

# CATEGORICAL FOUNDATIONS OF GRADIENT-BASED LEARNING

(CRUTTWELL, GAVRANOVIĆ, GHANI, WILSON, ZANASI)

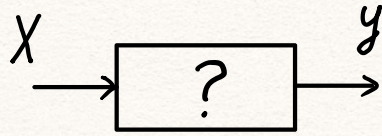
GOAL:

PROVIDE A CATEGORICAL FRAMEWORK

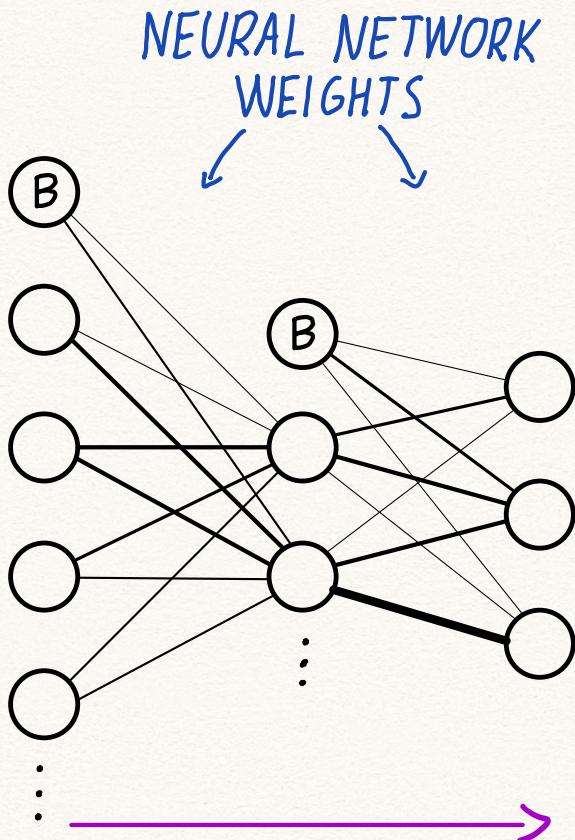
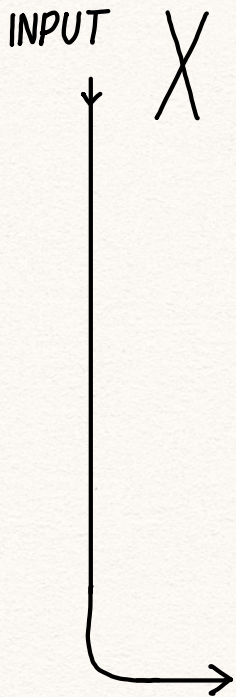
—

FOR DEEP LEARNING

# SUPERVISED LEARNING WITH NEURAL NETWORKS IN ONE SLIDE:

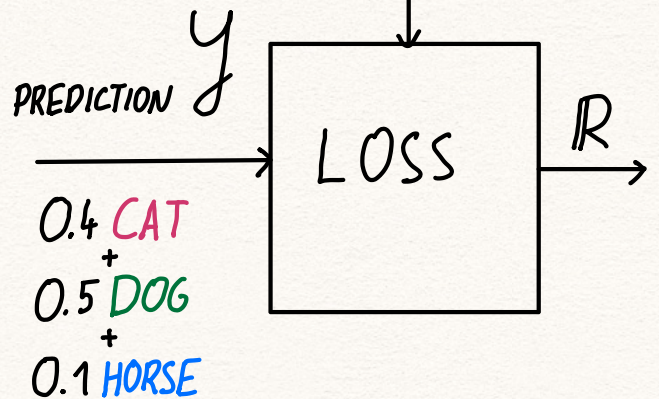


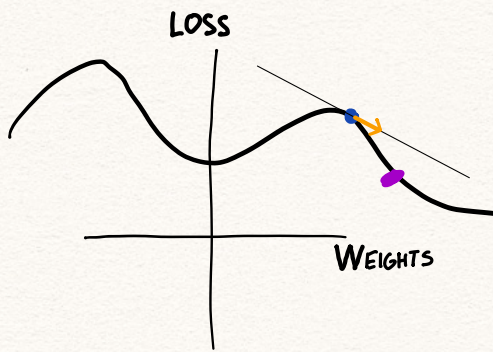
DATASET: List  $X, y$



1 CAT  
+  
0 DOG  
+  
0 HORSE

LABEL  $y$





# GRADIENT DESCENT ~ "OPTIMIZER"

- NN IS COMPUTATION PARAMETERIZED BY WEIGHTS
- BACKPROPAGATION OF CHANGES
- PARAMETER UPDATE - "OPTIMIZERS"

NEURAL NETWORKS



- LINEAR LAYER
- BIAS TERM
- ACTIVATION FUNCTION

THIS SIMPLE STORY PERMEATES DEEP LEARNING!

# PLAN FOR TODAY?

## TAKE A BIRD'S EYE VIEW OF NEURAL NETWORKS



- TRACE OUT THE INFORMATION FLOW ABOVE
  - PRECISELY WRITE DOWN ALL THE HIGH-LEVEL NOTIONS IN ISOLATION:
    - DIFFERENTIATION - REVERSE DERIVATIVE CATS.
    - BIDIRECTIONALITY - OPTICS/LENSES
    - PARAMETERIZATION - PARA
- AND STUDY THEIR INTERACTION.

## PARAMETERIZED OPTICS

AS A COMMON STRUCTURE BEHIND

- NEURAL NETWORKS
- LOSS FUNCTIONS
- OPTIMIZERS

• PAUL: CONCRETE EXAMPLES OF NEURAL NETWORKS

# DIFFERENTIATION

- CARTESIAN (FORWARD) DIFFERENTIAL CATEGORIES  
(Blute et al.)
- CARTESIAN REVERSE DIFFERENTIAL CATEGORIES (CRDC)  
(Cockett et al.)

## DEFINITION.

A CRDC  $\mathcal{C}$  is a Cartesian left-additive category which for every map

$$f: A \longrightarrow B$$

has a REVERSE DIFFERENTIAL COMBINATOR

$$R[f]: A \times B \longrightarrow A$$

(compare  
 $D[f]: A \times A \longrightarrow B$ )

subject to 7 axioms.

EXAMPLE. Smooth is a CRDC.  $\text{Poly}_{\mathbb{Z}_2}$  IS A CRDC.

EXAMPLE. Let  $\mathbb{R}^2 \xrightarrow{f} \mathbb{R}$  in Smooth.  
 $(x, y) \longmapsto x^2 + 3yx$

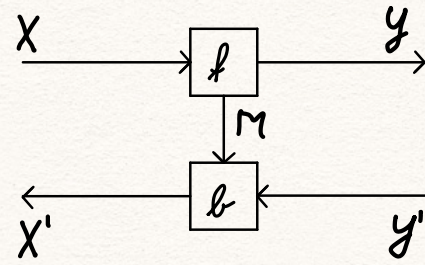
Then  $R[f]: \mathbb{R}^2 \times \mathbb{R} \longrightarrow \mathbb{R}^2$   
 $(x, y), w \longmapsto (2xw, 3xw)$

PLAN: STUDY CRDC'S THROUGH OPTICS/LENSES

# OPTICS / LENSES

DEFINITION. Let  $\mathcal{C}$  be a SMC. Category  $\text{Optic}(\mathcal{C})$ :

• Objects - pairs of objects  $\begin{pmatrix} X \\ X' \end{pmatrix}$  in  $\mathcal{C}$



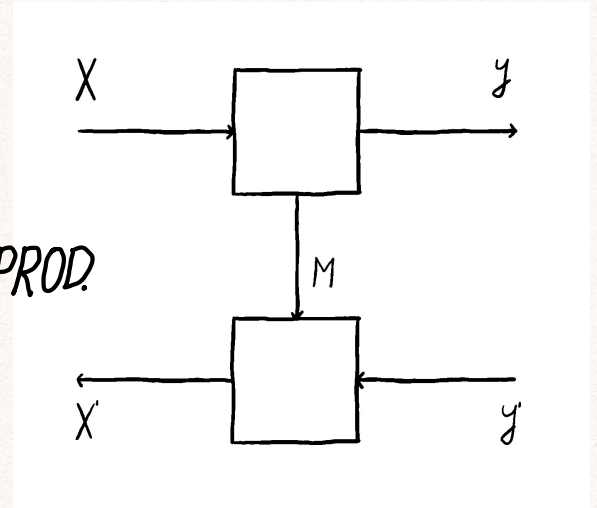
$$\text{Optic}(\mathcal{C}) \left( \begin{pmatrix} X \\ X' \end{pmatrix}, \begin{pmatrix} Y \\ Y' \end{pmatrix} \right) = \int^{m:e \leftarrow \text{COEND}} \mathcal{C}(X, Y \otimes M) \times \mathcal{C}(Y' \otimes M, X')$$

$$(M, f, b) \quad \begin{aligned} f: X &\longrightarrow Y \otimes M \\ b: Y' \otimes M &\longrightarrow X' \end{aligned}$$

PROP. If  $\mathcal{C}$  is Cartesian,

$$\int^{m:e} \mathcal{C}(X, M \times Y) \times \mathcal{C}(M \times Y', X') \cong \text{UNIV. PROPERTY OF PROD.}$$

$$\int^{m:e} \mathcal{C}(X, Y) \times \mathcal{C}(X, M) \times \mathcal{C}(M \times Y', X') \cong \text{YONEDA REDUCTION}$$

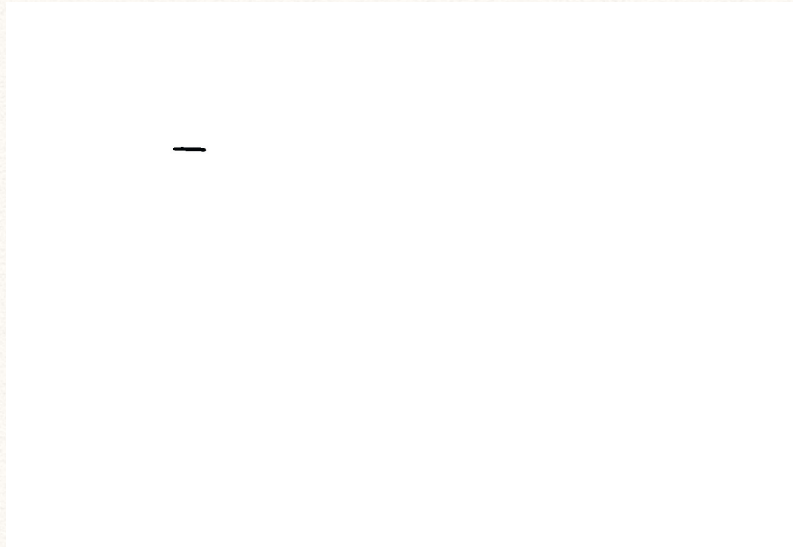


$$\int \mathcal{C}(X, Y) \times \mathcal{C}(X \times Y', X')$$

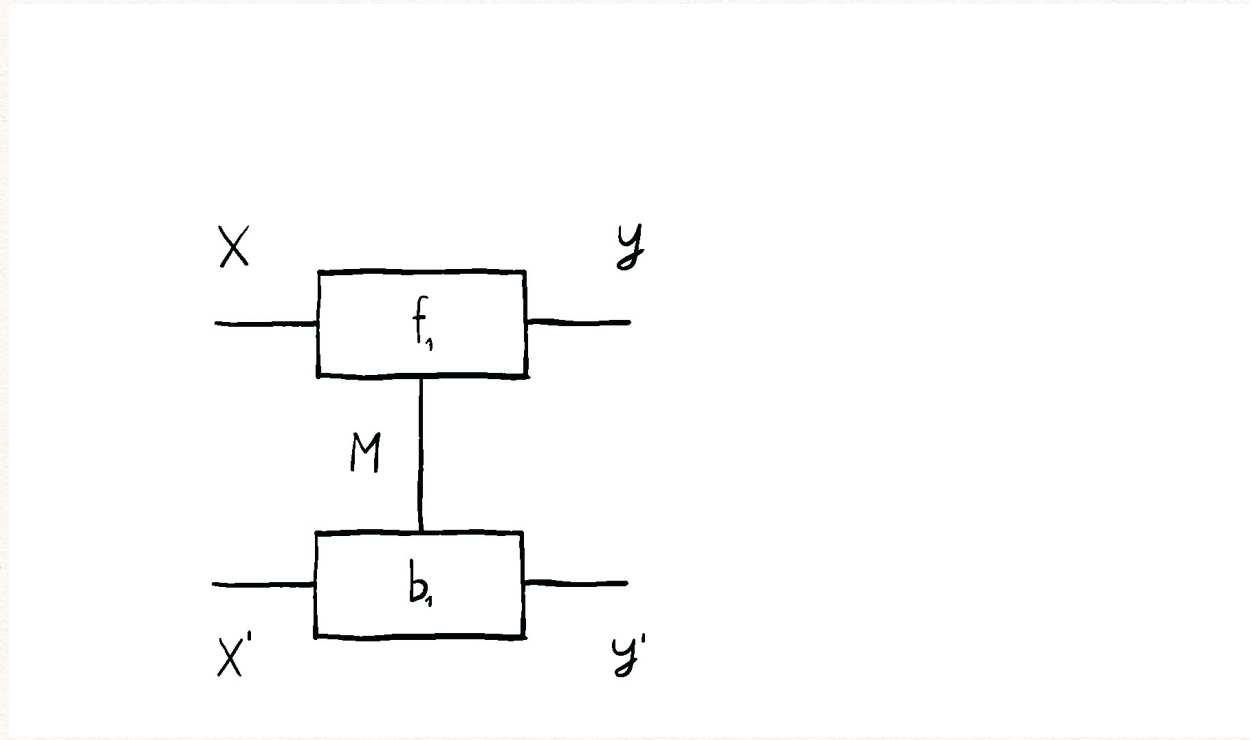
*get*                      *put*

then  $\text{Optic}(\mathcal{C}) \cong \text{Lens}(\mathcal{C})$

BIDIRECTIONAL INFORMATION FLOW



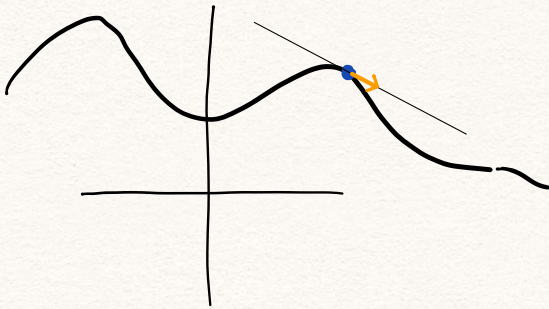
OPTICS CAN BE COMPOSED



PROPOSITION.  $\text{Optic}(\mathcal{C})$  is symmetric monoidal.

# EXAMPLE.

## GRADIENT DESCENT



$$P \times P' \xrightarrow{u} P$$

$$(p, \nabla p) \mapsto p - d \nabla p$$

is a lens, for  $\mathcal{E} := \text{Smooth}$

$$\begin{pmatrix} P \\ p \end{pmatrix} \xrightarrow{(id_p, u)} \begin{pmatrix} P \\ p' \end{pmatrix}$$



# EXAMPLE. STATEFUL OPTIMIZERS

- MOMENTUM,

INTERNAL STATE

$$\text{get: } P \times P \longrightarrow P$$

$$(w, p) \longmapsto p$$

$$\text{put: } P \times P \times P \longrightarrow P \times P$$

$$(w, p, \nabla p) \longmapsto (w', p - w')$$

where  $w' = \gamma w + \epsilon p'$

$$\begin{pmatrix} S \times P \\ S \times P \end{pmatrix} \longrightarrow \begin{pmatrix} P \\ P' \end{pmatrix}$$

- NESTEROV MOMENTUM

$$\text{get: } P \times P \longrightarrow P$$

$$(w, p) \longmapsto p - \gamma w$$

put - same as above

- ADAGRAD

- ADAM

...



BACK TO CRDC's:

$$\begin{array}{l} f: A \longrightarrow B \\ \text{R}[f]: A \times B \longrightarrow A \end{array} \quad \begin{array}{l} \sim \\ \sim \end{array} \quad \begin{array}{l} \text{'get' MAP OF A LENS} \\ \text{'put' MAP OF A LENS} \end{array}$$

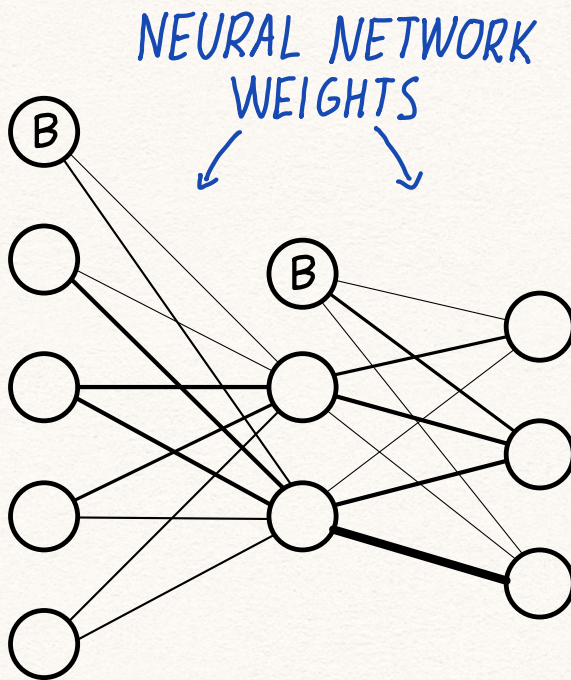
## PROPOSITION.

For each CRDC  $\mathcal{C}$  there is a symmetric monoidal functor

$$\begin{array}{ccc} \mathcal{C} & \xrightarrow{F} & \text{Lens}(\mathcal{C}) \cong \text{Optic}(\mathcal{C}) \\ A & \longmapsto & (A, A) \\ f \downarrow & & \downarrow (f, \text{R}[f]) \\ B & \longmapsto & (B, B) \end{array}$$

• THIS IS OUR FRAMEWORK FOR BACKPROPAGATION

# PARAMETERIZATION



Fix a SMC  $(\mathcal{C}, \otimes, I)$ .

DEF. Bicategory  $\text{Para}(\mathcal{C})$

Objects - objects of  $\mathcal{C}$

$$\text{Para}(\mathcal{C})(A, B) = \int_{P: \mathcal{C}}^{\text{op}} \mathcal{C}(P \otimes A, B)$$

CATEGORY OF ELEMENTS

$$A \xrightarrow{(P: \mathcal{C}, f: P \otimes A \rightarrow B)} B$$

$$A \begin{array}{c} \xrightarrow{(P, f)} \\ \Downarrow \tau \\ \xrightarrow{(Q, g)} \end{array} B$$

2-cells are reparameterizations: a 2-cell

is a map  $Q \xrightarrow{\tau} P$  such that

$$\begin{array}{ccc} Q \otimes A & \xrightarrow{\tau \otimes A} & P \otimes A \\ g \searrow & & \swarrow f \\ & B & \end{array} \text{ commutes.}$$

# EXAMPLE.

$(\text{Set}, x, 1)$

$\text{Para}(\text{Set})$

SETS AND  
PARAMETERIZED FUNCTIONS

$(\text{Smooth}, x, 1)$

$\text{Para}(\text{Smooth})$

EUCLIDEAN SPACES AND  
PARAMETERIZED SMOOTH  
FUNCTIONS

$(\text{Optic}(\mathcal{C}), \otimes, 1)$

$\text{Para}(\text{Optic}(\mathcal{C}))$

PAIRS OF OBJECTS AND  
PARAMETERIZED OPTICS

• • •

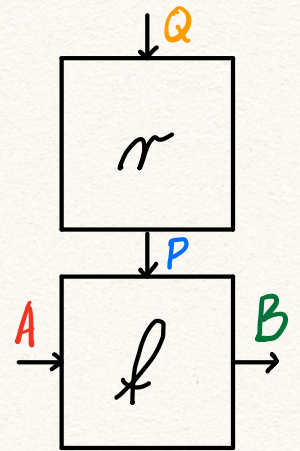
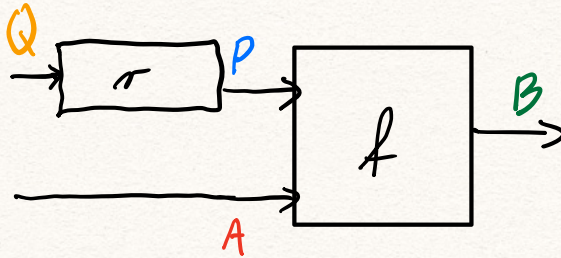
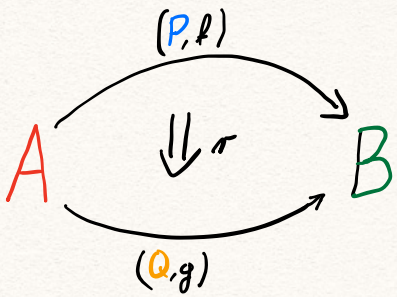
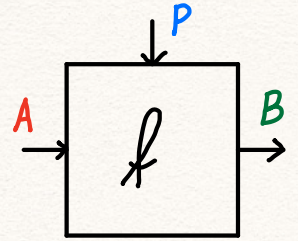
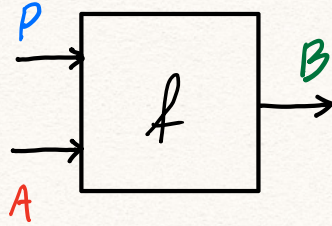
# GRAPHICAL LANGUAGE

TEXTUAL  
NOTATION

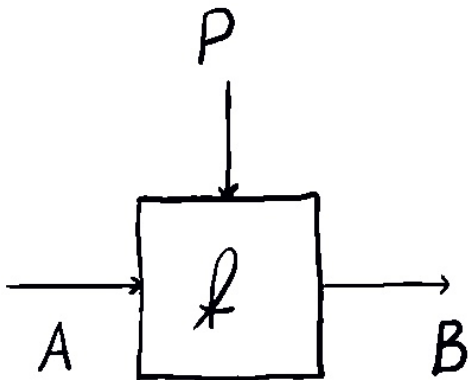
STANDARD  
STRING DIAGRAM

2D  
STRING DIAGRAM

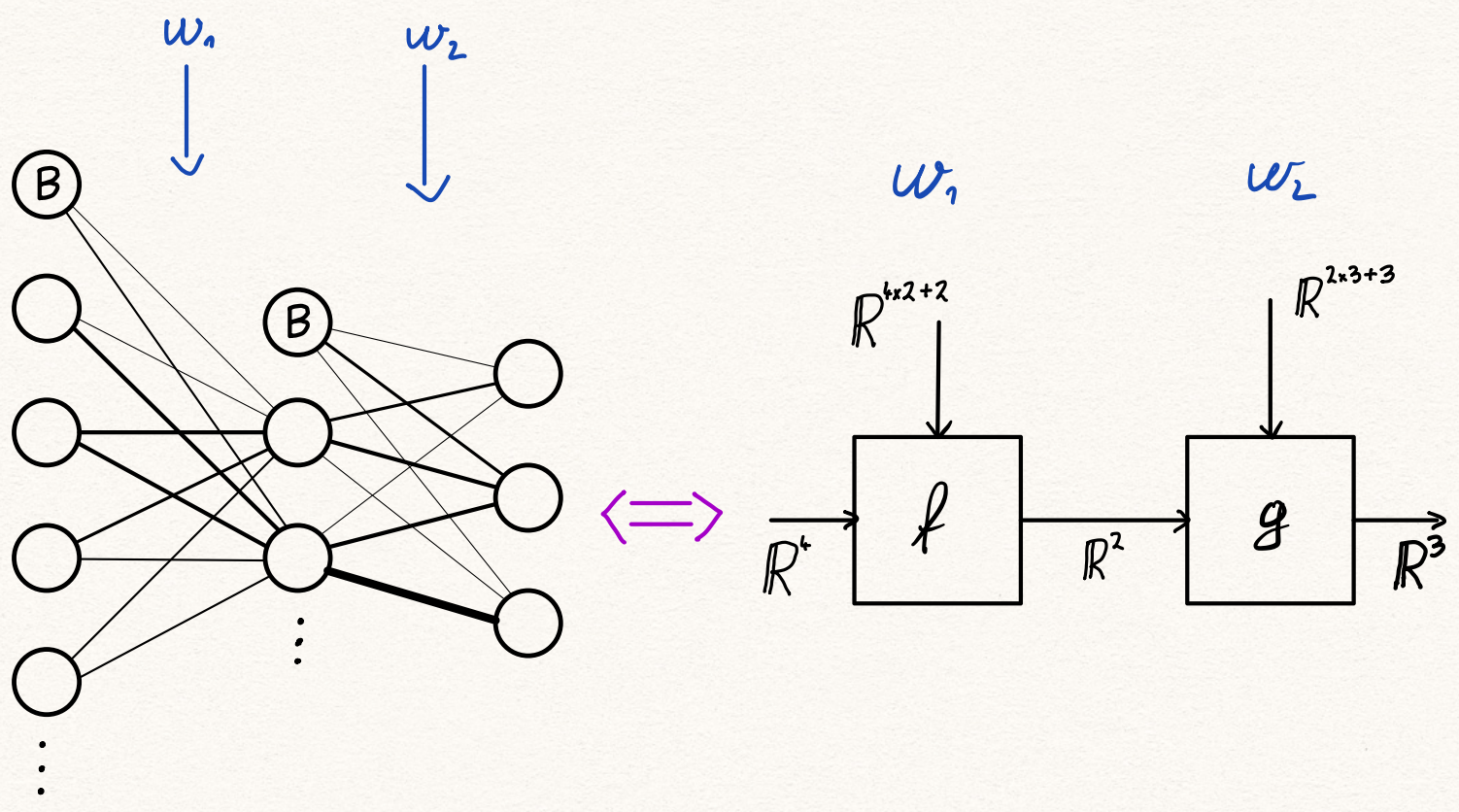
$$f: P \otimes A \longrightarrow B$$



HOW DOES COMPOSITION WORK?



# RECAP



Para IS NATURAL WITH RESPECT TO BASE CHANGE.

# DEFINITION.

Let  $G: \mathcal{C} \rightarrow \mathcal{D}$  be a symm. monoidal functor. We define

$$\begin{array}{ccc} \text{Para}(G): \text{Para}(\mathcal{C}) & \longrightarrow & \text{Para}(\mathcal{D}) \\ & & \\ & A \longmapsto & GA \\ & \downarrow (P, f) & \downarrow (GP, f') \\ & B \longmapsto & GB \end{array}$$

where  $f'$  is the composite

$$G(P) \otimes G(A) \xrightarrow{\mu_{P,A}} G(P \otimes A) \xrightarrow{G(f)} G(B)$$

+ MORE.

Para IS RICH IN CATEGORICAL STRUCTURE.

- Cokleisli category of a graded comonad
- Double category
- Actegorical Para

...

# PARAMETERIZED OPTICS

$$\mathcal{C} \longrightarrow \text{Optic}(\mathcal{C}) \longrightarrow \text{Para}(\text{Optic}(\mathcal{C}))$$

• Objects - Objects of  $\text{Optic}(\mathcal{C})$  - pairs  $\begin{pmatrix} X \\ X' \end{pmatrix}$  in  $\mathcal{C}$

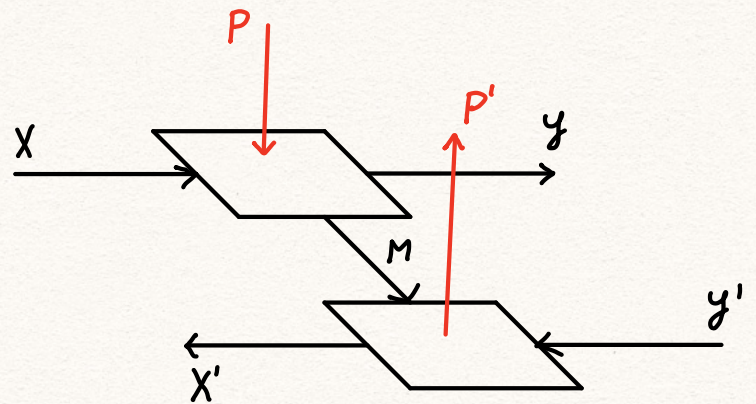
• Morphisms  $\begin{pmatrix} X \\ X' \end{pmatrix} \xrightarrow{((P), f)} \begin{pmatrix} Y \\ Y' \end{pmatrix}$  where  $f: \begin{pmatrix} P \otimes X \\ P' \otimes X' \end{pmatrix} \longrightarrow \begin{pmatrix} Y \\ Y' \end{pmatrix}$

$(M, f, b)$

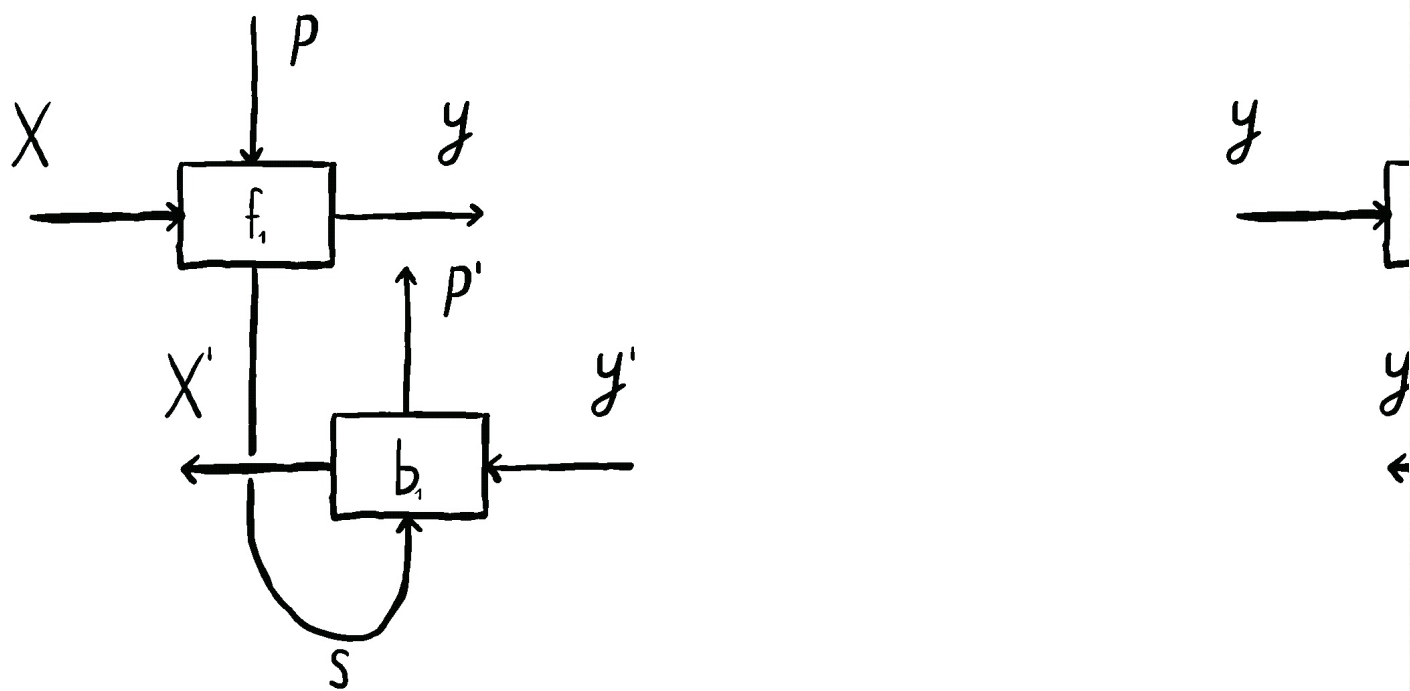
$M: \mathcal{C}$

$f: P \otimes X \longrightarrow Y \otimes M$

$b: Y' \otimes M \longrightarrow P' \otimes X'$



• WE CAN COMPOSE PARAMETERIZED OPTICS



• We automatically get two parameter ports

• A 2-cell  $(X, S) \xrightarrow{\tau} (y, R)$  is an optic

$$\begin{pmatrix} z \\ w \end{pmatrix} \xrightarrow{\tau} \begin{pmatrix} p \\ q \end{pmatrix}$$

## THEOREM.

GRADIENT DESCENT IS A 2-cell IN  $\text{Para}(\text{Optic}(\mathcal{C}))$ .

(Since it is a lens)





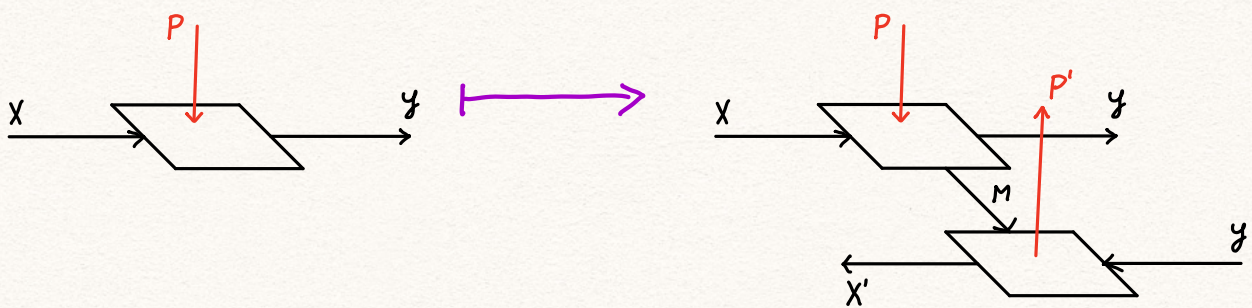
# THEOREM.

APPLYING  $\text{Para}$  TO THE CRDC FUNCTOR

$$\mathcal{C} \xrightarrow{F} \text{Optic}(\mathcal{C})$$

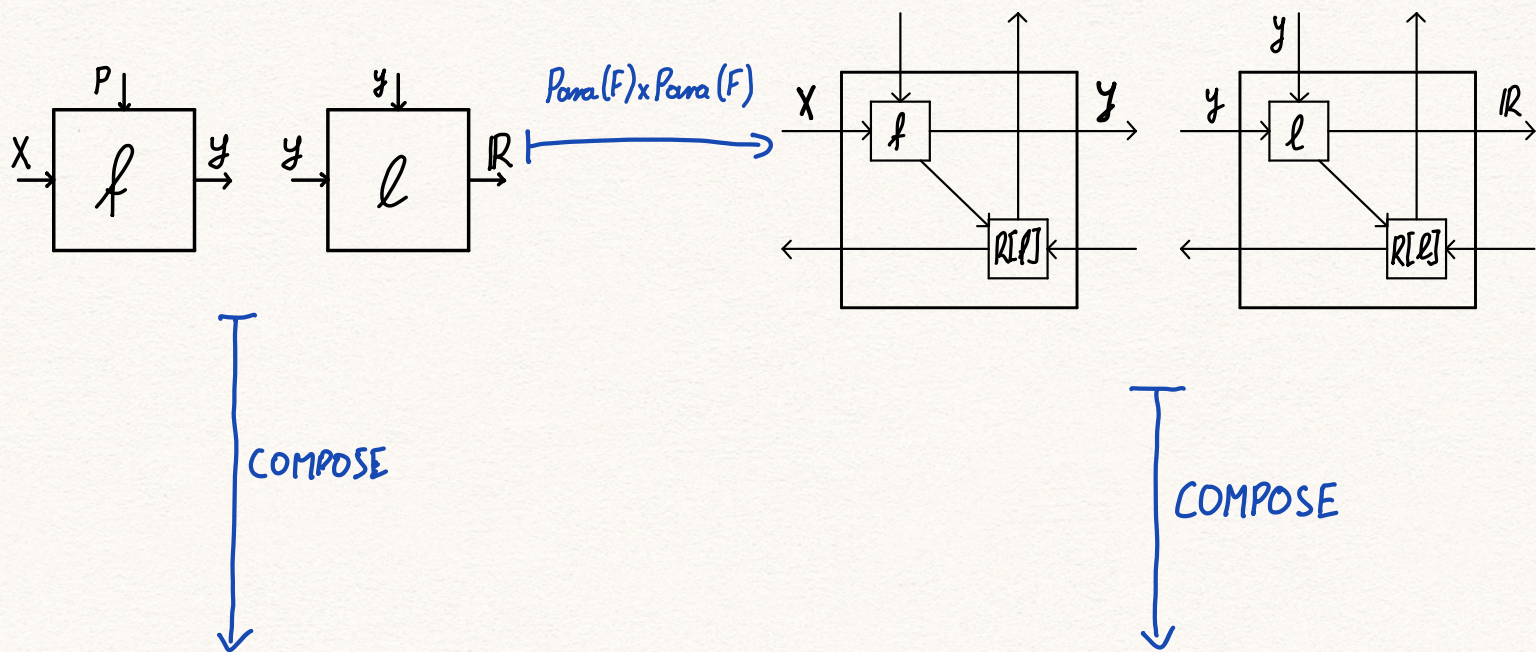
RESULTS IN A FUNCTOR

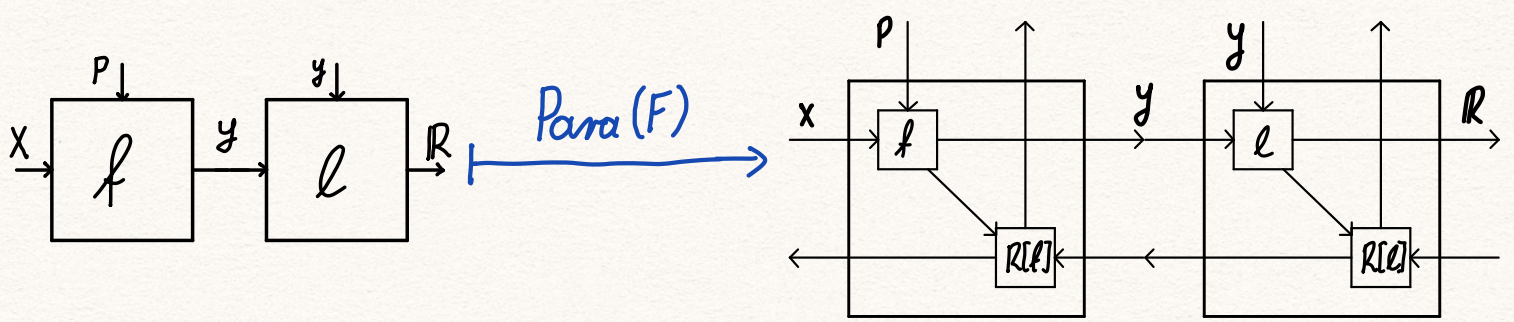
$$\text{Para}(\mathcal{C}) \xrightarrow{\text{Para}(F)} \text{Para}(\text{Optic}(\mathcal{C}))$$



• FUNCTORIALITY IS IMPORTANT!

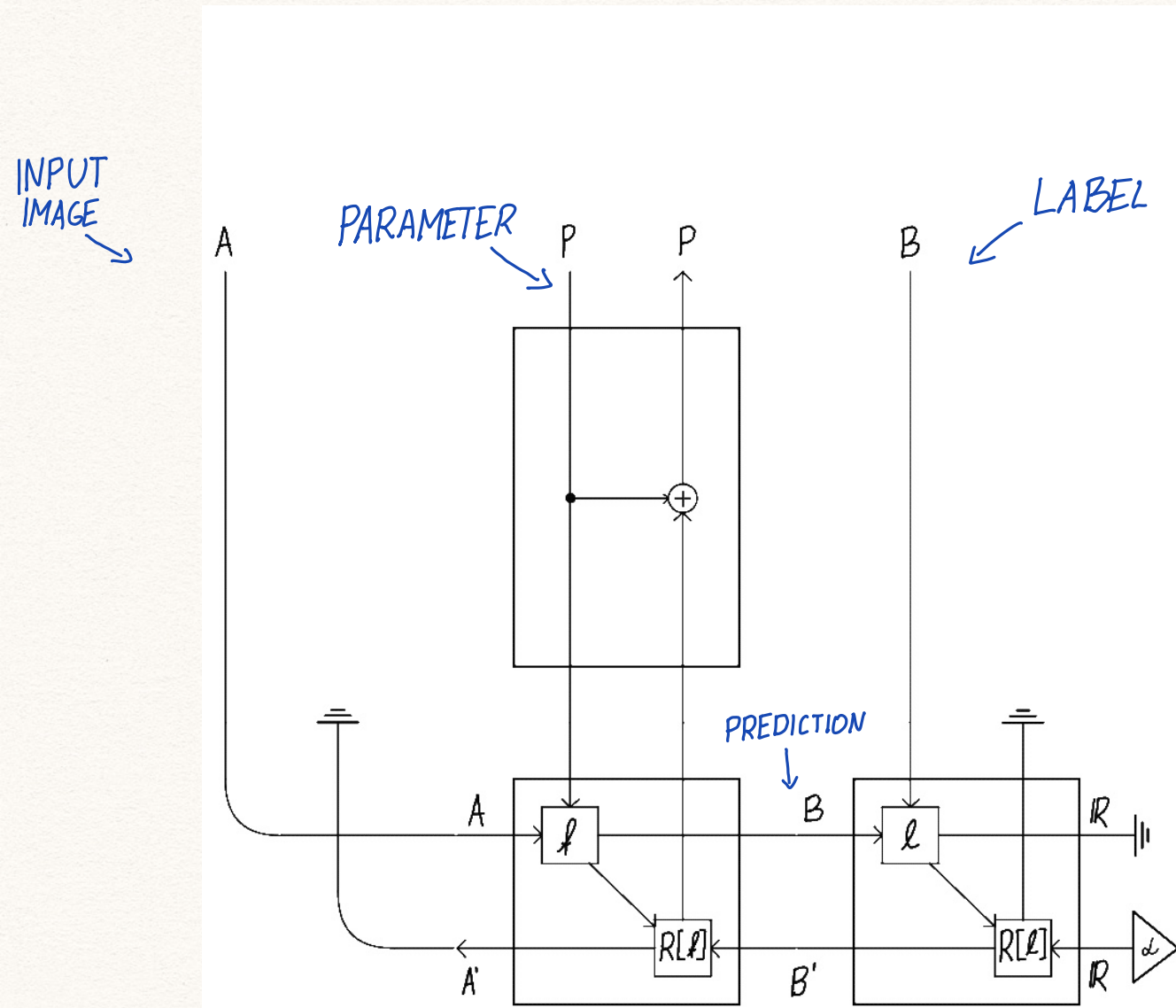
EXAMPLE. A NEURAL NETWORK + A LOSS FUNCTION



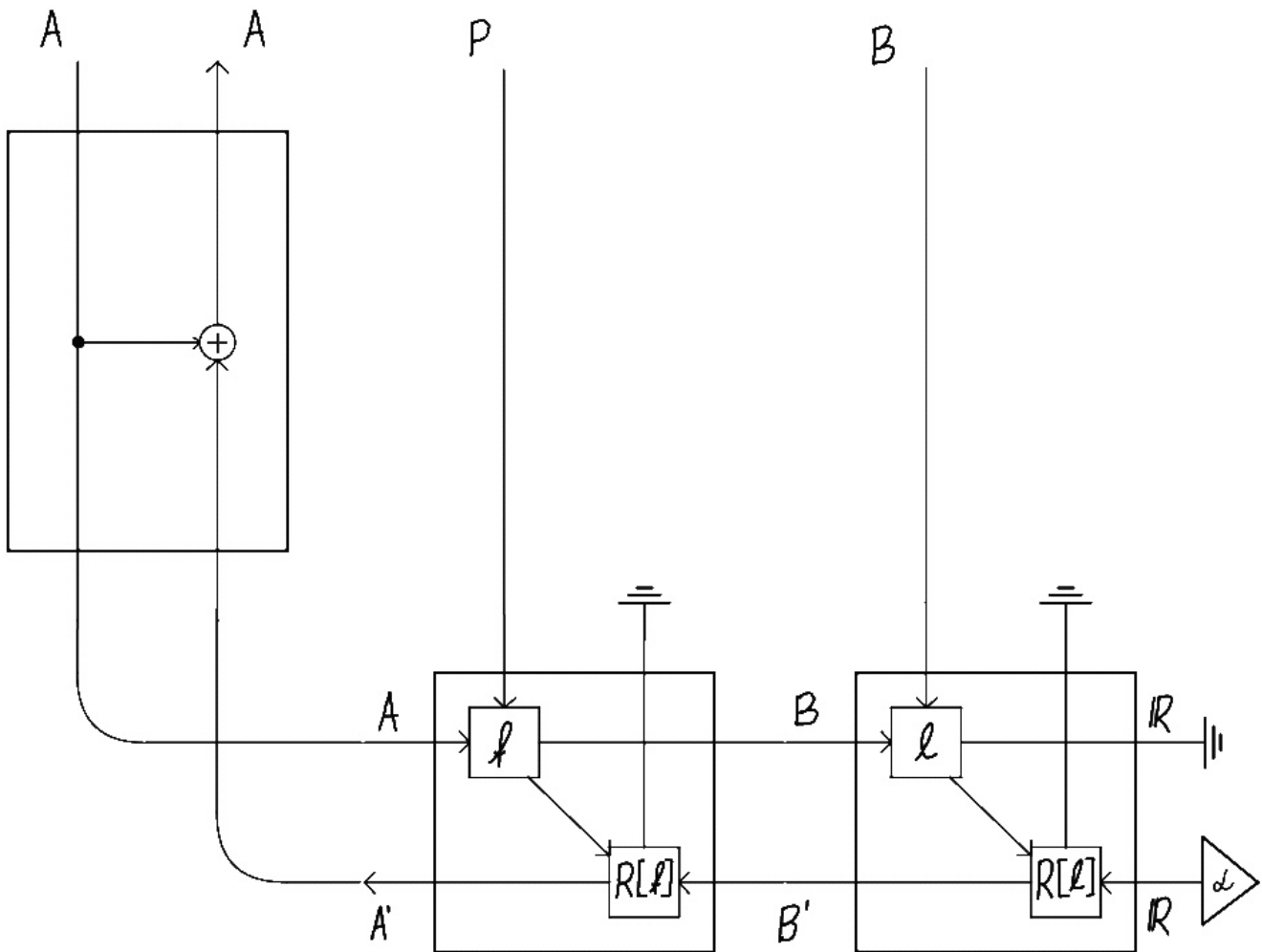


WE CAN PUT THE PIECES TOGETHER.

# SUPERVISED LEARNING



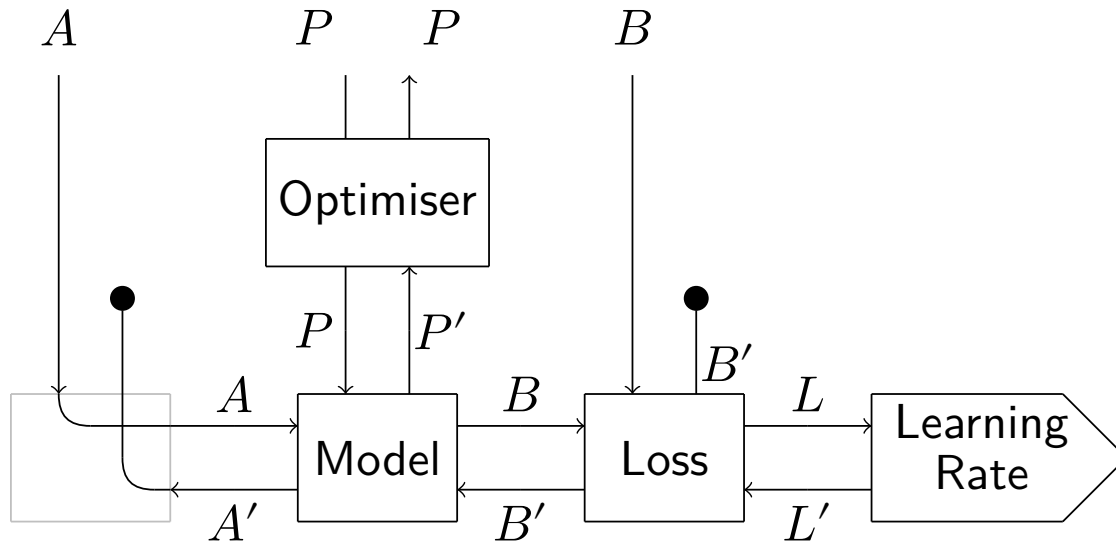
# DEEP DREAMING



# Categorical Foundations of Gradient-Based Learning

# How to Build a Neural Network out of Lenses

Using Para and Lens we get a high level picture



Now we'll see some examples of what can be plugged into each of these boxes.

# The Setting

- ▶ Each box in the diagram is a pair of maps
- ▶ Guiding example: Simple hidden layer neural network, basic gradient descent, MSE loss.
- ▶ We'll specify each pair of maps for each box
- ▶ Goal: you (roughly) understand how to translate this into code
- ▶ Implementation:  
[github.com/statusfailed/numeric-optics-python/](https://github.com/statusfailed/numeric-optics-python/)
  - ▶ examples include a convolutional image classifier for MNIST<sup>1</sup>

---

<sup>1</sup>Lecun et al., "Gradient-Based Learning Applied to Document Recognition."

# Supervised Learning

In supervised learning, we want to learn a map

$$f : A \rightarrow B$$

from a dataset of examples

$$(a, b) \in A \times B$$

Now, based on our beliefs about the structure of  $A$  and  $B$ , we design a *parametrised* map:

$$\text{model} : P \times A \rightarrow B$$

and we search for some  $\theta \in P$  such that  $\text{model}(\theta, -)$  best represents the data.

# Gradient-Based Learning

We want to use a datapoint  $(a, b) \in A \times B$  to improve  $\theta$ , so we need a map

$$??? : P \times A \times B \rightarrow P$$

The reverse derivative is almost what we want. For a map  $f : A \rightarrow B$ ,

$$R[f] : A \times B' \rightarrow A'$$

(while in an RDC  $A' = A$  and  $B' = B$ , it's useful think of the “primed” objects as representing **changes**)

So the reverse derivative of our model morphism has the following type:

$$R[\text{model}] : P \times A \times B' \rightarrow P' \times A'$$



# Updates, “Displacement” and Reverse Derivatives

This is not quite enough: we have two problems:

1. We have a “true” value  $b \in B$  and a “predicted” value  $\text{model}(\theta, a) \in B$  but we need a  $B'$
2. The reverse derivative gives us a  $P'$  and we want a  $P$

This is exactly what the update and loss lenses are for:

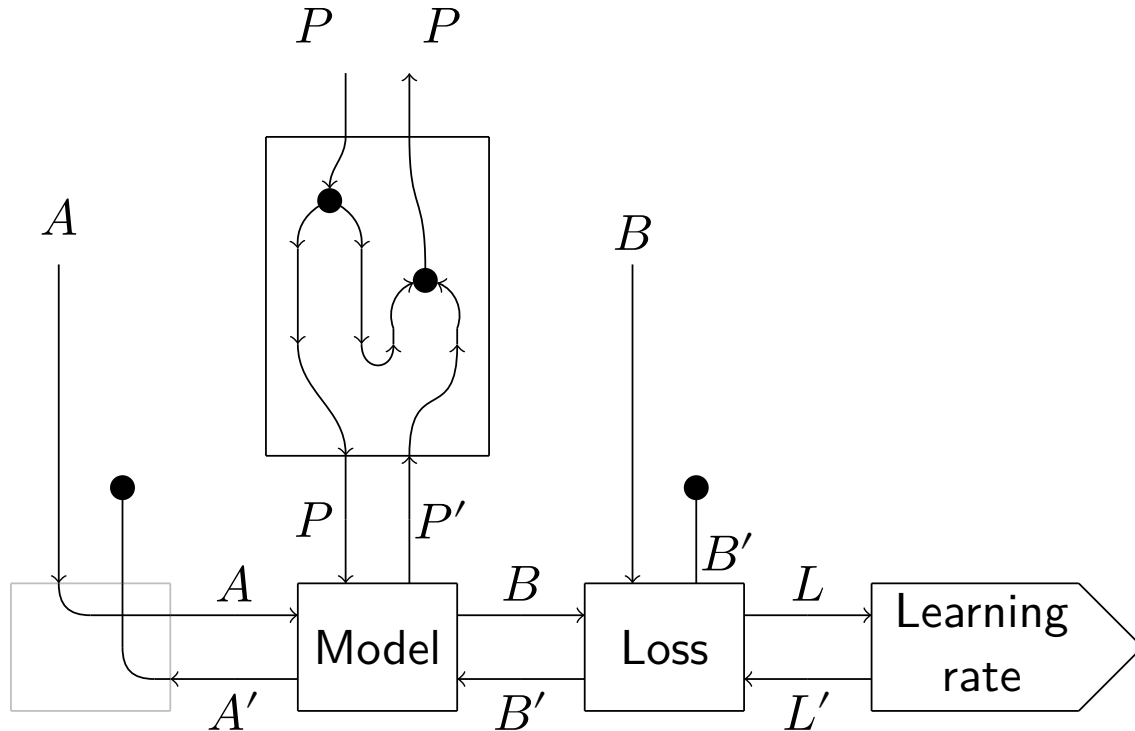
$$R[\text{model}] : P \times A \times B' \rightarrow P' \times A'$$

$$\text{loss}_{\text{put}} : B \times B \rightarrow B' \times B'$$

$$\text{update}_{\text{put}} : P \times P' \rightarrow P$$

# Updates

Updates are like “generalised addition”: add a vector to a point. The most obvious choice is just to add! That’s basic gradient descent:



where  $\rightarrow \bullet \leftarrow$  is copying and  $\rightarrow \bullet \rightarrow$  is addition

## Updates 2

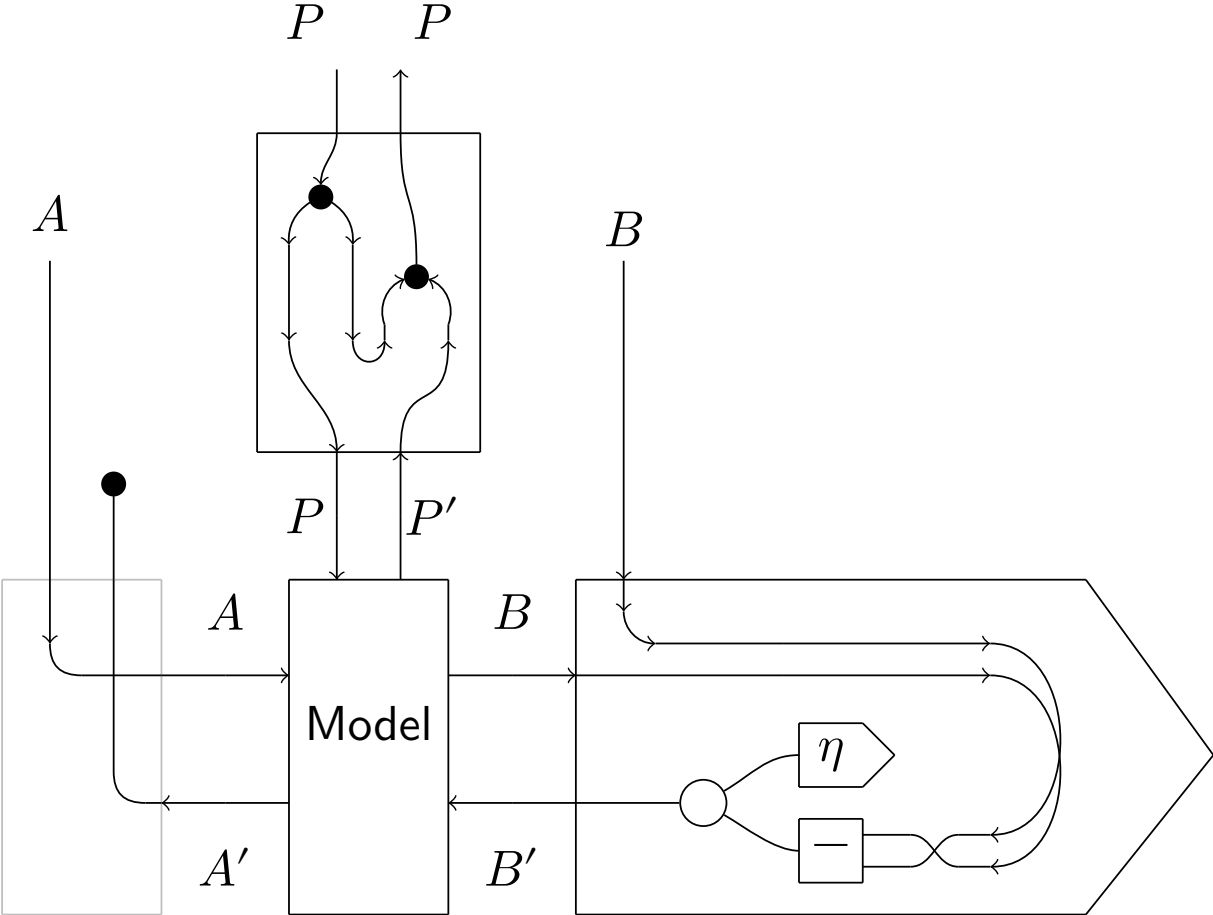
So basic gradient descent is comprised of this pair of maps:

$$\begin{aligned} \text{get} : P &\rightarrow P \\ \theta &\mapsto \theta \end{aligned}$$

$$\begin{aligned} \text{put} : P \times P' &\rightarrow P \\ \theta \quad \theta' &\mapsto \theta + \theta' \end{aligned}$$

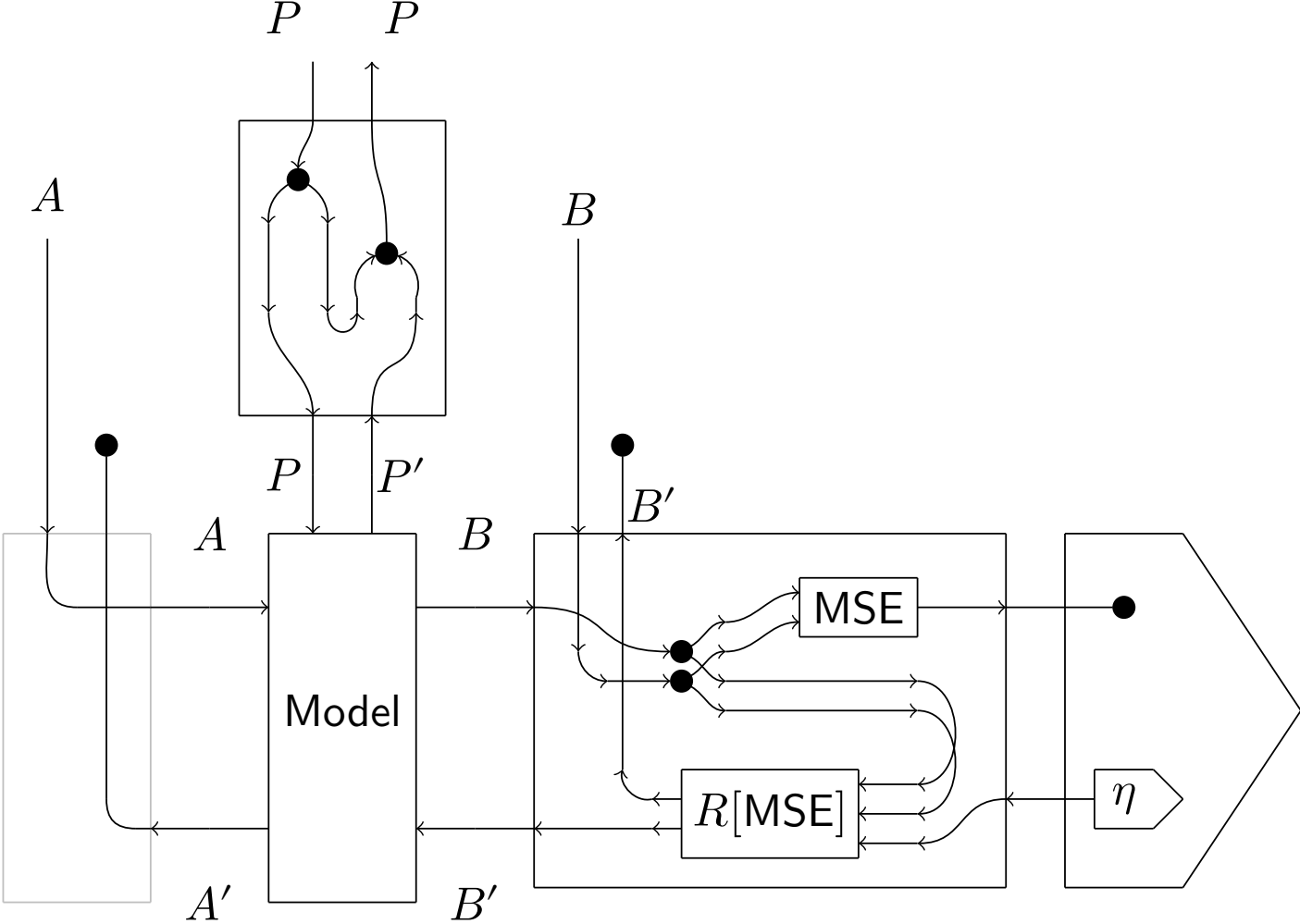
# Loss + Learning Rate

Simple choice is just to subtract:



# Loss + Learning Rate

This is just MSE Loss + fixed learning rate!



# Loss + Learning Rate

We can think of MSE loss as the parametrised lens with maps

$$\text{get} : B \times B \rightarrow \mathbb{R}$$

$$y \quad \hat{y} \mapsto \frac{1}{2n} \sum_i^n (y_i - \hat{y})^2$$

$$\text{put} : B \times B \times \mathbb{R} \rightarrow P$$

$$y \quad \hat{y} \quad l' \mapsto l'(\hat{y} - y)$$

And the fixed learning rate as

$$\text{get} : \mathbb{R} \rightarrow I$$

$$l \mapsto \langle \rangle$$

$$\text{put} : \mathbb{R} \times I \rightarrow \mathbb{R}$$

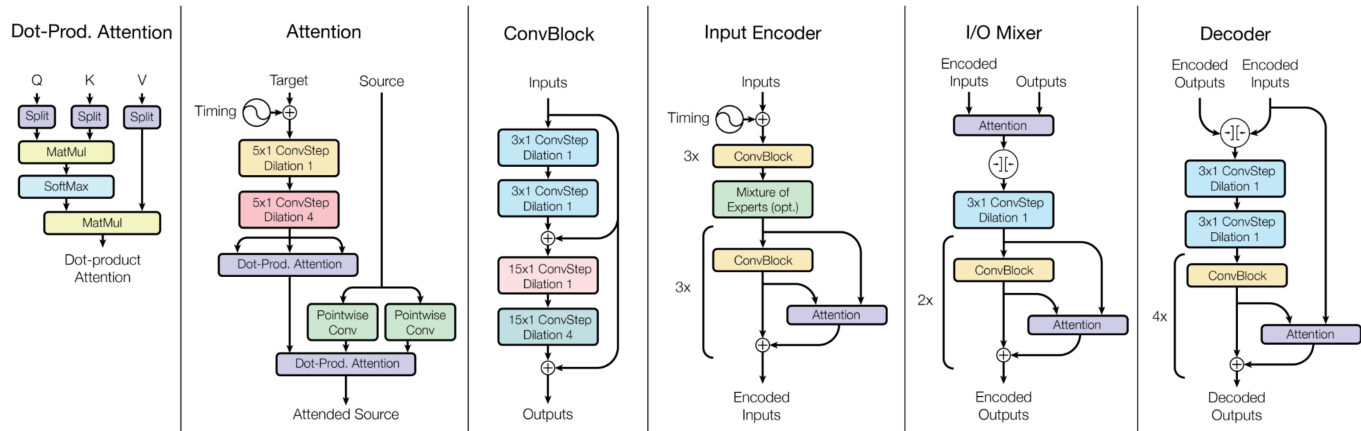
$$l \mapsto \eta$$

# Models, Architectures, and Layers

Two levels of detail in the model: “architecture” and “layers”.

- ▶ Architecture: the whole program as a collection of subroutines (a composition of parametrised lenses)
- ▶ Layer<sup>2</sup>: an individual subroutine (a parametrised lens / pair of maps)

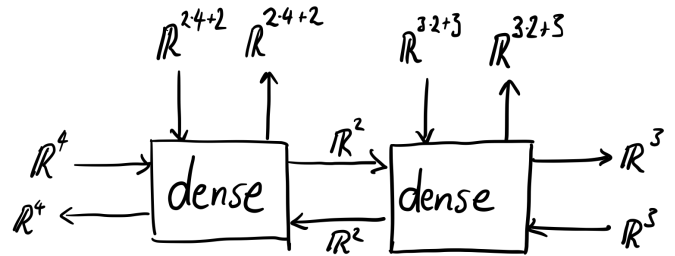
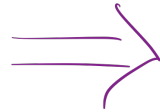
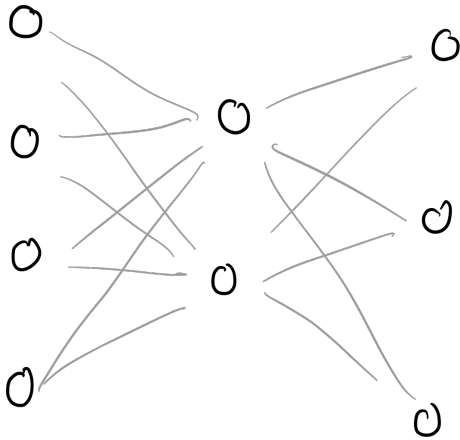
Example of a complicated architecture<sup>3</sup>:



<sup>2</sup>ambiguous terminology warning

<sup>3</sup>Kaiser et al., “One Model to Learn Them All.”

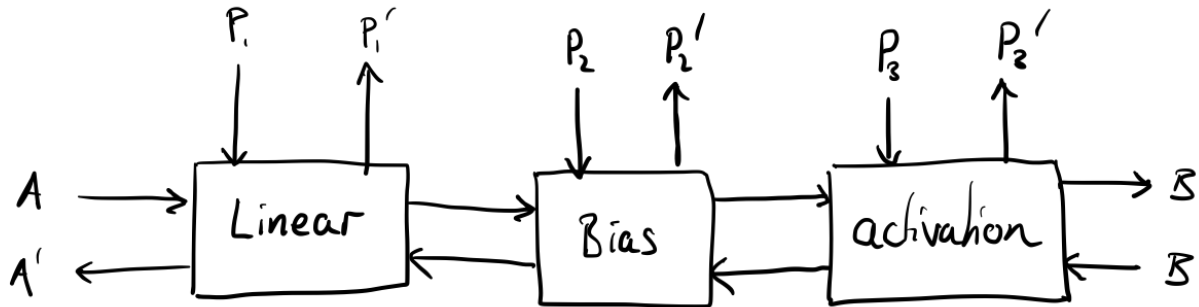
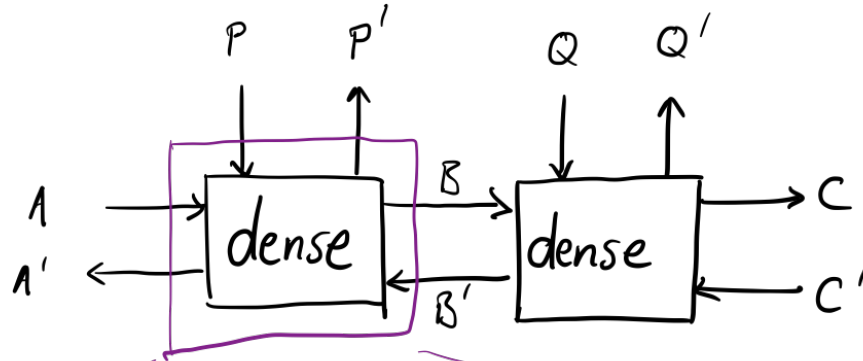
# The Old Ways



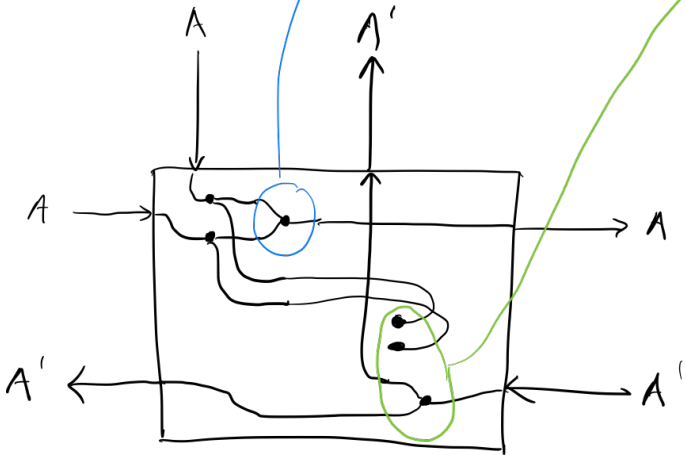
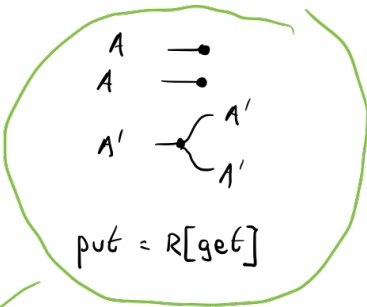
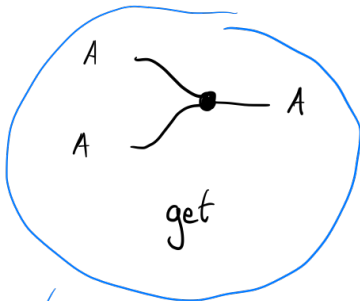


# Dense Layers

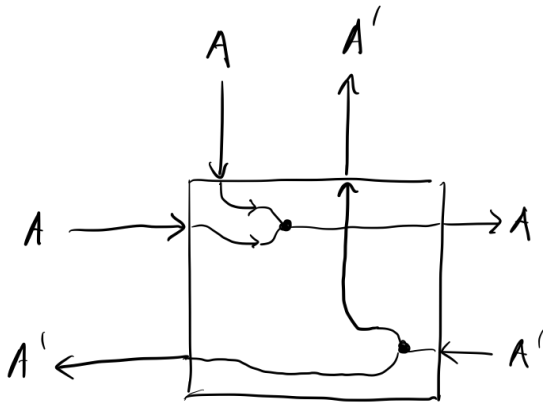
A simple hidden layer neural network is a composition of two dense layers. Let's unpack a dense layer and see what's inside...



# Bias Layers



Simplify  
⇒



# Linear Layers

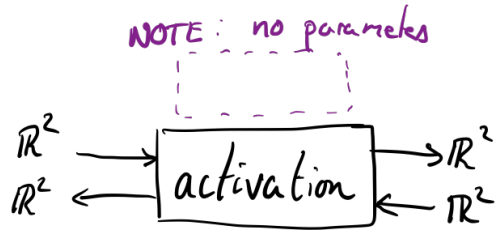
- ▶ Parameters  $P = \mathbb{R}^{b \cdot a}$  are the coefficients of a matrix
- ▶ Input  $A = \mathbb{R}^a$  is an  $a$ -dimensional vector
- ▶ Forward pass multiplies the matrix by the vector:

$$\begin{aligned} \text{get} &: \text{Mat}(A, B) \times \text{Vec}(A) \rightarrow \text{Vec}(B) \\ \text{get}(M, x) &\mapsto Mx \end{aligned}$$

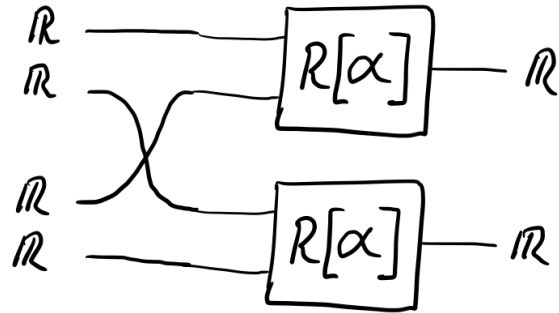
- ▶ Reverse pass does this (note that it typechecks!):

$$\begin{aligned} \text{put} &: \text{Mat}(A, B) \times \text{Vec}(A) \times \text{Vec}(B) \rightarrow \text{Mat}(A, B) \times \text{Vec}(A) \\ \text{put}(M, x, y) &\mapsto \langle y \otimes x, M^T y \rangle \end{aligned}$$

# Activation Layer



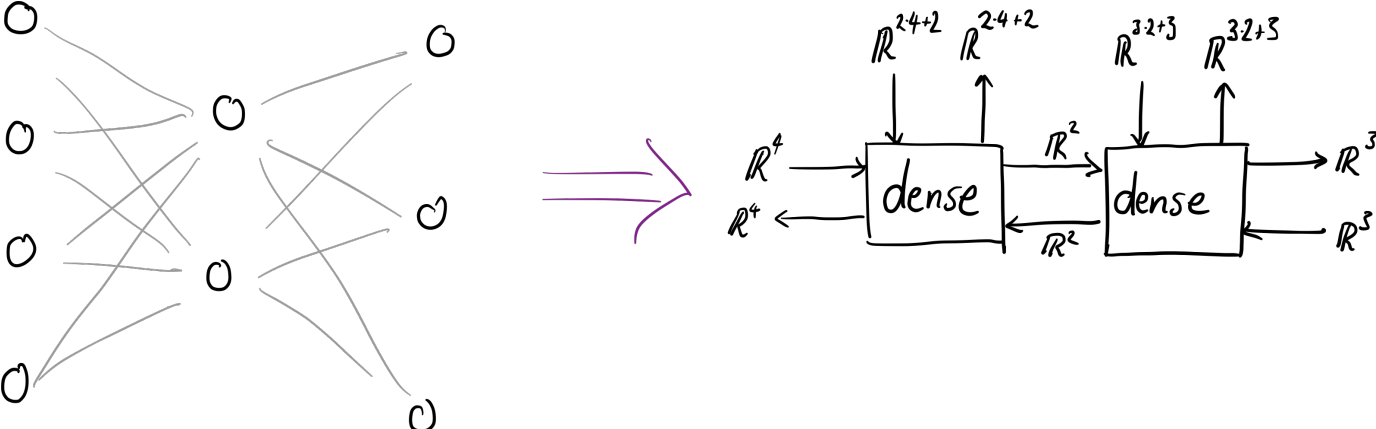
activation<sub>GET</sub>



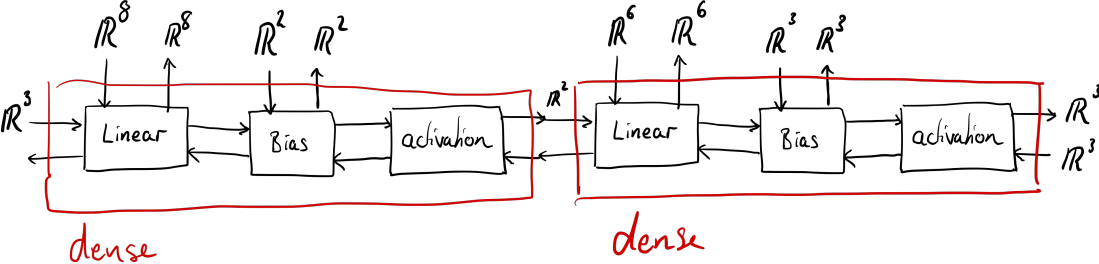
activation<sub>PUT</sub>

# Hidden Layer Neural Network

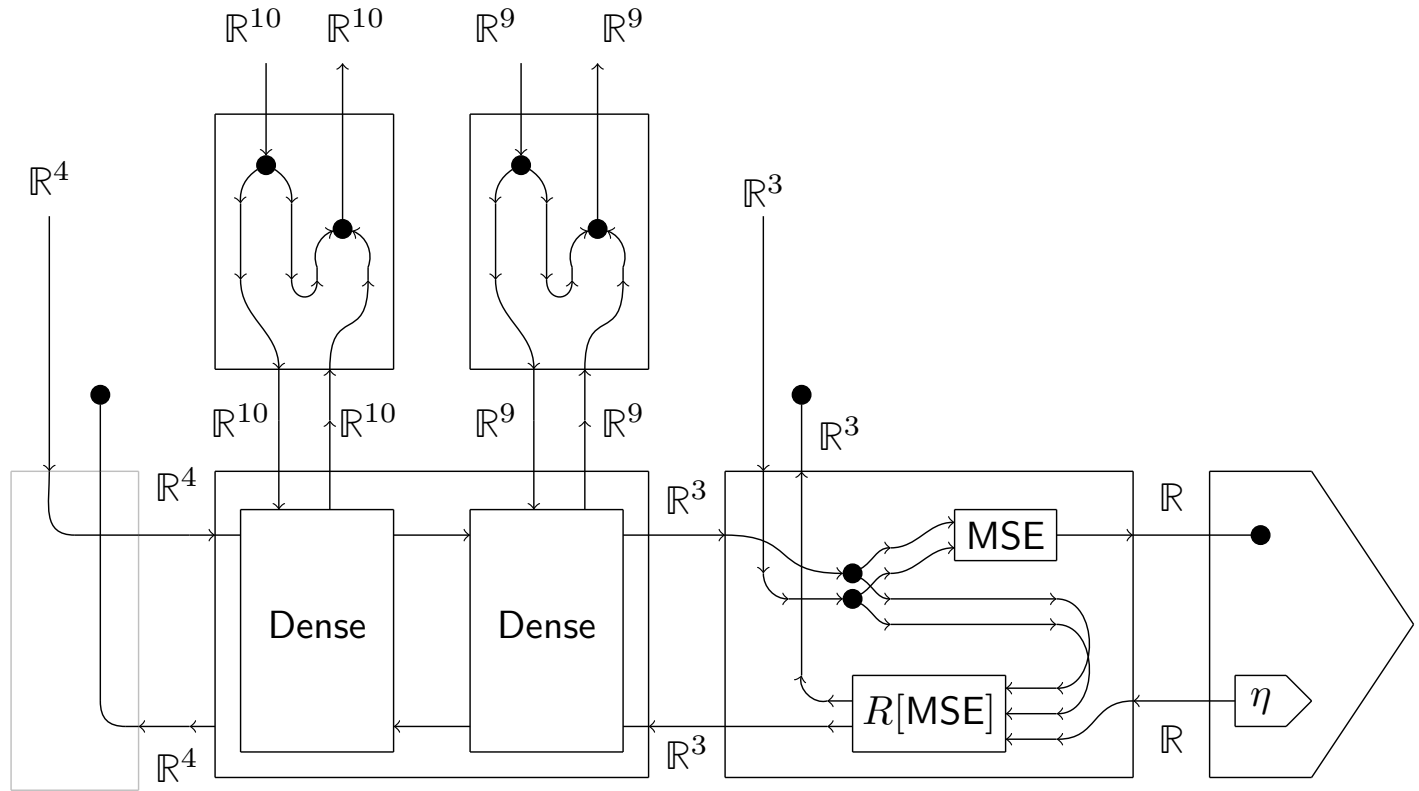
Returning to the "standard" picture of a neural network:



Expanding out "dense":



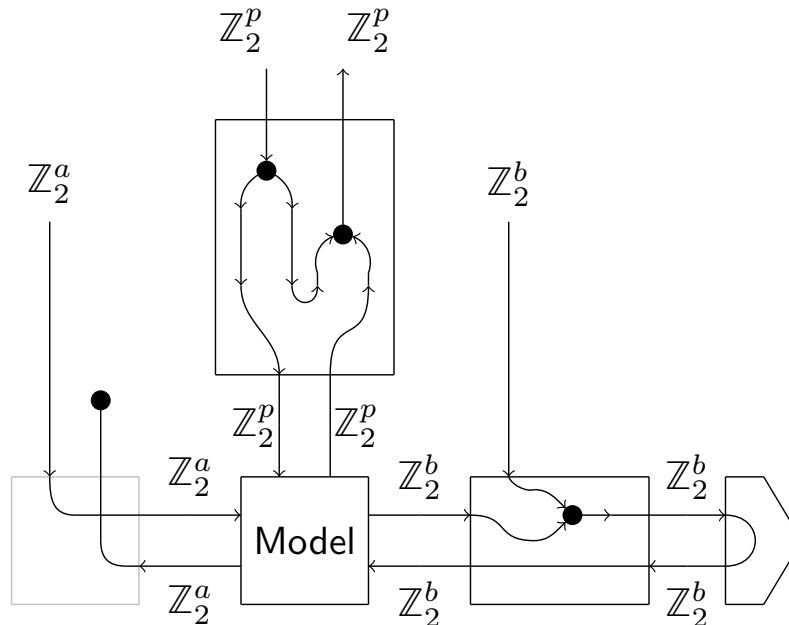
# A Hidden Layer Neural Network as a Parametrised Lens



... with MSE loss, basic gradient descent, and fixed error rate

# What else can we plug in?

- ▶ So far we've only seen neural networks, where objects are  $\mathbb{R}^n$  for  $n \in \mathbb{N}$ .
- ▶ We can do learning with boolean circuits too, as in Reverse Derivative Ascent<sup>4</sup>:



<sup>4</sup>Wilson and Zanasi, "Reverse Derivative Ascent."

Questions?



# Reverse Derivatives, Graphically

## CARTESIAN STRUCTURE

copy



$$x \mapsto \langle x, x \rangle$$

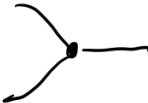
discard



$$x \mapsto \langle \rangle$$

## LEFT-ADDITIVE STRUCTURE

add



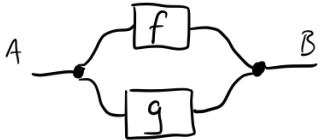
$$x_1, x_2 \mapsto \langle x_1 + x_2 \rangle$$

zero

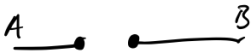


$$\langle \rangle \mapsto \langle 0 \rangle$$

Addition & zero maps



$$f + g : A \rightarrow B$$

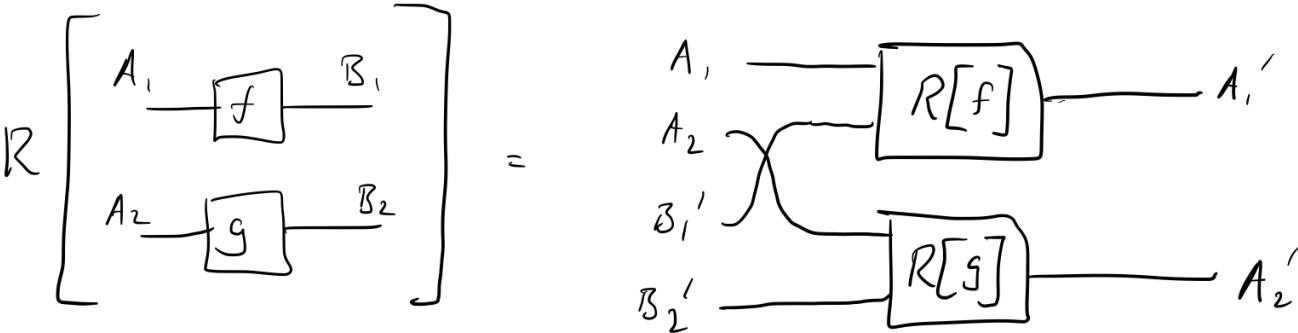
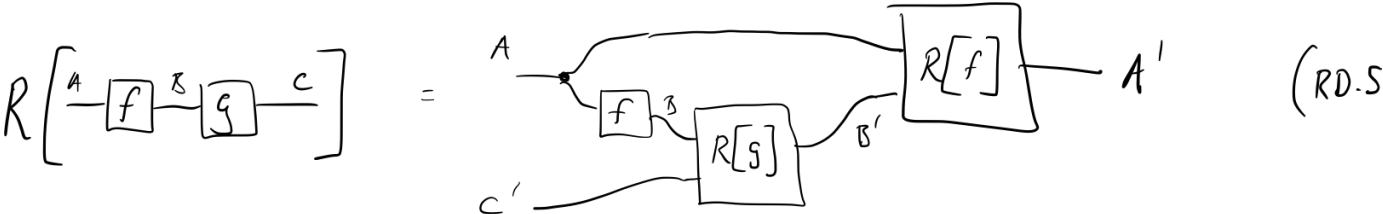


$$0 : A \rightarrow B$$



# Reverse Derivatives, Graphically

$$A \xrightarrow{f} B \quad \Rightarrow \quad A \times B' \xrightarrow{R[f]} A'$$



# References

- Kaiser, Lukasz, Aidan N. Gomez, Noam Shazeer, Ashish Vaswani, Niki Parmar, Llion Jones, and Jakob Uszkoreit. “One Model to Learn Them All,” 2017. <http://arxiv.org/abs/1706.05137>.
- Lecun, Yann, Léon Bottou, Yoshua Bengio, and Patrick Haffner. “Gradient-Based Learning Applied to Document Recognition.” In *Proceedings of the Ieee*, 2278–2324, 1998. <https://doi.org/10.1109/5.726791>.
- Wilson, Paul, and Fabio Zanasi. “Reverse Derivative Ascent: A Categorical Approach to Learning Boolean Circuits.” *Electronic Proceedings in Theoretical Computer Science* 333 (February 2021): 247–60. <https://doi.org/10.4204/eptcs.333.17>.