

Definition 1 (TRDS) Let $T(U)$ be the set of rooted labeled ordered trees immediate subtrees t_1, t_2, \dots, t_d . The mapping $g : T(U) \rightarrow M$ is defined by $h : M \rightarrow \mathbb{E}^q$ where $q \in \mathbb{N}$. Let $t \in T(U)$ be an arbitrary tree with $d \leq k$ of the “empty tree” and two functions f, h of the type $f : U \times M^k \rightarrow M$ and $(\text{TRDS}) (M, \text{nil}, f, h)$ consists of a state space $M \subseteq \mathbb{E}^m$, an encoding $\text{nil} \in M$ denote the vector concatenation operator. A tree-recursive dynamical system (where nodes carry labels of the type U) with maximum outdegree k and let \oplus ,

Recently, the “zoo” of analog models of computation and learning has been extended by some inhabitants that are capable to capture structured information (for an introduction based on a probabilistic framework see Frasconi et al. [2]). For the purpose of this paper we focus our attention to the class of tree-recursive dynamical systems (TRDS) which can be regarded as deterministic machines designed to map rooted labeled ordered trees to the Euclidean space.

This paper investigates whether and how RTT can be generalized – while conserving its appealing algorithmic properties – to calculate the gradient information for models operating on the domain of rooted labeled ordered trees. The answer is partly negative. It turns out that a postorder traversal of the tree has to be obeyed in order to keep the space consumption independent from the size of the input structures. By processing vertices in an inverse topological ordering the algorithm can also be applied on labeled directed acyclic graphs. However, we show that on this graph domain the memory consumption grows (in the worst case) linearly with the size of the input structure.

Abstract. Recently, the so-called *Backpropagation Through Structure* (BPTS) algorithm has to be favored to BPTT if long sequences are computed by continuous information. BPTS can be viewed as an extension of the well-known *Backpropagation Through Time* (BPTT) algorithm for discrete-time dynamical systems and sequence processing.

Andreas Küchler

Tree-Recursive Computation of Gradient Information for Structures

In this paper we investigate whether and how RTRL can be generalized – while conserving its appealing algorithmic properties – to operate for TRDs on the tree domain. We present the mathematical derivation of a RTRL-style gradient calculation for TRDS in Section 2 (for the special case of neural networks, see Sperduti and Starita [5]). The algorithmic perspective is elaborated in Section 3 and Section 4. Section 5 is to summarize the implications.

Recently, the well-known *Backpropagation Through Time* (BPTT) algorithm for discrete-time dynamical systems and sequence processing has been generalized to the so-called *Backpropagation Through Structure* (BPTS) scheme which allows to compute the first-order gradient for functional expressions where rooted labeled ordered trees (mapped by the FA) are embedded [3]. Both BPTT and BPTS are time-efficient algorithms, however their memory consumption grows linearly with the size of the input structures. The well-known *Real-time Recurrent Learning* (RTRL) algorithm works with constant space and is to be favored to (the functionally equivalent) BPTT if long sequences are going to be processed (we refer to Williams and Zipser [6]). Furthermore, it is applicable to an infinite stream of data, while BPTT is not.

The basic idea is to incorporate adaptivity (“learning”) into TDS models by choosing parameterized functions f , h and to estimate these parameters according to the given task and data. The learning task can often be adequately formulated as the optimization of an error measure E (in the parameter space) which depends somehow on the parameter vectors w_f , w_h (found in f , h) and on the output $\tilde{e}(t)$ that is computed by the given TRDS (under its current parameter settings) for trees from a given data set P . Gradient-based methods are usually taken into account to solve this optimization problem if the error measure is a continuous and differentiable function of the error.

The so-called *neural folding architecture* (FA) is an example for “neural networks” instances of the TRDS model (Goller and Krichel [3]). The mappings of weights and biases implemented by standard multi-layer neural networks. Several other neural network architectures (e.g., cascade correlation, neutral trees) have been used in a similar manner (Sperduti and Starita [5]). TRDS are very powerful models of computation. Hammer [4] proved the FA to be a universal approximator for mappings of the type $\Xi : T(\mathbb{R}^n) \rightarrow \mathbb{R}$.

where the operator root_ℓ extracts the root node of a tree and Δ its label information. Finally the mapping $\Xi : T(U) \hookrightarrow \mathbb{R}^d$ composed by the TRDS is defined as follows:

$$(I) \quad (\overbrace{l\mathfrak{m} \oplus \cdots \oplus l\mathfrak{m}}^{p-\mathfrak{q}} \oplus (p_7)b \oplus \cdots \oplus (z_7)b \oplus (\tau_7)b, ((7)mod 4)\chi)f = (7)b$$

tree-recursion

2 Tree-Recurisve Calculation of the Gradient

For reasons of simplicity we assume for a given TDS (M, null, f, h) with $M = \mathbb{R}^m$, $f : \mathbb{R}^{n+km} \rightarrow \mathbb{R}^m$, $h : \mathbb{R}^m \rightarrow \mathbb{R}^q$, $\text{null} \in \mathbb{R}^m$ the mappings f , h and the error function $E : U^f \times U_h \rightarrow \mathbb{R}$ to be continuous and differentiable on the whole parameter space $U^f \times U_h$. Let the data be given by $P = \{(s_1, E(s_1), (s_2, E(s_2), \dots, (s_p, E(s_p))\}$ where $s_i \in T(\mathbb{R}^n)$ and $E(s_i) \in \mathbb{R}^q$, $i = 1, 2, \dots, p$ denote the i -th function E : $T(\mathbb{R}^n) \rightarrow \mathbb{R}^q$ of the field Ξ computed by the TRDS under the current parameter assignment. Further, let w_f, w_h be the area $n_f = \dim(U^f)$, $n_h = \dim(U_h)$ -dimensional parameter vectors belonging to f, h . Without loss of generality the calculation of the gradient is performed by superimposed learning and the popular mean squared error function

$$\begin{array}{ll}
m \times f u \mathbb{W} \ni E & f^{\mathbf{M}} Q / ((\ell) x)^f Q =: ((\ell) x)^{f_i} E \\
\text{` } m \times u \mathbb{W} \ni (x) f \mathbf{F} & ((\ell) x)^{i+u(\ell-1)+u x^i f} f =: ((\ell) x)_{(x) f}^{f_i} f \\
b \times u \mathbb{W} \ni H & q^{\mathbf{M}} Q / ((\ell) \delta)^f q Q =: ((\ell) \delta)^{f_i} H \\
b \times u \mathbb{W} \ni q \mathbf{F} & ((\ell) \delta)^{i x^i f} q =: ((\ell) \delta)^{f_i} q f \\
m \times f u \mathbb{W} \ni G & f^{\mathbf{M}} Q / (\ell)^f g Q =: (\ell)^{f_i} G
\end{array}$$

Thus, the objective is to calculate $e_f(t) = \partial E(t) / \partial W_f$ and $e_h(t) = \partial E(t) / \partial W_h$. Following function matrices as compact representations of partial derivatives, following tree $t \in T(\mathbb{P}^n)$. We define the following function $e_f(t) \in \mathbb{P}^n$ and $e_h(t) \in \mathbb{P}^n$ for a given input tree $t \in T(\mathbb{P}^n)$. We define the following functions of partial derivatives of E with respect to W_f and W_h :

$$\zeta((?s)^{\ell}\Xi - (?s)^{\ell}\Xi) \sum_b^{\mathbb{I}} \sum_d^{\mathbb{I}} \frac{\zeta}{\mathbb{I}} = (?s)\mathcal{E} \sum_d^{\mathbb{I}} \frac{\zeta}{\mathbb{I}} = \mathcal{E}$$

For reasons of simplicity we assume for a given TDS (M, null, f, h) with $M = \mathbb{R}^m$, $f : \mathbb{R}^{n+km} \rightarrow \mathbb{R}^m$, $h : \mathbb{R}^m \rightarrow \mathbb{R}^q$, $\text{null} \in \mathbb{R}^m$ the mappings f , h and the error function $E : U^f \times U_h \rightarrow \mathbb{R}$ to be continuous and differentiable on the whole parameter space $U^f \times U_h$. Let the data be given by $P = \{(s_1, E(s_1), (s_2, E(s_2), \dots, (s_p, E(s_p))\}$ where $s_i \in T(\mathbb{R}^n)$ and $E(s_i) \in \mathbb{R}^q$, let $\Xi_j = \{(s_1, E(s_1), (s_2, E(s_2), \dots, (s_p, E(s_p))\}$ where $s_i \in T(\mathbb{R}^n)$ and $E(s_i) \in \mathbb{R}^q$, $j = 1, 2, \dots, q$ denote the j -th function $\Xi_j : T(\mathbb{R}^n) \rightarrow \mathbb{R}$ of the field Ξ computed by the TDS under the current parameter assignment. Further, let w_f, w_h be the area $n_f = \dim(U^f)$, $n_h = \dim(U_h)$ -dimensional parameter vectors belonging to f, h . Without loss of generality the calculation of the gradient is performed by superimposing learning and the popular mean squared error function

Equation 3 measures the total impact of the variables w_f on the target g , i.e. the direct effects ($F(x(t))$, via f) and the indirect effects (via the tree-recursive composition of f on the immediate subtrees). One can easily verify that the recursion is well-founded.

$$\underbrace{\mathbb{I}^{in} \oplus \cdots \oplus \mathbb{I}^{in}}_{p-3} \oplus (({}^p\mathbb{I})\mathbf{x})f \oplus \cdots \oplus (({}^1\mathbb{I})\mathbf{x})f \oplus ((\mathbb{I})too,r)\chi = (\mathbb{I})\mathbf{x}$$

$$((\mathbb{I})\mathbf{x})_{(\mathbb{I})f}\mathbf{f} \cdot ({}^i\mathbb{I})\mathfrak{G} \sum_p^{\mathbb{I}=?} + ((\mathbb{I})\mathbf{x})\mathbb{H} = (\mathbb{I})\mathbb{G}$$

Now, let $t \in T(\mathbb{N}^n)$ be an arbitrary tree with $d \leq k$ immediate subtrees t_1, t_2, \dots, t_d . The tree-recurstive dynamics (see Equation 1) can be utilized to directly evolve the partial derivative $G(t)$ to

$$(2) \quad ((\varphi)\mathbf{z} \cdot ((\varphi)\delta)_y \mathbf{f}) \cdot (\varphi)\mathbf{G} = (\varphi)_y \mathbf{G} \quad (\varphi)\mathbf{z} \cdot ((\varphi)\delta) \mathbf{H} = (\varphi)_y \mathbf{H}$$

Note that $h_{j,i} := \partial h_j / \partial x_i$ is a shorthand for the partial derivative of h_j with respect to its i -th argument, i.e., $J_f^{(r)}$ and \mathbf{J}_f are the transposed Jacobian matrices of the mappings f and h at the locations $\mathbf{x}(t) \in \mathbb{R}^{n+m}$ and $g(t) \in \mathbb{R}^m$. The value k denotes the maximum outdegree found in the tree domain $T(\mathbb{R}^n)$ and $\mathbf{z}(t) := (\Xi(t) - \Xi(t)) \in \mathbb{R}^k$ the difference between computed and desired output for a given input tree. By applying the well-known generalized chain rule for differentials on the the composition $\Xi = h \circ g$ (see Definition 1) we get

is source of potential errors.

²However, this requires additional evaluations of f and h in the environment of $\mathbf{x}(t)$ and

direct translation of Equation 3. b) For each node v of the input tree memory ALg. 1 (Lines 3–4) directly implements Equation 2, Lines 5–11 of ALg. 2 are a direct node under consideration have already been visited before. Furthermore, input tree will be touched exactly once. Secondly, all subtrees (if there are any) of the node under consideration that each node of the

Proof (Sketch) a) The postorder traversal guarantees that each node of the input tree is calculated to $\Theta(2kn_f m + n_h + a + n)$. b) The memory consumption is independent from the size of the input tree, i.e. is calculated to $\Theta(2kn_f m + n_h + a + n)$.

putted by Algorithm 1 (and Algorithm 2) for any given (compatible) pair of input tree t and target value $\Xi(t)$. This fact by interlacing these computations into one tree-recurisve process. the state $f(\mathbf{x}(t))$ are defined by tree-recurisve on t . Algorithm 2 makes use of both the computation of the gradient information $G(t)$ and the computation that and $u \cdot \mathbf{y} (O(1) \text{ space})$. By a careful analysis of Equation 3 one can observe that node v has to be augmented by two entries to keep the memory pointers $u \cdot G$ is returned by a claim operation. Furthermore, the data structure representing a tree lists (of claimed and released memory cells) and only a pointer to a released cell and RELEASE-MEMORY take $O(1)$ time if implemented as two double-linked binsed with an explicit memory management. The operators CLAIM-MEMORY and Algorithm 2 specifies a postorder left-to-right tree traversal which is com-

bined with a finite Taylor series expansion of the original mappings². The derivatives by a finite Taylor series expansion might be to numerically approximate the partial code. An alternative solution might be to numerically approximate the partial code. Together with the mappings f and h) compiled a priori into the program are (together with the mappings f and h) determined in an analytical way and target value $\Xi(t)$ as input. We assume that the matrices of partial derivatives ALgortithm 1 (TRGC) expects a pair of input tree $t \in T(\mathbb{N})$ and desired

1: $mem \rightarrow \text{ALLOCATE-MEMORY}(2k, (n_f \times m + m))$	Algorithm 1 TRGC($t, \Xi(t)$)
2: $(\mathbf{y}, G) \rightarrow \text{TRGC-POSTORDER-TREE-WALK}(root(t), mem)$	
3: $\mathbf{z} \rightarrow h(\mathbf{y}) - \Xi(t)$	
4: $e_f \rightarrow G \cdot (\mathbf{f}_h'(\mathbf{y}) \cdot \mathbf{z}) ; e_h \rightarrow H(\mathbf{y}) \cdot \mathbf{z}$	
5: FREE-MEMORY(mem) ; return (e_f, e_h)	

Let $u = |t|$ denote the size (number of nodes) of the given input tree. One can observe that the formulation of Equation 3 is not tail-recursiv. Thus, a direct algorithmic translation would generate a tree-recurisve process (at runtime) allocating worst case $\Theta(u_f m + km)$ memory, i.e. equivalent to store u times the matrix $G(t)$ and the vector $\mathbf{x}(t)$. Can this behavior be improved?

3 The Algorithmic Perspective

Proof (Sketch) Consider the RLDOAG built of $u = 2(n + 1)$ vertices $V = \{v_1, v_2, \dots, v_{2n+2}\}$ ($n \in \mathbb{N}$) and the edges (in adjacency description): $\text{succ}(v_i) = \{v_{2n-i}, v_{i+1}\}$ if $1 \leq i \leq n + 1$, $\text{succ}(v_i) = (v_{i+1}, v_{i+1})$ if $n + 2 \leq i \leq 2n + 1$,

Proposition 2 The memory consumption of Algorithm 3 grows linearly with the number of vertices in the given (worst case) input RLDOAG.

Again, an explicit memory management is used to achieve an optimal memory consumption. The data structure representing a vertex is augmented by a set of vertices that are incident to it. The operator INCIDENT returns this set. Memory kept for the successors of a vertex v (under consideration) is freed if all vertices incident to it are larger or equal according to the chosen topological ordering (lines 12, 13).

Again, an explicit memory management is used to achieve an optimal memory consumption. The data structure representing a vertex v is already calculated for the successors v_i of the vertex v to be visited next. All the state information $v_i.y = f(x)$ is already calculated for the successors and the state information $v_i.y = f(x)$ is guaranteed to be guaranteed for the successors v_i . An inverse topological ordering it is guaranteed that the gradient information $v_i.G$ is used to calculate the gradient information $v_i.y$. By processing the given vertices V in an different computation for RLDOAGs. Algorithm 3 describes a variant of tree-reversive gradient computation. Algorithm 3 applies a considerable speed-up in iteration (CSE) technique which in turn leads to a considerable elimination (CSE) to RLDOAGs by applying a common subexpression elimination (CSE). On the other hand, (a set of) purely symbolic trees (i.e. from the domain $T(\Sigma)$) can be “compressed” to RLDOAGs by applying a common subexpression elimination (CSE). Thus, the expressive power of TRDS models can be slightly enriched. Graphs (RLDOAGs) if a topological ordering is obeyed while visiting the vertices [3].

The BPTS algorithm is known to work on rooted labeled directed acyclic graphs (RLDOAGs) if a topological ordering is obeyed while visiting the vertices [3].

4 Extending the Tree Domain

will be claimed only once (see Alg. 2, line 4). The postorder traversal ensures that after each node has been processed the release operation is applied on each successor node. It can easily be shown that at most $2k$ memory cells will be claimed at any point of time. The worst case becomes apparent when a rightmost node (of k siblings) rooted itself to k subtrees is going to be visited. ■

Algorithm 2 TRGC-POSTORDER-TREE-WALK(v, mem)	
1: $d \rightarrow \text{OUTDEGREE}(v) ; x \rightarrow \chi(v)$	
2: for $i \rightarrow 1$ to d do	
3: TRGC-POSTORDER-TREE-WALK(v_i, mem)	
4: $(v.G, v.y) \rightarrow \text{CLAIM-MEMORY}(mem, v.y)$	
5: for $i \rightarrow 1$ to d do	
6: $x \rightarrow x \oplus v_i.y ; \text{RELEASE-MEMORY}(mem, v.y)$	
7: for $i \rightarrow d + 1$ to k do	
8: $x \rightarrow x \oplus null$	
9: $v.y \rightarrow f(x) ; v.G \rightarrow F(x)$	
10: for $i \rightarrow 1$ to d do	
11: $v.G \rightarrow v.G + v_i.G - f(v_i)(x) ; \text{RELEASE-MEMORY}(mem, v.G)$	
12: return $(v.y, v.G, G)$	

- [1] P. Baldi, Gradient Descent Learning Algorithm Overview: A General Dynamical Systems Perspective. *IEEE Trans. on Neural Networks*, 6(1):182–195, 1995.

[2] P. Frasconi, M. Gori, and A. Sperduti. A General Framework for Adaptive Pro-cessing of Data Structures. *IEEE Trans. on Neural Networks*, 9(5):768–786, 1998.

[3] C. Goller and A. Krichev. Learning Task-Dependent Distributed Representations by Backpropagation Through Structure. In *Proc. IEEE Int. Conf. on Neural Networks (ICNN'96)*, pp. 347–352, 1996.

[4] B. Hammer. On the Approximation Capability of Recurrent Neural Networks. In *Proc. Int. Symp. on Neural Computation (NC'98)*, pp. 512–518, 1998.

[5] A. Sperduti and A. Statta. Supervised Neural Networks for the Classification of Structures. *IEEE Trans. on Neural Networks*, 8(3):714–735, 1997.

[6] R. J. Williams and D. Zipser. Gradient-Based Learning Algorithms for Recurrent Networks and Their Computational Completeness. In *Backpropagation: Theory, Networks and Applications and Architectures*, pp. 433–486. LEA, 1994.

References

RTRD for sequence processing has the appealing property that the memory consumption is independent from the size of the input structure. This property can be conserved by tree-recurisive gradient computation for finite trees if a postorder traversal is obeyed, while it is lost in the domain of finite RLDAGs. It is hard to imagine a reasonable definition of TRDS and the TRGC algorithm on the domain of infinite trees (RLDAGs).

5 Conclusion

■ Else due to the linear chain structure of the edges (v_i, v_{i+1}) and $\text{succ}(v_i) = (\)$, due to the vertex v_{n+1} is reached. Thus, $O(n)$ blocks of memory equivalent to store the matrix G and the state y have to be allocated.

```

1: Choose a topological ordering <TOP on V
2: for each  $v \in V$  in inverse topological ordering <TOP do
   3:    $d \rightarrow$  OUTDEGREE( $v$ ) :  $x \rightarrow A(v)$ 
   4:   ( $v$ ,  $G$ ,  $a$ ,  $y$ )  $\rightarrow$  ALLOCATE-MEMORY( $n_f \times m + m$ )
   5:   for  $i \rightarrow 1$  to  $d$  do
   6:      $x \rightarrow x \oplus a^i y$ 
   7:   for  $i \rightarrow d+1$  to  $k$  do
   8:      $x \rightarrow x \oplus u_i l$ 
   9:    $a$ ,  $G \rightarrow F(x)$  ;  $a$ ,  $y \rightarrow f(x)$ 
  10:  for  $i \rightarrow 1$  to  $d$  do
  11:     $a$ ,  $G \rightarrow a^i G + a^i G \cdot f^{(i)}(x)$ 
  12:    if  $Aw \in INCIDENT(a^i) : (v < TOP(w)) \vee (a = w)$  then
  13:      FREE-MEMORY( $a^i$ ,  $G$ ,  $a^i$ ,  $y$ )
  14: return ( $a$ ,  $y$ ,  $a$ ,  $G$ ) where  $v$  is the root vertex in  $V$ 

```
