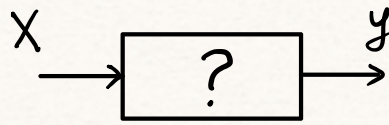# GENERALIZED OPEN LEARNERS (AGENTS)

**WHAT:** UNDERSTAND THE INFORMATION FLOW IN ARTIFICIAL NEURAL NETWORKS
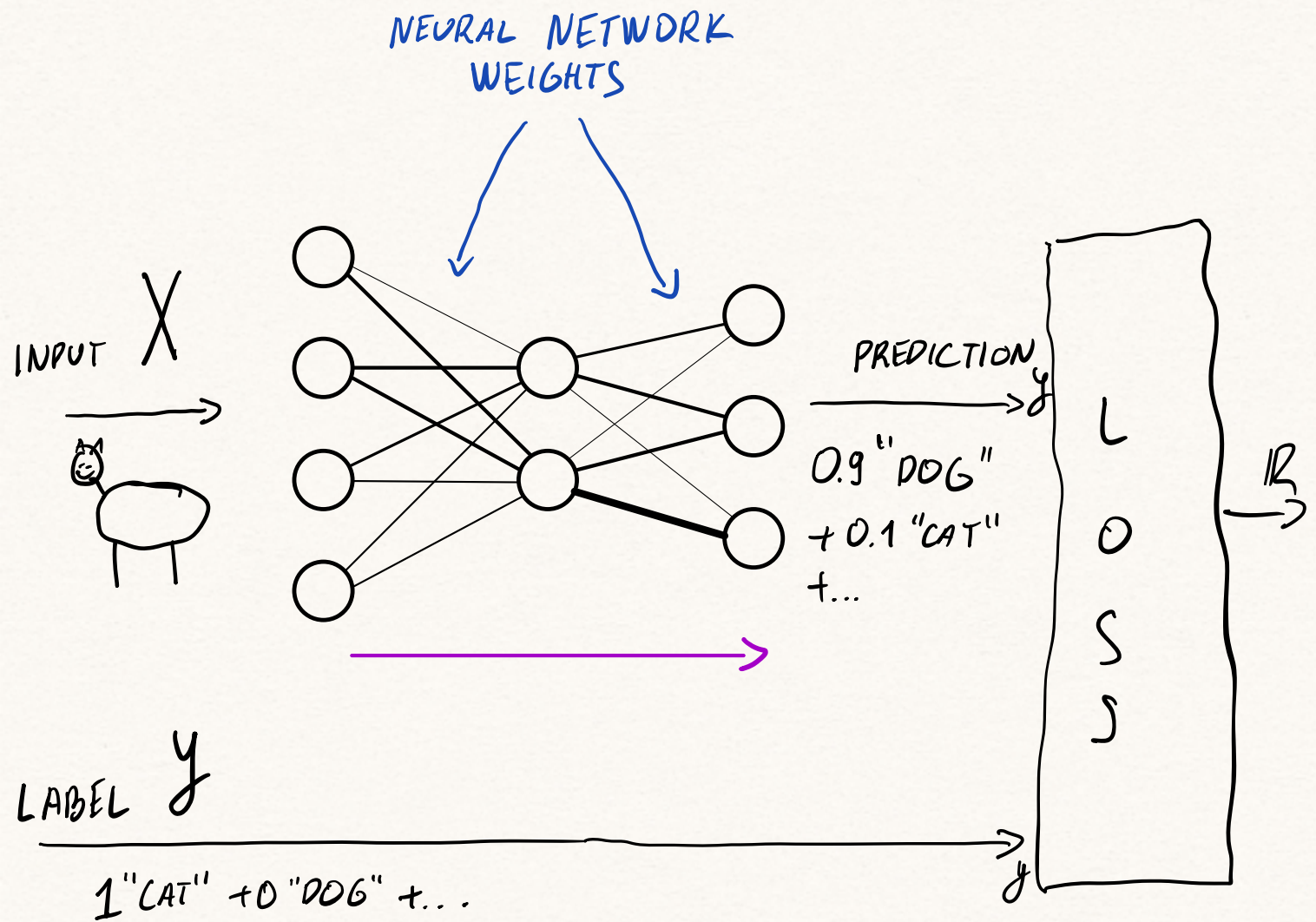
**WHY:** CONCRETE PROXY FOR UNDERSTANDING GENERAL CYBERNETICS SYSTEMS

**HOW:** TAKE PARAMETERIZATION AND BIDIRECTIONALITY SERIOUSLY
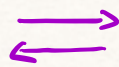
# SUPERVISED LEARNING WITH NEURAL NETWORKS IN ONE SLIDE

$$X \rightarrow \boxed{?} \rightarrow y$$

## DATASET: List $X \times Y$

NEURAL NETWORK WEIGHTS

INPUT $X$

PREDICTION $y$

0.9 "DOG"
+0.1 "CAT"
+...

L
O
S
S

$\mathbb{R}$

LABEL $Y$

$y$

1 "CAT" + 0 "DOG" + ...

- WEIGHTS CONTROL THE NN PERFORMANCE
- PROPAGATING CHANGES - „BACKPROPAGATION"

ONE STEP, THIS IS ITERATED

"PARAMETERIZATION" and "BIDIRECTIONALITY"

can be defined <u>separately</u> and then composed

- GENERALIZED OPEN LEARNERS (Agents)

- Abstract treatment of gradient-based learning
  - Backpropagation — ROCs
  - <u>Loss function</u>
  - <u>Update rule</u>

                              SIMPLER, PEDAGOGICAL

- Works on Euclidean spaces and Boolean circuits

- Defined optimizers (GRADIENT DESCENT, MOMENTUM, ADAM...)
                                              as lenses

- Optimizers are 2-cells in GEN. OPEN LEARNERS

# Game Plan:

- Parameterization
- Bidirectionality
- Parameterization + Bidirectionality
- Differentiation (How Do We Construct Learners?)
- How does learning work?
- Examples

# PARAMETERIZATION

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

**DEF.** $\mathsf{Para}(\mathcal{C})$

Objects – objects of $\mathcal{C}$

CATEGORY OF ELEMENTS

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

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

2-cells are reparameterizations: a 2-cell

$$A \overset{(P,f)}{\underset{(Q,g)}{\Downarrow r}} B$$

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

$$Q \otimes A \xrightarrow{r \otimes A} P \otimes A$$
$$g \searrow \quad \swarrow f$$
$$B$$

**EXAMPLE.**

$\mathsf{Para}(\mathsf{Set})$, $\mathsf{Para}(\mathsf{Smooth})$, $\mathsf{Para}(\mathsf{Optic}(\mathcal{C}))$, ...

$\underset{\mathsf{Para}}{\cong}$
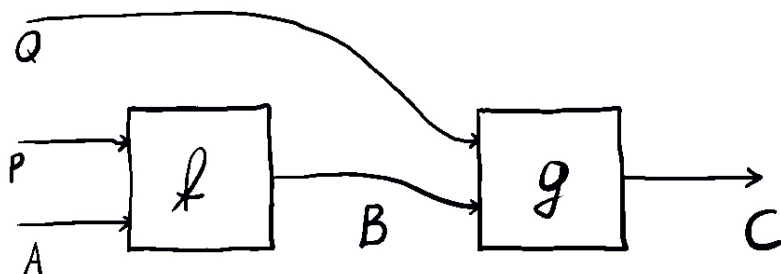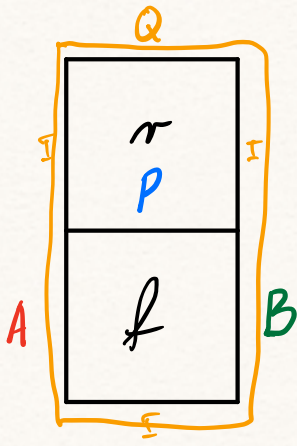
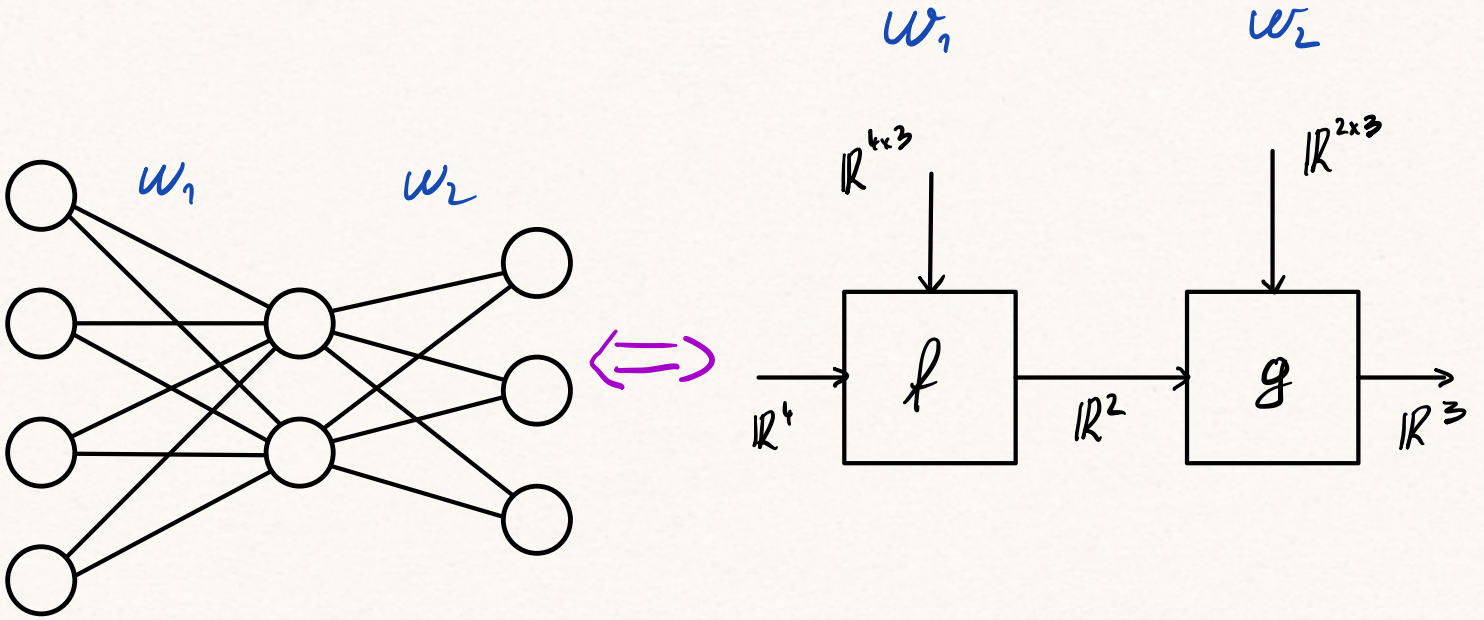# GRAPHICAL LANGUAGE

$f : P \otimes A \longrightarrow B$
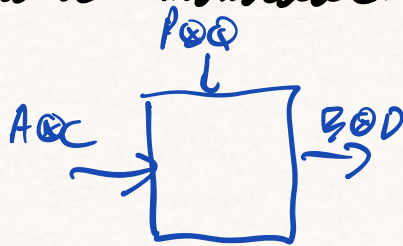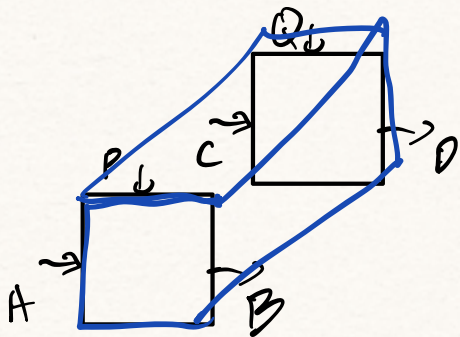






I DIDN'T TELL YOU HOW COMPOSITION WORKS'

- TILING LANGUAGE OF DOUBLE CATEGORIES (MYERS)

**PROP.** $Para(\mathcal{C})$ is symmetric monoidal.



$Para$ is natural w.r.t. base change.

**DEF.**

Let $G: \mathcal{C} \longrightarrow \mathcal{D}$ be a sym. monoidal functor. We

define $\quad Para(G): Para(\mathcal{C}) \longrightarrow Para(\mathcal{D})$ ⟵ — RELEVANT TO BACKPROPAGATION

$$
\begin{array}{ccc}
A & \longmapsto & G(A) \\
{\scriptstyle(P,\,f)}\Big\downarrow & & \Big\downarrow{\scriptstyle(G(P),\,f')} \\
B & \longmapsto & G(B)
\end{array}
$$

where $f'$ is the composite

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

**ALSO:** $Para$ is an endofunctor on SMC (also a monad) + a lot more structure
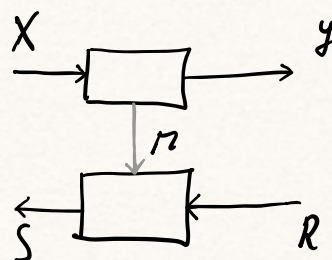
# BIDIRECTIONALITY
## (lenses, optics, dependent lenses, Cont, Poly, ... )

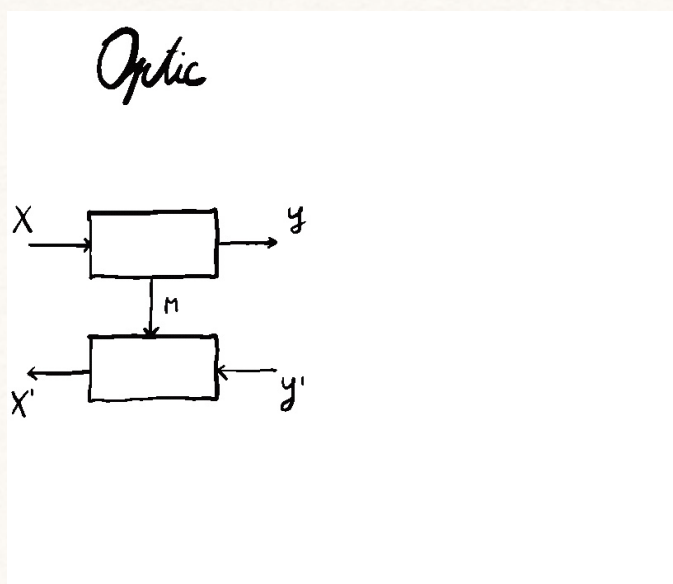- CURRENTLY FOCUSING ON OPTICS AS THE CANONICAL "BIDIRECTIONAL" STRUCTURE, OPEN TO AMENDMENTS

**DEF.** Category $\text{Optic}(\mathcal{C})$

- Objects – pairs of objects in $\mathcal{C}$ $\binom{X}{S}$

- $\text{Optic}(\mathcal{C})\binom{X}{S}, \binom{Y}{R} = \int^{M:\mathcal{C}} \mathcal{C}(X, M\otimes Y) \times \mathcal{C}(M\otimes R, S)$

COEND

"FROM THE OUTSIDE"

Optic

**PROP.** If $\mathcal{C}$ is Cartesian, then $\text{Lens}(\mathcal{C}) \cong \text{Optic}(\mathcal{C})$

Lens($\mathcal{C}$)

**PROP.** $\mathrm{Optic}(\mathcal{C})$ is symmetric monoidal.

## EXAMPLE. GRADIENT DESCENT



$$P \times P \xrightarrow{\;u\;} P$$
$$(p, \nabla p) \longmapsto p - \alpha \nabla p$$

is a lens, for $\mathcal{C} := \mathrm{Smooth}$

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

## EXAMPLE. OTHER OPTIMIZERS:
- MOMENTUM,

$$get: P \times P \longrightarrow P$$
$$(v, p) \longmapsto p$$

$$put: P \times P \times P \longrightarrow P \times P$$
$$(v, p, \nabla p) \longmapsto (v', p - v')$$
$$\text{where } v' = \gamma v + \varepsilon p'$$

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

- NESTEROV MOMENTUM

$$get: P \times P \longrightarrow P$$
$$(v, p) \longmapsto p - \gamma v$$

put - same as above

- ADAGRAD

- ADAM
...

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

- Objects — objects of $\text{Optic}(\mathcal{C})$ — pairs $\begin{pmatrix} X \\ S \end{pmatrix}$ in $\mathcal{C}$

- Morphisms $\quad \begin{pmatrix} X \\ S \end{pmatrix} \xrightarrow{\left(\binom{P}{Q}, f\right)} \begin{pmatrix} y \\ R \end{pmatrix}$ where $f: \begin{pmatrix} P \otimes X \\ Q \otimes S \end{pmatrix} \longrightarrow \begin{pmatrix} y \\ R \end{pmatrix}$

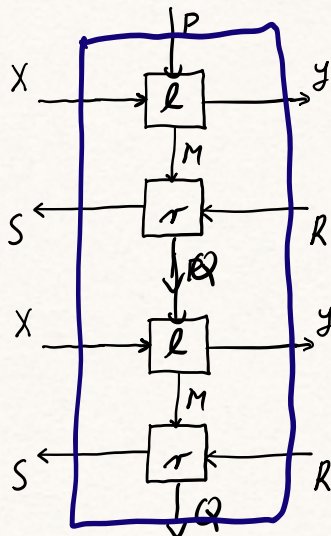$\begin{pmatrix} M, \ell, r \end{pmatrix}$

$M: \mathcal{C}$

$\ell: P \otimes X \longrightarrow M \otimes y$

$r: M \otimes R \longrightarrow Q \otimes S$

$$\begin{pmatrix} X \otimes X \\ S \otimes S \end{pmatrix} \longrightarrow \begin{pmatrix} y \otimes y \\ R \otimes R \end{pmatrix}$$
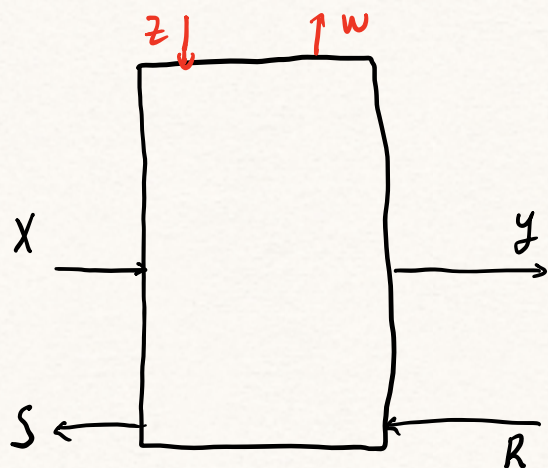


$\begin{pmatrix} x:X \\ S_x \end{pmatrix}$

$S: X \to \text{Set}$



- We automatically get two parameter ports

- A 2-cell $\begin{pmatrix} X \\ S \end{pmatrix} \Downarrow r \begin{pmatrix} y \\ R \end{pmatrix}$ is an optic

$\begin{pmatrix} P \\ Q, \ell \end{pmatrix}$

$\begin{pmatrix} Z \\ W, g \end{pmatrix}$

$\begin{pmatrix} Z \\ W \end{pmatrix} \xrightarrow{r} \begin{pmatrix} P \\ Q \end{pmatrix}$

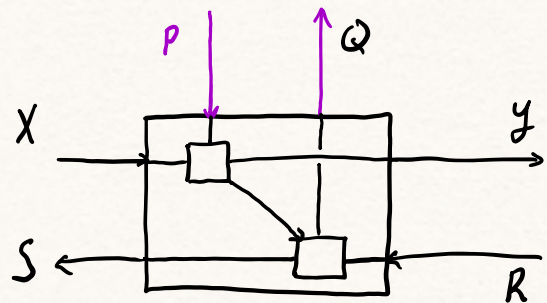**Definition II.1.** *Let A and B be sets. A* supervised learning algorithm, *or simply* learner, *A → B is a tuple* $(P, I, U, r)$ *where P is a set, and I, U, and r are functions of types:*

$$I: P \times A \to B,$$
$$U: P \times A \times B \to P,$$
$$r: P \times A \times B \to A.$$

# THEOREM.   $\text{Learn} \cong \text{Para}(\text{Optic}(\text{Set}))$

# HOW DO WE ACTUALLY CONSTRUCT SOME USEFUL LEARNERS (NEURAL NETWORKS)?

(THIS IS WHAT TENSORFLOW/PYTORCH/JAX ARE DOING)

# DIFFERENTIATION

· Reverse derivative categories (RDC).

**DEF.** A RDC $\mathcal{C}$ is a category which for every

$$A \xrightarrow{f} B \qquad \longleftarrow \text{the „get" map}$$

has a map

$$R[f] : A \times B \longrightarrow A \qquad \text{the „put" map of a lens}$$

subject to some conditions.

**EXAMPLE.** Smooth IS A RDC. BoolCirc IS A RDC

**EXAMPLE.** Let $\mathcal{C} := $ Smooth

Let $\mathbb{R} \xrightarrow{f} \mathbb{R}$ . Then $\mathbb{R} \times \mathbb{R} \xrightarrow{R[f]} \mathbb{R}$

$\qquad\qquad x \longmapsto x^2 \qquad\qquad (x, w) \longmapsto 2xw \qquad f'(x) \cdot w$

**PROP.** For each RDC $\mathcal{C}$ there is a symm. mon. functor

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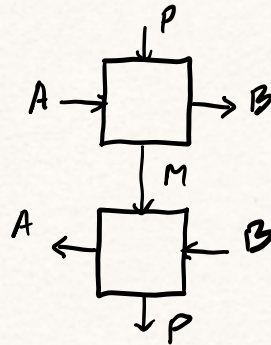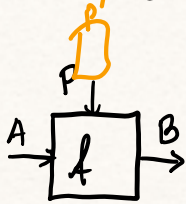\text{Lens}(\mathcal{C}) & \binom{A}{A} \xrightarrow{(f, R[f])} \binom{B}{B}
\end{array}
$$

FUNCTORIALITY OF F IS THE CHAIN RULE.

**THEOREM.** THE ACTION OF $Para$ ON THE FUNCTOR

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

RESULTS IN A FUNCTOR

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



**EXAMPLE.** Consider a linear layer $(\mathbb{R}^2, \ell): Para(Smooth)(\mathbb{R}, \mathbb{R})$
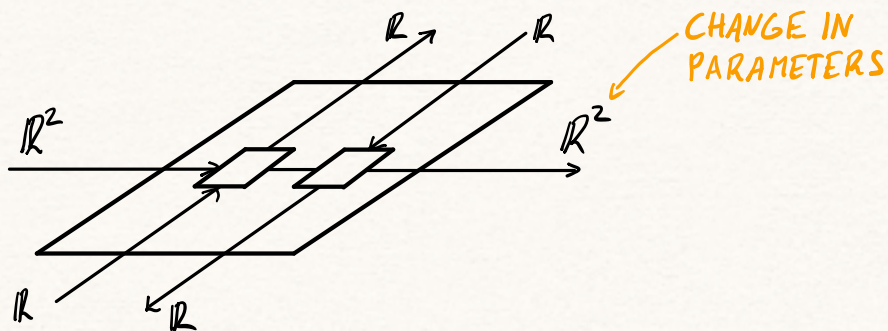
where

$$\mathbb{R}^2 \times \mathbb{R} \xrightarrow{\ell} \mathbb{R}$$
$$((w_0, w_1), x) \longmapsto w_0 + w_1 x$$

Then

$$\mathbb{R}^2 \times \mathbb{R} \times \mathbb{R} \xrightarrow{R[\ell]} \mathbb{R}^2 \times \mathbb{R}$$
$$((w_0, w_1), x, dy) \longmapsto ((dy, x dy), w_1 dy)$$



CHANGE IN PARAMETERS

# HOW DOES (DERIVATIVE BASED) (SUPERVISED) LEARNING HAPPEN?

## 1) UPDATE MAP

THM. Gradient descent is a reparameterization.

$$P \times P \xrightarrow{\;u\;} P$$
$$(p, \nabla p) \longmapsto p - \alpha Dp$$

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



## 2) LOSS FUNCTION



$$\ell : \mathcal{C}(P \times X, y)$$

$$\ell : \mathcal{C}(y \times y, \mathbb{R})$$

$$u : \text{Lens}(\mathcal{C}) \begin{pmatrix} P \\ P \end{pmatrix}, \begin{pmatrix} P \\ P \end{pmatrix}$$

# THEOREM:
- FIX A RDC $\mathcal{C}$

- FIX **LOSS FUNCTION** DATA    (MSE, SOFTMAX CROSS ENTROPY...)

- FIX **PARAMETER UPDATE** DATA (G.D., MOMENTUM, ADAGRAD, ADAM,...)

Then we can define a functor $Para(\mathcal{C}) \longrightarrow Para(Lens(\mathcal{C}))$
which creates a learner.

UPDATE

BACKWARD
PASS

LOSS
FUNCTION

UPDATE
MAP

$$\mathsf{Para}(\mathcal{C}) \xrightarrow{\ \mathsf{Para}(F)\ } \mathsf{Para}(\mathsf{Lens}(\mathcal{C})) \longrightarrow \mathsf{Para}(\mathsf{Lens}(\mathcal{C})) \longrightarrow \mathsf{Para}(\mathsf{Lens}(\mathcal{C}))$$



**Theorem III.2.** *Fix a real number $\varepsilon > 0$ and $e(x, y)\colon \mathbb{R} \times \mathbb{R} \to \mathbb{R}$ differentiable such that $\frac{\partial e}{\partial x}(x_0, -)\colon \mathbb{R} \to \mathbb{R}$ is invertible for each $x_0 \in \mathbb{R}$. Then we can define a faithful, injective-on-objects, strong symmetric monoidal functor*

$$L_{\varepsilon,e}\colon \mathsf{Para} \longrightarrow \mathsf{Learn} \cong \mathsf{Para}(\mathsf{Optic}(\mathsf{Set}))$$
$$\cong$$
$$\mathsf{Para}(\mathsf{Smooth})$$

*that sends each parametrised function $I\colon P \times A \to B$ to the learner $(P, I, U_I, r_I)$ defined by*
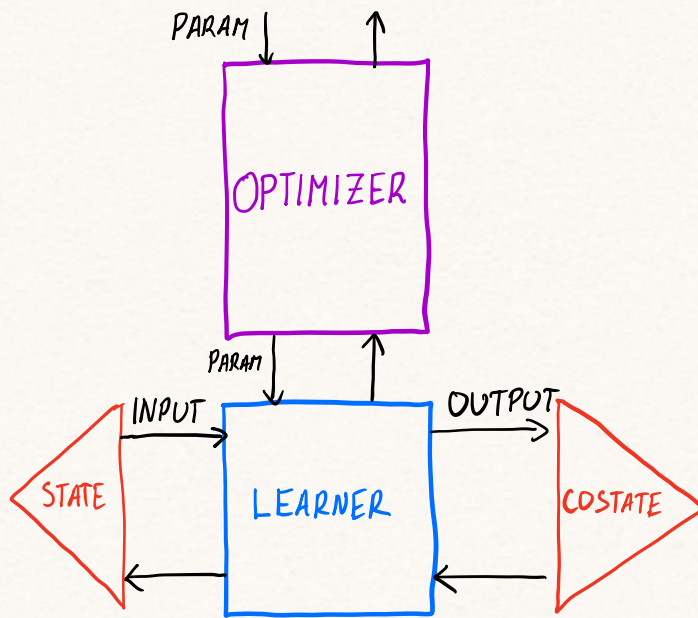
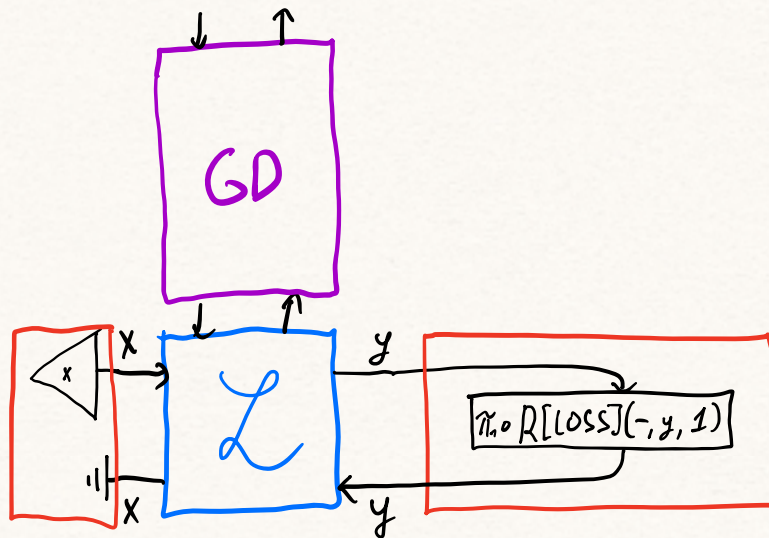$$U_I(p, a, b) := p - \varepsilon \nabla_p E_I(p, a, b)$$

*and*

$$r_I(p, a, b) := f_a\!\left(\nabla_a E_I(p, a, b)\right),$$

*where $E_I(p, a, b) := \sum_j e(I_j(p, a), b_j)$, and $f_a$ is component-wise application of the inverse to $\frac{\partial e}{\partial x}(a_i, -)$ for each $i$.*
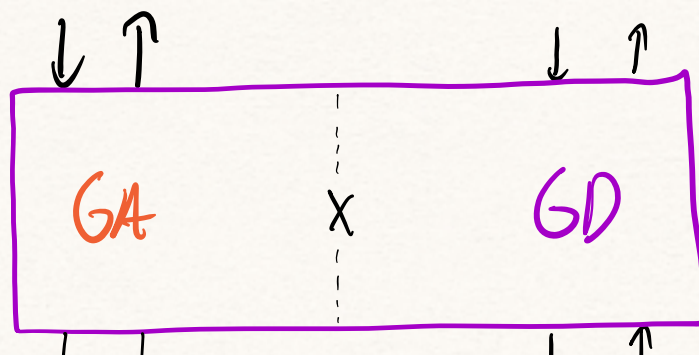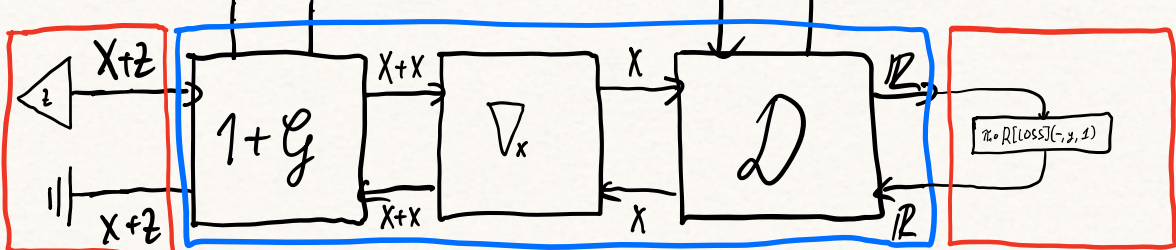
PARAM →
PARAM

**OPTIMIZER**
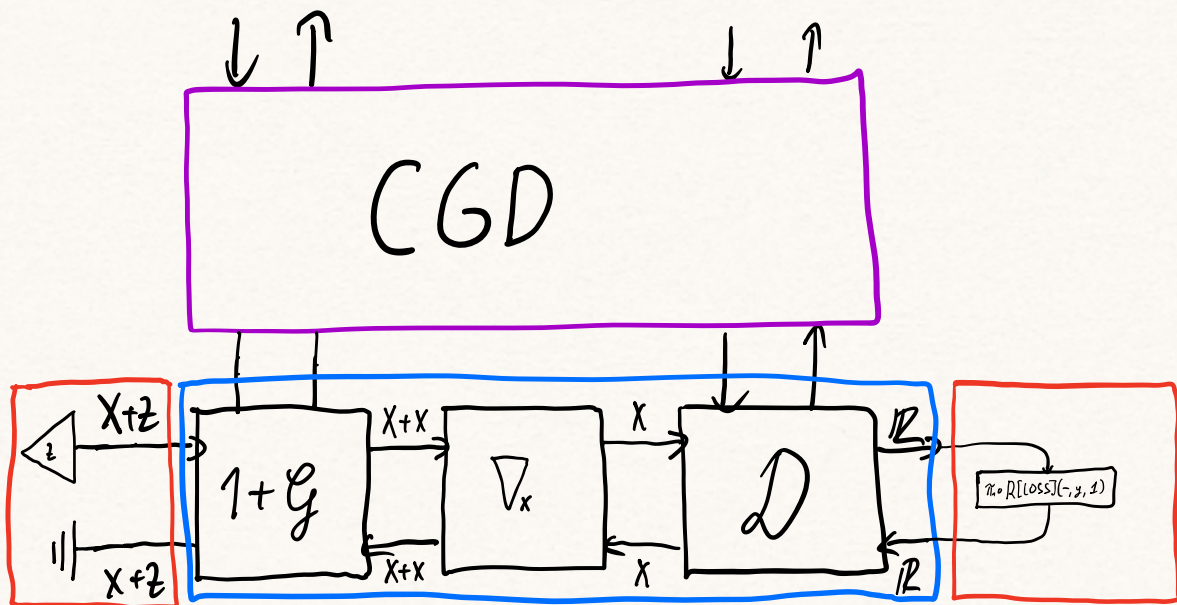
INPUT

**LEARNER**

OUTPUT

STATE

COSTATE

## Learning on Smooth



**GD**

$x$

$\mathcal{L}$

$y$

$\pi_{\circ} R[\text{LOSS}](-, y, 1)$

$x$

$y$

# GENERATIVE ADVERSARIAL NETWORKS



**GA**     $x$     **GD**

X+z

x+z

$1+G$

$\nabla_x$

$\mathcal{D}$

$X+X$

$\bar{X}+X$

X

IR

IR

$\pi_0 R[\text{LOSS}](\cdot, y, 1)$

# CGD



CGD

X+z

x+z

$1+G$

$\nabla_x$

$\mathcal{D}$

$X+X$

$\bar{X}+X$

X

X

IR

IR

$\pi_0 R[\text{LOSS}](\cdot, y, 1)$

# AUTOENCODERS

$\dim(Z) \ll \dim(X)$



P

Q

GD   X   GD

X

X

$e$

$d$

z

z
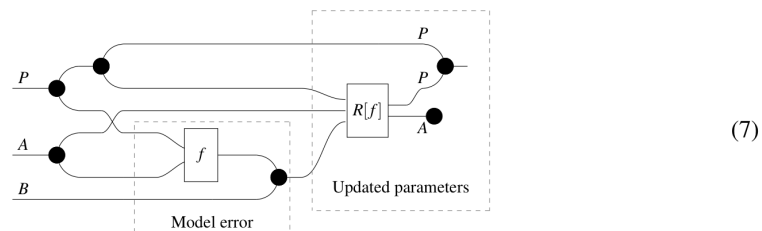
X

X

# Learning on BoolCirc

## 4 Reverse Derivative Ascent

### 4.1 Reverse Derivative Ascent Algorithm

We now introduce our machine learning algorithm, *reverse derivative ascent*. The definition refers to the category **BoolCirc**, as boolean circuits are our motivating example. However, our formulation makes sense in any reverse differential category.

We proceed in two parts: the inner 'step' of the algorithm, which we call `rdaStep`, and the outer 'iteration' of `rdaStep`, which is `rda`.

**Definition 20.** Let $f : p + a \to b$ be a boolean circuit in **BoolCirc**, thus computing a parametrised boolean function with $p$ parameters. We define $\text{rdaStep}_f$ as
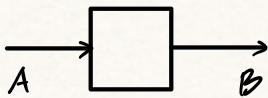


$$(7)$$

# TODO/FUTURE WORK

- DEPENDENT TYPES

- META-LEARNING

- AUTOMATA LEARNING

- OPEN GAMES

- LEARNING ITERATION/REPEATED GAMES

- OPEN DYNAMICAL SYSTEMS?

- PROBABILITY DISTRIBUTIONS?

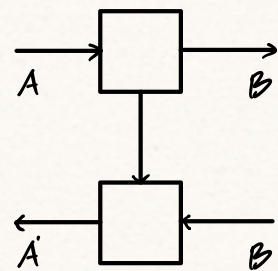- Para AND Optic ARE PRETTY GENERAL, WHAT ELSE CAN WE MODEL?

References:
- BACKPROP AS FUNCTOR (FONG, SPIVAK, TUYERAS)
- OPEN GAMES (GHANI, HEDGES,...)
- REVERSE DERIVATIVE CATEGORIES (COCKETT, CRUTTWELL, GALLAGHER...)
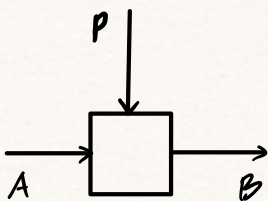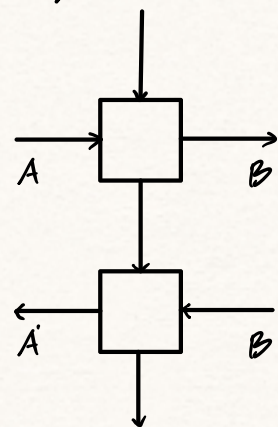- REVERSE DERIVATIVE ASCENT (ZANASI, WILSON)
- DIOPTICS (DARLYMPLE)

$c$



$A \longrightarrow \boxed{\phantom{x}} \longrightarrow B$

$Optic(c)$



$Para(c)$



$Para(Optic(c))$

# EXTRA: OPEN GAMES, COSTRATEGIES

DEF. (CONTEXT)

$$\bar{\mathcal{C}}(A,B) = \int^M \mathcal{C}(I, A \otimes M) \times \mathcal{C}(B \otimes M, I)$$



$$s : (\Sigma \to \Sigma') \leadsto P(\Sigma)$$