

Gating Revisited: Deep Multi-layer RNNs That Can Be Trained

Mehmet Ozgur Turkoglu¹, Stefano D’Aronco¹, Jan Dirk Wegner^{1,2}, and Konrad Schindler¹

¹EcoVision Lab - Photogrammetry and Remote Sensing, ETH Zurich, Switzerland

²Institute for Computational Science, University of Zurich, Switzerland

Abstract—We propose a new STACKable Recurrent cell (STAR) for recurrent neural networks (RNNs), which has fewer parameters than widely used LSTM [16] and GRU [10] while being more robust against vanishing or exploding gradients. Stacking recurrent units into deep architectures suffers from two major limitations: (i) many recurrent cells (e.g., LSTMs) are costly in terms of parameters and computation resources; and (ii) deep RNNs are prone to vanishing or exploding gradients during training. We investigate the training of multi-layer RNNs and examine the magnitude of the gradients as they propagate through the network in the “vertical” direction. We show that, depending on the structure of the basic recurrent unit, the gradients are systematically attenuated or amplified. Based on our analysis we design a new type of gated cell that better preserves gradient magnitude. We validate our design on a large number of sequence modelling tasks and demonstrate that the proposed STAR cell allows to build and train deeper recurrent architectures, ultimately leading to improved performance while being computationally more efficient.

Index Terms—Recurrent neural network, Deep RNN, Multi-layer RNN.



1 INTRODUCTION

Recurrent Neural Networks (RNN) have established themselves as a powerful tool for modelling sequential data. They have led to significant progress for a variety of applications, notably language processing and speech recognition [13], [38], [42].

The basic building block of an RNN is a computational *unit* (or *cell*) that combines two inputs: the data of the current time step in the sequence and the unit’s own output from the previous time step. While RNNs can in principle handle sequences of arbitrary and varying length, they are (in their basic form) challenged by long-term dependencies, since learning those would require the propagation of gradients over many time steps. To alleviate this limitation, *gated* architectures have been proposed, most prominently Long Short-Term Memory (LSTM) cells [16] and Gated Recurrent Units (GRU) [10]. They use gating mechanisms to store and propagate information over longer time intervals, thus mitigating the vanishing gradient problem.

In general, abstract features are often represented better by deeper architectures [5]. In the same way that multiple hidden layers can be stacked in traditional feed-forward networks, multiple recurrent cells can also be stacked on top of each other, i.e., the output (or the hidden state) of the lower cell is connected to the input of the next-higher cell, allowing for different dynamics. For instance one might expect low-level cues to vary more with lighting, whereas high-level representations might exhibit object-specific variations over time. Several works [11], [34], [49] have shown the ability of deeper recurrent architectures to extract more complex features from the input and make better predictions. However, such architectures are usually composed of just two or three layers because training deeper

recurrent architectures still presents an open problem. More specifically, deep RNNs suffer from two main shortcomings: (i) they are difficult to train because of gradient instability, i.e., the gradient either explodes or vanishes during training; and (ii) the large number of parameters contained in each single cell makes deep architectures extremely resource-intensive. Both issues restrict the practical use of deep RNNs and particularly their usage for image-like input data, which generally requires multiple convolutional layers to extract discriminative, abstract representations. Our work aims to address these weaknesses by designing a recurrent cell that, on the one hand, requires fewer parameters and, on the other hand, allows for stable gradient back-propagation during training; thus allowing for deeper architectures.

Contributions We present a detailed, theoretical analysis of how the gradient magnitude changes as it propagates through a cell in a deep RNN lattice. Our analysis offers a different perspective compared to existing literature about RNN gradients, as it focuses on the gradient flow across layers in depth direction, rather than the recurrent flow across time. We show that the two dimensions behave differently, i.e., the ability to preserve gradients in time direction does not necessarily mean that they are preserved across layers, too. We believe that the analysis in this paper contributes a further, complementary step towards a full understanding of gradient propagation in deep RNNs.

We leverage our analysis to design a new, lightweight gated cell, termed the STACKable Recurrent (STAR) unit. The STAR cell better preserves the gradient magnitude in the deep RNN lattice, while at the same time using fewer parameters than existing gated cells like LSTM [16] and GRU [10].

We compare deep recurrent architectures built from different cells in an extensive set of experiments with several

popular datasets. The results confirm our analysis: training very deep recurrent nets fails with most conventional units, whereas the proposed STAR unit allows for significantly deeper architectures. Moreover, our experiments show that the proposed cell outperforms alternative designs on several different tasks and datasets.

2 RELATED WORK

Vanishing or exploding gradients during training are a long-standing problem of recurrent (and other) neural networks [6], [15]. Perhaps the most effective measure to address them so far has been to introduce gating mechanisms in the RNN structure, as first proposed by [16] in the form of the LSTM (long short-term memory), and later by other architectures such as gated recurrent units [10].

Importantly, RNN training needs proper initialisation. In [14], [22] it has been shown that initialising the weight matrices with identity and orthogonal matrices can be useful to stabilise the training. This idea is further develop in [3], [45], where authors impose the orthogonality throughout the entire training to keep the amplification factor of the weight matrices close to unity, leading to a more stable gradient flow. Unfortunately, it has been shown [43] that such hard orthogonality constraints hurt the representation power of the model and in some cases even *destabilise* the optimisation.

Another line of work has studied ways to mitigate the vanishing gradient problem by introducing additional (skip) connections across time and/or layers. Authors in [7] have shown that skipping state updates in RNNs shrinks the effective computation graph and thereby helps to learn longer-range dependencies. Other works such as [19], [30] introduce a residual connection between LSTM layers; however, the performance improvements are limited. In [11] the authors propose a gated feedback RNN that extends the stacked RNN architecture with extra connections. An obvious disadvantage of such an architecture are the extra computations and memory costs of the additional connections. Moreover, the authors only report results for rather shallow networks up to 3 layers.

Many of the aforementioned works propose new RNN architectures by leveraging a gradient propagation analysis. However all of these studies, as well as other studies which specifically aim at modelling accurately gradient propagation in RNNs [3], [9], [28], overlook the propagation of the gradient along the "vertical" depth dimension. In this work we will employ similar gradient analysis techniques, but focus on the depth dimension of the network.

Despite the described efforts, it remains challenging to train deep RNNs. In [49] authors propose to combine LSTMs and highway networks [36] to form Recurrent Highway Networks (RHN) and train deeper architectures. RHN are popular and perform well on language modelling tasks, but they are still prone to exploding gradients, as illustrated in our experiments. Another solution to alleviate gradient instability in deep RNNs was recently proposed in [25]. The work suggests the use of a restricted RNN called IndRNN where all interactions are removed between neurons in the hidden state of a layer. This idea combined with the usage of batch normalization appears to greatly stabilize the gradient

propagation through layers at the cost of a much lower representation power per layer. Such feature hinders IndRNN ability to achieve high performance for complex problems such as satellite image sequence classification or other computer vision tasks. In these tasks it is very important to merge information from neighboring pixels to increase the receptive field of the network so that the model has the ability to represent long-range spatial dependencies. Since IndRNN has no interaction between neurons it is difficult to achieve good spatio-temporal modeling effectively.

To process image sequence data, computer vision systems often rely on Convolutional LSTMs [46]. But while very deep CNNs are very effective and now standard [21], [35], stacks made of more than a few convLSTMs do not train well. Moreover, the computational cost increase rather quickly due to the large numbers of parameters in each LSTM cell. In practice, shallow versions are preferred, for instance [26] use a single layer for action recognition, and [48] use two layers to recognise hand gestures (combined with a deeper feature extractor without recursion).

3 BACKGROUND AND PROBLEM STATEMENT

In this section we revisit the mathematics of RNNs with particular emphasis on the gradient propagation. We will then leverage this analysis to design a more stable recurrent cell, which is described in Sec. 4. A RNN cell is a non-linear transformation that maps the input signal x_t at time t and the hidden state of the previous time step $t-1$ to the current hidden state h_t :

$$h_t = f(x_t, h_{t-1}, W), \quad (1)$$

with W the trainable parameters of the cell. The input sequences have an overall length of T , which can vary. It depends on the task whether the final state h_T , the complete sequence of states $\{h_t\}$, or a single sequence label (typically defined as the average $\frac{1}{T} \sum_t h_t$) are the desired target prediction for which loss \mathcal{L} is computed. Learning amounts to fitting W to minimise the loss, usually with stochastic gradient descent.

When stacking multiple RNN cells on top of each other, the hidden state of the lower level $l-1$ is passed on as input to the next-higher level l (Fig. 1). In mathematical terms this corresponds to the recurrence relation

$$h_t^l = f(h_t^{l-1}, h_{t-1}^l, w). \quad (2)$$

Temporal unfolding leads to a two-dimensional lattice with depth L and length T (Fig. 1), the forward pass runs from left to right and from bottom to top. Gradients flow in opposite direction: at each cell the gradient w.r.t. the loss arrives at the output gate and is used to compute the gradient w.r.t. (i) the weights, (ii) the input, and (iii) the previous hidden state. The latter two gradients are then propagated through the respective gates to the preceding cells in time and depth. In the following, we investigate how the magnitude of these gradients changes across the lattice. The analysis, backed up by numerical simulations, shows that common RNN cells are biased towards attenuating or amplifying the gradients and thus prone to destabilising the training of deep recurrent networks.

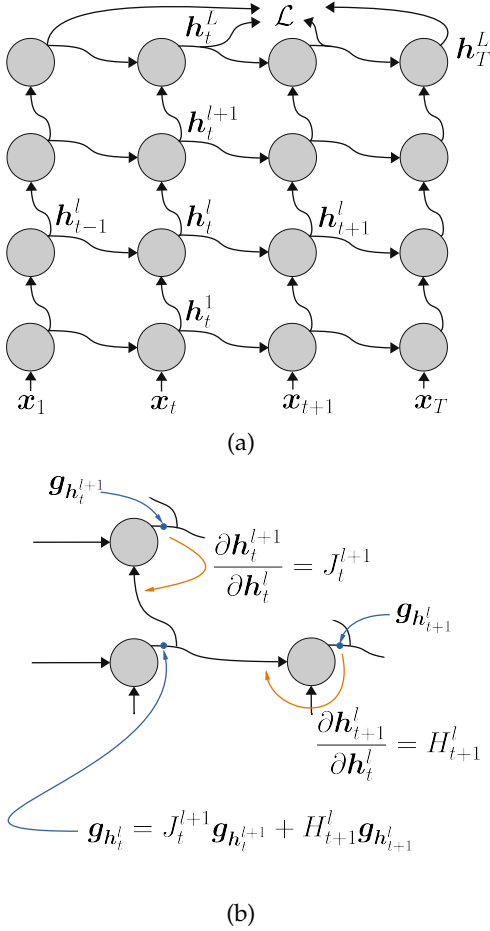


Fig. 1: (a) General structure of an unfolded deep RNN (b) Detail of the gradient backpropagation in the two dimensional lattice.

3.1 Gradient Magnitudes

The gradient w.r.t. the trainable weights at a single cell in the lattice is

$$\mathbf{g}_w = \frac{\partial \mathbf{h}_t^l}{\partial w} \mathbf{g}_{h_t^l}, \quad (3)$$

where $\frac{\partial \mathbf{h}_t^l}{\partial w}$ denotes the Jacobian matrix and $\mathbf{g}_{h_t^l}$ is a column vector containing the partial derivatives of the loss w.r.t. the cell's output (hidden) state. From the equation, it becomes apparent that the Jacobian acts as a "gain matrix" on the gradients, and should on average preserve their magnitude to prevent them from vanishing or exploding. We obtain the recurrence for propagation by expanding the gradient $\mathbf{g}_{h_t^l}$

$$\mathbf{g}_{h_t^l} = \frac{\partial \mathbf{h}_t^{l+1}}{\partial \mathbf{h}_t^l} \mathbf{g}_{h_{t+1}^l} + \frac{\partial \mathbf{h}_t^{l+1}}{\partial \mathbf{h}_t^l} \mathbf{g}_{h_{t+1}^l} = J_t^{l+1} \mathbf{g}_{h_{t+1}^l} + H_{t+1}^l \mathbf{g}_{h_{t+1}^l}, \quad (4)$$

with J_t^l the Jacobian w.r.t. the input and H_t^l the Jacobian w.r.t. the hidden state. Ideally we would like the gradient magnitude $\|\mathbf{g}_{h_t^l}\|_2$ to remain stable for arbitrary l and t . Characterising that magnitude completely is difficult because correlations may exist between $\mathbf{g}_{h_t^{l+1}}$ and $\mathbf{g}_{h_{t+1}^l}$ for instance, due to weight sharing. Nonetheless, it is evident that the two Jacobians J_t^{l+1} and H_{t+1}^l play a fundamental role: if their singular values are small, they will attenuate the gradients and cause them to vanish sooner or later. If their

singular values are large, they will amplify the gradients and make them explode.¹

In the following, we analyse the behaviour of the two matrices for two widely used RNN cells. We first consider the most simple RNN cell, hereinafter called Vanilla RNN (vRNN). Its recurrence equation reads

$$\mathbf{h}_t^l = \tanh(\mathbf{W}_x \mathbf{h}_t^{l-1} + \mathbf{W}_h \mathbf{h}_{t-1}^l + \mathbf{b}) \quad (5)$$

from which we get the two Jacobians

$$J_t^l = D_{\tanh(\mathbf{W}_x \mathbf{h}_t^{l-1} + \mathbf{W}_h \mathbf{h}_{t-1}^l + \mathbf{b})} \mathbf{W}_x \quad (6)$$

$$H_t^l = D_{\tanh(\mathbf{W}_x \mathbf{h}_t^{l-1} + \mathbf{W}_h \mathbf{h}_{t-1}^l + \mathbf{b})} \mathbf{W}_h \quad (7)$$

where D_x denotes a diagonal matrix with the elements of vector x as diagonal entries. Ideally, we would like to know the expected values of the two matrices' singular values. Unfortunately, there is no easy way to derive a closed-form analytical expressions for them, but we can compute them for a fixed, representative point. The most natural and illustrative choice is to set $\mathbf{h}_t^{l-1} = \mathbf{h}_{t-1}^l = 0$ because (i) in practice, RNNs' initial hidden states are set to $\mathbf{h}_0^l = 0$ (like in our experiments), and (ii) it is a stable and attracting fixed point so if the hidden state is perturbed around this point, it tends to return to its initial point. Note that this does not mean that the hidden state of the network is always equal to zero. The goal of the assumption is only simply to fix the value of the hidden state, in order to analyse the gradient propagation. We further choose weight matrices \mathbf{W}_h and \mathbf{W}_x with average singular value equal to one and $\mathbf{b} = 0$ (different popular initialisation strategies, such as orthogonal and identity matrices, are aligned with this assumption). Moreover, according to [18], [24], in the limit of a very wide network the parameters tend to stay close to their initial values, as a result the assumptions made are still legitimate during training (see the Appendix for empirical evidence). Since the derivative $\tanh'(0) = 1$, the average singular values of all matrices in Eq. (7) are equal to 1 in this configuration.

We expect to obtain a gradient $\mathbf{g}_{h_t^l}$ with a larger magnitude by combining the contributions of $\mathbf{g}_{h_{t+1}^l}$ and $\mathbf{g}_{h_{t+1}^{l+1}}$. To obtain a more precise estimate of the resulting gradient we should take into account the correlation between the two terms. However, if we examine two extreme cases (i) there is no or very small correlation between two gradient contributions, and (ii) they are highly (positively) correlated. The scaling factors of vRNN for the gradient are 1.414 and 2. respectively. Therefore, regardless of the correlation between the two terms, the gradient of vRNN is systematically growing while it propagates back in time and through layers. A deep network made of vRNN cells with orthogonal or identity initialisation can thus be expected to suffer, especially in the initial training phase, from exploding gradients as we move towards shallower layers and further back in time. To validate this assumption, we set up a toy example of a deep vRNN and compute the average gradient magnitude w.r.t. the network parameters for each cell in the unfolded network. For the numerical simulation we initialise all the

1. A subtle point is that sometimes large gradients are the precursor of vanishing gradients, if the associated large parameter updates cause the non-linearities to saturate.

$$J_t^l = D_{\tanh(c_t^l)} D_{(o_t^l)'} W_{xo} + D_{\tanh(c_t^l)'} D_{o_t^l} (D_{c_{t-1}^l} D_{(f_t^l)'} W_{xf} + D_{z_t^l} D_{(i_t^l)'} W_{xi} + D_{i_t^l} D_{(z_t^l)'} W_{xz}) \quad (8)$$

$$H_t^l = D_{\tanh(c_t^l)} D_{(o_t^l)'} W_{ho} + D_{\tanh(c_t^l)'} D_{o_t^l} (D_{c_{t-1}^l} D_{(f_t^l)'} W_{hf} + D_{z_t^l} D_{(i_t^l)'} W_{hi} + D_{i_t^l} D_{(z_t^l)'} W_{hz}) \quad (9)$$

hidden states and biases to 0, and chose random orthogonal matrices for the weights. Input sequences are generated with the random process $x_t = \alpha x_{t-1} + (1 - \alpha)z$, where $z \sim \mathcal{N}(0, 1)$ and the correlation factor $\alpha = 0.5$ (the choice of the correlation factor does not seem to qualitatively affect the results). Figure 2 depicts average gradient magnitudes over 10K runs with different weight initialisations and input sequences. As expected, the magnitude grows rapidly towards the earlier and shallower part of the network.

We perform a similar analysis for the classical LSTM cell [16]. The recurrent equations of the LSTM cell are the following:

$$i_t^l = \sigma(W_{xi} h_t^{l-1} + W_{hi} h_{t-1}^l + b_i) \quad (10)$$

$$f_t^l = \sigma(W_{xf} h_t^{l-1} + W_{hf} h_{t-1}^l + b_f) \quad (11)$$

$$o_t^l = \sigma(W_{xo} h_t^{l-1} + W_{ho} h_{t-1}^l + b_o) \quad (12)$$

$$z_t^l = \tanh(W_{xz} h_t^{l-1} + W_{hz} h_{t-1}^l + b_z) \quad (13)$$

$$c_t^l = f_t^l \circ c_{t-1}^l + i_t^l \circ z_t^l \quad (14)$$

$$h_t^l = o_t^l \circ \tanh(c_t^l), \quad (15)$$

where i, f, o are the input, forget, and output gate activations, respectively, c is the cell state. The expressions of the Jacobians are reported in Eqs. (8,9) where D_x again denotes a diagonal matrix with the elements of vector x as diagonal entries. The equations are slightly more complicated, but are still amenable to the same type of analysis. We again choose the same exemplary conditions as for the vRNN above, i.e., hidden states and biases equal to zero and orthogonal weight matrices. By substituting the numerical values in the aforementioned equations, we can see that the sigmoid function causes the expected singular value of the two Jacobians to drop to 0.25. Contrary to the vRNN cell, we expect that even the two Jacobians combined will produce an attenuation factor well below 1 (considering the same two extreme cases, i.e., uncorrelated and highly correlated, the value is 0.354 and 0.5, respectively) such that the gradient magnitude will decline and eventually vanish. We point out that LSTM cells have a second hidden state, the so-called ‘‘cell state’’. The cell state only propagates along the time dimension and not across layers, which makes the overall effect of the corresponding gradients more difficult to analyse. However, for the same reason one would, in a first approximation, expect that the cell state mainly influences the gradients in the time direction, but cannot help the flow through the layers. Again the numerical simulation results support our hypothesis as can be seen in Fig. 2. The LSTM gradients propagate relatively well backward through time, but vanish quickly towards shallower layers.

In summary, the gradient propagation behaves differently in time and depth directions. When considering the latter we need to take into consideration the gradient of the output w.r.t. the input state, too, and not exclusively consider the gradient w.r.t. the previous hidden state. Moreover, we need to take into account that the output of each cell is connected to *two* cells rather than one adjacent cell. Note that

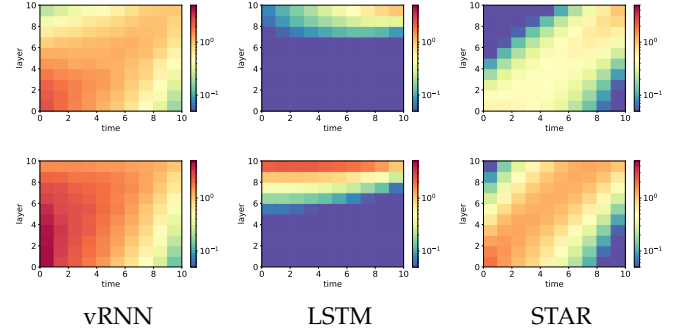


Fig. 2: Mean value of gradient magnitude with respect to the parameters for different RNN units. *top row*: loss $\mathcal{L}(h_T^L)$ only on final prediction. *bottom row*: loss $\mathcal{L}(h_1^L \dots h_T^L)$ over all time steps. As the gradients flow back through time and layers, for a network of vanilla RNN units they get amplified; for LSTM units they get attenuated; whereas the proposed STAR unit approximately preserves their magnitude. See the Appendix for the results with the real data.

this analysis is valid both when the loss is computed only using the final state T , and when all states are used (Fig. 2). In the latter case, we simply need to sum the contribution of all the separate losses. Usually, parameters are shared among different times t in RNNs, but not among different layers. If parameters are shared among different time steps, gradients accumulate row-wise (Fig. 2) increasing the gradient magnitude w.r.t. the parameters. This, however, is not true in the vertical direction as weights are not shared. As a consequence, it is particularly important to ensure that the gradient magnitude is preserved between adjacent layers.

4 THE STAR UNIT

Building upon the previous analysis, we introduce a novel RNN cell designed to avoid vanishing or exploding gradients while reducing the number of parameters. We start from the Jacobian matrix of the LSTM cell and investigate what design features are responsible for such low singular values. We see in Eq. (9) that every multiplication with \tanh non-linearity ($D_{\tanh(\cdot)}$), gating functions ($D_{\sigma(\cdot)}$), and with their derivatives can only ever decrease the singular values of W , since all those terms are always < 1 . The effect is particularly pronounced for the sigmoid and its derivative, $|\sigma'(\cdot)| \leq 0.25$ and $\mathbb{E}[\sigma(x)] = 0.5$ for zero-mean, symmetric distribution of x . In particular, the output gate o_t^l is a sigmoid and plays a major role in shrinking the overall gradients, as it multiplicatively affects all parts of both Jacobians. As a first measure, we thus propose to remove the output gate, which leads to h_t^l and c_t^l carrying the same information (the hidden state becomes an element-wise non-linear transformation of the cell state). To avoid this duplication and further simplify the design, we transfer

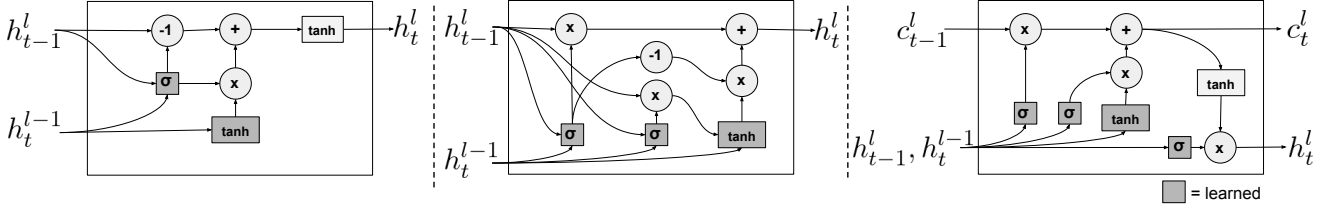


Fig. 3: RNN cell structures: STAR, GRU and LSTM, respectively.

the tanh non-linearity to the hidden state and remove the cell state altogether.

As a final modification, we also remove the input gate i_t^l from the architecture and couple it with the forget gate. We observed in detailed simulations that the input gate harms performance of deeper networks. This finding is consistent with the theory: for an LSTM cell with only the output gate removed, the Jacobians H_t^l, J_t^l will on average have singular values 1, respectively 0.5 (under the same conditions of Sec. 3). This suggests exploding gradients, which we indeed observe in numerical simulations. Moreover, signal propagation is less stable: state values can easily saturate if the two gates that control flow into the memory go out of sync. The gate structure of RHN [49] is similar to that configuration, and does empirically suffer from exploding, then vanishing, gradient (Fig. 4b).

More formally, our proposed STAR cell in the l -th layer takes the input \mathbf{h}_t^{l-1} (in the first layer, \mathbf{x}_t) at time t and non-linearly projects it to the space where the hidden vector \mathbf{h}^l lives, equation 16. Furthermore, the previous hidden state and the new input are combined into the gating variable \mathbf{k}_t^l (equation 17). \mathbf{k}_t^l is our analogue of the forget gate and controls how information from the previous hidden state and the new input are combined into a new hidden state. The complete dynamics of the STAR unit is given by the expressions

$$\mathbf{z}_t^l = \tanh(\mathbf{W}_z \mathbf{h}_t^{l-1} + \mathbf{b}_z) \quad (16)$$

$$\mathbf{k}_t^l = \sigma(\mathbf{W}_x \mathbf{h}_t^{l-1} + \mathbf{W}_h \mathbf{h}_{t-1}^l + \mathbf{b}_k) \quad (17)$$

$$\mathbf{h}_t^l = \tanh((1 - \mathbf{k}_t^l) \circ \mathbf{h}_{t-1}^l + \mathbf{k}_t^l \circ \mathbf{z}_t^l). \quad (18)$$

The Jacobian matrices for the STAR cell can be computed similarly to how it is done for the vRNN and LSTM (see the Appendix). In this case each of the two Jacobians has average singular values equal to 0.5. In the same two extreme cases previously considered, the scaling factor for the gradient becomes 0.707 and 1, respectively. Even if the gradient decays in the first case (worst case scenario, no correlation between two gradient contributions), it does so more slowly compared to LSTM. In the second case, the gradient can propagate without decaying or amplifying which is the ideal scenario. Empirically we have observed that, for an arbitrary STAR cell in the grid, these two terms are highly positively correlated, leading, ultimately, to a gradient scaling factor close to one. We repeat the same numerical simulations as above for the STAR cell, and find that it indeed maintains healthy gradient magnitudes throughout most of the deep RNN (Fig. 2). Finally, we point out that our proposed STAR architecture requires significantly less memory than most alternative designs. For the same input and hidden state

size, STAR has a 50%, respectively and 60% smaller memory footprint than GRU and LSTM. In the next section, we experimentally validate on real datasets that deep RNNs built from STAR units can be trained to a significantly greater depth while performing on par or better than state-of-the-art despite having fewer parameters.

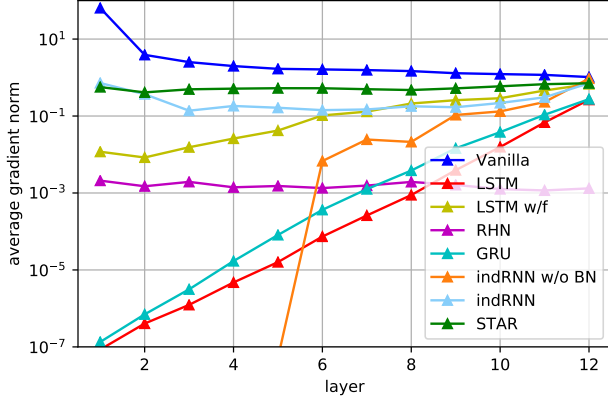
5 EXPERIMENTS

We evaluate the performance of several well-known RNN cells as well as that of the proposed STAR cell on different sequence modelling tasks with ten different datasets: sequential versions of MNIST [23], the adding [16], and copy memory [6] problems, music modeling [2, 29], character-level language modeling [27], which are a common testbeds for recurrent networks; three different remote sensing datasets, where time series of intensities observed in satellite images shall be classified into different agricultural crops [31, 32, 33]; and Jester [1] for hand gesture recognition. We use convolutional layers for gesture recognition and pixel-wise crop classification, whereas we employ conventional fully connected layers for the other tasks. The recurrent units we compare include the vRNN, the LSTM, the LSTM with only a forget gate [40], the GRU, the RHN [49], IndRNN [25], temporal convolution network (TCN) [4], transformer [41], and the proposed STAR. The experimental protocol is similar for all tasks: For each model variant, we train multiple versions with different depth (number of layers) and the best performing one is picked. Classification performance is measured by the rate of correct predictions (top-1 accuracy) for classification tasks, bits per character for character-level language modeling task and negative log-likelihood (NLL) for the rest of the tasks. Throughout the different experiments, we use orthogonal initialisation for weight matrices of RNNs. Training and network details for each experiment can be found in the Appendix.²

5.1 Pixel-by-pixel MNIST

We flatten all 28×28 grey-scale images of handwritten digits of the MNIST dataset [23] into 784×1 vectors, and the 784 values are sequentially presented to the RNN. The models' task is to predict the digit after having seen all pixels. The second task, pMNIST, is more challenging. Before flattening the images pixels are shuffled with a fixed random permutation, turning correlations between spatially close pixels into non-local long-range dependencies. As a consequence, the

²Code and trained models (in Tensorflow), as well as code for the simulations (in PyTorch), are available online: https://github.com/0zgur0/STAR_Network.



(a) gradient norm versus layer

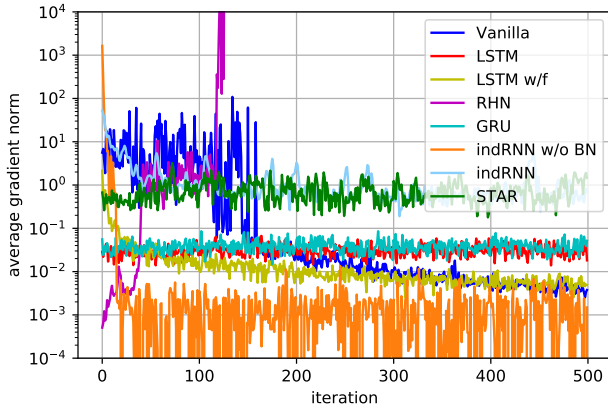
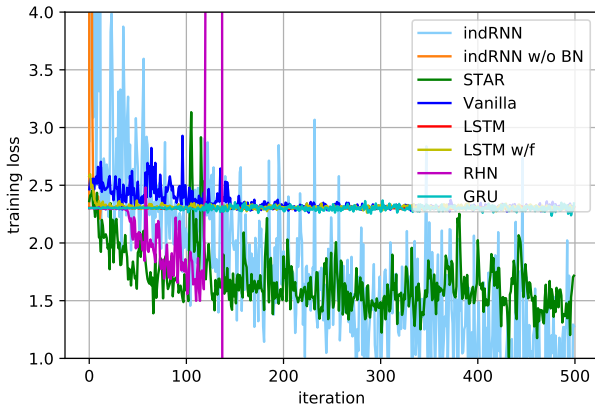
(b) gradient norm versus iteration, 1st epoch(c) training loss versus iteration, 1st epoch

Fig. 4: Gradient magnitudes of pix-by-pix MNIST. (a) Mean gradient norm per layer at the start of training. (b) Evolution of gradient norm during 1st training epoch. (c) Loss during 1st epoch.

model needs to remember dependencies between distant parts of the sequence to classify the digit correctly. Fig. 4a shows the average gradient norms per layer at the start of training for 12-layer networks built from different RNN cells. Propagation through the network increases the gradients for the vRNN and shrinks them for the LSTM. As the optimisation proceeds, we find that STAR and IndRNN remain stable, whereas all other units see a rapid decline of the gradients already within the first epoch, except for RHN, where the gradients explode, see Fig. 4b. Consequently, STAR and IndRNN are the only units for which a 12-layer model can be trained, as also confirmed by the evolution of the training loss, Fig. 4c.

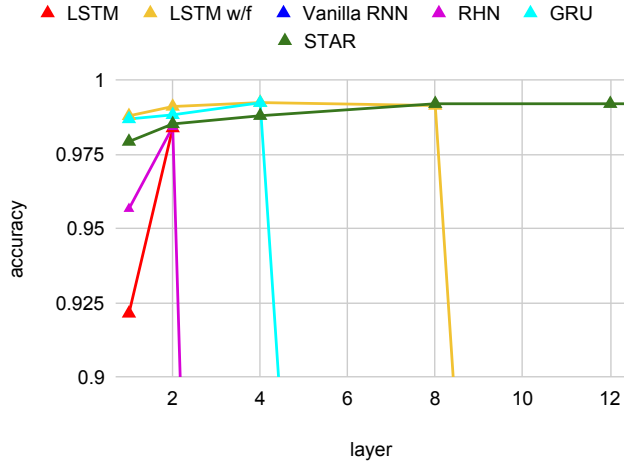
IndRNN's gradient propagation through layers is also stable even though not as good as STAR's. However, IndRNN strongly relies on Batch Normalization (BN) [17] for stable gradient propagation through layers while STAR does not require BN. If we remove the BNs between consecutive layers in IndRNN (denoted IndRNN w/o BN), its gradient propagation through layers and iterations becomes very unstable (see Fig. 4a and 4b). Indeed, IndRNN cannot be trained in those cases. It does not only fail for deeper, 12-layer setups applied to sequential MNIST, but also for shallower designs. Apart from increasing the computation overhead, general drawbacks of IndRNN's dependency on BNs are: (i) slow convergence during training and (ii) poor performance during inference if batch size is small (see the Appendix for further quantitative analysis).

Fig. 5 confirms that stacking into deeper architectures does benefit RNNs (except for vRNN); but it increases the risk of a catastrophic training failure. STAR is significantly more robust in that respect and can be trained up to >20 layers. On the comparatively easy and saturated MNIST data, the performance is comparable to a successfully trained LSTM (at depth 2-8 layers, LSTM training sometimes catastrophically fails; the displayed accuracies are averaged only over successful training runs).

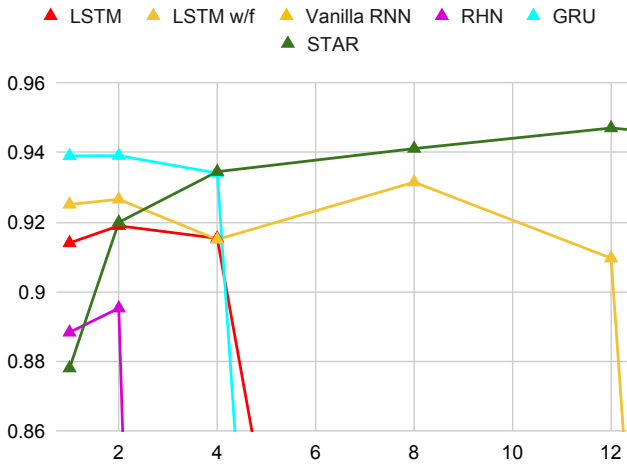
In 1 we show that our STAR cell mostly outperforms existing methods. As STAR is specifically designed to improve gradient propagation in the vertical direction, we conduct one additional experiment with a hybrid architecture: we use LSTM with a forget gate (which achieves good performance on the MNIST dataset in the one layer case) as first layer of the network and we stack seven layers of STAR cells on top. Such a design increases the capacity of the first layer without endangering gradient propagation. This further improves accuracy for both MNIST and pMNIST, leading to on par performance across both tasks with the best state-of-the-art methods BN-LSTM [12] and IndRNN [25]. Both methods employ Batch Normalization [17] inside the cells to improve the performance wrt to the simpler form of LSTM and IndRNN. We tested a version of the STAR cell which used BN and also in this case the modification lead to some performance improvements. This modification, however, is rather general and independent of the cell architecture, as it can be added to most of the other existing methods.

5.2 Adding Problem / Copy Memory

The adding problem [16] and the copy memory [6] are common benchmarks to evaluate whether a network is able



(a) MNIST



(b) pMNIST

Fig. 5: Accuracy results for pixel-by-pixel MNIST tasks.

to learn long-term memory. In the adding problem, two sequences of length T are taken as input: the first sequence consists of independent samples in range $(0, 1)$, while the second sequence is a binary vector with two entries set to 1 and the rest 0. The goal is to sum the two entries of the first sequence indicated by the two entries of 1 in the second sequence. In copy memory task, the input sequence is of length $T + 20$. The first 10 values in the sequence are chosen randomly among the digits $\{1, \dots, 8\}$, the sequence is then followed by T zeros, the last 11 entries are filled with the digit 9 (the first 9 is a delimiter). The goal is to generate an output of the same length that is zero everywhere except for the last 10 values after the delimiter, where the model is expected to repeat the 10 values encountered at the beginning of the input sequence. We perform experiments with two different sequence lengths, $T = 200$ and $T = 1000$, using different RNNs with the same number of parameters (70K). The results are shown in Fig. 6, 7. The vRNN is unable to perform long-term memorization, whereas LSTM has issue with longer sequences ($T = 1000$). In contrast, both STAR and GRU, can learn long-term memory even

Method	MNIST	pMNIST	units
vRNN (1 layer)	24.3%	44.0%	128
LSTM (2 layers)	98.4%	91.9%	128
GRU (2 layers)	98.8%	93.9%	128
RHN (2 layers)	98.4%	89.5%	128
iRNN [22]	97.0%	82.0%	100
uRNN [3]	95.1%	91.4%	512
FC uRNN [45]	96.9%	94.1%	512
Soft ortho [43]	94.1%	91.4%	128
AntisymRNN [8]	98.8%	93.1%	128
IndRNN [25]	99.0%	96.0%	128
BN-LSTM [12]	99.0%	95.4%	100
sTANH-RNN [47]	98.1%	94.0%	128
STAR (8 layers)	99.2%	94.1%	128
STAR (12 layers)	99.2%	94.7%	128
LSTM w/ f STAR (8 layers)	99.4%	95.4%	128

TABLE 1: Performance comparison for pixel-by-pixel MNIST tasks. Our best performing configuration **bold underlined**, top performers state-of-the-art **bold**.

when the sequences are very long. An advantage of STAR in this case is its faster convergence.

5.3 TUM & BreizhCrops Time Series Classification

We evaluate model performance on a more realistic sequence modelling problem, where the aim is to classify agricultural crop types using sequences of satellite images. In this case, time-series modelling captures phenological evidence, i.e. different crops have different growing patterns over the season. For the TUM dataset, the input is a time series of 26 multi-spectral Sentinel-2A satellite images with a ground resolution of 10m collected over a 102 km x 42 km area north of Munich, Germany between December 2015 and August 2016 [31]. We use patches of 3x3 pixels recorded in 6 spectral channels and flattened into 54x1 vectors as input. For the BreizhCrop dataset, the input is a time series of 45 multi-spectral Sentinel-2A satellite images with a ground resolution of 10m collected from 580k field parcels in the Region of Brittany, France of the season 2017. The input is 4 spectral channels (R, G, B, NIR) [33]. In the first task, only TUM dataset is used. The vectors are sequentially presented to the RNN model, which outputs a prediction at every time step (note that for this task the correct answer can sometimes be "cloud", "snow", "cloud shadow" or "water", which are easier to recognise than many crops). STAR outperforms all baselines, and it is again more suitable for stacking into deep architectures (Fig. 8). In the second task, both datasets are used. The goal is a single-step prediction i.e., the model predicts a crop type after the entire sequence is presented. STAR significantly outperforms all the baselines including TCN and the recently proposed method, IndRNN [25] (Tab. 2). Note that IndRNN also aims to build deep multi-layer RNNs. The performance gain is stronger in the BreizhCrop datasets. This is probably because the sequence is longer and the depth of the network helps to capture more complex dependencies in the data.

5.4 Music Modeling

JSB Chorales [2] is a polyphonic music dataset consisting of the entire corpus of 382 four-part harmonized chorales

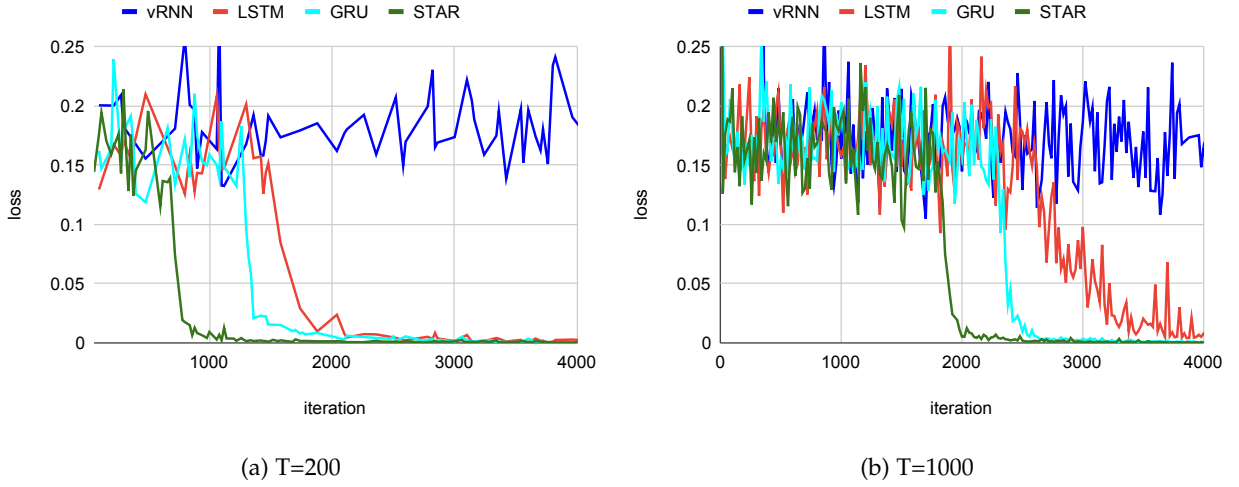


Fig. 6: Performance comparison for adding problem.

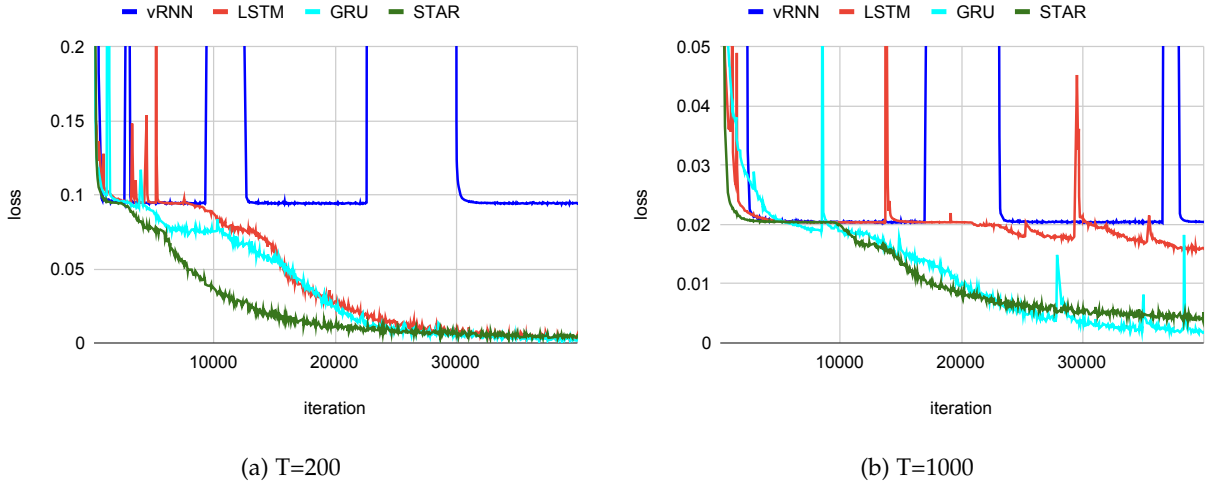


Fig. 7: Performance comparison for copy memory task.

Method	TUM		BreizhCrops	
	Acc	#params	Acc	#params
vRNN (2 layers)	84.4%	45k	38.1%	55k
LSTM (2 layers)	85.9%	170k	60.1%	210k
LSTM w/f (2 layers)	85.7%	90k	58.3%	105k
GRU (2/4 layers)	85.7%	130k	66.4%	350k
RHN (2 layers)	85.6%	100k	-	-
IndRNN (4 layers)	85.5%	90k	56.8%	105k
TCN (2 layers)	84.9%	300k	61.5%	360k
STAR (4 layers)	87.7%	130k	68.2%	170k
STAR (6 layers)	87.6%	210k	69.6%	270k

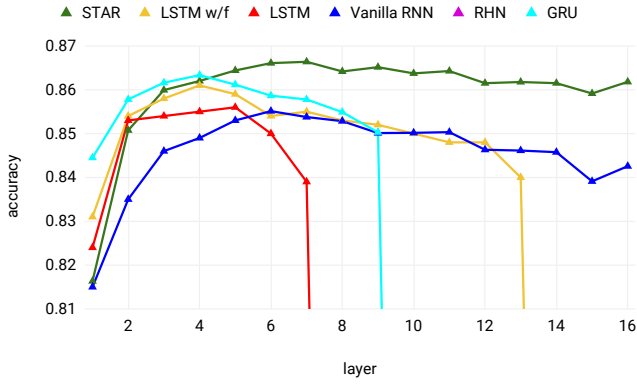
TABLE 2: Performance comparison for time series crop classification.

by J. S. Bach. Each input is a sequence of chord elements. Each element is an 88-bit binary code that corresponds to the 88 keys of a piano, with 1 indicating a key pressed at a given time. Piano-Midi [29] is a classical piano MIDI archive that consists of 130 pieces by various composers. These datasets have been used in several previous works to

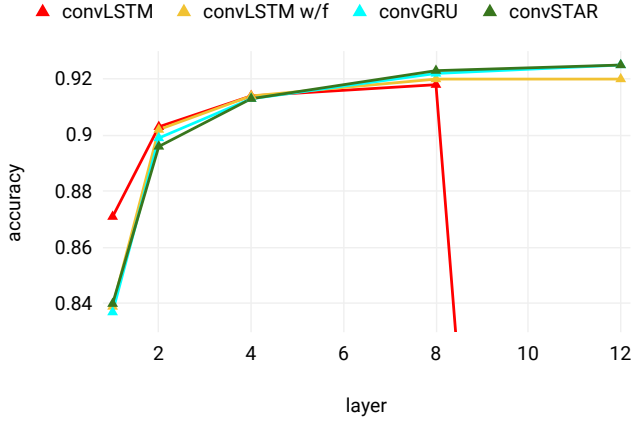
Method	JSB Chorales		Piano-Midi	
	NLL	#params	NLL	#params
vRNN [37]	8.72	40k	7.65	140k
LSTM [37]	8.51	650k	7.84	480k
GRU [37]	8.53	640k	7.62	690k
diagRNN [37]	8.14	420k	7.48	360k
TCN (2 layers) [4]	8.10	300k	-	-
STAR (2 layers)	8.13	360k	7.40	480k
STAR (4 layers)	8.09	830k	-	-

TABLE 3: Performance comparison for music task. The performance is measured in terms of negative log-likelihood (NLL).

investigate the ability of RNNs to represent music [10], [37]. The performance on both tasks is measured in terms of per-frame negative log-likelihood (NLL) on a test set. We follow the exact same experimental setup described in [37]. STAR works better than all tested RNN baselines, and performs on par with TCN (see Tab. 3).



(a) TUM, per time-step labels



(b) Jester, single label / sequence

Fig. 8: Time series classification. (a) Crop classes. (b) Hand gestures (convolutional RNNs).

5.5 Character-level Language Modeling

For this task we used the PennTreebank (PTB) [27]. When used as a character-level language corpus, PTB contains 5,059K characters for training, 396K for validation, and 446K for testing, with an alphabet size of 50. The goal is to predict the next character given the preceding context. We follow the exact same experimental setup as [4]. The performance is measured in terms of bits per character (BPC, i.e. average cross entropy over the alphabet) on the test set. On this task STAR outperforms all baselines, including Transformer and TCN (see Tab. 4).

5.6 Hand-gesture recognition from video

We also evaluate STAR on sequences of images, using *convolutional* layers. We analyse performance of STAR versus state-of-the-art on gesture recognition from video and pixel-wise crop classification. The 20BN-Jester dataset V1 [1] is a large collection of densely-labelled short video clips, where each clip contains a predefined hand gesture performed by a worker in front of a laptop camera or webcam. In total, the dataset includes 148'094 RGB video files of 27 types of gestures (see Fig. 11). The task is to classify which gesture is seen in a video. 32 consecutive frames of size 112×112 pixels are sequentially presented to the convolutional RNN.

Method	BPC	#params
vRNN [4]	1.48	3M
LSTM (2 layers) [4]	1.36	3M
GRU [4]	1.37	3M
IndRNN (6 layers)*	1.42	3M
TCN (3 layers) [4]	1.31	3M
Transformer (3 layers) [44]	1.45	-
STAR (6 layers)	1.30	3M

TABLE 4: Performance comparison for PennTreebank character-level language modeling. The performance is measured in terms of bits per character (BPC). *We run this experiment as designed in [4]’s experimental setup with a limited number of parameters to allow for a fair comparison. Note that [25] reports a better result, but uses many more model parameters.

Method	Accuracy	#params
convLSTM (8 layers)	91.8 %	2.2M
convLSTM w/f (8 layers)	92.0 %	1.1M
convGRU (12 layers)	92.5 %	2.5M
convSTAR (8 layers)	92.3 %	0.8M
convSTAR (12 layers)	92.5 %	1.2M
convLSTM convSTAR (8 layers)	92.7 %	0.9M

TABLE 5: Performance comparison for the gesture recognition task (Jester).

At the end, the model again predicts a gesture class via an averaging layer over all time steps. The outcome for convolutional RNNs is coherent with the previous results, see Fig. 8b, Tab. 5. Going deeper improves the performance of all four tested convRNNs. The improvement is strongest for convolutional STAR, and the best performance is reached with a deep model (12 layers). In summary, the results confirm both our intuitions that depth is particularly useful for convolutional RNNs, and that STAR is more suitable for deeper architectures, where it achieves higher performance with better memory efficiency. We note that in the shallow 1-2 layer setting the conventional LSTM performs a slightly better than the three others, likely due to its larger capacity. Lastly, we conduct the same additional experiment with the hybrid architecture as we do for MNIST tasks. We stack seven layers of STAR on top of one layer of LSTM. This further improves the results and achieves 92.7% accuracy (compared this to eight LSTM layers, which achieve only 91.8% accuracy with about twice as many parameters).

5.7 TUM image series pixel-wise classification

In another experiment with convolutional RNNs, we classify crops pixel-wise (and thus use convolutional layers) using a dataset [32] (TUM) containing Sentinel-2A optical satellite image sequences (RGB and NIR at 10 m ground sampling distance) accompanied by ground-truth land cover maps. Each satellite image sequence contains 30 images of size 48×48 px collected in 2016 within a $102 \text{ km} \times 42 \text{ km}$ region north of Munich, Germany (see Fig. 12). We compare pixel-wise classification accuracy for a network with a fixed depth of four layers and for four different basic recurrent cells LSTM, LSTM with only a forget gate, GRU, and the proposed STAR cell (Tab. 6). Moreover

Method	Acc	#params	#compute
bi-convGRU (1 layer) [32]	89.7%	6.2M	46bn
convLSTM (4 layers)	90.6%	292k	2.7bn
convLSTM w/f (4 layers)	89.6%	161k	1.5bn
convGRU (4 layers)	90.1%	227k	2.1bn
convSTAR (4 layers)	91.9%	124k	1.1bn

TABLE 6: Performance comparison for TUM pixel-wise image classification task.

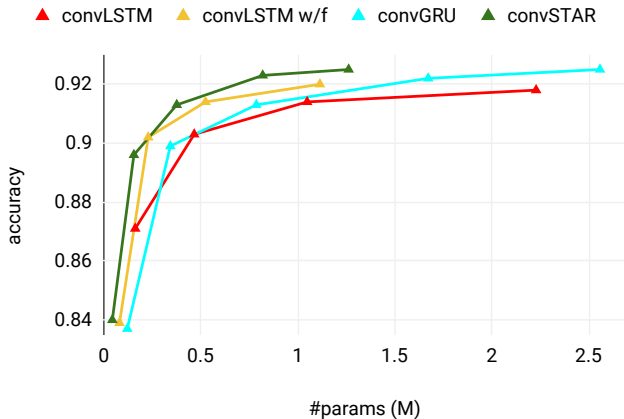


Fig. 9: Accuracy versus number of model parameters for the gesture recognition task (Jester).

we include the performance obtained in [32] using a bi-directional convolutional GRU with a single layer. Our STAR cell outperforms all other methods (Tab. 6) while requiring less memory and being computationally less costly.

5.8 Computational Resources and Training Time

Last, we compare to widely used recurrent units LSTM and GRU in terms of parameter efficiency and training time for the convolutional version used in gesture recognition. We plot performance versus number of parameters (Fig. 9) STAR outperforms LSTM and performs on par with GRU, but requires only half the number of parameters. We plot accuracy on the validation dataset versus training time for different recurrent units for the gesture recognition task in Fig. 10. STAR does not only require significantly less parameters but can also be trained much faster: the validation accuracy on the dataset after 8 hours is comparable to the best validation achieved by the LSTM and the GRU after 20 hours of training.

6 CONCLUSION

We have proposed STAR, a novel stackable recurrent cell type that it is specifically designed to be employed in deep recurrent architectures. A theoretical analysis and associated numerical simulations indicated that widely used standard RNN cells like LSTM and GRU do not preserve gradient

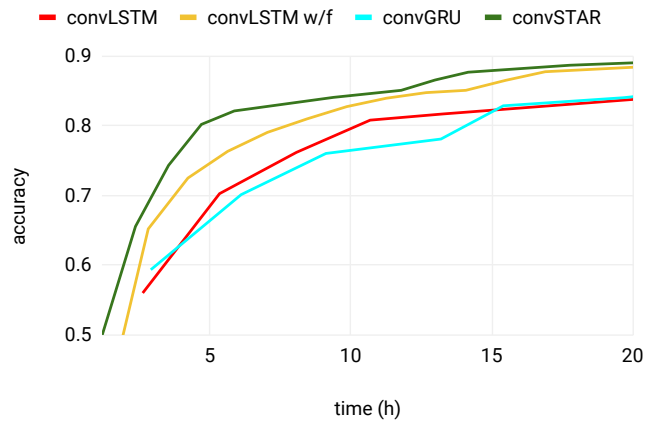


Fig. 10: Test accuracy versus training time for the gesture recognition task (Jester), 4 layers networks.

magnitudes in the “vertical” direction during backpropagation. As the depth of the network grows, the risk of either exploding or vanishing gradients increases. We leveraged this analysis to design a novel cell that better preserves the gradient magnitude between two adjacent layers, is better suited for deep architectures, and requires fewer parameters than other widely used recurrent units. An extensive experimental evaluation on several publicly available datasets confirms that STAR units can be stacked into deeper architectures and in many cases performs better than state-of-the-art architectures.

We see two main directions for future work. On the one hand, it would be worthwhile to develop a more formal and thorough mathematical analysis of the gradient flow, and perhaps even derive rigorous bounds for specific cell types, that could, in turn, inform the network design. On the other hand, it appears promising to investigate whether the analysis of the gradient flows could serve as a basis for better initialisation schemes to compensate the systematic influences of the cells structure, e.g., gating functions, in the training of deep RNNs.

ACKNOWLEDGMENTS

We thank the Swiss Federal Office for Agriculture (FOAG) for partially funding this Research project through the Deep-Field Project.

REFERENCES

- [1] The 20bn-jester dataset v1. <https://20bn.com/datasets/jester>.
- [2] Moray Allan and Christopher Williams. Harmonising chorales by probabilistic inference. In *Advances in neural information processing systems*, 2005.
- [3] Martin Arjovsky, Amar Shah, and Yoshua Bengio. Unitary evolution recurrent neural networks. In *ICML*, 2016.
- [4] Shaojie Bai, J Zico Kolter, and Vladlen Koltun. An empirical evaluation of generic convolutional and recurrent networks for sequence modeling. *arXiv preprint arXiv:1803.01271*, 2018.
- [5] Yoshua Bengio. Learning deep architectures for AI. *Foundations and trends® in Machine Learning*, 2(1):1–127, 2009.
- [6] Yoshua Bengio, Patrice Simard, Paolo Frasconi, et al. Learning long-term dependencies with gradient descent is difficult. *IEEE TNN*, 5(2):157–166, 1994.

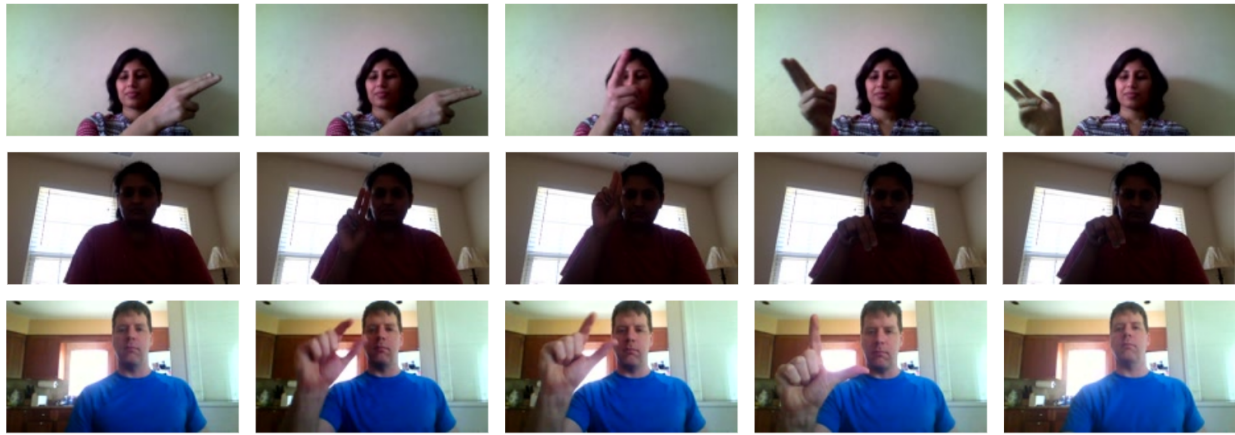


Fig. 11: Example frames of the Jester dataset. Columns show 1st, 8th, 16th, 24th, 32nd frames, respectively. First row: Sliding two fingers right. Second row: Sliding two fingers down. Third row: Zooming in with two fingers.



Fig. 12: Example satellite images of the TUM dataset. Each row shows randomly sampled images (in order, only R, G and B channels) from a satellite image time-series. The last column shows the ground-truth where different colors correspond to different crop types.

- [7] Víctor Campos, Brendan Jou, Xavier Giró-i Nieto, Jordi Torres, and Shih-Fu Chang. Skip RNN: Learning to skip state updates in recurrent neural networks. In *ICLR*, 2018.
- [8] Bo Chang, Minmin Chen, Eldad Haber, and Ed H Chi. AntisymmetricRNN: A dynamical system view on recurrent neural networks. In *ICLR*, 2019.
- [9] Minmin Chen, Jeffrey Pennington, and Samuel S Schoenholz. Dynamical isometry and a mean field theory of rnns: Gating enables signal propagation in recurrent neural networks. In *ICML*, 2018.
- [10] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. In *NIPS Workshop*, 2014.
- [11] Junyoung Chung, Caglar Gulcehre, Kyunghyun Cho, and Yoshua Bengio. Gated feedback recurrent neural networks. In *ICML*, 2015.
- [12] T. Cooijmans, N. Ballas, C. Laurent, C. Gulcehre, and A. Courville. Recurrent batch normalization. In *ICLR*, 2017.
- [13] Alex Graves, Abdel-rahman Mohamed, and Geoffrey Hinton. Speech recognition with deep recurrent neural networks. In *ICASSP*, 2013.
- [14] Mikael Henaff, Arthur Szlam, and Yann LeCun. Recurrent orthogonal networks and long-memory tasks. In *ICML*, 2016.
- [15] Sepp Hochreiter. Untersuchungen zu Dynamischen Neuronalen Netzen. *Diploma Thesis, Technische Universität München*, 91(1), 1991.
- [16] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997.
- [17] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015.
- [18] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in neural information processing systems*, 2018.
- [19] Jaeyoung Kim, Mostafa El-Khamy, and Jungwon Lee. Residual LSTM: Design of a deep recurrent architecture for distant speech recognition. In *Interspeech*, 2017.
- [20] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2014.
- [21] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. ImageNet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, 2012.
- [22] Quoc V Le, Navdeep Jaitly, and Geoffrey E Hinton. A simple way to initialize recurrent networks of rectified linear units. *arXiv*

- preprint *arXiv:1504.00941*, 2015.
- [23] Yann LeCun, Léon Bottou, Yoshua Bengio, Patrick Haffner, et al. Gradient-based learning applied to document recognition. *Proc. IEEE*, 86(11):2278–2324, 1998.
- [24] Jaehoon Lee, Lechao Xiao, Samuel Schoenholz, Yasaman Bahri, Roman Novak, Jascha Sohl-Dickstein, and Jeffrey Pennington. Wide neural networks of any depth evolve as linear models under gradient descent. In *Advances in neural information processing systems*, 2019.
- [25] Shuai Li, Wanqing Li, Chris Cook, Ce Zhu, and Yanbo Gao. Independently recurrent neural network (indrnn): Building a longer and deeper rnn. In *CVPR*, 2018.
- [26] Zhenyang Li, Kirill Gavriluk, Efstratios Gavves, Mihir Jain, and Cees GM Snoek. VideoLSTM convolves, attends and flows for action recognition. *CVIU*, 166:41–50, 2018.
- [27] Mitchell Marcus, Beatrice Santorini, and Mary Ann Marcinkiewicz. Building a large annotated corpus of english: The penn treebank. 1993.
- [28] Zakaria Mhammedi, Andrew Helicar, Ashfaqur Rahman, and James Bailey. Efficient orthogonal parametrisation of recurrent neural networks using householder reflections. In *ICML*, 2017.
- [29] Graham E Poliner and Daniel PW Ellis. A discriminative model for polyphonic piano transcription. *EURASIP Journal on Advances in Signal Processing*, 2007:1–9, 2006.
- [30] Sabeek Pradhan and Shayne Longpre. Exploring the depths of recurrent neural networks with stochastic residual learning, 2016.
- [31] Marc Rußwurm and Marco Körner. Temporal vegetation modelling using long short-term memory networks for crop identification from medium-resolution multi-spectral satellite images. In *CVPR Workshops*, 2017.
- [32] Marc Rußwurm and Marco Körner. Multi-temporal land cover classification with sequential recurrent encoders. *ISPRS International Journal of Geo-Information*, 7(4):129, 2018.
- [33] Marc Rußwurm, Sébastien Lefèvre, and Marco Körner. Breizhcrops: A satellite time series dataset for crop type identification. In *ICML Workshop*, 2019.
- [34] Amir Shahrudiy, Jun Liu, Tian-Tsong Ng, and Gang Wang. Ntu rgb+ d: A large scale dataset for 3d human activity analysis. In *CVPR*, 2016.
- [35] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015.
- [36] Rupesh Kumar Srivastava, Klaus Greff, and Jürgen Schmidhuber. Highway networks. In *ICML Workshop*, 2015.
- [37] Y Cem Subakan and Paris Smaragdis. Diagonal rnns in symbolic music modeling. In *2017 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*. IEEE, 2017.
- [38] Ilya Sutskever, Oriol Vinyals, and Quoc V Le. Sequence to sequence learning with neural networks. In *Advances in neural information processing systems*, 2014.
- [39] Corentin Tallec and Yann Ollivier. Can recurrent neural networks warp time? In *ICLR*, 2018.
- [40] Jos Van Der Westhuizen and Joan Lasenby. The unreasonable effectiveness of the forget gate. *arXiv preprint arXiv:1804.04849*, 2018.
- [41] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, 2017.
- [42] Oriol Vinyals and Quoc Le. A neural conversational model. *arXiv preprint arXiv:1506.05869*, 2015.
- [43] Eugene Vorontsov, Chiheb Trabelsi, Samuel Kadoury, and Chris Pal. On orthogonality and learning recurrent networks with long term dependencies. In *ICML*, 2017.
- [44] Zhiwei Wang, Yao Ma, Zitao Liu, and Jiliang Tang. R-transformer: Recurrent neural network enhanced transformer. *arXiv preprint arXiv:1907.05572*, 2019.
- [45] Scott Wisdom, Thomas Powers, John Hershey, Jonathan Le Roux, and Les Atlas. Full-capacity unitary recurrent neural networks. In *Advances in neural information processing systems*, 2016.
- [46] Shi Xingjian, Zhourong Chen, Hao Wang, Dit-Yan Yeung, Wai-Kin Wong, and Wang-chun Woo. Convolutional LSTM network: A machine learning approach for precipitation nowcasting. In *Advances in neural information processing systems*, 2015.
- [47] et al. Zhang, Saizheng. Architectural complexity measures of recurrent neural networks. In *Advances in neural information processing systems*, 2016.
- [48] Liang Zhang, Guangming Zhu, Lin Mei, Peiyi Shen, Syed Afaq Ali Shah, and Mohammed Bennamoun. Attention in convolutional LSTM for gesture recognition. In *Advances in neural information processing systems*, 2018.

- [49] Julian Georg Zilly, Rupesh Kumar Srivastava, Jan Koutník, and Jürgen Schmidhuber. Recurrent highway networks. In *ICML*, 2017.



Mehmet Ozgur Turkoglu received his BSc degrees in both electrical engineering and physics from Bogazici University in 2016. He studied a master's in electrical engineering with a specialization in computer vision at the University of Twente. He is a PhD candidate in the EcoVision group at ETH Zürich since 2018. His research interests include computer vision, deep learning and their applications to remote sensing data. He is particularly interested in deep sequence modeling of time-series data.



Stefano D'Aronco received his BS and MS degrees in electronic engineering from the Università degli studi di Udine, in 2010 and 2013 respectively. He then joined the Signal Processing Laboratory (LTS4) in 2014 as a PhD student under the supervision of Prof. Pascal Frossard. He received his PhD in Electrical Engineering from École Polytechnique Fédérale de Lausanne in 2018. He is Postdoctoral researcher in the EcoVision group at ETH Zürich since 2018. His research interests include several machine learning topics, such as Bayesian inference method and deep learning, with particular emphasis on applications related to remote sensing an environmental monitoring.



Jan Dirk Wegner is associate professor at University of Zurich and head of the EcoVision Lab at ETH Zurich. Jan was PostDoc (2012-2016) and senior scientist (2017-2020) in the Photogrammetry and Remote Sensing group at ETH Zurich after completing his PhD (with distinction) at Leibniz Universität Hannover in 2011. He was granted multiple awards, among others an ETH Postdoctoral fellowship and the science award of the German Geodetic Commission. Jan was selected for the WEF Young Scientist Class 2020 as one of the 25 best researchers world-wide under the age of 40 committed to integrating scientific knowledge into society for the public good. Jan is chair of the ISPRS II/WG 6 "Large-scale machine learning for geospatial data analysis" and organizer of the CVPR EarthVision workshops.



Konrad Schindler (M'05–SM'12) received the Diplomingenieur (M.Tech.) degree from Vienna University of Technology, Vienna, Austria, in 1999, and the Ph.D. degree from Graz University of Technology, Graz, Austria, in 2003. He was a Photogrammetric Engineer in the private industry and held researcher positions at Graz University of Technology, Monash University, Melbourne, VIC, Australia, and ETH Zürich, Zürich, Switzerland. He was an Assistant Professor of Image Understanding with TU Darmstadt, Darmstadt, Germany, in 2009. Since 2010, he has been a Tenured Professor of Photogrammetry and Remote Sensing with ETH Zürich. His research interests include computer vision, photogrammetry, and remote sensing.

APPENDIX

A.1 RNN Cells Dynamics

In the following, we provide more detailed insights about the updating rules of the tested cell types.

Vanilla RNN update rule:

$$h_t^l = \tanh(W_x h_t^{l-1} + W_h h_{t-1}^l + b) \quad (19)$$

LSTM update rule:

$$i_t^l = \sigma(W_{xi} h_t^{l-1} + W_{hi} h_{t-1}^l + b_i) \quad (20)$$

$$f_t^l = \sigma(W_{xf} h_t^{l-1} + W_{hf} h_{t-1}^l + b_f) \quad (21)$$

$$o_t^l = \sigma(W_{xo} h_t^{l-1} + W_{ho} h_{t-1}^l + b_o) \quad (22)$$

$$z_t^l = \tanh(W_{xz} h_t^{l-1} + W_{hz} h_{t-1}^l + b_z) \quad (23)$$

$$c_t^l = f_t^l \circ c_{t-1}^l + i_t^l \circ z_t^l \quad (24)$$

$$h_t^l = o_t^l \circ \tanh(c_t^l). \quad (25)$$

LSTM with only forget gate, update rule:

$$f_t^l = \sigma(W_{xf} h_t^{l-1} + W_{hf} h_{t-1}^l + b_f) \quad (26)$$

$$z_t^l = \tanh(W_{xz} h_t^{l-1} + W_{hz} h_{t-1}^l + b_z) \quad (27)$$

$$h_t^l = \tanh(f_t^l \circ h_{t-1}^l + (1 - f_t^l) \circ z_t^l) \quad (28)$$

GRU update rule:

$$z_t^l = \sigma(W_{xz} h_t^{l-1} + W_{hz} h_{t-1}^l + b_z) \quad (29)$$

$$r_t^l = \sigma(W_{xr} h_t^{l-1} + W_{hr} h_{t-1}^l + b_r) \quad (30)$$

$$h_t^l = (1 - z_t^l) \circ h_{t-1}^l + z_t^l \circ \tanh(W_{xh} h_t^{l-1} + W_{hh}(r_t^l \circ h_{t-1}^l) + b_h) \quad (31)$$

STAR Jacobians:

$$J_t^l = D_{\tanh(h_{t-1}^l + k_t^l \circ (z_t^l - h_{t-1}^l))}' \quad (32)$$

$$\cdot (D_{z_t^l - h_{t-1}^l} D_{(k_t^l)'} W_x + D_{k_t^l} D_{(z_t^l)'} W_z)$$

$$H_t^l = D_{\tanh(h_{t-1}^l + k_t^l \circ (z_t^l - h_{t-1}^l))}' \quad (33)$$

$$\cdot (I + D_{z_t^l - h_{t-1}^l} D_{(k_t^l)'} W_h - D_{k_t^l})$$

Convolutional STAR: We briefly describe the convolutional version of our proposed cell. The main difference is matrix multiplications now become convolutional operations. The dynamics of the convSTAR cell is given in the following equations.

$$K_t^l = \sigma(W_x * H_t^{l-1} + W_h * H_{t-1}^l + B_K) \quad (34)$$

$$Z_t^l = \tanh(W_z * H_t^{l-1} + B_z) \quad (35)$$

$$H_t^l = \tanh(H_{t-1}^l + K_t^l \circ (Z_t^l - H_{t-1}^l)) \quad (36)$$

A.2 Further Numerical Gradient Propagation Analysis

In this section, we extend the numerical simulations of the gradient propagation in the unfolded recurrent neural network to two further cell architectures, namely the GRU [10] and the LSTM with only forget gate for the synthetic dataset (Sec. A.2.1); and for the real dataset, MNIST (Sec. A.2.2).

A.2.1 Synthetic Dataset

The setup of the numerical simulations is the same as the one described in Section 3. As can be seen from Fig. 13 the GRU and the LSTM with only forget gate mitigate the attenuation of gradients to some degree. However, we observe that the corresponding standard deviations are much higher, i.e., the gradient norm greatly varies across different runs, see Fig. 14. We found that the gradients within a single run oscillate a lot more, for both LSTMw/f and GRU, and make training unstable which is undesirable. Moreover, the gradient magnitudes evolve very differently for different initial values, meaning that the training is less robust against fluctuations of the random initialisation.

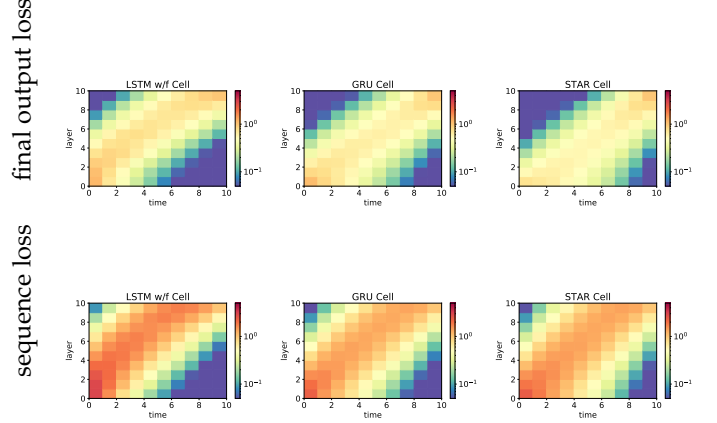


Fig. 13: Mean gradient magnitude w.r.t. the parameters for LSTM with only forget gate, GRU, and the proposed STAR cell. *top row*: loss $\mathcal{L}(h_T^L)$ only on final prediction. *bottom row*: loss $\mathcal{L}(h_1^L \dots h_T^L)$ over all time steps.

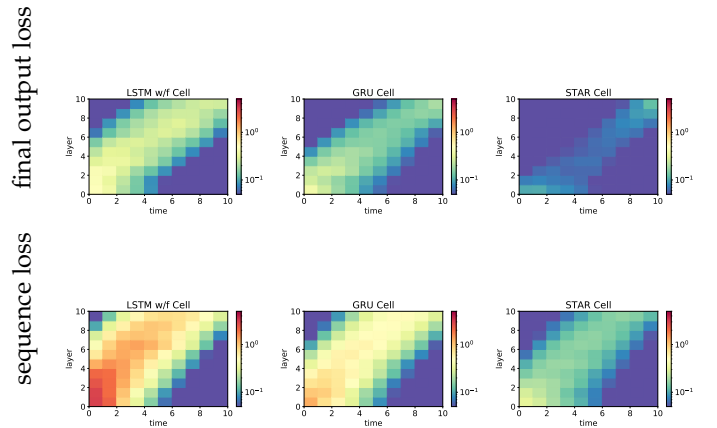


Fig. 14: Mean-normalised standard deviation of gradient magnitude for LSTM with only forget gate, GRU, and the proposed STAR cell. *top row*: loss $\mathcal{L}(h_T^L)$ only on final prediction. *bottom row*: loss $\mathcal{L}(h_1^L \dots h_T^L)$ over all time steps.

A.2.2 MNIST Dataset

In this section, we perform the same numerical analysis conducted before but using MNIST as input data. The goal is to verify whether during the first epoch the gradient propagation behaves in the same way as for the synthetic

dataset. First, in Fig 17 and Fig 18, we plot the evolution of the Hilbert-Schmidt norm (also called Frobenius norm) normalized by the square root of the hidden state size and the average hidden state value, respectively. The experiments are conducted using the proposed STAR method with MNIST as input data, the figures show the evolution of the norms and hidden states for the different layers of the recurrent network. The plots show the validity of our assumptions, during the initial training phase and with orthogonal matrix initialization the norm of the matrices is close to one, which translates to singular values close to one. The mean values of the hidden states are instead close to zero, which is consistent with the analysis conducted in Section 3. We reiterate that the plot shows how the *average value* of the hidden state remains close to 0 during training. This does *not* say that the hidden state is always $\equiv 0$.

Additionally, we show the gradient propagation in the two-dimensional lattice, as done in Fig. 13, for different cell types with MNIST as input data. We create 12-by-784 lattices with twelve layers RNNs. RNN weights are initialized the same way in the real experiments except the forget bias of the LSTM which is set to one (popular initialization scheme for the LSTM) due to numerical instability with the chrono method [39].

In Fig. 15 we can see that cells show similar behavior for the MNIST dataset. Even though on average STAR and GRU signal propagation looks fine, gradients within a single run oscillate a lot more for GRU (see Fig. 16) as seen in the previous numerical simulation (see Fig. 14).

A.3 Training details

We provide more details about training procedures for the experimental analysis in the main paper in this section.

A.3.1 Pixel-by-pixel MNIST

Following [39], chrono initialisation is applied for the bias term of k , b_k . The basic idea is that k should not be too large; such that the memory h can be retained over longer time intervals. The same initialisation is used for the input and forget bias of the LSTM and the RHN and for the forget bias of the LSTMw/f and the GRU. For the final prediction, a feedforward layer with softmax activation converts the hidden state to a class label. The numbers of hidden units in the RNN layers are set to 128. All networks are trained for 100 epochs with batch size 100, using the Adam optimizer [20] with learning rate 0.001, $\beta_1 = 0.9$ and $\beta_2 = 0.999$.

A.3.2 TUM time series classification

We use the same training procedure as described in the previous section for pixel-by-pixel MNIST. Again, a feed-forward layer is appended to the RNN output to obtain a prediction. The numbers of hidden units in the RNN layers is set to 128. All networks are trained for 30 epochs with batch size 500, using Adam [20] with learning rate 0.001, $\beta_1 = 0.9$ and $\beta_2 = 0.999$.

A.3.3 BreizhCrop time series classification

A feedforward layer is appended to the RNN output to obtain a prediction. The numbers of hidden units in the RNN layers is set to 128. All networks are trained for 30

epochs with batch size 1024, using Adam [20] with learning rate 0.001 and $\beta_1 = 0.9$ and $\beta_2 = 0.999$. The learning rate scheduler of [41] is used with 10 warm-up steps.

A.3.4 Adding problem / Copy memory

Following [39], chrono initialisation is applied for the bias term of k , b_k . The same initialisation is used for the input and forget bias of the LSTM and for the forget bias of the GRU. The number of hidden units is set to 128 for STAR and LSTM, 150 for GRU and 256 for vRNN. 2-layer STAR is used to have same number of parameters. Networks are trained using Adam [20] with learning rate 0.001, $\beta_1 = 0.9$ and $\beta_2 = 0.999$.

A.3.5 Music modeling

We follow the exact same experimental setup described in [37]. Baseline results are taken from [37]. The input sequence length is set to 200. STAR is trained for 500 iterations with batch size 1, using RMSProp. Dropout with keep probability 0.8 is applied. Other hyper-parameters (number of layer, momentum etc.) are searched as described in [37].

A.3.6 Character-level language modeling

We follow the exact same experimental setup described in [4]. Results for vRNN, LSTM, GRU and TCN are directly taken from [4]. The input sequence length is set to 400. The number of hidden units is set to 410 for STAR; therefore, the total number of parameters for 6-layers STAR makes 3M. STAR is trained for 50 epochs with batch size 32, using Adam [20] with learning rate 0.001, $\beta_1 = 0.9$ and $\beta_2 = 0.999$. The learning rate is decayed when the validation performance is no longer improved. Gradient clipping with 1 is applied. For IndRNN, the input sequence length is set to 50 because it performs poorly if set to 400. We set the number of hidden units to 660; therefore, the total number of parameters for 6-layers IndRNN is 3M and we train it for 100 epochs. Note that we took the IndRNN implementation from https://github.com/Sunnydreamrain/IndRNN_pytorch.

A.3.7 Hand-gesture recognition from video

All convolutional kernels are of size 3×3 . Each convolutional RNN layer has 64 filters. A shallow CNN is used to convert the hidden state to a label, with 4 layers that have filter depths 128, 128, 256 and 256, respectively. All models are trained with stochastic gradient descent (SGD) with momentum ($\beta = 0.9$). The batch size is set to 8, the learning rate starts at 0.001 and decays polynomially to 0.000001 over a total of 30 epochs. L_2 -regularisation with weight 0.00005 is applied to all parameters.

A.3.8 TUM image series pixel-wise classification

All convolutional kernels are of size 3×3 . Each convolutional RNN layer has 32 filters. A shallow CNN is used to convert the hidden state to a label, with 2 layers that have filter depths 64. All models are fitted with Adam [20]. The batch size is set to 1, the learning rate starts at 0.001 and decays polynomially to 0.000001 over a total of 25 epochs.

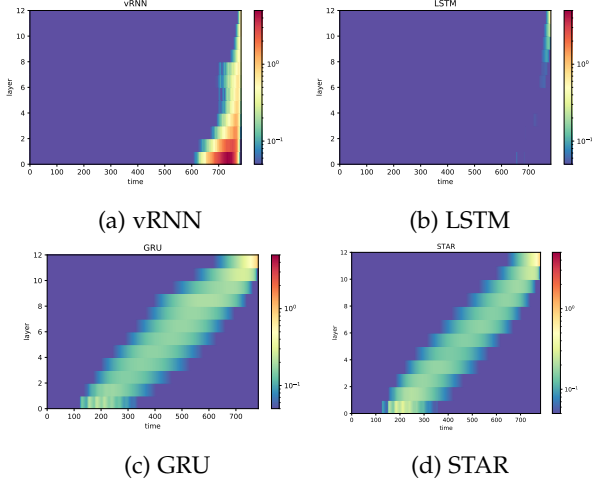


Fig. 15: Mean gradient magnitude w.r.t. the parameters for vRNN, LSTM, GRU, and the proposed STAR cell for MNIST dataset. Loss $\mathcal{L}(h_T^L)$ only on final prediction.

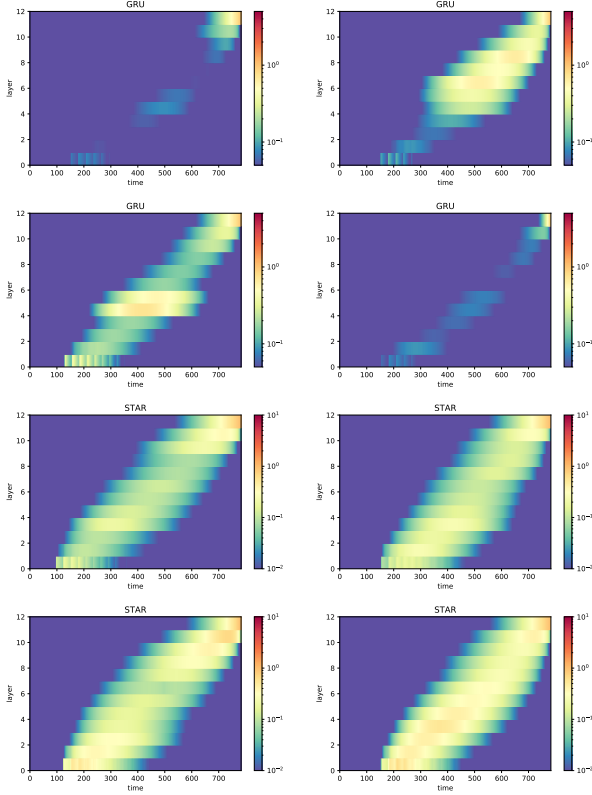


Fig. 16: Gradient magnitude comparison within a single run for MNIST dataset. *top two rows*: GRU samples. *bottom two rows*: STAR samples. Samples are randomly picked.

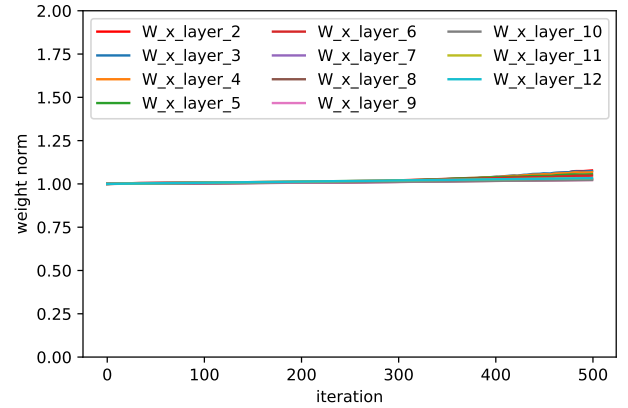
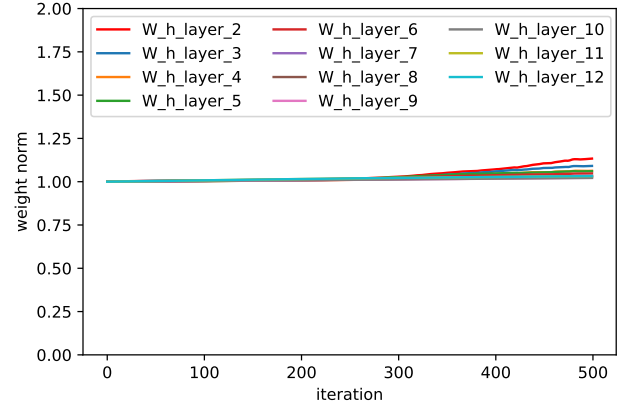
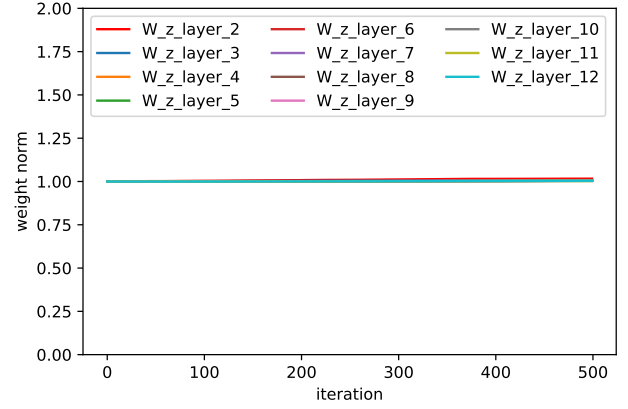


Fig. 17: Weight matrix norms of pix-by-pix MNIST during 1st epoch, the Hilbert–Schmidt norm, $\|A_{mxm}\| = \sqrt{\text{Tr}(AA^T)}$, divided by \sqrt{m} . Different curves correspond different layers.

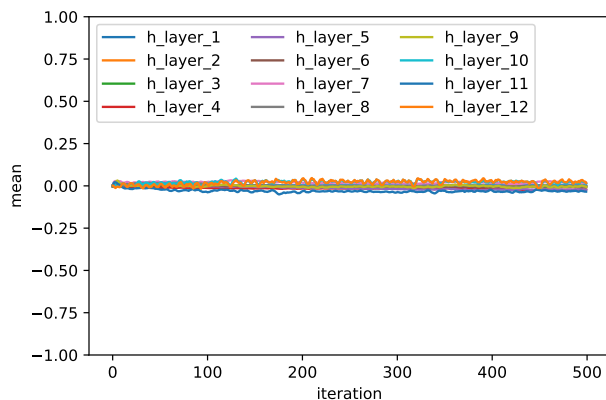


Fig. 18: Mean hidden state vector, $\mathbb{E}_{t,n}[\mathbf{h}^l]$ of pix-by-pix MNIST during 1st epoch. Different curves correspond different layers.

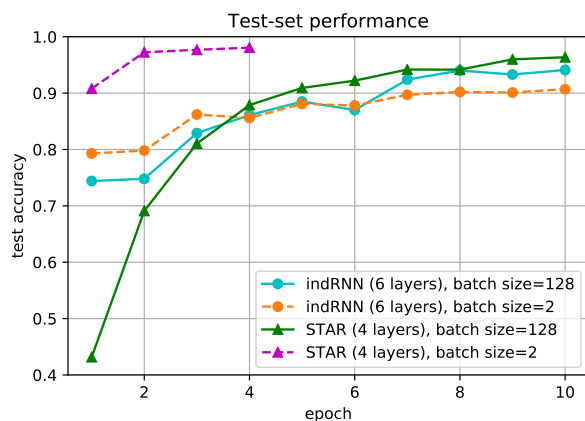


Fig. 19: Performance comparison for different batch sizes on the sequential MNIST task. If using a batch size of 128, both STAR and IndRNN converge to a solution; IndRNN is faster at the beginning but STAR eventually achieves better performance. IndRNN becomes very slow to train for a batch size of 2 (64 \times more steps per epoch) and it cannot achieve the same test performance as with the standard batch size (128). In contrast, STAR does not encounter these problems and clearly performs superior.