

Differentiation of causal functions

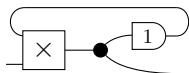
Shin-ya Katsumata and David Sprunger*
National Institute of Informatics
ERATO MMSD

IU Logic Colloquium
March 21, 2019

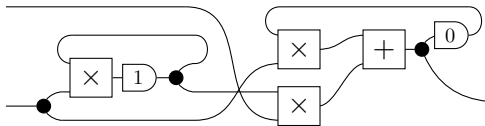


Where we're heading

I would like to convince you the derivative of this function:



Is this:



How we get there

- 1 Causal functions
- 2 Computation sequences
 - Main ideas
 - Sanity check
- 3 Delayed trace
- 4 Neural networks
- 5 Cartesian differential categories
- 6 Recap and future directions

Causal functions on sequences

A^ω is the set of A -valued infinite sequences. The entries of $\sigma \in A^\omega$ are $[\sigma]_k = \sigma_k \in A$ for $k \in \mathbb{N}$, so $\sigma = (\sigma_0, \sigma_1, \dots, \sigma_k, \dots)$.

Causal functions on sequences

A^ω is the set of A -valued infinite sequences. The entries of $\sigma \in A^\omega$ are $[\sigma]_k = \sigma_k \in A$ for $k \in \mathbb{N}$, so $\sigma = (\sigma_0, \sigma_1, \dots, \sigma_k, \dots)$.

Slicing extracts a finite list from an infinite sequence:

$$[\cdot]_{i:j} : \sigma \mapsto (\sigma_i, \sigma_{i+1}, \dots, \sigma_j)$$

Causal functions on sequences

A^ω is the set of A -valued infinite sequences. The entries of $\sigma \in A^\omega$ are $[\sigma]_k = \sigma_k \in A$ for $k \in \mathbb{N}$, so $\sigma = (\sigma_0, \sigma_1, \dots, \sigma_k, \dots)$.

Slicing extracts a finite list from an infinite sequence:

$$[\cdot]_{i:j} : \sigma \mapsto (\sigma_i, \sigma_{i+1}, \dots, \sigma_j)$$

Definition

A *causal* function $f : A^\omega \rightarrow B^\omega$ satisfies

$$\forall \sigma, \tau \in A^\omega, \forall k \in \mathbb{N}, [\sigma]_{0:k} = [\tau]_{0:k} \rightarrow [f(\sigma)]_{0:k} = [f(\tau)]_{0:k}$$

Intuitively, the first k outputs of f only depend on the first k inputs.

Examples

1. $\bar{b} : 1^\omega \rightarrow B^\omega$ defined by $[\bar{b}]_k = b$ is causal.

Examples

1. $\bar{b} : 1^\omega \rightarrow B^\omega$ defined by $[\bar{b}]_k = b$ is causal.
2. If $g : A \rightarrow B$, then $\text{map}(g) : A^\omega \rightarrow B^\omega$ defined by $[\text{map}(g)(\sigma)]_k = g(\sigma_k)$ is causal.

Examples

1. $\bar{b} : 1^\omega \rightarrow B^\omega$ defined by $[\bar{b}]_k = b$ is causal.
2. If $g : A \rightarrow B$, then $\text{map}(g) : A^\omega \rightarrow B^\omega$ defined by $[\text{map}(g)(\sigma)]_k = g(\sigma_k)$ is causal.
3. $\text{runProd} : \mathbb{R}^\omega \rightarrow \mathbb{R}^\omega$ defined by $[\text{runProd}(\sigma)]_k = \prod_{i=0}^k \sigma_i$ is causal.

Examples

1. $\bar{b} : 1^\omega \rightarrow B^\omega$ defined by $[\bar{b}]_k = b$ is causal.
2. If $g : A \rightarrow B$, then $\text{map}(g) : A^\omega \rightarrow B^\omega$ defined by $[\text{map}(g)(\sigma)]_k = g(\sigma_k)$ is causal.
3. $\text{runProd} : \mathbb{R}^\omega \rightarrow \mathbb{R}^\omega$ defined by $[\text{runProd}(\sigma)]_k = \prod_{i=0}^k \sigma_i$ is causal.
4. $\text{zip} : (A \times A)^\omega \rightarrow A^\omega$ defined by
$$[\text{zip}(\sigma, \tau)]_k = \begin{cases} \sigma_{k/2} & \text{if } k \equiv 0 \pmod{2} \\ \tau_{(k-1)/2} & \text{if } k \equiv 1 \pmod{2} \end{cases}$$
 is causal.

Examples

1. $\bar{b} : 1^\omega \rightarrow B^\omega$ defined by $[\bar{b}]_k = b$ is causal.
2. If $g : A \rightarrow B$, then $\text{map}(g) : A^\omega \rightarrow B^\omega$ defined by $[\text{map}(g)(\sigma)]_k = g(\sigma_k)$ is causal.
3. $\text{runProd} : \mathbb{R}^\omega \rightarrow \mathbb{R}^\omega$ defined by $[\text{runProd}(\sigma)]_k = \prod_{i=0}^k \sigma_i$ is causal.
4. $\text{zip} : (A \times A)^\omega \rightarrow A^\omega$ defined by
$$[\text{zip}(\sigma, \tau)]_k = \begin{cases} \sigma_{k/2} & \text{if } k \equiv 0 \pmod{2} \\ \tau_{(k-1)/2} & \text{if } k \equiv 1 \pmod{2} \end{cases}$$
 is causal.
5. $\text{t1} : A^\omega \rightarrow A^\omega$ defined by $[\text{t1}(\sigma)]_k = \sigma_{k+1}$ is **not causal**, since output k is input $k + 1$.

Background

Causal functions appear all over computer science:

- 1 Mealy machines [Mealy, 1955 and Eilenberg, 1974]
- 2 Coalgebra [Rutten, 2006]
- 3 Synchronous digital circuits [Leiserson and Saxe, 1986]
- 4 Signal flow graphs [Bonchi, Sobociński and Zanasi, 2014 and Milius, 2010]
- 5 Recurrent neural networks [Rumelhart, 1986 and Hochreiter and Schmidhuber, 1997]

(If you have examples from mathematics where this class of functions is important, let me know!)

Automata theory / coalgebraic approach

Mealy machines with input alphabet A and output alphabet B have behaviors much like causal functions $A^\omega \rightarrow B^\omega$.

Automata theory / coalgebraic approach

Mealy machines with input alphabet A and output alphabet B have behaviors much like causal functions $A^\omega \rightarrow B^\omega$.

Indeed, causal functions $A^\omega \rightarrow B^\omega$ carry a final coalgebra for the functor $M_{A,B}(X) = (B \times X)^A$. [Rutten, 2006]

Automata theory / coalgebraic approach

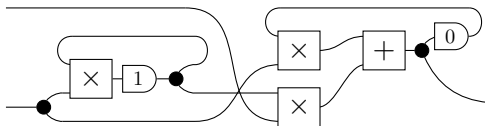
Mealy machines with input alphabet A and output alphabet B have behaviors much like causal functions $A^\omega \rightarrow B^\omega$.

Indeed, causal functions $A^\omega \rightarrow B^\omega$ carry a final coalgebra for the functor $M_{A,B}(X) = (B \times X)^A$. [Rutten, 2006]

Thus one way to specify a causal function is to specify a state in a Mealy machine. This is the basis of the *syntactic (GSOS) method* in stream differential equations. [Hansen, Kupke and Rutten, 2017]

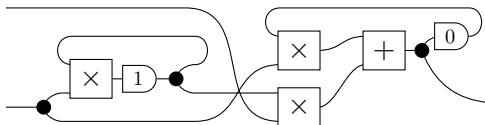
Digital circuit / signal flow graph approach

Circuit diagram treatments of causal functions are more explicit about how state is stored, how it is updated, and how outputs are produced. For example:

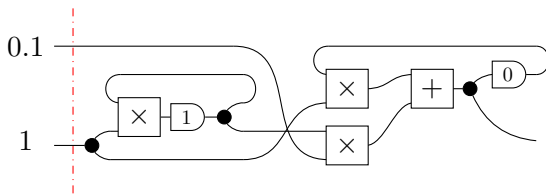


Digital circuit / signal flow graph approach

Circuit diagram treatments of causal functions are more explicit about how state is stored, how it is updated, and how outputs are produced. For example:

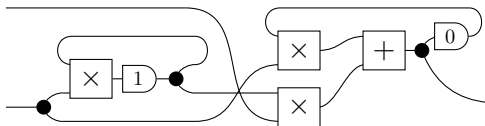


Let's evaluate this diagram once:

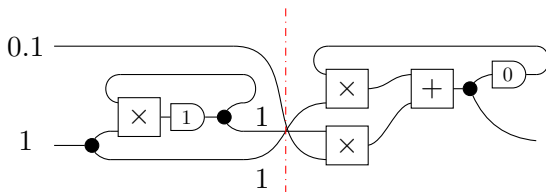


Digital circuit / signal flow graph approach

Circuit diagram treatments of causal functions are more explicit about how state is stored, how it is updated, and how outputs are produced. For example:

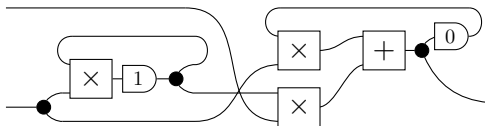


Let's evaluate this diagram once:

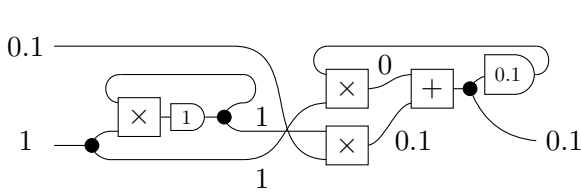


Digital circuit / signal flow graph approach

Circuit diagram treatments of causal functions are more explicit about how state is stored, how it is updated, and how outputs are produced. For example:



Let's evaluate this diagram once:



Recurrent neural networks

RNNs are special instances of the circuit diagram approach with the requirement that the functions involved are *differentiable*.

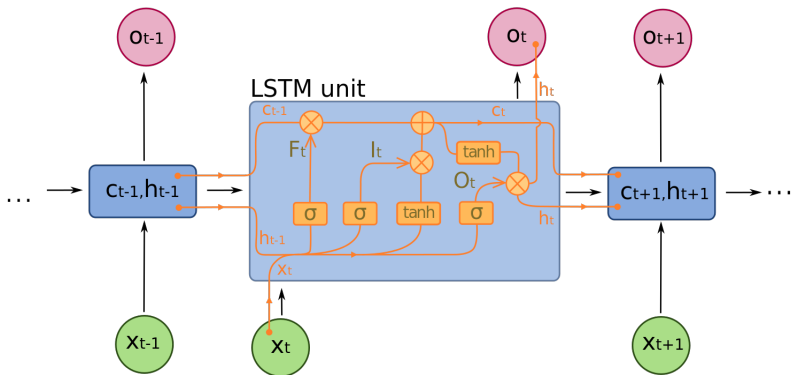


Figure: François Deloche - Own work, CC BY-SA 4.0, [wikimedia](#)

Finite approximations approach

Another perspective on causal functions is to treat them as a limit of functions on finite lists. This is based on the following observation:

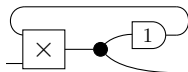
Theorem

The following are equivalent:

- 1 a causal function $f : A^\omega \rightarrow B^\omega$,
- 2 an infinite sequence of functions $u_k : A^{k+1} \rightarrow B$, and
- 3 an infinite sequence of functions $t_k : A^{k+1} \rightarrow B^{k+1}$ satisfying $[t_{k+1}(\vec{a})]_{0:k} = t_k([\vec{a}]_{0:k})$ for all $\vec{a} \in A^{k+2}$.

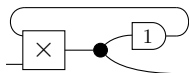
Proof Idea. u_k gives the $k + 1^{\text{st}}$ entry in the sequence produced by f , and t_k provides the first $k + 1$ entries.

What causal function is represented by this circuit?

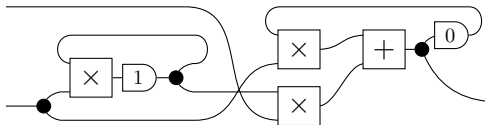


Where we're heading, update I

I would like to convince you the derivative of the running product:



Is this:



Outline

- 1 Causal functions
- 2 Computation sequences
 - Main ideas
 - Sanity check
- 3 Delayed trace
- 4 Neural networks
- 5 Cartesian differential categories
- 6 Recap and future directions

Main ideas

We start from a (strictified Cartesian) category \mathbb{C} whose morphisms represent simple computations. Our goal is to define an extended category whose morphisms are causal functions.

Main ideas

We start from a (strictified Cartesian) category \mathbb{C} whose morphisms represent simple computations. Our goal is to define an extended category whose morphisms are causal functions.

There are three steps:

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

Step 1: Single computation steps

We separate inputs and outputs of a \mathbb{C} -morphism into two classes: *values*, exchanged with the environment, and *states*, exchanged with next/previous steps of the computation.

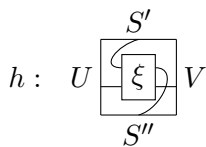
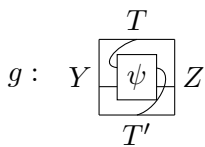
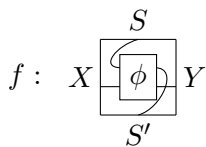
Definition

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

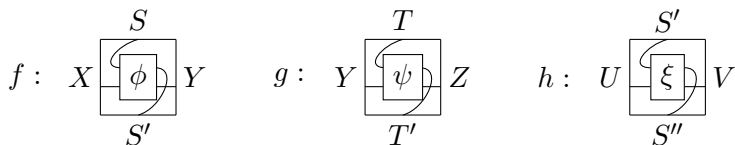
Step 1: Single computation steps

We can draw a single step like so:

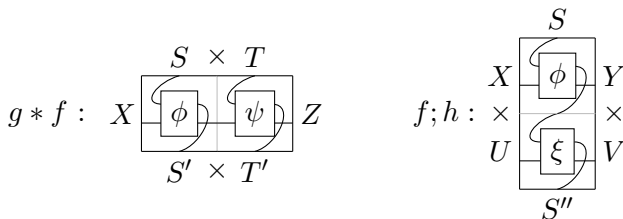


Step 1: Single computation steps

We can draw a single step like so:

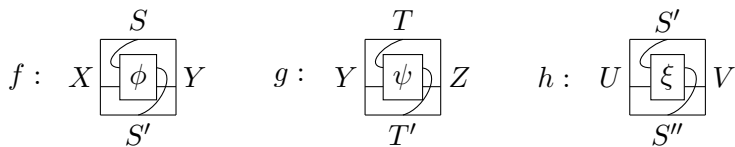


Computation steps can be composed in two different ways:

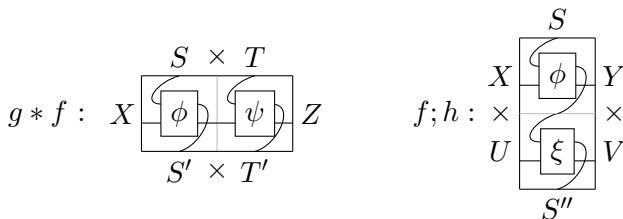


Step 1: Single computation steps

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 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

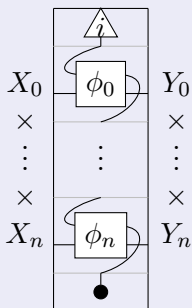
$i : 1 \xrightarrow[S_0]{1} 1$. We say \mathbf{f} takes sequences of type $[X_k] = (X_0, X_1, \dots)$ to sequences of type $[Y_k] = (Y_0, Y_1, \dots)$.

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

Step 3: Comparing computation sequences

Definition

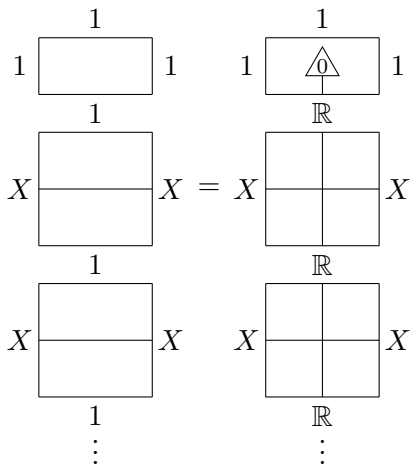
The n th truncation of a computation sequence is the morphism of the vertical composite of the first $n + 1$ steps:



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 $n \in \mathbb{N}$.

Step 3: Extensionally equivalent identity computations



Step 3: Comparing computation sequences

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 check

Theorem

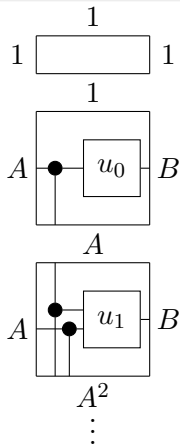
The homset $\text{St}(\text{Set})([A], [B])$ is in 1-1 correspondence with the set of causal functions from A^ω to B^ω .

Sanity check

Theorem

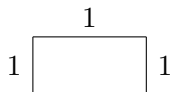
The homset $\text{St}(\text{Set})([A], [B])$ is in 1-1 correspondence with the set of causal functions from A^ω to B^ω .

Idea. A causal $f : A^\omega \rightarrow B^\omega$ determines a unique sequence of functions $u_k : A^{k+1} \rightarrow B$. Form a computation sequence as follows:




Quiz

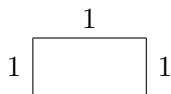
How would you implement \boxed{i} as a computation sequence?



⋮

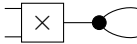
Quiz

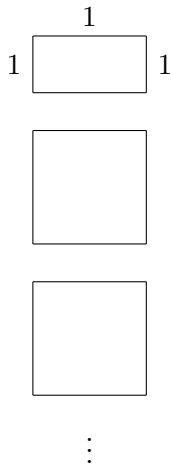
How would you implement
 as a computation
sequence?

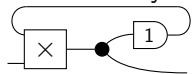


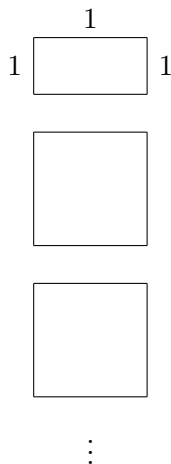
⋮

Quiz

How would you implement  as a computation sequence?



How would you implement  as a computation sequence?



Outline

- 1 Causal functions
- 2 Computation sequences
 - Main ideas
 - Sanity check
- 3 Delayed trace**
- 4 Neural networks
- 5 Cartesian differential categories
- 6 Recap and future directions

Delayed trace

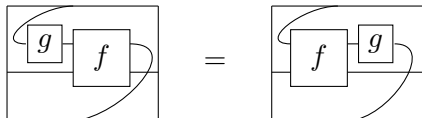
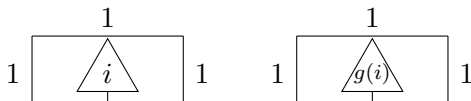
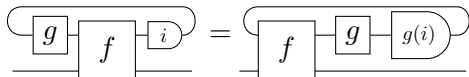
This loop-with-delay-gate is a trace-like operation.

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

It satisfies **most** of the the trace axioms [Joyal, Street, Verity, 1996] but misses two: yanking and dinaturality. For monoidal trace, those are

The diagram shows two equations. The first equation shows a yanking move: a box with two wires entering from the bottom and two wires exiting from the top, with a loop on the top wire, is equal to a simple crossing of the two wires. The second equation shows a dinaturality move: a box with two wires entering from the bottom and two wires exiting from the top, with a loop on the top wire, is equal to a box with two wires entering from the bottom and two wires exiting from the top, with a loop on the top wire, but the box is split into two boxes, f and g , with the loop on the top wire of g .

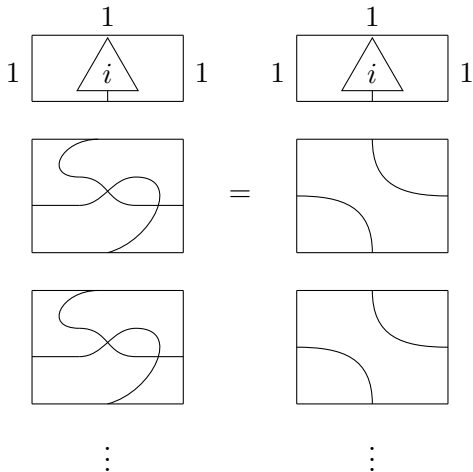
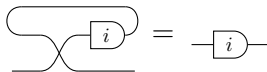
Dinaturality \rightarrow retiming



\vdots

\vdots

Yanking \rightarrow delay



Advantages of delayed trace

Often, people formalizing circuits in category theory use these steps:

- 1 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 use these steps:

- 1 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.

Our idea is 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 use these steps:

- 1 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.

Our idea is to:

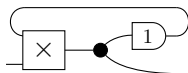
- 1 add a delayed trace to the category,
- 2 recover registers by delay-tracing symmetry

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

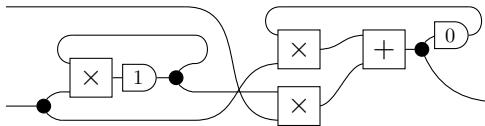
We are also interested in looking at delayed traces abstractly. . .

Where we're heading, update II

I would like to convince you the derivative of the running product:



Is this:



Now we know how to represent these circuits as computation sequences.

Outline

- 1 Causal functions
- 2 Computation sequences
 - Main ideas
 - Sanity check
- 3 Delayed trace
- 4 Neural networks**
- 5 Cartesian differential categories
- 6 Recap and future directions

(Feedforward) neural networks

A *neural network* looks like $\begin{matrix} \theta \\ x \end{matrix} \rightarrow \boxed{N} \rightarrow y$. We assume its input and output spaces are real Euclidean spaces \mathbb{R}^n .

Training a neural network (with data \hat{x}_i, \hat{y}_i) means finding θ^* so that:

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

(Feedforward) neural networks

A *neural network* looks like $\begin{matrix} \theta \\ x \end{matrix} \rightarrow \boxed{N} \rightarrow y$. We assume its input and output spaces are real Euclidean spaces \mathbb{R}^n .

Training a neural network (with data \hat{x}_i, \hat{y}_i) means finding θ^* 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 θ . This allows us to make smart updates to θ^* .

(Feedforward) neural networks

A *neural network* looks like $\begin{matrix} \theta \\ x \end{matrix} \rightarrow \boxed{N} \rightarrow y$. We assume its input and output spaces are real Euclidean spaces \mathbb{R}^n .

Training a neural network (with data \hat{x}_i, \hat{y}_i) means finding θ^* so that:

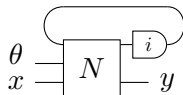
$$\begin{matrix} \triangle \theta^* \\ \hat{x}_i \end{matrix} \rightarrow \boxed{N} \rightarrow y_i \approx \hat{y}_i$$

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 θ . This allows us to make smart updates to θ^* .

Various algorithms can be used to compute these partial derivatives, including *backpropagation* [Rumelhart, 1986].

(Feedback/recurrent) neural networks

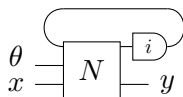
Recurrent neural networks (RNNs) are designed to process sequences of inputs using some internal state:



It's less clear what the derivatives of these functions are.

(Feedback/recurrent) neural networks

Recurrent neural networks (RNNs) are designed to process sequences of inputs using some internal state:

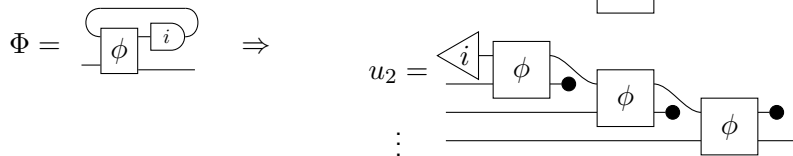


It's less clear what the derivatives of these functions are.

The most common training method for RNNs relies on a gradient-finding algorithm called *backpropagation through time*.

Backpropagation through time

First, we *unroll* the network:



Then, whenever the derivative of Φ is needed at an input of length $k + 1$, the derivative of u_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

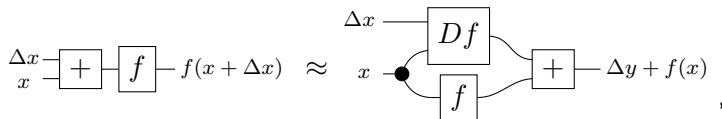
- 1 Causal functions
- 2 Computation sequences
 - Main ideas
 - Sanity check
- 3 Delayed trace
- 4 Neural networks
- 5 Cartesian differential categories**
- 6 Recap and future directions

Derivatives, diagrammatically

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 .

Derivatives, diagrammatically

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



Derivatives, diagrammatically

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

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

(CD2) $\triangleleft 0 \rightarrow Df \rightarrow = \bullet \rightarrow \triangleleft 0$, and

Derivatives, diagrammatically

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

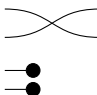
$$\frac{\Delta x}{x} \rightarrow + \rightarrow f \rightarrow f(x + \Delta x) \approx \begin{array}{c} \Delta x \\ x \end{array} \rightarrow \begin{array}{c} Df \\ f \end{array} \rightarrow + \rightarrow \Delta y + f(x)$$

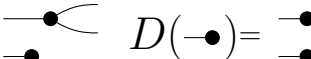
(CD2) $\triangleleft 0 \rightarrow Df \rightarrow = \bullet \triangleleft 0$, and

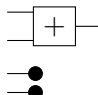
(CD3) $\rightarrow + \rightarrow Df \rightarrow = \begin{array}{c} Df \\ Df \end{array} \rightarrow +$.

Derivatives of basic functions

The derivatives of some basic functions include:

$$(CD1) \quad D(\text{---}) = \text{---} \bullet \quad D(\text{---} \times \text{---}) = \text{---} \bullet \bullet$$


$$D(\text{---} \bullet \text{---}) = \text{---} \bullet \text{---} \quad D(\text{---} \bullet) = \text{---} \bullet \bullet$$


$$D(\triangleleft 0 \text{---}) = \triangleleft 0 \text{---} \quad D(\text{---} \square + \text{---}) = \text{---} \bullet \bullet$$


Properties of derivatives, diagrammatically

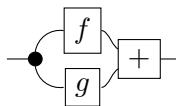
We can also express the chain rule and “parallel rule”:

$$(CD4) \quad D(-\boxed{f}-\boxed{g}-) = \begin{array}{c} \text{---} \\ | \\ \boxed{Df} \\ | \\ \bullet \\ / \quad \backslash \\ \boxed{f} \quad \boxed{Dg} \\ \backslash \quad / \\ \text{---} \end{array}$$

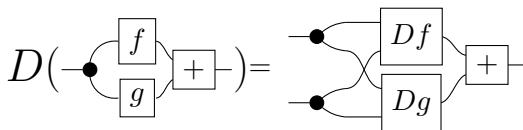
$$(CD5) \quad D\left(\begin{array}{c} \boxed{f} \\ \boxed{g} \end{array}\right) = \begin{array}{c} \text{---} \\ | \\ \boxed{Df} \\ | \\ \text{---} \\ \backslash \quad / \\ \boxed{Dg} \\ | \\ \text{---} \end{array}$$

Sum rule, diagrammatically

The sum rule is a consequence of the axioms we have seen so far. Here's a diagram for the sum of two functions:



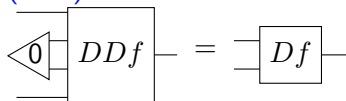
The rules we've seen allows us to compute:



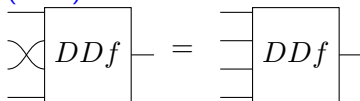
Properties of higher-order derivatives

It's also possible to express properties of higher order derivatives, but it's not worth the time to explain them in much detail.

(CD6)



(CD7)



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

A *Cartesian differential category* is a

- Cartesian category
- with left additive structure (providing 0 and $+$)
- an operation sending $f : X \rightarrow Y$ to $Df : X \times X \rightarrow Y$
- which satisfies CD1-7.

This is a categorical representation of the major algebraic properties of derivatives.

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

A *Cartesian differential category* is a

- Cartesian category
- with left additive structure (providing 0 and $+$)
- an operation sending $f : X \rightarrow Y$ to $Df : X \times X \rightarrow Y$
- which satisfies CD1-7.

This is a categorical representation of the major algebraic properties of derivatives.

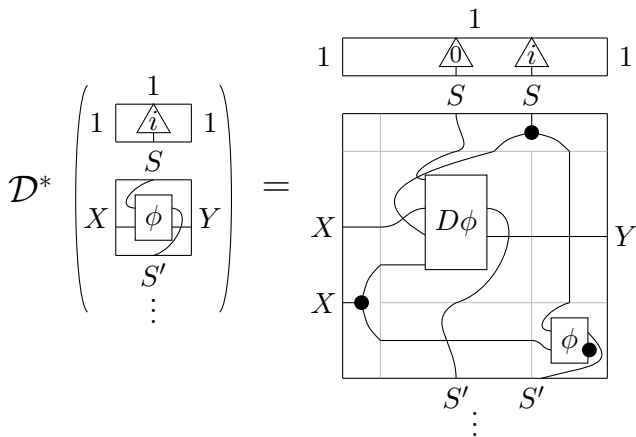
The canonical example of a Cartesian differential category is \mathbf{Euc}_∞ , where objects are \mathbb{R}^n for all $n \in \mathbb{N}$ and morphisms are smooth functions.

Key contribution: differential operator lifts

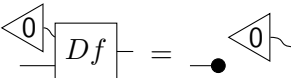
Theorem (Katsumata, S)

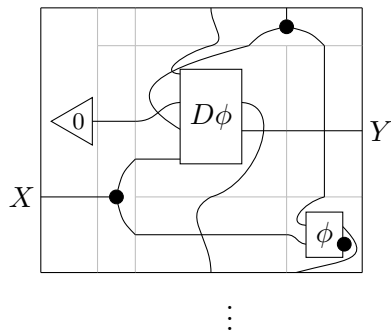
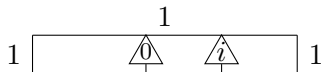
If \mathbb{C} is a Cartesian differential category, so is $\text{St}(\mathbb{C})$.

The Cartesian differential operator on $\text{St}(\mathbb{C})$ is defined as follows:

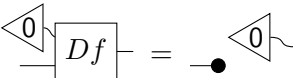


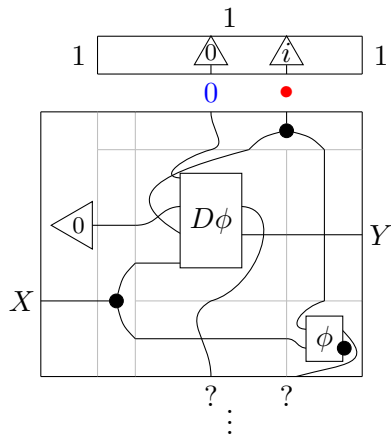
Key contribution: differential operator lifts

Proof idea. For CD2: 

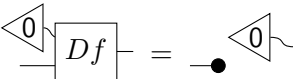


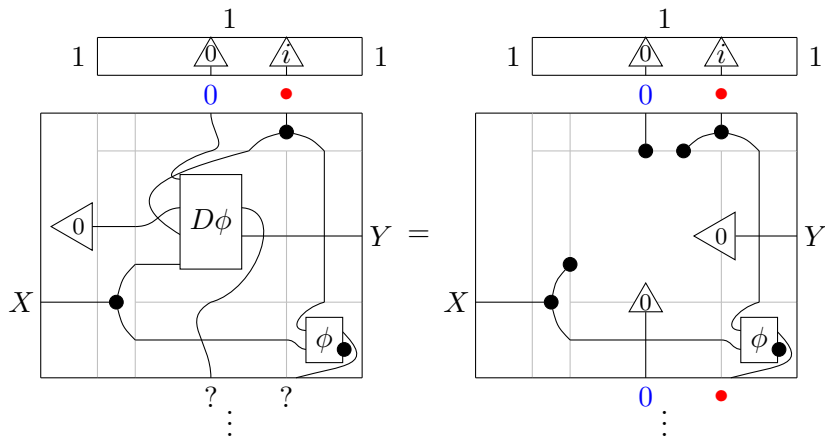
Key contribution: differential operator lifts

Proof idea. For CD2: 

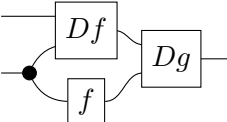


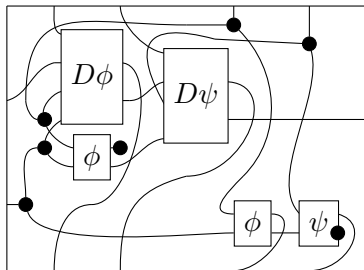
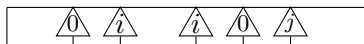
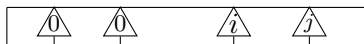
Key contribution: differential operator lifts

Proof idea. For CD2: 

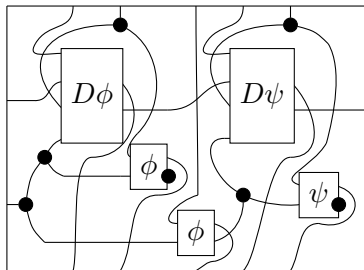


Key contribution: differential operator lifts

Proof idea. For CD4: $D(-\boxed{f}-\boxed{g}-) =$ 



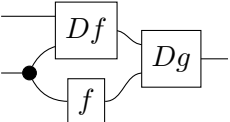
=

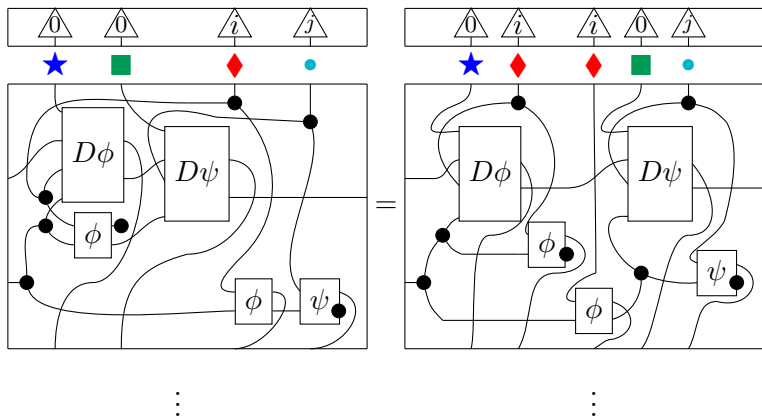


⋮

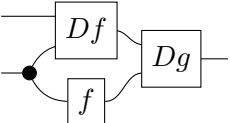
⋮

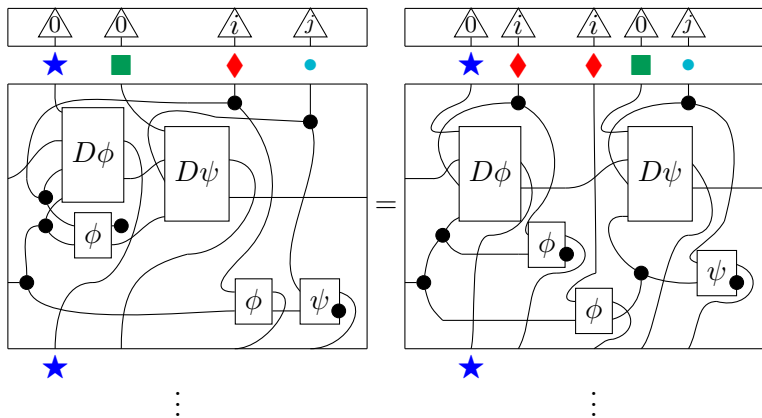
Key contribution: differential operator lifts

Proof idea. For CD4: $D(-\boxed{f}-\boxed{g}-) =$ 

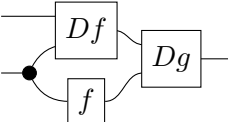


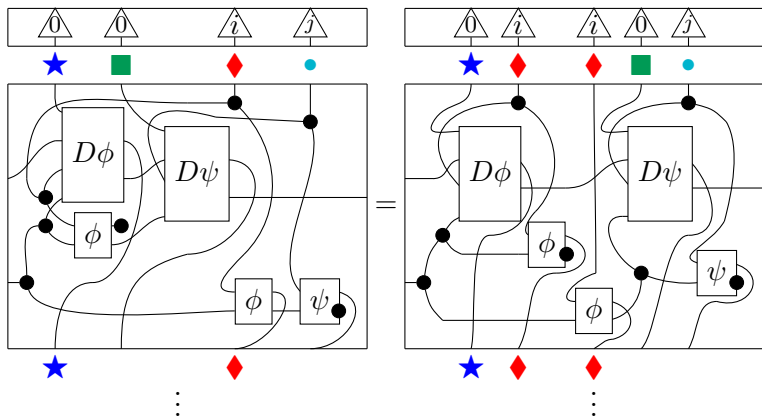
Key contribution: differential operator lifts

Proof idea. For CD4: $D(-\boxed{f}-\boxed{g}-) =$ 

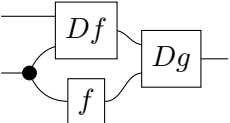


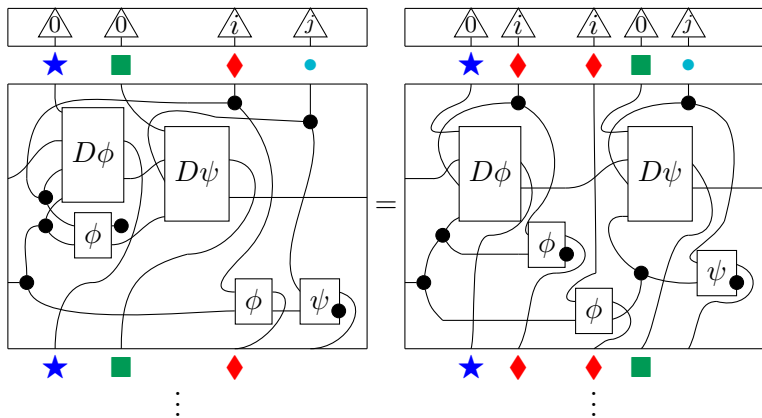
Key contribution: differential operator lifts

Proof idea. For CD4: $D(-\boxed{f}-\boxed{g}-) =$ 

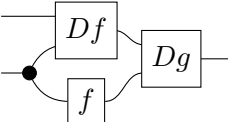


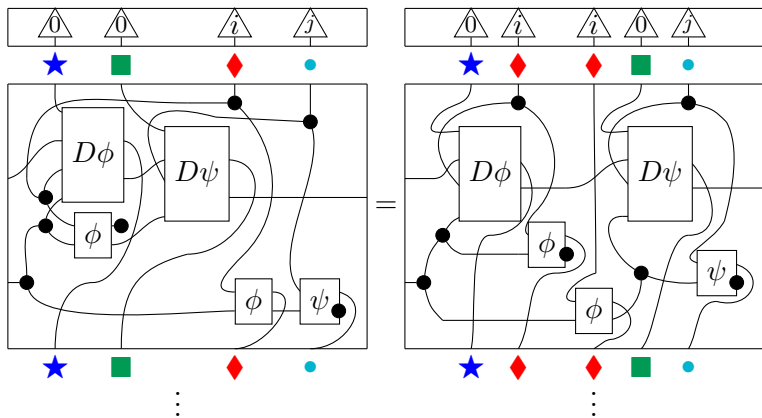
Key contribution: differential operator lifts

Proof idea. For CD4: $D(-\boxed{f}-\boxed{g}-) =$ 

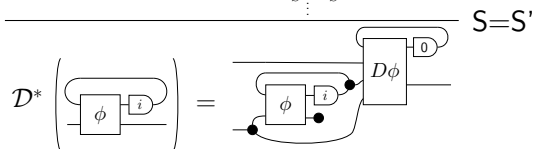
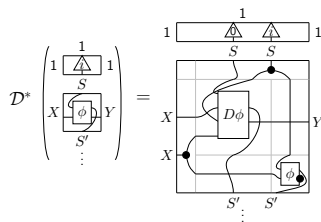


Key contribution: differential operator lifts

Proof idea. For CD4: $D(-\boxed{f}-\boxed{g}-) =$ 



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 $(i, [\phi])$.*

Where we're heading, update III

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

Outline

- 1 Causal functions
- 2 Computation sequences
 - Main ideas
 - Sanity check
- 3 Delayed trace
- 4 Neural networks
- 5 Cartesian differential categories
- 6 Recap and future directions

Recap

We gave a mechanism for defining causal functions based on computation sequences. We organized these with a construction on Cartesian categories, $\text{St}(-)$. These categories have a *delayed trace*.

We showed that if \mathbb{C} is a Cartesian differential category, so is $\text{St}(\mathbb{C})$. Taking $\mathbb{C} = \mathbf{Euc}_\infty$, we can find derivatives of recurrent neural networks. The differential operator on $\text{St}(\mathbf{Euc}_\infty)$ matches the derivatives taken in backpropagation through time; the fact that it is a differential operator obtains many properties for BPTT.

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!