networks with different resolutions in a multi-scale architecture (Sec. 3.3.2).

3.3.1 Adaptive Scaling Ratios

Dynamic resolution can be achieved by scaling features with adaptive ratios. For example, a small sub-network is first executed to predict a scale distribution of faces on the face detection task, then the input images are adaptively zoomed, so that all the faces fall in a suitable range for recognition [164]. A plug-in module is used by [165] to predict the stride for the first convolution block in each ResNet stage, producing features with dynamic resolution.

3.3.2 Dynamic Resolution in Multi-scale Architectures

An alternative approach to achieving dynamic resolution is building multiple sub-networks in a parallel [166] or cascading [32] way. These sub-networks with different feature resolutions are selectively activated conditioned on

the input during inference. For instance, Elastic [166] realizes a *soft* selection from multiple branches at every layer, where each branch performs a downsample-convolution-upsample procedure with an independent scaling ratio. To practically avoid redundant computation, a *hard* selection is realized by [32], which allows each sample to conditionally activate sub-networks that process feature representations with resolution from low to high (see Fig. 3 (c) in Sec. 2.1.1).

4 TEMPORAL-WISE DYNAMIC NETWORKS

Apart from the spatial dimension (Sec. 3), adaptive computation could also be performed along the temporal dimension of sequential data, such as texts (Sec. 4.1) and videos (Sec. 4.2). In general, network efficiency can be improved by dynamically allocating less/no computation to the inputs at unimportant temporal locations.

4.1 RNN-based Dynamic Text Processing

Traditional RNNs mostly follow a static inference paradigm, i.e. input tokens are read sequentially to update a hidden state at each time step, which could be written as

$$\mathbf{h}_t = \mathcal{F}(\mathbf{x}_t, \mathbf{h}_{t-1}), t = 1, 2, \dots, T. \quad (11)$$

Such a static inference paradigm induces significant redundant computation, as different tokens usually have different contributions to the downstream tasks. A type of dynamic RNN is developed for allocating appropriate computational cost at each step. Some learn to “*skim*” unimportant tokens by dynamic update of hidden states (Sec. 4.1.1), and others conduct *adaptive reading* to avoid processing task-irrelevant tokens. Specifically, such adaptive reading can be achieved by *early exiting* (Sec. 4.1.2) or *dynamic jumping* (Sec. 4.1.3).

4.1.1 Dynamic Update of Hidden States

Since not all the tokens are essential for capturing the task-relevant information in a sequence, dynamic RNNs can be built to adaptively update their hidden states at each time step. Less informative tokens will be coarsely *skimmed*, i.e. the states are updated with cheaper computation.

1) Skipping the update. For unimportant inputs at certain temporal locations, dynamic models can learn to entirely skip the update of hidden states (see Fig. 9 (a)), i.e.

$$\mathbf{h}_t = \alpha_t \mathcal{F}(\mathbf{x}_t, \mathbf{h}_{t-1}) + (1 - \alpha_t) \mathbf{h}_{t-1}, \alpha_t \in \{0, 1\}. \quad (12)$$

For instance, Skip-RNN [167] updates a controlling signal in every step to determine whether to update or *copy* the hidden state from the previous step. An extra agent is adopted by Structural-Jump-LSTM [168] to make the skipping decision conditioned on the previous state and the current input. Without training the RNNs and the controllers jointly as in [167] and [168], a predictor is trained in [169] to estimate whether each input will make a “significant change” on the hidden state. The update is identified worthy to be executed only when the predicted change is greater than a threshold.

2) Coarse update. As directly skipping the update may be too aggressive, dynamic models could also update the hidden states with adaptively allocated operations. In specific, a network can adapt its architecture in every step, i.e.

$$\mathbf{h}_t = \mathcal{F}_t(\mathbf{x}_t, \mathbf{h}_{t-1}), t = 1, 2, \dots, T, \quad (13)$$

where \mathcal{F}_t is determined based on the input \mathbf{x}_t . One implementation is selecting a subset of dimensions of the hidden state to calculate, and copying the remaining from the

previous step [170], [171], as shown in Fig. 9 (b). To achieve the partial update, a subset of rows in weight matrices of the RNN is dynamically activated in [170], while Skim-RNN [171] makes a choice between two independent RNNs.

When the hidden states are generated by a multi-layer network, the update could be interrupted at an intermediate layer based on an accumulated halting score [11].

To summarize, a coarse update can be realized by with data-dependent network *depth* [11] or *width* [170], [171].

3) Selective updates in hierarchical RNNs. Considering the intrinsic hierarchical structure of texts (e.g. sentence-word-character), researchers have developed hierarchical RNNs to encode the temporal dependencies with different timescales using a dynamic update mechanism [172], [173]. During inference, the RNNs at higher levels will selectively update their states conditioned on the output of low-level ones (see Fig. 9 (c)). For example, when a character-level model in [172] detects that the input satisfies certain conditions, it will “*flush*” (reset) its states and feed them to a word-level network. Similar operations have also been realized by a gating module on question answering tasks [173].

4.1.2 Temporally Early Exiting in RNNs

Despite that the dynamic RNNs in Sec. 4.1.1 are able to update their states with data-dependent computational costs at each step, all the tokens still must be read, leading to inefficiency in scenarios where the task-relevant results can be obtained before reading the entire sequence.

Ideally, an efficient model should adaptively stop reading before the last step T in Eq. 11 is reached, once the captured information is satisfactory to solve the task. For instance, reasoning network (ReasonNet) [63] terminates its reading procedure when sufficient evidence has been found for question answering. Similarly, early stopping is implemented for sentence-level [174] and paragraph-level [65] text classification, respectively. Note that the approaches discussed here focus on making early predictions with respect to the *temporal* dimension of sequential input, rather than along the *depth* dimension of networks as in Sec. 2.1.1.

4.1.3 Jumping in Texts

Although early exiting in Sec. 4.1.2 can largely reduce redundant computation, all the tokens must still be fed to the model one by one. More aggressively, dynamic RNNs could further learn to decide “*where to read*” by strategically skipping some tokens without reading them, and directly jumping to an arbitrary temporal location (see Fig. 9 (d)).

Such dynamic jumping, together with early exiting, is realized in [175] and [64]. Specifically, LSTM-Jump [175] implements an auxiliary unit to predict the jumping stride within a defined range, and the reading process ends when the unit outputs zero. The model in [64] first decides whether to stop at each step. If not, it will further choose to re-read the current input, or to skip a flexible number of words. Moreover, structural information is exploited by Structural-Jump-LSTM [168], which utilizes an agent to decide whether to jump to the next punctuation. Apart from looking ahead, LSTM-Shuttle [176] also allows backward jumping to supplement the missed history information.

4.2 Temporal-wise Dynamic Video Recognition

For video recognition, where a video could be seen as a sequential input of frames, temporal-wise dynamic networks

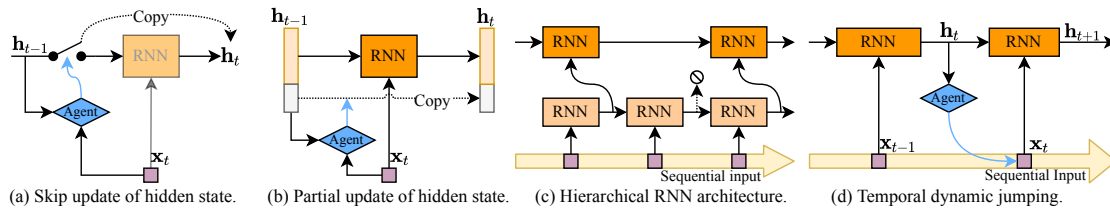


Fig. 9. Temporally adaptive inference. The first three approaches dynamically allocate computation in each step by (a) skipping the update, (b) partially updating the state, or (c) conditional computation in a hierarchical structure. The agent in (d) decides where to read in the next step.

are designed to allocate adaptive computational resources for different frames. This can generally be achieved by two approaches: 1) dynamically updating the hidden states in each time step of recurrent models (Sec. 4.2.1), and 2) performing adaptive *pre-sampling* for key frames (Sec. 4.2.2).

4.2.1 Video Recognition with Dynamic RNNs

Video recognition is often conducted via a recurrent procedure, where the video frames are first encoded by a 2D CNN, and the obtained frame features are fed to an RNN sequentially for updating its hidden state. Similar to the approaches introduced in Sec. 4.1, RNN-based adaptive video recognition is typically realized by 1) treating unimportant frames with relatively cheap computation ("*glimpse*") [177], [178]; 2) *early exiting* [61], [62]; and 3) performing dynamic *jumping* to decide "where to see" [61], [179], [180], [181].

1) Dynamic update of hidden states. To reduce redundant computation at each time step, LiteEval [177] makes a choice between two LSTMs with different computational costs. ActionSpotter [178] decides whether to update the hidden state according to each input frame. AdaFuse [182] selectively reuses certain feature channels from the previous step to efficiently make use of historical information. Recent work has also proposed to adaptively decide the numerical precision [183] or modalities [184], [185] when processing the sequential input frames. Such a *glimpse* procedure (i.e. allocating cheap operations to unimportant frames) is similar to the aforementioned text *skimming* [167], [168].

2) Temporally early exiting. Humans are able to comprehend the contents easily before watching an entire video. Such early stopping is also implemented in dynamic networks to make predictions only based on a portion of video frames [61], [62], [186]. Together with the *temporal* dimension, the model in [62] further achieves early exiting from the aspect of network *depth* as discussed in Sec. 2.1.1.

3) Jumping in videos. Considering encoding those unimportant frames with a CNN still requires considerable computation, a more efficient solution could be dynamically skipping some frames without watching them. Existing arts [179], [180], [187] typically learn to predict the location that the network should jump to at each time step. Furthermore, both early stopping and dynamic jumping are allowed in [61], where the jumping stride is limited in a discrete range. Adaptive frame (AdaFrame) [181] generates a continuous scalar within the range of $[0, 1]$ as the relative location.

4.2.2 Dynamic Key Frame Sampling

Rather than processing video frames recurrently as in Sec. 4.2.1, another line of work first performs an adaptive *pre-sampling* procedure, and then makes prediction by processing the selected subset of key frames or clips.

1) Temporal attention is a common technique for networks to focus on salient frames. For face recognition, neural ag-

gregation network [22] uses *soft* attention to adaptively aggregate frame features. To improve the inference efficiency, *hard* attention is realized to remove unimportant frames iteratively with RL for efficient video face verification [188].

2) Sampling module is also a prevalent option for dynamically selecting the key frames/clips in a video. For example, the frames are first sampled uniformly in [189], [190], and discrete decisions are made for each selected frame to go forward or backward step by step. As for clip-level sampling, SCSample [191] is designed based on a trained classifier to find the most informative clips for prediction. Moreover, dynamic sampling network (DSN) [192] segments each video into multiple sections, and a sampling module with shared weights across the sections is exploited to sample one clip from each section.

Adjusting multiple factors of deep models simultaneously has attracted researches in both static [193], [194] and dynamic networks [195], [196], [197], [198]. For example, together with *temporal-wise* frame sampling, *spatially* adaptive computation can be achieved by spatial [196]/temporal [199] resolution adaptation and patch selection [197], [200]. It would be promising to exploit the redundancy in both *input data* and *network structure* for further improving the efficiency of deep networks.

5 INFERENCE AND TRAINING

In previous sections, we have reviewed three different types of dynamic networks (sample-wise (Sec. 2), spatial-wise (Sec. 3) and temporal-wise (Sec. 4)). It can be observed that making data-dependent *decisions* at the inference stage is essential to achieve high efficiency and effectiveness. Moreover, *training* dynamic models is usually more challenging than optimizing static networks.

Note that since parameter adaptation (Sec. 2.2) could be conveniently achieved by differentiable operations, models with dynamic parameters [13], [20], [119] can be directly trained by stochastic gradient descent (SGD) without specific techniques. Therefore, in this section we mainly focus on discrete decision making (Sec. 5.1) and its training strategies (Sec. 5.2), which are absent in most static models.

5.1 Decision Making of Dynamic Networks

As described above, dynamic networks are capable of making data-dependent decisions during inference to transform their architectures, parameters, or to select salient spatial/temporal locations in the input. Here we summarize three commonly seen decision making schemes as follows.

5.1.1 Confidence-based Criteria

Many dynamic networks [12], [32], [45] are able to output "easy" samples at early exits if a certain confidence-based criterion is satisfied. These methods generally require estimating the confidence of intermediate predictions, which is

compared to a predefined threshold for decision making. In classification tasks, the confidence is usually represented by the maximum element of the *SoftMax* output [12], [32]. Alternative criteria include the entropy [45], [58] and the score margin [49]. On NLP tasks, a *model patience* is proposed in [60]: when the predictions for one sample stay unchanged after a number of classifiers, the inference procedure stops.

In addition, the halting score in [11], [33], [35], [36] could also be viewed as confidence for whether the current feature could be output to the next time step or calculation stage.

Empirically, the confidence-based criteria are easy to implement, and generally require no specific training techniques. A trade-off between accuracy and efficiency is controlled by manipulating the thresholds, which are usually tuned on a validation dataset. Note that the *overconfidence* issue in deep models [201], [202] might affect the effectiveness of such decision paradigm, when the incorrectly classified samples could obtain a high confidence at early exits.

5.1.2 Policy Networks

It is a common option to build an additional policy network learning to adapt the network topology based on different samples. Specifically, each input sample is first processed by the policy network, whose output directly determines which parts of the main network should be activated. For example, BlockDrop [71] and GaterNet [90] use a policy network to adaptively decide the *depth* and *width* of a backbone network. More generally, dynamic routing in a *SuperNet* can also be controlled by a policy network [106].

One possible limitation of this scheme is that the architectures and the training process of some policy networks are developed for a specific backbone [71], [90], and may not be easily adapted to different architectures.

5.1.3 Gating Functions

Gating function is a general and flexible approach to decision making in dynamic networks. It can be conveniently adopted as a plug-in module at arbitrary locations in any backbone network. During inference, each module is responsible for controlling the local inference graph of a layer or block. The gating functions take in intermediate features and efficiently produce binary-valued gate vectors to decide: 1) which channels need to be activated [15], [85], [86], [87], [88] *width*, 2) which layers need to be skipped [47], [48], [92], [93], 3) which paths should be selected in a *SuperNet* [107], or 4) what locations of the input should be allocated computations [146], [147], [148], [182].

Compared to the aforementioned decision policies, the gating functions demonstrate notable generality and applicability. However, due to their lack of differentiability, these gating functions usually need specific training techniques, which will be introduced in the following Sec. 5.2.

5.2 Training of Dynamic Networks

Besides architecture design, training is also essential for dynamic networks. Here we summarize the existing training strategies for dynamic models from the perspectives of objectives and optimization.

5.2.1 Training Objectives for Efficient Inference

1) Training multi-exit networks. We first notice that early-exiting dynamic networks [12], [32] are generally trained

by minimizing a weighted cumulative loss of intermediate classifiers. One challenge for training such models is the joint optimization of multiple classifiers, which may interfere with each other. MSDNet [12] alleviates the problem through its special architecture design. Several improved training techniques [56] are proposed for multi-exit networks, including a gradient equilibrium algorithm to stabilize the training process, and a bi-directional knowledge transfer approach to boost the collaboration of classifiers. For temporal-wise early exiting, the training of the policy network in FrameExit [186] is supervised by pseudo labels.

2) Encouraging sparsity. Many dynamic networks adapt their inference procedure by conditionally activating their computational units [47], [87] or strategically sampling locations from the input [148]. Training these models without additional constraints would result in superfluous computational redundancy, as a network could tend to activate all the candidate units for minimizing the task-specific loss.

The overall objective function for restraining such redundancy are typically written as $\mathcal{L} = \mathcal{L}_{\text{task}} + \gamma \mathcal{L}_{\text{sparse}}$, where γ is the hyper-parameter balancing the two items for the trade-off between accuracy and efficiency. In real-world applications, the second item can be designed based on the gate/mask values of candidate units (e.g. channels [86], [87], layers [47], [48] or spatial locations [148]). Specifically, one may set a target activation rate [48], [86] or minimizing the \mathcal{L}_1 norm of the gates/masks [148]. It is also practical to directly optimize a resource-aware loss (e.g. FLOPs) [92], [107], [147], which can be estimated according to the input and output feature dimension for every candidate unit.

3) Others. Note that extra loss items are mostly designed for but not limited to improving efficiency. Take [162] as an example, the model progressively focuses on a selected region, and is trained with an additional *inter-scale pairwise ranking loss* for proposing more discriminative regions. Moreover, knowledge distilling is utilized to boost the co-training of multiple sub-networks in [84] and [56].

5.2.2 Optimization of Non-differentiable Functions

A variety of dynamic networks contain non-differentiable functions that make discrete decisions to modify their architectures or sampling spatial/temporal locations from the input. These functions can not be trained directly with back-propagation. Therefore, specific techniques are studied to enable the end-to-end training as follows.

1) Gradient estimation is proposed to approximate the gradients for those non-differentiable functions and enable back-propagation. In [72], [172], straight-through estimator (STE) is exploited to heuristically copy the gradient with respect to the stochastic output directly as an estimator of the gradient with respect to the *Sigmoid* argument.

2) Reparameterization is also a popular technique to optimize the discrete decision functions. For instance, the gating functions controlling the network width [86] or depth [48] can both be trained with *Gumbel SoftMax* [259], [260], which is also used for pixel-level dynamic convolution [147], [148]. An alternative technique is *Improved SemHash* [261] adopted in [88] and [90] to train their hard gating modules.

Note that although these reparameterization techniques enable joint optimizing dynamic models together with gating modules in an end-to-end fashion, they usually lead to

TABLE 3

Applications of Dynamic Networks. For the type column, Sa, Sp and Te stand for sample-wise, spatial-wise and temporal-wise respectively.

Fields	Data	Type	Subfields & references
Computer Vision	Image	Sa	Object detection (face [40], [203], [204], facial point [205], pedestrian [206], general [33], [207], [208], [209], [210]) Image segmentation [107], [211], [212], Super resolution [213], Style transfer [214], Coarse-to-fine classification [215]
		Sa & Sp	Image segmentation [34], [129], [146], [148], [150], [154], [156], [216], [217], [218], [219], [220], Image-to-image translation [221], Object detection [110], [111], [147], [148], [164], Semantic image synthesis [222], [223], [224], Image denoising [225], Fine-grained classification [158], [162], [226], [227] Eye tracking [158], Super resolution [151], [153], [228]
		Sa & Sp & Te	General classification [39], [159], [161], Multi-object classification [229], [230], Fine-grained classification [160]
	Video	Sa	Multi-task learning (human action recognition and frame prediction) [231]
		Sa & Te	Classification (action recognition) [61], [177], [181], [189], [190], [191], [192], [196], [232], Semantic segmentation [233] Video face recognition [22], [188], Action detection [179], [180], Action spotting [178], [187]
		Sa & Sp & Te	Classification [196], [197], Frame interpolation [234], [235], Super resolution [236], Video deblurring [237], [238], Action prediction [239]
Natural Language Processing	Text	Sa & Sp	3D Shape classification and segmentation, 3D scene segmentation [240], 3D semantic scene completion [241]
		Sa	Neural language inference, Text classification, Paraphrase similarity matching, and Sentiment analysis [59], [60]
		Sa & Te	Language modeling [11], [16], [118], [170], [172], Machine translation [16], [35], [36], Classification [64], [65], [174], Sentiment analysis [168], [169], [171], [175], [176], Question answering [35], [63], [168], [171], [173]
Cross-Field	Image captioning [130], [242], Video captioning [243], [244], Visual question answering [123], [124], [245], Multi-modal sentiment analysis [246], [247]		
Others	Time series forecasting [248], [249], [250], Link prediction [251], Recommendation system [77], [252], [253], [254] Graph classification [121], Document classification [156], [255], [256], [257], Stereo confidence estimation [258]		

a longer training process for the decision functions to converge into a stable situation [144]. Moreover, the model performance might be sensitive to some extra hyper-parameters (e.g. temperature in *Gumbel SoftMax*), which might also increase the training cost for these dynamic networks.

3) Reinforcement learning (RL) is widely exploited for training non-differentiable decision functions. In specific, the backbones are trained by standard SGD, while the agents (either policy networks in Sec. 5.1.2 or gating functions in Sec. 5.1.3) are trained with RL to take discrete actions for dynamic inference graphs [15], [47], [71] or spatial/temporal sampling strategies [39], [190].

One challenge for RL-based training is the design of reward functions, which is important to the accuracy-efficiency tradeoff of dynamic models. Commonly seen reward signals are usually constructed to minimize a penalty item of the computational cost [15], [47]. Moreover, the training could be costly due to a multi-stage procedure: a pre-training process may be required for the backbone networks before the optimization of decision [71] or sampling [39] modules, and joint finetuning may be indispensable finally.

6 APPLICATION OF DYNAMIC NETWORKS

In this section, we summarize the applications of dynamic networks. Representative methods are listed in Table 3 based on the input data modality.

For image recognition, most dynamic CNNs are designed to conduct *sample-wise* or *spatial-wise* adaptive inference on classification tasks, and many inference paradigms can be generalized to other applications. Note that as mentioned in Sec. 3.2, the object recognition could be formulated as a sequential decision problem [39], [160]. By allowing early exiting in these approaches, *temporally* adaptive inference procedure could also be enabled.

For text data, reducing its intrinsic temporal redundancy has attracted great research interests, and the inference paradigm of *temporal-wise* dynamic RNNs (see Sec. 4.1) is also general enough to process audios [262]. Based on large language models such as Transformer [6] and BERT [7], adaptive depths [57], [58], [59], [60] are extensively studied to reduce redundant computation in network architectures.

For video-related tasks, the three types of dynamic inference can be implemented simultaneously [160], [197], [234],

[235]. However, for the networks that do not process videos recurrently, e.g. 3D CNNs [263], [264], [265], most of them still follow a static inference scheme. Few researches have been committed to building dynamic 3D CNNs [195], which might be an interesting future research direction.

Moreover, dynamic networks (especially the attention mechanism) have also been applied to dynamically fuse the features from different modalities in some multi-modal learning tasks, e.g. RGB-D image segmentation [212] and image/video captioning [130], [242], [243], [244].

Finally, dynamic networks have also been exploited to tackle some fundamental problems in deep learning. For example, multi-exit models can be used to: 1) alleviate the *overthinking* issue while reducing the overall computation [50], [266]; 2) perform *long-tailed classification* [267] by inducing early exiting in the training stage; and 3) improve the model *robustness* [268]. For another example, the idea of dynamic routing is implemented for: 1) reducing the *training cost* under a multi-task setting [269] and 2) finding the optimal fine-tuning strategy for per example in *transfer learning* [270].

7 CHALLENGES AND FUTURE DIRECTIONS

Though recent years have witnessed significant progress in the research of dynamic neural networks, there still exist many open problems that are worth exploring. In this section, we summarize a few challenges together with possible future directions in this field.

7.1 Theories for Dynamic Networks

Despite the success of dynamic neural networks, relatively few researches has been committed to analyze them from the theoretical perspective. In fact, theories for a deep understanding of current dynamic learning models and further improving them in principled ways are highly valuable. Notably, it has been proven that a dynamic network with an adaptive width can preserve the representation power of an unsparsified model [79]. However, there are more theoretical problems that are fundamental for dynamic networks. Here we list several of them as follows.

1) Optimal decision in dynamic networks. An essential operation in most dynamic networks (especially those designed for improving computational efficiency) is making data-dependent decisions, e.g., determining whether a module should be evaluated or skipped. Existing solutions either

use confidence-based criteria, or introduce policy networks and gating functions. Although being effective in practice (as mentioned in Sec. 5), they may not be optimal and lack theoretical justifications. Take early exiting as an example, the current heuristic methods [12], [32] might face the issues of overconfidence, high sensitivity for threshold setting and poor transferability. As for policy networks or gating modules, runtime decisions can be made based on a learned function. However, they often introduce extra computations, and usually require a long and unstable training procedure. Therefore, principled approaches with theoretical guarantees for decision function design in dynamic networks is a valuable research topic.

2) Generalization issues. In a dynamic model, a sub-network might be activated for a set of test samples that are not uniformly sampled from the data distribution, e.g., smaller sub-networks tend to handle “easy” samples, while larger sub-networks are used for “hard” inputs [12]. This brings a divergence between the training data distribution and that of the inference stage, and thus violates the common *i.i.d.* assumption in classical machine learning. Therefore, it would be interesting to develop new theories to analyze the generalization properties of dynamic networks under such distribution mismatch. Note that transfer learning also aims to address the issue of distributional shift at test time, but the samples of the target domain are assumed to be accessible in advance. In contrast, for dynamic models, the test distribution is not available until the training process is finished, when the network architecture and parameters are finalized. This poses greater challenges than analyzing the generalization issues in transfer learning.

7.2 Architecture Design for Dynamic Networks

Architecture design has been proven to be essential for deep networks. Existing researches on architectural innovations are mainly proposed for static models [4], [5], [27], while relatively few are dedicated to developing architectures specially for dynamic networks. It is expected that architectures developed specifically for dynamic networks may further improve their effectiveness and efficiency. For example, the interference among multiple classifiers in an early-exiting network could be mitigated by a carefully designed multi-scale architecture with dense connections [12].

Possible research direction include designing dynamic network structures either by hand (as in [12], [32], [35], [67]), or by leveraging the NAS techniques (as in [83], [106]). Moreover, considering the popularity of Transformers [138], recent work has proposed dynamic vision Transformers with adaptive early exiting [271] or token sparsification [272], [273]. Developing a dynamic version of this family of models could also be an interesting direction.

Note that the research on dynamic networks differs from a seemingly close topic, i.e. model compression [28], [29], [31]. One common goal of them is improving the network efficiency with minimal accuracy drop. However, model compression may focus on reducing the *size* of deep networks, while dynamic networks pay more attention to the *computation*, even at the price of slightly *increasing* model size [15], [47]. Moreover, model compression typically adopts pruning [28] or quantization [29] techniques to

produce compact *static* models, which treat all the inputs in the same way. In contrast, dynamic networks perform data-dependent computation on different inputs, which can effectively reduce the intrinsic redundancy in static models.

7.3 Applicability for More Diverse Tasks

Many existing dynamic networks (e.g., most of the sample-wise adaptive networks) are designed specially for classification tasks, and cannot be applied to other vision tasks such as object detection and semantic segmentation. The difficulty arises from the fact that for these tasks there is no simple criterion to assert whether an input image is easy or hard, as it usually contains multiple objects and pixels that have different levels of difficulty. Although many efforts, e.g., spatially adaptive models [33], [39], [148] and soft attention based models [13], [20], [21], have been made to address this issue, it remains a challenging problem to develop a unified and elegant dynamic network that can serve as an off-the-shelf backbone for a variety of tasks.

7.4 Gap between Theoretical & Practical Efficiency

The current deep learning hardware and libraries are mostly optimized for static models, and they may not be friendly to dynamic networks. Therefore, we usually observe that the practical runtime of dynamic models lags behind the theoretical efficiency. For example, some spatially adaptive networks involve sparse computation, which is known to be inefficient on modern computing devices due to the memory access bottleneck [148]. A recent line of work focuses on the codesign of algorithm and hardware for accelerating deep models on platforms with more flexibility such as FPGA [274]. Many input-dependent operations, including pixel-level dynamic computation [114], [275], [276], adaptive channel pruning [277], [278] and early exiting [279], have also been tailored together with hardware for further improving their practical efficiency. It is an interesting research direction to simultaneously optimize the algorithm, hardware and deep learning libraries to harvest the theoretical efficiency gains of dynamic networks.

In addition, a data-dependent inference procedure, especially for the dynamic *architectures*, usually requires a model to handle input samples sequentially, which also poses challenge for parallel computation. Although inference with batches has been enabled for early-exiting networks [271], the conflict between adaptive computational graph and parallel computation still exists for other types of dynamic architectures. This issue is mitigated in the scenario of mobile/edge computing, where the input signal by itself is sequential and the computing hardware is less powerful than high-end platforms. However, designing dynamic networks that are more compatible with existing hardware and software is still a valuable and challenging topic.

7.5 Robustness Against Adversarial Attack

Dynamic models may provide new perspectives for the research on adversarial robustness of deep neural networks. For example, recent work [268] has leveraged the multi-exit structure to improve the robustness against adversarial attacks. Moreover, traditional attacks are usually aimed at causing *misclassification*. For dynamic networks, it is possible to launch attacks on *efficiency* [280], [281]. Specifically, by

adjusting the objective function of the adversarial attack, input-adaptive models could be fooled to activate all their intermediate layers [280] or yielding confusing predictions at early exits [281] even for “easy” samples. It has also been observed that the commonly used adversarial training is not effective to defend such attacks. The robustness of dynamic network is an interesting yet understudied topic.

7.6 Interpretability

Dynamic networks inherit the black-box nature of deep learning models, and thus also invite research on interpreting their working mechanism. What is special here is that the adaptive inference paradigm, e.g., spatial/temporal adaptiveness, conforms well with that of the human visual system, and may provide new possibilities for making the model more transparent to humans. In a dynamic network, it is usually convenient to analyze which part of the model is activated for a given input or to locate which part of the input the model mostly relies on in making its prediction. It is expected that the research on dynamic network will inspire new work on the interpretability of deep learning.

ACKNOWLEDGMENTS

This work is supported in part by the National Science and Technology Major Project of the Ministry of Science and Technology of China under Grants 2018AAA0100701, the National Natural Science Foundation of China under Grants 61906106 and 62022048, the Institute for Guo Qiang of Tsinghua University and Beijing Academy of Artificial Intelligence.

REFERENCES

- [1] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *NeurIPS*, 2012.
- [2] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015.
- [3] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *CVPR*, 2015.
- [4] Kaiping He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [5] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *CVPR*, 2017.
- [6] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017.
- [7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *ACL*, 2019.
- [8] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. In *NeurIPS*, 2020.
- [9] Barret Zoph and Quoc V Le. Neural architecture search with reinforcement learning. In *ICLR*, 2017.
- [10] Hanxiao Liu, Karen Simonyan, and Yiming Yang. DARTS: Differentiable Architecture Search. In *ICLR*, 2018.
- [11] Alex Graves. Adaptive computation time for recurrent neural networks. *arXiv preprint arXiv:1603.08983*, 2016.
- [12] Gao Huang, Danlu Chen, Tianhong Li, Felix Wu, Laurens van der Maaten, and Kilian Weinberger. Multi-scale dense networks for resource efficient image classification. In *ICLR*, 2018.
- [13] Brandon Yang, Gabriel Bender, Quoc V Le, and Jiquan Ngiam. Condconv: Conditionally parameterized convolutions for efficient inference. In *NeurIPS*, 2019.
- [14] Sara Sabour, Nicholas Frosst, and Geoffrey E Hinton. Dynamic routing between capsules. In *NeurIPS*, 2017.
- [15] Ji Lin, Yongming Rao, Jiwen Lu, and Jie Zhou. Runtime neural pruning. In *NeurIPS*, 2017.
- [16] Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarczyk, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. In *ICLR*, 2017.
- [17] Luca Bertinetto, João F Henriques, Jack Valmadre, Philip HS Torr, and Andrea Vedaldi. Learning feed-forward one-shot learners. In *NeurIPS*, 2016.
- [18] Xin Wang, Fisher Yu, Ruth Wang, Trevor Darrell, and Joseph E Gonzalez. Tafe-net: Task-aware feature embeddings for low shot learning. In *CVPR*, 2019.
- [19] Yinpeng Chen, Xiyang Dai, Mengchen Liu, Dongdong Chen, Lu Yuan, and Zicheng Liu. Dynamic convolution: Attention over convolution kernels. In *CVPR*, 2020.
- [20] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *CVPR*, 2018.
- [21] Sanghyun Woo, Jongchan Park, Joon-Young Lee, and In So Kweon. Cbam: Convolutional block attention module. In *ECCV*, 2018.
- [22] Jiaolong Yang, Peiran Ren, Dongqing Zhang, Dong Chen, Fang Wen, Hongdong Li, and Gang Hua. Neural aggregation network for video face recognition. In *CVPR*, 2017.
- [23] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- [24] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015.
- [25] Yulin Wang, Xuran Pan, Shiji Song, Hong Zhang, Gao Huang, and Cheng Wu. Implicit semantic data augmentation for deep networks. In *NeurIPS*, 2019.
- [26] Ekin D Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasudevan, and Quoc V Le. Autoaugment: Learning augmentation strategies from data. In *CVPR*, 2019.
- [27] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.
- [28] Gao Huang, Shichen Liu, Laurens Van der Maaten, and Kilian Q Weinberger. Condensenet: An efficient densenet using learned group convolutions. In *CVPR*, 2018.
- [29] Itay Hubara, Matthieu Courbariaux, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. Binarized neural networks. In *NeurIPS*, 2016.
- [30] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. In *NeurIPS Workshop*, 2014.
- [31] Max Jaderberg, Andrea Vedaldi, and Andrew Zisserman. Speeding up convolutional neural networks with low rank expansions. In *BMVC*, 2014.
- [32] Le Yang, Yizeng Han, Xi Chen, Shiji Song, Jifeng Dai, and Gao Huang. Resolution Adaptive Networks for Efficient Inference. In *CVPR*, 2020.
- [33] Michael Figurnov, Maxwell D Collins, Yukun Zhu, Li Zhang, Jonathan Huang, Dmitry Vetrov, and Ruslan Salakhutdinov. Spatially adaptive computation time for residual networks. In *CVPR*, 2017.
- [34] Xiaoxiao Li, Ziwei Liu, Ping Luo, Chen Change Loy, and Xiaoou Tang. Not all pixels are equal: Difficulty-aware semantic segmentation via deep layer cascade. In *CVPR*, 2017.
- [35] Mostafa Dehghani, Stephan Gouws, Oriol Vinyals, Jakob Uszkoreit, and Łukasz Kaiser. Universal Transformers. In *ICLR*, 2019.
- [36] Maha Elbayad, Jiatao Gu, Edouard Grave, and Michael Auli. Depth-Adaptive Transformer. In *ICLR*, 2020.
- [37] David H Hubel and Torsten N Wiesel. Receptive fields, binocular interaction and functional architecture in the cat’s visual cortex. *The Journal of physiology*, 1962.
- [38] Akira Murata, Vittorio Gallesse, Giuseppe Luppino, Masakazu Kaseda, and Hideo Sakata. Selectivity for the shape, size, and orientation of objects for grasping in neurons of monkey parietal area aip. *Journal of neurophysiology*, 2000.
- [39] Yulin Wang, Kangchen Lv, Rui Huang, Shiji Song, Le Yang, and Gao Huang. Glance and focus: a dynamic approach to reducing spatial redundancy in image classification. In *NeurIPS*, 2020.
- [40] Paul Viola and Michael J. Jones. Robust real-time face detection. *IJCV*, 2004.

- [41] Robert A Jacobs, Michael I Jordan, Steven J Nowlan, and Geoffrey E Hinton. Adaptive mixtures of local experts. *Neural computation*, 1991.
- [42] Wolfgang Maass. Networks of spiking neurons: the third generation of neural network models. *Neural networks*, 1997.
- [43] Eugene M Izhikevich. Simple model of spiking neurons. *TNN*, 2003.
- [44] Zhuang Liu, Jianguo Li, Zhiqiang Shen, Gao Huang, Shoumeng Yan, and Changshui Zhang. Learning efficient convolutional networks through network slimming. In *ICCV*, 2017.
- [45] Surat Teerapittayanon, Bradley McDanel, and Hsiang-Tsung Kung. Branchynet: Fast inference via early exiting from deep neural networks. In *ICPR*, 2016.
- [46] Tolga Bolukbasi, Joseph Wang, Ofer Dekel, and Venkatesh Saligrama. Adaptive neural networks for efficient inference. In *ICML*, 2017.
- [47] Xin Wang, Fisher Yu, Zi-Yi Dou, Trevor Darrell, and Joseph E Gonzalez. Skipnet: Learning dynamic routing in convolutional networks. In *ECCV*, 2018.
- [48] Andreas Veit and Serge Belongie. Convolutional networks with adaptive inference graphs. In *ECCV*, 2018.
- [49] Eunhyeok Park, Dongyoung Kim, Soobeom Kim, Yong-Deok Kim, Gunhee Kim, Sungroh Yoon, and Sungjoo Yoo. Big/little deep neural network for ultra low power inference. In *CODES+ISSS*, 2015.
- [50] Xin Wang, Yujia Luo, Daniel Crankshaw, Alexey Tumanov, Fisher Yu, and Joseph E Gonzalez. Idk cascades: Fast deep learning by learning not to overthink. In *AUAI*, 2017.
- [51] Sam Leroux, Steven Bohez, Elias De Coninck, Tim Verbelen, Bert Vankeirsbilck, Pieter Simoons, and Bart Dhoedt. The cascading neural network: building the internet of smart things. *KAIS*, 2017.
- [52] Jiaqi Guan, Yang Liu, Qiang Liu, and Jian Peng. Energy-efficient amortized inference with cascaded deep classifiers. In *IJCAI*, 2018.
- [53] Xin Dai, Xiangnan Kong, and Tian Guo. Epnet: Learning to exit with flexible multi-branch network. In *CIKM*, 2020.
- [54] Mason McGill and Pietro Perona. Deciding how to decide: Dynamic routing in artificial neural networks. In *ICML*, 2017.
- [55] Zequn Jie, Peng Sun, Xin Li, Jiashi Feng, and Wei Liu. Anytime recognition with routing convolutional networks. *TPAMI*, 2019.
- [56] Hao Li, Hong Zhang, Xiaojuan Qi, Ruigang Yang, and Gao Huang. Improved techniques for training adaptive deep networks. In *ICCV*, 2019.
- [57] Weijie Liu, Peng Zhou, Zhiruo Wang, Zhe Zhao, Haotang Deng, and QI JU. FastBERT: a Self-distilling BERT with Adaptive Inference Time. In *ACL*, 2020.
- [58] Ji Xin, Raphael Tang, Jaejun Lee, Yaoliang Yu, and Jimmy Lin. DeeBERT: Dynamic Early Exiting for Accelerating BERT Inference. In *ACL*, 2020.
- [59] Roy Schwartz, Gabriel Stanovsky, Swabha Swayamdipta, Jesse Dodge, and Noah A. Smith. The Right Tool for the Job: Matching Model and Instance Complexities. In *ACL*, 2020.
- [60] Wangchunshu Zhou, Canwen Xu, Tao Ge, Julian McAuley, Ke Xu, and Furu Wei. BERT Loses Patience: Fast and Robust Inference with Early Exit. In *NeurIPS*, 2020.
- [61] Hehe Fan, Zhongwen Xu, Linchao Zhu, Chenggang Yan, Jianjun Ge, and Yi Yang. Watching a small portion could be as good as watching all: Towards efficient video classification. In *IJCAI*, 2018.
- [62] Wenhao Wu, Dongliang He, Xiao Tan, Shifeng Chen, Yi Yang, and Shilei Wen. Dynamic Inference: A New Approach Toward Efficient Video Action Recognition. In *CVPR Workshop*, 2020.
- [63] Yelong Shen, Po-Sen Huang, Jianfeng Gao, and Weizhu Chen. Reasonet: Learning to stop reading in machine comprehension. In *KDD*, 2017.
- [64] Keyi Yu, Yang Liu, Alexander G. Schwing, and Jian Peng. Fast and accurate text classification: Skimming, rereading and early stopping. In *ICLR Workshop*, 2018.
- [65] Xianggen Liu, Lili Mou, Haotian Cui, Zhengdong Lu, and Sen Song. Finding decision jumps in text classification. *Neurocomputing*, 2020.
- [66] Sam Leroux, Pavlo Molchanov, Pieter Simoons, Bart Dhoedt, Thomas Breuel, and Jan Kautz. IamNN: Iterative and Adaptive Mobile Neural Network for Efficient Image Classification. In *ICML Workshop*, 2018.
- [67] Qiushan Guo, Zhipeng Yu, Yichao Wu, Ding Liang, Haoyu Qin, and Junjie Yan. Dynamic recursive neural network. In *CVPR*, 2019.
- [68] Haichao Yu, Haoxiang Li, Honghui Shi, Thomas S Huang, and Gang Hua. Any-precision deep neural networks. In *AAAI*, 2021.
- [69] Qing Jin, Linjie Yang, and Zhenyu Liao. Adabits: Neural network quantization with adaptive bit-widths. In *CVPR*, 2020.
- [70] Jianghao Shen, Yonggan Fu, Yue Wang, Pengfei Xu, Zhangyang Wang, and Yingyan Lin. Fractional skipping: Towards finer-grained dynamic cnn inference. In *AAAI*, 2020.
- [71] Zuxuan Wu, Tushar Nagarajan, Abhishek Kumar, Steven Rennie, Larry S Davis, Kristen Grauman, and Rogerio Feris. Blockdrop: Dynamic inference paths in residual networks. In *CVPR*, 2018.
- [72] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*, 2013.
- [73] Kyunghyun Cho and Yoshua Bengio. Exponentially increasing the capacity-to-computation ratio for conditional computation in deep learning. *arXiv preprint arXiv:1406.7362*, 2014.
- [74] Emmanuel Bengio, Pierre-Luc Bacon, Joelle Pineau, and Doina Precup. Conditional computation in neural networks for faster models. *ICLR Workshop*, 2016.
- [75] Andrew Davis and Itamar Arel. Low-rank approximations for conditional feedforward computation in deep neural networks. *arXiv preprint arXiv:1312.4461*, 2013.
- [76] David Eigen, Marc' Aurelio Ranzato, and Ilya Sutskever. Learning factored representations in a deep mixture of experts. In *ICLR Workshop*, 2013.
- [77] Jiaqi Ma, Zhe Zhao, Xinyang Yi, Jilin Chen, Lichan Hong, and Ed H Chi. Modeling task relationships in multi-task learning with multi-gate mixture-of-experts. In *KDD*, 2018.
- [78] Ravi Teja Mullapudi, William R Mark, Noam Shazeer, and Kayvon Fatahalian. Hydranets: Specialized dynamic architectures for efficient inference. In *CVPR*, 2018.
- [79] Xin Wang, Fisher Yu, Lisa Dunlap, Yi-An Ma, Ruth Wang, Azalia Mirhoseini, Trevor Darrell, and Joseph E. Gonzalez. Deep mixture of experts via shallow embedding. In *UAI*, 2020.
- [80] Shaofeng Cai, Yao Shu, and Wei Wang. Dynamic routing networks. In *WACV*, 2021.
- [81] William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *arXiv preprint arXiv:2101.03961*, 2021.
- [82] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 1997.
- [83] Zhihang Yuan, Bingzhe Wu, Zheng Liang, Shiwan Zhao, Weichen Bi, and Guangyu Sun. S2dnas: Transforming static cnn model for dynamic inference via neural architecture search. In *ECCV*, 2020.
- [84] Weizhe Hua, Yuan Zhou, Christopher M De Sa, Zhiru Zhang, and G Edward Suh. Channel gating neural networks. In *NeurIPS*, 2019.
- [85] Xitong Gao, Yiren Zhao, Łukasz Dudziak, Robert Mullins, and Cheng zhong Xu. Dynamic channel pruning: Feature boosting and suppression. In *ICLR*, 2019.
- [86] Charles Herrmann, Richard Strong Bowen, and Ramin Zabih. Channel selection using gumbel softmax. In *ECCV*, 2020.
- [87] Babak Ehteshami Bejnordi, Tijmen Blankevoort, and Max Welling. Batch-shaping for learning conditional channel gated networks. In *ICLR*, 2020.
- [88] Jinting Chen, Zhaocheng Zhu, Cheng Li, and Yuming Zhao. Self-adaptive network pruning. In *ICONIP*, 2019.
- [89] Changlin Li, Guangrun Wang, Bing Wang, Xiaodan Liang, Zhihui Li, and Xiaojun Chang. Dynamic slimmable network. In *CVPR*, 2021.
- [90] Zhourong Chen, Yang Li, Samy Bengio, and Si Si. You look twice: Gaternet for dynamic filter selection in cnns. In *CVPR*, 2019.
- [91] Chuanjian Liu, Yunhe Wang, Kai Han, Chunjing Xu, and Chang Xu. Learning instance-wise sparsity for accelerating deep models. In *IJCAI*, 2019.
- [92] Yue Wang, Jianghao Shen, Ting-Kuei Hu, Pengfei Xu, Tan Nguyen, Richard G. Baraniuk, Zhangyang Wang, and Yingyan Lin. Dual dynamic inference: Enabling more efficient, adaptive and controllable deep inference. *JSTSP*, 2020.
- [93] Wenhan Xia, Hongxu Yin, Xiaoliang Dai, and Niraj K Jha. Fully dynamic inference with deep neural networks. *IEEE Transactions on Emerging Topics in Computing*, 2021.
- [94] Ali Ehteshami Bejnordi and Ralf Krestel. Dynamic channel and layer gating in convolutional neural networks. In *KI*, 2020.

- [95] Geoffrey E Hinton, Sara Sabour, and Nicholas Frosst. Matrix capsules with EM routing. In *ICLR*, 2018.
- [96] Augustus Odena, Dieterich Lawson, and Christopher Olah. Changing model behavior at test-time using reinforcement learning. In *ICLR Workshop*, 2017.
- [97] Lanlan Liu and Jia Deng. Dynamic deep neural networks: Optimizing accuracy-efficiency trade-offs by selective execution. In *AAAI*, 2018.
- [98] Samuel Rota Buló and Peter Kotschieder. Neural decision forests for semantic image labelling. In *CVPR*, 2014.
- [99] Peter Kotschieder, Madalina Fiterau, Antonio Criminisi, and Samuel Rota Buló. Deep neural decision forests. In *ICCV*, 2015.
- [100] Nicholas Frosst and Geoffrey Hinton. Distilling a neural network into a soft decision tree. *arXiv preprint arXiv:1711.09784*, 2017.
- [101] Thomas M Hehn, Julian FP Kooij, and Fred A Hamprecht. End-to-end learning of decision trees and forests. *IJCV*, 2019.
- [102] Hussein Hazimeh, Natalia Ponomareva, Petros Mol, Zhenyu Tan, and Rahul Mazumder. The tree ensemble layer: Differentiability meets conditional computation. In *ICML*, 2020.
- [103] Zhicheng Yan, Hao Zhang, Robinson Piramuthu, Vignesh Jagadeesh, Dennis DeCoste, Wei Di, and Yizhou Yu. Hd-cnn: hierarchical deep convolutional neural networks for large scale visual recognition. In *ICCV*, 2015.
- [104] Yani Ioannou, Duncan Robertson, Darko Zikic, Peter Kotschieder, Jamie Shotton, Matthew Brown, and Antonio Criminisi. Decision forests, convolutional networks and the models in-between. *arXiv preprint arXiv:1603.01250*, 2016.
- [105] Ryutaro Tanno, Kai Arulkumaran, Daniel Alexander, Antonio Criminisi, and Aditya Nori. Adaptive neural trees. In *ICML*, 2019.
- [106] An-Chieh Cheng, Chieh Hubert Lin, Da-Cheng Juan, Wei Wei, and Min Sun. Instanas: Instance-aware neural architecture search. In *AAAI*, 2020.
- [107] Yanwei Li, Lin Song, Yukang Chen, Zeming Li, Xiangyu Zhang, Xingang Wang, and Jian Sun. Learning Dynamic Routing for Semantic Segmentation. In *CVPR*, 2020.
- [108] Adam W. Harley, Konstantinos G. Derpanis, and Iasonas Kokkinos. Segmentation-aware convolutional networks using local attention masks. In *ICCV*, 2017.
- [109] Hang Su, Varun Jampani, Deqing Sun, Orazio Gallo, Erik Learned-Miller, and Jan Kautz. Pixel-adaptive convolutional neural networks. In *CVPR*, 2019.
- [110] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. Deformable convolutional networks. In *ICCV*, 2017.
- [111] Xizhou Zhu, Han Hu, Stephen Lin, and Jifeng Dai. Deformable convnets v2: More deformable, better results. In *CVPR*, 2019.
- [112] Hang Gao, Xizhou Zhu, Stephen Lin, and Jifeng Dai. Deformable Kernels: Adapting Effective Receptive Fields for Object Deformation. In *ICLR*, 2019.
- [113] Siyuan Shan, Yang Li, and Junier B Oliva. Meta-neighborhoods. *NeurIPS*, 2020.
- [114] Qijing Huang, Dequan Wang, Zhen Dong, Yizhao Gao, Yaohui Cai, Tian Li, Bichen Wu, Kurt Keutzer, and John Wawrzynnek. Codenet: Efficient deployment of input-adaptive object detection on embedded fpgas. In *FPGA*, 2021.
- [115] Misha Denil, Babak Shakibi, Laurent Dinh, Marc’Aurelio Ranzato, and Nando De Freitas. Predicting parameters in deep learning. In *NeurIPS*, 2013.
- [116] Jürgen Schmidhuber. Learning to control fast-weight memories: An alternative to dynamic recurrent networks. *Neural Computation*, 1992.
- [117] Xu Jia, Bert De Brabandere, Tinne Tuytelaars, and Luc V Gool. Dynamic filter networks. In *NeurIPS*, 2016.
- [118] David Ha, Andrew Dai, and Quoc V Le. Hypernetworks. In *ICLR*, 2016.
- [119] Ningning Ma, Xiangyu Zhang, Jiawei Huang, and Jian Sun. WeightNet: Revisiting the Design Space of Weight Networks. In *ECCV*, 2020.
- [120] Irwan Bello. Lambdanetworks: Modeling long-range interactions without attention. In *ICLR*, 2021.
- [121] Martin Simonovsky and Nikos Komodakis. Dynamic Edge-Conditioned Filters in Convolutional Neural Networks on Graphs. In *CVPR*, 2017.
- [122] Di Kang, Debarun Dhar, and Antoni Chan. Incorporating side information by adaptive convolution. In *NeurIPS*, 2017.
- [123] Harm de Vries, Florian Strub, Jérémie Mary, Hugo Larochelle, Olivier Pietquin, and Aaron Courville. Modulating early visual processing by language. In *NeurIPS*, 2017.
- [124] Ethan Perez, Florian Strub, Harm De Vries, Vincent Dumoulin, and Aaron Courville. Film: Visual reasoning with a general conditioning layer. In *AAAI*, 2018.
- [125] Hyunjae Lee, Hyo-Eun Kim, and Hyeonseob Nam. Srm: A style-based recalibration module for convolutional neural networks. In *ICCV*, 2019.
- [126] Qilong Wang, Banggu Wu, Pengfei Zhu, Peihua Li, Wangmeng Zuo, and Qinghua Hu. ECA-net: Efficient channel attention for deep convolutional neural networks. In *CVPR*, 2020.
- [127] Jingda Guo, Xu Ma, Andrew Sansom, Mara McGuire, Andrew Kalaani, Qi Chen, Sihai Tang, Qing Yang, and Song Fu. Spanet: Spatial Pyramid Attention Network for Enhanced Image Recognition. In *ICME*, 2020.
- [128] Fei Wang, Mengqing Jiang, Chen Qian, Shuo Yang, Cheng Li, Honggang Zhang, Xiaogang Wang, and Xiaoou Tang. Residual attention network for image classification. In *CVPR*, 2017.
- [129] Abhijit Guha Roy, Nassir Navab, and Christian Wachinger. Concurrent spatial and channel ‘squeeze & excitation’ in fully convolutional networks. In *MICCAI*, 2018.
- [130] Long Chen, Hanwang Zhang, Jun Xiao, Liqiang Nie, Jian Shao, Wei Liu, and Tat-Seng Chua. Sca-cnn: Spatial and channel-wise attention in convolutional networks for image captioning. In *CVPR*, 2017.
- [131] Jie Hu, Li Shen, Samuel Albanie, Gang Sun, and Andrea Vedaldi. Gather-excite: Exploiting feature context in convolutional neural networks. In *NeurIPS*, 2018.
- [132] Yinpeng Chen, Xiyang Dai, Mengchen Liu, Dongdong Chen, Lu Yuan, and Zicheng Liu. Dynamic relu. In *ECCV*, 2020.
- [133] Ningning Ma, Xiangyu Zhang, and Jian Sun. Funnel activation for visual recognition. In *ECCV*, 2020.
- [134] Xiang Li, Wenhai Wang, Xiaolin Hu, and Jian Yang. Selective kernel networks. In *CVPR*, 2019.
- [135] Shenlong Wang, Linjie Luo, Ning Zhang, and Li-Jia Li. Autoscaler: Scale-attention networks for visual correspondence. In *BMVC*, 2017.
- [136] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *CVPR*, 2018.
- [137] Kaiyu Yue, Ming Sun, Yuchen Yuan, Feng Zhou, Errui Ding, and Fuxin Xu. Compact generalized non-local network. In *NeurIPS*, 2018.
- [138] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021.
- [139] Sneha Chaudhari, Gungor Polatkan, Rohan Ramanath, and Varun Mithal. An attentive survey of attention models. *TIST*, 2021.
- [140] Xizhou Zhu, Dazhi Cheng, Zheng Zhang, Stephen Lin, and Jifeng Dai. An empirical study of spatial attention mechanisms in deep networks. In *ICCV*, 2019.
- [141] Salman Khan, Muzammal Naseer, Munawar Hayat, Syed Waqas Zamir, Fahad Shahbaz Khan, and Mubarak Shah. Transformers in vision: A survey. *arXiv preprint arXiv:2101.01169*, 2021.
- [142] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization. In *CVPR*, 2016.
- [143] Mengye Ren, Andrei Pokrovsky, Bin Yang, and Raquel Urtasun. SBNet: Sparse Blocks Network for Fast Inference. *CVPR*, 2018.
- [144] Xuanyi Dong, Junshi Huang, Yi Yang, and Shuicheng Yan. More is less: A more complicated network with less inference complexity. In *CVPR*, 2017.
- [145] Shijie Cao, Lingxiao Ma, Wencong Xiao, Chen Zhang, Yunxin Liu, Lintao Zhang, Lanshun Nie, and Zhi Yang. Seernet: Predicting convolutional neural network feature-map sparsity through low-bit quantization. In *CVPR*, 2019.
- [146] Shu Kong and Charles Fowlkes. Pixel-wise attentional gating for scene parsing. In *WACV*, 2019.
- [147] Thomas Verelst and Tinne Tuytelaars. Dynamic Convolutions: Exploiting Spatial Sparsity for Faster Inference. In *CVPR*, 2020.
- [148] Zhenda Xie, Zheng Zhang, Xizhou Zhu, Gao Huang, and Stephen Lin. Spatially Adaptive Inference with Stochastic Feature Sampling and Interpolation. In *ECCV*, 2020.

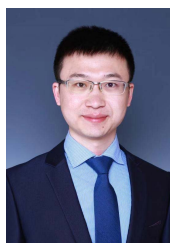
- [149] Amjad Almahairi, Nicolas Ballas, Tim Cooijmans, Yin Zheng, Hugo Larochelle, and Aaron Courville. Dynamic capacity networks. In *ICML*, 2016.
- [150] Alexander Kirillov, Yuxin Wu, Kaiming He, and Ross Girshick. Pointrend: Image segmentation as rendering. In *CVPR*, 2020.
- [151] Aritra Bhowmik, Suprosanna Shit, and Chandra Sekhar Seelamantula. Training-free, single-image super-resolution using a dynamic convolutional network. *IEEE Signal Processing Letters*, 2017.
- [152] Jialin Wu, Dai Li, Yu Yang, Chandrajit Bajaj, and Xiangyang Ji. Dynamic filtering with large sampling field for convnets. In *ECCV*, 2018.
- [153] Xuecai Hu, Haoyuan Mu, Xiangyu Zhang, Zilei Wang, Tieniu Tan, and Jian Sun. Meta-SR: A magnification-arbitrary network for super-resolution. In *CVPR*, 2019.
- [154] Jiaqi Wang, Kai Chen, Rui Xu, Ziwei Liu, Chen Change Loy, and Dahua Lin. CARAFE: Content-Aware ReAssembly of FEatures. In *ICCV*, 2019.
- [155] Jin Chen, Xijun Wang, Zichao Guo, Xiangyu Zhang, and Jian Sun. Dynamic region-aware convolution. In *CVPR*, 2021.
- [156] Guangrun Wang, Keze Wang, and Liang Lin. Adaptively connected neural networks. In *CVPR*, 2019.
- [157] Max Jaderberg, Karen Simonyan, and Andrew Zisserman. Spatial transformer networks. In *NeurIPS*, 2015.
- [158] Adria Recasens, Petr Kellnhofer, Simon Stent, Wojciech Matusik, and Antonio Torralba. Learning to zoom: a saliency-based sampling layer for neural networks. In *ECCV*, 2018.
- [159] Volodymyr Mnih, Nicolas Heess, and Alex Graves. Recurrent models of visual attention. In *NeurIPS*, 2014.
- [160] Zhichao Li, Yi Yang, Xiao Liu, Feng Zhou, Shilei Wen, and Wei Xu. Dynamic computational time for visual attention. In *ICCV Workshop*, 2017.
- [161] Amir Rosenfeld and Shimon Ullman. Visual concept recognition and localization via iterative introspection. In *ACCV*, 2016.
- [162] Jianlong Fu, Heliang Zheng, and Tao Mei. Look closer to see better: Recurrent attention convolutional neural network for fine-grained image recognition. In *CVPR*, 2017.
- [163] Jean-Baptiste Cordonnier, Aravindh Mahendran, Alexey Dosovitskiy, Dirk Weissenborn, Jakob Uszkoreit, and Thomas Unterthiner. Differentiable patch selection for image recognition. In *CVPR*, 2021.
- [164] Zekun Hao, Yu Liu, Hongwei Qin, Junjie Yan, Xiu Li, and Xiaolin Hu. Scale-aware face detection. In *CVPR*, 2017.
- [165] Zerui Yang, Yuhui Xu, Wenrui Dai, and Hongkai Xiong. Dynamic-stride-net: deep convolutional neural network with dynamic stride. In *SPIE Optoelectronic Imaging and Multimedia Technology*, 2019.
- [166] Huiyu Wang, Aniruddha Kembhavi, Ali Farhadi, Alan L. Yuille, and Mohammad Rastegari. Elastic: Improving cnns with dynamic scaling policies. In *CVPR*, 2019.
- [167] Víctor Campos, Brendan Jou, Xavier Giró-I-Nieto, Jordi Torres, and Shih Fu Chang. Skip RNN: Learning to skip state updates in recurrent neural networks. In *ICLR*, 2018.
- [168] Christian Hansen, Casper Hansen, Stephen Alstrup, Jakob Grue Simonsen, and Christina Lioma. Neural Speed Reading with Structural-Jump-LSTM. In *ICLR*, 2019.
- [169] Jin Tao, Urmish Thakker, Ganesh Dasika, and Jesse Beu. Skipping RNN State Updates without Retraining the Original Model. In *SenSys-ML*, 2019.
- [170] Yacine Jernite, Edouard Grave, Armand Joulin, and Tomas Mikolov. Variable computation in recurrent neural networks. In *ICLR*, 2017.
- [171] Minjoon Seo, Sewon Min, Ali Farhadi, and Hannaneh Hajishirzi. Neural Speed Reading via Skim-RNN. In *ICLR*, 2018.
- [172] Junyoung Chung, Sungjin Ahn, and Yoshua Bengio. Hierarchical multiscale recurrent neural networks. In *ICLR*, 2017.
- [173] Nan Rosemary Ke, Konrad Złotna, Alessandro Sordoni, Zhouhan Lin, Adam Trischler, Yoshua Bengio, Joelle Pineau, Laurent Charlin, and Christopher Pal. Focused Hierarchical RNNs for Conditional Sequence Processing. In *ICML*, 2018.
- [174] Zhengjie Huang, Zi Ye, Shuangyin Li, and Rong Pan. Length adaptive recurrent model for text classification. In *CIKM*, 2017.
- [175] Adams Wei Yu, Hongrae Lee, and Quoc Le. Learning to Skim Text. In *ACL*, 2017.
- [176] Tsu-Jui Fu and Wei-Yun Ma. Speed Reading: Learning to Read ForBackward via Shuttle. In *EMNLP*, 2018.
- [177] Zuxuan Wu, Caiming Xiong, Yu-Gang Jiang, and Larry S. Davis. Liteval: A coarse-to-fine framework for resource efficient video recognition. In *NeurIPS*, 2019.
- [178] Guillaume Vaudaux-Ruth, Adrien Chan-Hon-Tong, and Catherine Achard. Actionspotter: Deep reinforcement learning framework for temporal action spotting in videos. In *ICPR*, 2020.
- [179] Serena Yeung, Olga Russakovsky, Greg Mori, and Li Fei-Fei. End-to-end learning of action detection from frame glimpses in videos. In *CVPR*, 2016.
- [180] Yu-Chuan Su and Kristen Grauman. Leaving some stones unturned: dynamic feature prioritization for activity detection in streaming video. In *ECCV*, 2016.
- [181] Zuxuan Wu, Caiming Xiong, Chih-Yao Ma, Richard Socher, and Larry S. Davis. AdaFrame: Adaptive Frame Selection for Fast Video Recognition. In *CVPR*, 2019.
- [182] Yue Meng, Rameswar Panda, Chung-Ching Lin, Prasanna Sattigeri, Leonid Karlinsky, Kate Saenko, Aude Oliva, and Rogerio Feris. Adafuse: Adaptive temporal fusion network for efficient action recognition. In *ICLR*, 2021.
- [183] Ximeng Sun, Rameswar Panda, Chun-Fu Chen, Aude Oliva, Rogerio Feris, and Kate Saenko. Dynamic network quantization for efficient video inference. In *ICCV*, 2021.
- [184] Zejia Weng, Zuxuan Wu, Hengduo Li, and Yu-Gang Jiang. Hms: Hierarchical modality selection for efficient video recognition. *arXiv preprint arXiv:2104.09760*, 2021.
- [185] Rameswar Panda, Chun-Fu Chen, Quanfu Fan, Ximeng Sun, Kate Saenko, Aude Oliva, and Rogerio Feris. Adamml: Adaptive multi-modal learning for efficient video recognition. *arXiv preprint arXiv:2105.05165*, 2021.
- [186] Amir Ghodrati, Babak Ehteshami Bejnordi, and Amirhossein Habibi. Frameexit: Conditional early exiting for efficient video recognition. In *CVPR*, 2021.
- [187] Humam Alwassel, Fabian Caba Heilbron, and Bernard Ghanem. Action search: Spotting actions in videos and its application to temporal action localization. In *ECCV*, 2018.
- [188] Yongming Rao, Jiwen Lu, and Jie Zhou. Attention-aware deep reinforcement learning for video face recognition. In *ICCV*, 2017.
- [189] Yansong Tang, Yi Tian, Jiwen Lu, Peiyang Li, and Jie Zhou. Deep Progressive Reinforcement Learning for Skeleton-Based Action Recognition. In *CVPR*, 2018.
- [190] Wenhao Wu, Dongliang He, Xiao Tan, Shifeng Chen, and Shilei Wen. Multi-agent reinforcement learning based frame sampling for effective untrimmed video recognition. In *ICCV*, 2019.
- [191] Bruno Korbar, Du Tran, and Lorenzo Torresani. Scsampler: Sampling salient clips from video for efficient action recognition. In *ICCV*, 2019.
- [192] Yin-Dong Zheng, Zhaoyang Liu, Tong Lu, and Limin Wang. Dynamic Sampling Networks for Efficient Action Recognition in Videos. *TIP*, 2020.
- [193] Kai Han, Yunhe Wang, Qiulin Zhang, Wei Zhang, Chunjing Xu, and Tong Zhang. Model rubik's cube: Twisting resolution, depth and width for tinynets. *NeurIPS*, 2020.
- [194] Linxi Fan, Shyamal Buch, Guanzhi Wang, Ryan Cao, Yuke Zhu, Juan Carlos Nieves, and Li Fei-Fei. Rubiksnet: Learnable 3d-shift for efficient video action recognition. In *ECCV*, 2020.
- [195] Hengduo Li, Zuxuan Wu, Abhinav Shrivastava, and Larry S. Davis. 2d or not 2d? adaptive 3d convolution selection for efficient video recognition. In *CVPR*, 2021.
- [196] Yue Meng, Chung-Ching Lin, Rameswar Panda, Prasanna Sattigeri, Leonid Karlinsky, Aude Oliva, Kate Saenko, and Rogerio Feris. Ar-net: Adaptive frame resolution for efficient action recognition. In *ECCV*, 2020.
- [197] Yulin Wang, Zhaoxi Chen, Haojun Jiang, Shiji Song, Yizeng Han, and Gao Huang. Adaptive focus for efficient video recognition. *ICCV*, 2021.
- [198] Bowen Pan, Rameswar Panda, Camilo Fosco, Chung-Ching Lin, Alex Andonian, Yue Meng, Kate Saenko, Aude Oliva, and Rogerio Feris. Va-red²: Video adaptive redundancy reduction. In *ICLR*, 2021.
- [199] Mohsen Fayyaz, Emad Bahrami, Ali Diba, Mehdi Noroozi, Ehsan Adeli, Luc Van Gool, and Jurgen Gall. 3d cnns with adaptive temporal feature resolutions. In *CVPR*, 2021.
- [200] Thomas Verelst and Tinne Tuytelaars. Blockcopy: High-resolution video processing with block-sparse feature propagation and online policies. In *ICCV*, 2021.
- [201] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. In *ICML*, 2017.

- [202] Matthias Hein, Maksym Andriushchenko, and Julian Bitterwolf. Why relu networks yield high-confidence predictions far away from the training data and how to mitigate the problem. In *CVPR*, 2019.
- [203] Henry A Rowley, Shumeet Baluja, and Takeo Kanade. Neural network-based face detection. *TPAMI*, 1998.
- [204] Haoxiang Li, Zhe Lin, Xiaohui Shen, Jonathan Brandt, and Gang Hua. A convolutional neural network cascade for face detection. In *CVPR*, 2015.
- [205] Yi Sun, Xiaogang Wang, and Xiaoou Tang. Deep convolutional network cascade for facial point detection. In *CVPR*, 2013.
- [206] Anelia Angelova, Alex Krizhevsky, Vincent Vanhoucke, Abhijit Ogale, and Dave Ferguson. Real-Time Pedestrian Detection with Deep Network Cascades. In *BMVC*, 2015.
- [207] Fan Yang, Wongun Choi, and Yuanqing Lin. Exploit all the layers: Fast and accurate cnn object detector with scale dependent pooling and cascaded rejection classifiers. In *CVPR*, 2016.
- [208] Hong-Yu Zhou, Bin-Bin Gao, and Jianxin Wu. Adaptive feeding: Achieving fast and accurate detections by adaptively combining object detectors. In *ICCV*, 2017.
- [209] Tong Yang, Xiangyu Zhang, Zeming Li, Wenqiang Zhang, and Jian Sun. Metaanchor: Learning to detect objects with customized anchors. In *NeurIPS*, 2018.
- [210] Chunlin Chen and Qiang Ling. Adaptive Convolution for Object Detection. *IEEE Transactions on Multimedia*, 2019.
- [211] Hiroki Tokunaga, Yuki Teramoto, Akihiko Yoshizawa, and Ryoma Bise. Adaptive weighting multi-field-of-view cnn for semantic segmentation in pathology. In *CVPR*, 2019.
- [212] Yikai Wang, Wenbing Huang, Fuchun Sun, Tingyang Xu, Yu Rong, and Junzhou Huang. Deep multimodal fusion by channel exchanging. In *NeurIPS*, 2020.
- [213] Gernot Riegler, Samuel Schuster, Matthias Ruther, and Horst Bischof. Conditioned regression models for non-blind single image super-resolution. In *ICCV*, 2015.
- [214] Falong Shen, Shuicheng Yan, and Gang Zeng. Neural style transfer via meta networks. In *CVPR*, 2018.
- [215] Yu-Gang Jiang, Changmao Cheng, Hangyu Lin, and Yanwei Fu. Learning layer-skippable inference network. *TIP*, 2020.
- [216] Junjun He, Zhongying Deng, and Yu Qiao. Dynamic multi-scale filters for semantic segmentation. In *ICCV*, 2019.
- [217] Dmitrii Marin, Zijian He, Peter Vajda, Priyam Chatterjee, Sam Tsai, Fei Yang, and Yuri Boykov. Efficient segmentation: Learning downsampling near semantic boundaries. In *ICCV*, 2019.
- [218] Jun Li, Yongjun Chen, Lei Cai, Ian Davidson, and Shuiwang Ji. Dense transformer networks for brain electron microscopy image segmentation. In *IJCAI*, 2019.
- [219] Fei Wu, Feng Chen, Xiao-Yuan Jing, Chang-Hui Hu, Qi Ge, and Yimu Ji. Dynamic attention network for semantic segmentation. *Neurocomputing*, 2020.
- [220] Zilong Zhong, Zhong Qiu Lin, Rene Bidart, Xiaodan Hu, Ibrahim Ben Daya, Zhifeng Li, Wei-Shi Zheng, Jonathan Li, and Alexander Wong. Squeeze-and-Attention Networks for Semantic Segmentation. In *CVPR*, 2020.
- [221] Xun Huang, Ming-Yu Liu, Serge Belongie, and Jan Kautz. Multimodal unsupervised image-to-image translation. In *ECCV*, 2018.
- [222] Xihui Liu, Guojun Yin, Jing Shao, Xiaogang Wang, and hongsheng Li. Learning to Predict Layout-to-image Conditional Convolutions for Semantic Image Synthesis. In *NeurIPS*, 2019.
- [223] Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. Semantic image synthesis with spatially-adaptive normalization. In *CVPR*, 2019.
- [224] Peihao Zhu, Rameen Abdal, Yipeng Qin, and Peter Wonka. SEAN: Image Synthesis with Semantic Region-Adaptive Normalization. In *CVPR*, 2020.
- [225] Meng Chang, Qi Li, Huajun Feng, and Zhihai Xu. Spatial-adaptive network for single image denoising. In *ECCV*, 2020.
- [226] Tianjun Xiao, Yichong Xu, Kuiyuan Yang, Jiaying Zhang, Yuxin Peng, and Zheng Zhang. The application of two-level attention models in deep convolutional neural network for fine-grained image classification. In *CVPR*, 2015.
- [227] Heliang Zheng, Jianlong Fu, Tao Mei, and Jiebo Luo. Learning multi-attention convolutional neural network for fine-grained image recognition. In *ICCV*, 2017.
- [228] Wanjie Sun and Zhenzhong Chen. Learned image downscaling for upscaling using content adaptive resampler. *TIP*, 2020.
- [229] Jimmy Ba, Volodymyr Mnih, and Koray Kavukcuoglu. Multiple object recognition with visual attention. In *ICLR*, 2015.
- [230] SM Ali Eslami, Nicolas Heess, Theophane Weber, Yuval Tassa, David Szepesvari, and Geoffrey E. Hinton. Attend, infer, repeat: Fast scene understanding with generative models. In *NeurIPS*, 2016.
- [231] Ali Diba, Vivek Sharma, Luc Van Gool, and Rainer Stiefelhausen. Dynamonet: Dynamic action and motion network. In *ICCV*, 2019.
- [232] Ruohan Gao, Tae-Hyun Oh, Kristen Grauman, and Lorenzo Torresani. Listen to look: Action recognition by previewing audio. In *CVPR*, 2020.
- [233] Yu-Syuan Xu, Tsu-Jui Fu, Hsuan-Kung Yang, and Chun-Yi Lee. Dynamic video segmentation network. In *CVPR*, 2018.
- [234] Simon Niklaus, Long Mai, and Feng Liu. Video frame interpolation via adaptive separable convolution. In *ICCV*, 2017.
- [235] Simon Niklaus, Long Mai, and Feng Liu. Video frame interpolation via adaptive convolution. In *CVPR*, 2017.
- [236] Younghyun Jo, Seoung Wug Oh, Jaeyeon Kang, and Seon Joo Kim. Deep video super-resolution network using dynamic upsampling filters without explicit motion compensation. In *CVPR*, 2018.
- [237] Tae Hyun Kim, Kyoung Mu Lee, Bernhard Scholkopf, and Michael Hirsch. Online video deblurring via dynamic temporal blending network. In *CVPR*, 2017.
- [238] Shangchen Zhou, Jiawei Zhang, Jinshan Pan, Haozhe Xie, Wangmeng Zuo, and Jimmy Ren. Spatio-temporal filter adaptive network for video deblurring. In *ICCV*, 2019.
- [239] Lei Chen, Jiwen Lu, Zhanjie Song, and Jie Zhou. Part-activated deep reinforcement learning for action prediction. In *ECCV*, 2018.
- [240] Hugues Thomas, Charles R. Qi, Jean-Emmanuel Deschaud, Beatriz Marcotegeui, François Goulette, and Leonidas J. Guibas. Kpconv: Flexible and deformable convolution for point clouds. In *ICCV*, 2019.
- [241] Jie Li, Kai Han, Peng Wang, Yu Liu, and Xia Yuan. Anisotropic convolutional networks for 3d semantic scene completion. In *CVPR*, 2020.
- [242] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *ICML*, 2015.
- [243] Chiori Hori, Takaaki Hori, Teng-Yok Lee, Ziming Zhang, Bret Harsham, John R Hershey, Tim K Marks, and Kazuhiko Sumi. Attention-based multimodal fusion for video description. In *ICCV*, 2017.
- [244] Chen Sun, Austin Myers, Carl Vondrick, Kevin Murphy, and Cordelia Schmid. Videobert: A joint model for video and language representation learning. In *ICCV*, 2019.
- [245] Peng Gao, Hongsheng Li, Shuang Li, Pan Lu, Yikang Li, Steven CH Hoi, and Xiaogang Wang. Question-guided hybrid convolution for visual question answering. In *ECCV*, 2018.
- [246] AmirAli Bagher Zadeh, Paul Pu Liang, Soujanya Poria, Erik Cambria, and Louis-Philippe Morency. Multimodal language analysis in the wild: Cmu-mosei dataset and interpretable dynamic fusion graph. In *ACL*, 2018.
- [247] Wasifur Rahman, Md Kamrul Hasan, Sangwu Lee, Amir Zadeh, Chengfeng Mao, Louis-Philippe Morency, and Ehsan Hoque. Integrating multimodal information in large pretrained transformers. In *ACL*, 2020.
- [248] Yagmur Gizem Cinar, Hamid Mirisae, Parantapa Goswami, Eric Gaussier, Ali Ait-Bachir, and Vadim Strijov. Position-based content attention for time series forecasting with sequence-to-sequence rnnns. In *ICONIP*, 2017.
- [249] Chenyou Fan, Yuze Zhang, Yi Pan, Xiaoyue Li, Chi Zhang, Rong Yuan, Di Wu, Wensheng Wang, Jian Pei, and Heng Huang. Multi-horizon time series forecasting with temporal attention learning. In *KDD*, 2019.
- [250] Xiaoyong Jin, Yu-Xiang Wang, and Xifeng Yan. Inter-series attention model for covid-19 forecasting. In *SDM*, 2021.
- [251] Xiaotian Jiang, Quan Wang, and Bin Wang. Adaptive convolution for multi-relational learning. In *NAACL*, 2019.
- [252] Weiping Song, Zhiping Xiao, Yifan Wang, Laurent Charlin, Ming Zhang, and Jian Tang. Session-based social recommendation via dynamic graph attention networks. In *WSDM*, 2019.
- [253] Weiping Song, Chence Shi, Zhiping Xiao, Zhijian Duan, Yewen Xu, Ming Zhang, and Jian Tang. AutoInt: Automatic feature interaction learning via self-attentive neural networks. In *CIKM*, 2019.

- [254] Zhenhua Huang, Xin Xu, Honghao Zhu, and MengChu Zhou. An efficient group recommendation model with multiattention-based neural networks. *IEEE TNNLS*, 2020.
- [255] Giannis Nikolentzos, Antoine Tixier, and Michalis Vazirgiannis. Message passing attention networks for document understanding. In *AAAI*, 2020.
- [256] Gihyeon Choi, Shinhyeok Oh, and Harksoo Kim. Improving document-level sentiment classification using importance of sentences. *Entropy*, 2020.
- [257] Haopeng Zhang and Jiawei Zhang. Text graph transformer for document classification. In *EMNLP*, 2020.
- [258] Sunok Kim, Seungryong Kim, Dongbo Min, and Kwanghoon Sohn. Laf-net: Locally adaptive fusion networks for stereo confidence estimation. In *CVPR*, 2019.
- [259] Emil Julius Gumbel. Statistical theory of extreme values and some practical applications. *NBS Applied Mathematics Series*, 1954.
- [260] Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. In *ICLR*, 2017.
- [261] Łukasz Kaiser and Samy Bengio. Discrete autoencoders for sequence models. *arXiv preprint arXiv:1801.09797*, 2018.
- [262] Raffaele Tavarone and Leonardo Badino. Conditional-Computation-Based Recurrent Neural Networks for Computationally Efficient Acoustic Modelling. In *Interspeech*, 2018.
- [263] Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. Learning spatiotemporal features with 3d convolutional networks. In *ICCV*, 2015.
- [264] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *CVPR*, 2017.
- [265] Dongliang He, Zhichao Zhou, Chuang Gan, Fu Li, Xiao Liu, Yandong Li, Limin Wang, and Shilei Wen. Stnet: Local and global spatial-temporal modeling for action recognition. In *AAAI*, 2019.
- [266] Yigitcan Kaya, Sanghyun Hong, and Tudor Dumitras. Shallow-deep networks: Understanding and mitigating network overthinking. In *ICML*, 2019.
- [267] Rahul Duggal, Scott Freitas, Sunny Dhamnani, Duen Horng Jimeng Sun, et al. Elf: An early-exiting framework for long-tailed classification. *arXiv preprint arXiv:2006.11979*, 2020.
- [268] Ting-Kuei Hu, Tianlong Chen, Haotao Wang, and Zhangyang Wang. Triple Wins: Boosting Accuracy, Robustness and Efficiency Together by Enabling Input-Adaptive Inference. In *ICLR*, 2020.
- [269] Clemens Rosenbaum, Tim Klinger, and Matthew Riemer. Routing networks: Adaptive selection of non-linear functions for multi-task learning. In *ICLR*, 2018.
- [270] Yunhui Guo, Honghui Shi, Abhishek Kumar, Kristen Grauman, Tajana Rosing, and Rogerio Feris. Spottune: Transfer learning through adaptive fine-tuning. In *CVPR*, 2019.
- [271] Yulin Wang, Rui Huang, Shiji Song, Zeyi Huang, and Gao Huang. Not all images are worth 16x16 words: Dynamic vision transformers with adaptive sequence length, 2021.
- [272] Yongming Rao, Wenliang Zhao, Benlin Liu, Jiwen Lu, Jie Zhou, and Cho-Jui Hsieh. Dynamicvit: Efficient vision transformers with dynamic token sparsification. *arXiv preprint arXiv:2106.02034*, 2021.
- [273] Bowen Pan, Yifan Jiang, Rameswar Panda, Zhangyang Wang, Rogerio Feris, and Aude Oliva. Ia-red²: Interpretability-aware redundancy reduction for vision transformers. *arXiv preprint arXiv:2106.12620*, 2021.
- [274] Yifan Yang, Qijing Huang, Bichen Wu, Tianjun Zhang, Liang Ma, Giulio Gambardella, Michaela Blott, Luciano Lavagno, Kees Vissers, John Wawrzynek, et al. Syntegy: Algorithm-hardware co-design for convnet accelerators on embedded fpgas. In *FPGA*, 2019.
- [275] Jorge Albericio, Patrick Judd, Tayler Hetherington, Tor Aamodt, Natalie Enright Jerger, and Andreas Moshovos. Cnvlutin: Ineffectual-Neuron-Free Deep Neural Network Computing. In *ISCA*, 2016.
- [276] Yingyan Lin, Charbel Sakr, Yongjune Kim, and Naresh Shanbhag. Predictivenet: An energy-efficient convolutional neural network via zero prediction. In *ISCA*, 2017.
- [277] Vahideh Akhlaghi, Amir Yazdanbakhsh, Kambiz Samadi, Rajesh K Gupta, and Hadi Esmaeilzadeh. Snapea: Predictive early activation for reducing computation in deep convolutional neural networks. In *ISCA*, 2018.
- [278] Weizhe Hua, Yuan Zhou, Christopher De Sa, Zhiru Zhang, and G Edward Suh. Boosting the performance of cnn accelerators with dynamic fine-grained channel gating. In *MICRO*, 2019.
- [279] Debdeep Paul, Jawar Singh, and Jimson Mathew. Hardware-software co-design approach for deep learning inference. In *ICSCC*, 2019.
- [280] Mirazul Haque, Anki Chauhan, Cong Liu, and Wei Yang. Ilfo: Adversarial attack on adaptive neural networks. In *CVPR*, 2020.
- [281] Sanghyun Hong, Yigitcan Kaya, Ionuț-Vlad Modoranu, and Tudor Dumitras. A panda? no, it's a sloth: Slowdown attacks on adaptive multi-exit neural network inference. In *ICLR*, 2021.



Yizeng Han received the B.S. degree from the Department of Automation, Tsinghua University, Beijing, China, in 2018. And he is currently pursuing the Ph.D. degree in control science and engineering with the Department of Automation, Institute of System Integration in Tsinghua University. His current research interests include computer vision and deep learning, especially in dynamic neural networks.



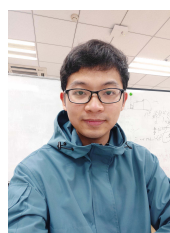
Gao Huang received the B.S. degree from the School of Automation Science and Electrical Engineering, Beihang University, Beijing, China, in 2009, and the Ph.D. degree from the Department of Automation, Tsinghua University, Beijing, China, in 2015. He was a Post-Doctoral Researcher with the Department of Computer Science, Cornell University, Ithaca, USA from 2015 to 2018. He is currently an assistant professor at the Department of Automation, Tsinghua University. His research interests include machine learning and computer vision.



Shiji Song (SM'17) received the Ph.D. degree in mathematics from the Department of Mathematics, Harbin Institute of Technology, Harbin, China, in 1996. He is currently a Professor with the Department of Automation, Tsinghua University, Beijing, China. He has authored over 180 research papers. His current research interests include pattern recognition, system modeling, optimization and control.



Le Yang received his B.S. degrees in the Department of Automation, Northwestern Polytechnical University in 2015, and the Ph.D. degree from the Department of Automation, Tsinghua University, Beijing, in 2021. He is currently an assistant professor at the School of Information and Communications Engineering, Xi'an Jiaotong University. His research interests mainly include machine learning and computer vision, especially in efficient deep learning and dynamic neural networks.



Honghui Wang received his B.S. degree from the Department of Automation, Tsinghua University, Beijing, China, in 2020. He is currently pursuing the Ph.D. degree in control science and engineering with the Department of Automation, Institute of System Integration in Tsinghua University. His current research interests include computer vision and deep learning.



Yulin Wang received his B.S. degree in Automation from Beihang University, Beijing, China, in 2019. He is currently pursuing the Ph.D. degree in the Department of Automation, Tsinghua University. He was a Visiting Student at UC, Berkeley, Berkeley, CA, USA, in 2018. His research interests include computer vision and deep learning.