

---

# Differentiate Everything with a Reversible Domain-Specific Language

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

Reverse-mode automatic differentiation (AD) suffers from the issue of having too much space overhead in order to trace back intermediate computational states for back propagation. The traditional method to trace back states is called checkpointing that stores intermediate states into a global stack and restore state through either stack pop or re-computing. The overhead of stack manipulations and re-computing makes the general purposed (or not tensor based) AD engines unable to meet many industrial needs. Instead of checkpointing, we propose to use reverse computing to trace back states by designing and implementing a reversible programming eDSL, where a program can be executed bi-directionally without implicit stack operations. The absence of implicit stack operations make the program compatible with existing compiler features, including utilizing existing optimization passes and compiling the code as GPU kernels. We implement AD for sparse matrix operations and some machine learning applications to show that the performance our framework has the state-of-the-art performance.

## 1 Introduction

Most popular automatic differentiation (AD) packages in the market, such as TensorFlow [Abadi et al. \(2015\)](#), Pytorch [Paszke et al. \(2017\)](#) and Flux [Innes et al. \(2018\)](#) implements reverse mode AD at the tensor level to meet the need of machine learning. Later, People in the scientific computing domain also find the powerfulness of these AD tools, people use these them to solve scientific problems such as seismic inversion [Zhu et al. \(2020\)](#), variational quantum circuits simulation [Bergholm et al. \(2018\)](#) and variational tensor network simulation [Liao et al. \(2019\)](#); [Roberts et al. \(2019\)](#). To meet the diversified need in these applications, one often needs to add manually defined backward rules, for example

1. In order to differentiate sparse matrix operations for Hamiltonian engineering [Hao Xie and Wang](#), we need to define backward rules for sparse matrix operations and dominant eigensolvers [Golub and Van Loan \(2012\)](#),
2. In tensor network algorithms to study the phase transition problem [Golub and Van Loan \(2012\)](#); [Liao et al. \(2019\)](#); [Seeger et al. \(2017\)](#); [Wan and Zhang \(2019\)](#); [Hubig \(2019\)](#), one needs defining backward rules for singular value decomposition (SVD) function and QR decomposition.

To avoid defining backward rules manually, one can also use a general purposed AD (GP-AD) packages like Tapenade [Hascoet and Pascual \(2013\)](#), OpenAD [Utke et al. \(2008\)](#) and Zygote [Innes \(2018\)](#); [Innes et al. \(2019\)](#) to differentiate a general program. These tools has been used in non-tensor based applications such as bundle adjustment [Shen and Dai \(2018\)](#) and earth system

simulation [Forget et al. \(2015\)](#). They read the source code from a user and generate code to compute the gradients. However, these packages have their own limitations too. In many practical applications, differentiating a program might do billions of computations. Frequent caching of data slows down the program significantly, while the memory usage will become a bottleneck as well. Moreover, implicit caching is not compatible with differentiating GPU kernels.

These needs call for a GP-AD framework that does not cache for users automatically. Hence we propose to implement the reverse mode AD on a reversible (domain-specific) programming language [Perumalla \(2013\)](#); [Frank \(2017\)](#). So that the intermediate states of a program can be traced backward with no extra effort. In a reversible languages, the memory allocation and deallocation are explicit, with which flexible time-space tradeoff is available. Meanwhile, one can also utilize the reversibility to reverse the program without space overhead, which is proven to significantly decrease the memory usage in unitary recurrent neural networks [MacKay et al. \(2018\)](#), normalizing flow [Dinh et al. \(2014\)](#), hyperparameter learning [Maclaurin et al. \(2015\)](#) and residual neural networks [Behrmann et al. \(2018\)](#). Using reversible programming language makes all these happen naturally without extra framework designs.

There have been many prototypes of reversible languages like Janus [Lutz \(1986\)](#), R (not the popular one) [Frank \(1997\)](#), Erlang [Lanese et al. \(2018\)](#) and object-oriented ROOPL [Haulund \(2017\)](#). In the past, the primary motivation to study reversible programming is to support reversible computing devices [Frank and Knight Jr \(1999\)](#) like adiabatic complementary metal-oxide-semiconductor (CMOS) [Koller and Athas \(1992\)](#), molecular mechanical computing system [Merkle et al. \(2018\)](#) and superconducting system [Likharev \(1977\)](#); [Semenov et al. \(2003\)](#), where a reversible computing device is more energy-efficient by the Landauer’s principle [Landauer \(1961\)](#). The above reversible programming languages are well defined so that can be compiled to future hardwares. However, the drawback is they can hardly be used in real scientific programming, since most of them do not have basic elements like floating point numbers, arrays and complex numbers. This motivates us to build a new embeded domain specific language (eDSL) in Julia [Bezanson et al. \(2012, 2017\)](#).

In this paper, we first introduce the language design of a reversible programming language and introduce our reversible eDSL NiLang in Sec. 2. In Sec. 3, we explain the implementation of automatic differentiation in this eDSL. In Sec. 4, we benchmark the performance of NiLang with other AD packages and explain why reversible programming AD is fast. In the appendix, we show the detailed language design of NiLang, show some examples used in the benchmark, discuss several important issues including the time-space tradeoff, reversible instructions and hardware, and finally, an outlook to some open problems to be solved.

## 2 Language design

### 2.1 An introduction to the reversible language design

#### 2.1.1 Reversible memory management

A distinct feature of reversible memory management is that the content of a variable must be known when it is deallocated. We denote the allocation of a zero emptied memory as  $x \leftarrow \emptyset$ , and the corresponding deallocation as  $x \rightarrow \emptyset$ . If a variable is allocated and deallocated in a local scope, we call it an ancilla. A variable can also be pushed to a stack and used later with a pop statement.

#### 2.1.2 Reversible control flows

One can define reversible `if`, `for` and `while` statements in a reversible program. Fig. 1 (a) shows the flow chart of executing the reversible `if` statement. There are two condition expressions in this chart, a precondition and a postcondition. The precondition decides which branch to enter in the forward execution. After executing the specific branch, the program checks the consistency between precondition and postcondition to make sure they are consistent. To reverse this statement, one can exchange the precondition and postcondition, and reverse the expressions in both branches. Fig. 1 (b) shows the flow chart of the reversible `while` statement. There are also two conditions expressions. Before executing the condition expressions, the program preassumes the postcondition is false. After each iteration, the program asserts the postcondition to be true. To reverse this statement, one can exchange the precondition and postcondition, and reverse the body statements. The reversible `for`

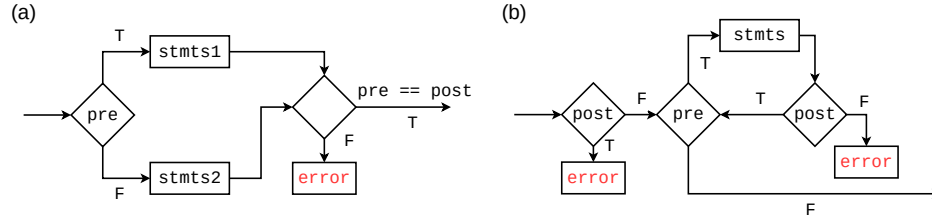


Figure 1: The flow chart for reversible (a) if statement and (b) while statement. “pre” and “post” represents precondition and postconditions respectively.

statement is similar to the irreversible one except that after execution, the program will assert the iterator to be unchanged. To reverse this statement, one can exchange `start` and `stop` and inverse the sign of `step`.

### 2.1.3 Reversible arithmetic instructions

Mathematically, a irreversible mapping  $y = f(\text{args} \dots)$  can be trivially transformed to its reversible form  $y \oplus= f(\text{args} \dots)$  or  $y \vee= f(\text{args} \dots)$  ( $\vee$  is the bitwise XOR), where  $y$  is a pre-empted variable. But in numeric computing with finite precision, this is not the whole story, because the reversibility of arithmetic instruction is closely related to the number system. For integer and fixed point number system,  $y \oplus= f(\text{args} \dots)$  and  $y \oplus= f(\text{args} \dots)$  are rigorously reversible. For logarithmic number system and tropical number system [Speyer and Sturmfels \(2009\)](#),  $y \oplus= f(\text{args} \dots)$  and  $y \oplus= f(\text{args} \dots)$  as reversible (not introducing the zero element). While for floating point numbers, none of the above operations are rigorously reversible. However, for convenience, we ignore the rounding errors in floating point  $+$  and  $-$  operations and treat them on equal footing with fixed point numbers in the following discussion. Besides the above operations, SWAP operation that exchanges the values of two variables is also widely used in reversible programming.

## 2.2 NiLang, a reversible eDSL suited for scientific computing

In the introduction, we introduced some reversible programming languages. These languages lacks essential components for scientific programming like arrays and complex numbers. Meanwhile most of them are designed as a stand alone language that can not be embeded in other machine learning frameworks. Hence we developped an embedded domain-specific language (eDSL) NiLang on top of the host language language Julia [Bezanson et al. \(2012, 2017\)](#). Julia is a popular language for scientific programming and machine learning. We choose Julia because the speed is our primary concern. Julia is a language with high abstraction, however, its clever design of type inference and just in time compiling make it has a C like speed. Also, it has rich features for meta-programming, and the package for pattern matching [MLStyle](#) allows us to define an eDSL compiler in less than 2000 lines. Comparing with a regular reversible programming language, NiLang features array operations, rich number systems including floating point number, complex number, fixed point number and logarithmic number. Besides the above nice features, it also has some "bad" features to meet the practical needs. For example, it views the floating point  $+$  and  $-$  operations as reversible. It also allows user to extend instruction sets and sometimes inserting external statements. These features are not compatible with future reversible hardware. By the time of writting, the version of NiLang is v0.7.2.

Let's start learning it by defining a reversible adder.

Listing 1: A reversible adder

```

122 @i function adder(y!::Real, x::Real)
      y! += x
    end

    @assert adder(2, 3) == (5, 3)
    @assert (~adder)(5, 3) == (2, 3)

```

Macro @i generates two functions that reversible to each other `adder` and `~adder`, each defines a mapping  $\mathbb{R}^2 \rightarrow \mathbb{R}^2$ . The ! after a symbol is a part of the name to indicate that a variable is changed. A reversible += instruction is always defined as `y += f(args...)`, where `f` is a mapping that allows to be irreversible, or just leave empty for identity mapping. We can easily check these two functions are reversible to each other. Then let's see a more advanced example of computing the complex valued log (a built in function).

Listing 2: Reversible complex valued log function `y += log(|x|) + iArg(x)`.

```

129 @i @inline function (:+=)(log)(y!::Complex{T}, x::Complex{T}) where T
      @routine begin
        n ← zero(T)
        n += abs(x)
      end
      y!.re += log(n)
      y!.im += angle(x)
      ~@routine
    end

```

Here, the macro @inline tells the compiler that this function can be inlined. `n ← zero(T)` is the ancilla allocation statement. One can input “←” and “→” by typing “\leftarrow[TAB KEY]” and “\rightarrow[TAB KEY]” respectively in a Julia editor or REPL. @routine and ~@routine are macros for computing and uncomputing, i.e. ~@routine means running the statement marked with @routine backwards. One can use the `begin ... end` statement to wrap multiple statements as one. NiLang view every field of a variable as mutable, so that the real part (`y!.re`) and imaginary (`y!.im`) of a complex number can also be changed directly.

We can want to apply this log function to an array, we can define

Listing 3: Applying the log function to an array.

```

138 @i function broadcasted_log!(y!::Array{Complex{T}, N}, x::Array{Complex{T}, N}) where T
      N ← min(length(x), length(y!))
      for i=1:N
        y![i] += log(x[i])
      end
    end

```

### 139 3 Reversible automatic differentiation

#### 140 3.1 First order gradient

141 If we inline all the instructions, the program would be like Listing. 4. The automatically generated  
 142 inverse program (i.e.  $(y, x) \rightarrow (y - \log(x), x)$ ) is like Listing. 5.

Listing 4: The expanded function body of Listing 3.

```

143 N ← min(length(x), length(y!))
    for i=1:N
      @routine begin
        nsq ← zero(T)
        n ← zero(T)
        nsq += x[i].re ^ 2
        nsq += x[i].im ^ 2
        n += sqrt(nsq)
      end
      y![i].re += log(n)
      y![i].im += atan(x[i].im, x[i].re)
    ~@routine
  end
  N → min(length(x), length(y!))

```

Listing 5: The inverse of Listing 4.

```

N ← min(length(x), length(y!))
for i=N:-1:1
  @routine begin
    nsq ← zero(T)
    n ← zero(T)
    nsq += x[i].re ^ 2
    nsq += x[i].im ^ 2
    n += sqrt(nsq)
  end
  y![i].re -= log(n)
  y![i].im -= atan(x[i].im, x[i].re)
~@routine
end
N → min(length(x), length(y!))

```

144 If we want to compute the adjoint of the complex valued  $(:+=)$  (log) function, one simply insert  
 145 the gradient code into the reversed code in Listing 5. As show in Listing 6, the original code is  
 146 highlighted with yellow background color, they now apply on the value field (.x) of the input value.  
 147 Along with these reversed code, we have inserted a bundle of extra code to update the gradient  
 148 field (.g). @zeros TYPE var1 var2... is the macro to allocate multiple ancillas. Since these  
 149 “allocated” variables are scalars, they do not really access the system memory. Its inverse operations  
 150 starts with ~@zeros returns zero emptied ancillas to the system.

Listing 6: Insert the gradient code into Listing 5.

```

151 N ← min(length(x), length(y!))
    for i=N:-1:1
      @routine begin
        nsq ← zero(GVar{T,T})
        n ← zero(GVar{T,T})

        gsqa ← zero(T)
        gsqa += x[i].re.x * 2
        x[i].re.g -= gsqa * nsq.g
        gsqa -= nsq.x * 2
        gsqa -= x[i].re.x * 2
        gsqa → zero(T)
        nsq.x += x[i].re.x ^ 2

        gsqb ← zero(T)
        gsqb += x[i].im.x * 2
        x[i].im.g -= gsqb * nsq.g
        gsqb -= x[i].im.x * 2
        gsqb → zero(T)
        nsq.x += x[i].im.x ^ 2

        @zeros T ra rb
        ra += sqrt(nsq.x)
        rb += 2 * ra
        nsq.g -= n.g / rb
        rb -= 2 * ra

        ra -= sqrt(nsq.x)
        ~@zeros T ra rb
        n.x += sqrt(nsq.x)
      end

      y![i].re.x -= log(n.x)
      n.g += y![i].re.g / n.x

      y![i].im.x -= atan(x[i].im.x, x[i].re.x)
      @zeros T xy2 jac_x jac_y
      xy2 += abs2(x[i].re.x)
      xy2 += abs2(x[i].im.x)
      jac_y += x[i].re.x / xy2
      jac_x += (-x[i].im.x) / xy2
      x[i].im.g += y![i].im.g * jac_y
      x[i].re.g += y![i].im.g * jac_x
      jac_x -= (-x[i].im.x) / xy2
      jac_y -= x[i].re.x / xy2
      xy2 -= abs2(x[i].im.x)
      xy2 -= abs2(x[i].re.x)
      ~@zeros T xy2 jac_x jac_y

    ~@routine
  end
  %

```

152 In really implementation, we utilize Julia’s multiple dispatch. And “insert” the gradient code by  
 153 overloading the basic instructions for the gradient wrapper type GVar. The same strategy has been  
 154 used in the ForwardDiff package in Julia. Thanks to the just in time compiling technology, the above  
 155 code does not run as long as it looks. Computing the gradient takes similar time as computing the  
 156 complex valued log lone with Julia’s builtin log function.

157 One does not need to define a similar function on  $(:+=)$  (log) because macro @i will generate it  
 158 automatically. Notice that taking inverse and computing gradients commute [McInerney \(2015\)](#).

### 159 3.2 Hessians

160 Combining the uncomputing program in NiLang with dual-numbers is a simple yet efficient way to  
 161 obtain Hessians. The dual number is the scalar type for computing gradients in the forward mode

AD, it wraps the original scalar with a extra gradient field. The gradient field of a dual number is updated automatically as the computation marches forward. By wrapping the elementary type with `Dual` defined in package `ForwardDiff` [Revels et al. \(2016\)](#) and throwing it into the gradient program defined in `NiLang`, one obtains one row/column of the Hessian matrix straightforward. We will show a benchmark in Sec. 4.2.

### 3.3 Complex numbers

To differentiate complex numbers, we re-implemented complex instructions reversibly. For example, with the reversible function defined in Listing. 2, we can differentiated complex valued log with no extra effort.

### 3.4 CUDA kernels

CUDA programming is playing a significant role in high-performance computing. In Julia, one can write GPU compatible functions in native Julia language with `KernelAbstractions` [Besard et al. \(2017\)](#). Since `NiLang` does not push variables into stack automatically for users, it is safe to write differentiable GPU kernels with `NiLang`. We will show this feature in the benchmarks of bundle adjustment (BA) in Sec. 4.3. Here, one should notice that the shared read in forward pass will become shared write in the backward pass, which may result in incorrect gradients. We will review this issue in the supplementary material.

## 4 Benchmarks

It is interesting to see how does our framework comparing with the state-of-the-art GP-AD frameworks, including source code transformation based `Tapenade` and `Zygote` and operator overloading based `ForwardDiff` and `ReverseDiff`. Since most DS-AD packages like famous `Tensorflow` and `Pytorch` are not designed for the following using cases, we do not benchmark those package. In the following benchmarks, the CPU device is Intel(R) Xeon(R) Gold 6230 CPU @ 2.10GHz, and the GPU device is Nvidia Titan V. For `NiLang` benchmarks, we have turned off the reversibility check off to achieve a better performance. Codes used in benchmarks could be found the in Examples section of the supplementary material.

### 4.1 Sparse matrices

We benchmarked the call, uncall and backward propagation time used for sparse matrix dot product and matrix multiplication. Here, we estimate the time for back propagating gradients rather than including both forward and backward, since `mul!` does not output a scalar as loss.

	dot	mul! (complex valued)
Julia-O	3.493e-04	8.005e-05
NiLang-O	4.675e-04	9.332e-05
NiLang-B	5.821e-04	2.214e-04

Table 1: Absolute runtimes in seconds for computing the objectives (O) and the backward pass (B) of sparse matrix operations. The matrix size is  $1000 \times 1000$ , and the element density is 0.05. The total time used in computing gradient can be estimated by summing “O” and “B”.

The time used for computing backward pass is approximately 1.5-3 times the Julia’s native forward pass. This is because the instruction length of differentiating basic arithmetic instructions is longer than pure computing.

## 4.2 Graph embedding problem

Since one can combine ForwardDiff and NiLang to obtain Hessians, it is interesting to see how much performance we can get in differentiating the graph embedding program. The problem definition could be found in the supplementary material.

$k$	2	4	6	8	10
Julia-O	4.477e-06	4.729e-06	4.959e-06	5.196e-06	5.567e-06
NiLang-O	7.173e-06	7.783e-06	8.558e-06	9.212e-06	1.002e-05
NiLang-U	7.453e-06	7.839e-06	8.464e-06	9.298e-06	1.054e-05
NiLang-G	1.509e-05	1.690e-05	1.872e-05	2.076e-05	2.266e-05
ReverseDiff-G	2.823e-05	4.582e-05	6.045e-05	7.651e-05	9.666e-05
ForwardDiff-G	1.518e-05	4.053e-05	6.732e-05	1.184e-04	1.701e-04
Zygote-G	5.315e-04	5.570e-04	5.811e-04	6.096e-04	6.396e-04
(NiLang+F)-H	4.528e-04	1.025e-03	1.740e-03	2.577e-03	3.558e-03
ForwardDiff-H	2.378e-04	2.380e-03	6.903e-03	1.967e-02	3.978e-02
(ReverseDiff+F)-H	1.966e-03	6.058e-03	1.225e-02	2.035e-02	3.140e-02

Table 2: Absolute times in seconds for computing the objectives (O), uncall objective (U), gradients (G) and Hessians (H) of the graph embedding program.  $k$  is the embedding dimension, the number of parameters is  $10k$ .

In Table 2, we show the the performance of different implementations by varying the dimension  $k$ . The number of parameters is  $10k$ . As the baseline, (a) shows the time for computing the function call. We have reversible and irreversible implementations, where the reversible program is slower than the irreversible native Julia program by a factor of  $\sim 2$  due to the uncomputing overhead. The reversible program shows the advantage of obtaining gradients when the dimension  $k \geq 3$ . The larger the number of inputs, the more advantage it shows due to the overhead proportional to input size in forward mode AD. The same reason applies to computing Hessians, where the combo of NiLang and ForwardDiff gives the best performance for  $k \geq 3$ .

## 4.3 Gaussian mixture model and bundle adjustment

We reproduced the benchmarks for Gaussian mixture model (GMM) and bundle adjustment (BA) in [Srajer et al. \(2018\)](#) by re-writing the programs in a reversible style. We show the results in Table 3 and Table 4. In our new benchmarks, we also rewrite the ForwardDiff program for a fair benchmark, this explains the difference between our results and the original benchmark. The Tapenade data is obtained by executing the docker file provided by the original benchmark, which provides a baseline for comparison.

# parameters	3.00e+1	3.30e+2	1.20e+3	3.30e+3	1.07e+4	2.15e+4	5.36e+4	4.29e+5
Julia-O	9.844e-03	1.166e-02	2.797e-01	9.745e-02	3.903e-02	7.476e-02	2.284e-01	3.593e+00
NiLang-O	3.655e-03	1.425e-02	1.040e-01	1.389e-01	7.388e-02	1.491e-01	4.176e-01	5.462e+00
Tapende-O	1.484e-03	3.747e-03	4.836e-02	3.578e-02	5.314e-02	1.069e-01	2.583e-01	2.200e+00
ForwardDiff-G	3.551e-02	1.673e+00	4.811e+01	1.599e+02	-	-	-	-
NiLang-G	9.102e-03	3.709e-02	2.830e-01	3.556e-01	6.652e-01	1.449e+00	3.590e+00	3.342e+01
Tapenade-G	5.484e-03	1.434e-02	2.205e-01	1.497e-01	4.396e-01	9.588e-01	2.586e+00	2.442e+01

Table 3: Absolute runtimes in seconds for computing the objective (O) and gradients (G) of GMM with 10k data points. “-” represents missing data due to not finishing the computing in limited time.

In the GMM benchmark, NiLang’s objective function has overhead comparing with irreversible programs in most cases. Except the uncomputing overhead, it is also because our naive reversible matrix-vector multiplication is much slower than the highly optimized BLAS function, where the matrix-vector multiplication is the bottleneck of the computation. The forward mode AD suffers from too large input dimension in the large number of parameters regime. Although ForwardDiff batches the gradient fields, the overhead proportional to input size still dominates. The source to



source AD framework Tapenade is faster than NiLang in all scales of input parameters, but the ratio between computing the gradients and the objective function are close.

# measurements	3.18e+4	2.04e+5	2.87e+5	5.64e+5	1.09e+6	4.75e+6	9.13e+6
Julia-O	2.020e-03	1.292e-02	1.812e-02	3.563e-02	6.904e-02	3.447e-01	6.671e-01
NiLang-O	2.708e-03	1.757e-02	2.438e-02	4.877e-02	9.536e-02	4.170e-01	8.020e-01
Tapenade-O	1.632e-03	1.056e-02	1.540e-02	2.927e-02	5.687e-02	2.481e-01	4.780e-01
ForwardDiff-J	6.579e-02	5.342e-01	7.369e-01	1.469e+00	2.878e+00	1.294e+01	2.648e+01
NiLang-J	1.651e-02	1.182e-01	1.668e-01	3.273e-01	6.375e-01	2.785e+00	5.535e+00
NiLang-J (GPU)	1.354e-04	4.329e-04	5.997e-04	1.735e-03	2.861e-03	1.021e-02	2.179e-02
Tapenade-J	1.940e-02	1.255e-01	1.769e-01	3.489e-01	6.720e-01	2.935e+00	6.027e+00

Table 4: Absolute runtimes in seconds for computing the objective (O) and Jacobians (J) in bundle adjustment.

In the BA benchmark, reverse mode AD shows slight advantage over ForwardDiff. The bottleneck of computing this large sparse Jacobian is computing the Jacobian of a elementary function with 15 input arguments and 2 output arguments, where input space is larger than the output space. In this instance, our reversible implementation is even faster than the source code transformation based AD framework Tapenade. With KernelAbstractions, we run our zero allocation reversible program on GPU, which provides a >200x speed up.

## Broader Impact

Our automatic differentiation in a reversible eDSL brings the field of reversible computing to the modern context. We believe it will be accepted by the public to meet current scientific automatic differentiation needs and aim for future energy-efficient reversible devices. For solving practical issues, in an unpublished paper, we have successfully differentiated a spin-glass solver to find the optimal configuration on a  $28 \times 28$  square lattice in a reasonable time. There are also some interesting applications like normalizing flow and bundle adjustment in the example folder of [NiLang](#) repository and [JuliaReverse](#) organization. For the future, energy consumption is an even more fundamental issue than computing time and memory. Current computing devices, including CPU, GPU, TPU, and NPU consume much energy, which will finally hit the "energy wall". We must get prepared for the technical evolution of reversible computing (quantum or classical), which may cost several orders less energy than current devices.

We also see some drawbacks to the current design. It requires the programmer to change to programing style rather than put effort into optimizing regular codes. It is not fully compatible with modern software stacks. Everything, including instruction sets and BLAS functions, should be redesigned to support reversible programming better. We put more potential issues and opportunities in the discussion section of the supplementary material. Solving these issues requires the participation of people from multiple fields.

## References

- M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser, M. Kudlur, J. Levenberg, D. Mané, R. Monga, S. Moore, D. Murray, C. Olah, M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viégas, O. Vinyals, P. Warden, M. Wattenberg, M. Wicke, Y. Yu, and X. Zheng, “[TensorFlow: Large-scale machine learning on heterogeneous systems](#),” (2015), software available from tensorflow.org.
- A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, in *NIPS Autodiff Workshop* (2017).
- M. Innes, E. Saba, K. Fischer, D. Gandhi, M. C. Rudilosso, N. M. Joy, T. Karmali, A. Pal, and V. Shah, “Fashionable modelling with flux,” (2018), [arXiv:1811.01457 \[cs.PL\]](#).



258 W. Zhu, K. Xu, E. Darve, and G. C. Beroza, “A general approach to seismic inversion with  
259 automatic differentiation,” (2020), [arXiv:2003.06027 \[physics.comp-ph\]](#) .

260 V. Bergholm, J. Izaac, M. Schuld, C. Gogolin, M. S. Alam, S. Ahmed, J. M. Arrazola, C. Blank,  
261 A. Delgado, S. Jahangiri, K. McKiernan, J. J. Meyer, Z. Niu, A. Száva, and N. Killoran,  
262 “PennyLane: Automatic differentiation of hybrid quantum-classical computations,” (2018),  
263 [arXiv:1811.04968 \[quant-ph\]](#) .

264 H.-J. Liao, J.-G. Liu, L. Wang, and T. Xiang, [Physical Review X](#) **9** (2019), 10.1103/physrevx.9.031041.

265 C. Roberts, A. Milsted, M. Ganahl, A. Zalcman, B. Fontaine, Y. Zou, J. Hidary, G. Vidal,  
266 and S. Leichenauer, “TensorNetwork: A library for physics and machine learning,” (2019),  
267 [arXiv:1905.01330 \[physics.comp-ph\]](#) .

268 J.-G. L. Hao Xie and L. Wang, [arXiv:2001.04121](#) .

269 G. H. Golub and C. F. Van Loan, *Matrix computations*, Vol. 3 (JHU press, 2012).

270 M. Seeger, A. Hetzel, Z. Dai, E. Meissner, and N. D. Lawrence, “Auto-differentiating linear  
271 algebra,” (2017), [arXiv:1710.08717 \[cs.MS\]](#) .

272 Z.-Q. Wan and S.-X. Zhang, “Automatic differentiation for complex valued svd,” (2019),  
273 [arXiv:1909.02659 \[math.NA\]](#) .

274 C. Hubig, “Use and implementation of autodifferentiation in tensor network methods with complex  
275 scalars,” (2019), [arXiv:1907.13422 \[cond-mat.str-el\]](#) .

276 L. Hascoet and V. Pascual, [ACM Transactions on Mathematical Software \(TOMS\)](#) **39**, 20 (2013).

277 J. Utke, U. Naumann, M. Fagan, N. Tallent, M. Strout, P. Heimbach, C. Hill, and C. Wunsch,  
278 [ACM Trans. Math. Softw.](#) **34** (2008), 10.1145/1377596.1377598.

279 M. Innes, “Don’t unroll adjoint: Differentiating ssa-form programs,” (2018),  
280 [arXiv:1810.07951 \[cs.PL\]](#) .

281 M. Innes, A. Edelman, K. Fischer, C. Rackauckas, E. Saba, V. B. Shah, and W. Tebbutt,  
282 [CoRR abs/1907.07587](#) (2019), [arXiv:1907.07587](#) .

283 Y. Shen and Y. Dai, [IEEE Access](#) **6**, 11146 (2018).

284 G. Forget, J.-M. Campin, P. Heimbach, C. N. Hill, R. M. Ponte, and C. Wunsch, (2015).

285 K. S. Perumalla, *Introduction to reversible computing* (Chapman and Hall/CRC, 2013).

286 M. P. Frank, [IEEE Spectrum](#) **54**, 32–37 (2017).

287 M. MacKay, P. Vicol, J. Ba, and R. B. Grosse, in  
288 [Advances in Neural Information Processing Systems 31](#), edited by S. Bengio, H. Wallach,  
289 H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (Curran Associates, Inc., 2018) pp.  
290 9029–9040.

291 L. Dinh, D. Krueger, and Y. Bengio, “Nice: Non-linear independent components estimation,”  
292 (2014), [arXiv:1410.8516 \[cs.LG\]](#) .

293 D. Maclaurin, D. Duvenaud, and R. Adams, in [Proceedings of the 32nd International Conference on Machine Learning](#),  
294 Proceedings of Machine Learning Research, Vol. 37, edited by F. Bach and D. Blei (PMLR,  
295 Lille, France, 2015) pp. 2113–2122.

296 J. Behrmann, D. Duvenaud, and J. Jacobsen, [CoRR abs/1811.00995](#) (2018), [arXiv:1811.00995](#) .

297 C. Lutz, “Janus: a time-reversible language,” (1986), *Letter to R. Landauer*.

298 M. P. Frank, *The R programming language and compiler*, Tech. Rep. (MIT Reversible Computing  
299 Project Memo, 1997).

300 I. Lanese, N. Nishida, A. Palacios, and G. Vidal, [Journal of Logical and Algebraic Methods in Programming](#) **100**, 71–97 (2018)

301 T. Haulund, “Design and implementation of a reversible object-oriented programming language,”  
302 (2017), [arXiv:1707.07845 \[cs.PL\]](#) .

303 M. P. Frank and T. F. Knight Jr, *Reversibility for efficient computing*, Ph.D. thesis, Massachusetts  
304 Institute of Technology, Dept. of Electrical Engineering and ... (1999).

305 J. G. Koller and W. C. Athas, in *Workshop on Physics and Computation* (1992) pp. 267–270.

306 R. C. Merkle, R. A. Freitas, T. Hogg, T. E. Moore, M. S. Moses, and J. Ryley,  
307 *Journal of Mechanisms and Robotics* **10** (2018), 10.1115/1.4041209.

308 K. Likharev, *IEEE Transactions on Magnetics* **13**, 242 (1977).

309 V. K. Semenov, G. V. Danilov, and D. V. Averin, *IEEE Transactions on Applied Superconductivity* **13**, 938 (2003).

310 R. Landauer, IBM journal of research and development **5**, 183 (1961).

311 J. Bezanson, S. Karpinski, V. B. Shah, and A. Edelman, arXiv preprint arXiv:1209.5145 (2012).

312 J. Bezanson, A. Edelman, S. Karpinski, and V. B. Shah, *SIAM Review* **59**, 65–98 (2017).

313 D. Speyer and B. Sturmfels, *Mathematics Magazine* **82**, 163 (2009).

314 A. McInerney, *First steps in differential geometry* (Springer, 2015).

315 J. Revels, M. Lubin, and T. Papamarkou, “Forward-mode automatic differentiation in julia,” (2016),  
316 [arXiv:1607.07892 \[cs.MS\]](#) .

317 T. Besard, C. Foket, and B. D. Sutter, *CoRR abs/1712.03112* (2017), [arXiv:1712.03112](#) .

318 F. Srajer, Z. Kukelova, and A. Fitzgibbon, *Optimization Methods and Software* **33**, 889 (2018).