

FreezeNet

Making Proof of Work Useful

Téo Bouvard

Shaon Rahman

Derek Berger

Abstract

1 Introduction

Proof of Work is a consensus mechanism intended to deter denial-of-service attacks. The idea was first presented by Dwork and Naor^[1] in 1993, and was later formalized by Jakobsson and Juels^[2] in 1999.

The concept behind Proof of Work is to propose a challenge which can be solved with a measurable difficulty, and which solution is easy to verify.

An example implementation of a Proof of Work protocol is the Hashcash algorithm^[3], where the challenge is to find an integer value which, when appended to some given data, produces a hash digest inferior to a given threshold value. The most efficient way to find a satisfying value is to try random values until the resulting hash is below the target value. The difficulty of this challenge grows exponentially as the target value decreases. However, verifying a solution is very fast, as the only thing to do is to verify that the hash of the concatenation of the data and the solution value yields a digest value inferior to the target value.

Derivatives of Hashcash-like algorithms are used as consensus mechanisms for more than 470 cryptocurrencies^[4]. Despite its robustness, Proof of Work is virtually useless. The computational power needed to generate a solution to each challenge is never used again once the solution has been verified. Moreover, this computational power is significant. According to the International Energy Agency^[5], Bitcoin mining alone uses approximately 20-80 TWh of electricity annually, which represent the annual domestic electricity consumption of countries such as Portugal^[6].

To make Proof of Work more useful, we try to replace the Hashcash-like challenges with neural

network training, which is also a computationally intensive task. Such a task would fulfill the requirements of a Proof of Work protocol : it is a challenge with a measurable difficulty having solutions which are fast to verify. Solutions to this challenge could be verified using common metrics such as accuracy or F1-score on test datasets.

However, this opens up a larger attack surface than Hashcash-like challenges. The main attacks we would have to prevent are transfer learning attacks. In this type of attack, an attacker could pass verification tests, using pre-trained networks which was fine-tuned for the given task. To prevent such attacks, we need to devise a scheme where an attacker would have no incentive in using a pre-trained network, because it would not give him any advantage in completing the task.

We present different approaches using neural network training as an alternative to Hashcash-like challenges in Proof of Work protocols. These approaches deter transfer learning attacks by embedding information in the network. This information is used as proof that the network was trained from scratch for a specific task, and can be verified along with the task metrics when a solution is submitted.

2 Experimental setup

We use the CIFAR-10 dataset, which is an image classification dataset containing 60,000 images belonging to 10 classes. The images are 32×32 pixels with 3 color channels. We use 50,000 images for training, and 10,000 images for testing.

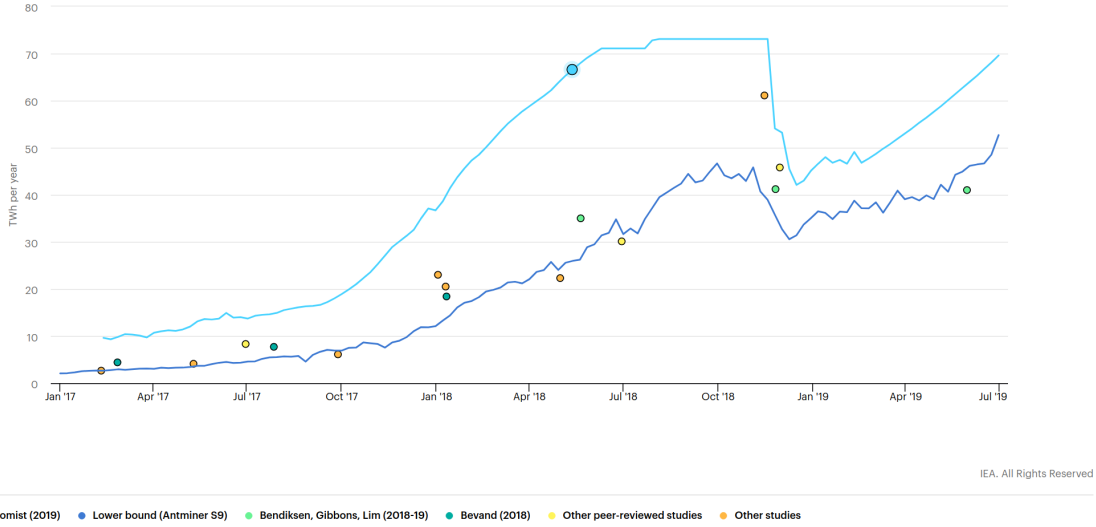


Figure 1: Bitcoin energy use estimates

3 White-box approach

In this approach, the main assumption is that we have access to the weights inside the model. We use this assumption to derive a watermarking procedure which imposes strict constraints on some weights of the model. This procedure is used to verify that a neural network was trained with a specific watermark. The watermark is derived from some data, using a cryptographic hash function and a Pseudo Random Number Generator. It can be thought of as a unique link tying a piece of data and a neural network together.

3.1 Creating a watermark

A watermark is created by hashing the data we want to embed in the network, and using the resulting digest as a seed for a Pseudo Random Number Generator. We then use this PRNG to generate a sequence of weights of a fixed size. This sequence of weights is the watermark, and it uniquely identifies the embedded data. As an example, we could use the bytes representation of a blockchain block as the embedding data.

3.2 Applying a watermark

Having generated the watermark, we now need an encoding function to apply it to a neural network.

Ideally, the encoding function should be structure agnostic, so that we could apply a watermark to a network without any constraints related to its architecture. The encoding function we use is to randomly assign the watermark’s weights to the model’s weights. The placement of the watermark weights is determined by indices drawn from the same PRNG used to generate the weights. This ties both the weights and their indices in the network to the embedded data, thus acting as a watermark. The watermark indices generated can be thought of as the indices of the model weights if they were sequentially flattened into a single one dimensional array. To apply the watermark, we simply replace the model weights at the watermark indices by the watermark weights.

3.3 Verifying a watermark

The verification process is similar to the watermarking process, except that we do not replace the weights but compute the difference between the watermark weights and the model weights. If this difference is below a given tolerance threshold, then the model is said to conform to the watermark. The tolerance is only introduced as an implementation side-effect, because of the intrinsic imprecision of floating point arithmetic. Theoretically, it would be sufficient to check for a strict equality between the watermark weights and the

model weights.

3.4 Training with a watermark

In order to use this watermarking procedure in a Proof of Work protocol, we simply apply the watermark to the network at the end of each training batch. Each time the backpropagation algorithm has updated the weights, we apply the watermark. This incidentally forces the network to learn with a strict constraint on certain weights.

3.5 Results

To conduct the following experiments, we use a simple custom baseline model comprising of 4 convolution blocks and 2 fully connected layers. The general structure of this model is presented in Figure 5.

Each convolution block consists of 2 successive convolutions followed by a max-pool and a dropout. The number of filters and the dropout ratio are functions of the convolution block index, as shown in Figure 6. This model achieves a test accuracy of 88.2% after 100 training epochs. More implementation details can be found in the notebook and the code.

The first result is that the network can still learn even when a significant proportion of its weights are constrained by the watermark. In Figure 2, experimentation show that even when the watermark size is 70% the number of weights in the model, the model reaches 65% accuracy after 100 epochs. Moreover, it seems that the learning curve is not yet reaching a plateau. This is a desired behaviour as it shows that the difficulty of training the network is proportional to the size of the watermark. Thus, the Proof of Work difficulty can be tuned by setting the desired size of watermark.

The second result is weak tampering resistance, where applying a watermark on a trained network decreases its predictive power according to an exponential decay with the watermark size. This shows that a plain model-reuse attack is not feasible if we set a watermark size large enough. In Figure 3, we show that once the watermark size represents more than 10% of the model weights, the model accuracy drops to a random classification score.

The third result is strong tampering resistance, where re-training a pre-trained model with a wa-

termark does not gives a significant advantage in solving the Proof of Work challenge. In Figure 4, experiments show that there is no incentive for an attacker to use transfer learning because honest training from scratch is as efficient. We see that in the first two-thirds of training a pre-trained model learns faster, but this advantage fades as the accuracy gets higher. Furthermore, this advantage is reduced by increasing the watermark size. For a watermark size of 70% of the model weights, honest training is nearly indiscernible from transfer attack, and the accuracy still reaches more than 60% after 100 epochs.

3.6 Possible Improvements

One possible improvement to this approach could be to apply the watermarked weights during the backpropagation, not after it has been performed. This would require lower-level access to the weights during the backpropagation algorithm, but it may lead to a faster learning curve. In our implementation, the weights are drawn from a uniform distribution in $[-1, 1]$, and the indices are drawn from a uniform distribution in all possible indices without replacement. This means that, for example, the probability of watermarking a weight in a middle hidden layer is the same as watermarking a weight in the output layer. It may be interesting to experiment with different strategies and distributions for generating these weights and indices. Another improvement would be to replace the two modules for hashing and generating weights by a single Deterministic Random Bit Generator derived from a hash function. In our case, we used SHA-256 as hash function to compute the seed, and initialized a Mersenne Twister PRNG with this seed. The alternative would be to use a Hash-DRBG^[7].

4 Black-box approach

5 Limitations

Unlike Hashcash challenges, using neural network training as Proof of Work adds a considerable overhead to the entire process. In practice, it requires the transfer of a dataset from a source to all workers in the network. Without introducing a way to generate these datasets in a decentralized manner, this process is inherently incompat-

ible with fully distributed blockchains. The overhead is also significant when workers present their solution, as they have to transfer all weights in the model leading to their solution. However we can greatly reduce the number of weights they have to transfer, as a proportion of these weights are constrained by the watermark and can be recomputed locally from the embedded data. This improvement might not be very significant a because the main overhead comes from the dataset transfer.

Appendices

References

- [1] Cynthia Dwork and Moni Naor. Pricing via processing or combatting junk mail. In Ernest F. Brickell, editor, *Advances in Cryptology — CRYPTO’ 92*, pages 139–147, Berlin, Heidelberg, 1993. Springer Berlin Heidelberg.
- [2] Markus Jakobsson and Ari Juels. *Proofs of Work and Bread Pudding Protocols(Extended Abstract)*, pages 258–272. Springer US, Boston, MA, 1999.
- [3] Adam Back. Hashcash - a denial of service counter-measure. Technical report, 2002.
- [4] CryptoSlate. Proof of work cryptocurrencies, 2020.
- [5] International Energy Agency. Bitcoin energy use estimates. 2019.
- [6] Enerdata. Electricity domestic consumption. 2020.
- [7] Elaine B. Barker and John M. Kelsey. Recommendation for random number generation using deterministic random bit generators. 2015.

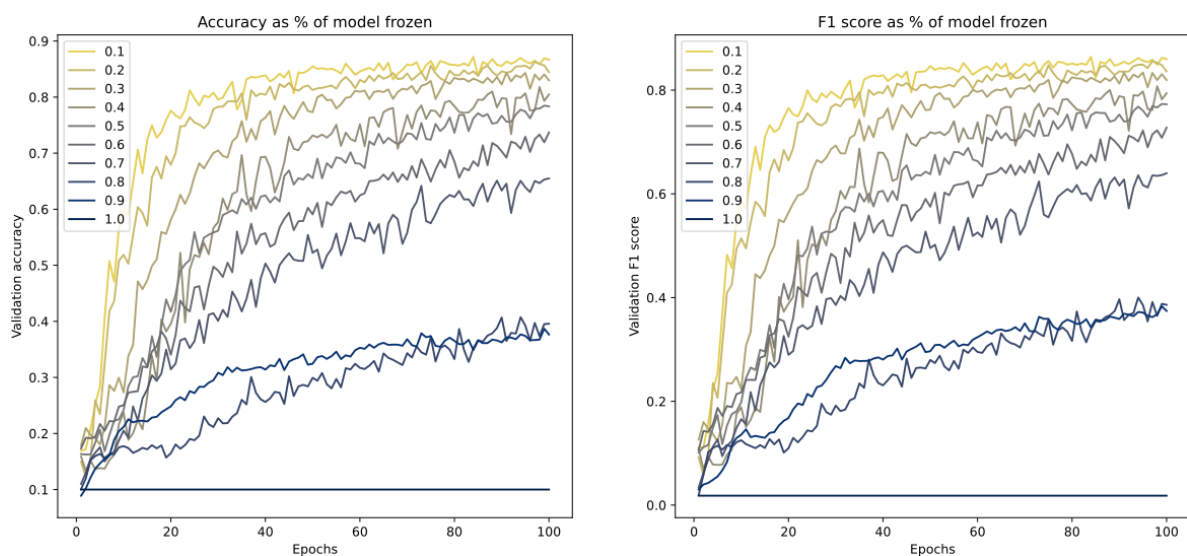


Figure 2: Training from scratch with a watermark

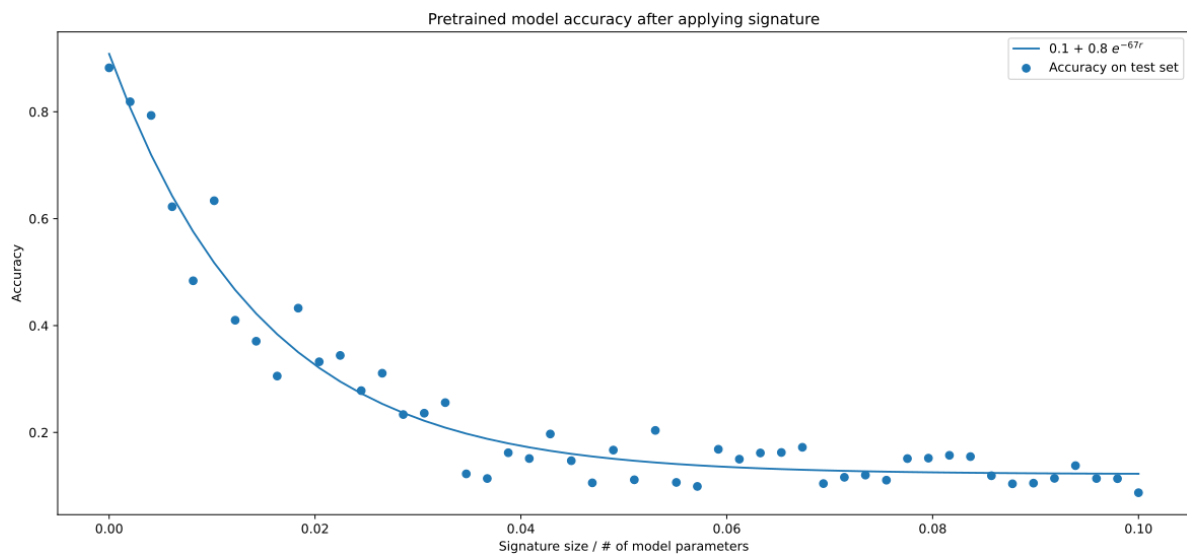


Figure 3: Pretrained model accuracy after watermarking

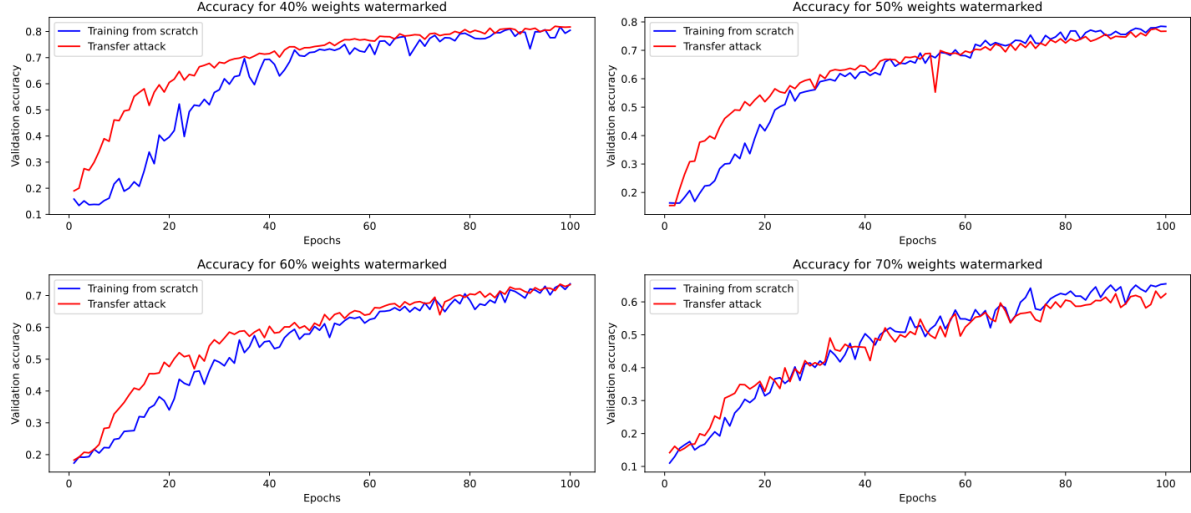


Figure 4: Transfer learning attack

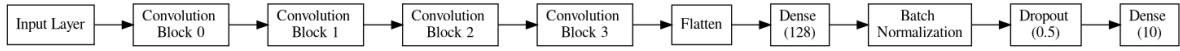


Figure 5: General model structure

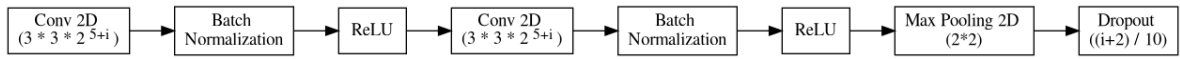


Figure 6: Convolution block