



Randomized Heuristics for Exploiting Jacobian Scarcity[☆]

Andrew Lyons^{1,2,*}, Ilya Safro², Jean Utke^{1,2}

Abstract

We describe a code transformation technique that, given code for a vector function F , produces code suitable for computing collections of Jacobian-vector products $F'(\mathbf{x})\mathbf{x}$ or Jacobian-transpose-vector products $F'(\mathbf{x})^T\bar{\mathbf{y}}$. Exploitation of scarcity — a measure of the degrees of freedom in the Jacobian matrix — means solving a combinatorial optimization problem that is believed to be hard. Our heuristics transform the computational graph for F , producing, in the form of a transformed graph G' , a representation of the Jacobian $F'(\mathbf{x})$ that is both concise and suitable for the evaluation of large collections of Jacobian-vector products or Jacobian-transpose-vector products. Our heuristics are randomized in nature and compare favorably in all cases with the best known heuristics.

Keywords: Automatic Differentiation, Scarcity

1. Introduction

The computation of Jacobian-vector products is a fundamental step in the context of science and engineering applications. Without loss of generality, Suppose a vector function $F : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is given as a straight-line evaluation procedure; real-life application codes often comprise such straight-line procedures. Thus, F may not represent the entire function of interest, but rather a small part that is executed many times. We are interested in algorithmically applying the chain rule, a technique known as automatic differentiation (AD), in order to obtain a new program that evaluates F along with some derivative information for F . Suppose, in particular, that we are interested in computing either a collection of p Jacobian-vector products $(F'(\mathbf{x})\mathbf{x}^i)_{i=1,\dots,p}$, or a

[☆]draft version hg:617fc2231dfd:142 compiled January 27, 2010

*Corresponding author

Email addresses: lyonsam@gmail.com (Andrew Lyons), safro@mcs.anl.gov (Ilya Safro), utke@mcs.anl.gov (Jean Utke)

URL: <http://www.mcs.anl.gov/~lyonsam/> (Andrew Lyons), <http://www.mcs.anl.gov/~safro/> (Ilya Safro), <http://www.mcs.anl.gov/~utke/> (Jean Utke)

¹Computation Institute, University of Chicago

²Mathematics and Computer Science Division, Argonne National Laboratory

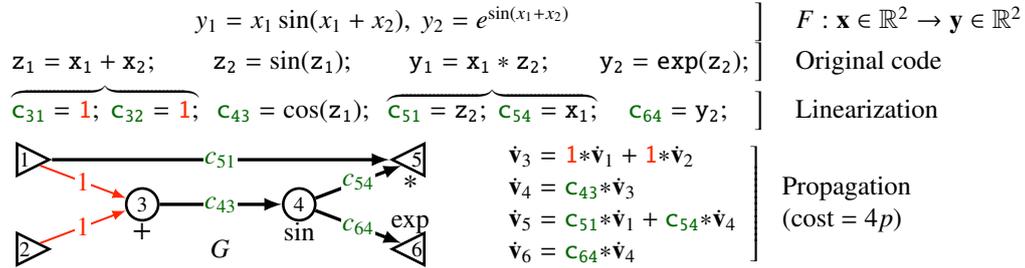


Figure 1: A function F (line 1) is given as a straight-line program (line 2). The edges of the computational graph G represent the direct dependencies among the program variables (which correspond to the vertices). The process of *linearization* (line 3), fundamental to AD, automatically produces code for evaluating the local partial derivatives $c_{ji} \equiv \partial v_j / \partial v_i$ at a small fixed cost.

collection of p Jacobian-transpose-vector products $(F'(\mathbf{x})^T \bar{\mathbf{y}}^i)_{i=1, \dots, p}$, where p is assumed to be sufficiently large. The vectors $\bar{\mathbf{x}}^i$ are directions in the domain; the vectors $\bar{\mathbf{y}}^i$ may be interpreted as weights. Due to symmetry, we are able to restrict our attention to the former without loss of generality. Our goal, therefore, will be to approximate the most efficient program for computing collections of Jacobian-vector products. The notion of Jacobian scarcity [1–3] generalizes the properties of sparsity and rank to capture a deficiency in the degrees of freedom of the Jacobian matrix. We describe new randomized heuristics that exploit scarcity for the optimized evaluation of collections of Jacobian-vector or Jacobian-transpose-vector products.

In the remainder of this section, we introduce the necessary definitions and concepts. In Section 3, we describe our heuristics. In Section 2, we provide experimental results.

Propagating derivatives. The basic concepts of AD are illustrated in Figure 1. Henceforth, we will assume that the linearized computational graph G is given, where every edge (i, j) is associated with either a unique variable representing the local partial derivative c_{ji} or a value in $\{1, -1\}$. Note that G is a directed acyclic graph (DAG). Traditional AD [3] prescribes the *forward mode* for evaluating Jacobian-vector products, whereby derivative values \dot{v}_j are *propagated* through G from the sources to the sinks by traversing the vertices in topological order. When a vertex j is visited, we compute $\dot{v}_j = \sum_{i \in P_j} c_{ji} * \dot{v}_i$, where P_j denotes the set of vertices that have outedges to j . $F'(\mathbf{x})\bar{X}$ can be evaluated either by propagating p directions $\bar{\mathbf{x}}^1, \dots, \bar{\mathbf{x}}^p$ separately (scalar mode), or by propagating $\bar{X} \in \mathbb{R}^{p \times n}$ in a single pass (vector mode). In terms of scalar multiplications, the total computational cost associated with each edge (i, j) is p if $c_{ji} \notin \{1, -1\}$ and 0 otherwise; this highlights the special significance of unit edges. Thus, in both scalar and vector modes, the cost of evaluating p Jacobian-vector products is $p|E^+(G)|$ scalar multiplications, where $E^+(G) \equiv \{(i, j) \in E \mid c_{ji} \notin \{1, -1\}\}$ denotes the set of *nonunit edges* in G .

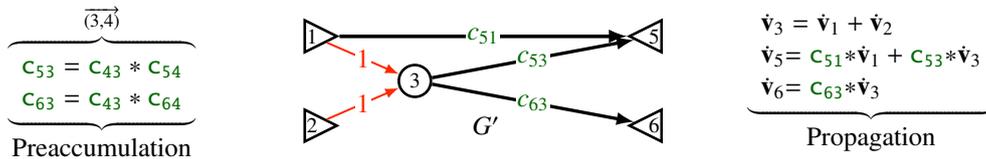


Figure 2: Partially preaccumulated Jacobian.

Graph transformations and preaccumulation. The essential property of G is that the entries of the Jacobian $F'(\mathbf{x})$ can be expressed as *Baur's formula*

$$\frac{\partial y_j}{\partial x_i} = \sum_{P \in [x_i \rightsquigarrow y_j]} \prod_{(k, \ell) \in P} c_{\ell k} \quad ,$$

where $[x_i \rightsquigarrow y_j]$ denotes the set of all paths from x_i to y_j in G . A sequence ρ of local graph transformations called *edge eliminations* allow us to transform G into the *remainder graph* $G'(\rho)$, which retains this property. Note that we do not consider other known types of transformations (normalizations, reroutings, etc. [4]) because of the caveats associated with them.

Front elimination of an edge (i, j) , denoted $\overrightarrow{(i, j)}$, entails updating $c_{\ell i+} = c_{ji} * c_{\ell j}$ for all successors ℓ of j . Similarly, *back elimination* of an edge (i, j) , denoted $\overleftarrow{(i, j)}$, entails updating $c_{jk+} = c_{ik} * c_{ji}$ for all predecessors k of i . If the elimination of an edge leaves an intermediate vertex with either no inedges or no outedges, then the vertex and all incident edges are removed from the graph. An edge elimination sequence ρ is *full* if $G'(\rho)$ contains no intermediate vertices. A full edge elimination sequence corresponds to fully accumulating $F'(\mathbf{x})$ as a matrix. Any edge elimination sequence that is not full is called *partial*.

For functions with scarce Jacobians, judicious choice of an edge elimination sequence can yield a remainder graph that has significantly fewer nonunit edges than G . Propagation can then be performed at a cost of $p|E^+(G'(\rho))| < p|E^+(G)|$ scalar multiplications. We distinguish between the propagation phase and the *preaccumulation* phase, which consists solely of edge eliminations and has a cost independent of p . When p is sufficiently large, the cost of the propagation phase dominates the computation. With this as our motivation, we focus on finding a sequence ρ of edge eliminations such that $|E^+(G'(\rho))|$ is minimized.

2. Randomized heuristics

A simple variant of the greedy heuristics described in [4] (hereafter called H_g) exploits scarcity by choosing the best possible edge elimination and keeping track of the best so far obtained complexity of the DAG. Here, we describe our experiments with two basic types of randomized local search methods, namely *Metropolis* [5] and *Simulated Annealing* (SA) [6]. Randomized local search methods are especially useful for hard combinatorial optimization problems about which little is known; successful use cases include problems such as TSP [7] VLSI design [8] and vertex elimination in AD [9, 10]. Our new randomized heuristics are compared with H_g when applied to the same set of sample codes and to a set of artificial DAGs. Among the heuristics we tested, a hybrid version of Metropolis produced the best results.

The edge elimination meta-graph. Consider a directed, Markov chain based dynamic *meta-graph* $\mathbf{G} = (\mathbf{V}, \mathbf{E})$ of all possible states G attains as it undergoes sequences of edge eliminations and their backtrackings along with a random walk process on \mathbf{G} . Each node $i \in \mathbf{V}$ corresponds to some state of G after a sequence of edge eliminations. The set \mathbf{E} of directed edges is partitioned into sets \mathbf{E}_s and \mathbf{E}_d for static and dynamic directed edges, respectively. A static directed edge $ij \in \mathbf{E}_s$ corresponds to the legal edge elimination that produces state j from i . A dynamic directed edge $ij \in \mathbf{E}_d$ represents a backward step (or backtracking) that is not an edge elimination. At any moment \mathbf{E}_d will contain only one edge. In other words, if at the k th step of a random walk, elimination ij was accepted, then at the $(k + 1)$ th step $ji \in \mathbf{E}_d$ will appear but the previous

backward edge will be removed. Note that, theoretically, at any state j many backward edges could be created since j can be reachable from more than one state. However, introduction of all backward edges can significantly increase the complexity of traversing the meta-graph, creating significant implementation difficulties. Thus, at any state (except the initial one) there will be only one backward edge. We denote by b_i the predecessor of i in the random walk over \mathbf{G} .

Let G_i denote the DAG corresponding to state i , $c(i)$ denote $|E^+(G_i)|$, and $N_i = N_i^+ \cup N_i^- \cup N_i^0$ denote the set of neighbors of i in \mathbf{G} , where

$$\begin{aligned} N_i^+ &= \{j \in \mathbf{V} \mid ij \in \mathbf{E} \text{ and } c(j) > c(i)\}; \\ N_i^- &= \{j \in \mathbf{V} \mid ij \in \mathbf{E} \text{ and } c(j) < c(i)\}; \\ N_i^0 &= \{j \in \mathbf{V} \mid ij \in \mathbf{E} \text{ and } c(j) = c(i)\}. \end{aligned}$$

The heuristics. We describe the classical version of Metropolis heuristic H_M in Algorithm 1 that we have started the computational experiments with.

Algorithm 1 Classical Metropolis algorithm

Require: initial graph G_0

```

1:  $i \leftarrow 0$ 
2: for  $k = 1, 2, \dots$  do
3:   while  $c(i)$  is sufficiently big do
4:     choose a random edge elimination  $ij$ 
5:     if  $c(j) \leq c(i)$  then
6:       accept  $ij$ 
7:     else
8:       accept  $ij$  with probability  $e^{-(c(j)-c(i))/T}$  for fixed  $T$ 
9:     end if
10:    if  $ij$  is accepted then
11:       $i \leftarrow j$ 
12:    end if
13:  end while
14: end for

```

The difference between classical Metropolis and SA lies in the choice of a temperature factor T . Instead of choosing a fixed T , a graduate cooling scheme for T is employed in SA. Carefully chosen (fixed and varying) schemes for T is a central issue of these algorithms. We refer the reader to [6] for a comprehensive background on these methods and to [9, 10] for example of using SA in automatic differentiation elimination problems. As discussed in [11], a Metropolis algorithm with the best temperature can outperform SA. The third heuristic H_h that we used, described in Algorithm 2, is a hybrid of Metropolis and a regular random walk.

3. Computational results

In this section, we describe numerical results obtained using Algorithm 2, which has been implemented as part of OpenAD [12]. We designed our numerical experiments with three types of computational graphs: (a) examples derived from applications; (b) a set of artificially generated single-expression-use (SEU) graphs [13]; and (c) random DAGs. We began with a series

Algorithm 2 Hybrid algorithm**Require:** initial graph G_0 , maximum number of steps K

```

1:  $i \leftarrow 0$ 
2: for  $k = 1, 2, \dots, K$  do
3:   if  $k$  is sufficiently big then
4:      $i \leftarrow 0$ 
5:   end if
6:   list all possible eliminations  $ij$ 
7:   if  $|N_i^- \setminus b_i| > 0$  then
8:     accept  $j \in N_i^+$  with normalized probability  $p_{ij} \propto e^{-(c(j)-c(i))/T}$ 
9:   else
10:    accept  $j \in N_i^+ \cup N_i^0$  with normalized probability  $p_{ij} \propto e^{-(c(j)-c(i)+1)/T}$ 
11:   end if
12:   if  $ij$  is accepted then
13:      $i \leftarrow j$ 
14:   end if
15: end for

```

of aggressive random walks over \mathbf{G} on SEU graph instances and randomly generated instances. After sufficiently many steps, the random walk is restarted, while keeping track of the best result achieved so far. Surprisingly, this trivial algorithm resulted in an improvement of 5-10% over H_g ; This provided the first indication that randomized local search can improve on H_g . In the real-life examples, however, this strategy was not able to beat H_g . The next stage of experiments consisted of designing the classical versions of Metropolis and SA. Independent of the aggressiveness of the gradual cooling scheme for T , both methods provided an improvement of up to 20% on the real-life instances and 10-15% on the artificial instances. The distribution of the results was proportional to the Gaussian which concentrated the most likely improvement on 12% over H_g on real-life instances. This series of experiments gave us an important observation regarding the steps that improve the current state i : *if there exist two eliminations ij and ik such that $ij, jk \in N_i^-$ and $c(j)$ and $c(k)$ are almost equal, the better of the two should not necessarily be accepted*. This observation guided the design of H_h which is the most successful of the heuristics. Note this issue cannot be addressed by any classical gradual cooling scheme, as such schemes work only on the elements of N_i^+ . Thus, no significant difference was observed between Metropolis and SA when different schemes were employed.

Our best computational results were obtained 2. The observed improvement on the real-life graphs was up to 35%. Examples of two experimental series for real-life graphs are presented in Figure 3. For every graph we ran 800 experiments. The results of H_g are 185 and 186 for the first and second graphs, respectively. The most interesting example can be observed in Figure 3(a), in which one can see that H_h almost separates two clusters of the solution quality. The maximum number of steps (K) (see algorithm 2) was $20|V|$ with 5 restarts (line 3) when k was reaching $4|V|$.

3.1. Interpretation of the Results

As is the case with other randomized heuristics one is curious if there is some structure in the problem that is essential for the improvement in the cost function. If such a structure could be found and characterized so that it may be recognized with relatively low computational

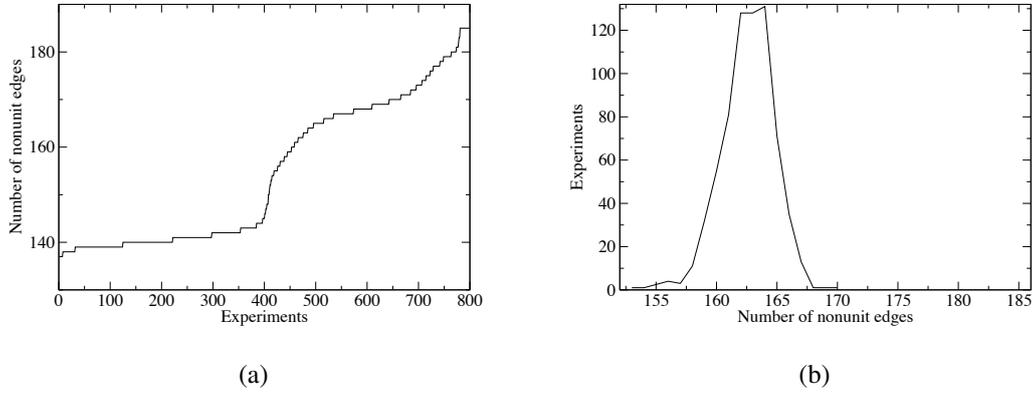


Figure 3: Results for H_h on two real-life instances: An example derived from a code for fluid dynamics (a); and an example from a complicated finite elements code.

complexity, then the heuristic could be modified to specifically search for and exploit it which then in turn would improve the chances of the thus modified randomized heuristic to find a good solution. The computational results show that (1) there can be a substantial improvement over H_g and (2) there can be a separation of the cost function values. We followed two paths to analyze the heuristic results.

One approach is to look at the *energy difference* δ_t of the t -th step of H_h relative to the step H_g would have taken. In other words, if i is the current state at the t -th step, ij is a transition accepted by H_h and ij^* is a transition that should minimize $c(k)$ over all possible transitions ik ,

$$\delta_t = c(j) - c(j^*).$$

There are 4 cases:

1. if there are targets with an improvement to the cost function we look among those targets for the energy δ_t of the randomized step vs. the deterministic step:
 - (a) if $\delta_t = 0$ then we randomize over the artificial order within the graph representation
 - (b) if $\delta_t > 0$ then we randomize wrt. the actual cost function
2. if there there are no targets with an improvement to the cost function we look among those targets with the analogous sub cases as described above:
 - (a) if $\delta_t = 0$ then we randomize over the artificial order within the graph representation
 - (b) if $\delta_t > 0$ then we randomize with respect to the actual cost function

For each randomized elimination sequence ρ we can then observe how many steps fall into one of the above four categories and also consider

$$\delta_\rho = \frac{\sum_{\delta_t \in \rho} \delta_t}{|\rho|}$$

as a (rough) measure of the distance of the given randomized heuristic from the deterministic heuristic as far as the actual cost function is concerned, where $|\rho|$ denotes the length of ρ . If we can find ρ with a substantial improvement of the cost function but all steps fall into categories 1(a) and 2(a), then the conclusion to be drawn is that the artificial order in the graph determines most of the cost, the suggested randomization over the energy would not be worth the effort, and one should break ties randomly. On the other hand, if there are no such sequences or even if there are only a few steps with a nonzero δ_s , then the suggested heuristic is justified.

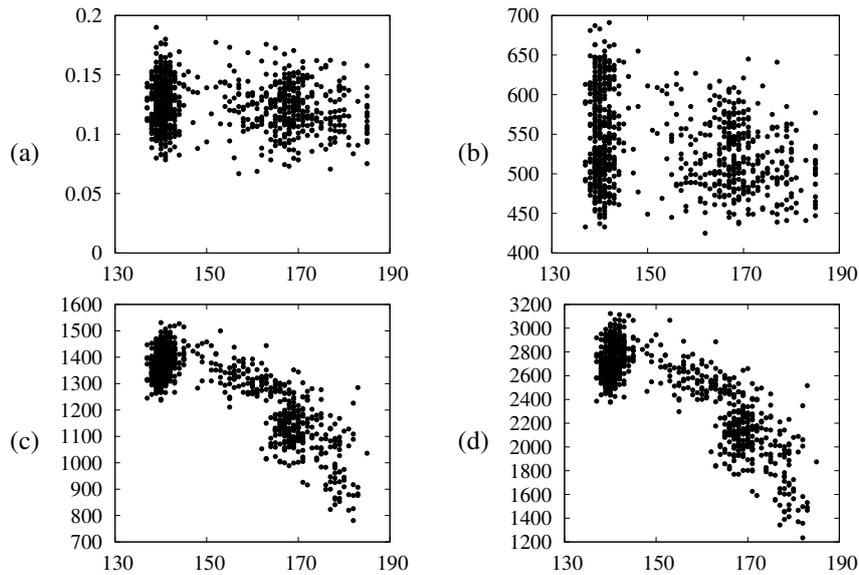


Figure 4: Plotted over the cost function values are δ_ρ in (a), $|\rho|$ in (b), vertex counts in compressed accumulation circuits in (c), edge counts in compressed accumulation circuits in (d).

The δ_ρ values shown in Figure 4(a) indicate that the cases 1(b) and 2(b) play a substantial role and therefore the suggested heuristic is effective. To gain insight into any structural properties, one will eventually have to look at the elimination sequences and compare them. This comparison is nontrivial because, as shown in Figure 4(b), the length of the sequences varies greatly. One might suspect that more elimination steps would be required to drive the cost down and therefore the best sequences would likely be longer than the others. We would like to point out that contrary to this plausible assumption, among the sequences with the lowest cost is also one that is the shortest (see the lower left datapoint in Figure 4(b)). This gives another indication that a structural property lies at the heart of the improvement.

The comparison of elimination sequences for the purpose of detecting structural properties is difficult for at least two reasons. First, the random sequences are of different length and, second, the sequences still embody a substantial artificial order that stems from tie breaking and is not at all relevant to the the question for structural properties.

The accumulation circuit. For a given edge elimination sequence ρ , The essential structure of the computation performed by the corresponding accumulation code is exhibited in the corresponding *accumulation circuit* $\Phi(\rho)$ (or just Φ where there is no ambiguity). $\Phi(\rho)$ is an arithmetic circuit whose leaves — the inputs to the circuit — correspond to the variables labeling the edges

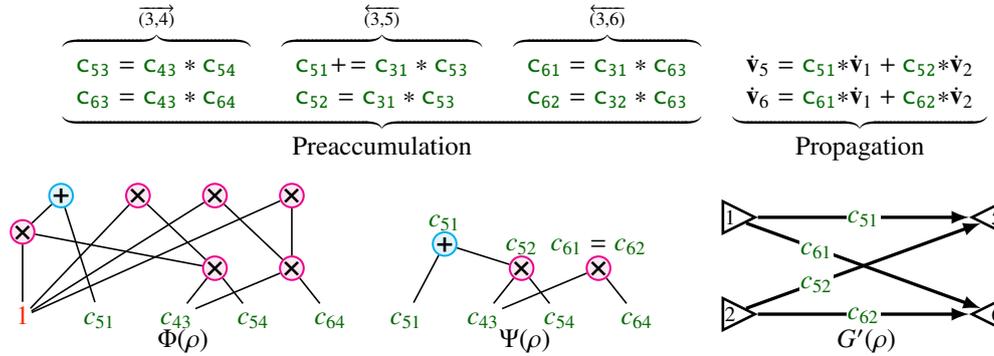


Figure 5: A full edge elimination sequence $\rho = ((3, 4), (5, 3), (6, 3))$ is shown along with the corresponding accumulation circuit $\Phi(\rho)$, compressed accumulation circuit $\Psi(\rho)$, and remainder graph $G'(\rho)$. Note the redundancy in the fully preaccumulated Jacobian. Any full edge elimination sequence will result in the same remainder graph G' , though different sequences generally imply different computational costs for the preaccumulation phase. However, sequences that are functionally identical may look quite different: For $\bar{\rho} = ((5, 4), (3, 4), (1, 2), (2, 3))$, we have $G'(\rho) = G'(\bar{\rho})$, $\Phi(\rho) = \Phi(\bar{\rho})$, and $\Psi(\rho) = \Psi(\bar{\rho})$.

$E(G)$. The internal nodes of $\Phi(\rho)$ all have exactly two predecessors and are labeled either $+$ (sum gates) or \times (product gates), where the label for a gate $\alpha \in \Phi$ is denoted o_α . In general, an accumulation circuit will have many outputs, each of which computes a value carried by an edge in the remainder graph G' . (Note that, for edge elimination, this is also true of some nodes in the accumulation circuit that are not maximal.) In general, the number of edge elimination sequences (partial or full) is much bigger than the number of accumulation circuits that can result from an edge elimination sequence. This gives us an equivalence relation where two edge elimination sequences ρ_1 and ρ_2 may satisfy $\Phi(\rho_1) = \Phi(\rho_2)$ in addition to satisfying $G'(\rho_1) = G'(\rho_2)$. Note that having $G'(\rho_1)$ equal to $G'(\rho_2)$ is necessary but not sufficient for two edge elimination sequences to be considered equivalent. *Compression* of the accumulation circuit, which establishes a kind of canonical form, allows for an even coarser equivalence relation.

The hope is that one may find a pair of sequences that has a substantial difference in the cost function yet the accumulation circuits are similar. The remaining difference then might point to a particular structure that triggers the difference in the cost. The following properties of the accumulation circuit guide the compression. All non-constant minimal vertices are distinct; all constant minimal vertices have either identical values or else are considered distinct; there can be non-maximal vertices referenced by the remainder graph; and all non-minimal vertices are either multiplication or addition operations. The circuit compression consists of the following steps. (1) collapse all vertices that are minimal, constant, and have identical values to a single representer vertex; (2) replace any constant valued subgraphs S that evaluate to exactly 1.0 and have a single outedge (S, j) by a new edge $(1.0, j)$ from the constant 1.0 representer vertex; (3) contract any edge (i, j) such that i is non-minimal, j is the only successor of i , $o_i = o_j$, and i is not referenced by the remainder graph; (4) collapse to i all non-minimal vertices j , if i and j have identical predecessor multi-sets³ and $o_i = o_j$.

The numbers shown in Figure 4(d) show that compression yields reductions between 611

³The accumulation circuit can have parallel paths; we must determine how often a given vertex is a predecessor.

and 1130. On the compressed accumulation circuits we can recursively build a signature $s_v = (o_v, c_v, \mathcal{V}_v, C_v)$ for each vertex v by considering its optional operation o_v or constant value c_v , a multiset \mathcal{V}_v of its dependencies on variable minimal vertices, and a multiset C_v of elements (c, o, C^*) representing operations with constant values.

- For all minimal v with constant value c we set $s(v) = (., c, \emptyset, \emptyset)$
- For all variable minimal v with edge label c_{ji} we set $s(v) = (., ., \{c_{ji}\}, \emptyset)$
- For all non-minimal v with direct predecessors \mathcal{P} :
 - $\mathcal{V}' = \bigcup_{p \in \mathcal{P}} \mathcal{V}_p$; $C' = \bigcup_{p \in \mathcal{P}} C_p$
 - if $\mathcal{V}' = C' = \emptyset$ then compute new constant c_v value from all predecessors and set $s(v) = (., c, \emptyset, \emptyset)$
 - if $\exists p \in \mathcal{P} : \mathcal{V}_p = C_p = \emptyset$, then compute new constant c_v value from the constants of those p and set $s(v) = (o_v, ., \mathcal{V}', \{(o_v, c_v, C')\})$ otherwise set $s(v) = (o_v, ., \mathcal{V}', C')$

The signatures can be built bottom up in the accumulation circuit and with the multisets lexicographically ordered and suitably represented as a string can be used to compare two vertices by string comparison. This has been implemented to enable comparisons between sequences from the cluster of good solutions and the cluster of solutions closer to the H_g result. Unfortunately, even in circuits of similar size we could match only less than half of the vertices. The alternative search for some vertices occurring only in circuits from the preferred cluster of solutions indicating a crucial step in the elimination has so far not produced tangible results and is subject to further research.

Acknowledgments

This work was supported in part by U.S. Department of Energy, under Contract DE-AC02-06CH11357.

- [1] A. Griewank, A mathematical view of automatic differentiation, in: Acta Numerica, Vol. 12, Cambridge University Press, 2003, pp. 321–398. doi:10.1017/S0962492902000132.
- [2] A. Griewank, O. Vogel, Analysis and exploitation of Jacobian scarcity, in: H. G. Bock, E. Kostina, H. X. Phu, R. Rannacher (Eds.), Modeling, Simulation and Optimization of Complex Processes, Springer, Berlin, 2005, pp. 149–164. doi:10.1007/3-540-27170-8_12.
- [3] A. Griewank, A. Walther, Evaluating Derivatives: Principles and Techniques of Algorithmic Differentiation, 2nd Edition, no. 105 in Other Titles in Applied Mathematics, SIAM, Philadelphia, PA, 2008.
- [4] A. Lyons, J. Utke, On the practical exploitation of scarcity, in: C. H. Bischof, H. M. Bücker, P. D. Hovland, U. Naumann, J. Utke (Eds.), Advances in Automatic Differentiation, Vol. 64 of Lecture Notes in Computational Science and Engineering, Springer, Berlin, 2008, pp. 103–114. doi:10.1007/978-3-540-68942-3_10.
- [5] N. Metropolis, A. W. Rosenbluth, M. N. Rosenbluth, A. H. Teller, E. Teller, Equation of state calculations by fast computing machines, The Journal of Chemical Physics 21 (6) (1953) 1087–1092. doi:10.1063/1.1699114.
- [6] S. Kirkpatrick, C. D. Gelatt, M. P. Vecchi, Optimization by simulated annealing, Science 220, 4598 (1983) 671–680.
- [7] J. J. Schneider, S. Kirkpatrick, Stochastic Optimization (Scientific Computation), Springer-Verlag New York, Inc., Secaucus, NJ, USA, 2006.
- [8] D. F. Wong, H. W. Leong, C. L. Liu, Simulated annealing for VLSI design, Kluwer Academic Publishers, Norwell, MA, USA, 1988.
- [9] U. Naumann, P. Gottschling, Prospects for simulated annealing in automatic differentiation, in: K. Steinhöfel (Ed.), Stochastic Algorithms: Foundations and Applications, no. 2264 in LNCS, Springer, Berlin, 2001, pp. 355–359. doi:10.1007/3-540-45322-9_9.

- [10] U. Naumann, P. Gottschling, Simulated annealing for optimal pivot selection in Jacobian accumulation, in: A. Albrecht, K. Steinhöfel (Eds.), *Stochastic Algorithms: Foundations and Applications*, no. 2827 in *Lecture Notes in Computer Science*, Springer, 2003, pp. 83–97. doi:10.1007/b13596.
- [11] M. Jerrum, A. Sinclair, M. C. Algorithms, *The markov chain monte carlo method: An approach to approximate counting and integration*, PWS Publishing, 1996, pp. 482–520.
- [12] J. Utke, U. Naumann, M. Fagan, N. Tallent, M. Strout, P. Heimbach, C. Hill, C. Wunsch, OpenAD/F: A modular, open-source tool for automatic differentiation of Fortran codes, *ACM Transactions on Mathematical Software* 34 (4) (2008) 18:1–18:36. doi:10.1145/1377596.1377598.
- [13] U. Naumann, Y. Hu, Optimal vertex elimination in single-expression-use graphs, *ACM Transactions on Mathematical Software* 35 (1) (2008) 1–20. doi:10.1145/1377603.1377605.

The submitted manuscript has been created by UChicago Argonne, LLC, Operator of Argonne National Laboratory (“Argonne”). Argonne, a U.S. Department of Energy Office of Science laboratory, is operated under Contract No. DE-AC02-06CH11357. The U.S. Government retains for itself, and others acting on its behalf, a paid-up nonexclusive, irrevocable worldwide license in said article to reproduce, prepare derivative works, distribute copies to the public, and perform publicly and display publicly, by or on behalf of the Government.