

Higher-order Reverse Automatic Differentiation with emphasis on the third-order.

R. Gower* · A. Gower†

June 22, 2013

Abstract

It is commonly assumed that calculating third order information is too expensive for most applications. But we show that the directional derivative of the Hessian ($D^3 f(x) \cdot d$) can be calculated at a cost proportional to that of a state-of-the-art method for calculating the Hessian matrix. We do this by first presenting a simple procedure for designing high order reverse methods and applying it to deduce several methods including a reverse method that calculates $D^3 f(x) \cdot d$. We have implemented this method taking into account symmetry and sparsity, and successfully calculated this derivative for functions with a million variables. These results indicate that the use of third order information in a general non-linear solver, such as Halley-Chebyshev methods, could be a practical alternative to Newton's method.

1 Introduction

Derivatives permeate mathematics and engineering right from the first steps of calculus, which together with the Taylor series expansion is a central tool in designing models and methods of modern mathematics. Despite this, successful methods for automatically calculating derivatives of n -dimensional functions is a relatively recent development. Perhaps most notably amongst recent methods is the advent of Automatic Differentiation (*AD*), which has the remarkable achievement of the “cheap gradient principle”, wherein the cost of evaluating the gradient is proportional to that of the underlying function [11]. This *AD* success is not only limited to the gradient, there also exists a number of efficient *AD* algorithms for calculating Jacobian [6, 10] and Hessian matrices [8, 7], that can accommodate for large dimensional sparse instances. The same success has not been observed in calculating higher order derivatives.

The assumed cost in calculating high-order derivatives drives the design of methods, typically favouring the use of lower-order methods. In the optimization community it is generally assumed that calculating any third-order information is too costly, so the design of methods revolves around using first

*School of Mathematics and Maxwell Institute for Mathematical Sciences The University of Edinburgh, e-mail: gowerrobert@gmail.com

†School of Mathematics, Statistics and Applied Mathematics, National University of Ireland Galway

and second order information. We will show that third-order information can be used at a cost proportional to the cost of calculating the Hessian. This has an immediate application in third-order nonlinear optimization methods such as the Chebyshev-Halley Family [14]. Furthermore, the need for higher order differentiation finds applications in calculating quadratures [16, 4], bifurcations and periodic orbits [13]. In the fields of numerical integration and solution of PDE's, a lot of attention has been given to refining and adapting meshes to then use first and second-order approximations over these meshes. An alternative paradigm would be to use fixed coarse meshes and higher approximations. With the capacity to efficiently calculate high-order derivatives, this approach could become competitive and lift the fundamental deterrent in higher-order methods.

Current methods for calculating derivatives of order three or higher in the AD community typically propagate univariate Taylor series [12] or repeatedly apply the tangent and adjoint operations [18]. In these methods, each element of the desired derivative is calculated separately. If AD has taught us anything it is that we should not treat elements of derivatives separately, for their computation can be highly interlaced. The cheap gradient principle illustrates this well, for calculating the elements of the gradient separately yields a time complexity of n times that of simultaneously calculating all entries. This same principle should be carried over to higher order methods, that is, be wary of overlapping calculations in individual elements. Another alternative for calculating high order derivatives is the use of forward differentiation [19]. The drawback of forward propagation is that it calculates the derivatives of all intermediate functions, in relation to the independent variables, even when these do not contribute to the desired end result. For these reasons, we look at calculating high-order derivatives as a whole and focus on reverse AD methods.

An efficient alternative to AD is that the end users hand code their derivatives. Though with the advent of evermore complicated models, this task is becoming increasingly error prone, difficult to write efficient code, and, let's face it, boring. This approach also rules out methods that use high order derivatives, for no one can expect the end user to code the total and directional derivatives of high order tensors.

The article flows as follows, first we develop algorithms that calculate derivatives in a more general setting, wherein our function is described as a sequence of compositions of maps, Section 2. We then use Griewank and Walther's [11] state-transformations in Section 3, to translate a composition of maps into an AD setting and an efficient implementation. Numerical tests are presented in Section 4, followed by our conclusions in Section 5.

2 Derivatives of Sequences of Maps

In preparation for the AD setting, we first develop algorithms for calculating derivatives of functions that can be broken into a composition of operators

$$F(x) = \Psi^\ell \circ \Psi^{\ell-1} \circ \dots \circ \Psi^1(x). \quad (1)$$

for Ψ^i 's of varying dimension: $\Psi^1(x) \in C^2(\mathbb{R}^n, \mathbb{R}^{m_1})$ and $\Psi^i(x) \in C^2(\mathbb{R}^{m_{i-1}}, \mathbb{R}^{m_i})$, each $m_i \in \mathbb{N}$ and for $i = 2, \dots, \ell$, so that $F : \mathbb{R}^n \rightarrow \mathbb{R}^{m_\ell}$. From this we define a functional $f(x) = y^T F(x)$, where $y \in \mathbb{R}^{m_\ell}$, and develop methods for calculating

the *gradient* $\nabla f(x) = y^T DF(x)$, the *Hessian* $D^2 f(x) = y^T D^2 F(x)$ and the *Tensor* $D^3 f(x) = y^T D^3 F(x)$.

For a given $d \in \mathbb{R}^n$, we also develop methods for the directional derivative $DF(x) \cdot d$, $D^2 F(x) \cdot d$, the Hessian-vector product $D^2 f(x) \cdot d = y^T D^2 F(x) \cdot d$ and the Tensor-vector product $D^3 f(x) \cdot d = y^T D^3 F(x) \cdot d$. Notation will be gradually introduced and clarified as is required, including the definition of the preceding directional derivatives.

2.1 First-Order Derivatives

Taking the derivative of F , equation (1), and recursively applying the chain rule, we get

$$y^T DF = y^T D\Psi^\ell D\Psi^{\ell-1} \dots D\Psi^1. \quad (2)$$

Note that $y^T DF$ is the transpose of the gradient $\nabla(y^T F)$. For simplicity's sake, the argument of each function is omitted in (2), but it should be noted that $D\Psi^i$ is evaluated at the argument $(\Psi^{i-1} \circ \dots \circ \Psi^1)(x)$, for each i from 1 to ℓ . If each of these arguments has been recorded, the gradient of $y^T F(x)$ can be calculated with what's called a *reverse sweep* in Algorithm (1). Reverse, for it transverses the maps from the last Ψ^ℓ to the first Ψ^1 , the opposite direction in which (1) is evaluated. The intermediate stages of the gradient calculation are accumulated in the vector \bar{v} , its dimension changing from one iteration to the next. This will be a recurring fact in the matrices and vectors used to store the intermediate phases of the archetype algorithms presented in this article.

Algorithm 1: Archetype Reverse Gradient.

initialization: $\bar{v} = y$
for $i = \ell, \dots, 1$ **do**
 $\bar{v}^T \leftarrow \bar{v}^T D\Psi^i$
end
Output: $y^T DF(x) = \bar{v}^T$

For a given direction $d \in \mathbb{R}^n$, we define the directional derivative of $F(x)$ as

$$\frac{d}{dt} F(x + td) = D_i F(x + td) d_i := DF(x + td) \cdot d, \quad (3)$$

where we have omitted the summation symbol for i , and instead, use Einstein notation where a repeated indexes implies summation over that index. We use this notation throughout the article unless otherwise stated. Again using the chain-rule and (1), we find

$$DF(x) \cdot d = D\Psi^\ell D\Psi^{\ell-1} \dots D\Psi^1 \cdot d.$$

This can be efficiently calculated using a *forward sweep* of the computational graph, detailed below.

2.2 Second-Order Derivatives

Here we develop a reverse algorithm for calculating the Hessian $D^2(y^T F(x))$. First we determine the Hessian for F as a composition of two maps, then we use induction to design a method for when F is a composition of ℓ maps.

Algorithm 2: Archetype 1st Order Directional Derivative.

initialization: $\dot{v}^0 = d$
for $i = 1, \dots, \ell$ **do**
 $\dot{v} \leftarrow D\Psi^i \dot{v}$
end
Output: $DF(x) \cdot d = \dot{v}^\ell$

For $F(X) = \Psi^2 \circ \Psi^1(x)$ and $\ell = 2$, we find the Hessian by differentiating in the j -th and k -th coordinate,

$$D_{jk}(y_i F_i) = (y_i D_{rs} \Psi_i^2) D_j \Psi_r^1 D_k \Psi_s^1 + (y_i D_r \Psi_i^2) D_{jk} \Psi_r^1, \quad (4)$$

where the arguments have been omitted. So the (j, k) component of the Hessian $[D^2(y^T F)]_{jk} = D_{jk}(y^T F)$. The higher the order of the derivative, the more messy and unclear component notation becomes. A way around this issue is to use a tensor notation

$$y^T D^2 F \cdot (v, w) := y_i D_{jk} F_i v_j w_k,$$

and

$$(y^T D^2 F \cdot w) \cdot v := y^T D^2 F \cdot (v, w), \quad (5)$$

for any vectors $v, w \in \mathbb{R}^n$, and in general,

$$[y^T D^2 F \cdot (\Delta, \square)]_{t_2 \dots t_q s_2 \dots s_p} := y_i D_{t_1 s_1} F_i \Delta_{t_1 t_2 \dots t_q} \square_{s_1 s_2 \dots s_p}, \quad (6)$$

and

$$(y^T D^2 F \cdot \square) \cdot \Delta := y^T D^2 F \cdot (\Delta, \square) \quad (7)$$

for any compatible Δ and \square . To use a matrix notation for a composition of maps can be aesthetically unpleasant. Using this tensor notation the Hessian of $y^T F$, see equation (4), becomes

$$y^T D^2 F = y^T D^2 \Psi^2 \cdot (D\Psi^1, D\Psi^1) + y^T D\Psi^2 \cdot D^2 \Psi^1. \quad (8)$$

We recursively use the identity (8) to design an algorithm that calculates the Hessian of a function $y^T F(x)$ composed of ℓ maps, as defined in equation (1).

Algorithm 3: Archetype Reverse Hessian.

initialization: $\bar{v} = y, W = 0$
for $i = \ell, \dots, 1$ **do**
 $W \leftarrow W \cdot (D\Psi^i, D\Psi^i)$
 $W \leftarrow W + \bar{v}^T D^2\Psi^i$
 $\bar{v}^T \leftarrow \bar{v}^T D\Psi^i$
end
Output: $y^T D^2F \leftarrow W, y^T DF \leftarrow \bar{v}^T$

Proof of Algorithm: We will use induction on the number of compositions ℓ . For $\ell = 1$ the output is $W = y^T D^2\Psi^1$. Now we suppose the Algorithm is correct for $m - 1$ map compositions, and use this assumption to show that for $\ell = m$ the output is $W = y^T D^2F$. Let

$$y^T X = y^T \Psi^m \circ \dots \circ \Psi^2,$$

so that $y^T F = y^T X \circ \Psi^1$. Then at the end of the iteration $i = 2$, by the chain rule, $\bar{v}^T = y^T DX$ and, by induction, $W = y^T D^2X$. This way, at termination, or after the iteration $i = 1$, we get

$$\begin{aligned}
W &= y^T D^2X \cdot (D\Psi^1, D\Psi^1) + y^T DX \cdot D^2\Psi^1 \\
&= y^T D^2(X \circ \Psi^1) \quad [\text{Equation (8)}] \\
&= y^T D^2F. \quad \blacksquare
\end{aligned}$$

Now we take a small detour to show how to calculate Hessian-vector products in a similar manner. We do this because it is an important component of graph-coloring based algorithms for calculating the Hessian [7] and its complexity is surprisingly the same as evaluating $y^T F$, the underlying functional [3]. Thus, analogously, we calculate the directional derivative of the gradient $y^T DF(x)$, for $\ell = 2$,

$$y^T D_{jk} F d_k = y^T D_{rs} \Psi^2 D_j \Psi_r^1 D_k \Psi_s^1 d_k + y^T D_r \Psi^2 D_{jk} \Psi_r^1 d_k, \quad (9)$$

or simply,

$$\boxed{y^T D^2F \cdot d = y^T D^2\Psi^2 \cdot (D\Psi^1, D\Psi^1 \cdot d) + y^T D\Psi^2 \cdot D^2\Psi^1 \cdot d}, \quad (10)$$

and use this recursively to calculate the directional derivative of $y^T DF(x)$ in Algorithm 4. This algorithm was first described in [3].

Proof of Algorithm: Let $y^T F$ be a composition of ℓ maps as in (1) and

$$X^m = y^T \Psi^\ell \circ \dots \circ \Psi^m,$$

so that $X^{m-1} = X^m \circ \Psi^{m-1}$. The first **for** loop simply accumulates the directional derivative $DF \cdot d$. For the second **for** loop, we use an induction hypothesis that at the end of the $i = m$ iteration $w = D^2 X^m \cdot \dot{v}^{m-1}$. The first iteration, $i = \ell$, the output is $w = y^T D^2 \Psi^\ell \cdot d = D^2 X^\ell \cdot \dot{v}^{\ell-1}$. Now suppose our hypothesis is true for $i = m + 1$, so that at the end of the $i = m + 1$ iteration, by the induction hypothesis,

$$w = D^2 X^{m+1} \cdot \dot{v}^m = D^2 X^{m+1} \cdot D\Psi^m \cdot \dot{v}^{m-1},$$

Algorithm 4: Archetype Gradient Directional Derivative

initialization: $\dot{v}^0 = d, \bar{v} = y \in \mathbb{R}^{m_\ell}, w = 0 \in \mathbb{R}^{m_\ell}$
for $i = 1, \dots, \ell$ **do**
 $\dot{v}^i \leftarrow D\Psi^i \cdot \dot{v}^{i-1}$
end
for $i = \ell, \dots, 1$ **do**
 $w \leftarrow w \cdot D\Psi^i$
 $w \leftarrow w + \bar{v}^T D^2\Psi^i \cdot \dot{v}^{i-1}$
 $\bar{v}^T \leftarrow \bar{v}^T D\Psi^i$
end
Output: $y^T D^2 F(x) \cdot d \leftarrow w, y^T DF \leftarrow \bar{v}^T$

and, by calculus,

$$\bar{v}^T = y^T D\Psi^\ell \dots D\Psi^{m+1} = DX^{m+1}.$$

Then for the next step, the $i = m$ iteration,

$$\begin{aligned}
w &\leftarrow w \cdot D\Psi^m + \bar{v}^T D^2\Psi^m \cdot \dot{v}^{m-1} \\
&= (D^2 X^{m+1} \cdot D\Psi^m \cdot \dot{v}^{m-1}) \cdot D\Psi^m + DX^{m+1} \cdot D^2\Psi^m \cdot \dot{v}^{m-1} \\
&= D^2 X^{m+1} \cdot (D\Psi^m, D\Psi^m \cdot \dot{v}^{m-1}) + DX^{m+1} \cdot D^2\Psi^m \cdot \dot{v}^{m-1} \\
&= D^2(X^{m+1} \circ \Psi^m) \cdot \dot{v}^{i-1} \quad [\text{Equation (10)}] \\
&= D^2 X^m \cdot \dot{v}^{m-1}.
\end{aligned}$$

Thus by induction we have proved that at the end of the $i = 1$ iteration,

$$w = D^2 X^1 \cdot \dot{v}^0 = D^2 X^1 \cdot d = y^T D^2 F(x) \cdot d. \quad \blacksquare$$

2.3 Third-Order Methods

Now we move on to the directional derivative of $y^T D^2 F(x)$, that is, the derivative of $y^T D^2 F(x + td)$ in t , where $d \in \mathbb{R}^n$, to get

$$\begin{aligned}
\frac{d}{dt} y^T D^2 F(x + td) &= y_i \frac{d}{dt} D_{jk} F_i(x + td) \\
&= y_i D_{jkm} F_i(x + td) d_m \\
&:= y^T D^3 F(x + td) \cdot d.
\end{aligned} \tag{11}$$

Here our tensor notation really facilitates working with third-order derivatives. Using matrix notation would lead to confusing equations and possibly deter intuition. The notation conventions from before carry over naturally to third-order derivatives, with

$$(y^T D^3 F \cdot (\Delta, \square, \diamond))_{t_2 \dots t_q s_2 \dots s_p l_2 \dots l_r} := y_i D^3 F_{t_1 s_1 l_1 \Delta_{t_1 \dots t_q} \square_{s_1 \dots s_p} \diamond_{l_1 \dots l_r}}, \tag{12}$$

and

$$y^T D^3 F \cdot (\Delta, \square, \diamond) = (y^T D^3 F \cdot \diamond) \cdot (\Delta, \square) = ((y^T D^3 F \cdot \diamond) \cdot \square) \cdot \Delta, \tag{13}$$

for any compatible \triangle , \square and \diamond . We begin by calculating the directional derivative of a composition of two maps $F = \Psi^2 \circ \Psi^1$,

$$\begin{aligned}
& \frac{d}{dt} (y^T D^2 F(x + dt)) \\
&= D (y^T D^2 \Psi^2 \cdot (D\Psi^1, D\Psi^1)) \cdot d + D ((y^T D\Psi^2) \cdot D^2 \Psi^1) \cdot d \\
&= (y^T D^3 \Psi^2 \cdot D\Psi^1 \cdot d) \cdot (D\Psi^1, D\Psi^1) + (y^T D^2 \Psi^2) \cdot (D\Psi^1, D^2 \Psi^1 \cdot d) \\
&\quad + (y^T D^2 \Psi^2) \cdot (D^2 \Psi^1 \cdot d, D\Psi^1) + (y^T D\Psi^2) \cdot D^3 \Psi^1 \cdot d \\
&\quad + (y^T D^2 \Psi^2 \cdot D\Psi^1 \cdot d) \cdot D^2 \Psi^1,
\end{aligned}$$

in conclusion, after some rearrangement,

$$\begin{aligned}
y^T \frac{d}{dt} D^2 F(x + dt) &= y^T D^3 \Psi^2 \cdot (D\Psi^1, D\Psi^1, D\Psi^1 \cdot d) + y^T D\Psi^2 \cdot D^3 \Psi^1 \cdot d \\
&\quad + y^T D^2 \Psi^2 \cdot ((D\Psi^1, D^2 \Psi^1 \cdot d) + (D^2 \Psi^1 \cdot d, D\Psi^1) + (D^2 \Psi^1, D\Psi^1 \cdot d)) \quad (14)
\end{aligned}$$

As usual, we have omitted all arguments to the maps. The above applied recursively gives us the Reverse Hessian Directional Derivative Algorithm 5, or *RevHedir* for short. To prove the correctness of **RevHedir**, we use induction based on $X^m = y^T \Psi^\ell \circ \dots \circ \Psi^m$, working from $m = \ell$ backwards towards $m = 1$ to calculate $y^T D^3 F(x) \cdot d$.

Algorithm 5: Archetype Reverse Hessian Directional Derivative (**RevHedir**)

initialization: $\dot{v}^1 = d, \bar{v} = y, W = Td = 0 \in \mathbb{R}^{m_\ell \times m_\ell}$
for $i = 1, \dots, \ell$ **do**
 $\dot{v}^i \leftarrow D\Psi^i \cdot \dot{v}^{i-1}$
end
for $i = \ell, \dots, 1$ **do**
 $Td \leftarrow Td \cdot (D\Psi^i, D\Psi^i)$
 $Td \leftarrow Td + W \cdot ((D\Psi^i, D^2 \Psi^i \cdot \dot{v}^{i-1}) + (D^2 \Psi^i \cdot \dot{v}^{i-1}, D\Psi^i))$
 $Td \leftarrow Td + W \cdot (D^2 \Psi^i, D\Psi^i \cdot \dot{v}^{i-1})$
 $Td \leftarrow Td + \bar{v}^T D^3 \Psi^i \cdot \dot{v}^{i-1}$
 $W \leftarrow W \cdot (D\Psi^i, D\Psi^i) + \bar{v}^T D^2 \Psi^i$
 $\bar{v}^T \leftarrow \bar{v}^T D\Psi^i$
end
Output: $y^T D^3 F(x) \cdot d \leftarrow Td, y^T D^2 F \leftarrow W, y^T DF \leftarrow \bar{v}^T$

Proof of Algorithm: Our induction hypothesis is that at the end of the $i = m$ iteration $Td = y^T D^3 X^m \cdot \dot{v}^{i-1}$. After the first iteration $i = \ell$, paying attention to the initialization of the variables, we have that $Td = \bar{v}^T D^3 \Psi^\ell \cdot \dot{v}^{\ell-1} = y^T D^3 X^\ell \cdot \dot{v}^{\ell-1}$. Now suppose the hypothesis is true for iterations up to $m + 1$, so that at the beginning of the $i = m$ iteration $Td = y^T D^3 X^{m+1} \cdot \dot{v}^m$. To prove the hypothesis we need the following results: at the end of the $i = m$ iteration

$$\bar{v}^T = y^T DX^m \quad \text{and} \quad W = y^T D^2 X^m, \quad (15)$$

both are demonstrated in the proof of Algorithm 10. Now we are equipt to examine Td at the end of the $i = m$ iteration,

$$\begin{aligned} Td \leftarrow Td \cdot (D\Psi^m, D\Psi^m) + W \cdot ((D\Psi^m, D^2\Psi^m \cdot \dot{v}^{m-1}) + (D^2\Psi^m \cdot \dot{v}^{m-1}, D\Psi^m)) \\ + W \cdot (D^2\Psi^m, D\Psi^m \cdot \dot{v}^{m-1}) + \bar{v}^T D^3\Psi^m \cdot \dot{v}^{m-1}, \end{aligned}$$

using the induction hypothesis followed by property (13) we get $Td \cdot (D\Psi^m, D\Psi^m) = y^T D^3 X^{m+1} \cdot \dot{v}^m \cdot (D\Psi^m, D\Psi^m) = y^T D^3 X^{m+1} \cdot (D\Psi^m, D\Psi^m, \dot{v}^m)$, and from the algorithm $\dot{v}^m = D\Psi^m \cdot \dot{v}^{m-1}$. Then using equations (15) to substitute W and \bar{v}^T we arrive at

$$\begin{aligned} Td = & y^T D^3 X^{m+1} \cdot (D\Psi^m, D\Psi^m, D\Psi^m \cdot \dot{v}^{m-1}) \\ & + y^T D^2 X^m \cdot ((D\Psi^m, D^2\Psi^m \cdot \dot{v}^{m-1}) + (D^2\Psi^m \cdot \dot{v}^{m-1}, D\Psi^m)) \\ & + y^T D^2 X^m \cdot (D^2\Psi^m, D\Psi^m \cdot \dot{v}^{m-1}) + (y^T D X^m) D^3\Psi^m \cdot \dot{v}^{m-1} \\ = & y^T D^3 X^m \cdot \dot{v}^{m-1} \quad [\text{Using equation 14}]. \end{aligned}$$

Finally, after iteration $i = 1$, we have

$$Td = y^T D^3 X^1 \cdot \dot{v}^0 = y^T D^3 F \cdot d. \quad \blacksquare$$

As is to be expected, in the computation of the Tensor-vector product, only 2-dimensional tensor arithmetic, or matrix arithmetic, is used, and it is not necessary to form a 3-dimensional tensor. On the other hand, calculating the entire $y^T D^3 F$ Tensor does involve 3-dimensional arithmetic. The final archetype algorithm we present is a reverse method for calculating the entire third-order Tensor $y^T D^3 F(x)$. We want an expression for the derivative such that

$$y^T \frac{d}{dt} D^2 F(x + td) = y^T D^3 F(x + td) \cdot d \quad (16)$$

for any d . From equation (14), we see that d is contracted with the last coordinate in every term except one. To account for this term, we need a *switching tensor* S such that

$$y^T D^2 \Psi^2 \cdot (D^2 \Psi^1 \cdot d, D\Psi^1) = y^T D^2 \Psi^2 \cdot (D^2 \Psi^1, D\Psi^1) \cdot S \cdot d,$$

in other words we define S as

$$S \cdot (v, w, z) = (v, z, w) \quad \text{or} \quad S_{abcijk} v_i w_j z_k = v_a z_b w_c \quad (17)$$

for any vectors v, w and z . This implies that S 's components are $S_{abcijk} = \delta_{ai} \delta_{cj} \delta_{bk}$, where $\delta_{nm} = 1$ if $n = m$ and 0 otherwise. Then for $F = \Psi^2 \circ \Psi^1$ we use equation (14) to reach

$$\begin{aligned} y^T D^3 F \cdot d = & (y^T D^3 \Psi^2 \cdot (D\Psi^1, D\Psi^1, D\Psi^1) + y^T D\Psi^2 \cdot D^3 \Psi^1 \\ & + y^T D^2 \Psi^2 \cdot ((D\Psi^1, D^2 \Psi^1) + (D^2 \Psi^1, D\Psi^1) \cdot S + (D^2 \Psi^1, D\Psi^1))) \cdot d. \end{aligned} \quad (18)$$

The above is true for all vectors d , thus we can remove d from both sides to arrive at our desired expression for $y^T D^3 F$. With this notation we have, as expected, $(y^T D^3 F)_{ijk} = y^T D_{ijk} F$. We can now use this result to build a recurrence for

Algorithm 6: Archetype Reverse Third Order Derivative

initialization: $\bar{v} = y, W = 0 \in \mathbb{R}^{m_\ell \times m_\ell}, T \in \mathbb{R}^{m_\ell \times m_\ell \times m_\ell}$
for $i = \ell, \dots, 1$ **do**
 $T \leftarrow T \cdot (D\Psi^i, D\Psi^i, D\Psi^i)$
 $T \leftarrow T + W \cdot ((D\Psi^i, D^2\Psi^i) + (D^2\Psi^i, D\Psi^i))$
 $T \leftarrow T + W \cdot (D^2\Psi^i, D\Psi^i) \cdot S + \bar{v}^T D^3\Psi^i$
 $W \leftarrow W \cdot (D\Psi^i, D\Psi^i) + \bar{v}^T D^2\Psi^i$
 $\bar{v}^T \leftarrow \bar{v}^T D\Psi^i$
end
Output: $y^T D^3F(x) \cdot d \leftarrow T, y^T D^2F \leftarrow W, y^T DF \leftarrow \bar{v}^T$

D^3X^m , defined by $X^m = y^T \Psi^\ell \circ \dots \circ \Psi^m$, working from $m = \ell$ backwards towards $m = 1$ to calculate $y^T D^3F(x) \cdot d$.

Proof of Algorithm: the demonstration of this algorithm can be carried out in an analogous fashion to the proof of Algorithm 5.

This notation, together with a closed expression for high-order derivatives of a composition of two maps, see [5], can be used to design algorithms of even higher-orders. Though this would require the presentation of a rather cumbersome notation. What we can extract from this generic formula in [5], is that the number of terms that need to be calculated grows combinatorially in the order of the derivative, thus posing a lasting computational challenge.

3 Implementing through State Transformations

When coding a function, the user would not commonly write a composition of maps such as the form used in the previous section, see equation (1). Instead users implement functions in a number of different ways. Automatic Differentiation (AD) packages standardize these hand written functions, through compiler tools and operator overloading, into an evaluation that fits the format of Algorithm 7. As an example, consider the function $f(x, y, z) = xy \sin(z)$, and its evaluation for a given (x, y, z) through the following list of commands

$$\begin{aligned}
 v_{-2} &= x \\
 v_{-1} &= y \\
 v_0 &= z \\
 v_1 &= v_{-2}v_{-1} \\
 v_2 &= \sin(v_0) \\
 v_3 &= v_2v_1.
 \end{aligned}$$

By naming the functions $\phi_1(v_{-2}, v_{-1}) := v_{-2}v_{-1}$, $\phi_2(v_0) := \sin(v_0)$ and $\phi_3(v_2, v_1) := v_2v_1$, this evaluation fits the format in Algorithm 7.

In general, each ϕ_i is an *elemental function* such as addition, multiplication, $\sin(\cdot)$, $\exp(\cdot)$, etc, which together with their derivatives are already coded in the AD package. In order, the algorithm first copies the *independent* variables x_i into internal *intermediate* variables v_{i-n} , for $i = 1, \dots, n$. Following convention, we use negative indexes for elements that relate to independent variables. For

consistency, we will shift all indexes of vectors and matrices by $-n$ from here on, e.g., the components of $x \in \mathbb{R}^n$ are x_{i-n} for $i = 1 \dots n$.

The next step in Algorithm 7, calculates the value v_1 that only depends on the intermediate variables v_{i-n} , for $i = 1, \dots, n$. In turn, the value v_2 may now depend on v_{i-n} , for $i = 1, \dots, n+1$, then v_3 may depend on v_{i-n} , for $i = 1, \dots, n+2$ and so on for all ℓ intermediate variables. Each v_i is calculated using only one elemental function ϕ_i . There is a dependency amongst the intermediate variables, for ϕ_i is evaluated at previously calculated intermediate variables. We say that j is a predecessor of i if v_j is a necessary argument of ϕ_i . Let $P(i)$ be the set of predecessors of i and $v_{P(i)}$ the vector of predecessors, thus $\phi_i(v_{P(i)}) = v_i$ and necessarily $j < i$ for any $j \in P(i)$. Analogously, $S(i)$ is the set of successors of i .

Algorithm 7: Function evaluation

Input: $v_{i-n} = x_i$, for $i = 1, \dots, n$
for $i = 1 \dots \ell$ **do**
 $v_i \leftarrow \phi_i(v_{P(i)})$
end
Output: $f(x) \leftarrow v_\ell$

We can bridge this algorithmic description of a function with that of compositions of maps (1) using Griewank and Walther's [11] state-transformations

$$\begin{aligned} \Phi^i: \mathbb{R}^{n+\ell} &\rightarrow \mathbb{R}^{n+\ell}, \\ v &\mapsto (v_{1-n}, \dots, v_{i-1}, \phi_i(v_{P(i)}), v_{i+1}, \dots, v_\ell)^T, \end{aligned} \quad (19)$$

for $i = 1, \dots, \ell$. In components,

$$\Phi_r^i(v) = v_r(1 - \delta_{ri}) + \delta_{ri}\phi_i(v_{P(i)}), \quad (20)$$

where here, and in the remainder of this article, we abandon Einstein's notation of repeated indexes, because having the limits of summation is useful for implementing. With this, the function $f(x)$ defined by Algorithm 7 can be written as

$$f(x) = e_{\ell+n}^T \Phi^\ell \circ \Phi^{\ell-1} \circ \dots \circ \Phi^1 \circ (P^T x), \quad (21)$$

where $e_{\ell+n}$ is the $(\ell+n)$ th canonical vector and P is the immersion matrix $[I \ 0]$ with $I \in \mathbb{R}^{n \times n}$ and $0 \in \mathbb{R}^{n \times (\ell-n)}$. The Jacobian of the i th state transformation Φ^i , in coordinates, is simply

$$D_j \Phi_r^i(v) = \delta_{rj}(1 - \delta_{ri}) + \delta_{ri} \frac{\partial \phi_i}{\partial v_j}(v_{P(i)}). \quad (22)$$

With the state-transforms and the structure of their derivatives, we look again at a few of the archetype algorithms in Section 2 and build a corresponding implementable version. Our final goal is to implement the `RevHedir` algorithm 5, for which we need the implementation of the reverse gradient and Hessian algorithms.

3.1 First-Order

To design an algorithm to calculate the gradient of $f(x)$, given in equation (21), we turn to the Archetype Reverse Gradient Algorithm 1 and identify¹ the Φ^i 's in place of the Ψ^i 's. Using (22) we find that $\bar{v}^T \leftarrow \bar{v}^T D\Phi^i$ becomes

$$\bar{v}_j \leftarrow \bar{v}_j(1 - \delta_{ij}) + \bar{v}_i \frac{\partial \phi_i}{\partial v_j}(v_{P(i)}) \quad \forall j \in \{1 - n, \dots, \ell\} \quad (23)$$

where \bar{v}_i is the i -th component of \bar{v} , also known as the i -th *adjoint* in the AD literature. Note that if $j \neq i$ in the above, then the above step will only alter \bar{v}_j if $j \in P(i)$. Otherwise if $j = i$, then this update is equivalent to setting $\bar{v}_i = 0$. We can disregard this update, as \bar{v}_i will not be used in subsequent iterations. This is because $i \notin P(m)$, for $m \leq i$. With these considerations, we arrive at the Algorithm 8, the component-wise version of Algorithm 1. Note how we have used the abbreviated operation $a+ = b$ to mean $a \leftarrow a+b$. Furthermore, the last step $\bar{v}^T \leftarrow \bar{v}^T P^T$ selects the adjoints corresponding to independent variables.

An abuse of notation that we will employ throughout, is that whenever we refer to \bar{v}_i in the body of the text, we are referring to the value of \bar{v}_i after iteration i of the Reverse Gradient algorithm has finished.

Algorithm 8: Reverse Gradient.

initialization: $\bar{v} = e_1 \in \mathbb{R}^{\ell+n}$
for $i = \ell, \dots, 1$ **do**
 for $j \in P(i)$ **do** $\bar{v}_j + = \bar{v}_i \partial \phi_i(v_{P(i)}) / \partial v_j$
end
Output: $\nabla f \leftarrow \bar{v}^T P^T = (\bar{v}_{1-n}, \dots, \bar{v}_0)^T$

Similarly, by using (22) again, each iteration i of the Archetype 1st Order Directional Derivative Algorithm 2, can be reduced to a coordinate form

$$\dot{v}_r \leftarrow (1 - \delta_{ri})\dot{v}_r + \delta_{ri} \sum_{j \in P(i)} \dot{v}_j \frac{\partial \phi_i}{\partial v_j}(v_{P(i)}),$$

where \dot{v}_j is the j -th component of \dot{v} . If $r \neq i$ in the above, then \dot{v}_r remains unchanged, while if $r = i$ then we have

$$\dot{v}_i \leftarrow \sum_{j \in P(i)} \dot{v}_j \frac{\partial \phi_i}{\partial v_j}(v_{P(i)}). \quad (24)$$

We implement this update by sweeping through the successors of each intermediate variable and incrementing a single term to the sum on the right-hand side of (24), see Algorithm 9. It is crucial to observe that the i -th component of \dot{v} will remain unaltered after the i -th iteration.

Again, when we refer to \dot{v}_i in the body of the text from this point on, we are referring to the value of \dot{v}_i after iteration i has finished in Algorithm 9.

¹Especially P^T would be Ψ^1 and Φ^i would be Ψ^{i+1} .

Algorithm 9: 1st Order Directional Derivative.

initialization: $\dot{v} = P^T d \in \mathbb{R}^{\ell+n}$
for $j = 1, \dots, \ell$ **do**
 for $i \in S(j)$ **do** $\dot{v}_i += \dot{v}_j \partial \phi_i(v_{P(i)}) / \partial v_j$
end
Output: $DF \cdot d = (\dot{v}_{1-n}, \dots, \dot{v}_0)^T$

3.2 Second-Order

Just by substituting Ψ^i s for Φ^i s in the Archetype Reverse Hessian, Algorithm 3, we can quickly reach a very efficient component-wise algorithm for calculating the Hessian of $f(x)$, given in equation (21). This component-wise algorithm is also known as `edge_pushing`, and has already been detailed in Gower and Mello [8]. Here we use a different notation which leads to a more concise presentation. Furthermore, the results below form part of the calculations needed for third order methods.

There are two steps of Algorithm 3 we must investigate, for we already know how to update \bar{v} from the above section. For these two steps, we need to substitute

$$D_{jk} \Phi_r^i(v) = \frac{\partial^2 \Phi_r^i}{\partial v_j \partial v_k}(v) = \delta_{ri} \frac{\partial^2 \phi_i}{\partial v_j \partial v_k}(v_{P(i)}), \quad (25)$$

and $D\Phi^i$, equation (22), in $W \leftarrow W \cdot (D\Phi^i, D\Phi^i) + \bar{v}^T D^2 \Phi^i$, resulting in

$$\begin{aligned} W_{jk} &\leftarrow \sum_{s,t=1-n}^{\ell} \frac{\partial \Phi_s^i}{\partial v_j} W_{st} \frac{\partial \Phi_t^i}{\partial v_k} + \sum_{s=1-n}^{\ell} \bar{v}_s \frac{\partial^2 \Phi_s^i}{\partial v_j \partial v_k} \\ &= (1 - \delta_{ji}) W_{jk} (1 - \delta_{ki}) + \frac{\partial \phi_i}{\partial v_j} W_{ii} \frac{\partial \phi_i}{\partial v_k} \\ &\quad + \frac{\partial \phi_i}{\partial v_j} W_{ik} (1 - \delta_{ki}) + (1 - \delta_{ji}) W_{ji} \frac{\partial \phi_i}{\partial v_k} \end{aligned} \quad (26)$$

$$+ \bar{v}_i \frac{\partial^2 \phi_i}{\partial v_j \partial v_k}. \quad (27)$$

Before translating these updates into an algorithm, we need a crucial result: at the beginning of iteration $i - 1$, the element W_{jk} is zero if $j \geq i$ for all k . We show this by using induction on the iterations of Algorithm 3. Note that W is initially set to zero, so the first step from (26) and (27) reduce to

$$W_{jk} \leftarrow \bar{v}_\ell \frac{\partial^2 \phi_\ell}{\partial v_j \partial v_k},$$

which is zero for $j = \ell$ because $\ell \notin P(\ell)$. Now we assume the induction hypothesis holds at the beginning of the iteration i , so that $W_{jk} = 0$ for $j \geq i + 1$. So letting $j \geq i + 1$ and executing the iteration i we get from the updates (26) and (27)

$$W_{jk} \leftarrow W_{jk} + W_{ji} \frac{\partial \phi_i}{\partial v_k},$$

because $j \notin P(i)$. Together with our hypothesis $W_{jk} = 0$ and $W_{ji} = 0$, we see that W_{jk} remains zero. While if $j = i$, then (26) and (27) sets $W_{jk} \leftarrow 0$ because

$i \notin P(i)$. Hence at the beginning of iteration $i - 1$ we have that $W_{jk} = 0$ for $j \geq i$ and this completes the induction.

Furthermore, W is symmetric at the beginning of iteration i because it is initialized to $W = 0$ and each iteration preserves symmetry. Consequentially, the only nonzero components W_{jk} appear when both $j, k \leq i$. We make use of this symmetry to avoid unnecessary calculations on symmetric counterparts. Let $W_{\{jk\}}$ denote both W_{jk} and W_{kj} . To accommodate for this symmetric representation, we perform (27) once for each pair $\{j, k\}$, as to opposed for every coordinate pair, see the **Creating** step in Algorithm 10.

The calculations in (26) are done by sweeping through the nonzero elements of W and then updating their contribution to the overall calculation.

Thus if $W_{\{ii\}} \neq 0$, looking to (26), this triggers the following increment

$$W_{\{jk\}}+ = \frac{\partial \phi_i}{\partial v_j} W_{\{ii\}} \frac{\partial \phi_i}{\partial v_k}.$$

Similarly to **Creating** step, the above should only be carried out for every pair $\{j, k\}$. While each nonzero off diagonal term W_{ik} and W_{ki} , for $k < i$, according to (26), has the effect of

$$W_{jk}+ = \frac{\partial \phi_i}{\partial v_j} W_{ik}, \quad (28)$$

$$W_{kj}+ = \frac{\partial \phi_i}{\partial v_j} W_{ki}. \quad (29)$$

It is redundant to update both symmetric elements, so we substitute both for just

$$W_{\{jk\}}+ = \frac{\partial \phi_i}{\partial v_j} W_{\{ik\}}.$$

Though we must take care when $j = k$, for according to (28) and (29), the two symmetric updates “double up” on the diagonal

$$W_{\{jj\}}+ = 2 \frac{\partial \phi_i}{\partial v_j} W_{\{ij\}}. \quad (30)$$

The operation (26) has been implemented with these above considerations in the **Pushing** step in Algorithm 10. The names of the steps **Creating** and **Pushing** are elusive to a graph interpretation [8].

3.3 Third-Order

The final algorithm that we translate to implementation is the Hessian directional derivative, the **RevHedir** Algorithm 5. This implementation has an immediate application in the Halley-Chebyshev class of third-order optimization methods, for at each step of these algorithms, such a directional derivative is required. For this reason we pay special attention to its implementation.

Identifying each Ψ^i with Φ^i , we address each of the five operations on the matrix Td in Algorithm 5 separately, pointing out how each one preserves the symmetry of Td and how to perform the component-wise calculations.

First, given that Td is symmetric, the *2D pushing* update

$$Td \leftarrow Td \cdot (D\Phi^i, D\Phi^i), \quad (31)$$

Algorithm 10: component-wise form of `edge_pushing`.

Input: Function evaluation \bar{v} , $x \in \mathbb{R}^n$.
initialization: $\bar{v} = e_{\ell+n} \in \mathbb{R}^{\ell+n}$, $W = 0 \in \mathbb{R}^{(\ell+n) \times (\ell+n)}$
for $i = \ell, \dots, 1$ **do**
 Pushing
 foreach $k \leq i$ such that $W_{\{ki\}} \neq 0$ **do**
 if $k < i$ **then**
 foreach $j \in P(i)$ **do**
 if $j = k$ **then**
 $W_{\{jj\}} += 2D_j \phi_i W_{\{ji\}}$
 else
 $W_{\{jk\}} += D_j \phi_i W_{\{ki\}}$
 end
 end
 else $k = i$
 foreach *unordered pair* $\{j, p\} \subset P(i)$ **do**
 $W_{\{jp\}} += D_p \phi_i D_j \phi_i W_{\{ii\}}$
 end
 end
 Creating
 foreach *unordered pair* $\{j, p\} \subset P(i)$ **do**
 $W_{\{jp\}} += \bar{v}_i D_{pj} \phi_i$
 end
 Adjoint
 foreach $j \in P(i)$ **do**
 $\bar{v}_j += \bar{v}_i D_j \phi_i$
 end
end
Output: $D^2 f = (W_{jk})_{1-n \leq j, k \leq 0}$

is exactly as detailed in (26) and the surrounding comments. While *3D creating*

$$Td \leftarrow Td + \bar{v}^T D^3 \Phi^i \cdot \dot{v}^{i-1},$$

can be written in coordinate form as

$$\begin{aligned} Td_{jk} &\leftarrow Td_{jk} + \sum_{r,p=1-n}^{\ell} \bar{v}_r D_{jkp} \Phi_r^i \dot{v}_p^{i-1} \\ &= Td_{jk} + \sum_{p \in P(i)} \bar{v}_i \frac{\partial^3 \phi_i}{\partial v_j \partial v_k \partial v_p} \dot{v}_p, \end{aligned} \quad (32)$$

where \dot{v}_p is the value given to \dot{v}_p after iteration p in Algorithm 9. Note that $\dot{v}_p^{i-1} = \dot{v}_p$ for $p \in P(i)$, because $p \leq i-1$, so on the iteration $i-1$ of Algorithm 9 the calculation of \dot{v}_p will already have been finalized. Another trick we employ is that, since the above calculation is performed on iteration i , we know that \bar{v}_i has already been calculated. These substitutions involving \bar{v}_i s and \dot{v}_i s will be carried out in the rest of the text with little or no comment. The update (32) also preserves the symmetry of Td .

To examine the update,

$$Td \leftarrow Td + W \cdot (D\Phi^i, D^2\Phi^i \cdot \dot{v}^{i-1}), \quad (33)$$

we use (22) and (25) to obtain the coordinate form

$$\begin{aligned} Td_{jk} &\leftarrow Td_{jk} + \sum_{r,s=1-n}^{\ell} W_{rs} \left(\delta_{rj}(1 - \delta_{ri}) + \delta_{ri} \frac{\partial \phi_i}{\partial v_j} \right) \delta_{si} \frac{\partial^2 \phi_i}{\partial v_k \partial v_p} \dot{v}_p \\ &= Td_{jk} + W_{ji}(1 - \delta_{ji}) \frac{\partial^2 \phi_i}{\partial v_k \partial v_p} \dot{v}_p + W_{ii} \frac{\partial \phi_i}{\partial v_j} \frac{\partial^2 \phi_i}{\partial v_k \partial v_p} \dot{v}_p. \end{aligned} \quad (34)$$

Upon inspection, the update

$$Td \leftarrow Td + W \cdot (D^2\Phi^i \cdot \dot{v}^{i-1}, D\Phi^i)$$

is the transpose of (34) due to the symmetry of W . So it can be written in coordinate form as

$$Td_{jk} \leftarrow Td_{jk} + W_{ik}(1 - \delta_{ki}) \frac{\partial^2 \phi_i}{\partial v_j \partial v_p} \dot{v}_p + W_{ii} \frac{\partial \phi_i}{\partial v_k} \frac{\partial^2 \phi_i}{\partial v_j \partial v_p} \dot{v}_p. \quad (35)$$

Thus update (35) together with (34) is equivalent to summing a symmetric matrix to Td , so the symmetry of Td is still preserved.

Last we translate

$$Td \leftarrow Td + W \cdot (D^2\Phi^i, D\Phi^i \cdot \dot{v}^{i-1}), \quad (36)$$

to its coordinate form

$$\begin{aligned} Td_{jk} &\leftarrow Td_{jk} + \sum_{r,s=1-n}^{\ell} W_{rs} \delta_{ri} D_{jk} \Phi_r^i \left(\delta_{sp}(1 - \delta_{si}) + \delta_{si} \frac{\partial \phi_i}{\partial v_p} \right) \dot{v}_p \\ &= Td_{jk} + W_{ip} \frac{\partial^2 \phi_i}{\partial v_j \partial v_k} (1 - \delta_{pi}) \dot{v}_p + W_{ii} \frac{\partial^2 \phi_i}{\partial v_j \partial v_k} D_p \phi_i \dot{v}_p. \end{aligned} \quad (37)$$

No change is affected by interchanging the indices j and k on the right-hand side of (37), so once again Td remains symmetric. For convenience of computing, we group updates (34), (35) and (37) into a set of updates called **2D Connecting**. The name indicating that these updates “connect” objects that contain second order derivative information.

More than just symmetric, through closer inspection of these operations, we see that the sparsity structure of Td is contained in that of W . This remains true even after execution, at which point $Td = D^3f(x) \cdot d$ and $W = D^2f(x)$ where, for each $j, k, p \in \{1 - n, \dots, 0\}$, we have

$$D_{jk}f(x) = 0 \implies D_{jkp}f(x)d_p = 0.$$

This fact should be explored when implementing the method, in that, the data structure of Td should imitate that of W .

3.3.1 Implementing Third-Order Directional Derivative

The matrices Td and W are symmetric, and based on the assumption that they will be sparse, we will represent them using a symmetric sparse data structure. Thus we now identify each pair (W_{jk}, W_{kj}) and (Td_{jk}, Td_{kj}) with the element $W_{\{jk\}}$ and $Td_{\{jk\}}$, respectively. Much like was done with `edge_pushing`, Algorithm 10, we will organize the computations by sweeping through all nonzero elements of $Td_{\{ik\}}$ and $W_{\{ik\}}$ and then updating their contribution to the overall calculation.

We must take care when updating our symmetric representation of Td , both for the 2D pushing update (31) and for the redundant symmetric counterparts (34) and (35) which “double-up” on the diagonal, much like in the `Pushing` operations of Algorithm 10.

Each operation (34), (35) and (37) depends on a diagonal element $W_{\{ii\}}$ and an off-diagonal element $W_{\{ik\}}$ of W , for $k \neq i$. Grouping together all terms that involve $W_{\{ii\}}$ we get the resulting update

$$Td_{\{jk\}}+ = W_{\{ii\}} \sum_{p \in P(i)} \dot{v}_p \left(\frac{\partial \phi_i}{\partial v_j} \frac{\partial^2 \phi_i}{\partial v_k \partial v_p} + \frac{\partial \phi_i}{\partial v_k} \frac{\partial^2 \phi_i}{\partial v_j \partial v_p} + \frac{\partial \phi_i}{\partial v_p} \frac{\partial^2 \phi_i}{\partial v_j \partial v_k} \right). \quad (38)$$

By appropriately renaming the indices in (34), (35) and (37), each nonzero off diagonal element $W_{\{ik\}}$ gives the updates (39), (40) and (41), respectively.

$$Td_{jk}+ = \sum_{p \in P(i)} \dot{v}_p \frac{\partial^2 \phi_i}{\partial v_j \partial v_p} W_{ik}, \quad \forall j \in P(i) \quad (39)$$

$$Td_{kj}+ = \sum_{p \in P(i)} \dot{v}_p \frac{\partial^2 \phi_i}{\partial v_j \partial v_p} W_{ki}, \quad \forall j \in P(i) \quad (40)$$

$$Td_{jp}+ = \sum_{p \in P(i)} \dot{v}_k \frac{\partial^2 \phi_i}{\partial v_j \partial v_p} W_{ik}, \quad \forall j \in P(i) \quad (41)$$

Note that (39) and (40) are symmetric updates, and when $j = k$ these two operations “double-up” resulting in the update

$$Td_{jj}+ = 2 \sum_{p \in P(i)} \dot{v}_p \frac{\partial^2 \phi_i}{\partial v_j \partial v_p} W_{ij}.$$

Passing to our symmetric notation, we dispense with (40) and account for this doubling effect in Algorithm 11. Finally we can eliminate redundant symmetric calculations performed in (41) by only performing this operation for each pair $\{j, p\}$. All these considerations relating to 2D **connecting** have been factored into our implementation of the **RevHedir** Algorithm 11.

Performing 3D **Creating** (32) using this symmetric representation is simply a matter of not repeating the obvious symmetric counterpart, but instead, performing these operations on $Td_{\{jk\}}$ once for each appropriate pair $\{j, k\}$, see 3D **Creating** in to Algorithm 11.

Algorithm 11: component-wise form of **RevHedir**.

Input: Function evaluation \bar{v} , $x \in \mathbb{R}^n$.
Initialization: $\bar{v}_{1-n} = \dots = \bar{v}_{\ell-1} = 0$, $\bar{v}_\ell = 1$, $W_{jk} = 0$, $Td_{\{jk\}} = 0$,
 $j < k \in \{1-n, \dots, \ell\}$
Calculate first order directional derivative \dot{v} using Algorithm (9)
for $i = \ell, \dots, 1$ **do**
 2D **Pushing** of Td , see **Pushing** in Algorithm 10
 2D **Connecting**
 foreach $p \in P(i)$, $\{j, k\} \subset P(i)$ **do**
 $Td_{\{jk\}} += W_{\{ii\}} \dot{v}_p (D_j \phi_i D_{kp} \phi_i + D_k \phi_i D_{jp} \phi_i + D_p \phi_i D_{jk} \phi_i)$
 end
 foreach $k < i$, $W_{\{ik\}} \neq 0$ **do**
 foreach $(j, p) \in P(i)^2$ **do**
 if $j = k$ **then**
 $Td_{\{kk\}} += 2W_{\{ik\}} \dot{v}_p D_{jp} \phi_i$
 end
 if $j \neq k$ **then**
 $Td_{\{jk\}} += W_{\{ik\}} \dot{v}_p D_{jp} \phi_i$
 end
 if $j \geq p$ **then**
 $Td_{\{jp\}} += W_{\{ik\}} \dot{v}_k D_{jp} \phi_i$
 end
 end
 end
 3D **Creating**
 foreach $p \in P(i)$, $\{j, k\} \subset P(i)$ **do**
 $Td_{\{jk\}} += \bar{v}_i D_{jkp} \phi_i \dot{v}_p$
 end
 Pushing and creating applied to W , see Algorithm 10
 Adjoint Iteration applied to \bar{v} , see Algorithm 8
end
Output: $(D^3 f(x) \cdot d)_{jk} = Td_{\{jk\}}$, $D^2 f(x)_{jk} = W_{\{jk\}}$
for each $j \leq k \in \{1-n, \dots, 0\}$.

4 Numerical experiment

We have implemented the `RevHedir` Algorithm 11 as an additional driver of ADOL-C, a well established automatic differentiation library coded in C and C++ [9]. We used version ADOL-C-2.4.0, the most recent available ². The tests were carried out on a personal laptop with 1.70GHz dual core processors Intel Core i5-3317U, 4GB of RAM, with the Ubuntu 13.0 operating system.

For those interested in replicating our implementation, we used a sparse undirected weighted graph data structure to represent the matrices W and Td . The data structure is an array of weighted neighbourhood sets, one for each node, where each neighbourhood set is a dynamic array that resizes when needed. Each neighbourhood set is maintained in order and the method used to insert or increment the weight of an edge is built around a binary search.

We have hand-picked sixteen problems from the CUTE collection [2], `augmagn` from [15], `toiqmerg` (Toint Quadratic Merging problem) and `chainros_trigexp` (Chained Rosenbrock function with Trigonometric and exponential constraints) from [17] for the experiments. We have also created a function

$$\text{heavey_band}(x, \text{band}) = \sum_{i=1}^{n-\text{band}} \sin \left(\sum_{j=1}^{\text{band}} x_{i+j} \right).$$

For our experiments, we tested `heavey_band(x, 20)`. The problems were selected based on the sparsity pattern of $D^3 f(x) \cdot d$, dimension scalability and sparsity. Our goal was to cover a variety of patterns, to easily change the dimension of the function and work with sparse matrices.

In Table 1, the ‘‘Pattern’’ column indicates the type of sparsity pattern: bandwidth³ of value x (B x), arrow, frame, number of diagonals (D x), or irregular pattern. The ‘‘nnz/n’’ column gives the number of nonzeros in $D^3 f(x) \cdot d$ over the dimension n , which serves as a measure of density. For each problem, we applied `RevHedir` and `edge_pushing` Algorithm 11 and 10 to the objective function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, with $x_i = i$ and $d_i = 1$, for $i = 1, \dots, n$, and give the runtime of each method for dimension $n = 10^6$ in Table 1. Note that all of these matrices are very sparse, partly due to the ‘‘thinning out’’ caused by the high order differentiation. This probably contributed to the relatively low runtime, for in these tests, the run-times have a 0.75 correlation with the density measure ‘‘nnz/n’’. This leads us to believe that the actual pattern is not a decisive factor in runtime.

We did not benchmark our results against an alternative algorithm for we could not find a known AD package that is capable of efficiently calculating such directional derivatives for such high dimensions. For small dimensions, we used the `tensor_eval` of ADOL-C to calculate the entire tensor using univariate forward Taylor series propagation[12]. Then we contract the resulting tensor with the vector d . This was useful to check that our implementation was correct, though it would struggle with dimensions over $n = 100$, thus not an appropriate comparison.

A remarkable feature of these tests, is that the time spent by `RevHedir` to calculate $D^3 f(x) \cdot d$ was, on average, 108% that of calculating $D^2 f(x)$. Thus,

²As checked May 28th, 2013

³The bandwidth of matrix $M = (m_{ij})$ is the maximum value of $2|i - j| + 1$ such that $m_{ij} \neq 0$.

name	Pattern	nnz/n	edge_pushing	RevHedir
cosine	B 3	3.0000	2.89	5.25
bc4	B 3	3.0000	3.93	7.87
cragglevy	B 3	2.9981	5.41	10.6
chainwood	B 3	1.4999	4.04	7.22
morebv	B 3	3.0000	4.57	9.44
scon1dls	B 3	0.7002	3.99	8.12
bdexp	B 5	0.0004	2.21	3.86
pspdoc	B 5	4.9999	3.05	5.97
augmlagn	5 × 5 diagonal blocks	4.9998	4.15	9.28
brybnd	B 11	12.9996	14.19	38.79
chainros_trigexp	B 3 + D 6	4.4999	6.51	12.87
toiqmerg	B 7	6.9998	4.33	8.89
arwhead	arrow	3.0000	3.63	6.78
nondquar	arrow + B 3	4.9999	2.9	5.61
sinquad	frame + diagonal	4.9999	5.12	10.01
bdqrtic	arrow + B 7	8.9998	8.98	19.62
noncvxu2	irregular	6.9998	4.95	9.55
ncvxqp3	irregular	6.9997	2.9	6.48
heavy_band	B 39	38.9995	20.74	61.27

Table 1: Description of problem set together with the execution time in seconds of `edge_push` and `RevHedir` for $n = 10^6$.

if the user is prepared to pay the price for calculating the Hessian, he could also gain some third order information for approximately the same cost. The code for these tests can be downloaded from the Edinburgh Research Group in Optimization website: <http://www.maths.ed.ac.uk/ERGO/>.

5 Conclusion

Our contribution boils down to a framework for designing high order reverse methods, and an efficient implementation of the directional derivative of the Hessian called `RevHedir`. The framework paves the way to obtaining a reverse method for all orders once and for all. Such an achievement could cause a paradigm shift in numerical method design, wherein, instead of increasing the number of steps or the mesh size, increasing the order of local approximations becomes conceivable. We have also shed light on existing AD methods, providing a concise proof of the `edge_pushing` [8] and the reverse gradient directional derivative [1] algorithms.

The novel algorithms 5 and 6 for calculating the third-order derivative and its contraction with a vector, respectively, fulfils what we set out to achieve: they accumulate the desired derivative “as a whole”, thus taking advantage of overlapping calculations amongst individual components. This is in contrast with what is currently being used, e.g., univariate Taylor expansions [12] and repeated tangent/adjoint operations [18]. These algorithms can also make use of the symmetry, as illustrated in our implementation of `RevHedir` Algorithm 11, wherein all operations are only carried out on a lower triangular matrix.

We implemented and tested the `RevHedir` with two noteworthy results. The first is its capacity to solve sparse problems of large dimension of up to a million variables. The second is how the time spent by `RevHedir` to calculate the directional derivative $D^3 f(x) \cdot d$ was very similar to that spent by `edge_pushing` to calculate the Hessian. We believe this is true in general and plan on confirming

this in future work through complexity analysis. Should this be confirmed, it would have an immediate consequence in the context of nonlinear optimization, in that the third-order Halley-Chebyshev methods could be used to solve large dimensional problems with an iteration cost proportional to that of Newton step. In more detail, at each step the Halley-Chebyshev methods require the Hessian matrix and its directional derivative. The descent direction is then calculated by solving the Newton system, and an additional system with the same sparsity pattern as the Newton system. If it is confirmed that solving these systems costs the same, in terms of complexity, then the cost of a Halley-Chebyshev iteration will be proportional to that of a Newton step. Though this comparison only holds if one uses these automatic differentiation procedures to calculate the derivatives in both methods.

The CUTE functions used to test both `edge_pushing` and `RevHedir` are rather limited, and further tests on real-world problems should be carried out. Also, complexity bounds need to be developed for both algorithms.

A current limitation of reverse AD procedures, such as the ones we have presented, is their issue with memory usage. All floating point values of the intermediate variables must be recorded on a forward sweep and kept for use in the reverse sweep. This can be a very substantial amount of memory, and can be prohibitive for large-scale functions [21]. As an example, when we used dimensions of $n = 10^7$, most of our above test cases exhausted the available memory on the personal laptop used. A possible solution to this, is to allow a trade off between run-time and memory usage by reversing only parts of the procedure at a time. This method is called checkpointing [21, 20].

References

- [1] J. ABATE, C. BISCHOF, L. ROH, AND A. CARLE, *Algorithms and design for a second-order automatic differentiation module*, in Proceedings of the 1997 International Symposium on Symbolic and Algebraic Computation, New York, 1997, ACM, pp. 149–155.
- [2] I. BONGARTZ, A. R. CONN, N. GOULD, AND P. L. TOINT, *CUTE: constrained and unconstrained testing environment*, ACM Transactions on Mathematical Software, 21 (1995), pp. 123–160.
- [3] B. CHRISTIANSON, *Automatic Hessians by reverse accumulation*, IMA Journal of Numerical Analysis, 12 (1992), pp. 135–150.
- [4] G. F. CORLISS, A. GRIEWANK, AND P. HENNEBERGER, *High-order stiff ODE solvers via automatic differentiation and rational prediction*, in Lecture Notes in Computer Science, Springer, 1997, pp. 114–125.
- [5] L. E. FRAENKEL, *Formulae for high derivatives of composite functions*, Mathematical Proceedings of the Cambridge Philosophical Society, 83 (2008), p. 159.
- [6] A. H. GEBREMEDHIN, F. MANNE, AND A. POTHEN, *What color is your Jacobian? Graph coloring for computing derivatives*, SIAM Review, 47 (2005), pp. 629–705.

- [7] A. H. GEBREMEDHIN, A. TARAFDAR, A. POTHEN, AND A. WALTHER, *Efficient computation of sparse Hessians using coloring and automatic differentiation*, *INFORMS Journal on Computing*, 21 (2009), pp. 209–223.
- [8] R. M. GOWER AND M. P. MELLO, *A new framework for the computation of Hessians*, *Optimization Methods and Software*, 27 (2012), pp. 251–273.
- [9] A. GRIEWANK, D. JUEDES, AND J. UTKE, *ADOL-C, a package for the automatic differentiation of algorithms written in C/C++*, *ACM Transactions on Mathematical Software*, 22 (1996), pp. 131–167.
- [10] A. GRIEWANK AND U. NAUMANN, *Accumulating Jacobians as chained sparse matrix products*, *Mathematical Programming*, 95 (2003), pp. 555–571.
- [11] A. GRIEWANK AND A. WALTHER, *Evaluating derivatives*, Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, second ed. ed., 2008.
- [12] A. GRIEWANK, A. WALTHER, AND J. UTKE, *Evaluating higher derivative tensors by forward propagation of univariate Taylor series*, *Mathematics of Computation*, 69 (2000), pp. 1117–1130.
- [13] J. GUCKENHEIMER AND B. MELOON, *Computing periodic orbits and their bifurcations with automatic differentiation*, *SIAM Journal on Scientific Computing*, 22 (2000), pp. 951–985.
- [14] J. GUTIÉRREZ AND M. HERNÁNDEZ, *A family of Chebyshev-Halley type methods in Banach spaces*, *Bulletin of the Australian Mathematical Society*, 55 (1997), p. 113.
- [15] W. HOCK AND K. SCHITTKOWSKI, *Test examples for nonlinear programming codes*, *Journal of Optimization Theory and Applications*, 30 (1980), pp. 127–129.
- [16] V. KARIWALA, *Automatic Differentiation-Based Quadrature Method of Moments for Solving Population Balance Equations*, *AIChE Journal*, 58 (2012), pp. 842–854.
- [17] LADISLAV LUKSAN JAN VLCEK, *Test problems for unconstrained optimization*, Tech. Report 897, Academy of Sciences of the Czech Republic, 2003.
- [18] U. NAUMANN, *The Art of Differentiating Computer Programs: An Introduction to Algorithmic Differentiation*, no. 24 in *Software, Environments, and Tools*, SIAM, Philadelphia, PA, 2012.
- [19] R. D. NEIDINGER, *An Efficient Method for the Numerical Evaluation of Partial Derivatives of Arbitrary Order*, *ACM Transactions on Mathematical Software*, 18 (1992), pp. 159–173.
- [20] J. STERNBERG AND A. GRIEWANK, *Reduction of Storage Requirement by Checkpointing for Time-Dependent Optimal Control Problems in ODEs*, in *Automatic Differentiation: Applications, Theory, and Implementations*, B. N. M. Bücker, G. Corliss, P. Hovland, U. Naumann, ed., no. 0, Springer, 1 ed., 2006, pp. 99–110.

- [21] A. WALTHER AND A. GRIEWANK, *Advantages of Binomial Checkpointing for Memory-reduced Adjoint Calculations*, Numerical Mathematics and Advanced Applications, (2004), pp. 834–843.