

Why we should develop reversible programming - the machine learning perspective

Jin-Guo Liu,¹ Tai-Ne Zhao,¹ and Lei Wang^{1,2,3,*}

¹*Institute of Physics, Chinese Academy of Sciences, Beijing 100190, China*

²*CAS Center for Excellence in Topological Quantum Computation,
University of Chinese Academy of Sciences, Beijing 100190, China*

³*Songshan Lake Materials Laboratory, Dongguan, Guangdong 523808, China*

This paper considers instruction level differential programming, i.e. knowing only the backward rule of basic instructions like +, -, * and /, differentiate a program with proper performance. We will review briefly why instruction level automatic differentiation is hard for current machine learning package even for a source to source automatic differentiation package. Then we propose a reversible Turing machine implementation to achieve instruction level automatic differentiation.

I. AUTOMATIC DIFFERENTIATION

There are two basic modes of automatic differentiation [?], the tangent mode[?] and the adjoint mode. Consider a multi-in (\vec{x}) multi-out (\vec{y}) function f , the tangent mode computes a column of its Jacobian $\frac{\partial \vec{y}}{\partial x_i}$ efficiently, where x_i is single input variable and \vec{y} is multiple output variables. Whereas the adjoint mode computes a row of Jacobian $\frac{\partial y_i}{\partial \vec{x}}$ efficiently.

Most popular automatic differentiation package implements the adjoint mode differentiation, because they are computational more efficient in most optimization applications, where the output loss is always a scalar. Implementing adjoint mode AD requires tracing a program and its intermediate state backward, which requires storing extra information

1. computational graph,
2. and intermediate result caching.

The computational graph is a DAG that stores function calls from inputs to results. Intermediate results are usually the input variable of a function, it is necessary to compute the adjoint of the function. In Pytorch [1] and Flux [?], every variable (tensor) has a tracker field that stores its parent (data and function that generate this variable) information in computational graph and intermediate state. TensorFlow [2] implements a static computational graph before actual computing happens. Source to source automatic differentiation package Zygote [3] use an intermediate representation SSA as the computational graph, so that it can back propagate over a native julia code. Still, intermediate caching is necessary.

Since every computational process is compiled to instructions, these instructions are a natural sequential computational graphs. These instructions are from a finite set of '+', '-', '*', '/', conditional jump statements et. al. With instruction level computational graph, we do not need to define primitives like `exp`, or even linear algebras functions like singular value decomposition [] and eigenvalue decomposition. where the manually derived backwards rule still faces the degenerate spectrum problem (gradients explodes), instruction level AD will return reasonable gradients. With instruction level AD,

people don't worry about inplace functions, which may be a huge problem in traditional approaches. We can back propagate over a quantum simulator, where all instructions are reversible two level unitaries (i.e. Jacobian rotation). We don't need extra effort to learn meta parameters. [] Neural ODE is much easier to design [4].

However, people don't use instructions computational graph for practical reasons. The cost of memorizing the computational graph and intermediate caching kills the performance for more than two orders (as we will show latter). A even more serious problem is the memory consumption for caching intermediate results increases linearly as time. Even in many traditional deep network like recurrent neural network and residual neural networks, where the depth is only several thousand, this memory cost can be a nightmare.

In this paper, we introduce a high performance instruction level AD by making a program time reversible. Making use of reversibility like information buffer [5] can reduce the memory allocations in recurrent neural network [6] and residual neural networks [7]. However, the use of reversibility in these cases are not general purposed. We develop a domain specific language (DSL) in Julia that implements reversible Turing machine. In the past, the reversible Turing machine is not a widely used computational model for having either polynomial overhead in computational time or additional memory cost that proportional to computational time. There has been some prototypes of reversible languages like Janus, where reversible control flows are introduced. Our DSL borrows the design of reversible control flow, meanwhile provides abstraction and memory management. With these additional features, differentiating over a general program requires less than 100 lines. In this paper, we show that in many useful applications, the computational overhead is just a constant factor. Only in some worst cases, it is equivalent to a traditional machine learning framework that cache every input.

* wanglei@iphy.ac.cn

A. Computing Jacobians, Hessians and more

B. First order gradient

Given a node $\vec{y} = f(\vec{x})$ in a computational graph, we can propagate the Jacobians in tangent mode like

$$J_{x_i}^{\vec{O}} = \sum_j J_{y_j}^{\vec{O}} J_{x_i}^{y_j} \quad (1)$$

and the adjoint mode

$$J_{\vec{I}}^{y_j} = \sum_i J_{x_i}^{y_j} J_{\vec{I}}^{x_i} \quad (2)$$

Here, \vec{I} is the inputs and \vec{O} is the outputs. This result can be obtained diagrammatically.

With reversible programming, f^{-1} is used for uncomputing intermediate results, $J_{\vec{x}}^{\vec{y}} \equiv (J_{\vec{y}}^{\vec{x}})^{-1}$ can also be obtained during uncomputing.

C. Second order gradient

D. Gradient on ancilla problem

Ancilla can also carry gradients during computation, sometimes these gradients can not be uncomputed even if their parents can be uncomputed regoriously. In these case, we simply “drop” the gradient field instead of raising an error. In this subsection, we prove doing this is safe, i.e. does not have effect on rest parts of program.

Ancilla is the fixed point of a function, which means

$$b, y \leftarrow f(x, a), \text{ where } b == a \quad (3)$$

$$\frac{\partial b}{\partial x} = 0$$

During the computation, the gradient field does not have any effect to the value. The key question is will the loss of gradient part in ancilla affect the reversibility of the gradient part of argument variables. The gradient of argument variable is defined as $\frac{\partial L}{\partial x} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial x} + \frac{\partial L}{\partial b} \frac{\partial b}{\partial x}$, where the second term vanish naturally.

E. Taylor propagation

The hessian back propagation rule is

$$H_{y_L, y'_L}^f = 0 \quad (4)$$

$$H_{y_{i-1}, y'_{i-1}}^f = J_{y_{i-1}}^{y_i} H_{y_i, y'_i}^f J_{y'_i}^{y'_i} + J_{y_i}^f H_{y_{i-1}, y'_{i-1}}^{y_i}$$

F. A new training scheme

We compute the output as

$$y = f(x), \quad (5)$$

and we have an expected output \hat{y} . In traditional machine learning, we define a loss and minimize it. However, this is not how human brain works. [?] Since the \hat{y} is different with y , the network start to “think”, if the output is \hat{y} , what should the input (including network parameters) be? So, we feed \hat{y} back to the output end of network, and inverse the tape, with the runtime information generated by the input image. The network will compute a different value to network parameters, then the network feel “confused”, and chose a value between old and new values.

Let’s look at the folloing example

```
function f(x::Ref, y::Ref)
    add!(y, x)
    exp!(x)
    neg!(x)
    exp!(y)
    add!(x, y)
    return x, y
end

function invf(x::Ref, y::Ref)
    sub!(x, y)
    log!(y)
    neg!(x)
    log!(x)
    sub!(y, x)
    return x, y
end

x, y = Ref(0.3), Ref(1.6)
for i=1:100
    x = Ref(0.3)
    f(x, y)
    x[] = 1.0
    invf(x, y)
end
```

$f(x, y)$ computes $e^{(y+x)} - e^x$ and stores the output in x , the target is to find a y that make the output x equal to the target value 1. After 200 steps training, y runs into the fixed point and x is equal to 1 upto machine precision.

This training is similar the recursive convergence method that widely used in mathematics and physics. However, in programming language but it has a loophole, it is vulnarable to loss of injectivity. For example, if the result is accumulated to another variable $z+ = x$, and take z as the loss, then the loss can not accumulated to the target value correctly since the inverse of above operation is $z- = x$ and x is unchanged in the inverse run. This is because the redundancy introduced in operation $z = z + e^x$, which keeps both x and z , the change of z shifting z directly is apparently a lazier approach to fit z with target value 1, to figure out the correct causality, corresponding change in output side to x and y are

required. This redundancy can be detected by computing the entanglement entropy (or mutual information) between two output datas. We call this type of invertibility inference, which should be avoided in this type of training. Then how shall we calculate $f(x) = e^x$ like functions, which are not partly invertible from the numeric perspective? Here we introduce the notion of "weak invertibility"

$$@asserty == 0 \quad (6)$$

$$f1(y)+ = \exp(f1(x)) \quad (7)$$

$$f1(x)- = \log(y) \quad (8)$$

We can keep x to ensure numerical invertibility, meanwhile erase most informations inside it. Since x is effectively 0, thus no redundancy (or entanglement entropy) in output, in this way, the above training is still valid.

G. Compactness of data

Measure of compactness by entanglement entropy. In invertible programming, we define the entanglement between output space garbage space as the compactness of information.

If the output information is compact, then we have $(x, x_a) \rightarrow (y, y_a)$ mutual information is

$$I(x, y) = S(x, y) - S(x) - S(y) \quad (9)$$

$$S(x, x_a) = S(y, y_a) \quad (10)$$

$$I(x, x_a) = 0 \quad (11)$$

$$I(y, y_a) = S(y, y_a) - S(y_a) \quad (12)$$

II. TIME SPACE TRADEOFF

Comparing with irreversible programming, reversible programming consumes more memory. Data have different life cycles, some are more persistent, like database, some are transient like an ancilla variable in register that live only in several clock cycle.

III. A INVERTIBLE PROGRAMMING STYLE

A. How does a general program work

In a modern programming language, functions are pushed to a global stack for scheduling. The memory layout of a function is consisted of input arguments, a function frame with informations like return address and saved memory segments, local variables and working stack. After each call, the function clears the input arguments, function frame, local variables and working stack and only stores the return value. In the invertible programming style, this kind of design pattern is no longer the best practise, the local information of a

function can not be emptied immediately after a function call to guarantee the invertibility.

Programming language Janus [8] is a prototype of reversible programming that can be compiled to invertible instructions.

B. π calculus

C. For Loop

D. Unitary Matrices

Recurrent networks with a unitary parametrized network ease the gradient exploding and vanishing problem [9–11]. Among different parametrization schemes, the most elegant one is Ref. 11, which parametrized the unitary matrix with two-level unitary operations, any unitary matrix of size $N \times N$ can be parametrized by $k = N(N - 1)/2$ two level unitary matrices [12]. All these two-level unitary matrices can be applied in $O(1)$ time as a two register instruction. Hence a real unitary matrix can be parametrized compactly by k rotation angles, each represents a rotation operations between datas in two target parameters.

```
@assert j==0, i==0, k==0, ip==0
loop N-1
  j + 1
  loop j
    i + 1
    k + 1
    ip + i + 1
    rot(x[i], x[ip], z[k])
    ip - 1 - i
  end
  i - j
end
```

Here, loop j means excuting the loop body for j times, it gurantes invertibility. Loop and branching can be implemented in a more regorious way [13] from instruction level, however it is still an open question how to implelement loops with instruction level invertibility without compiling technic. In the rotation instruction, $z[k]$ is the rotation angle, which represents rotating data in target registers by an angle of $\theta = z[k] * \pi$.

$$R(\theta) = \begin{bmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{bmatrix} \quad (13)$$

1. The gradient

Its backward rule is

$$\begin{aligned}\bar{\theta} &= \sum \frac{\partial R(\theta)}{\partial \theta} \odot (\bar{y}x^T) \\ &= \text{Tr} \left[\frac{\partial R(\theta)}{\partial \theta} \bar{y}x^T \right] \\ &= \text{Tr} \left[R \left(\frac{\pi}{2} - \theta \right) \bar{y}x^T \right]\end{aligned}\quad (14)$$

The resulting instruction is

```
-z[k]
z[k] + pi/2
rot(gx[i], gx[ip], theta[k])
gtheta + gx[i]*x[i]
gtheta + gx[ip]*x[ip]
theta[k] - pi/2
-theta[k]

rot(gx[i], gx[ip], -pi/2)
```

In this way, the gradient is defined in a closure form.

2. Rounding errors

$$a = a \cos(k\pi) - b \sin(k\pi) \quad (15)$$

$$b = a \sin(k\pi) + b \cos(k\pi) \quad (16)$$

$$\begin{aligned}\epsilon_a &= -a \sin(\theta)\theta\epsilon - b \cos(\theta)\theta\epsilon + a\epsilon \cos(\theta) - b\epsilon \sin(\theta) \\ &\sim \max(a, b)\epsilon\end{aligned}\quad (17)$$

The error accumulates linearly as the number of floating point operations, for a $N \times N$ unitary matrix, each number is operated N times. In a double precision computation, the rounding errors in unitary matrix multiplication will probably not have a substantial effect on resersibility.

3. Two level unitary for rectangular matrices

E. Invertible random numbers

IV. INSTRUCTIONS

Not only invertibility, but also the stability of gradient itself, requires reversible instruction support, otherwise invertibility can be easily ruined by rounding errors. Using information buffer in multiplication operations [5] in an approach to enforce invertibility in a memory efficient way. Invertibility has been studied in the cross field of computer science and physics a lot between 1980 and 2010. The motivation is saving energy. Carlin devised SRCL logic family from Pendulum instruction set architecture (PISA) [13] for a invertible programming device.

V. SPEUDO INVERSE

In many cases, for example the $\exp(x)$ operation, due to the lack of reversible multiplication instruction in a classical computer instruction set, we have to either use a list to bookmark intermediate results, or using uncomputing technic that cost polynomial or even exponential more times to limit the memory usage to a constant.

For these arithmetics, here we introduce a trick that computes these functions in constant memory with a constant uncomputing overhead.

Let's see the following Taylor serie implementation of $\exp(x)$ in a reversible fashion

```
using NiLang, NiLang.AD

@i function iexp(out!, x::T; atol::Float64=1e-14)
  where T
  @anc anc1::T
  @anc anc2::T
  @anc anc3::T
  @anc iplus::T

  out! += 1.0
  anc1 += 1.0
  while (val(anc1) > atol, !isapprox(iplus, 0.0))
    iplus += 1.0
    anc2 += anc1 * x
    anc3 += anc2 / iplus
    out! += anc3
    # pseudo inverse
    anc1 -= anc2 / x
    anc2 -= anc3 * iplus
    SWAP(anc1, anc3)
  end

  ~(while (val(anc1) > atol, !isapprox(iplus, 0.0))
    iplus += 1.0
    anc2 += anc1 * x
    anc3 += anc2 / iplus
    # pseudo inverse
    anc1 -= anc2 / x
    anc2 -= anc3 * iplus
    SWAP(anc1, anc3)
  end)
  anc1 -= 1.0
end
```

The two lines bellow the comment `# pseudo inverse` "uncompute" variables `anc1` and `anc2` to a value very close to zero. Notice `*` and `/` are not exactly dual to each other. It is reasonable to assume this inexact uncomputing will cause negligible error on final output, but harms reversibility. In the latter case, error accumulates in the whole program. In the second for loop inside the inverse notation `~`, we uncompute all ancilla bits rigorously. The while statement takes two conditions, the precondition and postcondition. Precondition `val(anc1) > atol` indicates when to break the forward pass and post condition `!isapprox(iplus, 0.0)` indicates when to break the backward pass.

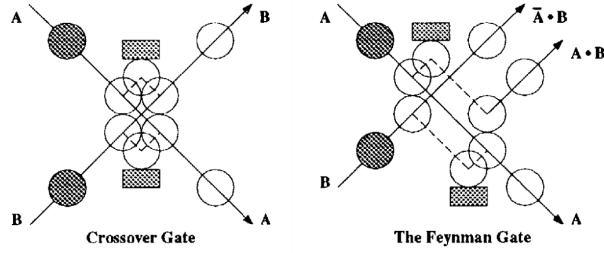


Figure 1. Two examples of BBM logic gates.

VI. HARDWARES

A. Traditional Invertible Computing Devices

A reversible computer, by no means refer to a computing device that every instruction or program is reversible. A better definition would be, a reversible computing device reserves the right to retract energy through uncomputing.

Like quantum computer, it should be able to reset to qubits 0. The reset operation, sometimes can be expensive. In a quantum device, this reset operation can be done by dissipation and decoherence (i.e. interact with environment and get thermalized). Alternatively, a immediate feedback loop can be introduced to a quantum device, where a classical computer is regarded as the dissipation source.

1. Pendulum Design

copy pasted begin

It is possible to design and fabricate a fully reversible processor using resources which are comparable to a conventional microprocessor. Existing CAD tools and silicon technology are sufficient for reversible computer design. Pendulum author demonstrated this by designing such a processor and fabricating and testing it in a commercially available CMOS process.

copy pasted end

2. Margolus's Billiard Ball Model Cellular Automaton (BBMCA)

The chip is known as Flattop, not very convenient to program, but is simple, universal, reversible, and scalable.

B. Quantum computers

Fully entangled quantum systems are not easy to prepare. For a quantum system, finding a reliable dissipation source is also not easy, this is why the time for resetting a qubit to 0 (or erasing a bit information) can be the bottleneck of a computing.

VII. DISCUSSION

One should notice the memory advantage of reversible programming to machine learning does comes from reversibility itself, but from a better data tracking strategy inspired from invertible programming. Normally, a reversible program is not as memory efficient as its irreversible counterpart due to the additional requirement of no information loss. A naive approach that keeping track of all informations will cost an additional space $O(T)$, where T stands for the execution time in a irreversible TM, the longer the program runs, the larger the memory usage is. This is exactly the approach to keeping reversibility in most machine learning packages in the market. The point it, an reversible Turing Machine is able to trade space with time. In some cases, it may cause polynomial overhead than its irreversible counterpart.

In the simplest g -segment trade off scheme [], it takes $Time(T) = T^{\log_g(4g-2)}$ and $Space(T) = (g-1)S \log_g T$. In practise, there are more practical trading off schemes that works much better in practise. [] Checkpointing [14].

ACKNOWLEDGMENTS

The authors are supported by the National Natural Science Foundation of China under the Grant No. 11774398, the Strategic Priority Research Program of Chinese Academy of Sciences Grant No. XDB28000000 and the research funding from Huawei Technologies under the Grant No. YBN2018095185.

-
- [1] 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).
 - [2] 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.

- [3] M. Innes, A. Edelman, K. Fischer, C. Rackauckas, E. Saba, V. B. Shah, and W. Tebbutt, [CoRR abs/1907.07587](#) (2019), [arXiv:1907.07587](#).
- [4] T. Q. Chen, Y. Rubanova, J. Bettencourt, and D. Duvenaud, [CoRR abs/1806.07366](#) (2018), [arXiv:1806.07366](#).
- [5] D. Maclaurin, D. Duvenaud, and R. Adams, in *Proceedings of the 32nd International Conference on Machine Learning*, Proceedings of Machine Learning Research, Vol. 37, edited by F. Bach and D. Blei (PMLR, Lille, France, 2015) pp. 2113–

2122.

- [6] M. MacKay, P. Vicol, J. Ba, and R. B. Grosse, in *Advances in Neural Information Processing Systems 31*, edited by S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (Curran Associates, Inc., 2018) pp. 9029–9040.
- [7] J. Behrmann, D. Duvenaud, and J. Jacobsen, *CoRR abs/1811.00995* (2018), arXiv:1811.00995.
- [8] C. Lutz, “Janus: a time-reversible language,” (1986), *Letter to R. Landauer*.
- [9] M. Arjovsky, A. Shah, and Y. Bengio, *CoRR abs/1511.06464* (2015), arXiv:1511.06464.
- [10] S. Wisdom, T. Powers, J. R. Hershey, J. L. Roux, and L. Atlas, “Full-capacity unitary recurrent neural networks,” (2016), arXiv:1611.00035 [stat.ML].
- [11] L. Jing, Y. Shen, T. Dubcek, J. Peurifoy, S. A. Skirlo, M. Tegmark, and M. Soljacic, *CoRR abs/1612.05231* (2016), arXiv:1612.05231.
- [12] C.-K. LI, R. ROBERTS, and X. YIN, *International Journal of Quantum Information* **11**, 1350015 (2013).
- [13] C. J. Vieri, *Reversible Computer Engineering and Architecture*, *Ph.D. thesis*, Cambridge, MA, USA (1999), aAI0800892.
- [14] T. Chen, B. Xu, C. Zhang, and C. Guestrin, *CoRR abs/1604.06174* (2016), arXiv:1604.06174.

1. NiLang Grammar

Appendix A: NiLang

- Terminals in quotes or lower case.
- ident, symbols
- num, numbers
- null, empty statement
- expr, native julia expression.
- expr::Int, native julia expression with integer return value.
- expr::Bool, native julia expression with boolean return value.
- { }* indiceates zero or more repetitions.
- [] indicates zero or one repetitions.

```

Function ::= 'function' Fname '(' Fargs [',' ident
          '...']; Fargs ')' 'where' '{' Symbols '}'
          Statements
          'end'

Fname ::= ident
        | '(' ident '::' ident ')'

Fargs ::= ident, Fargs
        | ident '::' ident, Fargs
        | ident '=' ident, Fargs
        | ident '::' ident '=' expr, Fargs
        | null

Symbols ::= ident, Symbols
          | ident

Statements ::= Ifstmt Statements
            | Whilestmt Statements
            | Forstmt Statements
            | Callstmt Statements
            | Instrstmt Statements
            | Revstmt Statements
            | @anc Statements
            | @routine Statements
            | @safe Statements
            | null

Ifstmt ::= 'if' '(' expr::Bool ',' expr::Bool ')'
          Statements
          ['else'
           Statements]
          'end'

Whilestmt ::= 'while' '(' expr::Bool ',' expr::Bool
              ')'
              Statements
              'end'

Forstmt ::= 'for' ident '=' expr::Int[':' expr::Int]
           ':' expr::Int
           Statements
           'end'

Callstmt ::= ident '(' Dataviews ')'

Instrstmt ::= ident '+=' ident ['(' Dataviews ')']
            | ident '-=' ident ['(' Dataviews ')']
            | ident '∇=' ident ['(' Dataviews ')']
            | ident '.+=' ident ['.' '(' Dataviews ')']
            | ident '.-=' ident ['.' '(' Dataviews ')']
            | ident '∇.=' ident ['.' '(' Dataviews
            ')']

Revstmt ::= '~' '(' Statements ')'

@routine ::= '@routine' ident 'begin'
            Statements
            'end'
            | '@routine' ident

@anc ::= '@anc' ident '=' expr
        | '@deanc' ident '=' expr

@safe ::= '@safe' expr

Dataviews ::= Dataview, Dataviews
            | null

Dataview ::= Dataview '[' expr::Int ']'
            | ident {'.' ident}*
            | ident '(' Dataview ')'
            | Constant

Constant ::= num | 'π'

```

Dataview is a special surjective mapping of parent data, e.g. a field of an object. The dataview can feedback to parent data with the `chfield` method, so that the modified object can generate desired dataview.

Appendix B: Julia based DSL implementation details

Macro `‘invfunc’` defines a invertible function, ancillas are binded to a function, since ancillas are uncomputed to 0 at the end of call, so that it can be used repeatedly in a function, it is like a class variable in a class, with no side effects.

Some variables can be uncomputed to 0, but we choose not to for performance reason. For example `infer!(argmax, i, x)` which computes the location of maximum value in x and store the output to i , if we uncompute it, it doubles additional computational time. Here, we trade off the memory with computation time. As a result, we must feed `imax` to the function, so that this variable can be manipulated in outer scopes.