
A BRIEF REVIEW OF HYPERNETWORKS IN DEEP LEARNING

Vinod Kumar Chauhan^{1*}, Jiandong Zhou¹, Ping Lu¹, Soheila Molaei¹ and David A. Clifton^{1,2}

¹Institute of Biomedical Engineering, University of Oxford,
Old Road Campus Research Building, Headington, Oxford OX3 7DQ, UK
²Oxford-Suzhou Institute of Advanced Research (OSCAR), Suzhou, China

June 13, 2023

ABSTRACT

Hypernetworks, or hypernets in short, are neural networks that generate weights for another neural network, known as the target network. They have emerged as a powerful deep learning technique that allows for greater flexibility, adaptability, faster training, information sharing, and model compression etc. Hypernets have shown promising results in a variety of deep learning problems, including continual learning, causal inference, transfer learning, weight pruning, uncertainty quantification, zero-shot learning, natural language processing, and reinforcement learning etc. Despite their success across different problem settings, currently, there is no review available to inform the researchers about the developments and help in utilizing hypernets. To fill this gap, we review the progress in hypernets. We present an illustrative example to train deep neural networks using hypernets and propose to categorize hypernets on five criteria that affect the design of hypernets as inputs, outputs, variability of inputs and outputs, and architecture of hypernets. We also review applications of hypernets across different deep learning problem settings. Finally, we discuss the challenges and future directions that remain under-explored in the field of hypernets. We believe that hypernetworks have the potential to revolutionize the field of deep learning. They offer a new way to design and train neural networks, and they have the potential to improve the performance of deep learning models on a variety of tasks. Through this review, we aim to inspire further advancements in deep learning through hypernetworks.

Keywords Hypernetworks · Deep learning · Neural Networks · Parameter generation · Weight generation

1 Introduction

Deep learning (DL) has revolutionized the field of artificial intelligence by enabling remarkable advancements in various domains, including computer vision [12], natural language processing [18], causal learning [11], and reinforcement learning [36] etc. Standard deep neural networks (DNNs) have proven to be powerful tools for learning complex representations from data. However, despite their success, standard DNNs remain restrictive in certain conditions. For example, weights of DNNs are fixed once the training process is completed [66]. This lack of adaptability restricts the flexibility of the models, making them less suitable for scenarios where dynamic adjustments are required [24]. DNNs, generally have a large number of weights and they need large amounts of data to optimize those weights [3]. This could be challenging in situations where large amounts of data are not available, e.g., in domains such as healthcare, the collection of sufficient data for rare diseases can be particularly challenging due to the limited number of patients available per year [69]. In addition to that, due to the large sizes of DNNs, it may not be feasible to train and deploy such models in the resource-constrained setting [35], e.g., deploying deep learning models on embedded systems with low-power processors, limited memory, and restricted energy supply poses constraints on model complexity and inference speed. Finally, uncertainty quantification (UQ) in the DNNs' predictions is essential as it provides a measure of confidence, enabling better decision-making in high-stakes applications [14]. The existing UQ techniques have several limitations, such as the need to train multiple models [1], and UQ is still considered an open problem [32].

*Corresponding author: Vinod Kumar Chauhan (vinod.kumar@eng.ox.ac.uk)

Similarly, domain adaptation, domain generalization, adversarial defence, neural style transfer, and neural architecture search etc. are important problems but far from being solved.

Hypernetworks (or hypernets in short) have emerged as a promising architectural paradigm to enhance the flexibility of DNNs. Hypernets are a class of neural networks that generate the weights/parameters of another neural network called target/main/primary network, where both the networks are trained in an end-to-end differentiable manner [24]. Hypernets complement the existing DNNs and provide a new framework to train DNNs, resulting in a new class of DNNs called as HyperDNNs (please refer to Section 2 for details). In addition to improving the flexibility, expressivity, and generalization of DNNs, hypernets offer several advantages, such as: (a) Parameter efficiency: Hypernets can have a smaller number of learnable weights than the target network, resulting in weight compression [74]. This can be particularly useful when working with limited resources or when dealing with high-dimensional data. (b) Transfer learning/information sharing: Hypernetworks can be used to generate the weights of DNNs for solving related tasks [15, 68], helping to improve performance by sharing knowledge learned from one task to another task. (c) Dynamic architectures: Hypernetworks can be used to generate the weights of a network with a dynamic architecture, where the number of layers or the structure of the network changes during training or inference. This can be particularly useful for tasks where the optimal network structure is not known beforehand [24]. (d) Faster training: Hypernetworks can be trained to generate the weights of a target network, allowing the target network to be trained faster [45]. (e) Uncertainty quantification: Hypernets can train uncertainty-aware DNNs by sampling multiple inputs from the noise distribution [33] or using dropout in the hypernets [15], as it generates multiple sets of weights for the main network, and thereby helping to estimate uncertainty in the model predictions which is vital for safety-critical applications, such as healthcare.

Ha et. al [25] coined the term hypernets (are also referred to as meta-networks or meta-models) and trained target network and hypernet in an end-to-end differentiable way, however, the concept of learnable context-dependent weights were discussed even earlier, such as *fast weights* in [57, 58] and HyperNEAT [62]. Our discussion on hypernets, focuses on neural networks generating weights for the target neural network, due to their popularity, expressiveness, and flexibility [13, 12]. Recently, hypernets have gained significant attention and have produced state-of-the-art results across several deep learning problems, such as ensemble learning [32], multitasking [65], neural architecture search [73], continual learning [68], weight pruning [40], Bayesian neural networks [17], generative models [17], hyperparameter optimization [41], information sharing [15], adversarial defence [63], and reinforcement learning (RL) [53] etc. (please refer to Section 4 for more details).

Despite the success of hypernets across different problem settings, to the best of our knowledge, there is no review of hypernets to guide the researchers. To fill this gap, we provide a brief review of hypernets in deep learning. We illustrate hypernets using an example and differentiate HyperDNNs from DNNs (Section 2). To facilitate better understanding and organization, we propose a systematic categorization of hypernets based on five distinct criteria, resulting in different classifications that consider factors such as input characteristics, output characteristics, variability of inputs and outputs, and the architecture of hypernets (Section 3). Furthermore, we offer a comprehensive overview of the diverse applications of hypernets in deep learning, spanning various problem settings (Section 4). By examining real-world applications, we aim to demonstrate the practical advantages and potential impact of hypernetworks. Finally, we discuss the challenges and future directions of hypernet research. This includes addressing initialization, stability and complexity concerns, as well as exploring avenues for enhancing the theoretical understanding and uncertainty quantification of DNNs etc. By providing a comprehensive review of hypernetworks, this paper aims to serve as a valuable resource for researchers and practitioners in the field. Through this review, we hope to inspire further advancements in deep learning by leveraging the potential of hypernets to develop more flexible, adaptive, and high-performing models.

Contributions: This review paper makes the following key contributions:

- To the best of our knowledge, we present the first review on hypernetworks in deep learning, which have shown impressive results across several deep learning problems.
- We propose to categorize hypernetworks based on five criteria, leading to different classifications of hypernets, such as based on inputs, outputs, variability of inputs and outputs, and architecture of hypernets.
- We present a comprehensive overview of applications of hypernetworks across different problem settings, such as uncertainty quantification, continual learning, causal inference, transfer learning, and federated learning etc., and summarize our review, as per our categorization, in a table (Table 2).
- Finally, we identify the challenges and future directions of hypernetwork research, including initialization, stability, scalability, and efficiency concerns, and the need for theoretical understandings and interpretability of hypernetworks. By highlighting these areas, we aim to inspire further advancements in hypernetworks and provide guidance for researchers interested in addressing these challenges.

The rest of the paper is organized as: Section 2 presents background about the hypernets, and Section 3 proposes categorization of hypernets. Applications of hypernets to different problems are presented in Section 4, and Section 5 discusses challenges and future research directions. Finally, the concluding remarks are discussed in Section 6.

2 Background

In this section, we discuss and differentiate the workings of standard Deep neural networks (DNNs) and DNNs trained with hypernets, referred to as HyperDNNs, using a generic example. Fig. 1 depicts the structural differences as well as gradient flows in DNN and HyperDNN. DNN and HyperDNN solve the same problem using exactly the same DNN architecture at the inference time. However, differences exist in how they are trained, i.e., how the gradients flow and the way optimized weights are obtained, as discussed below in details.

Let’s denote a dataset using $\{X, Y\}$ to solve a general task \mathcal{T} , where X is a matrix of features and Y is a vector of labels, and $x \in X$ denotes one data point and $y \in Y$ is the corresponding label. Let a DNN be denoted as a function $\mathcal{F}(X; \Theta)$, where X denotes the inputs, and Θ represents the weights of the DNN. During the forward pass, inputs $x \in X$ pass through the layers of \mathcal{F} to produce predictions $\hat{y} \in \hat{Y}$, and are used along with true labels $y \in Y$, to calculate an objective function to measure discrepancy between actual values and the values predicted by the model using a loss function $\mathcal{L}(Y, \hat{Y})$. During the backward pass, DNNs typically use back-propagation to propagate the error backward through the layers and help in calculating gradients of \mathcal{L} w.r.t. Θ . Then, optimization algorithms, such as Adam [30], use those gradients to update the weights. At the end of the training, we receive optimized weights Θ that are used in the DNN $\mathcal{F}(X; \Theta)$ to make predictions with the test data for solving task \mathcal{T} . So, in standard DNNs, Θ are learnable weights.

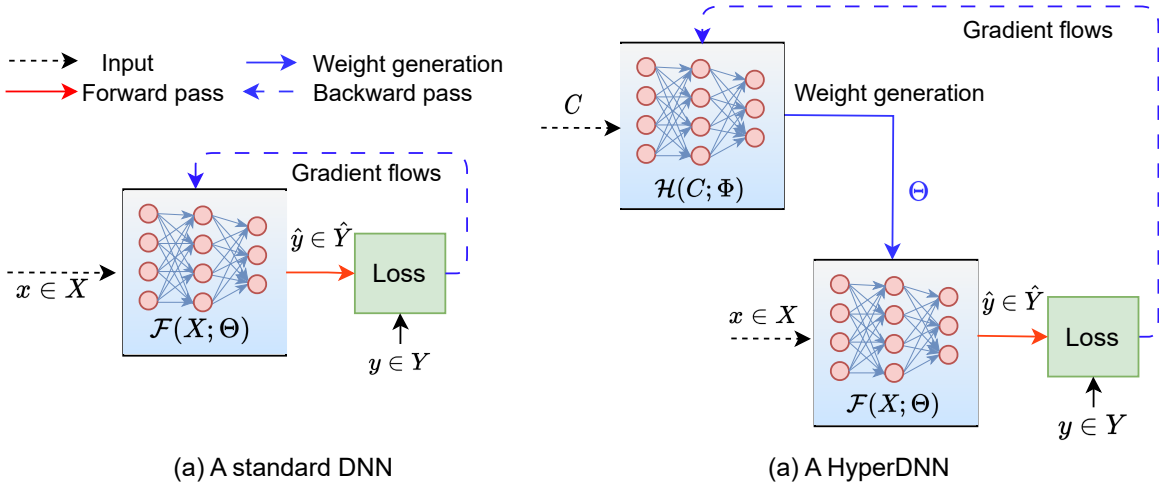


Figure 1: An overview of the architectures and gradient flows for a standard DNN $\mathcal{F}(X; \Theta)$ and the same DNN implemented with hypernets, referred to as HyperDNN $\mathcal{F}(X; \Theta) = \mathcal{F}(X; \mathcal{H}(C; \Phi))$. For DNN, gradients flow through the DNN, and DNN weights Θ are learned during the training. While for HyperDNN, gradients flow through the HyperDNN and hypernet weights Φ are learned during the training that produces DNN weights Θ as outputs.

Hypernets provide an alternative way of learning weights Θ of the DNN $\mathcal{F}(X; \Theta)$ to solve task \mathcal{T} , where Θ are not learned but rather generated by another DNN, i.e., with hypernets we solve the same task using the same DNN but with a different training framework. Let a hypernet be denoted as $\mathcal{H}(C; \Phi)$ which generates the task-specific weights of the DNN $\mathcal{F}(X; \Theta)$, where C is a task-specific context vector that acts as input to \mathcal{H} and Φ are weights of the hypernet \mathcal{H} . That is $\Theta = \mathcal{H}(C; \Phi)$ where Φ are the only learnable weights in the overall architecture. As discussed in detail in the following section, the context vector C can be generated from the data [2], sampled from noise distribution [33] or corresponds to task identity/embedding [4]. During the forward pass, a task-specific context vector C is passed to hypernet \mathcal{H} that generates weights Θ for the DNN \mathcal{F} , and now similar to standard DNN, an input $x \in X$ is passed through the DNN \mathcal{F} to predict the output Y , and loss is calculated as $\mathcal{L}(Y, \hat{Y})$. However, during the backward pass, the error is back-propagated through the hypernet \mathcal{H} and gradients of \mathcal{L} are calculated w.r.t. the weights of hypernet Φ . That is the learning algorithm optimizes Φ to generate Θ so that performance on the target task \mathcal{T} is optimized. At the test time, Θ generated from optimized hypernet \mathcal{H} are used in the DNN $\mathcal{F}(X; \Theta)$ to make predictions with the test data

for solving task \mathcal{T} . The optimization problems for the standard DNN and the HyperDNN can be written as given below (ignoring regularization terms for simplicity).

$$\text{DNN: } \min_{\Theta} \mathcal{F}(X; \Theta), \quad \text{HyperDNN: } \min_{\Phi} \mathcal{F}(X; \Theta) = \mathcal{F}(X; \mathcal{H}(C; \Phi)). \quad (1)$$

Thus, DNNs learn their weights directly from the data but in HyperDNN the weights of the hypernet are learned and the weights of the DNN are generated by hypernet.

3 Categorization of Hypernetworks

In this section, we propose to categorize the hypernetworks based on five criteria that affect the design of hypernets, as depicted in Fig. 2 as: (i) input-based, i.e., what kind of input is taken by the hypernetworks to generate the target neural network weights? (ii) output-based, i.e., how are the outputs, that is, the target weights generated? (iii) variability of inputs, i.e., are the inputs of hypernet fixed? (iv) variability of outputs, i.e., does the target network have a fixed number of weights? and (v) Architecture-based, i.e., what kind of architecture does hypernet use to generate the target weights? We discuss these in following subsections. One can categorize hypernets based on the architecture of the target network but that is not considered because hypernets mostly generate target weights independent of their architecture.

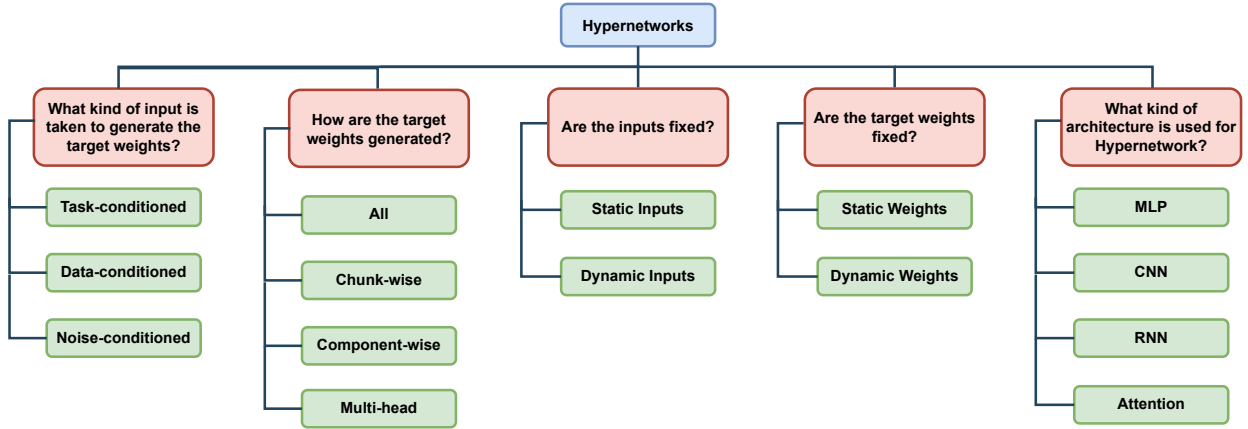


Figure 2: Proposed categorization of hypernets based on five criteria that affect the design of hypernets.

3.1 Input-based Hypernetworks

Hypernetworks take a context vector as an input and generate weights of the target DNN as output. Depending on what context vector is used, we can have the following types of hypernetworks.

Task-conditioned hypernetworks: These hypernetworks take task-specific information as input. The task information can be in the form of task identity/embedding, hyperparameters, architectures, or any other task-specific cues. The hypernetwork generates weights that are tailored to the specific task. This allows the hypernet to adapt its behavior accordingly and allows information sharing among the tasks, resulting in better performance on the tasks. For example [15] applied hypernets to solve treatment effects estimation problem in causal inference that uses an identity or embedding of potential outcome (PO) functions to generate weights corresponding to the PO function. The hypernetworks enabled dynamic end-to-end inter-treatment information sharing among treatment groups and helped to calculate reliable treatment estimates in observational studies with limited-size datasets. Similarly, task-conditioned hypernets have been used to solve other problems, such as multitasking [45], natural language processing (NLP) [24], and continual learning [68] etc.

Data-conditioned hypernetworks: These hypernetworks are conditioned on the data that the target network is being trained on. The hypernetwork generates weights based on the characteristics of the input data. This enables the neural network to dynamically adjust its behavior based on the specific input pattern or features, leading to more flexible and adaptive models, and resulting in better generalization to unseen data. For example, [2] applied hypernets for image editing where the input of hypernet is based on the input images and initial approximation of reconstruction to generate modulations to the weights of the pre-trained generator. Similarly, data-conditioned hypernets have been used to solve other problems, such as adversarial defence [63], knowledge graphs learning [5] and shape learning [39] etc.

Noise-conditioned hypernetworks: These hypernetworks are not conditioned on any input data or task cues, but rather on randomly sampled noise. This makes them more general-purpose and helps in predictive uncertainty quantification for DNNs, but it also means that they may not perform as well as task-conditioned or data-conditioned hypernetworks on multiple tasks or datasets. For example, [33] applied hypernetworks to approximate Bayesian inference in the DNNs and evaluated the approach for active learning, model uncertainty, regularization, and anomaly detection. Similarly, noise-conditioned hypernets have been used to solve other problems, such as manifold learning [17] and uncertainty quantification [52] etc.

These different types of conditioning enable hypernetworks to enhance the flexibility, adaptability, and performance of deep learning models in various contexts. The specific type of hypernetwork that is used will depend on the specific task or application.

3.2 Output-based Hypernetworks

Based on the outputs of hypernets, i.e., weight generation strategy, we can classify hypernetworks according to whether all weights are generated together or not. This classification of hypernetworks is important because it controls the scalability and complexity of the hypernetworks, as typically DNNs have a large number of weights, and producing all of them together can make the size of the last layer of hypernets very large. So, there are ways to manage the complexity of the hypernets that lead to different strategies of weight generation, as discussed below. It is possible to train HyperDNN with a lesser number of weights than the target DNN – this is called weight compression [74].

Table 1: Comparison of different weight generation strategies, i.e., output-based hypernetworks.

Type/ Feature	All	Component-wise	Chunk-wise	Split/Multi-head
Weight generation	Generates all target weights together	Generates target weights for one component at a time	Generates target weights in chunks	Complements all other weight generation strategies so can generate weights as any of the other
Wastage of generated weights	No	Yes, because weights are generated according to the weights of the largest layer	No	Depends on which base strategy it is applied
Are all weights of a layer generated together	Yes	Yes	No	Depends on which base strategy it is applied
The complexity of output space	Largest	Smaller than All	Simplest	Can further improve Chunk-wise
The complexity of input space	Simplest	More complex than All but simpler than Chunk-wise if the number of target layers is lower than the number of chunks	The most complex (assuming the number of chunks is greater than the number of target layers)	Does not have any effect on input space complexity

All Hypernetworks: All hypernetworks generate weights of the entire target DNN altogether. This approach does not lead to wastage of any weights and all the weights of each layer are generated together, unlike the other weight generation strategies. However, this weight generation approach is not suitable for large target networks because that can lead to complex hypernets. For example, [59, 22, 73] used all weight generation.

Chunk-wise Hypernetworks: Chunk-wise hypernetworks generate weights of the target network in chunks. This can lead to the wastage of some of the weights because the weights are generated as per the chunk size, which may not match the layer sizes. If the chunk size is smaller than the layer size, then all the weights of a layer may not be generated together. Moreover, these hypernets need additional embedding to distinguish different chunks and to produce specific weights for the chunks. However, overall chunk-wise hypernets lead to reducing complexity and improving the scalability of hypernets. For example, [15, 68, 15] used chunk-wise weight generation.

Component-wise Hypernetworks: Component-wise hypernetworks generate weights for each individual component (such as layer or channel) of the target model separately. This is helpful in generating specific weights because different layers or channels represent different features or patterns in the network. However, similar to the chunk-wise approach,

component-wise hypernets need an embedding for each component to distinguish among different components and produce weights specific to that component, and also to help to reduce the complexity and improve the scalability of hypernets. Since the weights are generated as per the size of the largest layer so this weight generation approach can lead to the wastage of weights in smaller layers. This strategy can be seen as a special case of a chunk-wise approach, where one chunk is equal to the size of one component. For example, [74, 2, 43] used component-wise weight generation.

Split/Multi-head Hypernetworks: Split/multi-head hypernetworks have multiple heads for producing weights and this weight generation approach can complement the other approaches. This simplifies the complexity and reduces the number of weights required in the last layer of the hypernets by the number of head times. This approach does not need additional embeddings and in general, does not lead to any wastage of weights, unlike component-wise and chunk-wise weight generation approaches. For example, [6, 53, 15] used split/multi-heads to produce target weights.

By classifying hypernetworks based on their weight generation strategy, we can control the scalability and complexity of the hypernetworks to some extent. Each type of weight generation strategy offers unique benefits and considerations based on the specific characteristics and requirements of the task at hand. The comparative study of characteristics of different weight generation approaches is summarized in Table 1.

3.3 Variability of Inputs

We can categorize hypernets based on the variability of the inputs. We have two classes, static inputs, where the inputs are predefined and are fixed, and dynamic inputs, where inputs change and generally are dependent on data on which the target network is trained. This can be seen as a super categorization over input-based hypernets where task-conditioned hypernets fall in the static inputs category while random-noise and data-conditioned hypernets fall in the dynamic category. Both the categories have their own advantages as static inputs help in information sharing [15], transfer learning [68], and are suitable where we have multiple tasks to solve [59]. On the other hand, dynamic inputs help hypernetworks to introduce a new level of adaptability by dynamically generating the weights of the target network. This dynamic weight generation enables hypernetworks to respond to input-dependent context and adjust their behavior accordingly. By generating network weights based on specific inputs, hypernetworks can capture intricate patterns and dependencies that may vary across different instances of data. This adaptability leads to enhanced model performance, especially in scenarios with complex and evolving data distributions [67]. Thus, dynamic input-based hypernets help in domain adaptation [67], density estimation [28] and knowledge graph learning [5] etc.

3.4 Variability of Outputs

When classifying hypernetworks based on the nature of the target network’s weights, we can categorize them into two types, static weights or dynamic weights, depending on whether the weights of the target network are fixed in size at inference time or not. That means for a dynamic weight-based hypernet, the target network is dynamic. These hypernetworks adaptively generate weights, allowing the target network to continuously adjust its behavior as per the inputs. The dynamic weights can be generated, mainly, in two situations, first when the hypernet architecture is dynamic, e.g., [24] used recurrent neural network (RNN) to propose HyperRNN based on non-shared weights. Second, the dynamic weights can be generated when the inputs are dynamic, i.e., hypernet adapts as per the input data, e.g., [39] applied convolutional neural network (CNN) based hypernet to generate dynamic weights for shape learning from an image of a shape. Similarly, [50, 37] also produce dynamic weights, while [49, 64] produce static weights.

3.5 Architecture-based Hypernetworks

In the categorization of hypernetworks based on their architectures, we can classify them into four major types: multi-layer perceptrons (MLPs), convolutional neural networks (CNNs), recurrent neural networks (RNNs), and attention-based networks, which have their standard differences. MLP hypernetworks employ a dense and fully-connected architecture, allowing every input neuron to connect with every output neuron. This architecture enables a comprehensive weight generation process by considering the entire input information, e.g., [15]. CNN hypernetworks, on the other hand, leverage convolutional layers to capture local patterns and spatial information. These hypernetworks excel in tasks involving spatial data, such as an image or video analysis, by extracting features from the input and generating weights or parameters accordingly, e.g., [47]. RNN hypernetworks incorporate recurrent connections in their architecture, facilitating feedback loops and sequential information processing. They dynamically generate weights or parameters based on previous states or inputs, making them well-suited for tasks involving sequential data, such as natural language processing or time series analysis, e.g., [24]. Attention-based hypernetworks incorporate attention mechanisms [66] into their architecture. By selectively focusing on relevant input features, these hypernetworks generate weights for the target network, allowing them to capture long-range dependencies and improve the quality of generated outputs, e.g., [67].

Each type of architecture has its own strengths and applicability, enabling hypernetworks to adapt and generate weights in a manner that aligns with the specific characteristics and demands of the target network and the data being processed.

Table 2: Important applications of hypernetworks and their categorization based on input: (i) task-conditioned, (ii) noise-conditioned, and (iii) data-conditioned; output, i.e., weight generation: (i) all, (ii) component-wise, (iii) chunk-wise, and (iv) split/multi-head; Input variability: (i) static inputs, and (ii) dynamic inputs; Output variability: (i) static weights, and (ii) dynamic weights; and architecture of hypernets (SN: Serial Number, Ref.: Reference, DL: Deep Learning).

SN	Ref.	DL Problem	Input				Output				Input Var.		Out Var.		Architecture	Code
			(i)	(ii)	(iii)	(i)	(ii)	(iii)	(iv)	(i)	(ii)	(i)	(ii)			
1	[24]	Image classification, NLP	✓				✓					✓	✓	✓	RNN, MLP	
2	[33]	Uncertainty quantification		✓		✓						✓	✓		MLP	
3	[63]	Adversarial defence			✓		✓					✓	✓		MLP	
4	[41]	Hyperparameter optimization	✓			✓						✓	✓		MLP	
5	[8]	Neural architecture search	✓					✓				✓	✓		CNN	Link
6	[49]	Spatio-temporal learning			✓	✓						✓	✓		MLP, CNN, RNN,	
7	[73]	Neural architecture search	✓			✓						✓	✓		MLP	
8	[17]	Manifold learning		✓			✓					✓	✓		CNN	
9	[52]	Uncertainty quantification		✓			✓		✓			✓	✓		GAN	
10	[40]	Weight pruning	✓			✓						✓	✓		MLP	Link
11	[5]	Knowledge graphs learning			✓		✓					✓	✓		MLP	Link
12	[39]	Shape learning			✓				✓			✓	✓		CNN	Link
13	[32]	Uncertainty quantification			✓	✓						✓	✓		MLP	
14	[31]	Image processing			✓				✓			✓	✓		CNN	
15	[68]	Continual learning, transfer learning	✓			✓	✓	✓		✓		✓	✓		MLP	Link
16	[74]	Few-shot learning	✓					✓				✓	✓		MLP	
17	[22]	Complexity of NN			✓	✓						✓	✓		MLP	
18	[37]	Weight pruning	✓				✓					✓	✓		MLP	Link
19	[50]	Neural architecture search	✓				✓					✓	✓		CNN	Link
20	[45]	Pareto-Front Learning (multi-tasking, fairness, image segmentation)	✓						✓	✓		✓	✓		MLP	Link
21	[59]	Federated Learning	✓			✓				✓		✓	✓		MLP	Link
22	[47]	Semantic segmentation			✓		✓		✓			✓	✓		CNN	Link
23	[43]	Multitasking, NLP, language model	✓				✓					✓	✓		MLP	Link
24	[56]	Reinforcement learning	✓			✓						✓	✓		MLP	Link
25	[29]	Continual RL	✓			✓						✓	✓		MLP	Link
26	[60]	Density estimation	✓				✓					✓	✓		MLP	
27	[44]	Neural image enhancement			✓	✓						✓	✓		MLP	
28	[27]	Continual learning	✓			✓						✓	✓		MLP	
29	[20]	Continual learning	✓					✓				✓	✓		MLP	Link
30	[34]	Adaptation of neural network architectures			✓	✓						✓	✓		MLP	
31	[46]	Network compression	✓			✓						✓	✓		MLP	
32	[7]	Internal learning (computer vision)			✓				✓			✓	✓		CNN	Link
33	[16]	Learning differential equations	✓			✓						✓	✓		MLP	
34	[2]	Image editing			✓		✓					✓	✓		CNN	Link
35	[67]	Domain adaptation, NLP	✓		✓	✓						✓	✓		Attention	Link
36	[48]	Autonomous driving			✓				✓			✓	✓		Attention, RNN, CNN,	
37	[71]	NLP	✓		✓		✓					✓	✓		MLP	Link
38	[51]	Domain generalization	✓			✓						✓	✓		MLP	
39	[61]	3D point cloud processing			✓	✓						✓	✓		MLP	
40	[19]	Image processing			✓	✓						✓	✓		CNN	Link
41	[72]	Few-shot learning			✓	✓						✓	✓		CNN	Link
42	[15]	Treatment effects estimation, Uncertainty quantification	✓			✓	✓	✓	✓	✓		✓	✓		MLP	
43	[6]	Meta-Reinforcement Learning	✓						✓	✓		✓	✓		MLP	
44	[53]	Zero-shot RL	✓						✓			✓	✓		MLP	Link
45	[64]	Sound representation			✓	✓						✓	✓		CNN	
46	[28]	Density estimation			✓	✓						✓	✓		CNN	
47	[9]	Quantum computing	✓			✓						✓	✓		MLP	Link
48	[55]	Neural style transfer			✓	✓						✓	✓		MLP	
49	[21]	Camera pose localization			✓		✓					✓	✓		Attention, MLP	Link
50	[70]	Knowledge distillation, visualization	✓			✓						✓	✓		MLP	

4 Applications of Hypernetworks

Hypernetworks have demonstrated their effectiveness and versatility across a wide range of domains and tasks in deep learning. In this section, we explore some of the important applications of hypernetworks and highlight their

contributions to advancing the state-of-the-art in these areas. We summarize the applications of hypernets as per our proposed categorization and also provide links to code repositories for the benefit of the researchers, wherever available, in Table 2.

Continual Learning: Continual learning, also known as lifelong learning or incremental learning, is a machine learning paradigm that focuses on the ability of a model to learn and adapt continuously over time, in a sequential manner, without forgetting previously learned knowledge. Unlike traditional batch learning, which assumes static and independent training and testing sets, continual learning deals with dynamic and non-stationary data distributions, where new data arrives incrementally and the model needs to adapt to these changes while retaining previously acquired knowledge. The challenge in continual learning lies in mitigating *catastrophic forgetting*, which refers to the tendency of a model to forget previously learned information when it is trained on new data. To address this, various strategies have been proposed, including regularization techniques, rehearsal methods, dynamic architectures, and parameter isolation. [68] applied task-conditioned hypernets to address the catastrophic forgetting issue and proposed a regularizer for rehearsing task-specific weight realizations rather than the data from previous tasks. They achieved state-of-the-art (SOTA) results on benchmarks and empirically showed that the task-conditioned hypernets have a long capacity to retain memories of previous tasks. Similarly, [29] also applied task-conditioned hypernets to continual learning in reinforcement learning (RL).

Federated Learning: Federated Learning is a decentralized approach to machine learning where the training process is distributed across multiple devices or edge devices, without the need to centralize data in a single location. In this paradigm, each device or edge node locally trains a model using its own data, and only the model updates, rather than the raw data, are shared and aggregated on a central server. This enables collaborative learning while preserving data privacy and security. It also reduces communication costs and latency, making it suitable for scenarios with limited bandwidth or intermittent connectivity. [59] applied hypernets to federated learning problems where a central hypernet is trained to generate the weights for the client models. This allows information sharing across different clients while making the hypernet size independent of communication cost, as hypernet weights are never transmitted. The hypernet-based federated learning achieved the SOTA results and also showed better generalization to new clients whose distributions were different than the existing clients.

Few-shot Learning: Few-shot learning is a sub-field of machine learning that focuses on training models to learn new concepts or tasks with only a limited number of training examples. Unlike traditional machine learning approaches that typically require large amounts of labeled data for each task, few-shot learning aims to generalize knowledge from a small support set of labeled examples to classify or recognize new instances. To address the practical difficulties of existing techniques to operate in high-dimensional parameter spaces with extremely limited-data settings, [54] applied data-conditioned hypernets. They also achieved SOTA results and showed that the proposed technique can capture uncertainty in the data. Similarly, [74] applied task-conditioned hypernets and utilized soft weight sharing for few-shot learning.

Manifold Learning: Manifold learning is a sub-field of machine learning that focuses on capturing the underlying structure or geometry of high-dimensional data in lower-dimensional representations or manifolds. It aims to uncover the intrinsic relationships and patterns within the data by mapping it to a lower-dimensional space, enabling better visualization, clustering, or classification. Hypernetworks can be utilized in the context of manifold learning to enhance the representation learning process. By generating weights or parameters for the target network based on the input, hypernetworks can adaptively learn a manifold that captures the intricate data structure [59]. [17] applied noise-conditioned hypernetworks to map latent vectors for generating target network weights that generalize mode connectivity in loss landscape to higher dimensional manifolds.

AutoML: AutoML, short for Automated Machine Learning, refers to the development of algorithms, systems, and tools that automate various aspects of the machine learning pipeline, e.g., neural architecture search (NAS) and automated hyperparameter optimization. [73] applied hypernetworks for NAS where they modeled neural architectures as graphs and used them as input to generate the target network weights. They achieved about 10 times faster results than the SOTA. [8] also applied hypernets for NAS by proposing a memory-read-write based mechanism. For hyperparameter optimization, [41] applied hypernets that take hyperparameters as input and generate optimal weights for the target network, and hence perform joint training for target network parameters and hyperparameters which are otherwise trained in nested optimization loops. The authors proved the efficacy of the proposed technique against the SOTA to train thousands of hyperparameters.

Pareto-front Learning: Pareto-front learning, also known as multi-objective optimization, is a technique that addresses problems with multiple conflicting objectives, e.g., multitasking has multiple tasks that may have conflicting gradients. It aims to find a set of solutions that represent the trade-off among different objectives, rather than a single optimal solution. In Pareto-front learning, the goal is to identify a set of solutions that cannot be improved in one objective without sacrificing performance in another objective. These solutions are referred to as Pareto-optimal or non-dominated

solutions and lie on the Pareto-front, which represents the best possible trade-off between objectives. [45] applied hypernets to learn the entire Pareto-front, which at inference time takes a preferential point on the Pareto-front and generates Pareto-front weights for the target network whose loss vector is in the direction of the ray. They showed that the proposed hypernets are computationally very efficient as compared with the SOTA and can scale to large models, such as ResNet18.

Reinforcement Learning: Reinforcement Learning (RL) focuses on training agents to make sequential decisions in an environment to maximize a cumulative reward. RL operates through an interaction loop where the agent takes actions, receives feedback in the form of rewards, and learns optimal policies through trial and error. Hypernets can be applied to reinforcement learning to enhance the learning process and improve the performance of RL agents. Hypernets can be used to dynamically generate or adapt network architectures, model parameters, or exploration strategies in RL agents. By using a hypernetwork, the RL agent can effectively learn to customize its internal representations or policies based on the specific characteristics of the environment or task. For example, [56] applied hypernets to generate the building blocks of RL, i.e., policy networks and Q-functions, rather than using MLPs. They showed faster training and improved performance on different algorithms for RL and in meta-RL. Similarly, task-conditioned hypernets were used in RL for generalization across tasks [6], continual RL [29], and zero-shot learning [53].

Domain Adaptation: Domain adaptation refers to the process of adapting a machine learning model trained on a source domain to perform well in a different target domain. It is a crucial challenge in machine learning when there is a shift or discrepancy between the distribution of the source and the target data. Hypernets can play a valuable role in domain adaptation by dynamically generating or adapting model parameters, architectures, or other components to effectively handle domain shifts. For example, [67] were the first to propose hypernets for domain adaptation. They used data-conditioned hypernets where examples from the target domains are used to generate weights for the target network.

Causal Inference: Causal inference is a field of study that focuses on understanding and estimating causal relationships between variables. It aims to uncover the cause-and-effect relationships within a system by leveraging observational or experimental data. Causal inference is particularly important when inferring the impact of treatments/ interventions/ policies on outcomes of interest. Recently, [15] were the first to apply hypernets to heterogeneous treatment effects (HTE) estimation problem where they applied task-conditioned hypernets using embedding of potential outcome (PO) functions to generate factual and counterfactual models. Based on soft weight sharing of hypernets, this work presents the first general mechanism to train HTE learners that enables end-to-end inter-treatment information sharing among the PO functions and helps to get reliable estimates, especially with limited-size observational data.

Natural Language Processing: Natural language processing (NLP) is a sub-field of artificial intelligence that focuses on the interaction between computers and human language. It involves various tasks, such as language generation, sentiment analysis, machine translation, and question answering, among others. In the context of NLP, hypernets can be used to generate or adapt neural network architectures, model hyperparameters, or attention mechanisms, and for transfer learning and domain adaptation. For example, [67] applied data-conditioned hypernet for out-of-distribution (OOD) generalization, and [43] applied task-conditioned hypernets to fine-tune the pre-trained language models by generating weights for the bottleneck adapters.

Computer Vision: Computer vision focuses on enabling computers to understand and interpret visual information from images or videos. Computer vision algorithms aim to replicate human visual perception by detecting and recognizing objects, understanding their spatial relationships, extracting features, and making sense of the visual scene. Some applications of hypernets in computer vision are: [24] first applied task-conditioned hypernets for image classification, [2, 44] applied data-conditioned hypernets, where image acts as input to hypernet, for image enhancement, and [52] applied noise-conditioned hypernets for image classification. Data-conditioned hypernets are also applied to semantic segmentation [47].

Uncertainty Quantification: Uncertainty quantification (UQ) is a critical aspect of deep learning and decision-making that involves estimating and understanding the uncertainty associated with model predictions or outcomes. It provides a measure of confidence or reliability in the predictions made by a model, particularly in situations where the model encounters unseen or uncertain data. Hypernets can be used for UQ as they can generate ensembles of models or adaptively generate weights as per the input data. For example, [33] proposed Bayesian hypernets that take random noise as input to produce distributions over the weights of the target network, and showed competitive performance for uncertainty. [52] also applied noise-conditioned hypernets for UQ and showed that the proposed technique provides a better estimate of uncertainty as compared to the ensemble learning technique. In addition, [15] used dropout in the task-conditioned hypernets to generate multiple sets of weights for the target network and thus helping to estimate uncertainty.

Adversarial Defence: Adversarial defense in deep learning refers to the techniques used to enhance the robustness and resilience of models against adversarial attacks. Adversarial attacks involve making carefully crafted perturbations

to input data in order to deceive or mislead deep learning models [42]. By incorporating hypernetworks, models can enhance their ability to detect and defend against adversarial attacks by dynamically generating or adapting their weights or architectures. For example, [63] generated data-dependent adaptive convolution kernels to improve the robustness of CNNs against adversarial attacks, and were successful in spontaneously detecting attacks generated by Gaussian noise, fast gradient sign methods, and black-box attack methods. [32, 52, 33] also found noise-conditioned hypernets robust to adversarial examples as compared with the SOTA.

Multitasking: Multitasking refers to the capability of a model to perform multiple tasks or learn multiple objectives simultaneously. It involves leveraging shared representations and parameters across different tasks to enhance learning efficiency and overall performance. Hypernets can be applied in the context of multitasking to facilitate the joint learning of multiple tasks by dynamically generating or adapting the model’s parameters or architectures. As already discussed above in overlapping deep learning applications, [43] applied task-conditioned hypernets that share knowledge across the tasks as well as generate task-specific models, and achieved benchmark results. [45] also studied task-conditioned hypernets for Pareto-front learning to address the conflicting gradients among different objectives and obtained impressive results on multitasking, including fairness and image segmentation.

These applications demonstrate the wide-ranging potential of hypernetworks in deep learning, enabling adaptive and task-specific parameter generation for improved model performance and generalization.

5 Challenges and Future Directions

Hypernetworks have shown great potential in enhancing deep learning models with increased flexibility, efficiency, and generalization. However, several challenges and opportunities for future research and development remain under-explored. In this section, we discuss some of the key challenges and propose potential directions for future exploration.

Initialization Challenge: The initialization challenge in hypernetworks refers to the difficulty of initializing the hypernetwork parameters effectively, as finding suitable initial values for the hypernetwork parameters is far from being resolved. One reason for the initialization challenge is that the weights of the target network which are generated at the output layer of hypernet may cause the generation of weights of a layer either separately or together with other layers. Due to this reason, hypernets fail to generate weights in the scale as the classical initialization techniques, such as Xavier [23] and Kaiming initialization [26]. The performance of the hypernetwork is highly influenced by the initial state of the target network and its parameters that are generated at the output layer of the hypernet. If the target network is poorly initialized, it can propagate errors or uncertainties to the hypernetwork, affecting its ability to generate or adapt parameters effectively. [10] were the first to discuss the challenge of initializing hypernets. They showed that classical techniques of initializing DNNs do not work well with hypernets, however, adaptive optimizers, such as Adam [30], can address the issue to some extent. The authors suggested initializing the hypernet weights in a manner that allows the target network weights to approximate the conventional initialization of DNNs. However, it is difficult to adopt this because the weights of the target network are typically generated together. Moreover, recently, [6] also showed that initialization challenge of hypernets occurs even in meta-RL and classical initialization techniques fail.

Complexity/Scalability: One of the primary challenges in hypernetworks is scalability and efficiency of hypernetwork-based models. As the size and complexity of target deep learning models increase, hypernetworks also become very complex, e.g., the size of the output layer is typically $m \times n$ where m is the number of neurons in the penultimate layer of hypernet and n is the number of weights in the target network. So, hypernets may not be suitable for large models unless appropriate weight-generation strategies are developed and used. Although, there are some approaches, such as split-head/multi-head [15] and chunk-wise weight generation [8] to manage the complexity of hypernets but it needs more research to address the scalability challenge and make hypernetworks more practical for real-world applications.

Numerical Stability: Numerical stability in hypernetworks refers to the ability of the model to maintain accurate and reliable computations throughout the training and inference process. Hypernets, like other neural networks, can encounter numerical stability issues [56]. One common numerical stability issue in hypernetworks is the vanishing or exploding gradients problem. During the training process, gradients can become extremely small or large, making it difficult for the model to effectively update the parameters. This can result in slow convergence or unstable training dynamics. To address numerical stability issues in hypernets, various techniques can be employed, such as careful initialization of the model’s parameters, the use of gradient clipping, which bounds the gradient values to prevent them from becoming too large, and different regularization techniques such as weight decay, dropout, and spectral norm [15] that help improve numerical stability by preventing overfitting and promoting smoother optimization. Furthermore, similar to standard DNNs, using appropriate activation functions, such as ReLU or Leaky ReLU, can help alleviate the vanishing gradient problem by providing non-linearities that allow for more effective gradient propagation. It is also

important to choose appropriate optimization algorithms that are known for their stability, such as Adam [30], which can handle the training dynamics of hypernetworks more effectively [10].

Theoretical Understanding: Theoretical analysis of hypernetworks involves studying their representational capacity, learning dynamics, and generalization properties. By understanding the theoretical foundations of hypernetworks, researchers can gain insights into the underlying principles that drive their effectiveness and explore new avenues for improving their performance. Just like DNNs, understanding the working of hypernets is far from being solved. Although, there are some works that provide theoretical insights into hypernets, e.g., [38] highlighted that infinitely wide hypernetworks may not converge to a global minimum using gradient descent, but convexity can be achieved by increasing the dimensionality of the hypernetwork’s output. [22] also studied the modularity of hypernets and showed that hypernets can be more efficient than the embedding-based method for mapping an input to a function. Intuitively, hypernets map an input to one point on a low-dimensional manifold for weights of target network [59] – theoretical insights into the connection between two can be very helpful. Thus, more research into the theoretical properties of hypernets will help to make them more popular and will also attract more research.

Uncertainty-aware Deep Learning: Uncertainty-aware neural networks allow for more reliable and robust predictions, especially in scenarios where uncertainty estimation is crucial, such as decision-making under uncertainty, safety-critical applications, or when working with limited or noisy data [1]. Despite the success of DNNs and the development of different uncertainty quantification techniques, it still remains an open problem to quantify the prediction uncertainty [32]. Hypernets have opened a new door to uncertainty quantification as noise-conditioned hypernets can generate distribution on target network weights and have been shown to have better uncertainties than the SOTA [33, 52]. Similarly, [15] used task-conditioned hypernets with dropout to generate multiple sets of weights for the target network. Further research into this can provide computationally efficient and effective techniques as compared with other techniques, such as ensemble methods, which need to train multiple models.

Interpretability Enhancement: It will be helpful for the community to develop methods for visualizing, analyzing, and explaining the task-specific weights generated by hypernetworks. This includes developing intuitive visualization methods, and feature relevance analysis techniques that provide deeper insights into the decision-making process of hypernetwork-based models.

Model Compression and Efficiency: Hypernetworks can aid in model compression and efficiency in some problem settings [74], where smaller hypernets are trained to generate larger target networks that can reduce the memory footprint and computational requirements of the model. This is particularly useful in resource-constrained environments where memory and computational resources are limited, and hypernets can be studied specifically for such settings.

Usage Guidelines: Hypernetworks add additional complexity to solving problems. As with hypernets, we have two neural networks rather than one in standard DNNs. Hypernets introduce additional hyperparameters related to the weight generation process, e.g., what kind of weight generation should be used and how many chunks should be used. Some research and guidelines are needed to guide the researchers through these choices, stressing the need for a comparative study of different approaches.

Thus, the field of hypernetworks in deep learning presents several challenges and opportunities for future research. The advancements in these areas will pave the way for the widespread adoption and effective utilization of hypernetworks in various domains of deep learning.

6 Conclusion

Hypernetworks have emerged as a promising approach to enhance deep learning models with increased flexibility, efficiency, generalization, uncertainty awareness, and information sharing etc. They have opened up new avenues for research and applications across various domains. In this paper, we have presented the first brief review of hypernetworks in the context of deep learning. We presented an illustrative example to discuss the working of hypernetworks, and proposed categorization of hypernets based on five criteria that influence the design of hypernets based on its inputs, outputs, variability of inputs and outputs, and based on the architecture of hypernets. We discussed some of the important applications of hypernets to different deep learning problems, including multitasking, continual learning, federated learning, causal inference, NLP, RL, computer vision etc. Finally, we also highlighted the challenges that need to be addressed in the future. These challenges include initialization, stability, scalability and efficiency, the need for theoretical insights etc. Future research should focus on tackling these challenges to further advance the field of hypernetworks and make them more accessible and practical for real-world applications.

Acknowledgements

This work was supported in part by the National Institute for Health Research (NIHR) Oxford Biomedical Research Centre (BRC) and in part by InnoHK Project Programme 3.2: Human Intelligence and AI Integration (HIAI) for the Prediction and Intervention of CVDs: Warning System at Hong Kong Centre for Cerebro-cardiovascular Health Engineering (COCHE). DAC was supported by an NIHR Research Professorship, an RAEng Research Chair, the InnoHK Hong Kong Centre for Cerebro-cardiovascular Health Engineering (COCHE), the NIHR Oxford Biomedical Research Centre (BRC), and the Pandemic Sciences Institute at the University of Oxford. The views expressed are those of the authors and not necessarily those of the NHS, the NIHR, the Department of Health, the InnoHK – ITC, or the University of Oxford.

Statements and Declarations

Competing Interests

The authors declare that they have no competing interests.

Data Availability

Data sharing not applicable to this article as no datasets were generated or analysed during the current study.

References

- [1] Abdar, M., Pourpanah, F., Hussain, S., Rezazadegan, D., Liu, L., Ghavamzadeh, M., Fieguth, P., Cao, X., Khosravi, A., Acharya, U. R., et al. (2021). A review of uncertainty quantification in deep learning: Techniques, applications and challenges. *Information Fusion*, 76:243–297.
- [2] Alaluf, Y., Tov, O., Mokady, R., Gal, R., and Bermano, A. (2022). Hyperstyle: Stylegan inversion with hypernetworks for real image editing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18511–18521.
- [3] Alzubaidi, L., Zhang, J., Humaidi, A. J., Al-Dujaili, A., Duan, Y., Al-Shamma, O., Santamaría, J., Fadhel, M. A., Al-Amidie, M., and Farhan, L. (2021). Review of deep learning: Concepts, cnn architectures, challenges, applications, future directions. *Journal of big Data*, 8:1–74.
- [4] Armstrong, J. and Clifton, D. (2021). Continual learning of longitudinal health records. *arXiv preprint arXiv:2112.11944*.
- [5] Balažević, I., Allen, C., and Hospedales, T. M. (2019). Hypernetwork knowledge graph embeddings. In *Artificial Neural Networks and Machine Learning—ICANN 2019: Workshop and Special Sessions: 28th International Conference on Artificial Neural Networks, Munich, Germany, September 17–19, 2019, Proceedings 28*, pages 553–565. Springer.
- [6] Beck, J., Jackson, M. T., Vuorio, R., and Whiteson, S. (2023). Hypernetworks in meta-reinforcement learning. In *Conference on Robot Learning*, pages 1478–1487. PMLR.
- [7] Bensadoun, R., Gur, S., Galanti, T., and Wolf, L. (2021). Meta internal learning. In Ranzato, M., Beygelzimer, A., Dauphin, Y., Liang, P., and Vaughan, J. W., editors, *Advances in Neural Information Processing Systems*, volume 34, pages 20645–20656. Curran Associates, Inc.
- [8] Brock, A., Lim, T., Ritchie, J., and Weston, N. (2018). SMASH: One-shot model architecture search through hypernetworks. In *International Conference on Learning Representations*.
- [9] Carrasquilla, J., Hibat-Allah, M., Inack, E., Makhzani, A., Neklyudov, K., Taylor, G. W., and Torlai, G. (2023). Quantum hypernetworks: Training binary neural networks in quantum superposition. *arXiv preprint arXiv:2301.08292*.
- [10] Chang, O., Flokas, L., and Lipson, H. (2020). Principled weight initialization for hypernetworks. In *International Conference on Learning Representations*.
- [11] Chauhan, V. K., Molaei, S., Tania, M. H., Thakur, A., Zhu, T., and Clifton, D. A. (2023a). Adversarial deconfounding in individualised treatment effects estimation. In *Proceedings of The 26th International Conference on Artificial Intelligence and Statistics*, volume 206, pages 837–849. PMLR.
- [12] Chauhan, V. K., Singh, S., and Sharma, A. (2023b). Hcr-net: A deep learning based script independent handwritten character recognition network. *arXiv preprint arXiv:2108.06663*.

- [13] Chauhan, V. K., Thakur, A., O’Donoghue, O., and Clifton, D. A. (2022a). Coper: Continuous patient state perceiver. In *2022 IEEE-EMBS International Conference on Biomedical and Health Informatics (BHI)*, pages 1–4. IEEE.
- [14] Chauhan, V. K., Thakur, A., O’Donoghue, O., Rohanian, O., and Clifton, D. A. (2022b). Continuous patient state attention models. *medRxiv*, pages 2022–12.
- [15] Chauhan, V. K., Zhou, J., Molaei, S., Ghosheh, G., and Clifton, D. A. (2023c). Dynamic inter-treatment information sharing for heterogeneous treatment effects estimation. *arXiv preprint arXiv:2305.15984*.
- [16] de Avila Belbute-Peres, F., fan Chen, Y., and Sha, F. (2021). HyperPINN: Learning parameterized differential equations with physics-informed hypernetworks. In *The Symbiosis of Deep Learning and Differential Equations*.
- [17] Deutsch, L., Nijkamp, E., and Yang, Y. (2019). A generative model for sampling high-performance and diverse weights for neural networks. *arXiv preprint arXiv:1905.02898*.
- [18] Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- [19] Dinh, T. M., Tran, A. T., Nguyen, R., and Hua, B.-S. (2022). Hyperinverter: Improving stylegan inversion via hypernetwork. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11389–11398.
- [20] Ehret, B., Henning, C., Cervera, M., Meulemans, A., Oswald, J. V., and Grewe, B. F. (2021). Continual learning in recurrent neural networks. In *International Conference on Learning Representations*.
- [21] Ferens, R. and Keller, Y. (2023). Hyperpose: Camera pose localization using attention hypernetworks. *arXiv preprint arXiv:2303.02610*.
- [22] Galanti, T. and Wolf, L. (2020). On the modularity of hypernetworks. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H., editors, *Advances in Neural Information Processing Systems*, volume 33, pages 10409–10419. Curran Associates, Inc.
- [23] Glorot, X. and Bengio, Y. (2010). Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 249–256. JMLR Workshop and Conference Proceedings.
- [24] Ha, D., Dai, A., and Le, Q. V. (2016). Hypernetworks. *arXiv preprint arXiv:1609.09106*.
- [25] Ha, D., Dai, A. M., and Le, Q. V. (2017). Hypernetworks. In *International Conference on Learning Representations*.
- [26] He, K., Zhang, X., Ren, S., and Sun, J. (2015). Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pages 1026–1034.
- [27] Henning, C., Cervera, M., D’Angelo, F., Oswald, J. V., Traber, R., Ehret, B., Kobayashi, S., Grewe, B. F., and Sacramento, J. (2021). Posterior meta-replay for continual learning. In Beygelzimer, A., Dauphin, Y., Liang, P., and Vaughan, J. W., editors, *Advances in Neural Information Processing Systems*.
- [28] Höfer, T., Kiefer, B., Messmer, M., and Zell, A. (2023). Hyperposepdf - hypernetworks predicting the probability distribution on so(3). In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 2369–2379.
- [29] Huang, Y., Xie, K., Bharadhwaj, H., and Shkurti, F. (2021). Continual model-based reinforcement learning with hypernetworks. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pages 799–805. IEEE.
- [30] Kingma, D. P. and Ba, J. (2014). Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- [31] Klocek, S., Maziarka, Ł., Wołczyk, M., Tabor, J., Nowak, J., and Śmieja, M. (2019). Hypernetwork functional image representation. In *Artificial Neural Networks and Machine Learning—ICANN 2019: Workshop and Special Sessions: 28th International Conference on Artificial Neural Networks, Munich, Germany, September 17–19, 2019, Proceedings 28*, pages 496–510. Springer.
- [32] Kristiadi, A., Däubener, S., and Fischer, A. (2019). Predictive uncertainty quantification with compound density networks. *arXiv preprint arXiv:1902.01080*.
- [33] Krueger, D., Huang, C.-W., Islam, R., Turner, R., Lacoste, A., and Courville, A. (2018). Bayesian hypernetworks.

- [34] Lamb, A., Saveliev, E., Li, Y., Tschitschek, S., Longden, C., Woodhead, S., Hernández-Lobato, J. M., Turner, R. E., Cameron, P., and Zhang, C. (2021). Contextual hypernetworks for novel feature adaptation. *arXiv preprint arXiv:2104.05860*.
- [35] Lee, E. and Lee, C.-Y. (2020). Neuronscale: Efficient scaling of neurons for resource-constrained deep neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [36] Li, Y. (2017). Deep reinforcement learning: An overview. *arXiv preprint arXiv:1701.07274*.
- [37] Li, Y., Gu, S., Zhang, K., Van Gool, L., and Timofte, R. (2020). Dhp: Differentiable meta pruning via hypernetworks. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part VIII 16*, pages 608–624. Springer.
- [38] Littwin, E., Galanti, T., Wolf, L., and Yang, G. (2020). On infinite-width hypernetworks. *Advances in neural information processing systems*, 33:13226–13237.
- [39] Littwin, G. and Wolf, L. (2019). Deep meta functionals for shape representation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1824–1833.
- [40] Liu, Z., Mu, H., Zhang, X., Guo, Z., Yang, X., Cheng, K.-T., and Sun, J. (2019). Metapruning: Meta learning for automatic neural network channel pruning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3296–3305.
- [41] Lorraine, J. and Duvenaud, D. (2018). Stochastic hyperparameter optimization through hypernetworks.
- [42] Madry, A., Makelov, A., Schmidt, L., Tsipras, D., and Vladu, A. (2017). Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*.
- [43] Mahabadi, R. K., Ruder, S., Dehghani, M., and Henderson, J. (2021). Parameter-efficient multi-task fine-tuning for transformers via shared hypernetworks. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 565–576.
- [44] Muller, L. K. (2021). Overparametrization of hypernetworks at fixed flop-count enables fast neural image enhancement. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, pages 284–293.
- [45] Navon, A., Shamsian, A., Fetaya, E., and Chechik, G. (2021). Learning the pareto front with hypernetworks. In *International Conference on Learning Representations*.
- [46] Nguyen, P., Tran, T., Le, K., Gupta, S., Rana, S., Nguyen, D., Nguyen, T., Ryan, S., and Venkatesh, S. (2021). Fast conditional network compression using bayesian hypernetworks. In Oliver, N., Pérez-Cruz, F., Kramer, S., Read, J., and Lozano, J. A., editors, *Machine Learning and Knowledge Discovery in Databases. Research Track*, pages 330–345, Cham. Springer International Publishing.
- [47] Nirkin, Y., Wolf, L., and Hassner, T. (2021). Hyperseg: Patch-wise hypernetwork for real-time semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4061–4070.
- [48] Oh, G. and Peng, H. (2022). Cvae-h: Conditionalizing variational autoencoders via hypernetworks and trajectory forecasting for autonomous driving. *arXiv preprint arXiv:2201.09874*.
- [49] Pan, Z., Liang, Y., Zhang, J., Yi, X., Yu, Y., and Zheng, Y. (2018). Hyperst-net: Hypernetworks for spatio-temporal forecasting. *arXiv preprint arXiv:1809.10889*.
- [50] Peng, H., Du, H., Yu, H., Li, Q., Liao, J., and Fu, J. (2020). Cream of the crop: Distilling prioritized paths for one-shot neural architecture search. *Advances in Neural Information Processing Systems*, 33:17955–17964.
- [51] Qu, J., Faney, T., Wang, Z., Gallinari, P., Yousef, S., and de Hemptinne, J.-C. (2022). Hmoe: Hypernetwork-based mixture of experts for domain generalization. *arXiv preprint arXiv:2211.08253*.
- [52] Ratzlaff, N. and Fuxin, L. (2019). Hypergan: A generative model for diverse, performant neural networks. In *International Conference on Machine Learning*, pages 5361–5369. PMLR.
- [53] Rezaei-Shoshtari, S., Morissette, C., Hogan, F. R., Dudek, G., and Meger, D. (2023). Hypernetworks for zero-shot transfer in reinforcement learning. *arXiv preprint arXiv:2211.15457*.
- [54] Rusu, A. A., Rao, D., Sygnowski, J., Vinyals, O., Pascanu, R., Osindero, S., and Hadsell, R. (2019). Meta-learning with latent embedding optimization. In *International Conference on Learning Representations*.
- [55] Ruta, D., Gilbert, A., Motiian, S., Faieta, B., Lin, Z., and Collomosse, J. (2023). Hypernst: Hyper-networks for neural style transfer. In Karlinsky, L., Michaeli, T., and Nishino, K., editors, *Computer Vision – ECCV 2022 Workshops*, pages 201–217, Cham. Springer Nature Switzerland.

- [56] Sarafian, E., Keynan, S., and Kraus, S. (2021). Recomposing the reinforcement learning building blocks with hypernetworks. In *International Conference on Machine Learning*, pages 9301–9312. PMLR.
- [57] Schmidhuber, J. (1992). Learning to control fast-weight memories: An alternative to dynamic recurrent networks. *Neural Computation*, 4(1):131–139.
- [58] Schmidhuber, J. (1993). A ‘self-referential’ weight matrix. In *ICANN’93: Proceedings of the International Conference on Artificial Neural Networks Amsterdam, The Netherlands 13–16 September 1993 3*, pages 446–450. Springer.
- [59] Shamsian, A., Navon, A., Fetaya, E., and Chechik, G. (2021). Personalized federated learning using hypernetworks. In *International Conference on Machine Learning*, pages 9489–9502. PMLR.
- [60] Shih, A., Sadigh, D., and Ermon, S. (2021). Hyperspns: Compact and expressive probabilistic circuits. In Ranzato, M., Beygelzimer, A., Dauphin, Y., Liang, P., and Vaughan, J. W., editors, *Advances in Neural Information Processing Systems*, volume 34, pages 8571–8582. Curran Associates, Inc.
- [61] Spurek, P., Zieba, M., Tabor, J., and Trzcinski, T. (2022). General hypernetwork framework for creating 3d point clouds. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(12):9995–10008.
- [62] Stanley, K. O., D’Ambrosio, D. B., and Gauci, J. (2009). A hypercube-based encoding for evolving large-scale neural networks. *Artificial life*, 15(2):185–212.
- [63] Sun, Z., Ozay, M., and Okatani, T. (2017). Hypernetworks with statistical filtering for defending adversarial examples. *arXiv preprint arXiv:1711.01791*.
- [64] Szatkowski, F., Piczak, K. J., Spurek, P., Tabor, J., and Trzcinski, T. (2022). Hypersound: Generating implicit neural representations of audio signals with hypernetworks. In *Sixth Workshop on Meta-Learning at the Conference on Neural Information Processing Systems*.
- [65] Tay, Y., Zhao, Z., Bahri, D., Metzler, D., and Juan, D.-C. (2021). Hypergrid transformers: Towards a single model for multiple tasks. In *International Conference on Learning Representations*.
- [66] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.
- [67] Volk, T., Ben-David, E., Amosy, O., Chechik, G., and Reichart, R. (2022). Example-based hypernetworks for out-of-distribution generalization. *arXiv preprint arXiv:2203.14276*.
- [68] von Oswald, J., Henning, C., Grewe, B. F., and Sacramento, J. (2020). Continual learning with hypernetworks. In *International Conference on Learning Representations*.
- [69] Wiens, J., Guttag, J., and Horvitz, E. (2014). A study in transfer learning: leveraging data from multiple hospitals to enhance hospital-specific predictions. *Journal of the American Medical Informatics Association*, 21(4):699–706.
- [70] Wu, Q., Bauer, D., Chen, Y., and Ma, K.-L. (2023). Hyperinr: A fast and predictive hypernetwork for implicit neural representations via knowledge distillation. *arXiv preprint arXiv:2304.04188*.
- [71] Wullach, T., Adler, A., and Minkov, E. (2022). Character-level hypernetworks for hate speech detection. *Expert Systems with Applications*, 205:117571.
- [72] Yin, L., Perez-Rua, J. M., and Liang, K. J. (2022). Sylph: A hypernetwork framework for incremental few-shot object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9035–9045.
- [73] Zhang, C., Ren, M., and Urtasun, R. (2019). Graph hypernetworks for neural architecture search. In *International Conference on Learning Representations*.
- [74] Zhao, D., Kobayashi, S., Sacramento, J., and von Oswald, J. (2020). Meta-learning via hypernetworks. In *4th Workshop on Meta-Learning at NeurIPS 2020 (MetaLearn 2020)*. NeurIPS.