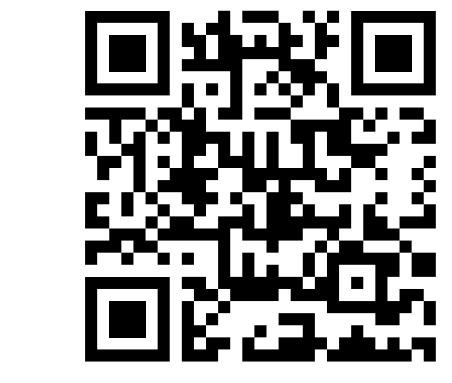


# Turn Logical Formulas into Loss Functions

TL;DR: Yivan Zhang<sup>1,2</sup> Masashi Sugiyama<sup>2,1</sup>

<sup>1</sup>The University of Tokyo <sup>2</sup>RIKEN AIP

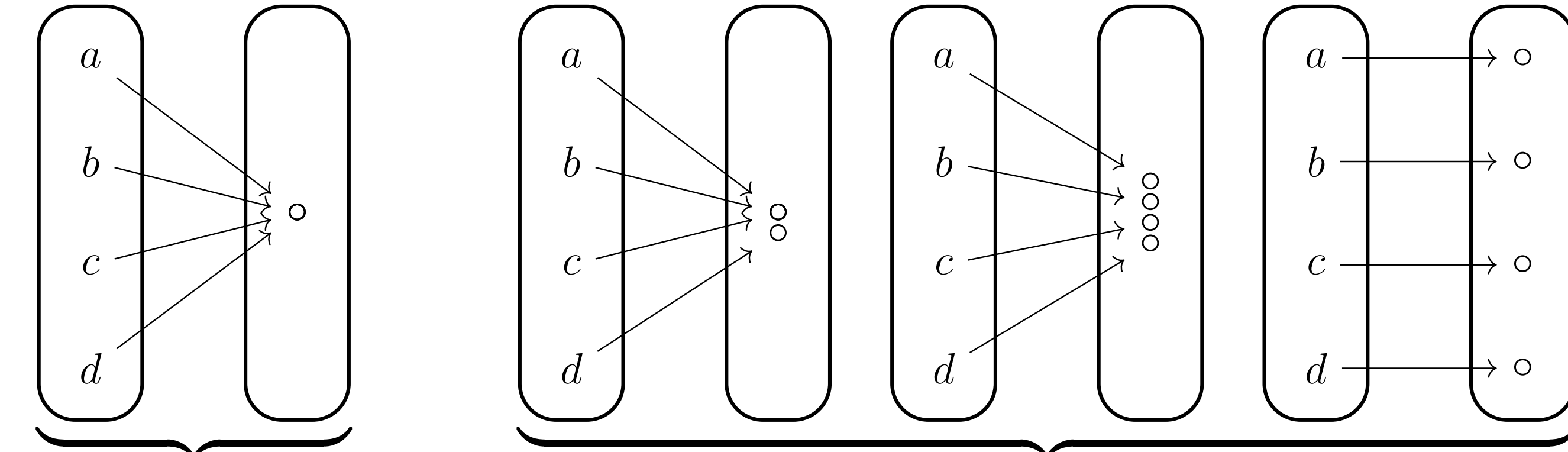


arXiv:2305.11512  
 https://yivan.xyz  
 yivan.xyz@gmail.com



## Motivation: measuring properties of functions

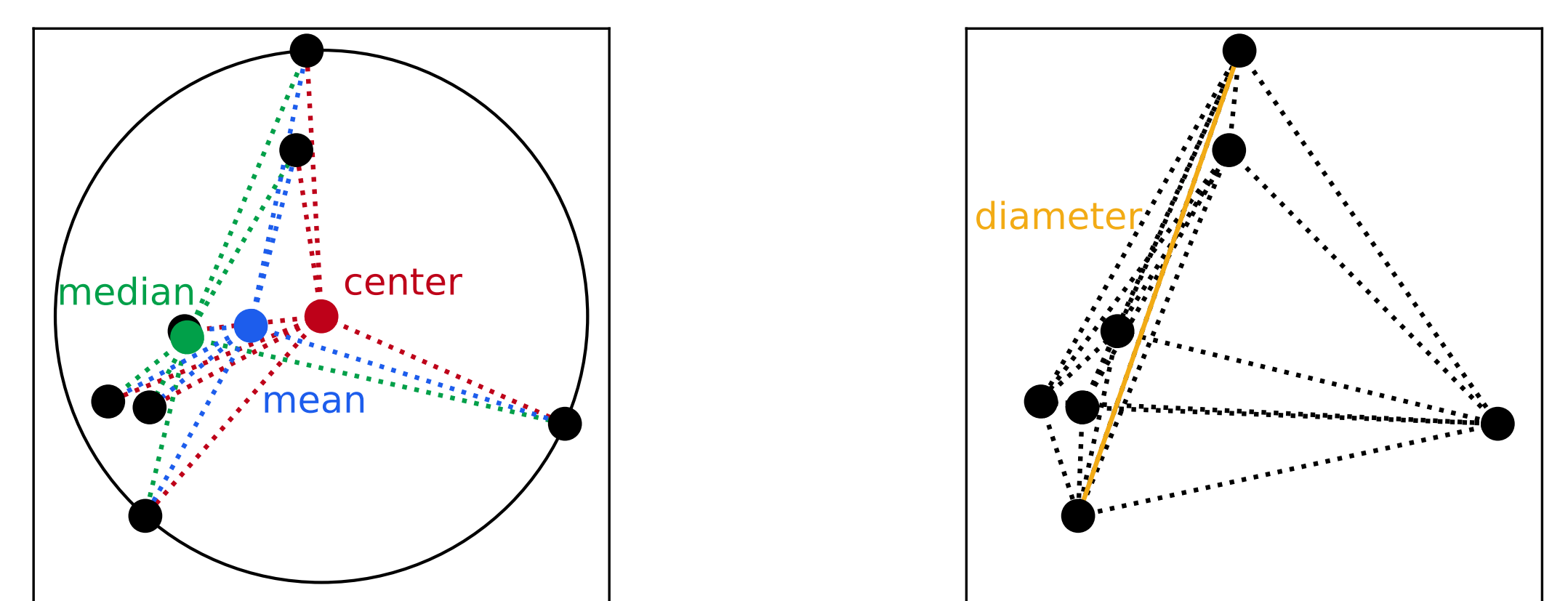
A *constant function* is a function that maps all inputs to the same output, but **“how constant” is a non-constant function?**



- ✓ a constant function
- ✗ non-constant functions
- Intuitively, we can measure how the outputs are distributed over the output space (**constancy** = 0 **deviation**).
- If we have a real-valued and preferably differentiable metric for the degree of constancy, we can *measure* the constancy of a function and *optimize* it using gradient descent.
- How can we measure properties of functions?**

## Idea: quantifying logical predicates

- We can derive metrics from the definition of properties.**
- The constancy of a function  $f : A \rightarrow B$  can be defined by (binary-valued) **logical predicates**:
  - (a)  $\exists b \in B. \forall a \in A. f(a) = b$  or (b)  $\forall a, a' \in A. f(a) = f(a')$
- From these definitions, we can derive their corresponding (real-valued) **quantitative metrics**:
  - (a)  $\inf_{b \in B} \text{agg}_{a \in A} d(f(a), b)$  or (b)  $\text{agg}_{a, a' \in A} d(f(a), f(a'))$



(a) how far the outputs are from a central point (b) how far the outputs are from each other

We can use them as *learning objectives* or *evaluation criteria*. They differ in terms of *computation cost* and *differentiability*.

## Background: logic and metric in machine learning

- Even in the *supervised learning* setting, we can view the **total cost**  $L(f, g) := \sum_{x \in X} \ell(f(x), g(x))$  as a measure of the **function equality** ( $f = g$ ) :=  $\forall x \in X. f(x) = g(x)$  between a learning model  $f : X \rightarrow Y$  and the ground-truth  $g : X \rightarrow Y$ .
- Can we extend this parallel to *representation learning*?

## Problem: disentangled representation learning

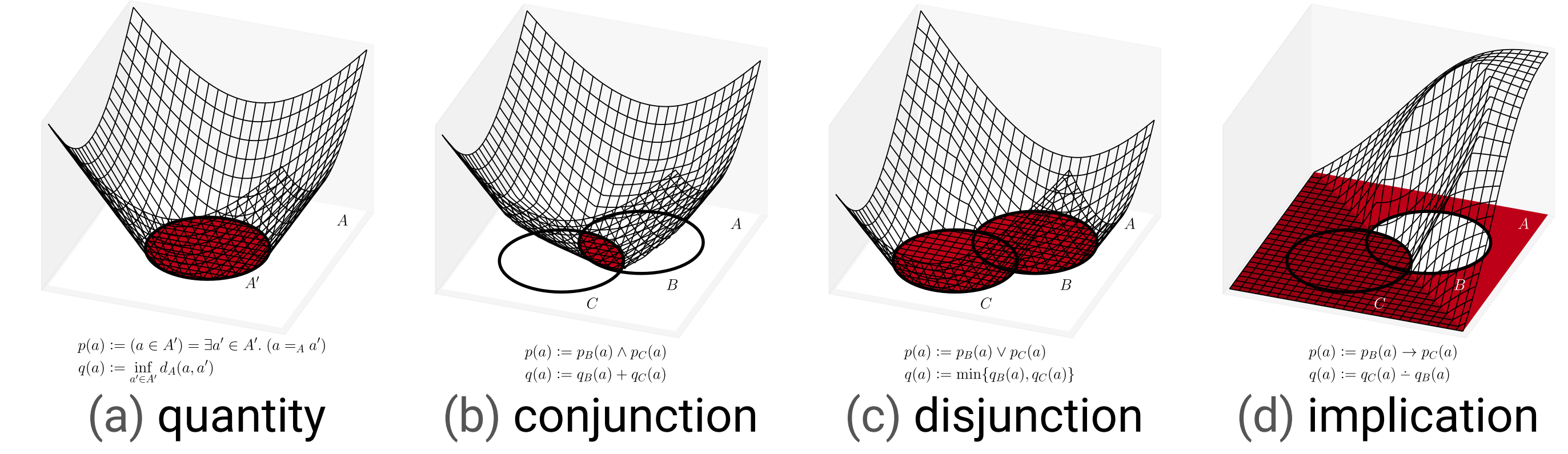
- No unified logical definition, many evaluation metrics [Carbonneau et al., 2022, Zhang and Sugiyama, 2023]
- Unclear what properties these metrics quantify
- Usually non-differentiable and computationally inefficient
- Unproven if a learning method can truly optimize a metric

## Enrichment: from logic to metric

Like logic, we can construct metrics *compositionally!*

Logic		Metric	
truth values	$\{\top, \perp\}$	real values	$[0, \infty]$
predicate equality	$A \xrightarrow{p} \{\top, \perp\}$	quantity	$A \xrightarrow{q} [0, \infty]$
conjunction	$(a = a')$	strict premetric	$d(a, a')$
disjunction	$\wedge$	addition	$+$
implication	$\vee$	minimum	$\min$
universal quantifier	$\rightarrow$	subtraction*	$-$
existential quantifier	$\forall$	aggregator**	$\nabla$
	$\exists$	infimum	$\inf$

\* truncated subtraction:  $b \dot{-} a := \max\{b - a, 0\}$   
 \*\* e.g., maximum, sum, mean, and mean square



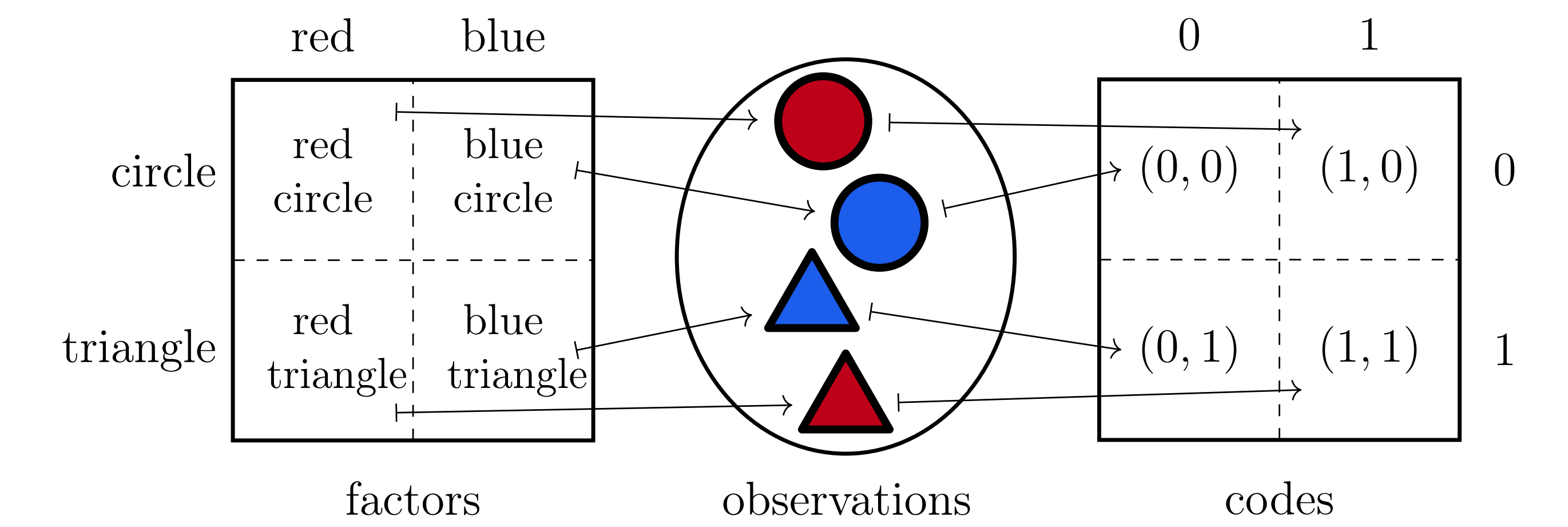
## Theory: (sub)homomorphisms from metric to logic

- Zero predicate:**  $\zeta : [0, \infty] \rightarrow \{\top, \perp\} := x \mapsto (x = 0)$
- Homomorphisms from metric to logic via  $\zeta$ :**

$$\begin{array}{ccc}
 A & \xrightarrow{\text{id}_A} & A \\
 q \downarrow & & \downarrow p \\
 [0, \infty] & \xrightarrow{\zeta} & \{\top, \perp\}
 \end{array}$$
 from quantity  $q$  to predicate  $p$
- Subhomomorphisms:** replace equality = by implication  $\rightarrow$ . **Truncated subtraction** is subhomomorphic to **implication**. No **continuous operation** is homomorphic to **implication**.
- Main theorem:** *If the components are (sub)homomorphic, so is the compound:*  $q(a) = 0$  equals (implies)  $p(a) = \top$ .
- Benefits:** (1) **no failure modes**; (2) **no hyperparameters**; (3) **no stochastic components**; (4) some are **differentiable**.

## Logical definitions of disentangled representations

### Modularity: reconstruct the product structure



$$Y_{\text{color}} \times Y_{\text{shape}} \xrightarrow{g} X \xrightarrow{f} Z_{\text{color}} \times Z_{\text{shape}}$$

$$m := f \circ g = m_{\text{color}} \times m_{\text{shape}}$$

**Example:**

$  \begin{cases}  (\text{red}, \text{circle}) \mapsto (2, 4) \\  (\text{red}, \text{triangle}) \mapsto (2, 6) \\  (\text{blue}, \text{circle}) \mapsto (6, 4) \\  (\text{blue}, \text{triangle}) \mapsto (6, 6)  \end{cases}  $	$  = \underbrace{\begin{cases} \text{red} \mapsto 2 \\ \text{blue} \mapsto 6 \end{cases}}_{m_{\text{color}}} \times \underbrace{\begin{cases} \text{circle} \mapsto 4 \\ \text{triangle} \mapsto 6 \end{cases}}_{m_{\text{shape}}}  $	$  \begin{cases}  (\text{red}, \text{circle}) \mapsto (1, 2) \\  (\text{red}, \text{triangle}) \mapsto (3, 4) \\  (\text{blue}, \text{circle}) \mapsto (5, 6) \\  (\text{blue}, \text{triangle}) \mapsto (7, 8)  \end{cases}  $
---	---	---

✓ a product function      ✗ not a product function

**Definition:**  
 $p_{\text{product}}(m) := \exists m_{1,1} : Y_1 \rightarrow Z_1. \exists m_{2,2} : Y_2 \rightarrow Z_2. m = m_{1,1} \times m_{2,2}$

## Quantitative metrics of disentangled representations

**Metric:**  
 $q_{\text{product}}(m) := \inf_{m_{1,1} \in [Y_1, Z_1]} \inf_{m_{2,2} \in [Y_2, Z_2]} d_{[Y, Z]}(m, m_{1,1} \times m_{2,2})$

**Instantiations:**

- $\text{mean}_{y_1 \in Y_1} \text{var}_{y_2 \in Y_2} m_1(y_1, y_2) + \text{mean}_{y_2 \in Y_2} \text{var}_{y_1 \in Y_1} m_2(y_1, y_2)$
- $\text{max}_{y_1 \in Y_1} \text{diam}_{y_2 \in Y_2} m_1(y_1, y_2) + \text{max}_{y_2 \in Y_2} \text{diam}_{y_1 \in Y_1} m_2(y_1, y_2)$

**Results:**

- We have derived fine-grained, efficient, and differentiable quantitative metrics for disentangled representations.
- We can quantify any logically defined properties!**

	Modularity					Informativeness					Existing metrics						
	Product approx.	Constancy	Retraction approx.	Contraction	Pair	Info.	Regressor	Info.	Info.	Info.	Info.	Info.	Info.	Info.			
entanglement rotation	0.44	0.75	0.96	0.19	0.82	0.76	0.96	0.99	0.44	0.78	0.89	0.83	0.18	0.28	0.28	1.00	
duplicate	0.24	0.43	0.67	0.06	0.56	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	0.59	1.00
complement	0.12	0.28	0.55	0.01	0.42	1.00	1.00	1.00	1.00	1.00	1.00	1.00	0.00	1.00	1.00	0.63	1.00
misalignment	0.22	0.44	0.74	0.05	0.58	1.00	1.00	1.00	1.00	1.00	1.00	1.00	0.00	1.00	1.00	1.00	1.00
redundancy	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	0.33	1.00	1.00	0.93	1.00
contraction	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	0.18	0.49	1.00	1.00	1.00	1.00	1.00	1.00	1.00
nonlinear	1.00	1.00	1.00	1.00	1.00	0.79	0.93	0.99	0.65	0.95	1.00	1.00	0.88	1.00	1.00	1.00	1.00
constant	1.00	1.00	1.00	1.00	1.00	0.42	0.76	0.90	0.18	0.48	0.33	0.33	0.00	0.00	0.00	0.00	0.00
random	0.22	0.48	0.78	0.05	0.61	0.42	0.76	0.90	0.22	0.83	0.34	0.33	0.00	0.00	0.00	0.00	0.04

**References**  
 Marc-André Carbonneau, Julian Zaidi, Jonathan Boilard, and Ghyslain Gagnon. Measuring disentanglement: A review of metrics. *IEEE Transactions on Neural Networks and Learning Systems*, 2022.  
 Yivan Zhang and Masashi Sugiyama. A category-theoretical meta-analysis of definitions of disentanglement. In *ICML*, 2023.