

# Recognizing Images with at most one Spike per Neuron

Christoph Stöckl<sup>1</sup>, Wolfgang Maass<sup>1,\*</sup>

January 22, 2020

<sup>1</sup>*Institute of Theoretical Computer Science, Graz University of Technology, Inffeldgasse 16b, Graz, Austria*

\* *To whom correspondence should be addressed; E-mail: maass@igi.tugraz.at.*

## Abstract

In order to port the performance of trained artificial neural networks (ANNs) to spiking neural networks (SNNs), which can be implemented in neuromorphic hardware with a drastically reduced energy consumption, an efficient ANN to SNN conversion is needed. Previous conversion schemes focused on the representation of the analog output of a rectified linear (ReLU) gate in the ANN by the firing rate of a spiking neuron. But this is not possible for other commonly used ANN gates, and it reduces the throughput even for ReLU gates. We introduce a new conversion method where a gate in the ANN, which can basically be of any type, is emulated by a small circuit of spiking neurons, with At Most One Spike (AMOS) per neuron. We show that this AMOS conversion improves the accuracy of SNNs for ImageNet from 74.60% to 80.97%, thereby bringing it within reach of the best available ANN accuracy (85.0%). The Top5 accuracy of SNNs is raised to 95.82%, getting even closer to the best Top5 performance of 97.2% for ANNs. In addition, AMOS conversion improves latency and throughput of spike-based image classification by several orders of magnitude. Hence these results suggest that SNNs provide a viable direction for developing highly energy efficient hardware for AI that combines high performance with versatility of applications.

## Introduction

Spiking neural networks (SNNs) are among the leading candidates to solve one of the major impediments of more widespread uses of modern AI: The energy consumption of the very large artificial neural networks (ANNs) that are needed. These ANNs need to be large, since they need to have a sufficiently large number of parameters in order to absorb enough information from the huge data sets on which they are trained, such as the 1.2 million images of ImageNet2012. Inference on these large ANNs is power hungry

(García-Martín et al., 2019), which impedes their deployment in mobile devices or autonomous vehicles. Spiking neurons have been in the focus of recent work on novel computing hardware for AI with a drastically reduced energy budget, because the giant SNN in the brain –consisting of about 100 billion neurons– consumes just 20W (Ling, 2001). Most spiking neuron models that are considered in neuromorphic hardware are in fact inspired by neurons in the brain. Their output is a train of stereotypical electrical pulses –called spikes. Hence their output is very different from the analog numbers that an ANN neuron produces as output.

But whereas large ANNs, trained with ever more sophisticated learning algorithms on giant data sets, approach –and sometimes exceed– human performance in several categories of intelligence, the performance of SNNs is lagging behind. One strategy for closing this gap is to design an ANN-to-SNN conversion that enables us to port the performance of a trained ANN into an SNN. The most common –and so far best performing– conversion method was based on the idea of (firing-) rate coding, where the analog output of an ANN unit is emulated by the firing rate of a spiking neuron (Rueckauer et al., 2017). This method can readily be used for ANN units that are based on the ReLU (rectified linear) activation function. It has produced impressive results for professional benchmark tasks such as ImageNet, but a significant gap to the accuracy, latency, and throughput of ANN solutions has thwarted its practical application. Problems with the timing and precision of resulting firing rates on higher levels of the resulting SNNs have been cited as possible reasons for the loss in accuracy of the SNN. In addition, the transmission of an analog value through a firing rate requires a fairly large number of time steps –typically in the order of 100, which reduces both latency and throughput for inference. An additional impediment for a rate-based ANN-to-SNN conversion emerged more recently in the form of better performing ANNs such as EfficientNet (Tan and Le, 2019). They employ the Swish function as activation function, defined by  $x \cdot \text{sigmoid}(x)$ , which contains more complex non linearities than the ReLU function. Furthermore, the Swish function also produces negative values that can not be readily encoded by the – necessarily non-negative – firing rate of a spiking neuron.

We introduce a new ANN-to-SNN conversion method that we call AMOS conversion because it requires At-Most-One-Spike (AMOS) per neuron. This method is obviously very different from rate-based conversions, and structurally more similar to temporal coding, where the delay of a single spike is used to encode an analog value. However temporal coding has turned out to be difficult to implement in a noise-robust and efficient manner in neuromorphic hardware. This arises from the difficulty to implement delays with sufficiently high precision without sacrificing latency or throughput of the SNN, and the difficulty to design spiking neurons that can efficiently process such temporal code (Maass and Natschläger, 1998), (Thorpe et al., 2001), (Rueckauer et al., 2017), (Kheradpisheh and Masquelier, 2019). In contrast, AMOS coding requires just on the order of  $\log N$  different delays for transmitting integers between 1 to  $N$ . Furthermore, these delays can be arranged in a data-independent manner that supports pipelining, so that a new image can be processed by the SNN at every time step.

We show that even the simplest type of spiking neuron, the McCulloch Pitts neuron

or threshold gate (McCulloch and Pitts, 1943) can efficiently compute with AMOS codes. Thus no temporal integration of information is needed for the spiking neuron model in order to efficiently emulate inference by ANNs. This simple version of a spiking neuron had previously already been used for image classification by SNNs in hardware (Esser et al., 2016). We will describe in the first subsection of Results the design of an AMOS unit that replaces the gate of an ANN –with basically any activation function—in this new ANN-to-SNN conversion. We then show that this conversion produces an SNN that carries out inference for classifying images from the full ImageNet2012 dataset with drastically improved accuracy, latency, and throughput. Whereas the design of the AMOS unit for the conversion of ANN-neurons with the Swish function requires training of its parameters, one can design an AMOS unit for the special case of the ReLU activation function explicitly: It reduces in this special case to an analog-to-digital conversion via binary coding. This will be made explicit in the last subsection of Results.

# 1 Results

## 1.1 Architecture of the AMOS unit

We present in Fig. 1B the architecture of an AMOS unit –consisting of  $K$  spiking neurons– that approximates a generic ANN gate with activation function  $f$  shown in Fig. 1A. Besides fixed delays (shown in green) the AMOS unit contains weight coefficients  $c_1, \dots, c_K, d_1, \dots, d_K, h_{ij}$  for  $i, j \in \{1, \dots, K\}$ , and thresholds  $T_1, \dots, T_K$  (shown in blue).

The case that the activation function  $f$  of an ANN gate outputs positive and negative numbers is of particular importance in view of the Swish function (see Fig. 5) that was introduced in (Zoph and Le, 2018) and used as activation function in EfficientNet (Tan and Le, 2019). It is defined by

$$\text{Swish}(x) = x \cdot \frac{1}{1 + e^{-x}}. \quad (1)$$

Neuron  $i$  in the AMOS unit outputs

$$z_i = \Theta(c_i \cdot x - H_i - T_i), \quad (2)$$

where  $\Theta$  is the Heaviside activation function defined by

$$\Theta(x) = \begin{cases} 0, & \text{if } x < 0 \\ 1, & \text{if } x \geq 0, \end{cases} \quad (3)$$

the coefficient  $c_i$  and the threshold  $T_i$  are trained parameters, and  $H_i$  is an inhibitory input defined by

$$H_i = \sum_{j=1}^{i-1} h_{ij} z_j \quad (4)$$

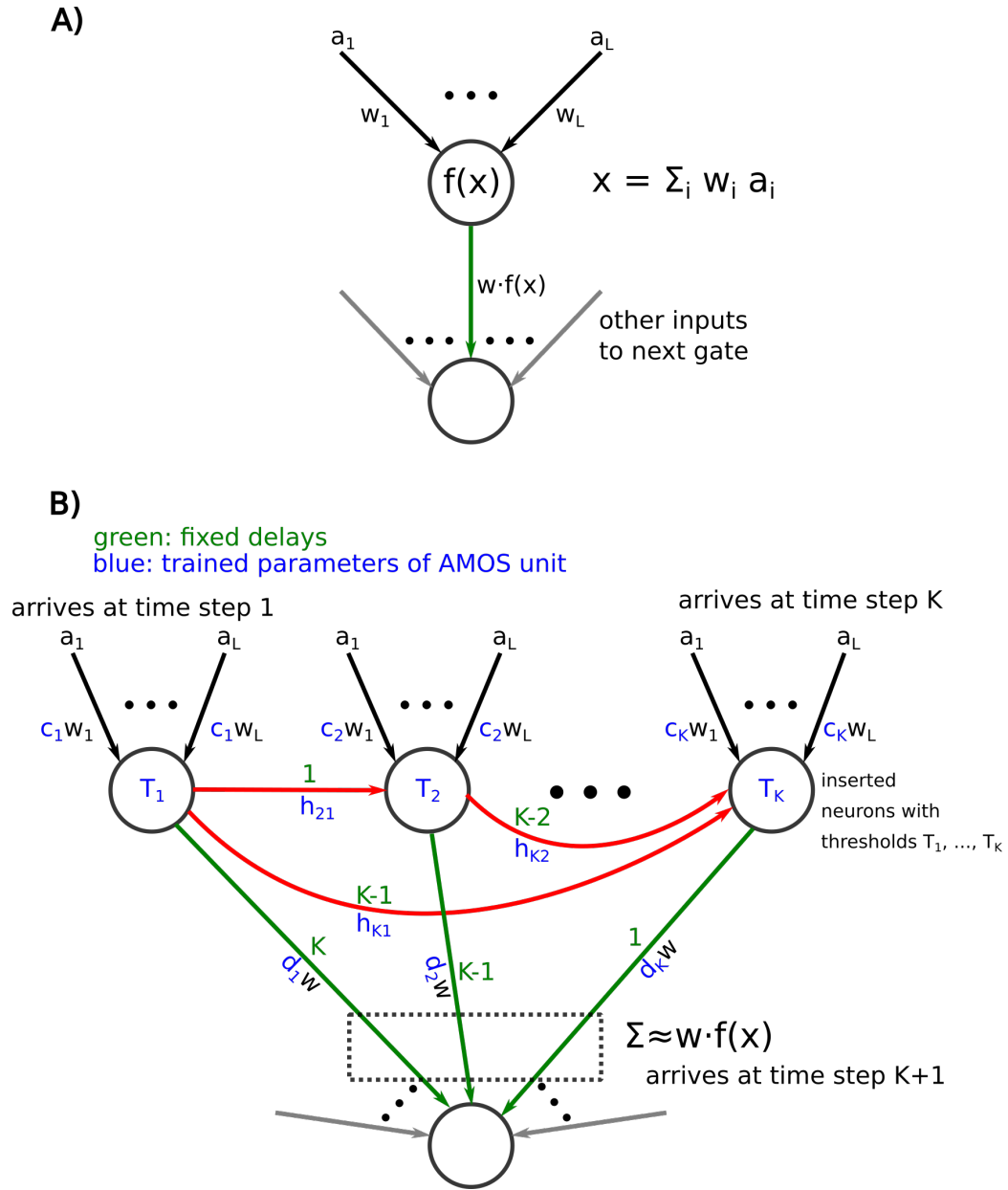
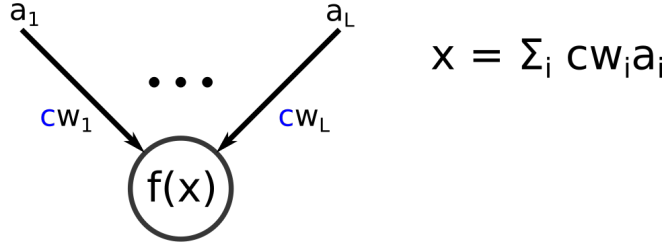


Figure 1: Architecture of an AMOS unit, consisting of  $K$  spiking neurons.  
**A)** An ANN gate with activation function  $f$  that is to be emulated by the AMOS unit.  
**B)** The AMOS unit receives the same inputs  $a_1, \dots, a_L$ , duplicated  $K$  times. It outputs after  $K + 1$  time steps an approximation of the value  $wf(x)$  which the ANN gate sends to subsequent gates.

Input motif from figure 1B:



generalized input motif for the case of two inputs:

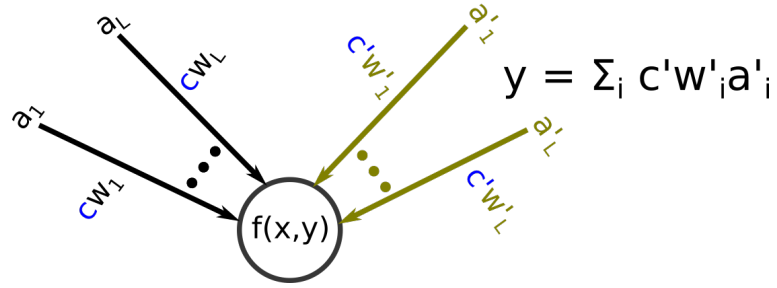


Figure 2: AMOS approximation of a function  $f(x, y)$  with two input values  $x, y$ , e.g.  $x \cdot y$

with trained negative weights  $h_{ij}$ . The output  $y$  of the AMOS unit, which is fed into subsequent AMOS units that emulate subsequent gates to which the ANN that it emulates is connected (only one subsequent gate is shown for simplicity in Fig. 1), can be written as

$$y = \sum_{i=1}^K wd_i z_j, \quad (5)$$

where the  $d_i$  are additional trained weights of the AMOS unit, and  $w$  is the weight of the corresponding connection from the ANN gate in Fig. 1A to the next ANN gate. Thus the computational function of the entire AMOS unit can be expressed as

$$\text{AMOS}(x) = y = \sum_{i=1}^K w \cdot d_i \cdot \Theta(c_i \cdot x - H_i - T_i). \quad (6)$$

For the case of functions  $f(x, y)$  with two input values  $x$  and  $y$ , one just needs to replace the motif of Fig. 1A in Fig. 1 by the motif indicated in Fig. 2B.

This changes equation 2 to:

$$z_i = \Theta(c_i \cdot x + c'_i \cdot x' - H_i - T_i) \quad (7)$$

## 1.2 Approximation of some common ANN gates by AMOS units

In the case of the ReLU activation function no training of the parameters of the AMOS unit from Fig. 1B is needed. If one defines its thresholds  $T_j$  by  $2^{K-j}$  and weights  $d_j, c_j, h_{ij}$  by  $c_j = 1, d_j = 2^{K-j}, h_{ij} = 2^{K-j}$ , an AMOS unit with  $K$  neurons produces an approximation  $\text{ReLU}(x)$  of the activation function  $\text{ReLU}(x)$  that deviates for  $x$  from 0 by at most  $\alpha 2^{-K}$  for an approximation in the interval from  $-\infty$  to  $\alpha$ , for any  $\alpha \in \mathbb{R}$ . The resulting approximation is plotted in Fig. 3.

For the case of other gate functions  $f$  we train the additional weights and thresholds of the AMOS unit through backpropagation, using a triangle-shaped pseudo-derivative in place of the non-existing derivative of the Heaviside function. Results of such approximations are plotted in Fig. 4 - 6 for the case of the sigmoid function, Swish function, and multiplication.

Note that the AMOS unit can be used in a pipelined manner, processing  $K$  different inputs  $x$  in its  $K$  time steps. Hence the resulting SNN can be used in a pipelined manner, processing a new network input at each time step. Hence its throughput is much better than that of SNNs that result from rate-based ANN-to-SNN conversions, such as for example (Rueckauer et al., 2017; Sengupta et al., 2019).

The number of neurons in the network is increased through the AMOS conversion by some factor. However a hardware implementation can reuse AMOS units for multiple time steps (since all AMOS units have the same internal parameters), thereby reducing the required size of the network at the cost of a corresponding reduction of the throughput.

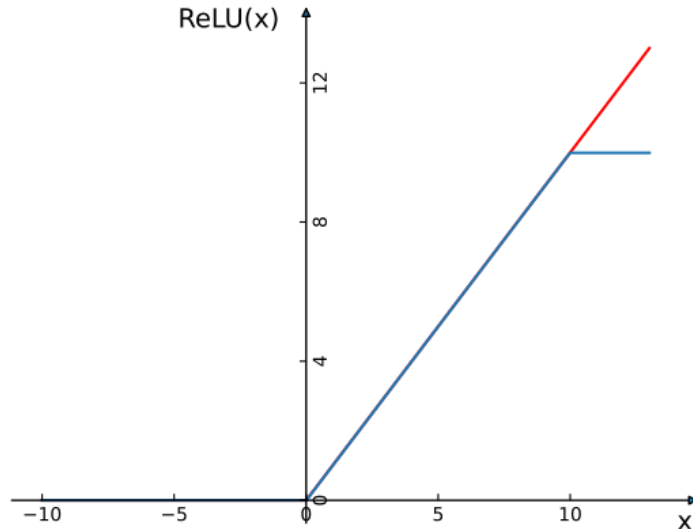


Figure 3: *AMOS approximation of the ReLU function,  $K = 10$ , (red: target function, blue: AMOS approximation)*

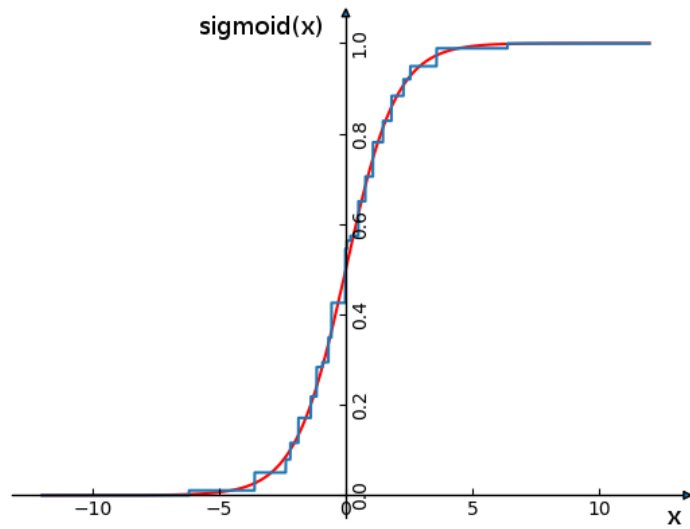


Figure 4: *AMOS approximation of the sigmoid function,  $K = 8$  (red: target function, blue: AMOS approximation)*

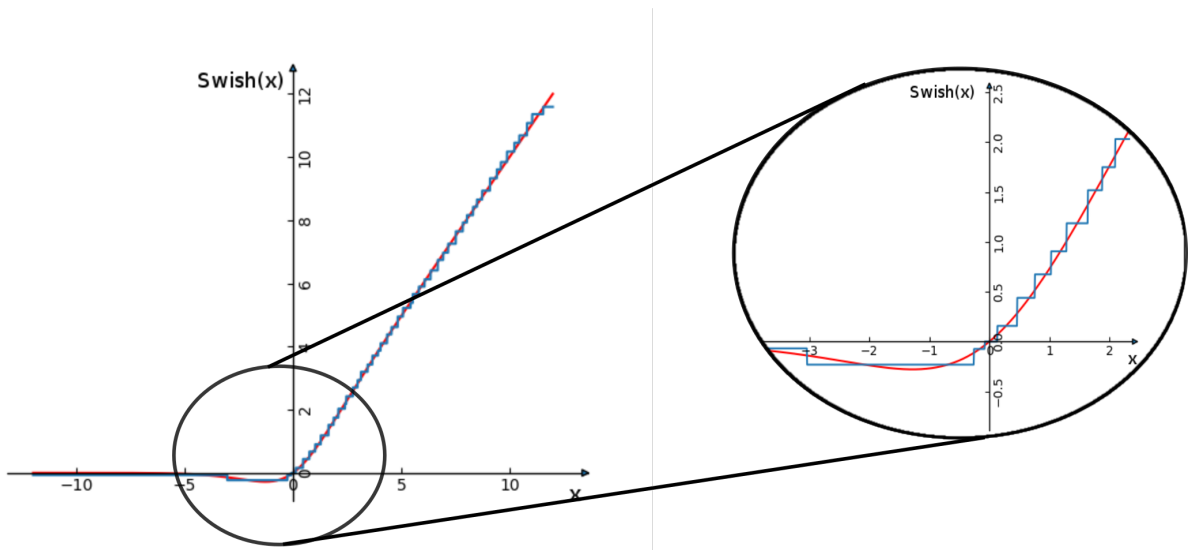


Figure 5: *AMOS approximation of the Swish function,  $K = 12$ , (red: target function, blue: AMOS approximation)*

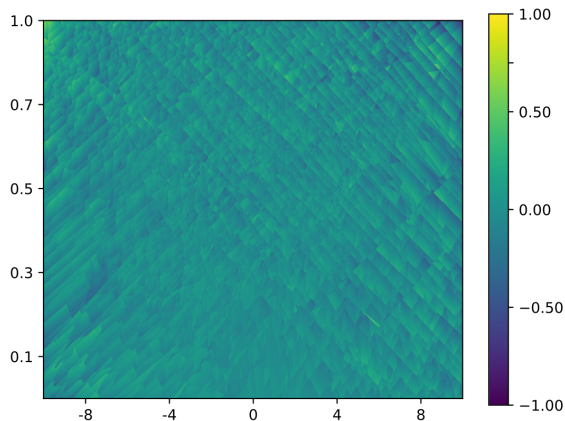


Figure 6: *Approximation error of the trained AMOS unit for multiplication with  $K = 40$  ( $MSF = 0.0102$ ).*

### 1.3 Application of the AMOS conversion to image classification with SNNs

The ImageNet data set (Russakovsky et al., 2015) has become the most popular benchmark for image classification in machine learning (we are using here the ImageNet2012 version). This data set consists of 1.281.167 training images and 50.000 test images (both RGB images of different sized), that are labeled by 1000 different categories. Classifying imaged from ImageNet is a nontrivial task even for a human, since this data set contains for example 59 categories for birds of different species and gender (Van Horn et al., 2015). This may explain why a relaxed performance measurement, where one records whether the target class is among the top 5 classifications that are proposed by the neural network ("Top5"), is typically much higher.

A new record in ANN performance for the classification of images from the ImageNet 2012 dataset was achieved by (Tan and Le, 2019). They introduced a new ANN family called EfficientNet, that achieved 84.4%. With a modified training procedure it can achieve 85% accuracy (Cubuk et al., 2019). The parameters of the trained network are publicly available\*.

This accuracy was achieved by the EfficientNet-B7 model that has 66M parameters. Besides a new scaling method for balancing network depth, width, and image resolution, they also introduced the Swish function as activation function, instead of ReLU. The Swish function emerged from automatic search for better performing activation functions (Zoph and Le, 2018). Other nonlinearities that are used in EfficientNet are sigmoids and multiplication. These occur in "Squeeze and Excitation Layers" (Hu et al., 2018) of EfficientNet,

\*<https://storage.googleapis.com/cloud-tpu-checkpoints/efficientnet/randaug/efficientnet-b7-randaug.tar.gz>



see Fig. 7. The main building block of the EfficientNet architecture is the mobile inverted bottleneck (Sandler et al., 2018), which uses depthwise separable convolutions. This type of convolutional layer uses neurons with linear activation functions. Although it would certainly be possible to approximate linear functions using AMOS conversion, we simply collapsed linear layers into the generation of the weighted sums that form the inputs to the next layers.

We evaluated the classification performance of the SNN that results from an application of the AMOS conversion described in the previous section to EfficientNet-B7. The resulting SNN achieved a classification performance of 80.97%, and 95.82% for Top5; see Table 1. The values for  $K$  used in the AMOS conversion are listed in Table 3 .

Model	# params	ANN	SNN	# layers	# neurons	# spikes
EfficientNet-B7	66M	85% (97.2%)	80.97% (95.82%)	4605	3110M	1799M
ResNet50	26M	75.22% (92.4%)	75.10% (92.36%)	500	96M	14M

Table 1: *Performance results for ANNs and the SNNs, that result from their AMOS conversion on ImageNet. Top5 accuracy is given in parentheses. #layers and #neurons refer to the SNN version of the network. Spike numbers refer to inference for a sample test image.*

Note that the resulting SNNs have an optimal throughput, since they can classify a new image at each time step. SNNs that result from a rate-based ANN conversion need up to 3000 time steps for reaching maximum accuracy, hence their throughput is by a corresponding factor smaller. The number of parameters is increased by the AMOS conversion only by a small additive term (see Table 3), since all AMOS units of the same type use the same parameters.

The accuracy of 75.22% for the ANN version of ResNet50 in Table 1 resulted from training a variant of ResNet50 where max pooling was replaced by average pooling, using the hyperparameters given in the TensorFlow repository. This accuracy is close to the best published performance of 76% for ResNet50 ANNs (Tan and Le, 2019, Table 2). Apart from max-pooling, ResNets use neither Swish nor sigmoid or multiplication as nonlinearities, just ReLU. This explains why the application of the AMOS conversion to ResNet yields SNNs whose Top1 and Top5 performance is almost indistinguishable from the ANN version. Note also that the resulting SNNs are only sparsely active, an activity regime that enhances the energy efficiency of some hardware implementations.

The best previous performance of an SNN on ImageNet was achieved by converting an Inception-v3 model (Szegedy et al., 2016) with a rate-based conversion scheme (Rueckauer et al., 2017). The reported test accuracy of the resulting SNN was 74.6%, where 550 time steps were used to simulate the model. Hence already the application of AMOS conversion to ResNet50 improves this result with regard to accuracy, and especially with regard to throughput. The AMOS conversion of EfficientNet-B7 yields an additional 5.87% accuracy improvement.

Model	ANN	SNN	# neurons	# spikes
ResNet50	92.99%	92.42%	4.751.369	647.245
ResNet20	91.58%	91.45%	1.884.160	261.779
ResNet14	90.49%	90.39%	1.310.720	190.107
ResNet8	87.22%	87.05%	737.280	103.753

Table 2: *SNN performances on CIFAR10 that result from AMOS conversions of ResNet models of different depths. (using  $K = 10$  in the AMOS units) #neurons refers to the SNN version.*

## 1.4 Results for the classification of images from the CIFAR10 data set

The results for the ANN versions of ResNet that are given in Table 2 are the outcome of training them with the hyperparameters given in the TensorFlow models repository. They are very close to the best results reported in the literature. The best ResNet on CIFAR10 is the ResNet110, where a test accuracy of 93.57% has been reported (He et al., 2016). Our ResNet50 achieves 92.99%, which is very close to the performance of the ResNet56 with 93.03%.

Spiking versions of ResNet20 have already been explored (Sengupta et al., 2019). Using a rate-based conversion scheme a performance of 87.46% was reported. Compared to these results, AMOS conversion yields a higher accuracy, while also using significantly fewer time steps, thereby reducing latency for inference. In addition, the throughput is reduced from a rather low value to the theoretically ideal value of one image per time step.

## 1.5 Trade-off between latency and network size of the SNNs

It is well-known that the extraction of bits from a weighted sum, as well as multiplication of binary numbers, can be carried out by threshold circuits –hence also by SNNs– with a small number of layers –typically 2 or 3– that does not depend on the bit length of the binary numbers involved, however at the cost of increasing the number of neurons to a low-degree polynomial of this bit length. A recent summary of such results is provided in section 3 of (Parekh et al., 2018). Hence one can replace the basic architectures of the AMOS units from Fig. 1 and 2 by the more shallow architectures that employ a larger number of spiking neurons. In order to make the resulting circuits applicable to larger input domains, one can train their parameters, similarly as for the basic architectures. Hence there exists a trade-off between the number of layers (latency) and the size of the SNN that can be achieved through an AMOS conversion, and several points on this trade-off curve can be reached through these modified AMOS conversions.

## 2 Methods

### 2.1 Squeeze and Excitation in EfficientNet

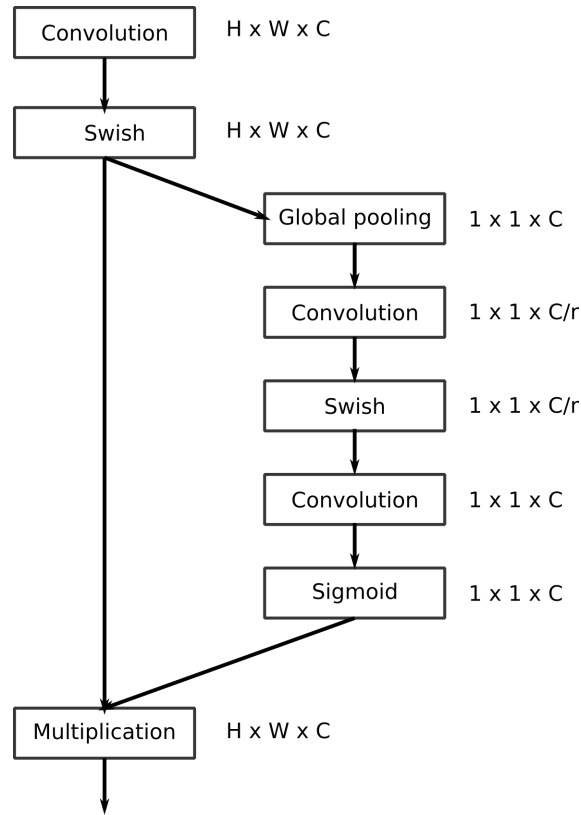


Figure 7: Squeeze and Excitation Layer

### 2.2 Number of parameters of AMOS units

Approximation	Number of parameters
Sigmoid (K=8)	52
ReLU (K= 10)	75
Swish (K=12)	102
Mult (K= 40)	940

Table 3: Number of parameters in the AMOS units that are used in this article

## Discussion

We have introduced a new method for converting ANNs to SNNs. Since the resulting SNN uses for inference at most one spike (AMOS) per neuron, AMOS conversion can be seen as

dual to the familiar rate-based conversion, since it uses space – i.e., more neurons – rather than time for replacing analog values by spikes. This conversion significantly improves accuracy, latency, and especially the throughput of the resulting SNN. Furthermore, since AMOS conversion can be applied to virtually any type of ANN gate, it demonstrates for the first time that SNNs are universal computing devices from the perspective of modern ANN-based machine learning and AI. For the classification of natural images, it raises the accuracy of SNNs for the full ImageNet 2012 dataset to 80.97% – and to 95.82% for Top5 – thereby bringing it into the range of the best ANN and human performance. This was achieved by converting EfficientNets, a recently proposed variant of ANNs that employ the Swish function as neural activation function, for which a rate-based conversion to SNNs is impossible. The resulting SNNs that achieve high performance are – like their ANN counterparts – very large. But one can sacrifice some of their accuracy by starting with a smaller ANN, or by reducing their perfect throughput of one image per step and reuse gates of a smaller SNN with online reconfigurable connections and parameters. Note that reuse of neurons would be more problematic for rate-based SNN computations, since each spiking neuron is occupied there during a fairly large and somewhat unpredictable number of time steps when emulating a computation step of an ANN gate. In contrast, AMOS conversion provides a tight bound on the number of time steps during which a spiking neuron is occupied. Hence it can also be used for converting recurrently connected ANNs to SNNs.

Altogether the proposed method for generating highly performant SNNs offers an opportunity to combine the computationally more efficient and functionally more powerful training of ANNs with the superior energy-efficiency of SNNs for inference. Note that one can also use the AMOS-converted SNN as initialization for subsequent direct training of the SNN for a more specific task. Altogether our results suggest that spike-based hardware may gain an edge in the competition for the development of drastically more energy efficient hardware for AI applications by combining energy efficiency and competitive performance with a versatility that optimized hardware for specific ANNs –such as a specific type of convolutional neural networks— cannot offer.

## Acknowledgements

We would like to thank Franz Scherr for helpful discussions. This research was partially supported by the Human Brain Project of the European Union (Grant agreement number 785907).

## References

- Cubuk, E. D., Zoph, B., Shlens, J., and Le, Q. V. (2019). RandAugment: Practical data augmentation with no separate search. *arXiv.org:1909.13719*.
- Esser, S. K., Merolla, P. A., Arthur, J. V., Cassidy, A. S., Appuswamy, R., Andreopoulos,

- A., Berg, D. J., McKinstry, J. L., Melano, T., Barch, D. R., Di Nolfo, C., Datta, P., Amir, A., Taba, B., Flickner, M. D., and Modha, D. S. (2016). Convolutional networks for fast, energy-efficient neuromorphic computing. *Proceedings of the National Academy of Sciences of the United States of America*, 113(41):11441–11446.
- García-Martín, E., Rodrigues, C. F., Riley, G., and Grahn, H. (2019). Estimation of energy consumption in machine learning. *Journal of Parallel and Distributed Computing*, 134:75–88.
- He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2016-Decem:770–778.
- Hu, J., Shen, L., and Sun, G. (2018). Squeeze-and-Excitation Networks. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 7132–7141.
- Kheradpisheh, S. R. and Masquelier, T. (2019). S4NN: temporal backpropagation for spiking neural networks with one spike per neuron.
- Ling, J. (2001). <https://hypertextbook.com/facts/2001/JacquelineLing.shtml>.
- Maass, W. and Natschläger, T. (1998). Emulation of hopfield networks with spiking neurons in temporal coding. In *Computational Neuroscience*, pages 221–226. Springer.
- McCulloch, W. and Pitts, W. (1943). A logical calculus of the ideas immanent in nervous activity. In: *Bulletin of Mathematical Biophysics*, Bd. 5, page 115–133.
- Parekh, O. D., Phillips, C. A., James, C. D., and Aimone, J. B. (2018). Constant-depth and subcubic-size threshold circuits for matrix multiplication. Technical report, Sandia National Lab.(SNL-NM), Albuquerque, NM (United States).
- Rueckauer, B., Lungu, I. A., Hu, Y., Pfeiffer, M., and Liu, S. C. (2017). Conversion of continuous-valued deep networks to efficient event-driven networks for image classification. *Frontiers in Neuroscience*, 11(DEC):1–12.
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A. C., and Fei-Fei, L. (2015). ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision*, 115(3):211–252.
- Sandler, M., Howard, A., Zhu, M., Zhmoginov, A., and Chen, L. C. (2018). MobileNetV2: Inverted Residuals and Linear Bottlenecks. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 4510–4520.

- Sengupta, A., Ye, Y., Wang, R., Liu, C., and Roy, K. (2019). Going Deeper in Spiking Neural Networks: VGG and Residual Architectures. *Frontiers in Neuroscience*, 13(1998):1–16.
- Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., and Wojna, Z. (2016). Rethinking the Inception Architecture for Computer Vision. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2016-Decem:2818–2826.
- Tan, M. and Le, Q. V. (2019). EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks.
- Thorpe, S., Delorme, A., and Rullen, R. (2001). Spike-based strategies for rapid processing. *Neural networks : the official journal of the International Neural Network Society*, 14:715–25.
- Van Horn, G., Branson, S., Farrell, R., Haber, S., Barry, J., Ipeirotis, P., Perona, P., and Belongie, S. (2015). Building a bird recognition app and large scale dataset with citizen scientists: The fine print in fine-grained dataset collection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 595–604.
- Zoph, B. and Le, Q. V. (2018). Searching for activation functions. *6th International Conference on Learning Representations, ICLR 2018 - Workshop Track Proceedings*, pages 1–13.