

Two Constraints to Scale Unbiased Recurrent Learning

Anonymous Authors¹

Abstract

Online scalable recurrent learning is challenging. Two popular gradient-based methods for recurrent learning are BPTT, and RTRL. BPTT looks at a complete sequence before computing gradients, and is unsuitable for online updates. RTRL can do online updates, but scales poorly with the number of parameters. In this paper, we propose two constraints that make RTRL scalable. We show that by either decomposing the network into independent modules, or learning a recurrent network incrementally, we can make RTRL scale linearly with the number of parameters. Both approaches result in different algorithms, that can be combined. We show the strengths and weaknesses of the proposed algorithms on a prediction learning benchmark inspired by animal learning.

1. Introduction

Structural credit-assignment — identifying how to change network parameters to improve predictions — is an essential component of learning in neural networks. Effective structural credit-assignment requires tracking the influence of parameters on future predictions. In recurrent networks, a parameter can influence the internal state of the network which, in turn, can affect a prediction made many steps in the future.

Back-Propagation Through Time (BPTT) (Werbos, 1988; Robinson and Fallside, 1987) is a popular algorithm for gradient-based structural credit-assignment in RNNs. BPTT extends the back-propagation algorithm for feed-forward networks — independently proposed by Werbos (1974) and Rumelhart *et al.* (1986) — to RNNs by storing network activations from prior steps, and repeatedly applying the chain-rule starting from the output of the network and ending at the activations at the beginning of the sequence. BPTT is unsuitable for online learning as it requires memory proportional to the length of the sequence. Moreover, it delays gradient computation until the end of the sequence. For online learning, this sequence can be never-ending or arbitrarily long.

RTRL—an alternative to BPTT—was proposed by Williams

and Zipser (1989). RTRL relies on forward-mode differentiation—using chain-rule to compute gradients in the direction of time—to compute gradients recursively. Unlike BPTT, RTRL does not delay gradient-computation until the final step. The memory requirement of RTRL also does not depend on the sequence length. As a result, it is a better candidate for real-time online learning. Unfortunately, RTRL requires maintaining the Jacobian $\frac{\partial h(t)}{\partial \theta}$ at every step, which requires $O(|h||\theta|)$ memory, where $|h|$ is the size of state of the network and $|\theta|$ is the number of total parameters. The Jacobian is recursively updated by applying chain rule as:

$$\frac{\partial h(t+1)}{\partial \theta} = \frac{\partial h(t+1)}{\partial h(t)} \frac{\partial h(t)}{\partial \theta} + \frac{\partial h(t+1)}{\partial \theta(t+1)},$$

which requires $O(|h|^2|\theta|)$ operations and scales poorly to large networks.

A promising direction to scale gradient-based credit-assignment to large networks is to approximate the gradient. Elman (1990) proposed to ignore the influence of parameters on future predictions completely for training RNNs. This resulted in a scalable but biased algorithm. Williams and Peng (1990) proposed a more general algorithm called Truncated BPTT (T-BPTT). T-BPTT tracks the influence of all parameters on predictions made up to k steps in the future. T-BPTT limits the BPTT computation to last k activation, and works well for a range of problems (Mikolov *et al.*, 2009, 2010; Sutskever, 2013 and Kapturowski *et al.*, 2018). Its main limitation is that the resultant gradient is blind to long-range dependencies. Mujika *et al.* (2018) showed that on a simple copy task, T-BPTT failed to learn dependencies beyond the truncation window. Tallec *et al.* (2017) demonstrated T-BPTT can even diverge when a parameter has a negative long-term effect on a target and a positive short-term effect.

RTRL can also be approximated to reduce its computational overhead. Ollivier *et al.* (2015) and Tallec *et al.* (2017) proposed NoBacktrack and UORO. Both of these algorithms provide stochastic unbiased estimates of the gradient and scale well. However, their estimates have high variance and require small step sizes for effective learning. Cooijmans and Martens (2019) and Menick *et al.* (2021) showed that, for practical problems, UORO does not perform well due to its high variance compared to other biased approximations.

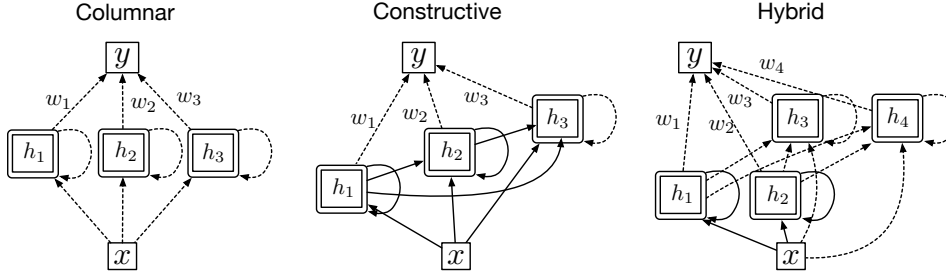


Figure 1. Three structures of recurrent neural networks that can be trained without truncation. Solid lines indicate parameters that are frozen, whereas dotted line parameters are being learned using gradients. Recurrent networks with a columnar structure can be trained end-to-end using gradients without any truncation, only requiring $O(n)$ operations and memory per step. However, they do not have hierarchical recurrent features—recurrent features made out of other recurrent features. Constructive networks have hierarchical recurrent features, however must be trained incrementally, one feature at a time, to prevent bias. Incremental learning is achieved by initializing all w_i to zero, and learning h_1, h_2 , and h_3 in order. Finally, columnar and constructive networks can be combined to get a hybrid network. The pairs (h_1, h_2) and (h_3, h_4) do not depend on each other, and can learn in parallel. However, (h_3, h_4) must be learned after (h_1, h_2) have been learned and fixed.

Menick *et al.* (2021) proposed an approximation to RTRL called SnAp- k . SnAp- k tracks the influence of a parameter on a state only if the parameter can influence the state within k steps. The bias introduced by SnAp- k is fundamentally different than the bias introduced by T-BPTT. SnAp-1 is computationally efficient but introduces significant bias. Laker k in SnAp- k reduces bias, but makes the algorithm computationally expensive.

2. Experiments

We evaluate our proposed methods on the sequential MNIST and the animal learning testbed (Rafiee *et al.* 2022). Animal learning testbed, on the other hand, is a sequential prediction problem. Certain patterns in the input—conditional stimuli (CS)—are followed by an unconditional stimuli (US). The goal of the agent is to predict the discounted future sum of US.

2.1. Sequential MNIST

In Sequential MNIST, the network sees a 28 x 28 MNIST image one row at a time and classifies the image into one of the ten categories—digits from 0 to 9.

3. Conclusion and Discussion

References

Bengio, Y., Bengio, S., and Cloutier, J. Learning a synaptic learning rule. Technical report, 1990.

Bengio, Y., Deleu, T., Rahaman, N., Ke, R., Lachapelle, S., Bilaniuk, O., Goyal, A., and Pal, C. A meta-transfer objective for learning to disentangle causal mechanisms. *ICLR*, 2020.

Cho, K., Van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., and Bengio, Y. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.

Cooijmans, T. and Martens, J. On the variance of unbiased online recurrent optimization. *arXiv preprint arXiv:1902.02405*, 2019.

Elman, J. L. Finding structure in time. *Cognitive science*, 1990.

Finn, C., Abbeel, P., and Levine, S. Model-agnostic meta-learning for fast adaptation of deep networks. In *ICML*, 2017.

Glorot, X. and Bengio, Y. Understanding the difficulty of training deep feedforward neural networks. In *International conference on artificial intelligence and statistics*. JMLR Workshop and Conference Proceedings, 2010.

Glorot, X., Bordes, A., and Bengio, Y. Deep sparse rectifier neural networks. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*. JMLR Workshop and Conference Proceedings, 2011.

Hochreiter, S. and Schmidhuber, J. Long short-term memory. *Neural computation*, 1997.

- Javed, K. and White, M. Meta-learning representations for continual learning. *NeurIPS*, 2019.
- Kapturowski, S., Ostrovski, G., Quan, J., Munos, R., and Dabney, W. Recurrent experience replay in distributed reinforcement learning. In *International conference on learning representations*, 2018.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *ICLR*, 2015.
- Li, Z., Zhou, F., Chen, F., and Li, H. Meta-sgd: Learning to learn quickly for few-shot learning. *arXiv:1707.09835*, 2017.
- Menick, J., Elsen, E., Evci, U., Osindero, S., Simonyan, K., and Graves, A. A practical sparse approximation for real time recurrent learning. *ICLR 2021*, 2020.
- Mikolov, T., Kopecky, J., Burget, L., Glembek, O., et al. Neural network based language models for highly inflective languages. In *International conference on acoustics, speech and signal processing*, 2009.
- Mikolov, T., Karafiát, M., Burget, L., Černocký, J., and Khudanpur, S. Recurrent neural network based language model. In *International speech communication association*, 2010.
- Mujika, A., Meier, F., and Steger, A. Approximating real-time recurrent learning with random kronecker factors. *Advances in Neural Information Processing Systems*, 2018.
- Ollivier, Y., Tallec, C., and Charpiat, G. Training recurrent networks online without backtracking. *arXiv preprint arXiv:1507.07680*, 2015.
- Rafiee, B., Abbas, Z., Ghiassian, S., Kumaraswamy, R., Sutton, R., Ludvig, E., and White, A. From eye-blinks to state construction: Diagnostic benchmarks for online representation learning. *Adaptive Behavior*, 2022.
- Robinson, A. and Fallside, F. *The utility driven dynamic error propagation network*. University of Cambridge Department of Engineering Cambridge, MA, 1987.
- Rumelhart, D. E., Hinton, G. E., and Williams, R. J. Learning representations by back-propagating errors. *Nature*, 1986.
- Santoro, A., Bartunov, S., Botvinick, M., Wierstra, D., and Lillicrap, T. Meta-learning with memory-augmented neural networks. In *International Conference on Machine Learning*, 2016.
- Schmidhuber, J. *Evolutionary principles in self-referential learning, or on learning how to learn: the meta-meta-... hook*. PhD thesis, Technische Universität München, 1987.
- Sutskever, I. *Training recurrent neural networks*. University of Toronto Toronto, Canada, 2013.
- Sutton, R. S. Adapting bias by gradient descent: An incremental version of delta-bar-delta. In *AAAI*. San Jose, CA, 1992.
- Tallec, C. and Ollivier, Y. Unbiased online recurrent optimization. *arXiv preprint arXiv:1702.05043*, 2017.
- Tange, O. Gnu parallel - the command-line power tool. *login: The USENIX Magazine*, 36(1):42–47, Feb 2011. doi: <http://dx.doi.org/10.5281/zenodo.16303>. URL <http://www.gnu.org/s/parallel>.
- Veeriah, V., Zhang, S., and Sutton, R. S. Crossprop: Learning representations by stochastic meta-gradient descent in neural networks. In *Machine Learning and Knowledge Discovery in Databases*, Cham, 2017.
- Werbos, P. Beyond regression:” new tools for prediction and analysis in the behavioral sciences. *Ph. D. dissertation, Harvard University*, 1974.
- Werbos, P. J. Generalization of backpropagation with application to a recurrent gas market model. *Neural networks*, 1988.
- Williams, R. J. and Peng, J. An efficient gradient-based algorithm for on-line training of recurrent network trajectories. *Neural computation*, 1990.
- Williams, R. J. and Zipser, D. A learning algorithm for continually running fully recurrent neural networks. *Neural computation*, 1989.

A. Forward-mode gradient computation of an LSTM cell

Here we derive the update equations for recursively computing the gradients of a single LSTM based recurrent column. Each column has a single hidden unit. Because all columns are identical, the same update equations can be used to learn the complete Columnar Network. The state of an LSTM column is updated using following equations:

$$i(t) = \sigma(W_i^T x_k(t) + u_i h(t-1) + b_i) \quad (1)$$

$$f(t) = \sigma(W_f^T x_k(t) + u_f h(t-1) + b_f) \quad (2)$$

$$o(t) = \sigma(W_o^T x_k(t) + u_o h(t-1) + b_o) \quad (3)$$

$$g(t) = \phi(W_g^T x_k(t) + u_g h(t-1) + b_g) \quad (4)$$

$$c(t) = f(t)c(t-1) + i(t)g(t) \quad (5)$$

$$h(t) = o(t)\phi(c(t)) \quad (6)$$

where σ and ϕ are the sigmoid and tanh activation functions, $h(t)$ is the state of the column at time t and $W_i^T x_k(t) = \sum_{k=1}^m W_{ik} x_k(t)$. The derivative of $\sigma(x)$ and $\phi(x)$ w.r.t to x are $\sigma(x)(1 - \sigma(x))$ and $(1 - \phi^2(x))$ respectively.

Let the length of input vector x be m . Then, W_i, W_f, W_o and W_g are vectors of length m whereas $u_i, b_i, u_f, b_f, u_o, b_o, u_g$ and b_g are scalars. We want to compute gradient of $h(t)$ with respect to all the parameters. We derive the update equations for $\frac{\partial h(t)}{\partial W_i}, \frac{\partial h(t)}{\partial u_i}, \frac{\partial h(t)}{\partial b_i}, \frac{\partial h(t)}{\partial W_f}, \frac{\partial h(t)}{\partial u_f}, \frac{\partial h(t)}{\partial b_f}, \frac{\partial h(t)}{\partial W_o}, \frac{\partial h(t)}{\partial u_o}, \frac{\partial h(t)}{\partial b_o}, \frac{\partial h(t)}{\partial W_g}, \frac{\partial h(t)}{\partial u_g}$, and $\frac{\partial h(t)}{\partial b_g}$ in the following sections.

A.1. $\frac{\partial h(t)}{\partial W_i}$

$W_i = (W_{i_1}, W_{i_2}, \dots, W_{i_m})$ is a vector of length m . Since all elements of W_i are symmetric, we show gradient derivation for W_{i_j} without loss of generality. Let

$$TH_{W_{i_j}}(t) := \frac{\partial h(t)}{\partial W_{i_j}} \quad (\text{By definition}) \quad (7)$$

$$TH_{W_{i_j}}(0) := 0 \quad (\text{By definition}) \quad (8)$$

$$TC_{W_{i_j}}(t) := \frac{\partial c(t)}{\partial W_{i_j}} \quad (\text{By definition}) \quad (9)$$

$$TC_{W_{i_j}}(0) := 0 \quad (\text{By definition}) \quad (10)$$

Then:

$$TH_{W_{i_j}}(t) = \frac{\partial}{\partial W_{i_j}} (o(t)\phi(c(t))) \quad \text{From equation 6 and definition 7}$$

$$= o(t) \frac{\partial \phi(c(t))}{\partial W_{i_j}} + \phi(c(t)) \frac{\partial o(t)}{\partial W_{i_j}} \quad \text{Product rule of differentiation}$$

$$= o(t)(1 - \phi^2(c(t))) \frac{\partial c(t)}{\partial W_{i_j}} + \phi(c(t)) \frac{\partial o(t)}{\partial W_{i_j}} \quad \text{Derivative of } \phi(x) \text{ is } (1 - \phi^2(x))$$

$$= o(t)(1 - \phi^2(c(t))) TC_{W_{i_j}}(t) + \phi(c(t)) \frac{\partial o(t)}{\partial W_{i_j}} \quad \text{From definition 9}$$

$$\frac{\partial o(t)}{\partial W_{i_j}} = \frac{\partial}{\partial W_{i_j}} \sigma(W_o^T x(t) + u_o h(t-1) + b_o) \quad \text{From equation 3}$$

$$= \sigma(y)(1 - \sigma(y)) u_o TH_{W_{i_j}}(t-1) \quad \text{Where } y \text{ equals } W_o^T x(t) + u_o h(t-1) + b_o$$

$$TC_{W_{i_j}}(t) = \frac{\partial c(t)}{\partial W_{i_j}} \quad \text{From definition 9}$$

$$= \frac{\partial}{\partial W_{i_j}} (f(t)c(t-1) + i(t)g(t)) \quad \text{From equation 5}$$

$$= f(t) TC_{W_{i_j}}(t-1) + c(t-1) \frac{\partial f(t)}{\partial W_{i_j}} + \frac{\partial}{\partial W_{i_j}} (i(t)g(t)) \quad \text{Product rule and definition 9}$$

$$= f(t)TC_{W_{i_j}}(t-1) + c(t-1)\frac{\partial f(t)}{\partial W_{i_j}} + i(t)\frac{\partial g(t)}{\partial W_{i_j}} + g(t)\frac{\partial i(t)}{\partial W_{i_j}} \quad \text{Product rule}$$

Where gradient of $g(t)$ w.r.t W_{i_j} is:

$$\begin{aligned} \frac{\partial g(t)}{\partial W_{i_j}} &= \frac{\partial}{\partial W_{i_j}} \phi(W_g^T x(t) + u_g h(t-1) + b_g) && \text{From equation 4} \\ &= (1 - \phi^2(y))u_g TH_{W_{i_j}}(t-1) && \text{Where } y \text{ equals } W_g^T x(t) + u_g h(t-1) + b_g \end{aligned}$$

, gradient of $f(t)$ w.r.t W_{i_j} is:

$$\begin{aligned} \frac{\partial f(t)}{\partial W_{i_j}} &= \frac{\partial}{\partial W_{i_j}} \sigma(W_f^T x(t) + u_f h(t-1) + b_f) && \text{From equation 2} \\ &= \sigma(y)(1 - \sigma(y))u_f TH_{W_{i_j}}(t-1) && \text{Where } y \text{ equals } W_f^T x(t) + u_f h(t-1) + b_f \end{aligned}$$

and gradient of $i(t)$ w.r.t W_{i_j} is:

$$\begin{aligned} \frac{\partial i(t)}{\partial W_{i_j}} &= \frac{\partial}{\partial W_{i_j}} \sigma(W_i^T x(t) + u_i h(t-1) + b_i) && \text{From equation 1} \\ &= \sigma(y)(1 - \sigma(y)) \left(x_j(t) + u_i TH_{W_{i_j}}(t-1) \right) && \text{Where } y \text{ equals } W_i^T x(t) + u_i h(t-1) + b_i \end{aligned}$$

The derivation shows that using two traces per parameter of W_i , it is possible to compute the gradient of $h(t)$ w.r.t W_i recursively. We provide the derivations for parameters u_i and b_i below. We skip the step-by-step derivations for the remaining parameters as they are similar.

A.2. $\frac{\partial h(t)}{\partial u_i}$

$$TH_{u_i}(t) := \frac{\partial h(t)}{\partial u_i} \quad (\text{By definition}) \quad (11)$$

$$TH_{u_i}(0) := 0 \quad (\text{By definition}) \quad (12)$$

$$TC_{u_i}(t) := \frac{\partial c(t)}{\partial u_i} \quad (\text{By definition}) \quad (13)$$

$$TC_{u_i}(0) := 0 \quad (\text{By definition}) \quad (14)$$

$$\begin{aligned} TH_{u_i}(t) &= \frac{\partial}{\partial u_i} (o(t)\phi(c(t))) && \text{From equation 6} \\ &= o(t)\frac{\partial \phi(c(t))}{\partial u_i} + \phi(c(t))\frac{\partial o(t)}{\partial u_i} && \text{Product rule} \\ &= o(t)(1 - \phi^2(c(t)))\frac{\partial c(t)}{\partial u_i} + \phi(c(t))\frac{\partial o(t)}{\partial u_i} && \text{Derivative of } \phi(x) \text{ is } 1 - \phi^2(x) \\ &= o(t)(1 - \phi^2(c(t)))TC_{u_i}(t) + \phi(c(t))\frac{\partial o(t)}{\partial u_i} && \text{Using definition 13} \\ \frac{\partial o(t)}{\partial u_i} &= \frac{\partial}{\partial u_i} \sigma(W_o^T x(t) + u_o h(t-1) + b_o) && \text{Using equations 3} \end{aligned}$$

$$= \sigma(x)(1 - \sigma(x))u_o TH_{u_i}(t - 1)$$

Where x equal $W_o^T x(t) + u_o h(t - 1) + b_o$

$$TC_{u_i}(t) = \frac{\partial c(t)}{\partial u_i}$$

Definition 13

$$= \frac{\partial}{\partial u_i}(f(t)c(t - 1) + i(t)g(t))$$

From equation 6

$$= f(t)TC_{u_i}(t - 1) + c(t - 1)\frac{\partial f(t)}{\partial u_i} + \frac{\partial}{\partial u_i}(i(t)g(t))$$

Product rule

$$= f(t)TC_{u_i}(t - 1) + c(t - 1)\frac{\partial f(t)}{\partial u_i} + i(t)\frac{\partial g(t)}{\partial u_i} + g(t)\frac{\partial i(t)}{\partial u_i}$$

Product rule

Gradient of $g(t)$ w.r.t u_i is:

$$\frac{\partial g(t)}{\partial u_i} = \frac{\partial}{\partial u_i}\phi(W_g^T x(t) + u_g h(t - 1) + b_g)$$

From equations 6

$$= (1 - \phi^2(y))u_g TH_{u_i}(t - 1)$$

Where y equals $W_g^T x(t) + u_g h(t - 1) + b_g$

, gradient of $f(t)$ w.r.t u_i is:

$$\frac{\partial f(t)}{\partial u_i} = \frac{\partial}{\partial u_i}\sigma(W_f^T x(t) + u_f h(t - 1) + b_f)$$

From equations 1

$$= \sigma(y)(1 - \sigma(y))u_f TH_{u_i}(t - 1)$$

Where y equals $W_f^T x(t) + u_f h(t - 1) + b_f$

and the gradient of $i(t)$ w.r.t u_i is

$$\frac{\partial i(t)}{\partial u_i} = \frac{\partial}{\partial u_i}\sigma(W_i^T x(t) + u_i h(t - 1) + b_i)$$

Using equations 1

$$= \sigma(y)(1 - \sigma(y))(h(t - 1) + u_i TH_{u_i}(t - 1))$$

Where y equals $W_i^T x(t) + u_i h(t - 1) + b_i$

A.3. $\frac{\partial h(t)}{\partial b_i}$

$$TH_{b_i}(t) := \frac{\partial h(t)}{\partial b_i}$$

(By definition)

(15)

$$TH_{b_i}(0) := 0$$

(By definition)

(16)

$$TC_{b_i}(t) := \frac{\partial c(t)}{\partial b_i}$$

(By definition)

(17)

$$TC_{b_i}(0) := 0$$

(By definition)

(18)

$$TH_{b_i}(t) = \frac{\partial}{\partial b_i}(o(t)\phi(c(t)))$$

From equation 6

$$= o(t)\frac{\partial \phi(c(t))}{\partial b_i} + \phi(c(t))\frac{\partial o(t)}{\partial b_i}$$

Product rule

$$= o(t)(1 - \phi^2(c(t)))\frac{\partial c(t)}{\partial b_i} + \phi(c(t))\frac{\partial o(t)}{\partial b_i}$$

Derivative of $\phi(x)$ is $1 - \phi^2(x)$

$$= o(t)(1 - \phi^2(c(t)))TC_{b_i}(t) + \phi(c(t))\frac{\partial o(t)}{\partial b_i}$$

From definition 17

$$\begin{aligned}
 \frac{\partial o(t)}{\partial b_i} &= \frac{\partial}{\partial b_i} \sigma(W_o^T x(t) + u_o h(t-1) + b_o) && \text{From equations 3} \\
 &= \sigma(y)(1 - \sigma(y))u_o TH_{b_i}(t-1) && \text{Where } y \text{ equal } W_o^T x(t) + u_o h(t-1) + b_o \\
 TC_{b_i}(t) &= \frac{\partial c(t)}{\partial b_i} && \text{From definition 17} \\
 &= \frac{\partial}{\partial b_i} (f(t)c(t-1) + i(t)g(t)) && \text{From equation 5} \\
 &= f(t)TC_{b_i}(t-1) + c(t-1)\frac{\partial f(t)}{\partial b_i} + \frac{\partial}{\partial b_i} i(t)g(t) && \text{Product rule} \\
 &= f(t)TC_{b_i}(t-1) + c(t-1)\frac{\partial f(t)}{\partial b_i} + i(t)\frac{\partial g(t)}{\partial b_i} + g(t)\frac{\partial i(t)}{\partial b_i} && \text{Product rule}
 \end{aligned}$$

Where gradient of $g(t)$ w.r.t b_i is:

$$\begin{aligned}
 \frac{\partial g(t)}{\partial b_i} &= \frac{\partial}{\partial b_i} \phi(W_g^T x(t) + u_g h(t-1) + b_g) && \text{From equation 4} \\
 &= (1 - \phi^2(y))u_g TH_{b_i}(t-1) && \text{Where } y \text{ equal } W_g^T x(t) + u_g h(t-1) + b_g
 \end{aligned}$$

, gradient of $f(t)$ w.r.t b_i is:

$$\begin{aligned}
 \frac{\partial f(t)}{\partial b_i} &= \frac{\partial}{\partial b_i} \sigma(W_f^T x(t) + u_f h(t-1) + b_f) && \text{From equation 2} \\
 &= \sigma(y)(1 - \sigma(y))u_f TH_{b_i}(t-1) && \text{Where } y \text{ equal } W_f^T x(t) + u_f h(t-1) + b_f
 \end{aligned}$$

and gradient of $i(t)$ w.r.t b_i is:

$$\begin{aligned}
 \frac{\partial i(t)}{\partial b_i} &= \frac{\partial}{\partial b_i} \sigma(W_i^T x(t) + u_i h(t-1) + b_i) && \text{From equation 1} \\
 &= \sigma(y)(1 - \sigma(y))(u_i TH_{b_i}(t-1) + 1) && \text{Where } y \text{ equal } W_i^T x(t) + u_i h(t-1) + b_i
 \end{aligned}$$

A.4. $\frac{\partial h(t)}{\partial W_{f_j}}$

The derivations for the remaining parameters is analogous to what previous derivations. The final equations are as follows.

$$\begin{aligned}
 \frac{\partial g(t)}{\partial W_{f_j}} &= (1 - \phi^2(y))(u_g TH_{W_{f_j}}(t-1)) \\
 \frac{\partial f(t)}{\partial W_{f_j}} &= \sigma(y)(1 - \sigma(y))(x_j + u_f TH_{W_{f_j}}(t-1)) \\
 \frac{\partial i(t)}{\partial W_{f_j}} &= \sigma(y)(1 - \sigma(y))(u_i TH_{W_{f_j}}(t-1)) \\
 \frac{\partial o(t)}{\partial W_{f_j}} &= \sigma(y)(1 - \sigma(y))(u_o TH_{W_{f_j}}(t-1)) \\
 TC_{W_{f_j}} &= f(t)TC_{f_j}(t-1) + c(t-1)\frac{\partial f(t)}{\partial b_i} + i(t)\frac{\partial g(t)}{\partial b_i} + g(t)\frac{\partial i(t)}{\partial b_i} \\
 TH_{W_{f_j}} &= o(t)(1 - \phi^2(c(t)))TC_{W_{f_j}}(t) + \phi(c(t))\frac{\partial o(t)}{\partial W_{ij}}
 \end{aligned} \tag{19}$$

A.5. $\frac{\partial h(t)}{\partial W_{o_j}}$

$$\begin{aligned}
 \frac{\partial g(t)}{\partial W_{o_j}} &= (1 - \phi^2(y))(u_g TH_{W_{o_j}}(t-1)) \\
 \frac{\partial f(t)}{\partial W_{o_j}} &= \sigma(y)(1 - \sigma(y))(u_f TH_{W_{o_j}}(t-1)) \\
 \frac{\partial i(t)}{\partial W_{o_j}} &= \sigma(y)(1 - \sigma(y))u_i TH_{W_{o_j}}(t-1) \\
 \frac{\partial o(t)}{\partial W_{o_j}} &= \sigma(x)(1 - \sigma(x))(x_j + u_o TH_{W_{o_j}}(t-1)) \\
 TC_{W_{o_j}} &= f(t)TC_{o_j}(t-1) + c(t-1)\frac{\partial f(t)}{\partial b_i} + i(t)\frac{\partial g(t)}{\partial b_i} + g(t)\frac{\partial i(t)}{\partial b_i} \\
 TH_{W_{o_j}} &= o(t)(1 - \phi^2(c(t)))TC_{W_{o_j}}(t) + \phi(c(t))\frac{\partial o(t)}{\partial W_{ij}}
 \end{aligned} \tag{20}$$

A.6. $\frac{\partial h(t)}{\partial W_{g_j}}$

$$\begin{aligned}
 \frac{\partial g(t)}{\partial W_{g_j}} &= (1 - \phi^2(y))(x_j + u_g TH_{W_{g_j}}(t-1)) \\
 \frac{\partial f(t)}{\partial W_{g_j}} &= \sigma(y)(1 - \sigma(y))(u_f TH_{W_{g_j}}(t-1)) \\
 \frac{\partial i(t)}{\partial W_{g_j}} &= \sigma(y)(1 - \sigma(y))(u_i TH_{W_{g_j}}(t-1)) \\
 \frac{\partial o(t)}{\partial W_{g_j}} &= \sigma(x)(1 - \sigma(x))(u_o TH_{W_{g_j}}(t-1)) \\
 TC_{W_{g_j}} &= f(t)TC_{g_j}(t-1) + c(t-1)\frac{\partial f(t)}{\partial b_i} + i(t)\frac{\partial g(t)}{\partial b_i} + g(t)\frac{\partial i(t)}{\partial b_i} \\
 TH_{W_{g_j}} &= o(t)(1 - \phi^2(c(t)))TC_{W_{g_j}}(t) + \phi(c(t))\frac{\partial o(t)}{\partial W_{ij}}
 \end{aligned} \tag{21}$$

A.7. $\frac{\partial h(t)}{\partial u_o}$

$$\begin{aligned}
 \frac{\partial g(t)}{\partial u_o} &= (1 - \phi^2(y))(u_g TH_{u_o}(t-1)) \\
 \frac{\partial f(t)}{\partial u_o} &= \sigma(y)(1 - \sigma(y))(u_f TH_{u_o}(t-1)) \\
 \frac{\partial i(t)}{\partial u_o} &= \sigma(y)(1 - \sigma(y))(u_i TH_{u_o}(t-1)) \\
 \frac{\partial o(t)}{\partial u_o} &= \sigma(x)(1 - \sigma(x))(u_o TH_{u_o}(t-1) + h(t-1)) \\
 TC_{u_o} &= f(t)TC_{i_j}(t-1) + c(t-1)\frac{\partial f(t)}{\partial b_i} + i(t)\frac{\partial g(t)}{\partial b_i} + g(t)\frac{\partial i(t)}{\partial b_i} \\
 TH_{u_o} &= o(t)(1 - \phi^2(c(t)))TC_{u_o}(t) + \phi(c(t))\frac{\partial o(t)}{\partial W_{ij}}
 \end{aligned} \tag{22}$$

A.8. $\frac{\partial h(t)}{\partial u_f}$

$$\begin{aligned}
 \frac{\partial g(t)}{\partial u_f} &= (1 - \phi^2(y))(u_g TH_{u_f}(t-1)) \\
 \frac{\partial f(t)}{\partial u_f} &= \sigma(y)(1 - \sigma(y))(u_f TH_{u_f}(t-1) + h(t-1)) \\
 \frac{\partial i(t)}{\partial u_f} &= \sigma(y)(1 - \sigma(y))(u_i TH_{u_f}(t-1)) \\
 \frac{\partial o(t)}{\partial u_f} &= \sigma(x)(1 - \sigma(x))(u_o TH_{u_f}(t-1)) \\
 TC_{u_f} &= f(t)TC_{i_j}(t-1) + c(t-1)\frac{\partial f(t)}{\partial b_i} + i(t)\frac{\partial g(t)}{\partial b_i} + g(t)\frac{\partial i(t)}{\partial b_i} \\
 TH_{u_f} &= o(t)(1 - \phi^2(c(t)))TC_{u_f}(t) + \phi(c(t))\frac{\partial o(t)}{\partial W_{ij}}
 \end{aligned} \tag{23}$$

A.9. $\frac{\partial h(t)}{\partial u_g}$

$$\begin{aligned}
 \frac{\partial g(t)}{\partial u_g} &= (1 - \phi^2(y))(u_g TH_{u_g}(t-1) + h(t-1)) \\
 \frac{\partial f(t)}{\partial u_g} &= \sigma(y)(1 - \sigma(y))(u_f TH_{u_g}(t-1)) \\
 \frac{\partial i(t)}{\partial u_g} &= \sigma(y)(1 - \sigma(y))(u_i TH_{u_g}(t-1)) \\
 \frac{\partial o(t)}{\partial u_g} &= \sigma(x)(1 - \sigma(x))(u_o TH_{u_g}(t-1)) \\
 TC_{u_g} &= f(t)TC_{i_j}(t-1) + c(t-1)\frac{\partial f(t)}{\partial b_i} + i(t)\frac{\partial g(t)}{\partial b_i} + g(t)\frac{\partial i(t)}{\partial b_i} \\
 TH_{u_g} &= o(t)(1 - \phi^2(c(t)))TC_{u_g}(t) + \phi(c(t))\frac{\partial o(t)}{\partial W_{ij}}
 \end{aligned} \tag{24}$$

A.10. $\frac{\partial h(t)}{\partial b_g}$

$$\begin{aligned}
 \frac{\partial g(t)}{\partial b_g} &= (1 - \phi^2(y))(u_g TH_{b_g}(t-1) + 1) \\
 \frac{\partial f(t)}{\partial b_g} &= \sigma(y)(1 - \sigma(y))(u_f TH_{b_g}(t-1)) \\
 \frac{\partial i(t)}{\partial b_g} &= \sigma(y)(1 - \sigma(y))(u_i TH_{b_g}(t-1)) \\
 \frac{\partial o(t)}{\partial b_g} &= \sigma(x)(1 - \sigma(x))(u_o TH_{b_g}(t-1)) \\
 TC_{b_g} &= f(t)TC_{i_j}(t-1) + c(t-1)\frac{\partial f(t)}{\partial b_i} + i(t)\frac{\partial g(t)}{\partial b_i} + g(t)\frac{\partial i(t)}{\partial b_i} \\
 TH_{b_g} &= o(t)(1 - \phi^2(c(t)))TC_{b_g}(t) + \phi(c(t))\frac{\partial o(t)}{\partial W_{ij}}
 \end{aligned} \tag{25}$$

A.11. $\frac{\partial h(t)}{\partial b_f}$

$$\begin{aligned}
 \frac{\partial g(t)}{\partial b_f} &= (1 - \phi^2(y))(u_g TH_{b_f}(t-1)) \\
 \frac{\partial f(t)}{\partial b_f} &= \sigma(y)(1 - \sigma(y))(u_f TH_{b_f}(t-1) + 1) \\
 \frac{\partial i(t)}{\partial b_f} &= \sigma(y)(1 - \sigma(y))(u_i TH_{b_f}(t-1)) \\
 \frac{\partial o(t)}{\partial b_f} &= \sigma(x)(1 - \sigma(x))(u_o TH_{b_f}(t-1)) \\
 TC_{b_f} &= f(t)TC_{i_j}(t-1) + c(t-1)\frac{\partial f(t)}{\partial b_i} + i(t)\frac{\partial g(t)}{\partial b_i} + g(t)\frac{\partial i(t)}{\partial b_i} \\
 TH_{b_f} &= o(t)(1 - \phi^2(c(t)))TC_{b_f}(t) + \phi(c(t))\frac{\partial o(t)}{\partial W_{ij}}
 \end{aligned} \tag{26}$$

A.12. $\frac{\partial h(t)}{\partial b_o}$

$$\begin{aligned}
 \frac{\partial g(t)}{\partial b_o} &= (1 - \phi^2(y))(u_g TH_{b_o}(t-1)) \\
 \frac{\partial f(t)}{\partial b_o} &= \sigma(y)(1 - \sigma(y))(u_f TH_{b_o}(t-1)) \\
 \frac{\partial i(t)}{\partial b_o} &= \sigma(y)(1 - \sigma(y))(u_i TH_{b_o}(t-1)) \\
 \frac{\partial o(t)}{\partial b_o} &= \sigma(x)(1 - \sigma(x))(u_o TH_{b_o}(t-1) + 1) \\
 TC_{b_o} &= f(t)TC_{i_j}(t-1) + c(t-1)\frac{\partial f(t)}{\partial b_i} + i(t)\frac{\partial g(t)}{\partial b_i} + g(t)\frac{\partial i(t)}{\partial b_i} \\
 TH_{b_o} &= o(t)(1 - \phi^2(c(t)))TC_{b_o}(t) + \phi(c(t))\frac{\partial o(t)}{\partial W_{ij}}
 \end{aligned} \tag{27}$$