

Instruction level automatic differentiation on a reversible Turing machine

Jin-Guo Liu^{1,*} and Taine Zhao²

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

²*Department of Computer Science, University of Tsukuba*

This paper considers the instruction-level adjoint mode automatic differentiation (AD). Here, instruction-level means, given the backward rules of basic instructions like $+$, $-$, $*$ and $/$, one can differentiate an arbitrary program efficiently. In this paper, we review why instruction-level AD is hard for traditional machine learning frameworks and propose a solution to these problems by back-propagating a reversible Turing machine. Instruction level AD is a powerful tool to generate backward rules. We demonstrate its power by differentiating over the reversible implementations of exp function, and linear algebra functions unitary matrix multiplication and QR decomposition. It is also a promising direction towards solving the notorious memory wall problem in machine learning. Also, we discuss the challenges that we face towards rigorous reversible programming from the instruction and hardware perspective.

I. INTRODUCTION

Computing the gradients of a numeric model $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$ plays a crucial role in scientific computing. Automatic differentiation (AD) is the key technique to obtain gradients automatically by propagating gradients with the chain rule. The propagation of gradients can happen in two directions. Consider a program $\vec{y} = f(\vec{x})$, a tangent mode AD propagate the gradient forward and computes one column of its Jacobian $\frac{\partial \vec{y}}{\partial x_i}$ efficiently, where x_i is one of the input variables, whereas an adjoint mode AD propagates the gradients backward and computes one row of Jacobian $\frac{\partial y_i}{\partial \vec{x}}$ efficiently. [1] For this reason, the adjoint mode AD is preferred in variational applications, where the loss as output is always a scalar. However, implementing adjoint mode AD is harder than implementing its tangent mode counterpart, because it requires a program's intermediate information for back-propagation. This information includes

1. the computational graph,
2. and variables used in back propagation.

A computational graph is a directed acyclic graph (DAG) that records the relationship between data (edges) and functions (nodes). In Pytorch [2] and Flux [3], every variable has a tracker field that stores its parent information, i.e., the input data and function that generate this variable. TensorFlow [4] implements a static computational graph as a description of the program before actual computation happens. Source to source AD package Zygote [5, 6] uses an intermediate representation (IR) of a program, the static single assignment (SSA) form, as the computational graph. To cache intermediate states, it has to access the global storage.

Several limitations are observed in these AD implementations due to the recording and caching. First of all, these packages require a lot of primitive functions with programmer-defined backward rules. These backward rules are not necessary given the fact that, at the lowest level, these

primitive functions are compiled to a finite set of instructions including $+$, $-$, $*$, $/$, and conditional jump. By defining backward rules for these basic instructions, AD should work without extra effort. These machine learning packages do not use instructions as the computational graph, because the overhead of memorizing the computational graph and caching intermediate states is significant. The memory overhead for caching inputs of every instruction is a linear to the computing time. In many deep learning models like recurrent neural network [7] and residual neural networks [8], the depth can reach several thousand. The memory is often the bottleneck of these programs, which is known as the memory wall problem. [9] Secondly, inplace functions are not supported in these packages. This limitation comes the requirement of caching intermediate states. If a data is allowed to change inplace, then the cached data would not be safe anymore. Even if Zygote uses SSA as the computational graph, it is not trivial to handle inplace functions properly. On the other side, most functions in BLAS and LAPACK are implemented as inplace functions. As a result, all AD packages that use BLAS and LAPACK functions have to define their own backward rules for their non-inplace wrappers due to the lack of AD support to inplace functions. Thirdly, obtaining higher-order gradients are not efficient in these packages. For example, in most machine learning packages, people back-propagate the whole program of obtaining first-order gradients to obtain second-order gradients. The repeated use of back-propagation algorithm causes an exponential overhead concerning the order of gradients. A better approach is using Taylor propagation like in JAX [10]. However, Taylor propagation requires tedious work of implementing higher order backward rules for all primitives.

Our solution to these issues is making a program reversible. It is not the same as the workarounds in the machine learning field that use information buffer [11] and reversible activation functions [12, 13] to reduce the memory allocations in recurrent neural networks [14] and residual neural networks [15]. Our approach is general purposed. We develop NiLang, an embedded domain-specific language (eDSL) in Julia language [16, 17] that implements a reversible Turing machine (RTM). [18, 19]. This eDSL contains a macro that generates reversible functions, which is completely

* cacate0129@iphy.ac.cn

compatible with Julia ecosystem. It provides reversible control flows, instructions and memory managements, meanwhile provides multiple-dispatch based abstraction. With these features and less than 100 lines of code, we implemented the AD engine that can differentiate any thing, including linear algebra functions.

NiLang is not the first implementation of RTM. In history, there have been some prototypes of reversible languages like Janus [20], R (not the popular one) [21], Erlang [22] and object-oriented ROOPL [23]. These languages have reversible control flow that allows user to input an additional postcondition in control flows to help programs run backward. In the past, the primary motivation for making a program time-reversible is to support reversible devices. Reversible devices do not erase information hence do not have a lower bound of energy consumption by Landauer principle [24]. However, people show less and less interest to reversible programming since 15 years ago, because the energy efficiency of complementary metal–oxide–semiconductor (CMOS) devices are still two orders [19] above this lower bound, removing this lower bound is not an urgent problem yet. Our work breaks the information barrier between the machine learning community and the reversible programming community, and provides yet another strong motivation to develop reversible programming.

In this paper, we first introduce the language design of NiLang in Sec. II. In Sec. III, we explain the algorithm to back-propagation Jacobians and Hessians in this eDSL. In Sec. IV, we show several examples including Fibonacci number, \exp function, unitary matrix multiplication and QR decomposition [25]. We show how to generate first order and second order backward rules for these functions. In Sec. V, we discuss several important issues, how time-space tradeoff works, reversible instructions and hardware, and finally, an outlook to some open problems to be solved. In the appendix, we show the grammar of NiLang and a gradient free self-consistent training strategy.

II. LANGUAGE DESIGN

A. Introductions to reversible language design

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 information like the return address and saved memory segments, local variables, and working stack. After the call, the function clears these runtime information, only stores the return value. In the reversible programming style, this kind of design pattern is no longer the best practice. One can not discard input variables and local variables easily after a function call, since discarding information may ruin reversibility. Hence reversible functions are very differential from irreversible ones from multiple perspectives.

First of all, the memory management in a reversible language is special. A variable can not be assigned or discarded directly in a reversible language. There are two basic types

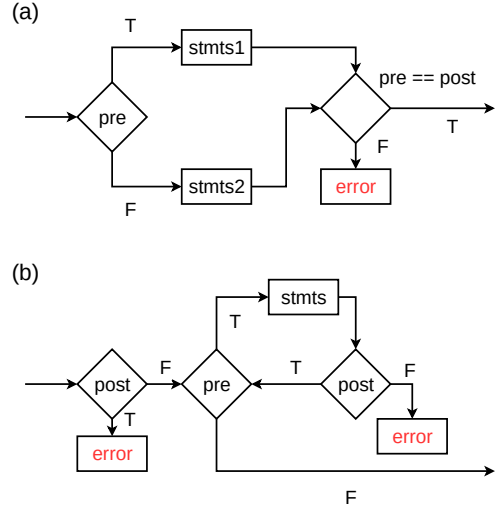


Figure 1. Flow chart for reversible (a) if statement and (b) while statement. “stmts”, “stmts1” and “stmts2” are statements, statements in true branch and statements in false branch respectively. “pre” and “post” are precondition and postconditions respectively.

of memory management in traditional reversible languages, the ancilla and the stack. In Janus, ancillas are handled by a statement pair of `local` and `delocal`. `local` initializes a variable to a specific value, while `delocal` deallocates that variable with a preassumed value. Stack operations includes `push` and `pop` statements, it is similar to a traditional stack except it zero-clears the variable after pushing and presupposes the variable being zero-cleared before popping. To summarize, reversible program deallocates a variable only when its state is known.

Secondly, the control flows are special. The `if` statement is shown in Fig. 1 (a), the program enters the branch specified by precondition. After executing this branch, the program checks the consistency of precondition and postcondition to make sure they are the same. In the reverse pass, the program enters the branch specified by the postcondition. For the `while` statement shown in Fig. 1 (b), before entering, the program preassumes the postcondition is false. After each iteration, the program asserts the postcondition to be true. The inverse function exchanges the precondition and postcondition. The definition of reversible `for` statement is similar to irreversible ones. Before entering the loop, the program stores variables start, step, and stop locally. After executing the loop, the program checks the values of these variables to make sure they are not changed. The reverse program exchanges start and stop and inverse the sign of step.

Lastly, the basic arithmetic and boolean instructions are also different. Every instruction has a unique inverse that can undo the changes. For example $z += x * y$ is legal reversible instruction for integers, with its inverse $z -= x * y$, while $x *= y$ is not. In the following discussion, we assume $x += y$ is reversible for floating point numbers although it is not true considering rounding error. A not very well known fact is $+$, $-$, $*$, $/$ instructions can be implemented

reversibly for both integers and floating point numbers, it is only a matter of decision. We will come to this point later in Sec. VB

B. NiLang’s Reversible IR

In our paper, we want use a reversible program as the computational graph. To obtain the gradient, we need to insert the code of obtaining gradients into the reversed program. The reversible language design should have related abstraction power to provide polymorphism on the same instruction sequence. This motivates us to design a new reversible language that can fit this task.

The grammar of NiLang is shown in Appendix A. Julia language’s meta programming and its package MLStyle [26] allow user to define an eDSL to explain the reversible statements conveniently. The type inference and just in time compiling can remove most overheads in our eDSL, providing a comparable performance with respect to its irreversible counterpart. Most importantly, the multiple dispatch in Julia provides the polymorphism that allowing us to obtain the gradients on the reversed program. The main feature of NiLang is contained in a single macro `@i`. The following example shows a minimal working example of compiling a NiLang function to native julia function.

```
julia> using NiLangCore, MacroTools

julia> macroexpand(Main, :(@i function f(x, y)
    SWAP(x, y)
end)) |> prettify

quote
  $(Expr(:meta, :doc))
  function $(Expr(:where, :(f(x, y))))
    gaur = SWAP(x, y)
    x = (NiLangCore.wrap_tuple(gaur))[1]
    y = (NiLangCore.wrap_tuple(gaur))[2]
    return (x, y)
  end
  if typeof(f) != typeof(~f)

    function $(Expr(:where, :(( # $ TODO: remove this comment
      mongoose::typeof(~f))(x, y)))
      mandrill = (~SWAP)(x, y)
      x = (NiLangCore.wrap_tuple(mandrill))[1]
      y = (NiLangCore.wrap_tuple(mandrill))[2]
      return (x, y)
    end
  end
  if ! (NiLangCore._hasmethod1(
    NiLangCore.isreversible, typeof(f)))
    NiLangCore.isreversible(::typeof(f)) = true
  end
end
```

Macro `@i` generates two functions `f` and `~f`, where `~f` is an callable of type `Inv{typeof(f)}`. `NiLangCore.wrap_tuple` is used to unify output data types, it will wrap any non-tuple variable to a tuple. The output of an instruction or a function call will be assigned

back to input variables, this requires the length output variables same as input variables. In NiLang, this is always true because input variables of a function will be returned automatically, i.e. a reversible function defined by NiLang is a mapping between the same set of variables. At last, `isreversible` is binded to `f` to mark it as reversible.

To design such an eDSL, we first introduce a reversible IR that plays a central role in NiLang. In this IR, a statement can be an instruction, a function call, a control flow, a macro call, or the inverse statement `~`. Any statement in this IR has a unique inverse as shown in Table I.

statement	inverse
<code><f>(<args>...)</code>	<code>(~<f>)(<args>...)</code>
<code><y!> += <f>(<args>...)</code>	<code><y!> -= <f>(<args>...)</code>
<code><y!> .+= <f>(<args>...)</code>	<code><y!> .-= <f>(<args>...)</code>
<code><y!> ∇= <f>(<args>...)</code>	<code><y!> ∇= <f>(<args>...)</code>
<code><y!> .∇= <f>(<args>...)</code>	<code><y!> .∇= <f>(<args>...)</code>
<code><a> ← <expr></code>	<code><a> → <expr></code>
<code>(<T1> => <T2>)(<x>)</code>	<code>(<T2> => <T1>)(<x>)</code>
<code>begin</code> <code><stmts></code> <code>end</code>	<code>begin</code> <code>~(<stmts>)</code> <code>end</code>
<code>if (<pre>, <post>)</code> <code><stmts1></code> <code>else</code> <code><stmts2></code> <code>end</code>	<code>if (<post>, <pre>)</code> <code>~(<stmts1>)</code> <code>else</code> <code>~(<stmts2>)</code> <code>end</code>
<code>while (<pre>, <post>)</code> <code><stmts></code> <code>end</code>	<code>while (<post>, <pre>)</code> <code>~(<stmts>)</code> <code>end</code>
<code>for <i>=<m>:<s>:<n></code> <code><stmts></code> <code>end</code>	<code>for <i>=<m>:-<s>:<n></code> <code>~(<stmts>)</code> <code>end</code>
<code>@safe <expr> x</code>	<code>@safe <expr></code>

Table I. A collection of reversible statements.

The conditional expression in a `if` or a `while` statements is a two-element tuple that consists of a precondition and a postcondition. `←` and `→` plays the same role as `local` and `delocal` in Janus, they can be input by typing “\leftarrow” followed by a Tab key in a Julia editor or REPL. `a ← <expr>` binds variable `a` to an initial value specified by `<expr>`. Its inverse `a → <expr>` deallocates the variable `a`. Before deallocating the variable, the program checks that the its value is the same as the value of `<expr>`, otherwise throws an `InvertibilityError`. `←` and `→` must appear in pairs inside a function call, a `while` statement, or a `for` statement. A function call in NiLang supports Julia’s broadcasting magic, a function defined on scalars can be broadcasted to vectors by adding a “.” after the function call. The `@safe` macro allows users to use external statements that do not break reversibility. For example, one can use `@safe @show var` for debugging. There is not stack operations in our language design because we haven’t find a use case that

requires a stack yet.

C. Compiling

The interpretation of a reversible function consists three stages.

The first stage preprocess human inputs to a reversible IR. The preprocessor expands the symbol \sim in the postcondition field of `if` statement to the precondition, adds missing ancilla deallocation operations to ensure \leftarrow and \rightarrow statements appear in pairs and expands `@routine` macro. Here, `@routine <stmt>` records a statement. In preprocessing stage, `~@routine` will be replaced to `~<stmt>` for uncomputing. We will use macro extensively in the "compute-copy-uncompute" design pattern of reversible programming. The following example preprocess an `if` statement to the reversible version.

```
julia> prettify(@code_preprocess if (pre, ~)
      z += x * y
      end)
:(if (pre, pre)
    z += x * y
  else
  end)
```

Here, the symbol \sim means *same as precondition*, it will be expanded at this stage.

The second stage generates the reversed IR according to table Table I.

```
julia> prettify(@code_reverse if (pre, post)
      z += x * y
      else
      z += x / y
      end)
:(if (post, pre)
    z -= x * y
  else
    z -= x / y
  end)
```

The third stage is translating this IR and its inverse to native Julia code. At the end of a function definition, it attaches a return statement that returns variables in the argument list. Now the function is ready to execute on the host language. The following example shows how to compile an `if` statement.

```
julia> prettify(@code_interpret if (pre, post)
      z += x * y
      else
      z += x / y
      end)
quote
  bat = pre
  if bat
    @assignback (PlusEq(*))(z, x, y)
  else
    @assignback (PlusEq(/))(z, x, y)
  end
  @invcheck post bat
end
```

The compiler translates the instruction according to Table II and adds `@assignback` before each instruction and function call statement. The macro `@assignback` assigns the output of a function to the argument list of a function, which will be explained in detail in the next subsection. At the end of the code `@invcheck post bat` is added to check the consistency between preconditions and postconditions to ensure reversibility. This statement will throw an error if target variables `bat` and `post` are not "equal" to each other up to the rounding error.

D. Types and Dataviews

So far, the language design is not too different from a traditional reversible language. To implement the adjoint mode AD, we introduce types and dataviews. The type that used in the reversible context is just a normal Julia type with extra requirements that its constructor is reversible. The inverse of a constructor is called a "destructor", which unpacks data and deallocates derived fields. Data packing and unpacking are implemented in following ancilla statements.

```
x ← new{TX, TG}(x, g)
x → new{TX, TG}(x, g)
```

Ancilla operations with new function is treated specially, because `new` keeps all information in input argument list. The first statement is equivalent to says, do `x = newTX, TG(x, g)` just like a normal ancilla allocation, but do not deallocate `g` because the information of `g` is transferred to `x`. Its inverse the second statement allocates a new ancilla `g`. The following example shows how to define as reversible type `GVar`.

```

julia> using NiLangCore

julia> @i struct GVar{T,GT} <: IWrapper{T}
    x::T
    g::GT

    function GVar{T,GT}(x::T, g::GT) where {T,GT}
        new{T,GT}(x, g)
    end
    function GVar(x::T, g::GT) where {T,GT}
        new{T,GT}(x, g)
    end
    @i function GVar(x::T) where T
        g ← zero(x)
        x ← new{T,T}(x, g)
    end
    @i function GVar(x::AbstractArray)
        GVar.(x)
    end
end

julia> GVar{0.5}
GVar{Float64,Float64}(0.5, 0.0)

julia> (~GVar)(GVar{0.5})
0.5

julia> (~GVar)(GVar([0.5, 0.6]))
2-element Array{Float64,1}:
 0.5
 0.6

```

GVar is a reversible type that copied from the AD submodule of NiLang. It is used to store gradient information of a variable. Here, we only pick two of its constructors. We put @i macro before both struct and function statements. The ones before functions mark reversible functions, while the one before struct keyword moves the definitions of GVar functions to the outside of this type definition.

The cast between two types can be defined conveniently with the macro @icast.

```

julia> @pure_wrapper A

julia> @icast A(x) => GVar(x, g) begin
    g ← zero(x)
    g += identity(1)
end

julia> x = A(0.5)
A(0.5)

julia> @instr (A=>GVar)(x)
GVar{Float64,Float64}(0.5, 1.0)

julia> @instr (GVar=>A)(x)
A(0.5)

```

Here, we first define a simple reversible wrapper type A with macro @pure_wrapper, then define the cast rule between A type and GVar type. The function body is a reversible

program that transforms x to (x, g). The compiler appends a default constructor DVar(xx, gg) at the end of program to instantiate a new object as the return value. Its inverse that converts an object of type GVar to type A is automatically generated by reversing the above statements.

To access and manipulate the fields of a data type, especially an immutable type, we introduce the dataview. A dataview of variable can be the variable itself, a field of its view, an array element of its view, or a bijective mapping of its view. Let us first consider the following example.

```

julia> arr = [GVar(3.0), GVar(1.0)]
2-element Array{GVar{Float64,Float64},1}:
 GVar{Float64,Float64}(3.0, 0.0)
 GVar{Float64,Float64}(1.0, 0.0)

julia> x, y = 1.0, 2.0
(1.0, 2.0)

julia> @instr -arr[2].g += x * y
2.0

julia> arr
2-element Array{GVar{Float64,Float64},1}:
 GVar{Float64,Float64}(3.0, 0.0)
 GVar{Float64,Float64}(1.0, -2.0)

```

-arr[2].g, x and y are all dataviews. In Julia language, -grad(arr[2]) += x * y statement will throw a syntax error because the function call “-” can not be assigned, and GVar is an immutable type. In our eDSL, we wish it works because a memory cell is assumed to be modifiable in our eDSL. The secret of how it works lies in the macro @instr, it translates the above statement to

```

1 res = (PlusEq{*})(-arr[2].g, x, y)
2 arr[2] = chfield(arr[2], Val{:g},
3   chfield(arr[2].g, -, res[1]))
4 x = res[2]
5 y = res[3]

```

The first line PlusEq{*})(-arr[2].g, x, y) computes the output, which is a tuple of length 3. At lines 2-3, chfield(x, Val{:g}, val) modifies the g field of x and chfield(x, -, res[1]) returns -res[1]. Here, modifying a field requires the default constructor of a type is not overwritten. The assignments at lines 4 and 5 are straightforward.

III. AUTOMATIC DIFFERENTIATION

A. First order gradient

Given a node $\vec{y} = f(\vec{x})$ in a computational graph, the adjoint mode AD propagates the Jacobians in the reversed direction

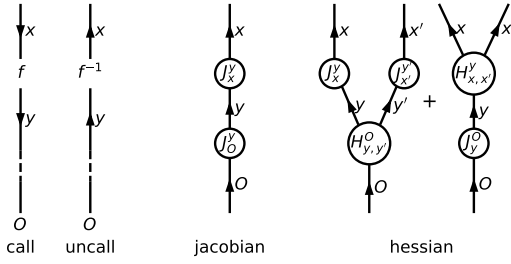


Figure 2. Adjoint rules for Jacobians and Hessians in tensor network language.

like

$$\begin{aligned} J_{O'}^O &= \delta_{O,O'}, \\ J_x^O &= J_y^O J_x^y, \end{aligned} \quad (1)$$

where O represents the outputs of the program, $J_{x/y}^O$ is the gradient to be propagated, and J_x^y is the local Jacobian matrix. The Einstein's notation is used here so that duplicated indices are summed over. This back-propagation rule can be rewritten in the language of tensor networks [27], as shown in Fig. 2, where a circle is a tensor, dangling edges are unpaired labels and connected edges are paired labels.

In reversible programming with multiple-dispatch, we implement the adjoint mode AD as

Algorithm 1: Reversible programming AD

Result: $\text{grad}(\vec{x}_g)$
 let iloss be the index of loss variable in \vec{x}
 $\vec{y} = f(\vec{x})$
 $\vec{y}_g = \text{GVar}(\vec{y})$
 $\text{grad}(\vec{y}_g[\text{iloss}]) += 1.0$
 $\vec{x}_g = f^{-1}(\vec{y}_g)$

Here, GVar is a reversible type. The constructor attaches a zero gradient field to a variable. If the input is an array, GVar will be broadcasted to each array element. One can access the gradient field of a GVar instance through the `grad` dataview. Its inverse `~GVar` deallocates the gradient field safely and returns its value field. Here, "safely" means before deallocation, the program will check the gradient field to make sure its value is restored to 0. When an instruction meets a GVar , besides computing its value field $\text{value}(\vec{y}) = \text{instruct}(\text{value}(\vec{x}))$, it also updates the gradient field $\text{grad}(\vec{y}) = [J_{\vec{x}}^{\vec{y}}]^{-1} \text{grad}(\vec{x})$, where $[J_{\vec{x}}^{\vec{y}}]^{-1}$ is the Jacobian of instruct^{-1} . One can define this gradient function on either `instruct` or `instruct-1`. If one defines the backward rule on `instruct`, the compiler will generate the backward rule for its inverse `instruct-1` as the inverse function. This is doable because the inverse and adjoint operation commutes [28]. This implementation that changes the elementary data type is similar to the use of dual number in tangent mode AD [29]. In the following example, We bind the adjoint function of ROT to its reverse IROT by defining a new function that dispatch to GVar .

```
@i function IROT(a!::GVar, b!::GVar, theta::GVar)
    IROT(value(a!), value(b!), value(theta))
    NEG(value(theta))
    value(theta) -= identity(pi/2)
    ROT(grad(a!), grad(b!), value(theta))
    grad(theta) += value(a!) * grad(a!)
    grad(theta) += value(b!) * grad(b!)
    value(theta) += identity(pi/2)
    NEG(value(theta))
    ROT(grad(a!), grad(b!), pi/2)
end
```

The definition of ROT instruction could be found in Sec. B. This backward rule has been included in NiLang, one can check the gradients by typing in a Julia REPL

```
julia> using NiLang, NiLang.AD

julia> x, y, theta = GVar(0.5), GVar(0.6), GVar(0.9)
(GVar(0.5, 0.0), GVar(0.6, 0.0), GVar(0.9, 0.0))

julia> @instr grad(x) += identity(1.0)

julia> @instr ROT(x, y, theta)

julia> x, y, theta
(GVar(-0.1591911616411577, 0.6216099682706646),
 GVar(0.7646294357761403, 0.7833269096274833),
 GVar(0.8999999999999999, 0.6))
```

The implementation of Algorithm 1 is so short that we present the function definition as follows.

```
@i function (g::Grad)(args...; kwargs...)
    @safe @assert count(x -> x isa Loss, args) == 1
    iloss = 0
    @routine for i=1:length(args)
        if (tget(args,i) isa Loss, iloss==i)
            iloss += identity(i)
            (~Loss)(tget(args,i))
        end
    end

    g.f(args...; kwargs...)
    GVar.(args)
    grad(tget(args,iloss)) += identity(1.0)
    (~g.f)(args...; kwargs...)

    ~@routine
end
```

Here, we use `Loss` type to mark the loss variable. Input variables must contain exactly one `Loss`. This program first checks the input parameters and locates the loss variable as `iloss`. Then `Loss` unwraps the loss variable. After computing the forward pass and backward pass, `@routine` uncomputes the ancilla `iloss` and returns the location infor-

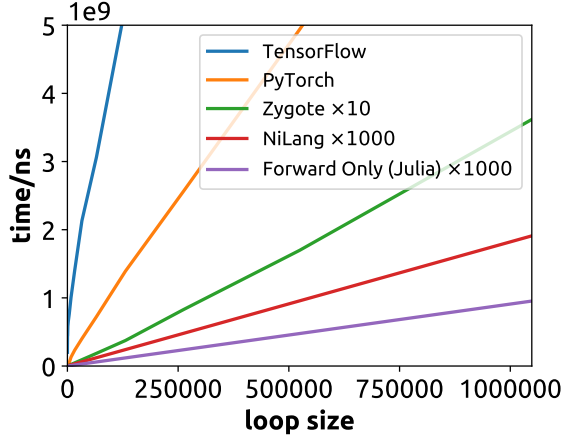


Figure 3. The time to obtain gradient as function of loop size. $\times n$ in legend represents a rescaling of time.

mation to the loss variable. `tget(args, i)` returns the i -th element of a tuple. Here, in order to avoid possible ambiguity in supporting array indexing, we forbid tuple indexing deliberately.

The overhead of using `GVar` type is negligible thanks to Julia’s multiple-dispatch and type inference. Let us consider a simple example that accumulates 1.0 to a target variable x for n times.

[JG: Grammarly here!]

```
julia> using NiLang, NiLang.AD, BenchmarkTools

julia> @i function prog(x, one, n::Int)
    for i=1:n
        x += identity(one)
    end
end

julia> @benchmark prog'(Loss(0.0), 1.0, 10000)
BenchmarkTools.Trial:
  memory estimate:  1.05 KiB
  allocs estimate:  39
  -----
  minimum time:     35.838 μs (0.00% GC)
  median time:      36.055 μs (0.00% GC)
  mean time:        36.483 μs (0.00% GC)
  maximum time:     185.973 μs (0.00% GC)
  -----
  samples:          10000
  evals/sample:     1
```

We implement the same function with TensorFlow, PyTorch and Zygote for comparison. The code could be found in our paper’s github repository [30]. Benchmark results on CPU Intel(R) Xeon(R) CPU E5-2680 v4 @ 2.40GHz are shown in Fig. 3. One can see that the NiLang implementation is unreasonably fast, it is approximately two times the forward pass written in native Julia code. Reversible programming is not always as fast as its irreversible counterparts. In

practical applications, a reversible program may have memory or computation overhead. We will discuss the details of time and space trade off in Sec. V A.

B. Second-order gradient

Second-order gradients can be obtained in two different approaches.

1. Back propagating first-order gradients

Back propagating the first-order gradients is the most widely used approach to obtain the second-order gradients. Suppose the function space is closed under gradient operation, one can obtain higher-order gradients by recursively differentiating lower order gradient functions without defining new backward rules.

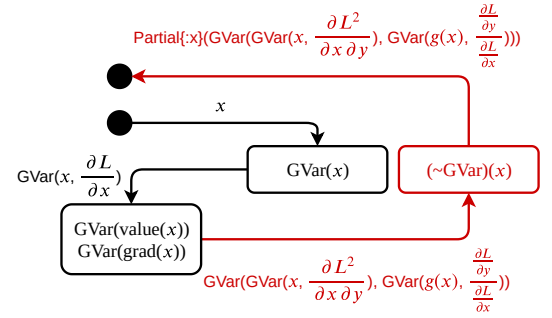


Figure 4. Obtaining the second-order gradient with the reverse differentiation approach. Black lines are computing gradients, red lines are back-propagating the process of obtaining the first-order gradients. Annotations on lines are data types used in the computation.

Fig. 4 shows the four passes in computing Hessian in this way. The first two passes (black lines) obtains the gradients. Before entering the third pass, the program wraps each field in `GVar` with another layer of `GVar`. Then we pick a variable x_i and add 1 to `grad(grad(x_i))` to compute the i -th row of Hessian. Before entering the final pass, the `~GVar` is called. It does not unwrap `GVar` directly because the second-order gradient fields may be nonzero in this case. Instead, we use `Partial{x}(·)` to take the x field of an instance without deallocating memory. By repeating the above process for different x_i , one can obtains the Hessian matrix.

2. Taylor propagation

A probably more efficient approach is back-propagating Hessians directly [31] using the relation

$$\begin{aligned} H_{O',O''}^O &= \mathbf{0}, \\ H_{x,x'}^O &= J_x^y H_{y,y'}^O J_{x'}^{y'} + J_y^O H_{x,x'}^y. \end{aligned} \quad (2)$$

Here, the Hessian tensor $H_{x,x'}^O$ is rank three, where the top index is often taken as a scalar and omitted. In tensor network language, the above equation can be represented as the right panel of Fig. 2. Hessian propagation is a special case of Taylor propagation. With respect to the order of gradients, Taylor propagation is exponentially more efficient in obtaining higher-order gradients than differentiating lower order gradients recursively. However, the exhausted support to Taylor propagation [10] requires much more effort than Jacobian propagation, this is why most AD packages choose the recursive approach. Instruction level automatic differentiation has the advantage of having very limited primitives. It is more flexible in obtaining higher-order gradients like Hessian. An example to obtain Hessians is provided in Sec. IV B.

C. Gradient on ancilla problem

An ancilla can also carry gradient during computation. As a result, even if an ancilla can be uncomputed rigorously in the original program, its GVar version can not be safely uncomputed. In these case, we simply “drop” the gradient field instead of raising an error. In this subsection, we prove the correctness of the following theorem

Theorem 1. *Dropping the gradient field of an ancilla when deallocating does not make the gradient function irreversible.*

Proof. Consider a reversible function $\vec{y}, b = f(\vec{x}, a)$, where a and b are the input and output values of an ancilla. The reversibility requires $b = a$ for any \vec{x} . So that

$$\frac{\partial b}{\partial \vec{x}} = \vec{0}. \quad (3)$$

Suppose in the backward pass, we discard the gradient field of b . Since the gradient fields are derived from the values of variables, they should not have any effect to the value fields. The rest is to show changing the value of $\text{grad}(b)$ does not result in a different $\text{grad}(\vec{x})$ in the backward pass. It can be seen from the back-propagation rule

$$\frac{\partial O}{\partial \vec{x}} = \frac{\partial O}{\partial \vec{y}} \frac{\partial \vec{y}}{\partial \vec{x}} + \frac{\partial O}{\partial b} \frac{\partial b}{\partial \vec{x}}, \quad (4)$$

where the second term with $\frac{\partial O}{\partial b}$ vanishes naturally. Hence the gradient field of an ancilla can be any value when entering a function, i.e. no need to cache the discarded value. \square

This theorem is very important to the AD in NiLang. Without it, ancillas and GVar will not fit with each other. On the other side, we emphasize that although the initial value of the gradient field of an ancilla can be randomly chosen, not having a gradient field at the beginning is a different story.

IV. EXAMPLES

A. Computing Fibonacci Numbers

The following is an example that everyone likes, computing Fibonacci number recursively.

```
using NiLang

@i function rfib(out!, n::T) where T
    n1 ← zero(T)
    n2 ← zero(T)
    @routine begin
        n1 += identity(n)
        n1 -= identity(1)
        n2 += identity(n)
        n2 -= identity(2)
    end
    if (value(n) <= 2, ~)
        out! += identity(1)
    else
        rfib(out!, n1)
        rfib(out!, n2)
    end
    ~@routine
end
```

Here, we use $x += \text{identity}(y)$ to represent $x += y$. The later is forbidden deliberately to avoid possible ambiguity between a function and a dataview. The time complexity of this recursive algorithm is exponential to input n . It is also possible to write a reversible linear time for loop algorithm.

A slightly non-trivial task is computing the first Fibonacci number that greater or equal to the second argument z , where a `while` statement is required.

```
@i function rfibn(n!, z)
    @safe @assert n! == 0
    out ← 0
    rfib(out, n!)
    while (out < z, n! != 0)
        ~rfib(out, n!)
        n! += identity(1)
        rfib(out, n!)
    end
    ~rfib(out, n!)
end
```

In this example, the postcondition $n!=0$ in the `while` statement is false before entering the loop, and becomes true in later iterations. In the reverse program, the `while` statement stops at $n==0$. If executed correctly, a user will see the following result.

```
julia> rfib(0, 10)
(55, 10)

julia> rfibn(0, 100)
(12, 100)

julia> (~rfibn)(rfibn(0, 100)...)
(0, 100)
```


B. exp function

An exp function can be computed using Taylor expansion

$$y+ = \sum_n \frac{x^n}{\text{factorial}(n)} \quad (5)$$

One can compute the accumulated item $s_n \equiv \frac{x^n}{\text{factorial}(n)}$ iteratively as $s_n = \frac{x s_{n-1}}{n}$. Considering the fact that product and division are considered as irreversible in NiLang, one can not deallocate s_{n-1} after computing s_n . This recursive computation mimics the famous pebble game [18]. However, there is no known constant memory and polynomial time solution to pebble game. Here the case is different. Notice $*$ and $/$ are arithmetically reversible to each other, we can “uncompute” previous state s_{n-1} by $s_{n-1} = \frac{n s_n}{x}$ approximately, and use the dirty ancilla in next iteration. The implementation is

```
using NiLang, NiLang.AD

@i function iexp(y!, x::T; atol::Float64=1e-14)
    where T
    anc1 ← zero(T)
    anc2 ← zero(T)
    anc3 ← zero(T)
    iplus ← 0
    expout ← zero(T)

    y! += identity(1.0)
    @routine begin
        anc1 += identity(1.0)
        while (value(anc1) > atol, iplus != 0)
            iplus += identity(1)
            anc2 += anc1 * x
            anc3 += anc2 / iplus
            expout += identity(anc3)
            # arithmetic uncompute
            anc1 -= anc2 / x
            anc2 -= anc3 * iplus
            SWAP(anc1, anc3)
        end
    end

    y! += identity(expout)

    ~@routine
end
```

Here, the definition of SWAP instruction can be found in Appendix B. The two lines bellow the comment “# arithmetic uncompute” uncompute variables anc1 and anc2 approximately, which is only arithmetically true. As a result, the final output is not exact due to the rounding error. On the other side, the reversibility is not affected since the inverse call at the last line of function uncomputes all ancillas rigorously.

To obtain gradients, one can wrap the variable y! with Loss type and feed it into iexp'

```
julia> y!, x = 0.0, 1.6
(0.0, 1.6)

julia> @instr iexp'(Loss(y!), x)

julia> grad(x)
4.9530324244260555
```

iexp' is a callable instance of type Grad{typeof(iexp)}. It wraps input variables with a gradient fields as outputs. This function itself is reversible and differentiable, one can back-propagate this function to obtain Hessians as introduced in Sec. III B 1. In NiLang, it is implemented as simple_hessian.

```
julia> y!, x = 0.0, 1.6
(0.0, 1.6)

julia> simple_hessian(iexp, (Loss(y!), x))
2×2 Array{Float64,2}:
 0.0  0.0
 0.0  4.95303
```

To obtain Hessians, we can also use the Taylor propagation approach as introduced in Sec. III B 2.

```
julia> y!, x = 0.0, 1.6
(0.0, 1.6)

julia> @instr iexp''(Loss(y!), x)

julia> collect_hessian()
2×2 Array{Float64,2}:
 0.0  0.0
 0.0  4.95303
```

iexp'' computes the second-order gradients. It wraps variables with type BeijingRing [32] in the backward pass. BeijingRing records Jacobians and Hessians for a variable, where Hessians are stored in a global storage. Whenever an n -th variable or ancilla is created, we push a ring of size $2n-1$ to a global tape. Whenever an ancilla is deallocated, we pop a ring from the top. The n -th ring stores $H_{i \leq n, n}$ and $H_{n, i < n}$. We didn't use the symmetry relation $H_{i,j} = H_{j,i}$ to save memory here in order to simplify the implementation of backward rules described in the right most panel of Fig. 2. The final result can be collected by calling collect_hessian(), it will read out the Hessian stored in the global storage.

C. QR decomposition

Let's consider a linear algebra function, the QR decomposition

```

using NiLang, NiLang.AD

@i function qr(Q, R, A::AbstractMatrix{T}) where T
    anc_norm ← zero(T)
    anc_dot ← zeros(T, size(A,2))
    ri ← zeros(T, size(A,1))
    for col = 1:size(A, 1)
        ri .+= identity.(A[:,col])
        for precol = 1:col-1
            dot(anc_dot[precol], Q[:,precol], ri)

            R[precol,col] += identity(anc_dot[precol])
            for row = 1:size(Q,1)
                ri[row] -= anc_dot[precol] * Q[row, precol]
            end
            norm2(anc_norm, ri)

            R[col, col] += anc_norm^0.5
            for row = 1:size(Q,1)
                Q[row,col] += ri[row] / R[col, col]
            end

            ~begin
                ri .+= identity.(A[:,col])
                for precol = 1:col-1
                    dot(anc_dot[precol], Q[:,precol], ri)
                    for row = 1:size(Q,1)
                        ri[row] -= anc_dot[precol] * Q[row, precol]
                    end
                end
                norm2(anc_norm, ri)
            end
        end
    end
end

```

Here, in order to avoid frequent uncomputing, we allocate ancillas `ri` and `anc_dot` as vectors. the expression in `~` is used to uncompute `ri`, `anc_dot` and `anc_norm`. `R[col, col] += anc_norm^0.5` is a ternary instruction, whose backward rule is defined in `NiLang`. The algorithm used to compute QR decomposition is very naive. It does not consider reorthogonalization. `idot` and `inorm2` are functions to compute dot product and vector norm. They are implemented as

```

@i function idot(out, v1::AbstractVector{T},
    v2) where T
    anc1 ← zero(T)
    for i = 1:length(v1)
        anc1 += identity(v1[i])
        CONJ(anc1)
        out += v1[i]*v2[i]
        CONJ(anc1)
        anc1 -= identity(v1[i])
    end
end

@i function inorm2(out, vec::AbstractVector{T})
    ) where T
    anc1 ← zero(T)
    for i = 1:length(vec)
        anc1 += identity(vec[i])
        CONJ(anc1)
        out += anc1*vec[i]
        CONJ(anc1)
        anc1 -= identity(vec[i])
    end
end

```

One can easily check the correctness of the gradient function

```

using Test
A = randn(4,4)
q = zero(A)
r = zero(A)

@i function test1(out, q, r, A)
    iqr(q, r, A)
    out += identity(q[1,2])
end

@i function test2(out, q, r, A)
    iqr(q, r, A)
    out += identity(r[1,2])
end

@test check_grad(test1, (Loss(0.0), q, r, A);
    atol=0.05, verbose=true)
@test check_grad(test2, (Loss(0.0), q, r, A);
    atol=0.05, verbose=true)

```

Here, the loss function `test1` and `test2` are defined as single elements in the output matrices. The `check_grad` function is a gradient checker function defined in module `NiLang.AD`. In this example, `GVar` is broadcasted to arrays. Thanks to Julia's just in time compiling, using `GVar` in an array does not have much overhead. It should be possible to define other reversible linear algebra functions too. We leave this future projects.

D. Unitary Matrices

Unitary matrices features uniform eigenvalues and reversibility. It is widely used as an approach to ease the gradient exploding and vanishing problem [33–35] and the memory wall problem [36]. One of the simplest way to parametrize a unitary matrix is representing a unitary matrix as a product of two-level unitary operations [35]. A real unitary matrix of size N can be parametrized compactly by $N(N - 1)/2$ rotation operations [37]

$$\text{ROT}(a!, b!, \theta) = \begin{bmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{bmatrix} \begin{bmatrix} a! \\ b! \end{bmatrix}, \quad (6)$$

where θ is the rotation angle, $a!$ and $b!$ are target registers.

```
using NiLang, NiLang.AD

@i function umm!(x!,  $\theta$ )
    @safe @assert length( $\theta$ ) ==
        length(x!)*(length(x!)-1)/2
    k ← 0
    for j=1:length(x!)
        for i=length(x!)-1:-1:j
            k += identity(1)
            ROT(x![i], x![i+1],  $\theta$ [k])
        end
    end
    k → length( $\theta$ )
end
```

Here, the ancilla k is deallocated manually by specifying its value, because we know the loop size is $N(N - 1)/2$. We define the test functions in order to check gradients.

```
julia> @i function isum(out!, x::Vector)
    for i=1:length(x)
        out! += identity(x[i])
    end
end

julia> @i function test!(out!, x!::Vector,  $\theta$ ::Vector)
    umm!(x!,  $\theta$ )
    isum(out!, x!)
end

julia> out, x,  $\theta$  = Loss(0.0), randn(4), randn(6);

julia> @instr test!(out, x,  $\theta$ )

julia> x
4-element Array{GVar{Float64,Float64},1}:
 GVar(1.220182125326287, 0.14540743042341095)
 GVar(2.1288634811475937, -1.3749962375499805)
 GVar(1.2696579252569677, 1.42868739498625)
 GVar(0.1083891125379283, 0.2170123344615735)

julia> @instr (~test!)(out, x,  $\theta$ )

julia> x
4-element Array{Float64,1}:
 1.220182125326287
 2.1288634811475933
 1.2696579252569677
 0.10838911253792821
```

In the above testing code, `test'` attaches a gradient field to each element of x . `~test'` is the inverse program that erase the gradient fields. Notably, this reversible implementation is inplace with zero memory allocation. In traditional automatic differentiation framework, one needs to implement special designs [12, 13] to describe this property.

V. DISCUSSION AND OUTLOOK

In this paper, we introduce automatic differentiation on a reversible Turing machine (RTM) and a Julia eDSL NiLang that simulates an RTM. We show a program on an RTM can be differentiated to any order reliably and efficiently without sophisticated designs to memorize computational graph and intermediate states.

In the following, we discussed some practical issues about reversible programming, and several future directions to go.

A. Time Space Tradeoff

In history, there has been many other interesting designs of reversible languages and instruction sets. However, current popular programming languages are all irreversible. Comparing with a irreversible Turing machine, an RTM has either a space overhead that proportional to computing time T or a computational overhead that sometimes can be exponential. The tradeoff between space and time is one of the most

important issue in the theory of RTM. In the simplest g-segment trade off scheme [38, 39], we have

$$Time(T) = \frac{T^{1+\epsilon}}{S^\epsilon}, \quad (7)$$

$$Space(T) = \epsilon 2^{1/\epsilon} (S + S \log \frac{T}{S}). \quad (8)$$

Here, T and S are the time and space usage on a irreversible Turing machine. ϵ is the control parameter. It is related to the g-segment trade off parameters by $g = k^n, \epsilon = \log_k(2k - 1)$ with $n \geq 1$ and $k \geq 1$. In this section, we try to convince the readers that the overhead of reversible computing is not as terrible as people thought.

First, let $\epsilon \rightarrow 0$, there is not overhead in time. The space used by a RTM is bounded the caching strategy used in a traditional machine learning package that memorizes every inputs of primitives. Memorizing the inputs always make a primitive reversible since it does not discard any information. For deep neural networks, people used checkpointing trick to trade time with space [40]. This trick is also widely used in reversible programming [18]. Reversible programming just provides more alternatives to trade time and space.

Second, many computational overheads come from of the irreversibility of $/=$ and $*=$ operations. This part is not fundamental because reversible floating point instructions have already been designed [41, 42]. Using reversible floating point instructions may significant decrease the computation time and memory usage of a RTM. Even in current stage, the overhead can be mitigated by "arithematic uncomputing" without sacrificing reversibility as shown in the `iexp` example. We will review this point in Sec. V B.

Thrid, clever compiling can remove most overheads. Often, when we define a new reversible function, we allocate some ancillas at the begining of the function and deallocate them through uncomputing at the end. The overhead comes from the uncomputing, in the worst case, the time used for uncomputing can be the same as the forward pass. In a hierarchical design, uncomputing can happend in every layer. To quantify the overhead of uncomputing, we introducing the concept

Definition 1 (program granularity). The logarithm of the ratio between the execution time of a reversible program and its irreversible counter part

$$\log_2 \frac{Time(T)}{T}. \quad (9)$$

The computing time increases exponentially as the granularity increases. A cleverer compilation of a program can reduce the granularity. Given plenty of space, all uncomputing can be merged to avoid repeated works since the uncomputing of ancillas can be executed in any level of hierarchy.

At last, making reversible programming an eDSL rather than an independent language allows flexible choices between reversibility and computational overhead. For example, in order to deallocate the gradient memory in a reversible language one has to uncompute the whole process of obtaining this gradient. In our eDSL, we can just deallocate the memory irreversibly, i.e. trade energy with time.

B. Instructions and Hardwares

So far, our eDSL is not really compiled to instructions, instead, it runs on a irreversible host Julia. In the future, it can be compiled to low level instructions and is execute on a reversible devices. For example, the control flow defined in this NiLang can be compiled to reversible instructions like reversible goto instruction, where the target instruction can be a `comefrom` instruction that specifying the postcondition. [43]

Arithmetic instructions should be redesigned to support better reversible programs. The major obstacle to exact reversibility programming is current floating point adders and multipliers used in our computing devices are not exactly reversible. There are proposals of reversible floating point adders and multipliers [41, 42, 44, 45] that introduces garbage bits to ensure reversibility, however this design with garbage qubits allocates new bits in each operaton, which is not too different from the information buffer approach [11]. With floating point numbers, rigorous reversible arithematic designs without using information buffer or garbage qubits is nearly impossible. Alternatives include fixed point numbers [46] and logarithmic numbers [47, 48], where logarithmic number system is reversible under $* =$ and $/ =$. With these infrastructures, a reversible program can be executed without suffering from the rounding error.

Reversible programming is not nessesarily related to reversible hardwares. Reversible programs is a subset of irreversible programs, hence can be simulated efficiently on traditional CMOS devices [43]. Reversible programming just provides an alternative to execute on an energy efficient reversible hardwares [49] like adiabatic CMOS [50], molecular mechanical computing system [51] and superconducting system [52, 53]. Reversible hardwares are not necessarily related to reversible gates such as the Toffoli gate and the Fredkin gate. Devices with the ability to recover signal energy can also be used to run reversible programs energy efficiently, which is known as the generalized reversible computing. [54, 55] In the following, we comment briefly on a special type of reversible device Quantum computer.

1. Quantum Computers

Building a universal quantum computer is difficult. The difficulty lies in the fact that, unlike a classical state, an unknown quantum state can not be copied. A quantum state in a environment suffers from decoherence, this underlines the simulation nature of quantum devices. Although there are proposals about quantum random access memory [56], they are difficult to implement, and are known to have many caveats [57]. Reversible computing does not enjoy the quantum advantage, nor the quantum disadvantages of non-cloning and decoherence. The reversibility of quantum computing comes from the fact that microscopic processes are unitary.

Given the fundamental limitations of quantum decoherence and non-cloning and the reversible nature of microscopic world. It is reasonable to have a reversible computing device

to bridge the gap between classical devices and universal quantum computing devices. By introducing entanglement little by little, we can accelerate some basic components in reversible computing. For example, quantum Fourier transformation provides an interesting alternative to the reversible adders and multipliers by introducing one additional CPHASE gate, even though adders and multipliers are classical functions. [58] The development of reversible compiling theory will also have profounding effect to quantum compiling.

C. Outlook

The reversible eDSL NiLang can be used to solve many existing scientific computing problems. First of all, it can be used to generate AD rules for existing machine learning packages. For example, in Zygote, a user sometimes need to define new primitives. If this primitive is written in NiLang, then the backward rules can be generated automatically. We can built reversible BLAS and LAPACK on top of NiLang. These functions are extensively used in physics [59, 60] and many other disciplines. However, manually derived backward rules for singular value decomposition and eigenvalue decomposition functions suffer from the gradient exploding problem [61–63]. Hopefully, the automatically generated backward rules do not have such problems.

Secondly, we can use it to break the memory wall problem. NiLang provides a systematic time-space trade off scheme through uncomputing. A successful related example is the memory efficient domain-specific AD engine in quantum simulator Yao [36]. This domain-specific AD engine is written in a reversible style and solved the memory bottleneck in variational quantum simulations. It also gives hitherto the best performance in differentiating quantum circuit parameters. Similarly, we can write memory efficient normalizing flow [64] with NiLang. Normalizing flow is a successful class of generative model in both computer vision [65] and quantum physics [66, 67]. Its building block bijector is required to be reversible and differentiable. We can use similar idea to differentiate reversible integrators [68, 69]. With reversible integrators, it should be possible to rewrite the control system in robotics [70] in a reversible style. In robotics, the control parameters are often floating point numbers rather than tensors. Writing a control program reversibly should boost the training performance a lot.

Thridly, reversibility is a resource for traning. For those

who are interested in non-gradient based training. In Appendix C, we provide a self-consistency training strategy for reversible programs.

Latstly, the reversible IR is a good starting point to study quantum compiling. Most quantum programming language preassumes a classical coprocessor and use classical control flows [71] in universal quantum computing. However, we believe reversible control flows are also very important to a universal quantum computer.

To solve the above problems better, NiLang can be improved from multiple perspectives.

- Add a compiling stage that decreases granularity and reduces uncomputing overheads.
- Better type inference. Current type inference assumes the variable types not changed by a function, which is not true. The type cast rule of a reversible function should be recorded somewhere to help type inference.
- Implement rigorous reversible floating point arithmetics on hardwares to let the reversibility free from rounding error.
- Show the advantage of NiLang in parallel computing in handle asynchronous computing [72] and debugging with bidirectional move [73]. It is interesting to see how NiLang combines with other parts of Julia ecosystem like CUDAnative [74] and Debugger.

These improvements needs participations of people from multiple fields.

VI. ACKNOWLEDGMENTS

Jin-Guo Liu thank Lei Wang for motivating the project with possible applications reversible integrator, normalizing flow and neural ODE. Xiu-Zhe Luo for discussion on the implementation details of source to source automatic differentiation, Shuo-Hui Li for helpful discussion on differential geometry. Damian Steiger for telling me the come from joke. Tong Liu and An-Qi Chen for helpful discussion on quantum adders and multipliers. Chris Rackauckas for helpful discussion on reversible integrators. Mike Innes for reviewing the comments about Zygote. Simon Byrne and Chen Zhao for helpful discussion on floating point and logarithmic numbers. 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.

[1] L. Hascoet and V. Pascual, *ACM Transactions on Mathematical Software (TOMS)* **39**, 20 (2013).
 [2] 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).

[3] 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 \[cs.PL\]](https://arxiv.org/abs/1811.01457).
 [4] 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. Jozefow-

- icz, 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.
- [5] M. Innes, “Don’t unroll adjoint: Differentiating ssa-form programs,” (2018), [arXiv:1810.07951 \[cs.PL\]](#).
- [6] M. Innes, A. Edelman, K. Fischer, C. Rackauckas, E. Saba, V. B. Shah, and W. Tebbutt, *CoRR* **abs/1907.07587** (2019), [arXiv:1907.07587](#).
- [7] Z. C. Lipton, J. Berkowitz, and C. Elkan, “A critical review of recurrent neural networks for sequence learning,” (2015), [arXiv:1506.00019 \[cs.LG\]](#).
- [8] K. He, X. Zhang, S. Ren, and J. Sun, *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2016), 10.1109/cvpr.2016.90.
- [9] “Breaking the Memory Wall: The AI Bottleneck,” <https://blog.semi.org/semi-news/breaking-the-memory-wall-the-ai-bottleneck>.
- [10] J. Bettencourt, M. J. Johnson, and D. Duvenaud, (2019).
- [11] 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.
- [12] A. N. Gomez, M. Ren, R. Urtasun, and R. B. Grosse, in *Advances in Neural Information Processing Systems 30*, edited by I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Curran Associates, Inc., 2017) pp. 2214–2224.
- [13] J.-H. Jacobsen, A. W. Smeulders, and E. Oyallon, in *International Conference on Learning Representations* (2018).
- [14] 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.
- [15] J. Behrmann, D. Duvenaud, and J. Jacobsen, *CoRR* **abs/1811.00995** (2018), [arXiv:1811.00995](#).
- [16] J. Bezanson, S. Karpinski, V. B. Shah, and A. Edelman, *arXiv preprint arXiv:1209.5145* (2012).
- [17] J. Bezanson, A. Edelman, S. Karpinski, and V. B. Shah, *SIAM Review* **59**, 65–98 (2017).
- [18] K. S. Perumalla, *Introduction to reversible computing* (Chapman and Hall/CRC, 2013).
- [19] M. P. Frank, *IEEE Spectrum* **54**, 32–37 (2017).
- [20] C. Lutz, “Janus: a time-reversible language,” (1986), *Letter to R. Landauer*.
- [21] M. P. Frank, *The R programming language and compiler*, Tech. Rep. (MIT Reversible Computing Project Memo, 1997).
- [22] I. Lanese, N. Nishida, A. Palacios, and G. Vidal, *Journal of Logical and Algebraic Methods in Programming* **100**, 71–97 (2018).
- [23] T. Haulund, “Design and implementation of a reversible object-oriented programming language,” (2017), [arXiv:1707.07845 \[cs.PL\]](#).
- [24] R. Landauer, *IBM journal of research and development* **5**, 183 (1961).
- [25] G. H. Golub and C. F. Van Loan, *Matrix computations*, Vol. 3 (JHU press, 2012).
- [26] “MLStyle.jl,” <https://github.com/thautwarm/MLStyle.jl>.
- [27] R. Orús, *Annals of Physics* **349**, 117–158 (2014).
- [28] A. McInerney, *First steps in differential geometry* (Springer, 2015).
- [29] J. Revels, M. Lubin, and T. Papamarkou, “Forward-mode automatic differentiation in julia,” (2016), [arXiv:1607.07892 \[cs.MS\]](#).
- [30] “Paper’s Github Repository,” <https://github.com/GiggleLiu/nilangpaper/tree/master/codes>.
- [31] J. Martens, I. Sutskever, and K. Swersky, “Estimating the hessian by back-propagating curvature,” (2012), [arXiv:1206.6464 \[cs.LG\]](#).
- [32] When people ask for the location in Beijing, they will start by asking which ring. We use the similar approach to locate the elements of Hessian matrix.
- [33] M. Arjovsky, A. Shah, and Y. Bengio, *CoRR* **abs/1511.06464** (2015), [arXiv:1511.06464](#).
- [34] S. Wisdom, T. Powers, J. R. Hershey, J. L. Roux, and L. Atlas, “Full-capacity unitary recurrent neural networks,” (2016), [arXiv:1611.00035 \[stat.ML\]](#).
- [35] L. Jing, Y. Shen, T. Dubcek, J. Peurifoy, S. A. Skirlo, M. Tegmark, and M. Soljacic, *CoRR* **abs/1612.05231** (2016), [arXiv:1612.05231](#).
- [36] X.-Z. Luo, J.-G. Liu, P. Zhang, and L. Wang, “Yao.jl: Extensible, efficient framework for quantum algorithm design,” (2019), [arXiv:1912.10877 \[quant-ph\]](#).
- [37] C.-K. LI, R. ROBERTS, and X. YIN, *International Journal of Quantum Information* **11**, 1350015 (2013).
- [38] C. H. Bennett, *SIAM Journal on Computing* **18**, 766 (1989).
- [39] R. Y. Levine and A. T. Sherman, *SIAM Journal on Computing* **19**, 673 (1990).
- [40] T. Chen, B. Xu, C. Zhang, and C. Guestrin, *CoRR* **abs/1604.06174** (2016), [arXiv:1604.06174](#).
- [41] M. Nachtigal, H. Thapliyal, and N. Ranganathan, in *10th IEEE International Conference on Nanotechnology* (2010) pp. 233–237.
- [42] M. Nachtigal, H. Thapliyal, and N. Ranganathan, in *2011 11th IEEE International Conference on Nanotechnology* (2011) pp. 451–456.
- [43] C. J. Vieri, *Reversible Computer Engineering and Architecture*, Ph.D. thesis, Cambridge, MA, USA (1999), aAI0800892.
- [44] T. D. Nguyen and R. V. Meter, “A space-efficient design for reversible floating point adder in quantum computing,” (2013), [arXiv:1306.3760 \[quant-ph\]](#).
- [45] T. Häner, M. Soeken, M. Roetteler, and K. M. Svore, “Quantum circuits for floating-point arithmetic,” (2018), [arXiv:1807.02023 \[quant-ph\]](#).
- [46] “FixedPointNumbers.jl,” <https://github.com/JuliaMath/FixedPointNumbers.jl>.
- [47] F. J. Taylor, R. Gill, J. Joseph, and J. Radke, *IEEE Transactions on Computers* **37**, 190 (1988).
- [48] “LogarithmicNumbers.jl,” <https://github.com/cjdoris/LogarithmicNumbers.jl>.
- [49] M. P. Frank and T. F. Knight Jr, *Reversibility for efficient computing*, Ph.D. thesis, Massachusetts Institute of Technology, Dept. of Electrical Engineering and ... (1999).
- [50] J. G. Koller and W. C. Athas, in *Workshop on Physics and Computation* (1992) pp. 267–270.
- [51] R. C. Merkle, R. A. Freitas, T. Hogg, T. E. Moore, M. S. Moses, and J. Ryley, *Journal of Mechanisms and Robotics* **10** (2018), 10.1115/1.4041209.
- [52] K. Likharev, *IEEE Transactions on Magnetics* **13**, 242 (1977).
- [53] V. K. Semenov, G. V. Danilov, and D. V. Averin, *IEEE Transactions on Applied Superconductivity* **13**, 938 (2003).
- [54] M. P. Frank, in *35th International Symposium on Multiple-Valued Logic (ISMVL’05)* (2005) pp. 168–185.

- [55] M. P. Frank, in *Reversible Computation*, edited by I. Phillips and H. Rahaman (Springer International Publishing, Cham, 2017) pp. 19–34.
- [56] V. Giovannetti, S. Lloyd, and L. Maccone, *Physical Review Letters* **100** (2008), 10.1103/physrevlett.100.160501.
- [57] S. Aaronson, *Nature Physics* **11**, 291 (2015).
- [58] L. Ruiz-Perez and J. C. Garcia-Escartin, *Quantum Information Processing* **16** (2017), 10.1007/s11128-017-1603-1.
- [59] J.-G. L. Hao Xie and L. Wang, *arXiv:2001.04121*.
- [60] H.-J. Liao, J.-G. Liu, L. Wang, and T. Xiang, *Physical Review X* **9** (2019), 10.1103/physrevx.9.031041.
- [61] M. Seeger, A. Hetzel, Z. Dai, E. Meissner, and N. D. Lawrence, “Auto-differentiating linear algebra,” (2017), *arXiv:1710.08717 [cs.MS]*.
- [62] Z.-Q. Wan and S.-X. Zhang, “Automatic differentiation for complex valued svd,” (2019), *arXiv:1909.02659 [math.NA]*.
- [63] C. Hubig, “Use and implementation of autodifferentiation in tensor network methods with complex scalars,” (2019), *arXiv:1907.13422 [cond-mat.str-el]*.
- [64] I. Kobyzev, S. Prince, and M. A. Brubaker, “Normalizing flows: An introduction and review of current methods,” (2019), *arXiv:1908.09257 [stat.ML]*.
- [65] D. P. Kingma and P. Dhariwal, 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. 10215–10224.
- [66] L. Dinh, J. Sohl-Dickstein, and S. Bengio, “Density estimation using real nvp,” (2016), *arXiv:1605.08803 [cs.LG]*.
- [67] S.-H. Li and L. Wang, *Physical Review Letters* **121** (2018), 10.1103/physrevlett.121.260601.
- [68] P. Hut, J. Makino, and S. McMillan, *The Astrophysical Journal* **443**, L93 (1995).
- [69] D. N. Laikov, *Theoretical Chemistry Accounts* **137** (2018), 10.1007/s00214-018-2344-7.
- [70] M. Gifftaler, M. Neunert, M. Stäuble, M. Frigerio, C. Semini, and J. Buchli, *Advanced Robotics* **31**, 1225–1237 (2017).
- [71] K. Svore, M. Roetteler, A. Geller, M. Troyer, J. Azariah, C. Granade, B. Heim, V. Kliuchnikov, M. Mykhailova, and A. Paz, *Proceedings of the Real World Domain Specific Languages Workshop 2018 on - RWDSL2018* (2018), 10.1145/3183895.3183901.
- [72] D. R. Jefferson, *ACM Transactions on Programming Languages and Systems (TOPLAS)* **7**, 404 (1985).
- [73] B. Boothe, *ACM SIGPLAN Notices* **35**, 299 (2000).
- [74] T. Besard, C. Foket, and B. D. Sutter, *CoRR* **abs/1712.03112** (2017), *arXiv:1712.03112*.
- [75] M. Bender, P.-H. Heenen, and P.-G. Reinhard, *Rev. Mod. Phys.* **75**, 121 (2003).

Appendix A: NiLang Grammar

To define a reversible function one can use “@i” plus a normal function definition like bellow

```

"""
docstring...
"""
@i function f(args..., kwargs...) where {...}
    <stmts>
end

```

where the definition of “<stmts>” are shown in the grammar on the next column. The following is a list of terminologies used in the definition of grammar

- *ident*, symbols
- *num*, numbers
- ϵ , empty statement
- *JuliaExpr*, native Julia expression
- [], zero or one repetitions.

Here, all *JuliaExpr* should be pure, otherwise the reversibility is not guaranteed. Dataview is a view of a data, it can be a bijective mapping of an object, an item of an array or a field of an object.

$\langle \text{Stmts} \rangle ::= \epsilon$
 $\quad | \langle \text{Stmt} \rangle$
 $\quad | \langle \text{Stmts} \rangle \langle \text{Stmt} \rangle$
 $\langle \text{Stmt} \rangle ::= \langle \text{BlockStmt} \rangle$
 $\quad | \langle \text{IfStmt} \rangle$
 $\quad | \langle \text{WhileStmt} \rangle$
 $\quad | \langle \text{ForStmt} \rangle$
 $\quad | \langle \text{InstrStmt} \rangle$
 $\quad | \langle \text{RevStmt} \rangle$
 $\quad | \langle \text{AncillaStmt} \rangle$
 $\quad | \langle \text{TypecastStmt} \rangle$
 $\quad | \langle @routine \rangle \langle \text{Stmt} \rangle$
 $\quad | \langle @safe \rangle \text{JuliaExpr}$
 $\quad | \langle \text{CallStmt} \rangle$
 $\langle \text{BlockStmt} \rangle ::= \text{begin} \langle \text{Stmts} \rangle \text{end}$
 $\langle \text{RevCond} \rangle ::= (\text{JuliaExpr} , \text{JuliaExpr})$
 $\langle \text{IfStmt} \rangle ::= \text{if} \langle \text{RevCond} \rangle \langle \text{Stmts} \rangle [\text{else} \langle \text{Stmts} \rangle] \text{end}$
 $\langle \text{WhileStmt} \rangle ::= \text{while} \langle \text{RevCond} \rangle \langle \text{Stmts} \rangle \text{end}$
 $\langle \text{Range} \rangle ::= \text{JuliaExpr} : \text{JuliaExpr} [: \text{JuliaExpr}]$
 $\langle \text{ForStmt} \rangle ::= \text{for } \text{ident} = \langle \text{Range} \rangle \langle \text{Stmts} \rangle \text{end}$
 $\langle \text{KwArg} \rangle ::= \text{ident} = \text{JuliaExpr}$
 $\langle \text{KwArgs} \rangle ::= [\langle \text{KwArgs} \rangle ,] \langle \text{KwArg} \rangle$
 $\langle \text{CallStmt} \rangle ::= \text{JuliaExpr} ([\langle \text{DataViews} \rangle] [: \langle \text{KwArgs} \rangle])$
 $\langle \text{Constant} \rangle ::= \text{num} \mid \pi$
 $\langle \text{InstrBinOp} \rangle ::= + = \mid - = \mid \forall =$
 $\langle \text{InstrTrailer} \rangle ::= [.] ([\langle \text{DataViews} \rangle])$
 $\langle \text{InstrStmt} \rangle ::= \langle \text{DataView} \rangle \langle \text{InstrBinOp} \rangle \text{ident} [\langle \text{InstrTrailer} \rangle]$
 $\langle \text{RevStmt} \rangle ::= \sim \langle \text{Stmt} \rangle$
 $\langle \text{AncillaStmt} \rangle ::= \text{ident} \leftarrow \text{JuliaExpr}$
 $\langle \text{TypecastStmt} \rangle ::= (\text{JuliaExpr} \Rightarrow \text{JuliaExpr}) (\text{ident})$
 $\langle @routine \rangle ::= @routine \text{ident} \langle \text{Stmt} \rangle$
 $\langle @safe \rangle ::= @safe \text{JuliaExpr}$
 $\langle \text{DataViews} \rangle ::= \epsilon$
 $\quad | \langle \text{DataView} \rangle$
 $\quad | \langle \text{DataViews} \rangle , \langle \text{DataView} \rangle$
 $\quad | \langle \text{DataViews} \rangle , \langle \text{DataView} \rangle \dots$
 $\langle \text{DataView} \rangle ::= \langle \text{DataView} \rangle [\text{JuliaExpr}]$
 $\quad | \langle \text{DataView} \rangle . \text{ident}$
 $\quad | \text{JuliaExpr} (\langle \text{DataView} \rangle)$
 $\quad | \langle \text{DataView} \rangle '$
 $\quad | - \langle \text{DataView} \rangle$
 $\quad | \langle \text{Constant} \rangle$
 $\quad | \text{ident}$

Appendix B: Instruction Table

The translation of instructions to Julia functions
The list of instructions implemented in NiLang

instruction	translated	symbol
$y += f(\text{args} \dots)$	$\text{PlusEq}(f)(\text{args} \dots)$	\oplus
$y -= f(\text{args} \dots)$	$\text{MinusEq}(f)(\text{args} \dots)$	\ominus
$y \forall = f(\text{args} \dots)$	$\text{XorEq}(f)(\text{args} \dots)$	\odot

Table II. Instructions and their interpretation in NiLang.

instruction	output
$\text{SWAP}(a, b)$	b, a
$\text{ROT}(a, b, \theta)$	$a \cos \theta - b \sin \theta, b \cos \theta + a \sin \theta, \theta$
$\text{IROT}(a, b, \theta)$	$a \cos \theta + b \sin \theta, b \cos \theta - a \sin \theta, \theta$
$y += a^b$	$y + a^b, a, b$
$y += \exp(x)$	$y + e^x, x$
$y += \log(x)$	$y + \log x, x$
$y += \sin(x)$	$y + \sin x, x$
$y += \cos(x)$	$y + \cos x, x$
$y += \text{abs}(x)$	$y + x , x$
$\text{NEG}(y)$	$-y$
$\text{CONJ}(y)$	y'

Table III. A collection of reversible instructions, “.” is the broadcasting operations in Julia.

Appendix C: Learn by consistency

Consider a training that with input \vec{x}^* and output \vec{y}^* , find a set of parameters \vec{p}_x that satisfy $\vec{y}^* = f(\vec{x}^*, \vec{p}_x)$. In traditional machine learning, we define a loss $\mathcal{L} = \text{dist}(\vec{y}^*, f(\vec{x}^*, \vec{p}_x))$ and minimize it with gradient $\frac{\partial \mathcal{L}}{\partial \vec{p}_x}$. This works only when the target function is locally differentiable.

Here we provide an alternative by making use of reversibility. We construct a reversible program $\vec{y}, \vec{p}_y = f_r(\vec{x}, \vec{p}_x)$, where \vec{p}_x and \vec{p}_y are “parameter” spaces on the input side and output side. The algorithm can be summarized as

Algorithm 2: Learn by consistency

Result: \vec{p}_x
Initialize \vec{x} to \vec{x}^* , parameter space \vec{p}_x to random.
if \vec{p}_y is *null* **then**
 $\quad \vec{x}, \vec{p}_x = f_r^{-1}(\vec{y}^*)$
else
 $\quad \vec{y}, \vec{p}_y = f_r(\vec{x}, \vec{p}_x)$
while $\vec{y} \neq \vec{y}^*$ **do**
 $\quad \vec{y} = \vec{y}^*$
 $\quad \vec{x}, \vec{p}_x = f_r^{-1}(\vec{y}, \vec{p}_y)$
 $\quad \vec{x} = \vec{x}^*$
 $\quad \vec{y}, \vec{p}_y = f_r(\vec{x}, \vec{p}_x)$

Here, $\text{parameter}(\cdot)$ is a function for taking the parameter space. This algorithm utilizes the self-consistency relation

$$\vec{p}_x^* = \text{parameter}(f_r^{-1}(\vec{y}^*, \text{parameter}(f_r(\vec{x}^*, \vec{p}_x^*)))), \quad (\text{C1})$$

Similar idea of training by consistency is used in self-consistent meanfield theory [75] in physics. Finding the

self-consistent relation is crucial to a self-consistency based training. Here, the reversibility provides a natural self-consistency relation. However, it is not a silver bullet, let's consider the following example

```
@i function f1(y!, x, p!)
    p! += identity(x)
    y! -= exp(x)
    y! += exp(p!)
end

@i function f2(y!, x!, p!)
    p! += identity(x!)
    y! -= exp(x!)
    x! -= log(-y!)
    y! += exp(p!)
end

function train(f)
    loss = Float64[]
    p = 1.6
    for i=1:100
        y!, x = 0.0, 0.3
        @instr f(y!, x, p)
        push!(loss, y!)
        y! = 1.0
        @instr (~f)(y!, x, p)
    end
    loss
end
```

Functions `f1` and `f2` computes $f(x, p) = e^{(p+x)} - e^x$ and stores the output in a new memory `y!`. The only difference is `f2` uncomputes x arithmetically. The task of training is to find a p that make the output value equal to target value 1. After 100 steps, `f2` runs into the fixed point with x equal to 1 upto machine precision. However, parameters in `f1` does change at all. The training of `f1` fails because this function actually computes $f1(y, x, p) = y + e^{(p+x)} - e^x, x, x+p$, where the training parameter p is completely determined by the parameter space on the output side $x \cup x + p$. As a result, shifting y directly is the only approach to satisfy the consistency relation. On the other side, $f2(y, x, p) = y + e^{(p+x)} - e^x, \tilde{0}, x + p$, the output parameters $\tilde{0} \cup x + p$ can not uniquely determine input parameters p and x . Here, we use $\tilde{0}$ to denote the zero with rounding error.

By viewing \vec{x} and parameters in \vec{p}_x as variables, we can

study the trainability from the information perspective.

Theorem 2. Only if the conditional entropy $S(\vec{y}|\vec{p}_y)$ is nonzero, algorithm 2 is trainable.

Proof. The above example reveals a fact that the training can not work when output parameters completely determines input parameters. In other words, if $S(\vec{p}_x|\vec{p}_y) = 0$, the training

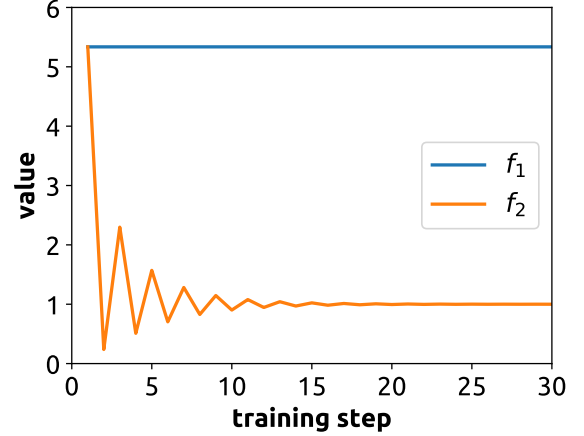


Figure 5. The output value `y!` as a function of self-consistent training step.

can not work.

$$\begin{aligned}
 S(\vec{p}_x|\vec{p}_y) &= S(\vec{p}_x \cup \vec{p}_y) - S(\vec{p}_y) \\
 &\leq S((\vec{p}_x \cup \vec{x}) \cup \vec{p}_y) - S(\vec{p}_y), \\
 &\leq S((\vec{p}_y \cup \vec{y}) \cup \vec{p}_y) - S(\vec{p}_y), \\
 &\leq S(\vec{y}|\vec{p}_y).
 \end{aligned} \tag{C2}$$

The third line uses the bijectivity $S(\vec{x} \cup \vec{p}_x) = S(\vec{y} \cup \vec{p}_y)$. This inequality shows that when the parameter space on the output side satisfies $S(\vec{y}|\vec{p}_y) = 0$, i.e. contains all information to determine the output field, the input parameters are also completely determined by this parameter space, hence training can not work. \square

In the above example, it corresponds to the case $S(e^{(x+y)} - e^x | x \cup x + y) = 0$ in `f1`. The solution is to remove the information redundancy in output parameter space through uncomputing as shown in `f2`.