

Exploring Randomly Wired Neural Networks for Image Recognition

Saining Xie Alexander Kirillov Ross Girshick Kaiming He

Facebook AI Research (FAIR)

Abstract

Neural networks for image recognition have evolved through extensive manual design from simple chain-like models to structures with multiple wiring paths. The success of ResNets [11] and DenseNets [16] is due in large part to their innovative wiring plans. Now, neural architecture search (NAS) studies are exploring the joint optimization of wiring and operation types, however, the space of possible wirings is constrained and still driven by manual design despite being searched. In this paper, we explore a more diverse set of connectivity patterns through the lens of randomly wired neural networks. To do this, we first define the concept of a stochastic network generator that encapsulates the entire network generation process. Encapsulation provides a unified view of NAS and randomly wired networks. Then, we use three classical random graph models to generate randomly wired graphs for networks. The results are surprising: several variants of these random generators yield network instances that have competitive accuracy on the ImageNet benchmark. These results suggest that new efforts focusing on designing better network generators may lead to new breakthroughs by exploring less constrained search spaces with more room for novel design.

1. Introduction

What we call deep learning today descends from the *connectionist* approach to cognitive science [38, 7]—a paradigm reflecting the hypothesis that *how computational networks are wired* is crucial for building intelligent machines. Echoing this perspective, recent advances in computer vision have been driven by moving from models with chain-like wiring [19, 53, 42, 43] to more elaborate connectivity patterns, e.g., ResNet [11] and DenseNet [16], that are effective in large part because of how they are wired.

Advancing this trend, neural architecture search (NAS) [55, 56] has emerged as a promising direction for jointly searching wiring patterns and which operations to perform. NAS methods focus on *search* [55, 56, 33, 26, 29, 27] while implicitly relying on an important—yet largely overlooked—component that we call a *network generator* (defined in §3.1). The NAS network generator defines a family of possible wiring patterns from which networks

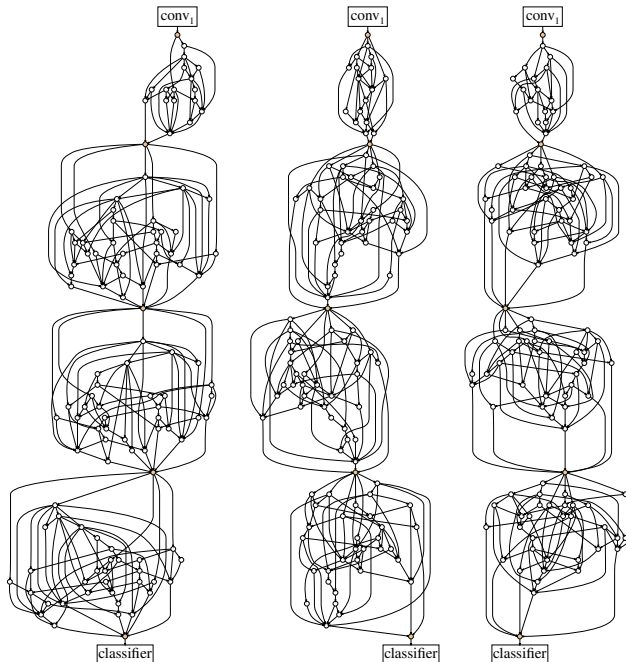


Figure 1. **Randomly wired neural networks** generated by the classical Watts-Strogatz (WS) [50] model: these three instances of random networks achieve (left-to-right) 79.1%, 79.1%, 79.0% classification accuracy on ImageNet under a similar computational budget to ResNet-50, which has 77.1% accuracy.

are sampled subject to a learnable probability distribution. However, like the wiring patterns in ResNet and DenseNet, the NAS network generator is *hand designed* and the space of allowed wiring patterns is constrained in a small subset of all possible graphs. Given this perspective, we ask: *What happens if we loosen this constraint and design novel network generators?*

We explore this question through the lens of *randomly wired neural networks* that are sampled from stochastic network generators, in which a human-designed *random process* defines generation. To reduce bias from us—the authors of this paper—on the generators, we use *three classical families of random graph models in graph theory* [51]: the Erdős-Rényi (ER) [6], Barabási-Albert (BA) [1], and Watts-Strogatz (WS) [50] models. To define complete networks, we convert a random graph into a directed acyclic graph (DAG) and apply a simple mapping from nodes to their functional roles (e.g., to the same type of convolution).

The results are surprising: several variants of these random generators yield networks with competitive accuracy on ImageNet [39]. The best generators, which use the WS model, produce multiple networks that outperform or are comparable to their fully manually designed counterparts and the networks found by various neural architecture search methods. We also observe that the variance of accuracy is low for different random networks produced by the same generator, yet there can be clear accuracy gaps between different generators. These observations suggest that the network generator *design* is important.

We note that these randomly wired networks are *not* “prior free” even though they are random. Many strong priors are in fact implicitly designed into the generator, including the choice of a particular rule and distribution to control the probability of wiring or not wiring certain nodes together. Each random graph model [6, 50, 1] has certain probabilistic behaviors such that sampled graphs likely exhibit certain properties (*e.g.*, WS is highly clustered [50]). Ultimately, the generator design determines a probabilistic distribution over networks, and as a result these networks tend to have certain properties. The generator design underlies the prior and thus should not be overlooked.

Our work explores a direction orthogonal to concurrent work on *random search* for NAS [23, 41]. These studies show that random search is competitive in “the NAS search space” [55, 56], *i.e.*, the “NAS network generator” in our perspective. Their results can be understood as showing that the prior induced by the NAS generator design tends to produce good models, similar to our observations. In contrast to [23, 41], our work goes beyond the design of established NAS generators and explores different random generator designs.

Finally, our work suggests a new transition from designing an individual network to *designing a network generator* may be possible, analogous to how our community have transitioned from designing features to designing a network that learns features. Rather than focusing primarily on search with a fixed generator, we suggest designing new network generators that produce new families of models for searching. The importance of the *designed* network generator (in NAS and elsewhere) also implies that machine learning has not been automated (*c.f.* “AutoML” [20])—the underlying human design and prior shift from network engineering to network generator engineering.

2. Related Work

Network wiring. Early recurrent and convolutional neural networks (RNNs and CNNs) [37, 21] use chain-like wiring patterns. LSTMs [14] use more sophisticated wiring to create a gating mechanism. Inception CNNs [45, 46, 44] concatenate multiple, irregular branching pathways, while ResNets [11] use $x + \mathcal{F}(x)$ as a regular wiring template;

DenseNets [16] use concatenation instead: $[x, \mathcal{F}(x)]$. The LSTM, Inception, ResNet, and DenseNet wiring patterns are effective in general, beyond any individual instantiation.

Neural architecture search (NAS). Zoph and Le [55] define a NAS search space and investigate reinforcement learning (RL) as an optimization algorithm. Recent research on NAS mainly focuses on optimization methods, including RL [55, 56], progressive [26], gradient-based [29, 27], weight-sharing [33], evolutionary [34], and random search [23, 41] methods. The search space in these NAS works, determined by the network generator implicit in [55], is largely unchanged in these works. While this is reasonable for comparing optimization methods, it inherently limits the set of feasible solutions.

Randomly wired machines. Pioneers of artificial intelligence were originally interested in randomly wired hardware and their implementation in computer programs (*i.e.*, artificial neural networks). In 1940s, Turing [47] suggested a concept of *unorganized machines*, which is a form of the earliest randomly connected neural networks. One of the first neural network learning machines, designed by Minsky [31] in 1950s and implemented using vacuum tubes, was randomly wired. In late 1950s the “Mark I Perceptron” visual recognition machine built by Rosenblatt [36] used an array of randomly connected photocells.

Relation to neuroscience. Turing [47] analogized the unorganized machines to an infant human’s cortex. Rosenblatt [36] pointed out that “*the physical connections of the nervous system ... are not identical from one organism to another*”, and “*at birth, the construction of the most important networks is largely random.*” Studies [50, 48] have observed that the neural network of a nematode (a worm) with about 300 neurons is a graph with *small-world* properties [18]. Random graph modeling has been used as a tool to study the neural networks of human brains [2, 4, 3].

Random graphs in graph theory. Random graphs are widely studied in graph theory [51]. Random graphs exhibit different probabilistic behaviors depending on the random process defined by the model (*e.g.*, [6, 1, 50]). The definition of the random graph model determines the prior knowledge encoded in the resulting graphs (*e.g.*, small-world [18]) and may connect them to naturally occurring phenomena. As a result, random graph models are an effective tool for modeling and analyzing real-world graphs, *e.g.*, social networks, world wide web, citation networks.

3. Methodology

We now introduce the concept of a network generator, which is the foundation of randomly wired neural networks.

3.1. Network Generators

We define a *network generator* as a mapping g from a parameter space Θ to a space of neural network architectures \mathcal{N} , $g: \Theta \mapsto \mathcal{N}$. For a given $\theta \in \Theta$, $g(\theta)$ returns a neural network instance $n \in \mathcal{N}$. The set \mathcal{N} is typically a family of related networks, for example, VGG nets [43], ResNets [11], or DenseNets [16].

The generator g determines, among other concerns, how the computational graph is wired. For example, in ResNets a generator produces a stack of blocks that compute $x + \mathcal{F}(x)$. The parameters θ specify the instantiated network and may contain diverse information. For example, in a ResNet generator, θ can specify the number of stages, number of residual blocks for each stage, depth/width/filter sizes, activation types, etc.

Intuitively, one may think of g as a function in a programming language, e.g. Python, that takes a list of arguments (corresponding to θ), and returns a network architecture. The network representation n returned by the generator is *symbolic*, meaning that it specifies the type of operations that are performed and the flow of data; it does *not* include values of network weights,¹ which are learned from data after a network is generated.

Stochastic network generators. The above network generator $g(\theta)$ performs a *deterministic* mapping: given the same θ , it always returns the same network architecture n . We can extend g to accept an additional argument s that is the *seed* of a pseudo-random number generator that is used internally by g . Given this seed, one can construct a (pseudo) random family of networks by calling $g(\theta, s)$ multiple times, keeping θ fixed but changing the value of $s = 1, 2, 3, \dots$. For a fixed value of θ , a uniform probability distribution over all possible seed values induces a (likely non-uniform) probability distribution over \mathcal{N} . We call generators of the form $g(\theta, s)$ *stochastic network generators*.

Before we discuss our method, we provide additional background by reinterpreting the work on NAS [55, 56] in the context of stochastic network generators.

NAS from the network generator perspective. The NAS methods of [55, 56] utilize an LSTM “controller” in the process of generating network architectures. But the LSTM is only part of the complete NAS network generator, which is in fact a stochastic network generator, as illustrated next.

The weight matrices of the LSTM are the parameters θ of the generator. The output of each LSTM time-step is a probability distribution conditioned on θ . Given this distribution and the seed s , each step samples a construction action (e.g., insert an operator, connect two nodes). The parameters θ of the LSTM controller, due to its probabilistic behavior, are optimized (searched for) by RL in [55, 56].

¹We use *parameters* to refer to network generator arguments and *weights* to refer to the learnable weights and biases of a generated network.

However, the LSTM is not the only component in the NAS network generator $g(\theta, s)$. *There are also hand-designed rules defined to map the sampled actions to a computational DAG, and these rules are also part of g .* Using the node/edge terminology in graph theory, for a NAS network in [56], if we map a combination operation (e.g., sum) to a node and a unary transformation (e.g., conv) to an edge (see the supplement), the rules of the NAS generator include, but are not limit to:

- A subgraph to be searched, called a cell [56], always accepts the activations of the output nodes from the 2 immediately preceding cells;
- Each cell contains 5 nodes that are wired to 2 and only 2 existing nodes, chosen by sampling from the probability distribution output by the LSTM;
- All nodes that have no output in a cell are concatenated by an extra node to form a valid DAG for the cell.

All of the generation rules, together with the choice of using an LSTM, and other hyper-parameters of the system (e.g., the number of nodes, say, 5), comprise the NAS network generator that produces a full DAG. It is also worth noticing that the view of “node as combination and edge as transformation” is not the only way to interpret a neural network as a graph, and so it is not the only way to turn a general graph into a neural network (we use a different mapping in §3.2).

Encapsulating the *complete* generation process, as we have illustrated, reveals which components are optimized and which are hard-coded. It now becomes explicit that the network space \mathcal{N} has been carefully restricted by hand-designed rules. For example, the rules listed above suggest that each of the 5 nodes in a cell always *has precisely input degree² 2 and output degree 1* (see the supplement). This does not cover all possible 5-(internal-)node graphs. It is in a highly restricted network space. Viewing NAS from the perspective of a network generator helps explain the recently demonstrated ineffectiveness of sophisticated optimization vs. random search [23, 41]: the manual design in the NAS network generator is a *strong* prior, which represents a meta-optimization beyond the search over θ (by RL, e.g.) and s (by random search).

3.2. Randomly Wired Neural Networks

Our analysis of NAS reveals that the network generator is hand-designed and encodes a prior from human knowledge. It is likely that the design of the network generator plays a considerable role—if so, current methods are short of achieving “AutoML” [20] and still involve significant human effort (c.f. “Our experiments show that Neural Architecture Search can design good models from *scratch*.” [55],

²In graph theory, “degree” is the number of edges connected to a node. We refer to “input/output degree” as that of input/output edges to a node.

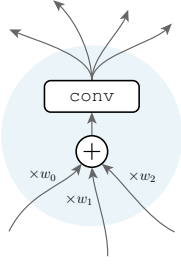


Figure 2. **Node operations** designed for our random graphs. Here we illustrate a node (blue circle) with 3 input edges and 4 output edges. The aggregation is done by weighted sum with learnable positive weights w_0, w_1, w_2 . The transformation is a ReLU-convolution-BN triplet, simply denoted as `conv`. The transformed data are sent out as 4 copies.

emphasis added). To investigate how important the generator design is, it is not sufficient to compare different optimizers (sophisticated or random) for the same NAS generator; it is necessary to study *new* network generators that are substantially different from the NAS generator.

This leads to our exploration of *randomly wired neural networks*. That is, we will define network generators that yield networks with random graphs, subject to different human-specific priors. To minimize the human bias from us—the authors of this paper—on the prior, we will use three *classical random graph models* in our study ([6, 1, 50]; §3.3). Our methodology for generating randomly wired networks involves the following concepts:

Generating general graphs. Our network generator starts by generating a general graph (in the sense of graph theory). It generates a set of nodes and edges that connect nodes, without restricting how the graphs correspond to neural networks. This allows us to freely use any general graph generator from graph theory (ER/BA/WS). Once a graph is obtained, it is mapped to a computable neural network.

The mapping from a general graph to neural network operations is in itself arbitrary, and thus also human-designed. We intentionally use a simple mapping, discussed next, so that we can focus on graph wiring patterns.

Edge operations. Assuming by construction that the graph is directed, we define that edges are data flow, *i.e.*, a directed edge sends data (a tensor) from one node to another node.

Node operations. A node in a directed graph may have some input edges and some output edges. We define the operations represented by one node (Figure 2) as:

- *Aggregation:* The input data (from one or more edges) to a node are combined via a weighted sum; the weights are learnable and positive.³

- *Transformation:* The aggregated data is processed by a transformation defined as a ReLU-convolution-BN triplet⁴ [12]. The same type of convolution is used for all nodes, *e.g.*, a 3×3 separable convolution⁵ by default.

³Applying sigmoid on unrestricted weights ensures they are positive.

⁴Instead of a triplet with a convolution followed by BN [17] then ReLU [32], we use the ReLU-convolution-BN triplet, as it means the aggregation (at the next nodes) can receive positive and negative activation, preventing the aggregated activation from being inflated in case of a large input degree.

⁵Various implementations of separable convolutions exist. We use the

- *Distribution:* The same copy of the transformed data is sent out by the output edges of the node.

These operations have some nice properties:

- (i) Additive aggregation (unlike concatenation) maintains the same number of output channels as input channels, and this prevents the convolution that follows from growing large in computation, which may increase the importance of nodes with large input degree simply because they increase computation, not because of how they are wired.

- (ii) The transformation should have the same number of output and input channels (unless switching stages; discussed later), to make sure the transformed data can be combined with the data from any other nodes. Fixing the channel count then keeps the FLOPs (floating-point operations) and parameter count unchanged for each node, regardless of its input and output degrees.

- (iii) Aggregation and distribution are almost parameter-free (except for a negligible number of parameters for weighted summation), regardless of input and output degrees. Also, given that every edge is parameter-free the overall FLOPs and parameter count of a graph are roughly proportional to the number of nodes, and nearly independent of the number of edges.

These properties nearly decouple FLOPs and parameter count from network wiring, *e.g.*, the deviation of FLOPs is typically $\pm 2\%$ among our random network instances or different generators. This enables the comparison of different graphs without inflating/deflating model complexity. Differences in task performance are therefore reflective of the properties of the wiring pattern.

Input and output nodes. Thus far, a general graph is not yet a valid neural network even given the edge/node operations, because it may have multiple input nodes (*i.e.*, those without any input edge) and multiple output nodes. It is desirable to have a single input and a single output for typical neural networks, *e.g.*, for image classification. We apply a simple post-processing step.

For a given general graph, we create a single extra node that is connected to all original input nodes. This is the unique input node that sends out the same copy of input data to all original input nodes. Similarly, we create a single extra node that is connected to all original output nodes. This is the unique output node; we have it compute the (unweighted) average from all original output nodes. These two nodes perform no convolution. When referring to the node count N , we exclude these two nodes.

Stages. With unique input and output nodes, it is sufficient for a graph to represent a valid neural network. Nevertheless, in image classification in particular, networks that

form of [5]: a 3×3 separable convolution is a 3×3 *depth-wise* convolution followed by a 1×1 convolution, with no non-linearity in between.

| stage | output | <i>small regime</i> | <i>regular regime</i> |
|-------------------|---------|---|----------------------------------|
| conv ₁ | 112×112 | 3×3 conv, $C/2$ | |
| conv ₂ | 56×56 | 3×3 conv, C | random wiring $N/2, C$ |
| conv ₃ | 28×28 | random wiring N, C | random wiring $N, 2C$ |
| conv ₄ | 14×14 | random wiring $N, 2C$ | random wiring $N, 4C$ |
| conv ₅ | 7×7 | random wiring $N, 4C$ | random wiring $N, 8C$ |
| classifier | 1×1 | 1×1 conv, 1280-d global average pool, 1000-d fc, softmax | |

Table 1. **RandWire architectures** for small and regular computation networks. A random graph is denoted by the node count (N) and channel count for each node (C). We use conv to denote a ReLU-Conv-BN triplet (expect conv₁ is Conv-BN). The input size is 224×224 pixels. The change of the output size implies a stride of 2 (omitted in table) in the convolutions that are right after the input of each stage.

maintain the full input resolution throughout are not desirable. It is common [19, 43, 11, 56] to divide a network into *stages* that progressively down-sample feature maps.

We use a simple strategy: the random graph generated above defines one stage. Analogous to the stages in a ResNet, *e.g.*, conv_{1,2,3,4,5} [11], our entire network consists of multiple stages. One random graph represents one stage, and it is connected to its preceding/succeeding stage by its unique input/output node. For all nodes that are directly connected to the input node, their transformations are modified to have a stride of 2. The channel count in a random graph is increased by $2 \times$ when going from one stage to the next stage, following [11].

Table 1 summarizes the randomly wired neural networks, referred to as **RandWire**, used in our experiments. They come in small and regular complexity regimes (more in §4). For conv₁ and/or conv₂ we use a single convolutional layer for simplicity with multiple random graphs following. The network ends with a classifier output (Table 1, last row). Figure 1 shows full computation graphs of three randomly wired network samples.

3.3. Random Graph Models

We now describe in brief the three classical random graph models used in our study. We emphasize that these random graph models are not proposed by this paper; we describe them for completeness. The three classical models all generate undirected graphs; we use a simple heuristic to turn them into DAGs (see the supplement).

Erdős-Rényi (ER). In the ER model [8, 6], with N nodes, an edge between two nodes is connected with probability P , independent of all other nodes and edges. This process is iterated for all pairs of nodes. The ER generation model has only a single parameter P , and is denoted as ER(P).

Any graph with N nodes has non-zero probability of being generated by the ER model, including graphs that are disconnected. However, a graph generated by ER(P) has high probability of being a single connected component if $P > \frac{\ln(N)}{N}$ [6]. This provides one example of an implicit bias introduced by a generator.

Barabási-Albert (BA). The BA model [1] generates a random graph by sequentially adding new nodes. The initial state is M nodes without any edges ($1 \leq M < N$). The method sequentially adds a new node with M new edges. For a node to be added, it will be connected to an existing node v with probability proportional to v ’s degree. The new node repeatedly adds non-duplicate edges in this way until it has M edges. Then this is iterated until the graph has N nodes. The BA generation model has only a single parameter M , and is denoted as BA(M).

Any graph generated by BA(M) has exactly $M \cdot (N - M)$ edges. So the set of all graphs generated by BA(M) is a *subset* of all possible N -node graphs—this gives one example on how an underlying prior can be introduced by the graph generator in spite of randomness.

Watts-Strogatz (WS). The WS model [50] was defined to generate small-world graphs [18]. Initially, the N nodes are regularly placed in a ring and each node is connected to its $K/2$ neighbors on both sides (K is an even number). Then, in a clockwise loop, for every node v , the edge that connects v to its clockwise i -th next node is *rewired* with probability P . “Rewiring” is defined as uniformly choosing a random node that is not v and that is not a duplicate edge. This loop is repeated $K/2$ times for $1 \leq i \leq K/2$. K and P are the only two parameters of the WS model, denoted as WS(K, P).

Any graph generated by WS(K, P) has exactly $N \cdot K$ edges. WS(K, P) only covers a small subset of all possible N -node graphs too, but this subset is different from the subset covered by BA. This provides an example on how a different underlying prior has been introduced.

3.4. Design and Optimization

Our randomly wired neural networks are generated by a stochastic network generator $g(\theta, s)$. The random graph parameters, namely, $P, M, (K, P)$ in ER, BA, WS respectively, are part of the parameters θ . The “optimization” of such a 1- or 2-parameter space is essentially done by *trial-and-error* by human designers, *e.g.*, by line/grid search. Conceptually, such “optimization” is not distinct from many other *designs* involved in our and other models (including NAS), *e.g.*, the number of nodes, stages, and filters.

Optimization can also be done by scanning the random seed s , which is an implementation of random search. Random search is possible for any stochastic network generator, including ours and NAS. But as we present by experiment,

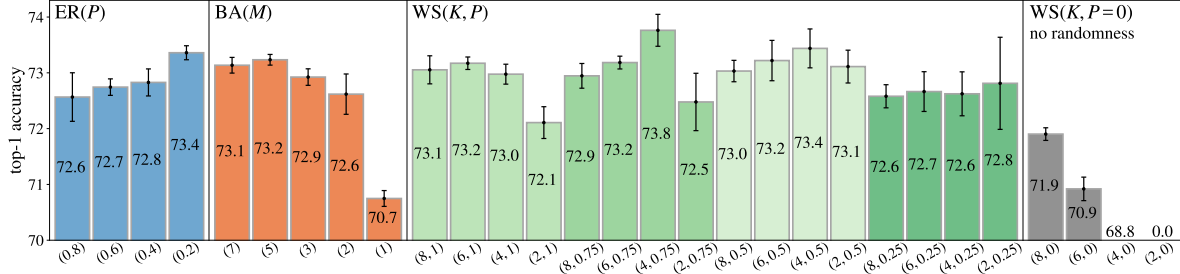


Figure 3. **Comparison on random graph generators: ER, BA, and WS** in the small computation regime. Each bar represents the results of a generator under a parameter setting for P , M , or (K, P) (tagged in x-axis). The results are ImageNet top-1 accuracy, shown as mean and standard deviation (std) over 5 random network instances sampled by a generator. At the rightmost, $WS(K, P=0)$ has no randomness.

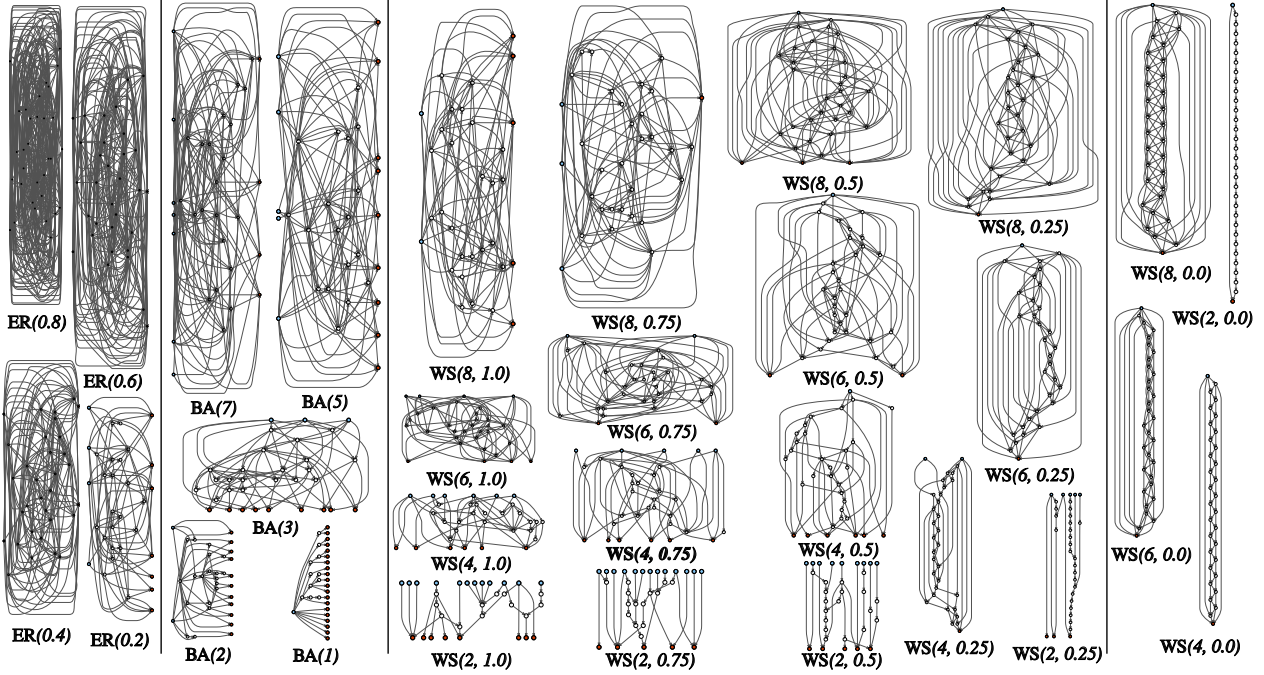


Figure 4. **Visualization of the random graphs generated by ER, BA, and WS**. Each plot represents one random graph instance sampled by the specified generator. The generators are those in Figure 3. The node count is $N=32$ for each graph. A blue/red node denotes an input/output node, to which an extra unique input/output node (not shown) will be added (see §3.2).

the accuracy variation of our networks is small for different seeds s , suggesting that the benefit of random search may be small. So we perform *no random search* and report *mean* accuracy of multiple random network instances. As such, our network generator has *minimal optimization* (1- or 2-parameter grid search) beyond their hand-coded *design*.

4. Experiments

We conduct experiments on the ImageNet 1000-class classification task [39]. We train on the training set with ~ 1.28 M images and test on the 50K validation images.

Architecture details. Our experiments span a small computation regime (e.g., MobileNet [15] and ShuffleNet [54]) and a regular computation regime (e.g., ResNet-50/101 [11]). RandWire nets in these regimes are in Table 1, where N nodes and C channels determine network complexity.

We set $N=32$, and then set C to the nearest integer such that target model complexity is met: $C=78$ in the small regime, and $C=109$ or 154 in the regular regime.

Random seeds. For each generator, we randomly sample 5 network instances (5 random seeds), train them from scratch, and evaluate accuracy for each instance. To emphasize that we perform *no random search* for each generator, we report the classification accuracy with “mean \pm std” for *all* 5 random seeds (i.e., we do *not* pick the best). We use the same seeds 1, ..., 5 for all experiments.

Implementation details. We train our networks for 100 epochs, unless noted. We use a half-period-cosine shaped learning rate decay [28, 16]. The initial learning rate is 0.1, the weight decay is $5e-5$, and the momentum is 0.9. We use label smoothing regularization [44] with a coefficient of 0.1. Other details of the training procedure are the same as [10].

4.1. Analysis Experiments

Random graph generators. Figure 3 compares the results of different generators in the small computation regime: each RandWire net has $\sim 580\text{M}$ FLOPs. Figure 4 visualizes one example graph for each generator. The graph generator is specified by the random graph model (ER/BA/WS) and its set of parameters: *e.g.*, ER(0.2). We observe:

All random generators provide decent accuracy over all 5 random network instances; none of them fails to converge. ER, BA, and WS all have certain settings that yield mean accuracy of $>73\%$, within a $<1\%$ gap from the best mean accuracy of 73.8% from WS(4, 0.75).

Moreover, the variation among the random network instances is low. Almost all random generators in Figure 3 have an standard deviation (std) of $0.2\sim 0.4\%$. As a comparison, training the same instance of a ResNet-50 multiple times has a typical std of $0.1\sim 0.2\%$ [10]. The observed low variance of our random generators suggests that even without random search (*i.e.*, picking the best from several random instances), it is likely that the accuracy of a network instance is close to the mean accuracy, subject to some noise.

On the other hand, different random generators may have a gap between their mean accuracies, *e.g.*, BA(1) has 70.7% accuracy and is $\sim 3\%$ lower than WS(4, 0.75). This suggests that random generator design, including the wiring priors (BA *vs.* WS) and generation parameters, plays an important role in the accuracy of sampled network instances.

Figure 3 also includes a set of *non-random* generators: WS($K, P=0$). “ $P=0$ ” means no random rewiring. Interestingly, the results of WS($K, P=0$) are all worse than their WS($K, P>0$) counterparts for any fixed K in Figure 3.

Graph damage. We explore *graph damage* by randomly removing one node or edge—an ablative setting inspired by [22, 49]. Formally, given a network instance *after training*, we remove one node or one edge from the graph and evaluate the validation accuracy *without any further training*.

When a node is removed, we evaluate the accuracy loss (Δ) *vs.* the output degree of that node (Figure 5, top). It is clear that ER, BA, and WS behave differently under such damage. For networks generated by WS, the mean degradation of accuracy is larger when the output degree of the removed node is higher. This implies that “hub” nodes in WS that send information to many nodes are influential.

When an edge is removed, we evaluate the accuracy loss *vs.* the input degree of this edge’s target node (Figure 5, bottom). If the input degree of an edge’s target node is smaller, removing this edge tends to change a larger portion of the target node’s inputs. This trend can be seen by the fact that the accuracy loss is generally decreasing along the x-axis in Figure 5 (bottom). The ER model is less sensitive to edge removal, possibly because in ER’s definition wiring of every edge is independent.

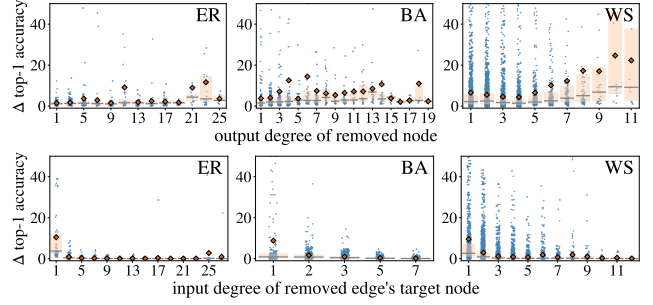


Figure 5. **Graph damage ablation.** We randomly *remove one node* (top) or *remove one edge* (bottom) from a graph after the network is trained, and evaluate the loss (Δ) in accuracy on ImageNet. From left to right are ER, BA, and WS generators. Red circle: *mean*; gray bar: *median*; orange box: *interquartile range*; blue dot: *an individual damaged instance*.

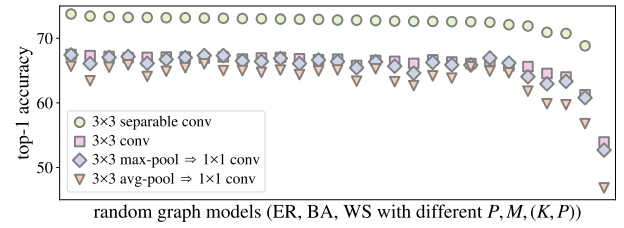


Figure 6. **Alternative node operations.** Each column is the mean accuracy of the same set of 5 random graphs equipped with different node operations, sorted by “ 3×3 separable conv” (from Figure 3). The generators roughly maintain their orders of accuracy.

Node operations. Thus far, all models in our experiment use a 3×3 separable convolution as the “conv” in Figure 2. Next we evaluate alternative choices. We consider: (i) 3×3 (regular) convolution, and (ii) 3×3 max-/average-pooling followed by a 1×1 convolution. We replace the transformation of *all nodes* with the specified alternative. We adjust the factor C to keep the complexity of all alternative networks.

Figure 6 shows the mean accuracy for each of the generators listed in Figure 3. Interestingly, almost all networks still converge to non-trivial results. Even “ 3×3 pool with 1×1 conv” performs similarly to “ 3×3 conv”. The network generators roughly maintain their accuracy ranking despite the operation replacement; in fact, the Pearson correlation between any two series in Figure 5 is $0.91\sim 0.98$. This suggests that the network wiring plays a role somewhat orthogonal to the role of the chosen operations.

4.2. Comparisons

Small computation regime. Table 2 compares our results in the *small computation regime*, a common setting studied in existing NAS papers. Instead of training for 100 epochs, here we train for 250 epochs following settings in [56, 34, 26, 27] for fair comparisons.

RandWire with WS(4, 0.75) has mean accuracy of 74.7% (with min 74.4% and max 75.0%). This result is better than or comparable to all existing hand-designed

| network | top-1 acc. | top-5 acc. | FLOPs (M) | params (M) |
|--------------------|----------------------------------|----------------------------------|---------------|---------------|
| MobileNet [15] | 70.6 | 89.5 | 569 | 4.2 |
| MobileNet v2 [40] | 74.7 | - | 585 | 6.9 |
| ShuffleNet [54] | 70.9 | 89.8 | 524 | ~5 |
| ShuffleNet v2 [30] | 73.7 | - | 524 | ~5 |
| NASNet-A [56] | 74.0 | 91.6 | 564 | 5.3 |
| NASNet-B [56] | 72.8 | 91.3 | 488 | 5.3 |
| NASNet-C [56] | 72.5 | 91.0 | 558 | 4.9 |
| Amoeba-A [34] | 74.5 | 92.0 | 555 | 5.1 |
| Amoeba-B [34] | 74.0 | 91.5 | 555 | 5.3 |
| Amoeba-C [34] | 75.7 | 92.4 | 570 | 6.4 |
| PNAS [26] | 74.2 | 91.9 | 588 | 5.1 |
| DARTS [27] | 73.1 | 91.0 | 595 | 4.9 |
| RandWire-WS | 74.7± 0.25 | 92.2± 0.15 | 583 ± 6.2 | 5.6 ± 0.1 |

Table 2. **ImageNet: small computation regime** (*i.e.*, <600M FLOPs). RandWire results are the mean accuracy (\pm std) of 5 random network instances, with WS(4, 0.75). Here we train for 250 epochs similar to [56, 34, 26, 27], for fair comparisons.

| network | top-1 acc. | top-5 acc. | FLOPs (B) | params (M) |
|--|----------------------------------|----------------------------------|----------------|-----------------|
| ResNet-50 [11] | 77.1 | 93.5 | 4.1 | 25.6 |
| ResNeXt-50 [52] | 78.4 | 94.0 | 4.2 | 25.0 |
| RandWire-WS, $C=109$ | 79.0± 0.17 | 94.4± 0.11 | 4.0 ± 0.09 | 31.9 ± 0.66 |
| ResNet-101 [11] | 78.8 | 94.4 | 7.8 | 44.6 |
| ResNeXt-101 [52] | 79.5 | 94.6 | 8.0 | 44.2 |
| RandWire-WS, $C=154$ | 80.1± 0.19 | 94.8± 0.18 | 7.9 ± 0.18 | 61.5 ± 1.32 |

Table 3. **ImageNet: regular computation regime** with FLOPs comparable to ResNet-50 (top) and to ResNet-101 (bottom). ResNeXt is the 32 \times 4 version [52]. RandWire is WS(4, 0.75).

| network | test size | epochs | top-1 acc. | top-5 acc. | FLOPs (B) | params (M) |
|--------------------|------------------|--------|-----------------|-----------------|-----------------|-----------------|
| NASNet-A [56] | 331 ² | >250 | 82.7 | 96.2 | 23.8 | 88.9 |
| Amoeba-B [34] | 331 ² | >250 | 82.3 | 96.1 | 22.3 | 84.0 |
| Amoeba-A [34] | 331 ² | >250 | 82.8 | 96.1 | 23.1 | 86.7 |
| PNASNet-5 [26] | 331 ² | >250 | 82.9 | 96.2 | 25.0 | 86.1 |
| RandWire-WS | 320 ² | 100 | 81.6 ± 0.13 | 95.6 ± 0.07 | 16.0 ± 0.36 | 61.5 ± 1.32 |

Table 4. **ImageNet: large computation regime.** Our networks are the same as in Table 3 ($C=154$), but we evaluate on 320 \times 320 images instead of 224 \times 224. Ours are only trained for 100 epochs.

wiring (MobileNet/ShuffleNet) and NAS-based results, except for AmoebaNet-C [34]. The *mean* accuracy achieved by RandWire is a competitive result, especially considering that we perform no random search in our random generators, and that we use a single operation type for all nodes.

Regular computation regime. Next we compare the RandWire networks with ResNet-50/101 [11] under similar FLOPs. In this regime, we use a regularization method inspired by our edge removal analysis: for each training mini-batch, we randomly remove one edge whose target node has input degree > 1 with probability of 0.1. This regularization is similar to DropPath adopted in NAS [56]. We train with a weight decay of 1e-5 and a Dropout [13] rate of 0.2 in the classifier fc layer. Other settings are the same as the small computation regime. We train the ResNet/ResNeXt competitors using the recipe of [10], but with the cosine

| backbone | AP | AP ₅₀ | AP ₇₅ | AP _S | AP _M | AP _L |
|--|-------------|------------------|------------------|-----------------|-----------------|-----------------|
| ResNet-50 [11] | 37.1 | 58.8 | 39.7 | 21.9 | 40.8 | 47.6 |
| ResNeXt-50 [52] | 38.2 | 60.5 | 41.3 | 23.0 | 41.5 | 48.8 |
| RandWire-WS, $C=109$ | 39.9 | 61.9 | 43.3 | 23.6 | 43.5 | 52.7 |
| ResNet-101 [11] | 39.8 | 61.7 | 43.3 | 23.7 | 43.9 | 51.7 |
| ResNeXt-101 [52] | 40.7 | 62.9 | 44.5 | 24.4 | 44.8 | 52.7 |
| RandWire-WS, $C=154$ | 41.1 | 63.1 | 44.6 | 24.6 | 45.1 | 53.0 |

Table 5. **COCO object detection** results fine-tuned from the networks in Table 3, reported on the val2017 set. The backbone networks have comparable FLOPs to ResNet-50 or ResNet-101.

schedule and label smoothing, for fair comparisons.

Table 3 compares RandWire with ResNet and ResNeXt under similar FLOPs as ResNet-50/101. Our mean accuracies are respectively 1.9% and 1.3% higher than ResNet-50 and ResNet-101, and are 0.6% higher than the ResNeXt counterparts. Both ResNe(X)t and RandWire can be thought of as hand-designed, but ResNe(X)t is based on designed wiring patterns, while RandWire uses a designed stochastic generator. These results illustrate different roles that manual design can play.

Larger computation. For completeness, we compare with the most accurate NAS-based networks, which use more computation. For simplicity, we use *the same trained networks* as in Table 3, but only increase the test image size to 320 \times 320 without retraining. Table 4 compares the results.

Our networks have mean accuracy 0.7%~1.3% lower than the most accurate NAS results, but ours use only $\sim 2/3$ FLOPs and $\sim 3/4$ parameters. Our networks are trained for 100 epochs and not on the target image size, *vs.* the NAS methods which use >250 epochs and train on the target 331 \times 331 size. Our model has no search on operations, unlike NAS. These gaps will be explored in future work.

COCO object detection. Finally, we report the transferability results by fine-tuning the networks for COCO object detection [25]. We use Faster R-CNN [35] with FPN [24] as the object detector. Our fine-tuning is based on 1 \times setting of the publicly available Detectron [9]. We simply replace the backbones with those in Table 3 (regular regime).

Table 5 compares the object detection results. A trend is observed similar to that in the ImageNet experiments in Table 3. These results indicate that the features learned by our randomly wired networks can also transfer.

5. Conclusion

We explored randomly wired neural networks driven by three classical random graph models from graph theory. The results were surprising: the mean accuracy of these models is competitive with hand-designed *and* optimized models from recent work on neural architecture search. Our exploration was enabled by the novel concept of a *network generator*. We hope that future work exploring new generator designs may yield new, powerful networks designs.

A. Appendix

Mapping a NAS cell to a graph. If one maps a combining op (e.g., addition, concatenation) to a node, and a unary transformation (e.g., 3×3 conv, 5×5 conv, identity) to an edge (Figure 7, right), then all cells in the NAS search space share this property: *internal nodes all have precisely input degree 2 and output degree 1*. This is an implicit prior induced by the design.

The mapping from the NAS cell to a graph is not unique. One may map both combining and unary transformations to nodes, and data flow to edges (Figure 7, left). The above property on the NAS search space can be instead described as: *internal merging nodes all have precisely input degree 2 and output degree 1*.

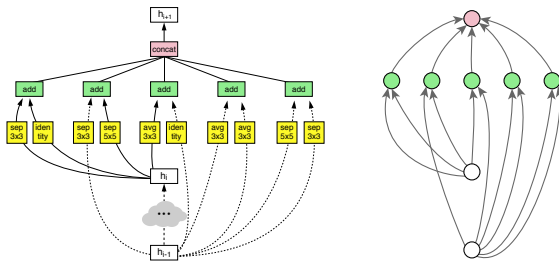


Figure 7. Mapping a NAS cell (left, credit: [56]) to a graph (right).

Converting undirected graphs into DAGs. ER, BA, and WS models generate random undirected graphs. We convert them to DAGs using a simple heuristic: we assign indices to all nodes in a graph, and set the direction of every edge as pointing from the smaller-index node to the larger-index one. This heuristic ensures that there is no cycle in the resulted directed graph. The node indexing strategies for the models are — ER: indices are assigned in a random order; BA: the initial M nodes are assigned indices 1 to M , and all other nodes are indexed following their order of adding to the graph; WS: indices are assigned sequentially in the clockwise order.

References

- [1] R. Albert and A.-L. Barabási. Statistical mechanics of complex networks. *Reviews of modern physics*, 74(1):47, 2002. 1, 2, 4, 5
- [2] D. S. Bassett and E. Bullmore. Small-world brain networks. *The neuroscientist*, 12(6):512–523, 2006. 2
- [3] D. S. Bassett and O. Sporns. Network neuroscience. *Nature neuroscience*, 20(3):353, 2017. 2
- [4] E. Bullmore and O. Sporns. Complex brain networks: graph theoretical analysis of structural and functional systems. *Nature reviews neuroscience*, 10(3):186, 2009. 2
- [5] F. Chollet. Xception: Deep learning with depthwise separable convolutions. In *CVPR*, 2017. 4
- [6] P. Erdős and A. Rényi. On the evolution of random graphs. *Publ. Math. Inst. Hung. Acad. Sci*, 5(1):17–60, 1960. 1, 2, 4, 5
- [7] J. A. Fodor and Z. W. Pylyshyn. Connectionism and cognitive architecture: A critical analysis. *Cognition*, 28(1-2):3–71, 1988. 1
- [8] E. N. Gilbert. Random graphs. *The Annals of Mathematical Statistics*, 30(4):1141–1144, 12 1959. 5
- [9] R. Girshick, I. Radosavovic, G. Gkioxari, P. Dollár, and K. He. Detectron, 2018. 8
- [10] P. Goyal, P. Dollár, R. Girshick, P. Noordhuis, L. Wesolowski, A. Kyrola, A. Tulloch, Y. Jia, and K. He. Accurate, large minibatch SGD: Training ImageNet in 1 hour. *arXiv:1706.02677*, 2017. 6, 7, 8
- [11] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 1, 2, 3, 5, 6, 8
- [12] K. He, X. Zhang, S. Ren, and J. Sun. Identity mappings in deep residual networks. In *ECCV*, 2016. 4
- [13] G. E. Hinton, N. Srivastava, A. Krizhevsky, I. Sutskever, and R. R. Salakhutdinov. Improving neural networks by preventing co-adaptation of feature detectors. *arXiv:1207.0580*, 2012. 8
- [14] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 1997. 2
- [15] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam. MobileNets: Efficient convolutional neural networks for mobile vision applications. *arXiv:1704.04861*, 2017. 6, 8
- [16] G. Huang, Z. Liu, L. van der Maaten, and K. Q. Weinberger. Densely connected convolutional networks. In *CVPR*, 2017. 1, 2, 3, 6
- [17] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015. 4
- [18] M. Kochen. *The Small world*. Ablex Pub., 1989. 2, 5
- [19] A. Krizhevsky, I. Sutskever, and G. Hinton. ImageNet classification with deep convolutional neural networks. In *NIPS*, 2012. 1, 5
- [20] Q. Le and B. Zoph. Using machine learning to explore neural network architecture. <https://ai.googleblog.com/2017/05/using-machine-learning-to-explore.html>, 2017. 2, 3
- [21] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel. Backpropagation applied to handwritten zip code recognition. *Neural computation*, 1989. 2
- [22] Y. LeCun, J. S. Denker, and S. A. Solla. Optimal brain damage. In *Advances in neural information processing systems*, 1990. 7
- [23] L. Li and A. Talwalkar. Random search and reproducibility for neural architecture search. *arXiv:1902.07638*, 2019. 2, 3
- [24] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie. Feature pyramid networks for object detection. In *CVPR*, 2017. 8
- [25] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft COCO: Common objects in context. In *ECCV*. 2014. 8

- [26] C. Liu, B. Zoph, M. Neumann, J. Shlens, W. Hua, L.-J. Li, L. Fei-Fei, A. Yuille, J. Huang, and K. Murphy. Progressive neural architecture search. In *ECCV*, 2018. 1, 2, 7, 8
- [27] H. Liu, K. Simonyan, and Y. Yang. DARTS: Differentiable architecture search. In *ICLR*, 2019. 1, 2, 7, 8
- [28] I. Loshchilov and F. Hutter. SGDR: Stochastic gradient descent with warm restarts. In *ICLR*, 2017. 6
- [29] R. Luo, F. Tian, T. Qin, E. Chen, and T.-Y. Liu. Neural architecture optimization. In *NIPS*, 2018. 1, 2
- [30] N. Ma, X. Zhang, H.-T. Zheng, and J. Sun. Shufflenet v2: Practical guidelines for efficient cnn architecture design. In *ECCV*, 2018. 8
- [31] M. L. Minsky. *Theory of neural-analog reinforcement systems and its application to the brain model problem*. Princeton University., 1954. 2
- [32] V. Nair and G. E. Hinton. Rectified linear units improve restricted boltzmann machines. In *ICML*, 2010. 4
- [33] H. Pham, M. Y. Guan, B. Zoph, Q. V. Le, and J. Dean. Efficient neural architecture search via parameter sharing. In *ICML*, 2018. 1, 2
- [34] E. Real, A. Aggarwal, Y. Huang, and Q. V. Le. Regularized evolution for image classifier architecture search. *arXiv:1802.01548*, 2018. 2, 7, 8
- [35] S. Ren, K. He, R. Girshick, and J. Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. In *NIPS*, 2015. 8
- [36] F. Rosenblatt. The perceptron: a probabilistic model for information storage and organization in the brain. *Psychological review*, 65(6):386, 1958. 2
- [37] D. E. Rumelhart, G. E. Hinton, and R. J. Williams. Learning representations by back-propagating errors. *Nature*, 1986. 2
- [38] D. E. Rumelhart and J. L. McClelland. Parallel distributed processing: Explorations in the microstructure of cognition. 1986. 1
- [39] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *IJCV*, 2015. 2, 6
- [40] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *CVPR*, 2018. 8
- [41] C. Scuto, K. Yu, M. Jaggi, C. Musat, and M. Salzmann. Evaluating the search phase of neural architecture search. *arXiv:1902.08142*, 2019. 2, 3
- [42] P. Sermanet, D. Eigen, X. Zhang, M. Mathieu, R. Fergus, and Y. LeCun. Overfeat: Integrated recognition, localization and detection using convolutional networks. In *ICLR*, 2014. 1
- [43] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015. 1, 3, 5
- [44] C. Szegedy, S. Ioffe, and V. Vanhoucke. Inception-v4, inception-resnet and the impact of residual connections on learning. In *ICLR Workshop*, 2016. 2, 6
- [45] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *CVPR*, 2015. 2
- [46] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna. Rethinking the inception architecture for computer vision. In *CVPR*, 2016. 2
- [47] A. M. Turing. Intelligent machinery. 1948. 2
- [48] L. R. Varshney, B. L. Chen, E. Paniagua, D. H. Hall, and D. B. Chklovskii. Structural properties of the caenorhabditis elegans neuronal network. *PLoS computational biology*, 7(2), 2011. 2
- [49] A. Veit, M. Wilber, and S. Belongie. Residual networks behave like ensembles of relatively shallow network. In *NIPS*, 2016. 7
- [50] D. J. Watts and S. H. Strogatz. Collective dynamics of small-world networks. *Nature*, 393(6684):440, 1998. 1, 2, 4, 5
- [51] D. B. West et al. *Introduction to graph theory*, volume 2. Prentice hall Upper Saddle River, NJ, 1996. 1, 2
- [52] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He. Aggregated residual transformations for deep neural networks. In *CVPR*, 2017. 8
- [53] M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional neural networks. In *ECCV*, 2014. 1
- [54] X. Zhang, X. Zhou, M. Lin, and J. Sun. ShuffleNet: An extremely efficient convolutional neural network for mobile devices. In *CVPR*, 2018. 6, 8
- [55] B. Zoph and Q. V. Le. Neural architecture search with reinforcement learning. In *ICML*, 2017. 1, 2, 3
- [56] B. Zoph, V. Vasudevan, J. Shlens, and Q. V. Le. Learning transferable architectures for scalable image recognition. In *CVPR*, 2018. 1, 2, 3, 5, 7, 8, 9