

# Differentiation of stateful functions

Shin-ya Katsumata and David Sprunger\*

ERATO MMSD Colloquium  
February 7, 2019

# Outline

- 1 Notation: Functions as diagrams
- 2 Related work: function unrolling and BPTT
- 3 Main goals
- 4 Causal function formalization in category theory
  - Main ideas
  - Sanity checks
- 5 Delayed trace
- 6 Cartesian differential structure
- 7 Recap and future directions

# Outline

- 1 Notation: Functions as diagrams
- 2 Related work: function unrolling and BPTT
- 3 Main goals
- 4 Causal function formalization in category theory
  - Main ideas
  - Sanity checks
- 5 Delayed trace
- 6 Cartesian differential structure
- 7 Recap and future directions

## Ordinary functions as diagrams

We depict functions as boxes with input/output wires. For example, multiplication  $\times : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$  and copying  $\Delta_{\mathbb{R}} : \mathbb{R} \rightarrow \mathbb{R} \times \mathbb{R}$  are:

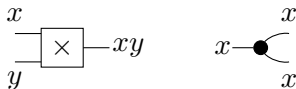


## Ordinary functions as diagrams

We depict functions as boxes with input/output wires. For example, multiplication  $\times : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$  and copying  $\Delta_{\mathbb{R}} : \mathbb{R} \rightarrow \mathbb{R} \times \mathbb{R}$  are:



We may also put values on the wires for scratchwork, so

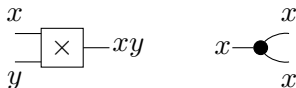


## Ordinary functions as diagrams

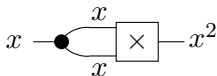
We depict functions as boxes with input/output wires. For example, multiplication  $\times : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$  and copying  $\Delta_{\mathbb{R}} : \mathbb{R} \rightarrow \mathbb{R} \times \mathbb{R}$  are:



We may also put values on the wires for scratchwork, so



We depict sequential composition by:

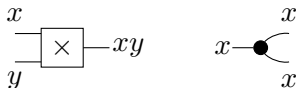


## Ordinary functions as diagrams

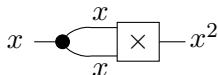
We depict functions as boxes with input/output wires. For example, multiplication  $\times : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$  and copying  $\Delta_{\mathbb{R}} : \mathbb{R} \rightarrow \mathbb{R} \times \mathbb{R}$  are:



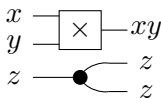
We may also put values on the wires for scratchwork, so



We depict sequential composition by:



We can also compose functions in parallel:



# Functions on sequences, part I

Ordinary functions can be upgraded to functions on sequences:

Input

Output



2

3

⋮



# Functions on sequences, part I

Ordinary functions can be upgraded to functions on sequences:

Input

Output



2

3

⋮

## Functions on sequences, part I

Ordinary functions can be upgraded to functions on sequences:

Input                      Output

1

1



3

⋮

## Functions on sequences, part I

Ordinary functions can be upgraded to functions on sequences:

Input                      Output

1

1

2



4

3

⋮

# Functions on sequences, part I

Ordinary functions can be upgraded to functions on sequences:

<u>Input</u>	<u>Output</u>
--------------	---------------

1	1
---	---

2	4
---	---



⋮

# Functions on sequences, part I

Ordinary functions can be upgraded to functions on sequences:

Input                      Output

1                                      1

2                                      4

3  9

⋮

# Functions on sequences, part I

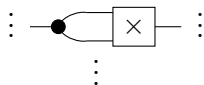
Ordinary functions can be upgraded to functions on sequences:

Input                      Output

1                              1

2                              4

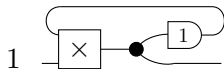
3                              9



## Functions on sequences, part II

We also use state and feedback mechanisms for sequence functions:

Input Output



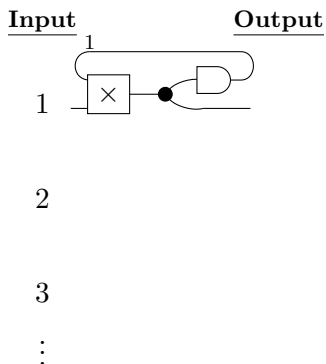
2

3

⋮

## Functions on sequences, part II

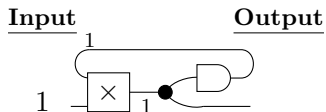
We also use state and feedback mechanisms for sequence functions:





## Functions on sequences, part II

We also use state and feedback mechanisms for sequence functions:

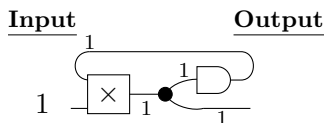


2

3

## Functions on sequences, part II

We also use state and feedback mechanisms for sequence functions:



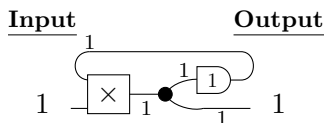
2

3

⋮

## Functions on sequences, part II

We also use state and feedback mechanisms for sequence functions:



2

3

⋮

## Functions on sequences, part II

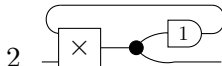
We also use state and feedback mechanisms for sequence functions:

Input

Output

1

1



2

3

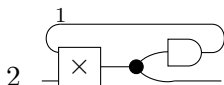
$\vdots$

## Functions on sequences, part II

We also use state and feedback mechanisms for sequence functions:

Input                      Output

1                                      1



3

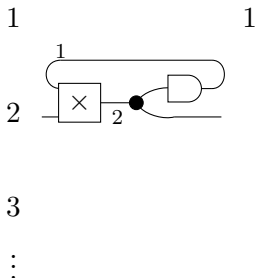
⋮

## Functions on sequences, part II

We also use state and feedback mechanisms for sequence functions:

Input

Output

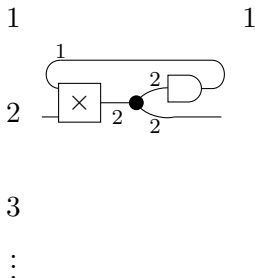


## Functions on sequences, part II

We also use state and feedback mechanisms for sequence functions:

Input

Output

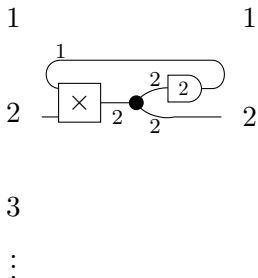


## Functions on sequences, part II

We also use state and feedback mechanisms for sequence functions:

Input

Output





## Functions on sequences, part II

We also use state and feedback mechanisms for sequence functions:

Input                      Output

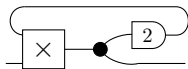
1

1

2

2

3



⋮

## Functions on sequences, part II

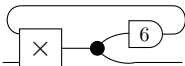
We also use state and feedback mechanisms for sequence functions:

Input                      Output

1                                      1

2                                      2

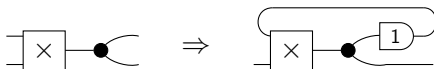
3                                      6



⋮                                      ⋮

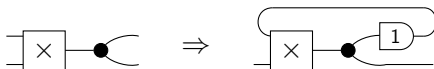
## What is this diagram?

What do diagrams like this, or operational semantics shown previously, remind you of?



# What is this diagram?

What do diagrams like this, or operational semantics shown previously, remind you of?



*\*ahem\**

Metamathematical  
Transfer

Meta-theoretician

What's happening here?

... uniform & comprehensive construction

$T + e \rightarrow T[e]$

$T_1 + e_1 \rightarrow T_1[e_1]$

$T_2 + e_2 \rightarrow T_2[e_2]$

$T_3 + e_3 \rightarrow T_3[e_3]$

FM  
techniques

new  
concerns

heterogenized  
techniques

# Outline

- 1 Notation: Functions as diagrams
- 2 Related work: function unrolling and BPTT**
- 3 Main goals
- 4 Causal function formalization in category theory
  - Main ideas
  - Sanity checks
- 5 Delayed trace
- 6 Cartesian differential structure
- 7 Recap and future directions

# Backpropagation & backpropagation through time

*Backpropagation*: an efficient implementation for computing derivatives.

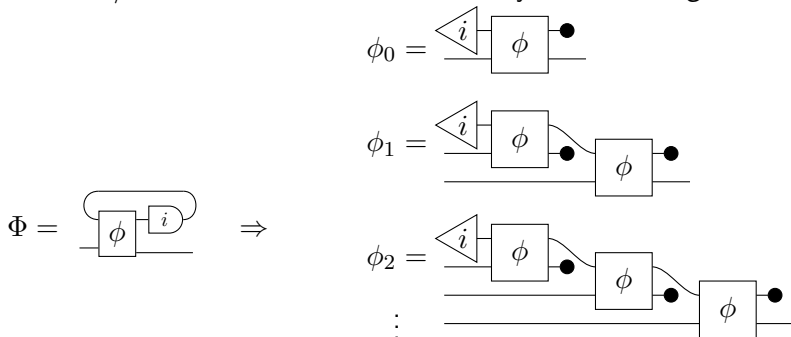
# Backpropagation & backpropagation through time

*Backpropagation*: an efficient implementation for computing derivatives. Works only for stateless/feedforward networks.

# Backpropagation & backpropagation through time

*Backpropagation*: an efficient implementation for computing derivatives. Works only for stateless/feedforward networks.

Stateful/recurrent networks are trained by first *unrolling*:

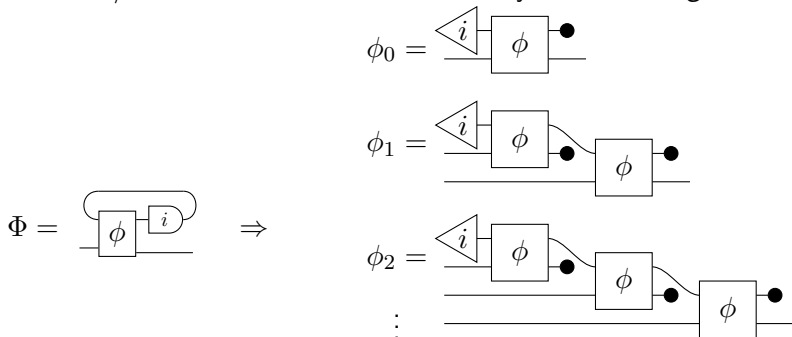




# Backpropagation & backpropagation through time

*Backpropagation*: an efficient implementation for computing derivatives. Works only for stateless/feedforward networks.

Stateful/recurrent networks are trained by first *unrolling*:



*Backpropagation through time (BPTT)*: Whenever the derivative of  $\Phi$  is needed at an input of length  $k + 1$ , the derivative of  $\phi_k$  is supplied instead.

**Does BPTT make sense,  
or is it just a hack?**

# Does BPTT make sense, or is it just a hack?

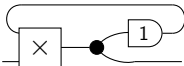
“Making sense” means having the usual properties of derivatives:

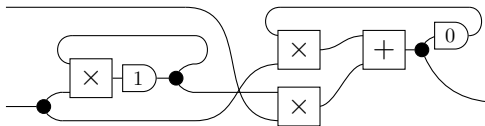
- 1 a sum rule,
  - 2 a chain rule,
  - 3 being linear when evaluated at any base point,
  - 4 symmetry of mixed partial derivatives,
- ⋮

# Outline

- 1 Notation: Functions as diagrams
- 2 Related work: function unrolling and BPTT
- 3 Main goals**
- 4 Causal function formalization in category theory
  - Main ideas
  - Sanity checks
- 5 Delayed trace
- 6 Cartesian differential structure
- 7 Recap and future directions

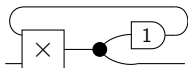
# Goals of this talk (specific to general)

A. The derivative of  is this monster

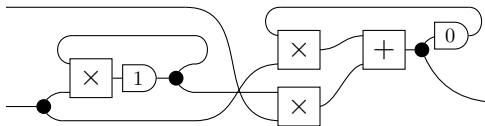


## Goals of this talk (specific to general)

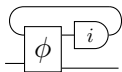
A. The derivative of



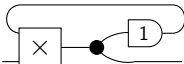
is this monster

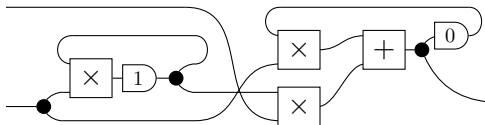


B. There is a general rule for derivatives of stateful functions

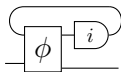


## Goals of this talk (specific to general)

A. The derivative of  is this monster

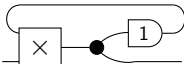


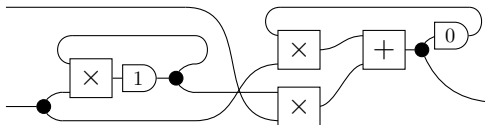
B. There is a general rule for derivatives of stateful functions



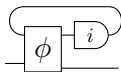
C. We understand many properties of this rule.

## Goals of this talk (specific to general)

A. The derivative of  is this monster



B. There is a general rule for derivatives of stateful functions



C. We understand many properties of this rule.

⋮

Z. Profit?? Get paper??



Goal Z & C: How do we get  $\mathbb{Y}\mathbb{Y}\mathbb{Y}$ ?

Goal Z & C: How do we get ¥¥¥?

**Answer:** Use it for machine learning.

## Goal Z & C: How do we get $\mathbb{Y}\mathbb{Y}\mathbb{Y}$ ?

**Answer:** Use it for machine learning.

A *neural network* looks like  $\begin{matrix} \theta \\ x \end{matrix} \rightarrow \boxed{N} \rightarrow y$ . Training a neural network means finding  $\theta^*$  so that:

$$\begin{matrix} \triangle \theta^* \\ \hat{x}_i \end{matrix} \rightarrow \boxed{N} \rightarrow y_i \approx \hat{y}_i$$

## Goal Z & C: How do we get $\mathbb{Y}\mathbb{Y}\mathbb{Y}$ ?

**Answer:** Use it for machine learning.

A *neural network* looks like  $\begin{matrix} \theta \\ x \end{matrix} \rightarrow \boxed{N} \rightarrow y$ . Training a neural network means finding  $\theta^*$  so that:

$$\begin{matrix} \triangle \theta^* \\ \hat{x}_i \end{matrix} \rightarrow \boxed{N} \rightarrow y_i \approx \hat{y}_i$$

*Gradient-based training* algorithms are based on the insight that  $\frac{\partial N}{\partial \theta}$  is a good approximation for the change in  $y$  that results from a small change in  $\theta$ . This allows us to make smart updates to  $\theta^*$ .

## Goal Z & C: How do we get $\forall\forall\forall$ ?

**Answer:** Use it for machine learning.

A *neural network* looks like  $\theta \begin{matrix} \text{---} \\ \text{---} \end{matrix} \boxed{N} \text{---} y$ . Training a neural network means finding  $\theta^*$  so that:

$$\hat{x}_i \begin{matrix} \triangle \theta^* \\ \text{---} \end{matrix} \boxed{N} \text{---} y_i \approx \hat{y}_i$$

*Gradient-based training* algorithms are based on the insight that  $\frac{\partial N}{\partial \theta}$  is a good approximation for the change in  $y$  that results from a small change in  $\theta$ . This allows us to make smart updates to  $\theta^*$ .

**Upshot:** If we can find (partial) derivatives of  $\theta \begin{matrix} \text{---} \\ \text{---} \end{matrix} \boxed{N} \begin{matrix} \text{---} \\ \text{---} \end{matrix} y$ , we can train recurrent neural networks. Using **properties** of these derivatives, maybe we can do it more efficiently than before.

## Goal B: Derivatives in diagram form

The derivative of  $x \text{ --- } \boxed{f} \text{ --- } y$  is  $\frac{\Delta x}{x} \text{ --- } \boxed{Df} \text{ --- } \Delta y$ , which gives the best linear approximation to  $\phi$  at a base point  $x$ .

## Goal B: Derivatives in diagram form

The derivative of  $x \rightarrow f \rightarrow y$  is  $\frac{\Delta x}{x} \rightarrow Df \rightarrow \Delta y$ , which gives the best linear approximation to  $f$  at a base point  $x$ . That is

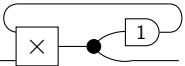
①  $\frac{\Delta x}{x} \rightarrow \boxed{+} \rightarrow \boxed{f} \rightarrow f(x + \Delta x) \approx \begin{array}{c} \Delta x \\ x \end{array} \rightarrow \begin{array}{c} \boxed{Df} \\ \boxed{f} \end{array} \rightarrow \boxed{+} \rightarrow \Delta y + f(x)$







## Goal A: Does the monster approximate output change?

Let  $\Phi =$   . Take  $\mathbf{x} = (1, 1, 1, \dots)$ , and  $\Delta \mathbf{x} = (0.1, 0.1, 0.1, \dots)$ . Then

$$\Phi(\mathbf{x}) = \Phi(1, 1, 1, \dots) = (1, 1, 1, \dots),$$

$$\Phi(\mathbf{x} + \Delta \mathbf{x}) = \Phi(1.1, 1.1, 1.1, \dots) = (1.1, 1.21, 1.331, \dots)$$

What is  $D\Phi((0.1, 1), (0.1, 1), (0.1, 1), \dots)$ ?

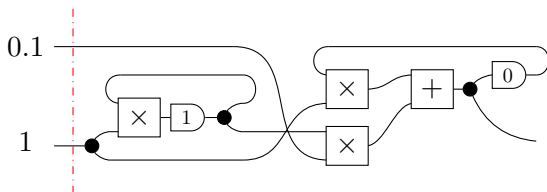
## Goal A: Does the monster approximate output change?

Let  $\Phi = \text{[Diagram]}$ . Take  $\mathbf{x} = (1, 1, 1, \dots)$ , and  $\Delta \mathbf{x} = (0.1, 0.1, 0.1, \dots)$ . Then

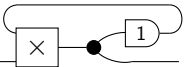
$$\Phi(\mathbf{x}) = \Phi(1, 1, 1, \dots) = (1, 1, 1, \dots),$$

$$\Phi(\mathbf{x} + \Delta \mathbf{x}) = \Phi(1.1, 1.1, 1.1, \dots) = (1.1, 1.21, 1.331, \dots)$$

What is  $D\Phi(\mathbf{0.1}, \mathbf{1}), (0.1, 1), (0.1, 1), \dots$ ?



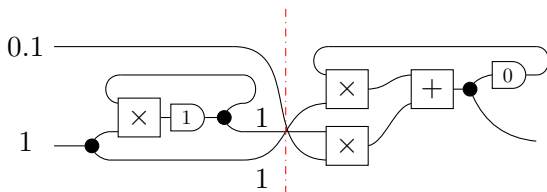
## Goal A: Does the monster approximate output change?

Let  $\Phi =$  . Take  $\mathbf{x} = (1, 1, 1, \dots)$ , and  $\Delta \mathbf{x} = (0.1, 0.1, 0.1, \dots)$ . Then

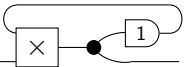
$$\Phi(\mathbf{x}) = \Phi(1, 1, 1, \dots) = (1, 1, 1, \dots),$$

$$\Phi(\mathbf{x} + \Delta \mathbf{x}) = \Phi(1.1, 1.1, 1.1, \dots) = (1.1, 1.21, 1.331, \dots)$$

What is  $D\Phi(\mathbf{0.1}, \mathbf{1}), (0.1, 1), (0.1, 1), \dots$ ?



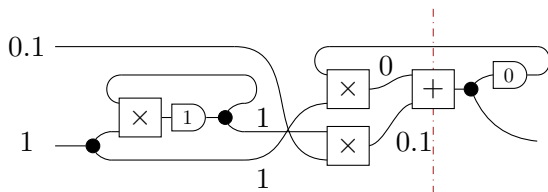
## Goal A: Does the monster approximate output change?

Let  $\Phi =$  . Take  $\mathbf{x} = (1, 1, 1, \dots)$ , and  $\Delta \mathbf{x} = (0.1, 0.1, 0.1, \dots)$ . Then

$$\Phi(\mathbf{x}) = \Phi(1, 1, 1, \dots) = (1, 1, 1, \dots),$$

$$\Phi(\mathbf{x} + \Delta \mathbf{x}) = \Phi(1.1, 1.1, 1.1, \dots) = (1.1, 1.21, 1.331, \dots)$$

What is  $D\Phi(\mathbf{0.1}, \mathbf{1}), (0.1, 1), (0.1, 1), \dots$ ?



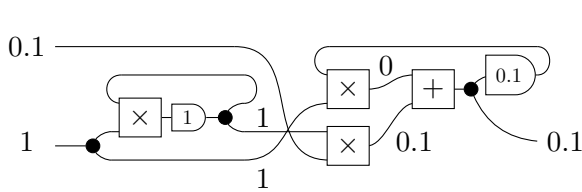
## Goal A: Does the monster approximate output change?

Let  $\Phi = \text{[Diagram]}$ . Take  $\mathbf{x} = (1, 1, 1, \dots)$ , and  $\Delta \mathbf{x} = (0.1, 0.1, 0.1, \dots)$ . Then

$$\Phi(\mathbf{x}) = \Phi(1, 1, 1, \dots) = (1, 1, 1, \dots),$$

$$\Phi(\mathbf{x} + \Delta \mathbf{x}) = \Phi(1.1, 1.1, 1.1, \dots) = (1.1, 1.21, 1.331, \dots)$$

$$D\Phi((\mathbf{0.1}, \mathbf{1}), (0.1, 1), (0.1, 1), \dots) = (\mathbf{0.1}, \dots).$$



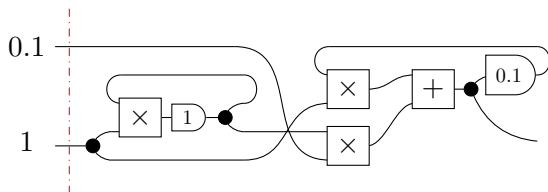
## Goal A: Does the monster approximate output change?

Let  $\Phi = \text{[Diagram]}$ . Take  $\mathbf{x} = (1, 1, 1, \dots)$ , and  $\Delta \mathbf{x} = (0.1, 0.1, 0.1, \dots)$ . Then

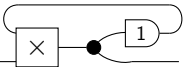
$$\Phi(\mathbf{x}) = \Phi(1, 1, 1, \dots) = (1, 1, 1, \dots),$$

$$\Phi(\mathbf{x} + \Delta \mathbf{x}) = \Phi(1.1, 1.1, 1.1, \dots) = (1.1, 1.21, 1.331, \dots)$$

$$D\Phi((0.1, 1), (\mathbf{0.1}, \mathbf{1}), (0.1, 1), \dots) = (0.1, \dots).$$



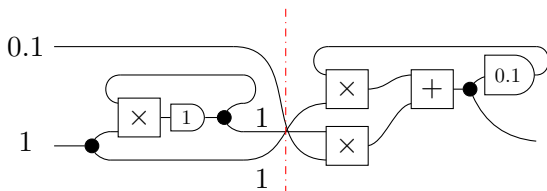
## Goal A: Does the monster approximate output change?

Let  $\Phi =$  . Take  $\mathbf{x} = (1, 1, 1, \dots)$ , and  $\Delta \mathbf{x} = (0.1, 0.1, 0.1, \dots)$ . Then

$$\Phi(\mathbf{x}) = \Phi(1, 1, 1, \dots) = (1, 1, 1, \dots),$$

$$\Phi(\mathbf{x} + \Delta \mathbf{x}) = \Phi(1.1, 1.1, 1.1, \dots) = (1.1, 1.21, 1.331, \dots)$$

$$D\Phi((0.1, 1), (\mathbf{0.1}, \mathbf{1}), (0.1, 1), \dots) = (0.1, \dots).$$





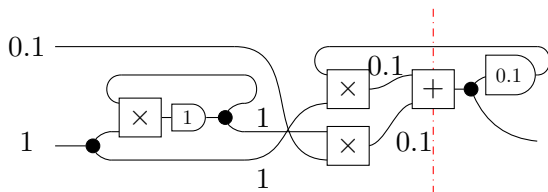
## Goal A: Does the monster approximate output change?

Let  $\Phi = \text{[Diagram]}$ . Take  $\mathbf{x} = (1, 1, 1, \dots)$ , and  $\Delta \mathbf{x} = (0.1, 0.1, 0.1, \dots)$ . Then

$$\Phi(\mathbf{x}) = \Phi(1, 1, 1, \dots) = (1, 1, 1, \dots),$$

$$\Phi(\mathbf{x} + \Delta \mathbf{x}) = \Phi(1.1, 1.1, 1.1, \dots) = (1.1, 1.21, 1.331, \dots)$$

$$D\Phi((0.1, 1), (\mathbf{0.1}, \mathbf{1}), (0.1, 1), \dots) = (0.1, \dots).$$



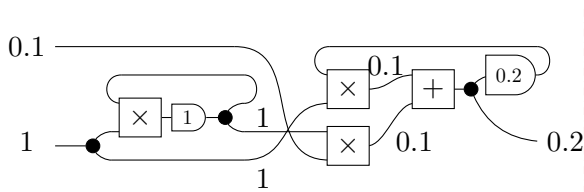
## Goal A: Does the monster approximate output change?

Let  $\Phi = \text{[Diagram: a box with } \times \text{, a box with } 1 \text{, and a black dot with a loop labeled } 1 \text{]} \text{. Take } \mathbf{x} = (1, 1, 1, \dots)$ , and  $\Delta \mathbf{x} = (0.1, 0.1, 0.1, \dots)$ . Then

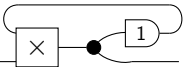
$$\Phi(\mathbf{x}) = \Phi(1, 1, 1, \dots) = (1, 1, 1, \dots),$$

$$\Phi(\mathbf{x} + \Delta \mathbf{x}) = \Phi(1.1, 1.1, 1.1, \dots) = (1.1, 1.21, 1.331, \dots)$$

$$D\Phi((0.1, 1), (\mathbf{0.1}, \mathbf{1}), (0.1, 1), \dots) = (0.1, \mathbf{0.2}, \dots).$$



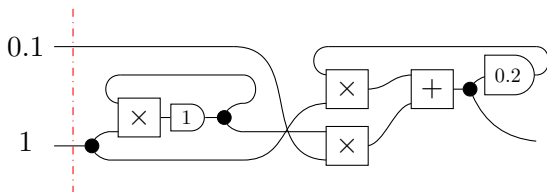
## Goal A: Does the monster approximate output change?

Let  $\Phi =$  . Take  $\mathbf{x} = (1, 1, 1, \dots)$ , and  $\Delta \mathbf{x} = (0.1, 0.1, 0.1, \dots)$ . Then

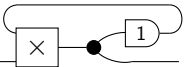
$$\Phi(\mathbf{x}) = \Phi(1, 1, 1, \dots) = (1, 1, 1, \dots),$$

$$\Phi(\mathbf{x} + \Delta \mathbf{x}) = \Phi(1.1, 1.1, 1.1, \dots) = (1.1, 1.21, 1.331, \dots)$$

$$D\Phi((0.1, 1), (0.1, 1), (\mathbf{0.1}, \mathbf{1}), \dots) = (0.1, 0.2, \dots).$$



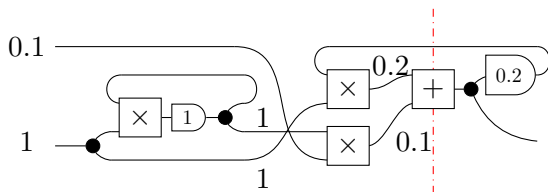
## Goal A: Does the monster approximate output change?

Let  $\Phi =$  . Take  $\mathbf{x} = (1, 1, 1, \dots)$ , and  $\Delta \mathbf{x} = (0.1, 0.1, 0.1, \dots)$ . Then

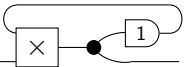
$$\Phi(\mathbf{x}) = \Phi(1, 1, 1, \dots) = (1, 1, 1, \dots),$$

$$\Phi(\mathbf{x} + \Delta \mathbf{x}) = \Phi(1.1, 1.1, 1.1, \dots) = (1.1, 1.21, 1.331, \dots)$$

$$D\Phi((0.1, 1), (0.1, 1), (\mathbf{0.1}, \mathbf{1}), \dots) = (0.1, 0.2, \dots).$$



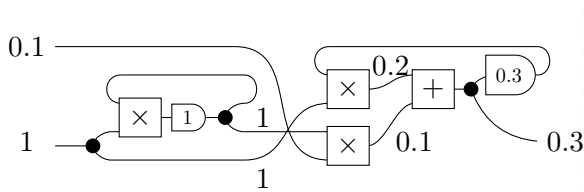
## Goal A: Does the monster approximate output change?

Let  $\Phi =$  . Take  $\mathbf{x} = (1, 1, 1, \dots)$ , and  $\Delta \mathbf{x} = (0.1, 0.1, 0.1, \dots)$ . Then

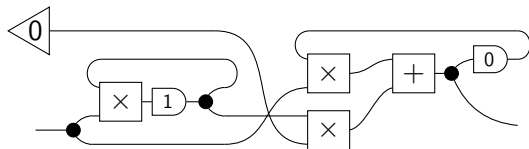
$$\Phi(\mathbf{x}) = \Phi(1, 1, 1, \dots) = (1, 1, 1, \dots),$$

$$\Phi(\mathbf{x} + \Delta \mathbf{x}) = \Phi(1.1, 1.1, 1.1, \dots) = (1.1, 1.21, 1.331, \dots)$$

$$D\Phi((0.1, 1), (0.1, 1), (\mathbf{0.1}, \mathbf{1}), \dots) = (0.1, 0.2, \mathbf{0.3}, \dots).$$

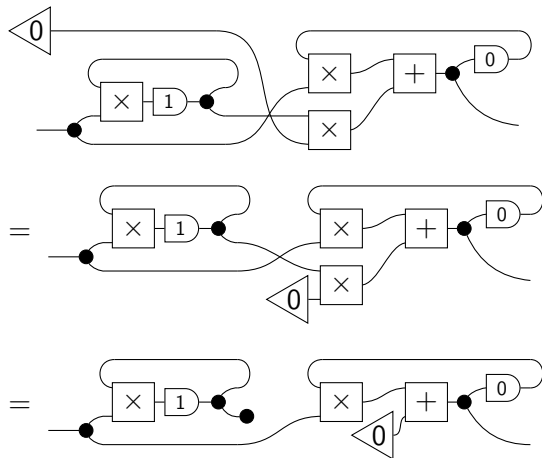


Goal A: Is  $D\Phi(0, \mathbf{x}) = 0$ ?



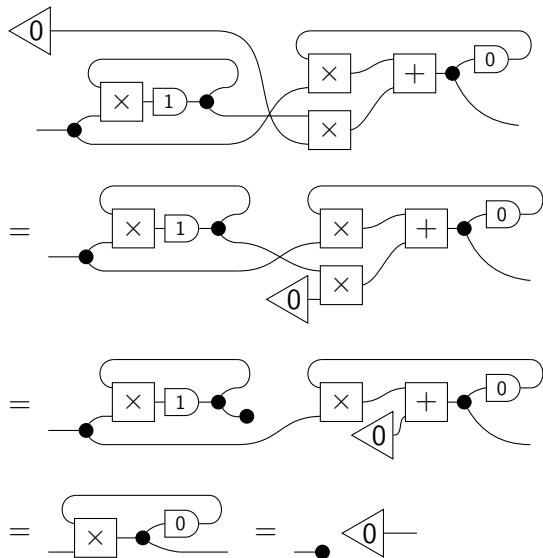


Goal A: Is  $D\Phi(0, \mathbf{x}) = 0$ ?





Goal A: Is  $D\Phi(0, \mathbf{x}) = 0$ ?



# Outline

- 1 Notation: Functions as diagrams
- 2 Related work: function unrolling and BPTT
- 3 Main goals
- 4 Causal function formalization in category theory**
  - Main ideas
  - Sanity checks
- 5 Delayed trace
- 6 Cartesian differential structure
- 7 Recap and future directions

## Main ideas

We imagine that we are starting from a (strictified Cartesian) category  $\mathbb{C}$  whose morphisms represent single computations, which we aim to extend to computations on sequences.

## Main ideas

We imagine that we are starting from a (strictified Cartesian) category  $\mathbb{C}$  whose morphisms represent single computations, which we aim to extend to computations on sequences.

We want to capture *causal* functions on sequences, meaning the  $n$ th element of the output sequence depends only on the first  $n$  elements of the input sequence.

## Main ideas

We imagine that we are starting from a (strictified Cartesian) category  $\mathbb{C}$  whose morphisms represent single computations, which we aim to extend to computations on sequences.

We want to capture *causal* functions on sequences, meaning the  $n$ th element of the output sequence depends only on the first  $n$  elements of the input sequence.

There are three steps:

- 1 describe a single step of the sequence computation, including a mechanism for sending and receiving state,
- 2 chain these single steps together to get a full sequence, and
- 3 quotient these computations by their observable behaviour.

## Step 1: Single computation steps

We separate inputs and outputs of a  $\mathbb{C}$ -morphism into two types: *values*, exchanged with the environment, and *states*, received from (sent to) the previous (next) step of the computation.

## Step 1: Single computation steps

We separate inputs and outputs of a  $\mathbb{C}$ -morphism into two types: *values*, exchanged with the environment, and *states*, received from (sent to) the previous (next) step of the computation.

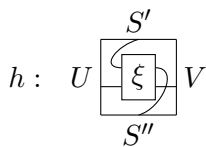
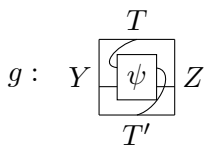
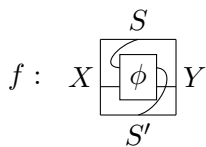
### Definition

A *computation step*  $f : X \xrightarrow[S']{S} Y$  is four objects and a morphism:

- 1  $S$  — the input state received from the previous step
- 2  $X$  — the input value received from the environment
- 3  $S'$  — the output state sent to the next step
- 4  $Y$  — the output value sent to the environment
- 5  $\phi : S \times X \rightarrow S' \times Y$  — the  $\mathbb{C}$ -morphism doing the computation

## Step 1: Single computation steps

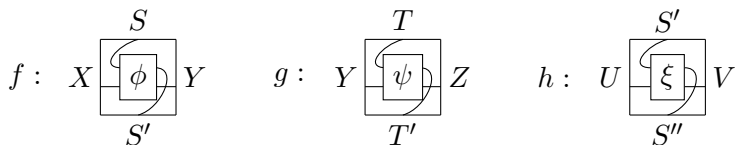
We can draw a single step like so:



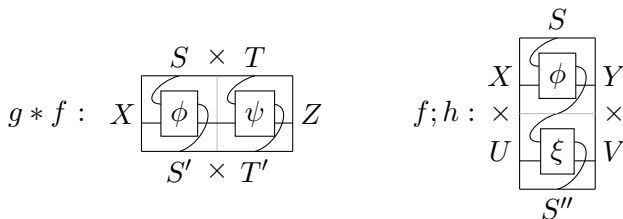


## Step 1: Single computation steps

We can draw a single step like so:

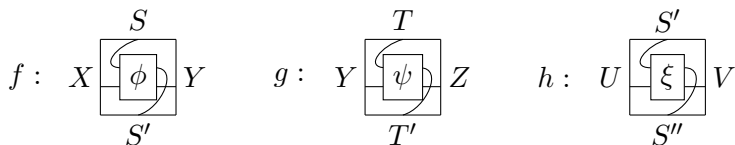


Computation steps can be composed in two different ways:

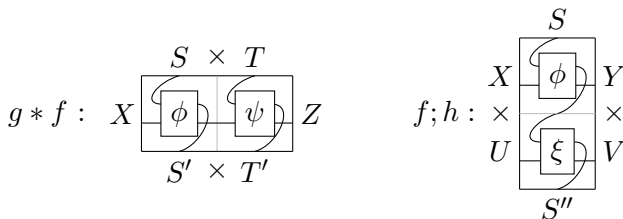


## Step 1: Single computation steps

We can draw a single step like so:



Computation steps can be composed in two different ways:



psst ... double category

## Step 2: Computation sequences

Next, we chain together infinitely many of these computation steps to compute a whole sequence of outputs.

### Definition

A *computation sequence*  $\mathbf{f} = (i, [f_k])$  is an infinite sequence of computation steps  $f_k : X_k \xrightarrow[S_{k+1}]{S_k} Y_k$ , and an initial state  $i : 1 \xrightarrow[S_0]{1} 1$ .

We say  $\mathbf{f}$  takes sequences of type  $[X_k] = (X_0, X_1, \dots, X_n, \dots)$  to sequences of type  $[Y_k] = (Y_0, Y_1, \dots, Y_n, \dots)$ .

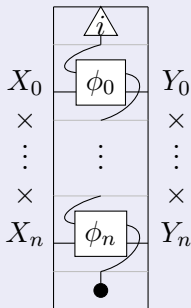
Note that  $f_k$  and  $f_{k+1}$  are vertically composable, as are  $i$  and  $f_0$ .

## Step 3: Comparing computation sequences

Two computation sequences might have different state spaces and still compute the same function. We need to make some identifications. . .

### Definition

The  $n$ th *truncation* of a computation sequence is the morphism of the vertical composite of the first  $n + 1$  steps:



## Step 3: Comparing computation sequences

### Definition

Two computation sequences  $\mathbf{f}, \mathbf{g} : [X_k] \rightarrow [Y_k]$  are *extensionally equivalent* means they have the same  $n$ th truncation for all  $n \in \mathbb{N}$ .

### Definition

A *stateful (sequence) function* is an extensional equivalence class of computation sequences.

### Definition

If  $\mathbb{C}$  is a (strict) Cartesian category, then its *stateful extension* is a category  $\text{St}(\mathbb{C})$  where

- objects are  $[X_k]$ , i.e. infinite sequences of objects in  $\mathbb{C}$  and
- morphisms are stateful functions  $\mathbf{f} : [X_k] \rightarrow [Y_k]$ .

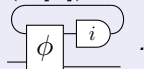
# Sanity checks

## Theorem

*The homset  $\text{St}(\text{Set})([A], [B])$  is in 1-1 correspondence with the set of causal functions from  $A^\omega$  to  $B^\omega$ .*

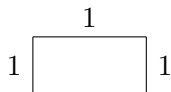
## Theorem

*The sequence of  $n$ th truncations of  $(i, [\phi])$ , each projected to their last component, is the unrolling of*




## Fill-in-the-blank

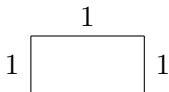
How would you implement  $\boxed{i}$  as a computation sequence?



⋮

## Fill-in-the-blank

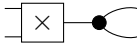
How would you implement  
 as a computation  
sequence?

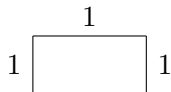


⋮

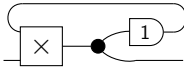


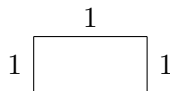
# Fill-in-the-blank

How would you implement  as a computation sequence?



⋮

How would you implement  as a computation sequence?



⋮

# Outline

- 1 Notation: Functions as diagrams
- 2 Related work: function unrolling and BPTT
- 3 Main goals
- 4 Causal function formalization in category theory
  - Main ideas
  - Sanity checks
- 5 Delayed trace**
- 6 Cartesian differential structure
- 7 Recap and future directions

## Delayed trace

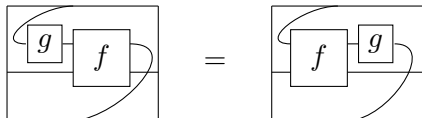
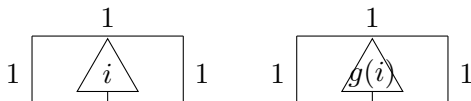
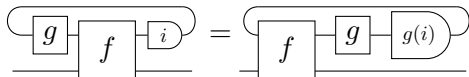
This loop-with-delay-gate is a trace-like operation.

$$\frac{\phi : S \times X \rightarrow S \times Y \quad \text{---} \boxed{\phi} \text{---}}{dtr_i^S(\phi) : X \rightarrow Y \quad \text{---} \boxed{\phi} \text{---} \text{---} \boxed{i} \text{---}}$$

It satisfies **most** of the the trace axioms (source & target naturality, vanishing 1 and  $\times$ , and superposing) but misses two: yanking and dinaturality. For regular trace, those are

$$\text{---} \text{---} \text{---} \text{---} \text{---} = \text{---} \text{---} \quad \text{---} \boxed{g} \boxed{f} \text{---} \text{---} = \text{---} \boxed{f} \boxed{g} \text{---} \text{---}$$

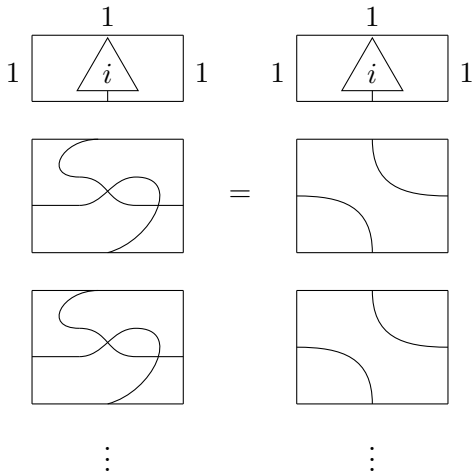
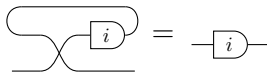
# Dinaturality $\rightarrow$ retiming



⋮

⋮

# Yanking $\rightarrow$ delay



## Advantages of delayed trace

Often, people formalizing circuits in category theory do these:

- 1 (freely) add registers to the category,
- 2 add trace to the category,
- 3 restrict to a subcategory such that all traces are taken on positions guarded by a register.

## Advantages of delayed trace

Often, people formalizing circuits in category theory do these:

- 1 (freely) add registers to the category,
- 2 add trace to the category,
- 3 restrict to a subcategory such that all traces are taken on positions guarded by a register.

We suggest instead to:

- 1 add a delayed trace to the category,
- 2 recover registers by delay-tracing symmetry

This ensures all loops are guarded (without a syntactic restriction!) and requires less structure.

## Advantages of delayed trace

Often, people formalizing circuits in category theory do these:

- 1 (freely) add registers to the category,
- 2 add trace to the category,
- 3 restrict to a subcategory such that all traces are taken on positions guarded by a register.

We suggest instead to:

- 1 add a delayed trace to the category,
- 2 recover registers by delay-tracing symmetry

This ensures all loops are guarded (without a syntactic restriction!) and requires less structure.

We are also interested in looking at delayed traces abstractly. . .



# Outline

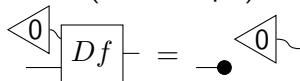
- 1 Notation: Functions as diagrams
- 2 Related work: function unrolling and BPTT
- 3 Main goals
- 4 Causal function formalization in category theory
  - Main ideas
  - Sanity checks
- 5 Delayed trace
- 6 Cartesian differential structure**
- 7 Recap and future directions

# Cartesian differential categories [Blute, Cockett, Seely '09]

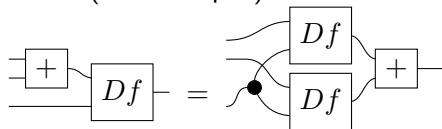
A *Cartesian differential category* has a differential operation on morphisms sending  $f : X \rightarrow Y$  to  $Df : X \times X \rightarrow Y$ , satisfying:

**CD1.**  $Ds = s \times !_{\text{dom}(s)}$  for  $s \in \{\text{id}_X, \sigma_{X,Y}, !_X, \Delta_X, 0_X, +_X\}$ .

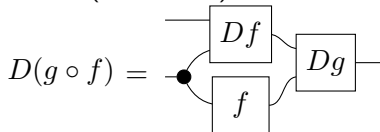
**CD2.** (linear map I)



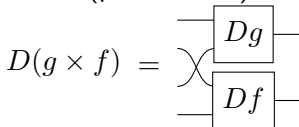
**CD3.** (linear map II)



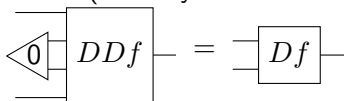
**CD4.** (chain rule)



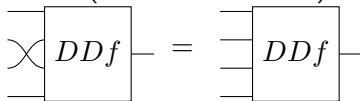
**CD5.** (parallel rule)



**CD6.** (linearity of 2<sup>nd</sup> deriv.)



**CD7.** (Schwartz theorem)



## $\text{Euc}_\infty$ is Cartesian differential

The differential operator on Euclidean spaces and smooth maps sends  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  to  $Df : (\Delta x, x) \mapsto Jf|_x \times \Delta x$ .

## $\text{Euc}_\infty$ is Cartesian differential

The differential operator on Euclidean spaces and smooth maps sends  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  to  $Df : (\Delta x, x) \mapsto Jf|_x \times \Delta x$ .

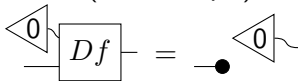
**CD1.**  $Ds(\Delta x, x) = s(\Delta x)$  for  $s \in \{\text{id}_X, \sigma_{X,Y}, !_X, \Delta_X, 0_X, +_X\}$ .

## $\text{Euc}_\infty$ is Cartesian differential

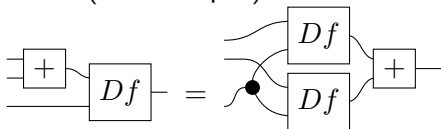
The differential operator on Euclidean spaces and smooth maps sends  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  to  $Df : (\Delta x, x) \mapsto Jf|_x \times \Delta x$ .

**CD1.**  $Ds(\Delta x, x) = s(\Delta x)$  for  $s \in \{\text{id}_X, \sigma_{X,Y}, !_X, \Delta_X, 0_X, +_X\}$ .

**CD2.** (linear map I)



**CD3.** (linear map II)

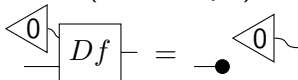


# $\text{Euc}_\infty$ is Cartesian differential

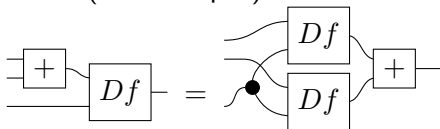
The differential operator on Euclidean spaces and smooth maps sends  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  to  $Df : (\Delta x, x) \mapsto Jf|_x \times \Delta x$ .

**CD1.**  $Ds(\Delta x, x) = s(\Delta x)$  for  $s \in \{\text{id}_X, \sigma_{X,Y}, !_X, \Delta_X, 0_X, +_X\}$ .

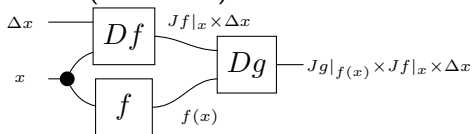
**CD2.** (linear map I)



**CD3.** (linear map II)



**CD4.** (chain rule)

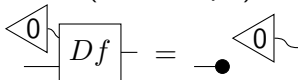


# $\text{Euc}_\infty$ is Cartesian differential

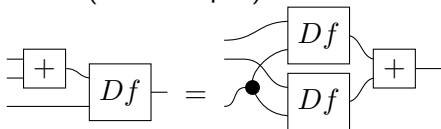
The differential operator on Euclidean spaces and smooth maps sends  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  to  $Df : (\Delta x, x) \mapsto Jf|_x \times \Delta x$ .

**CD1.**  $Ds(\Delta x, x) = s(\Delta x)$  for  $s \in \{\text{id}_X, \sigma_{X,Y}, !_X, \Delta_X, 0_X, +_X\}$ .

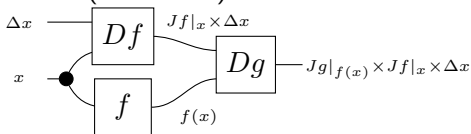
**CD2.** (linear map I)



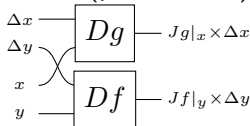
**CD3.** (linear map II)



**CD4.** (chain rule)



**CD5.** (parallel rule)



## Second derivative axioms in $\mathbf{Euc}_\infty$

For  $f : \mathbb{R} \rightarrow \mathbb{R}$ ,  $Df(a, b) = f'(b) \cdot a$ .



## Second derivative axioms in $\mathbf{Euc}_\infty$

For  $f : \mathbb{R} \rightarrow \mathbb{R}$ ,  $Df(a, b) = f'(b) \cdot a$ .

$$\begin{aligned} DDf(\Delta a, \Delta b, a, b) &= \frac{\partial}{\partial a}(f'(b) \cdot a) \cdot \Delta a + \frac{\partial}{\partial b}(f'(b) \cdot a) \cdot \Delta b \\ &= f'(b) \cdot \Delta a + f''(b) \cdot a \cdot \Delta b \end{aligned}$$

## Second derivative axioms in $\mathbf{Euc}_\infty$

For  $f : \mathbb{R} \rightarrow \mathbb{R}$ ,  $Df(a, b) = f'(b) \cdot a$ .

$$\begin{aligned} DDf(\Delta a, \Delta b, a, b) &= \frac{\partial}{\partial a}(f'(b) \cdot a) \cdot \Delta a + \frac{\partial}{\partial b}(f'(b) \cdot a) \cdot \Delta b \\ &= f'(b) \cdot \Delta a + f''(b) \cdot a \cdot \Delta b \end{aligned}$$

Note two things:

- 6 If  $a = \Delta b = 0$ , the expression is  $f'(b) \cdot \Delta a = Df(\Delta a, b)$ .
- 7 The expression is the same if  $a$  and  $\Delta b$  are swapped.

## Second derivative axioms in $\mathbf{Euc}_\infty$

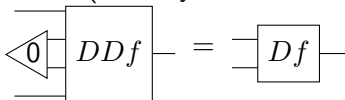
For  $f : \mathbb{R} \rightarrow \mathbb{R}$ ,  $Df(a, b) = f'(b) \cdot a$ .

$$\begin{aligned} DDf(\Delta a, \Delta b, a, b) &= \frac{\partial}{\partial a}(f'(b) \cdot a) \cdot \Delta a + \frac{\partial}{\partial b}(f'(b) \cdot a) \cdot \Delta b \\ &= f'(b) \cdot \Delta a + f''(b) \cdot a \cdot \Delta b \end{aligned}$$

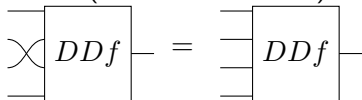
Note two things:

- 6 If  $a = \Delta b = 0$ , the expression is  $f'(b) \cdot \Delta a = Df(\Delta a, b)$ .
- 7 The expression is the same if  $a$  and  $\Delta b$  are swapped.

**CD6.** (linearity of 2<sup>nd</sup> deriv.)

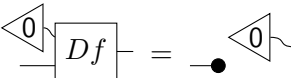


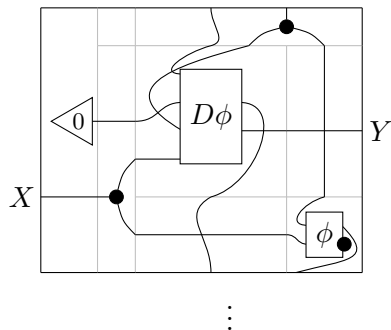
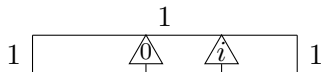
**CD7.** (Schwartz theorem)



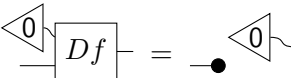


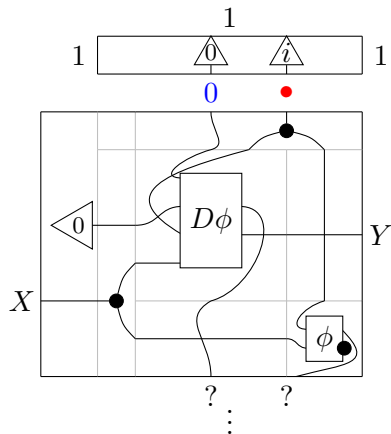
# Key contribution: differential operator lifts

**Proof idea.** For CD2: 

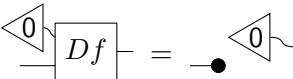


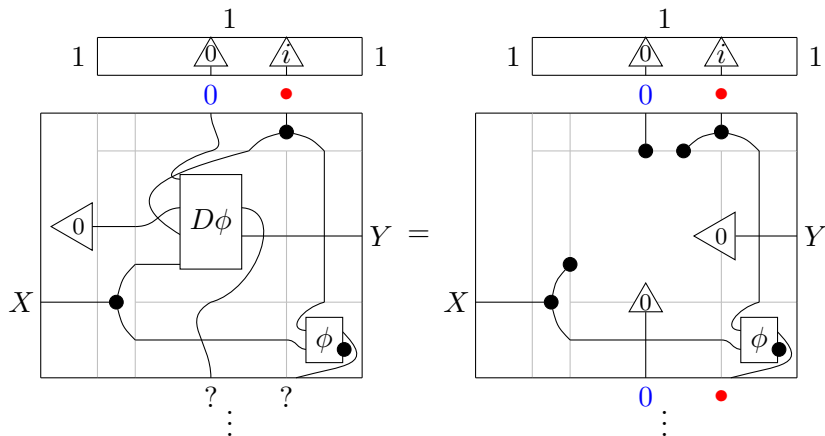
# Key contribution: differential operator lifts

**Proof idea.** For CD2: 



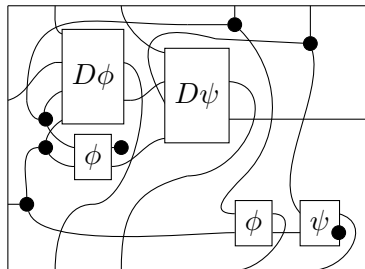
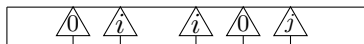
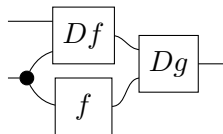
# Key contribution: differential operator lifts

**Proof idea.** For CD2: 

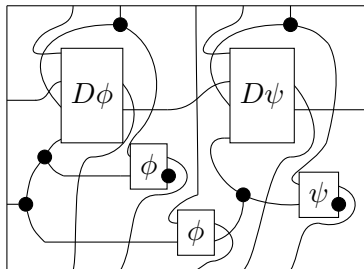


# Key contribution: differential operator lifts

**Proof idea.** For CD4:  $D(g \circ f) =$



=



⋮

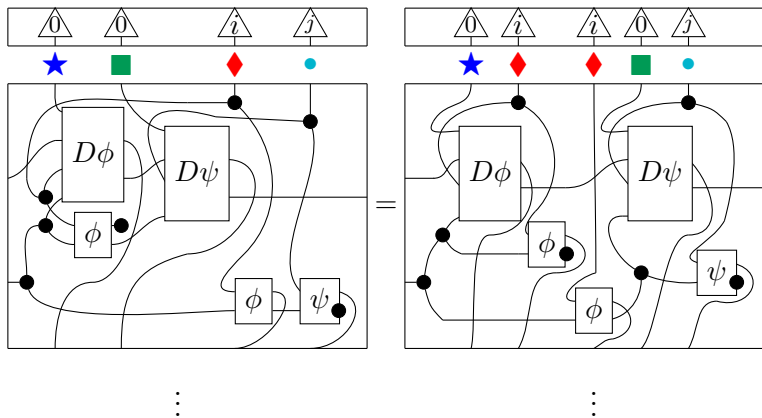
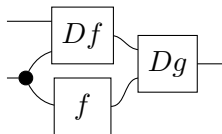
⋮



# Key contribution: differential operator lifts

**Proof idea.** For CD4:

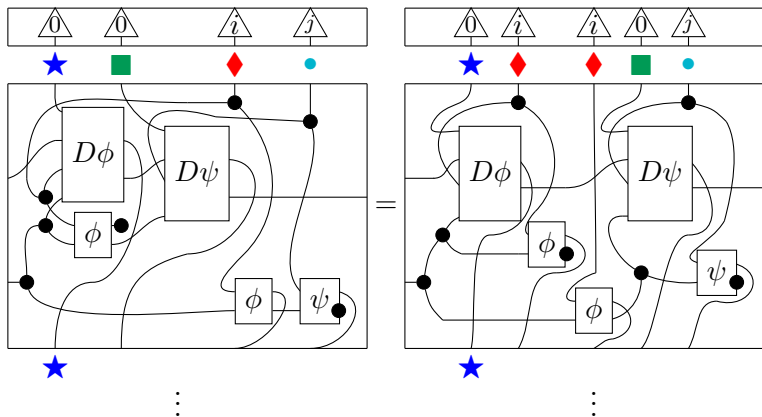
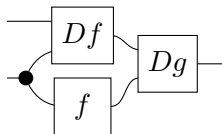
$$D(g \circ f) =$$



# Key contribution: differential operator lifts

**Proof idea.** For CD4:

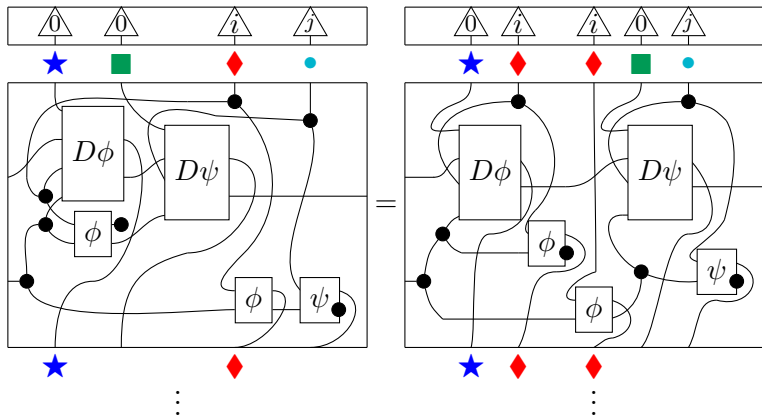
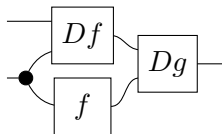
$$D(g \circ f) =$$



# Key contribution: differential operator lifts

**Proof idea.** For CD4:

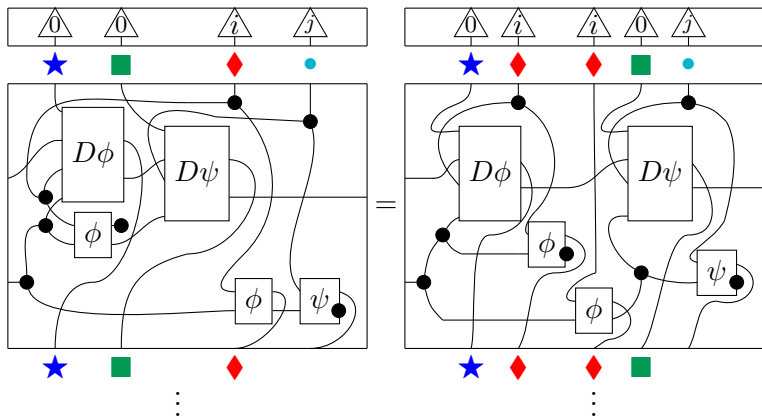
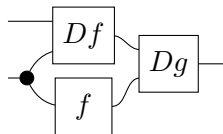
$$D(g \circ f) =$$



# Key contribution: differential operator lifts

**Proof idea.** For CD4:

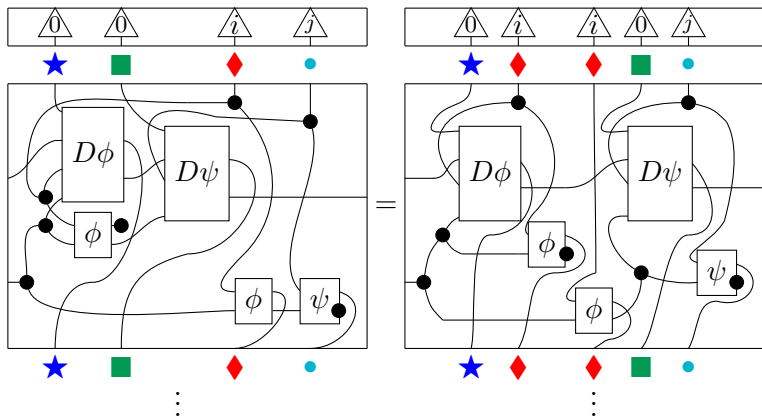
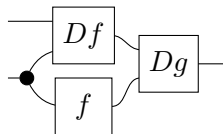
$$D(g \circ f) =$$



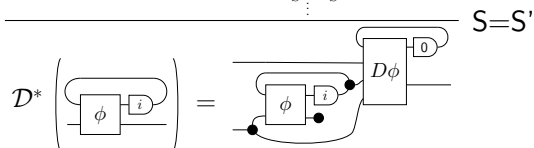
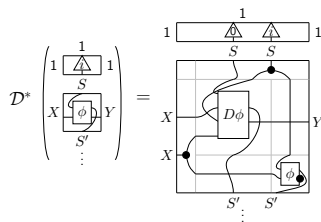
# Key contribution: differential operator lifts

**Proof idea.** For CD4:

$$D(g \circ f) =$$



# Differentiating circuits



## Theorem ( $\mathcal{D}^*$ matches BPTT)

*The unrolling of  $\mathcal{D}^*((i, [\phi]))$  is the component-wise application of  $D$  to the unrolling of  $(i, [\phi])$ .*



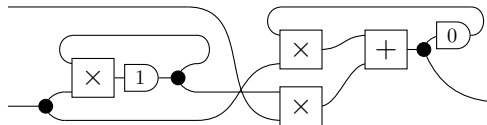
What if  $\phi = \Delta_{\mathbb{R}} \circ \times$ ?

$$D\left(\begin{array}{c} \text{---} \\ | \\ \boxed{\times} \\ | \\ \text{---} \end{array}\right) = \begin{array}{c} \text{---} \\ | \\ \boxed{\times} \\ | \\ \text{---} \\ | \\ \boxed{\times} \\ | \\ \text{---} \end{array} \rightarrow \boxed{+} \text{---}, \text{ so } D\left(\begin{array}{c} \text{---} \\ | \\ \boxed{\times} \\ | \\ \bullet \\ \text{---} \end{array}\right) = \begin{array}{c} \text{---} \\ | \\ \boxed{\times} \\ | \\ \text{---} \\ | \\ \boxed{\times} \\ | \\ \text{---} \end{array} \rightarrow \boxed{+} \bullet$$

Then

$$\mathcal{D}^*\left(\begin{array}{c} \text{---} \\ | \\ \boxed{\phi} \\ | \\ \text{---} \end{array}\right) = \begin{array}{c} \text{---} \\ | \\ \boxed{\phi} \\ | \\ \text{---} \end{array} \rightarrow \begin{array}{c} \text{---} \\ | \\ \boxed{\phi} \\ | \\ \text{---} \end{array} \rightarrow \boxed{D\phi} \begin{array}{c} \text{---} \\ | \\ \text{---} \end{array}$$

becomes



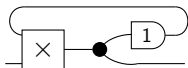


# Outline

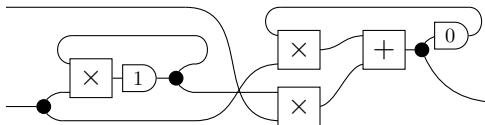
- 1 Notation: Functions as diagrams
- 2 Related work: function unrolling and BPTT
- 3 Main goals
- 4 Causal function formalization in category theory
  - Main ideas
  - Sanity checks
- 5 Delayed trace
- 6 Cartesian differential structure
- 7 Recap and future directions

## Goals of this talk, again

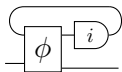
A. The derivative of



is this monster



B. There is a general rule for derivatives of stateful functions



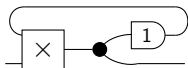
C. We understand many properties of this rule.

⋮

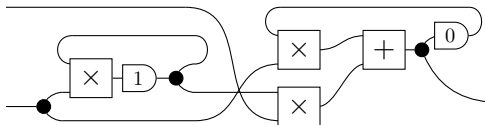
Z. Profit?? Get paper??

## Goals of this talk, again

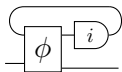
A. The derivative of



is this monster



B. There is a general rule for derivatives of stateful functions



C. We understand many properties of this rule.

⋮

Z. Profit?? Get paper?? Write papers??

## Future questions (roughly mathematical to practical)

- 1 Categorical properties of  $\text{St}(-)$ ?
- 2  $\mathcal{D}^*$  and derivatives in sequence spaces?
- 3 Bisimulations and extensional equality?
- 4 Basic results for delayed trace categories?

## Future questions (roughly mathematical to practical)

- 1 Categorical properties of  $\text{St}(-)$ ?
- 2  $\mathcal{D}^*$  and derivatives in sequence spaces?
- 3 Bisimulations and extensional equality?
- 4 Basic results for delayed trace categories?
- 5 Axiomatization of stateful function diagrams?
- 6 \* Explicit Jacobians using closed structure?
- 7 Probabilistic/nondeterministic causal functions?
- 8 \* Partial functions with differential restriction categories?

# Future questions (roughly mathematical to practical)

- 1 Categorical properties of  $\text{St}(-)$ ?
- 2  $\mathcal{D}^*$  and derivatives in sequence spaces?
- 3 Bisimulations and extensional equality?
- 4 Basic results for delayed trace categories?
- 5 Axiomatization of stateful function diagrams?
- 6 \* Explicit Jacobians using closed structure?
- 7 Probabilistic/nondeterministic causal functions?
- 8 \* Partial functions with differential restriction categories?
- 9 Iterations I: RNN parameters in  $\text{St}(\text{St}(\mathbb{C}))$ ?
- 10 Iterations II: RNN hyperparameters in  $\text{St}(\text{St}(\text{St}(\mathbb{C})))$ ?

## Future questions (roughly mathematical to practical)

- 1 Categorical properties of  $\text{St}(-)$ ?
- 2  $\mathcal{D}^*$  and derivatives in sequence spaces?
- 3 Bisimulations and extensional equality?
- 4 Basic results for delayed trace categories?
- 5 Axiomatization of stateful function diagrams?
- 6 \* Explicit Jacobians using closed structure?
- 7 Probabilistic/nondeterministic causal functions?
- 8 \* Partial functions with differential restriction categories?
- 9 Iterations I: RNN parameters in  $\text{St}(\text{St}(\mathbb{C}))$ ?
- 10 Iterations II: RNN hyperparameters in  $\text{St}(\text{St}(\text{St}(\mathbb{C})))$ ?
- 11 \* Measuring complexity with string diagrams?
- 12 What happens when  $\phi \in \{LSTM, GRU, \dots\}$ ?
- 13 Implementation/plugin to neural net library?

Thanks!