

Juvenile state hypothesis: What we can learn from lottery ticket hypothesis researches?

Di Zhang

University of Science and Technology of China
trotsky@mail.ustc.edu.cn

Abstract

The proposition of lottery ticket hypothesis revealed the relationship between network structure and initialization parameters and the learning potential of neural networks. The original lottery ticket hypothesis performs pruning and weight resetting after training convergence, exposing it to the problem of forgotten learning knowledge and potential high cost of training. Therefore, we propose a strategy that combines the idea of neural network structure search with a pruning algorithm to alleviate this problem. This algorithm searches and extends the network structure on existing winning ticket sub-network to producing new winning ticket recursively. This allows the training and pruning process to continue without compromising performance. A new winning ticket sub-network with deeper network structure, better generalization ability and better test performance can be obtained in this recursive manner. This method can solve: the difficulty of training or performance degradation of the sub-networks after pruning, the forgetting of the weights of the original lottery ticket hypothesis and the difficulty of generating winning ticket sub-network when the final network structure is not given. We validate this strategy on the MNIST and CIFAR-10 datasets. And after relating it to similar biological phenomena and relevant lottery ticket hypothesis studies in recent years, we will further propose a new hypothesis to discuss which factors that can keep a network juvenile, i.e., those possible factors that influence the learning potential or generalization performance of a neural network during training.

Introduction

After the proposition of the lottery ticket hypothesis (Frankle and Carbin 2019), more effective methods for neural network pruning appeared, and at the same time, new methods for neural network architect searching (NAS) were emerging. NAS algorithm searching for the network’s architectural parameters to construct the mapping function from the network’s input to the network’s output (Ren et al. 2021). Network pruning reduces the total number of the network’s parameters by removing those parameters of lesser importance from the network weight parameters (LeCun, Denker,

and Solla 1990). The main difference between NAS and network printing is that the former is based on the network’s architecture, while the latter is based on the network’s parameters. To some extent, they applied similar procedures to neural networks since network pruning can cut unnecessary links between neurons by eliminating the network’s unimportant parameters and finding a better architecture for the network.

Therefore, we can consider the pruning algorithm as a network structure optimization operator. The pruning algorithm optimizes the inter-neuron connection structure at a finer granularity when the neural architecture search algorithm for the network layer width and the path of forwarding propagation flow.

However, in practice, pruning algorithms are usually used to prune neural networks trained close to convergence to optimize their inference speed and memory efficiency (Liang et al. 2021), and many researchers point out that continuing training on the neural network after pruning could cause performance degradation on the validation set, and that training from scratch is hard using pruned structures (Li et al. 2016; Han et al. 2015). To address the latter issue, Frankle et al. proposed the lottery ticket hypothesis (LTH) that the initialization is a key point for one neural network and its sub-networks to get better performance on the validation set during a similar training process. By pruning the neural networks close to the training fit, the unmasked network weights are re-trained after re-adopting the parameters from the original initialization, and the process is repeated until obtaining a sub-network with the required amount of parameters and accuracy.

Even so, in large neural networks, the overhead of the pruning-reset-retraining cycle is unaffordable, and Frankle et al. then tried using weights after trained for iterations to reset the pruned network (Frankle et al. 2019; Renda, Frankle, and Carbin 2020) instead of the totally original initial weights. But the recovery of the original parameters and weights means forgetting of most knowledge learned in training, even some researcher assumed that pruned architecture keeps the memory of the network’s inductive bias (Frankle and Carbin 2019; Cohen and Shashua 2016).

But after trying to add new trainable layers and parameters to pruned near-convergent neural networks, also known as *winning tickets*, we found that the tested network could

continue training until convergence nearly at a higher test accuracy.

Therefore, we proposed a new iterative pruning algorithm that is based on a generalized inference of LTH. This pruning algorithm combines a simple network architecture search Process and pruning-retraining cycle into a pruning-growing-retraining approach to shape a deep and thin neural network to meet the requirement of parameter’s amount and test accuracy.

Then we further generalize the lottery ticket hypothesis into recursive lottery ticket hypothesis, and by this hypothesis and other inferences of lottery ticket hypotheses, we will further explore the juvenile state of neural networks beyond the luckiness of initialization to try to explain what determines the learning potential and convergence speed of neural networks.

Approach

Lottery Ticket Hypothesis

First, we would have a review of the naive Lottery ticket hypothesis (LTH).

Notion We follow the description of the lottery ticket hypothesis by Frankle et al. (Frankle and Carbin 2019). but with an additional symbol θ to represent the architecture of the network and shaping the vector space of weight parameters. Given a feed-forward neural network f with architecture θ , weights ω and taking x as input.

Furthermore, the weight parameter ω defined a set of parameters ω_i for each layer i . The architecture θ defines the number of neurons, width w in each layer i and the number of layers L , as well as the path of forwarding propagation flow especially inter-neuron connections through mask matrix m_i in each layer i .

Formal Definitions At start, the network weight parameters ω are randomly initialized with $\omega^0 \sim \mathcal{D}_\omega$ to adapt the structure θ^0 of the network. Then, the network $f(x; \theta^0, \omega^0)$ trained on a dataset \mathbf{D} and optimized by a gradient descent algorithm optimizer \mathcal{L} close to convergence.

Then the network $f(x; \theta^0, \omega^1)$ is pruned by removing the parameters of lesser importance from the network weight parameters ω , and produce mask matrices m_i for each layer i to define inter-neuron connections between layers.

The cycle of training and pruning could execute iteratively until the final network $f(x; \theta^*, \omega^*)$ achieve the expected test accuracy Acc^* and compression ratio r_{comp}^* of the number of parameters or meet the maximum iterations j , and produce a mask matrix m .

The network $f(x; \theta^*, \omega^*)$ is said to be a *lottery ticket* network if the following conditions are satisfied:

1. $\frac{\|\omega^*\|_0}{\|\omega^0\|_0} \leq r_{comp}^*$
2. Existing a (θ', ω') with the same architecture of θ^* and the same number of parameters of ω^* and a policy p , follow this policy p , we can training network $f(x; \theta', \omega')$ close to convergence, and achieve a new test accuracy Acc' close enough to Acc^* .

So, We take the identifying method which frankle et al. used in their work (Frankle and Carbin 2019) as the weight

reset policy, its work process can be described by the following pseudo-code:

Algorithm 1: Weight reset policy

- 1 **Given** $f(x; \theta^0, \omega^0)$
 - Result:** *winning ticket* $f(x; \theta', \omega')$
 - 2 initialize $\omega^0 \sim \mathcal{D}_\omega$ to adapt the structure θ^0
 - 3 record (θ^0, ω^0)
 - 4 training and pruning $f(x; \theta^0, \omega^0)$ for j iterations, arriving at $f(x; \theta^*, \omega^*)$
 - 5 $\omega' \leftarrow m \odot \omega^0$, reset ω^* to ω^0 according to the mask matrix m , and denote as ω'
 - 6 produce *winning ticket* $f(x; \theta', \omega')$
 - 7 **(Optional)** Verifying the *winning ticket* $f(x; \theta', \omega')$
 - 8 training $f(x; \theta', \omega')$ for j iterations and get a new test accuracy Acc'
 - 9 **If** $\frac{\|\omega'\|_0}{\|\omega^0\|_0} \leq r_{comp}^*$ **And** Acc' is closing enough to Acc^*
 - 10 **Then** $f(x; \theta', \omega')$ is a *winning ticket*
 - 11 **End**
-

With this policy, we can solve the accuracy degradation problem when we continue to train the pruned neural network.

That is, by resetting the network weights to the original initialization values, we can retrain the network with a smaller number of parameters and keep the final test accuracy with a little degradation or even better compared to the network with a larger number of parameters.

Recursive Lottery Ticket Hypothesis

Then, we propose the recursive lottery ticket hypothesis (RLTH), which is a generalization of the former naive lottery ticket hypothesis allows us to obtain a new *winning ticket* with more trainable parameters from an existing *winning ticket* by a recursive way like pruning-growing-retraining cycle.

Up to this point, we have used the pruning algorithm \mathcal{L}_{prun} only for optimizing the number of parameters of a trained neural network with a fixed structure θ .

In the following part, we will go a step further and use the pruning algorithm \mathcal{L}_{prun} as a fine-grained structure optimization operator to get a sub-network with optimal structure from original network. And try to extend the existed *winning ticket* with new neural network layers and trainable parameters, and obtain a new *winning ticket* recursively by a similar training-pruning process.

Under the weight reset strategy, the network weights must be reset if we want continue training on the pruned structure without compromising the test accuracy, but it implies discarding and forgetting the learned knowledge while preserving only a little learned inductive bias preserved by the structural information.

Assuming that we have already obtained a *winning ticket* sub-network by the original lottery ticket hypothesis, we can make following assumptions to induce the recursive lottery ticket hypothesis, which is a generalized inference of original lottery ticket hypothesis.

Inference 1: When neural network $f(x; \theta^*, \omega^*)$ is a *winning ticket* obtained following a policy p from neural network $f(x; \theta^0, \omega^0)$, and m is a mask matrix, a *winning ticket* structure parameterized by θ^* is a sub-network of structure parameterized by θ^0 , and on their corresponding graphical representations G^*, G^0 , Also meet that G^* is a *spanning subgraph* of G^0 .

Then we can construct a new father graph of that *winning ticket* by add new nodes and links to existing graph through specific network layers with new trainable parameters.

Inference 2: Given a *winning ticket* network $f(x; \theta^*, \omega^*)$, there existing a **specific** network which structure parameterized by θ' with weights ω' .

by adding that, we can construct a new neural network $f(x; \theta^1, \omega^1)$, where

$$\begin{cases} \theta^1 = \theta^* \cup \theta' \\ \omega^1 = \omega^* \cup \omega' \end{cases}$$

and according to the original LTH, neural network $f(x; \theta^*, \omega^*)$ can also be a *winning ticket* sub-network of $f(x; \theta^1, \omega^1)$, only if there exist a policy p and a pruning algorithm \mathcal{L} could train and prune $f(x; \theta^1, \omega^1)$ into $f(x; \theta^*, \omega^*)$.

Further more, we can also learn a simple fact from original LTH, that one neural network could have more than one *winning ticket* sub-networks. Since our goal is to get a new *winning ticket* sub-network, if our training policy and pruning algorithm meet the requirements fortunately, we can relax the premise of "specific extending layer" of Inference 2. Then the old *winning ticket* sub-network can be extended with a new randomly initialized layer, and thus obtain a new *winning ticket* through above training-pruning cycle.

So in summary, we can draw a generalized conclusion from the original lottery ticket hypothesis.

inference 3 Given a *winning ticket* neural networks $f(x; \theta^*, \omega^*)$ and a randomly initialized extending layer which structure parameterized by θ' with weights ω' , there exist a policy p and a pruning algorithm \mathcal{L} could train and prune its extended network $f(x; \theta^1, \omega^1)$ and produce a new *winning ticket* sub-network.

That is the Recursive lottery ticket hypothesis (RLTH). With this conclusion we can start training with a simple small neural network and expand it after getting a *winning ticket* sub-network. And repeat for several times, eventually getting a deeper *winning ticket* recursively without forgetting of learned knowledge during the training process.

Formally, the basic training policy for the recursive lottery ticket hypothesis can be outlined in the following form.

This approach allows the neural network to grow like an oyster, growing to a deeper structure and keeping the structurally optimal with a proper parameters amount at the same time.

Experiment

RLTH on Dense Networks

First, we test the weight reset policy of the original LTH and the basic structure growing policy of RLTH on the MNIST

Algorithm 2: Basic structure growing policy

```

1 Given  $f(x; \theta^0, \omega^0)$ 
   Result: winning ticket  $f(x; \theta', \omega')$ 
2 initialize  $\omega^0 \sim \mathcal{D}_\omega$  to adapt the structure  $\theta^0$ 
3 while sub-network can't meet requirements do
4   training and pruning  $f(x; \theta^0, \omega^0)$  for  $j$ 
   iterations, arriving at  $f(x; \theta^*, \omega^*)$ 
5   initialize a new layer  $\omega^1 \sim \mathcal{D}_\omega$  to adapt the
   structure  $\theta^1$ 
6    $\theta^1 \leftarrow \theta^* \cup \theta', \omega^1 \leftarrow \omega^* \cup \omega'$ 
7   construct a new sub-network  $f(x; \theta', \omega')$ 
8 end
9 produce winning ticket  $f(x; \theta', \omega')$ 
10 (Optional) Verifying the winning ticket  $f(x; \theta', \omega')$ 
11 training  $f(x; \theta', \omega')$  for  $j$  iterations and get a new test
   accuracy  $Acc'$ 
12 If  $\frac{\|\omega'\|}{\|\omega^0\|} \leq r_{comp}^*$  And  $Acc'$  is closing enough to  $Acc^*$ 
13 Then  $f(x; \theta', \omega')$  is a winning ticket
14 End

```

dataset (Lecun et al. 1998). The experiment was executed on Google Colab platform.

Then we iteratively training neural networks for 5 times to the same parameter amount by former two policies without any hyperparameter tunings, this process is a circle of training, obtain *winning ticket* and training again.

The neural network to be trained with weight reset policy iteratively is a random-initialized dense network. The network have a input layer project the input image to a vector of 128 elements. And those 5 hidden layers have been set to a linear layer with a shape of (128,128) and a ReLU activation layer. Finally, the output layer produce the category probability of the output image with a Softmax activation layer.

And The neural network to be trained with structure growing policy have the same input and output layer and have no hidden layers at begin. When iteratively training, hidden layers would be insert to the network via so-called structure growing operation. To simplify the structure growing operation, we insert a linear layer with a shape of (128,128) and a ReLU activation layer per iteration. And to avoid the interruption of the training process, the inserted layer was initialized with a identity matrix for weight and a zero matrix for bias.

Notion that, in more complicated scenarios, the structure growing operation can be implement more properly by a sophisticated methods like other NAS algorithms.

For the convenience of reproduction, The other settings of the experiment are as follows. All random seeds set to 42. Optimizer of networks was Adam (Kingma and Ba 2017) with a initial learning rate of 0.001. And the loss function is using categorical cross entropy. The pruning algorithm is Level pruner with a sparsity ratio of 95% from NNI toolbox (Microsoft 2021).

The loss and accuracy changes during training are shown in the Figure 1 and Figure 2, different color referring to dif-

ferent training iterations.

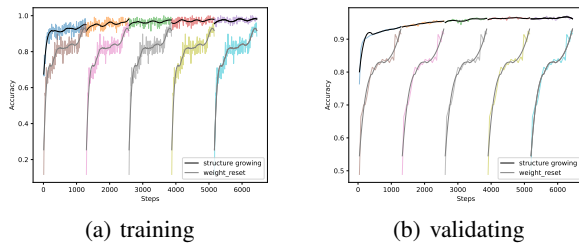


Figure 1: Accuracy on MNIST Dataset

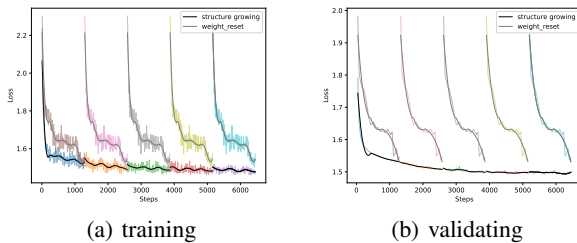


Figure 2: Loss on MNIST Dataset

RLTH on ResNet and Convolution Networks

This part of experiment was executed on CIFAR-10 Dataset (Krizhevsky, Hinton et al. 2009). To fitting the multi-channel image data of this dataset. The input and output layer was different from above. And training process was iteratively executed for 7 times.

The neural network to be trained with weight reset policy iteratively is a random-initialized residual convolution neural network. The network has an input layer with input channel of 3 and output channel of 8, and kernel size is 7. Thus those 7 hidden layers have been set to basicblocks of ResNet (He et al. 2016) with input channel of 8 and output channel of 8. Finally, the output layer produces the category probability of the output image with a Softmax activation layer.

And the neural network to be trained with structure growing policy has the same input and output layer and has no hidden layers at the beginning. When iteratively training, hidden layers would be inserted into the network via so-called structure growing operation, too. We insert a basicblock of ResNet with input channel of 8 and output channel of 8 per iteration, the convolution kernels of the inserted layer were initialized with an identity matrix for weight.

The other settings of the experiment are the same with above except the sparsity ratio changed to 80% and the pruning algorithms only compress hidden layers of convolution basicblocks.

The loss and accuracy changes during training are shown in Figure 3 and Figure 4, different colors refer to different training iterations.

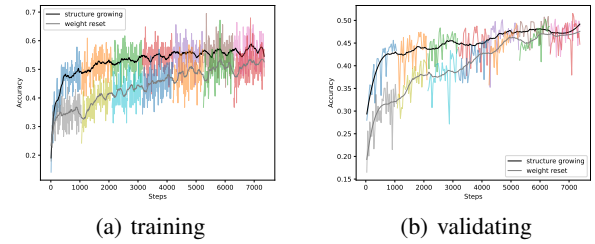


Figure 3: Accuracy on CIFAR-10 Dataset

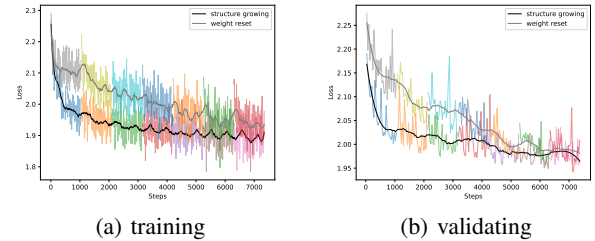


Figure 4: Loss on CIFAR-10 Dataset

From the above two parts of the experiments, we can see that although the weight reset policy avoids the training difficulties of the pruned neural network, its forgetting of knowledge is irrecoverable, and the forgetting slows down the training process and the final test performance when the *winning ticket* sub-networks are repeatedly acquired and trained. In contrast, the structural growth strategy avoids the training difficulties of the pruned neural network by adding new trainable parameters after the *winning ticket* sub-networks are acquired, and ensuring the transfer of the learned knowledge at the same time. This policy is useful for accelerating convergence and reducing overhead in the early training of large neural networks.

Discussion

Juvenile state of neural networks

Based on the above experiments, we can see that the weight reset policy and the structure growing policy are similar in that both strategies play a similar role in the training of neural networks and the acquisition of *winning ticket*. To clarify this role, we need to go back to the early starting point of Frankle's work, i.e., to overcome the training difficulties and test performance degradation faced by the re-training of pruned neural networks. And the weight reset policy addresses this difficulty by re-empowering the pruned neural networks to obtain the speed of convergence and test performance comparable to the original neural networks just after initialization.

Therefore, we can make a qualitative hypothesis that when the initialization method and network structure are appropriate, the original initialized neural network has a property that has been weakened during the pruning or improper

training, and it is this property that affects the so-called learning potential, also known as the ability of the network to converge at a better test accuracy with an appropriate training policy.

As for this property, when applied to the learning process of neural networks, it has a high degree of similarity to the learning process of natural animals, so we can introduce a bio-behavioral phenomena, juvenile state, to define this property, where the learning skills and proficiency of young animals progress significantly faster than those of adults (Delacour 1994; Boesch et al. 2019). The main manifestation of this in humans is in language learning, where young children who have not acquired any language skills learn their native language significantly faster than adults who learn the same language as a foreign language (Snow 1972; Leung, Tunkel, and Yurovsky 2021). For canines, there is a golden period of training, with puppies around half a year old acquiring tricks and skills much faster than adult dogs after sexual maturity (of Veterinary Medicine Vienna 2016; Wallis et al. 2016).

For neural networks, the juvenile state is the state where the model has good learning potential and fast convergence speed. When the neural network is in the juvenile state or the juvenile state has good properties, the neural network can converge at a good test performance level by a fast speed with appropriate learning strategies and optimizers.

In contrast, the juvenile state or property of juvenile of a neural network is weakened during an inappropriate training process or model pruning, which is most intuitively manifested by the difficulty of convergence and the degradation of test performance when the neural network continues training after pruning.

winning tickets and juvenile state

Returning to the lottery ticket hypothesis, we can consider that the *winning ticket* sub-networks are the sub-networks with better learning potential among the different sub-networks obtained from the same original neural network. Thus the weight reset policy and structure growing policy, to some extent, exerted a rejuvenation of the sub-networks by modifying the sub-networks structure or weights to restore the learning potential of the sub-networks that were damaged during the pruning process.

Therefore, based on the above hypothesis and analysis, we can consider that the process of obtaining a *winning ticket* sub-networks is to obtain a sub-networks with smaller structure and smaller parameters amount from a neural network who was already converged. And at the same time, the juvenile properties of this sub-network would also close to the original just initialized neural network. A good *winning ticket* sub-networks has a good juvenile properties to ensure it can converge as fast as the original just initialized neural network with a smaller parameters amount and a comparable test performance after convergence.

Factors affecting Juvenile state

In the above analysis, we have covered some behaviors that are harmful or beneficial to the neural network juvenile state

and properties. In the following, we would discuss those factors that may affect the neural network's juvenile state.

First, the weights and initialization of the neural network. The most important contribution of the work of Frankle et al. (Frankle and Carbin 2019) is that the convergence speed and testing performance of the network are highly correlated with the initialization process of the network, and the weight resetting policy restores the property of juvenile and learning potential of the network precisely by resetting the network to the initial weights after pruning.

Second, the structure of the neural network. The impact of the structure of the neural network on the neural network juvenile properties is reflected in the common sense. When a model structure is not suitable for training data, the neural network will face the situation of training difficult or degraded test performance even from the beginning of training.

Thus, If the number of model parameters is too large, it would make the model difficult to train and slow to converge. While if the number of parameters is too small, it would make the neural network easy to fall into the local optima and cause a overfitting of the model.

In contrast, experiments on the structural growth strategy show that appropriate overparameterization of structure is helpful for improving the property of juvenile of the model.

Third, the training strategy. In our usual sense, It contains the loss function, the optimization algorithm and its hyperparameter settings, such as learning rate and momentum, as well as other settings such as batch size and data set sampling methods appeared in deep learning practices.

Specifically, just for the gradient descent method, from a viewpoint of the *loss function surface landscape* (Li et al. 2017), the projection points (or the center of gravity of multiple projection points) of the same (or the same batch) dataset sample points on the loss surface, during the iterative optimization process, are also in constant motion, and ideally most of the projection points can reach a local elevation minimum along the fastest descending line, i.e., the model converges on a local optimal solution. And models can leave the local optima by the momentum methods to reach a better approximate solution.

The proper loss function and proper training policy provided a good *loss function surface landscape*, which enables faster convergence, avoidance of local optima traps and better test performance of the neural network in the above process, which corresponds to the juvenile state assumption.

Finally, better test scores tend to reflect better generalization performance for different *winning ticket* sub-networks within the same dataset. on the other hand, the work of Ibrahim et al. (Alabdulmohsin et al. 2021; Morcos et al. 2019) shown the existence of *winning ticket* sub-networks with the ability to generalize across datasets and training policies. So, we can assume that the learning potential of a neural network and juvenile state properties on datasets also affects the neural network's generalization ability.

Conclusion

In this paper,we introduced the recursive lottery ticket hypothesis by generalizing the lottery ticket hypothesis for

neural networks. By recursively acquiring the *winning ticket* sub-networks and continuing the training after add new trainable parameters, we can avoid the training difficulty of pruned network and the knowledge forgetting of the lottery ticket hypothesis does.

From the perspective of connectionism, the recursive lottery ticket hypothesis uses the pruning algorithm as a fine-grained optimizer of the link structure between neuron layers. The structure growth operation allows the introduction of more complex neural architecture search algorithms to above procedure, thus a bridge between structural search and structural optimization have been constructed. Finally, we further propose the juvenile state hypothesis by linking previous lottery ticket hypothesis studies with similar phenomena in biological behavior, and finally we made a qualitative analysis of those factors may affecting the juvenile state of neural networks.

Future works

Due to hardware and time constraints, the datasets and pruning algorithms involved in the experiments and the architectural search operations in this paper are relatively simple and naive. In the near future, we will verify that hypothesis on more complex problems and introduce more complex algorithms for the structural operations of neural networks.

In addition, another terminology closely related to juvenile state in biology is Neoteny (Shea 1989; Price 1999), also known as juvenile state's continuation, which is a state in which an organism can maintain its physical and mental young and strong learning ability for a long time even last for life, similar to continuous learning in machine learning that can combat with catastrophic forgettings (Delange et al. 2021), and it is one of the main directions of our future work.

References

- Alabdulmohsin, I.; Markeeva, L.; Keysers, D.; and Tolstikhin, I. 2021. A Generalized Lottery Ticket Hypothesis. *arXiv:2107.06825 [cs]*. ArXiv: 2107.06825.
- Boesch, C.; Bombjaková, D.; Meier, A.; and Mundry, R. 2019. Learning curves and teaching when acquiring nut-cracking in humans and chimpanzees. *Scientific Reports*, 9(1): 1515. Bandiera_abtest: a Cc.license.type: cc-by Cg.type: Nature Research Journals Number: 1 Primary_atype: Research Publisher: Nature Publishing Group Subject.term: Behavioural ecology; Biological anthropology Subject.term.id: behavioural-ecology; biological-anthropology.
- Cohen, N.; and Shashua, A. 2016. Inductive bias of deep convolutional networks through pooling geometry. Edition: arXiv preprint arXiv: 1605.06743.
- Delacour, J. 1994. *The Memory System of the Brain*. World Scientific. ISBN 978-981-02-1021-2. Google-Books-ID: CCtqDQAAQBAJ.
- Delange, M.; Aljundi, R.; Masana, M.; Parisot, S.; Jia, X.; Leonardis, A.; Slabaugh, G.; and Tuytelaars, T. 2021. A continual learning survey: Defying forgetting in classification tasks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 1–1.
- Frankle, J.; and Carbin, M. 2019. The Lottery Ticket Hypothesis: Finding Sparse, Trainable Neural Networks. *arXiv:1803.03635 [cs]*. ArXiv: 1803.03635.
- Frankle, J.; Dziugaite, G. K.; Roy, D. M.; and Carbin, M. 2019. Stabilizing the lottery ticket hypothesis. *arXiv preprint arXiv:1903.01611*.
- Han, S.; Pool, J.; Tran, J.; and Dally, W. J. 2015. Learning both Weights and Connections for Efficient Neural Networks. *CoRR*, abs/1506.02626.
- He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 770–778.
- Kingma, D. P.; and Ba, J. 2017. Adam: A Method for Stochastic Optimization. arXiv:1412.6980.
- Krizhevsky, A.; Hinton, G.; et al. 2009. Learning multiple layers of features from tiny images.
- Lecun, Y.; Bottou, L.; Bengio, Y.; and Haffner, P. 1998. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11): 2278–2324.
- LeCun, Y.; Denker, J. S.; and Solla, S. A. 1990. Optimal brain damage. In *Advances in neural information processing systems*, 598–605.
- Leung, A.; Tunkel, A.; and Yurovsky, D. 2021. Parents Fine-Tune Their Speech to Children's Vocabulary Knowledge. *Psychological Science*, 32(7): 975–984. Publisher: SAGE Publications Inc.
- Li, H.; Kadav, A.; Durdanovic, I.; Samet, H.; and Graf, H. P. 2016. Pruning Filters for Efficient ConvNets. *CoRR*, abs/1608.08710.
- Li, H.; Xu, Z.; Taylor, G.; and Goldstein, T. 2017. Visualizing the Loss Landscape of Neural Nets. *CoRR*, abs/1712.09913.
- Liang, T.; Glossner, J.; Wang, L.; Shi, S.; and Zhang, X. 2021. Pruning and quantization for deep neural network acceleration: A survey. *Neurocomputing*, 461: 370–403.
- Microsoft. 2021. microsoft/nni. Original-date: 2018-06-01T05:51:44Z.
- Morcos, A. S.; Yu, H.; Paganini, M.; and Tian, Y. 2019. One ticket to win them all: generalizing lottery ticket initializations across datasets and optimizers. *arXiv preprint arXiv:1906.02773*.
- of Veterinary Medicine Vienna, U. 2016. You can teach an old dog new tricks, but younger dogs learn faster.
- Price, E. O. 1999. Behavioral development in animals undergoing domestication. *Applied Animal Behaviour Science*, 65(3): 245–271.
- Ren, P.; Xiao, Y.; Chang, X.; Huang, P.-Y.; Li, Z.; Chen, X.; and Wang, X. 2021. A comprehensive survey of neural architecture search: Challenges and solutions. *ACM Computing Surveys (CSUR)*, 54(4): 1–34.
- Renda, A.; Frankle, J.; and Carbin, M. 2020. Comparing rewinding and fine-tuning in neural network pruning. *arXiv preprint arXiv:2003.02389*.

Shea, B. T. 1989. Heterochrony in human evolution: The case for neoteny reconsidered. *American Journal of Physical Anthropology*, 32(S10): 69–101. [_eprint: https://onlinelibrary.wiley.com/doi/pdf/10.1002/ajpa.1330320505](https://onlinelibrary.wiley.com/doi/pdf/10.1002/ajpa.1330320505).

Snow, C. E. 1972. Mothers' Speech to Children Learning Language. *Child Development*, 43(2): 549–565. Publisher: [Wiley, Society for Research in Child Development].

Wallis, L. J.; Virányi, Z.; Müller, C. A.; Serisier, S.; Huber, L.; and Range, F. 2016. Aging effects on discrimination learning, logical reasoning and memory in pet dogs. *AGE*, 38(1): 6.