



# Mathematical Foundation of Interpretable Equivariant Surrogate Models

Jacopo Joy Colombini<sup>1</sup>(✉), Filippo Bonchi<sup>2</sup>, Francesco Giannini<sup>1</sup>,  
Fosca Giannotti<sup>1</sup>, Roberto Pellungrini<sup>1</sup>, and Patrizio Frosini<sup>2</sup>

<sup>1</sup> Scuola Normale Superiore, Pisa, Italy  
{Filippo.Bonchi, Patrizio.Frosini}@sns.it

<sup>2</sup> Università di Pisa, Pisa, Italy  
{JacopoJoy.Colombini, Francesco.Giannini, Fosca.Giannotti,  
Roberto.Pellungrini}@unipi.it

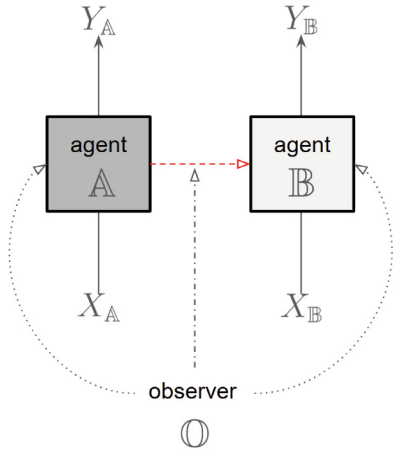
**Abstract.** This paper introduces a rigorous mathematical framework for neural network explainability, and more broadly for the explainability of equivariant operators called Group Equivariant Operators (GEOs), based on Group Equivariant Non-Expansive Operators (GENEOs) transformations. The central concept involves quantifying the distance between GEOs by measuring the non-commutativity of specific diagrams. Additionally, the paper proposes a definition of interpretability of GEOs according to a complexity measure that can be defined according to each user’s preferences. Moreover, we explore the formal properties of this framework and show how it can be applied in classical machine learning scenarios, like image classification with convolutional neural networks.

**Keywords:** Mathematical Foundation of XAI · XAI metrics · Equivariant Neural Networks

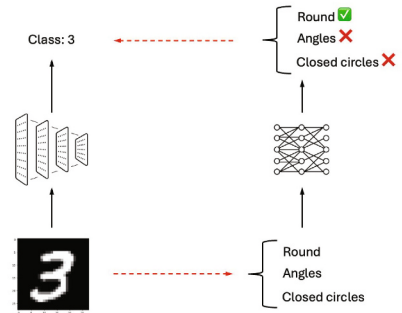
## 1 Introduction

What is an “*explanation*”? An explanation can be seen as a combination of elementary blocks, much like a sentence is formed by words, a formula by symbols, or a proof by axioms and lemmas. The key question is when such a combination effectively explains a phenomenon. Notably, the quality of an explanation is observer-dependent—what is clear to a scientist may be incomprehensible to a philosopher or a child. In our approach, an explanation of a phenomenon  $P$  is convenient for an observer  $\mathbb{O}$  if (i)  $\mathbb{O}$  finds it comfortable, meaning the building blocks are easy to manipulate, and (ii) it is convincing, meaning  $\mathbb{O}$  perceives  $P$  and the explanation as sufficiently close. We contextualize this perspective by assuming that the phenomenon is an AI agent, viewed as an operator, thus saying that the action of an agent  $\mathbb{A}$  is explained by another agent  $\mathbb{B}$  from the perspective of an observer  $\mathbb{O}$  if:

1.  $\mathbb{O}$  perceives  $\mathbb{B}$  as close to  $\mathbb{A}$ ;
2.  $\mathbb{O}$  perceives  $\mathbb{B}$  as less complex than  $\mathbb{A}$ .



(a) An agent  $\mathbb{B}$  explains an agent  $\mathbb{A}$  for an observer  $\mathbb{O}$ :  $\mathbb{O}$  perceives  $\mathbb{B}$  as close to  $\mathbb{A}$ , and  $\mathbb{B}$  as less complex than  $\mathbb{A}$



(b) A practical example: an image classification task can be surrogated by a model based on textual concepts

**Fig. 1.** Representation of interpretable surrogate models exemplified on MNIST.

This is represented in Fig. 1 where we show this concept with an example. Notwithstanding the fact that observer  $\mathbb{O}$  has the right to choose subjective criteria to measure how good  $\mathbb{B}$  is to approximate  $\mathbb{A}$  and their complexities, this paper introduces a mathematical framework for these measurements.

The growing use of complex neural networks in critical applications demands both high performance and transparency in decision-making. While AI interpretability and explainability have advanced, a rigorous mathematical framework for defining and comparing explanations is still lacking [33]. Recent efforts to formalize explanations [21] and interpretable models [29] do not provide practical guidelines for designing or training explainable models, nor do they incorporate the notion of an observer within the theory. Moreover, researchers emphasize the importance of Group Equivariant Operators (GEOs) in machine learning [3, 14, 19, 41], as they integrate prior knowledge and enhance neural network design control [5]. While standard neural networks are universal approximators [17], this typically increases complexity. However, no existing XAI technique addresses explaining an equivariant model using another equivariant model.

This paper addresses this gap by introducing a framework for learning interpretable surrogate models of a black-box and defining a measure of interpretability based on an observer’s subjective preferences. Given the importance of equivariant operators, our XAI framework is built on the theory of GEOs and Group Equivariant Non-Expansive Operators (GENEOs) [18]. GEOs, a broader class than standard neural networks, are well-suited for processing data with inherent symmetries. Indeed, equivariant networks, such as convolutional [20] and

graph neural networks [36], have proven effective across different tasks [35]. Using GENE-based transformations, we develop a theory for learning surrogate models of a given GEO by minimizing algebraic diagram commutation errors. The learned surrogate model can either perform a task or approximate a black-box model’s predictions while optimizing interpretability based on an observer’s perception of complexity, while allowing different observers to have distinct interpretability preferences for the same model architecture.

**Contributions.** Our contributions can be summarized as follows.

- Introduction of a mathematical framework to define interpretable surrogate models, where interpretability depends on a specific observer.
- Definition of a distance between GEOs using diagram non-commutativity, providing a quantitative method for model comparison and training.
- Formal definition of GEOs’ complexity to assess model interpretability.
- We show empirically that these metrics enable training of more interpretable models, usable for direct task-solving or as surrogates for black-box models.

The paper is organized as follows. Section 2 recalls basics from different Mathematics areas that we use to define our metrics in Sect. 3. Section 4 shows how these metrics are used in practice to define a learning problem for an interpretable surrogate model. We show how the proposed framework can be used via an experimental evaluation in Sect. 5. Finally, Sect. 6 comments on related work and Sect. 7 draws conclusions and remarks on future work. The Appendix contains additional material and all proofs.

## 2 Mathematical Preliminaries

The framework proposed in this paper is founded on mathematical structures studied in various fields, such as geometry and category theory. Metric spaces and groups are used to define GE(NE)Os, while categories to compose them.

### 2.1 Perception Spaces and GE(NE)Os

Recall that a *pseudo-metric space* is a pair  $(X, d)$  where  $X$  is a set and  $d: X \times X \rightarrow [0, \infty]$  is a pseudo-metric, namely a function such that, for all  $x, y, z \in X$ ,

$$(R) d(x, x) = 0, \quad (S) d(x, y) = d(y, x), \quad (T) d(x, z) \leq d(x, y) + d(y, z).$$

A *metric*  $d$  is a pseudo-metric that additionally satisfies  $d(x, y) = 0 \implies x = y$ .  $d: X \times X \rightarrow [0, \infty]$  is a *hemi-metric* if it only satisfies (R) and (T). We use the informal term *distance* to refer to either metrics, pseudo-metrics or hemi-metrics.

A *group*  $\mathbf{G} = (G, \circ, \text{id}_G)$  consists of a set  $G$ , an associative operation  $\circ: G \times G \rightarrow G$  having a unit element  $\text{id}_G \in G$  such that, for all  $g \in G$ , there exists  $g^{-1} \in G$  satisfying  $g \circ g^{-1} = g^{-1} \circ g = \text{id}_G$ . A *group homomorphism*  $T: (G, \circ_G, \text{id}_G) \rightarrow$

$(K, \circ_K, \text{id}_K)$  is a function  $T: G \rightarrow K$  such that, for all  $g_1, g_2 \in G$ ,  $T(g_1 \circ_G g_2) = T(g_1) \circ_K T(g_2)$ . Given a group  $(G, \circ, \text{id}_G)$  and a set  $X$ , a *group left action* is a function  $*$ :  $G \times X \rightarrow X$  such that, for all  $x \in X$  and  $g_1, g_2 \in G$ ,

$$\text{id}_G * x = x \quad \text{and} \quad (g_1 \circ g_2) * x = g_1 * (g_2 * x).$$

With these ingredients, we can now illustrate the notions of perception space GEO and GENE0. We refer the interested reader to [6, 18] and [1, 8, 10] for a more extensive description of GENE0s and their applications.

**Definition 1.** An (extended) perception space  $(X, d_X, \mathbf{G}, *)$ , shortly  $(X, \mathbf{G})$ , consists of a pseudo-metric space  $(X, d_X)$ , a group  $\mathbf{G}$ , a left group action  $*$ :  $G \times X \rightarrow X$  such that, for all  $x_1, x_2 \in X$  and every  $g \in G$ ,

$$d_X(g * x_1, g * x_2) = d_X(x_1, x_2).$$

*Example 1.*  $(X, \mathbf{G})$ , with  $\mathbf{G}$  the group of rotations of  $0^\circ, 90^\circ, 180^\circ, 270^\circ$ , and  $X$  a set of images closed under the actions of  $\mathbf{G}$ , is a perception space.

Notice that in any perception space, one can define a pseudo-metric over the group  $\mathbf{G}$  by fixing  $d_G(g_1, g_2) := \sup_{x \in X} d_X(g_1 * x, g_2 * x)$  for any  $g_1, g_2 \in G$ . With this definition, one can easily show that  $\mathbf{G}$  is a topological group and that the action  $*$  is continuous (see Proposition 2 in Appendix).

**Definition 2.** Let  $(X, G), (Y, K)$  be two (extended) perception spaces,  $f: X \rightarrow Y$  and  $t: G \rightarrow K$  a group homomorphism. We say that  $(f, t)$  is an (extended) group equivariant operator (GEO) if  $g(g * x) = t(g) * f(x)$  for every  $x \in X, g \in G$ .  $(f, t)$  is said an (extended) group equivariant non-expansive operator (GENEO) in case it is a GEO and it is also non-expansive, i.e.,

1.  $d_Y(f(x_1), f(x_2)) \leq d_X(x_1, x_2)$  for every  $x_1, x_2 \in X$ ,
2.  $d_K(t(g_1), t(g_2)) \leq d_G(g_1, g_2)$  for every  $g_1, g_2 \in G$ .

The previous extended definitions generalize original perception pairs, GEOs, and GENE0s beyond data represented as functions. We simply refer to them as perception space, GEO, and GENE0. With slight abuse of notation, we use  $d_{\text{dt}}$  for the metric  $d_X$  on the set of data, and  $d_{\text{gr}}$  for the metric  $d_G$  on the group  $G$ , relying on context to specify the perception space  $(X, G)$  under consideration.

*Example 2 (Neural Networks as GEOs).* Neural networks are a special case of GEOs, with different architectures equivariant to specific groups. Convolutional Neural Networks (CNNs) are equivariant to translations, while Graph Neural Networks (GNNs) respect graph permutations. Although standard Multi-Layer Perceptrons are not typically equivariant, they can be viewed as GEOs on the trivial group  $\mathbf{1}$ , containing only the neutral element.

*Example 3.* Let  $X_\alpha$  be the set of all subsets  $\mathbb{R}^3$  and the group  $\mathbf{G}_\alpha$  the group of all translations in  $\mathbb{R}^3$ , and let  $\tau_{(x,y,z)}$  represent the translations by  $(x, y, z)$ . Similarly define  $X_\beta$  and  $\mathbf{G}_\beta$  in  $\mathbb{R}^2$  with  $\tau_{(x,y)}$  translating by  $(x, y)$ . A GENE0  $(f, t)$  can be defined where  $f(x)$  gives the shadow (orthogonal projection) of  $x$  in  $X_\beta$  and the homomorphism  $t: \mathbf{G}_\alpha \rightarrow \mathbf{G}_\beta$  is given by  $t(\tau_{(x,y,z)}) = \tau_{(x,y)}$  for projections onto the  $xy$ -plane. Similarly, defining  $t(\tau_{(x,y,z)}) = \tau_{(y,z)}$  gives a GENE0 for projections onto the  $yz$ -plane.

### 2.2 A Categorical Algebra of GEOs

We introduce a simple language to specify combinations of GEOs. Our proposal relies on the algebra of monoidal categories (CD-categories [12]) that enjoy an intuitive –but formal– graphical representation by means of string diagrams [38].

*Syntax.* We fix a set  $\mathcal{S}$  of basic sorts and we consider the set  $\mathcal{S}^*$  of words over  $\mathcal{S}$ : we write 1 for the empty word and  $U \otimes V$ , or just  $UV$ , for the concatenation of any two words  $U, V \in \mathcal{S}^*$ . Moreover, we fix a set  $\Gamma$  of operator symbols and two functions  $ar, coar: \Gamma \rightarrow \mathcal{S}^*$ . For an operator symbol  $g \in \Gamma$ ,  $ar(g)$  represents its arity, intuitively the types of its input and  $coar(g)$  its coarity, intuitively its output. The tuple  $(\mathcal{S}, \Gamma, ar, coar)$ , shortly  $\Gamma$ , is what is called in categorical jargon a *monoidal signature*.

We consider terms generated by the following context-free grammar

$$c ::= \begin{array}{c} A_1 \\ \vdots \\ A_n \end{array} \left[ \begin{array}{c} B_1 \\ \vdots \\ B_m \end{array} \right] g \quad \left[ \begin{array}{c} \phantom{A} \\ \vdots \\ \phantom{A} \end{array} \right] \quad A \text{---} A \quad \begin{array}{c} A \\ \diagup \\ B \end{array} \begin{array}{c} B \\ \diagdown \\ A \end{array} \quad \begin{array}{c} A \\ \bullet \\ A \end{array} \quad A \text{---} \bullet \quad \begin{array}{c} c_1 \circ c_2 \\ \mid \\ c_1 \otimes c_2 \end{array}$$

where  $A, B, A_i, B_i$  are sorts in  $\mathcal{S}$  and  $g$  is a symbol in  $\Gamma$  with arity  $A_1 \otimes \dots \otimes A_n$  and coarity  $B_1 \otimes \dots \otimes B_m$ . Terms of our grammar can be thought of as circuits where information flows from left to right: the wires on the left represent the input ports, those on the right the outputs; the labels on the wires specify the types of the ports. The input type of a term is the word in  $\mathcal{S}^*$  obtained by reading from top to bottom the labels on the input ports; Similarly for the output. The circuit  $\begin{array}{c} A_1 \\ \vdots \\ A_n \end{array} \left[ \begin{array}{c} B_1 \\ \vdots \\ B_m \end{array} \right] g$  takes  $n$  inputs of type  $A_1, \dots, A_n$  and produce  $m$  outputs of type  $B_1, \dots, B_m$ ;  $\left[ \begin{array}{c} \phantom{A} \\ \vdots \\ \phantom{A} \end{array} \right]$  is the empty circuit with no inputs and no output;  $A \text{---} A$

is the wire where information of type  $A$  flows from left to right;  $\begin{array}{c} A \\ \diagup \\ B \end{array} \begin{array}{c} B \\ \diagdown \\ A \end{array}$  allows for crossing of wires;  $\begin{array}{c} A \\ \bullet \\ A \end{array}$  receives some information of type  $A$  and emit two copies as outputs;  $A \text{---} \bullet$  receives an information of type  $A$  and discards it. For arbitrary circuits  $c_1$  and  $c_2$ ,  $c_1 \circ c_2$  and  $c_1 \otimes c_2$  represent, respectively their sequential and parallel composition drawn as

$$\begin{array}{c} A_1 \\ \vdots \\ A_n \end{array} \left[ \begin{array}{c} B_1 \\ \vdots \\ B_m \end{array} \right] c_2 \left[ \begin{array}{c} C_1 \\ \vdots \\ C_o \end{array} \right] c_1 \quad \text{and} \quad \begin{array}{c} A_1 \\ \vdots \\ A_n \\ C_1 \\ \vdots \\ C_j \end{array} \left[ \begin{array}{c} B_1 \\ \vdots \\ B_m \\ D_1 \\ \vdots \\ D_k \end{array} \right] \begin{array}{c} c_1 \\ c_2 \end{array}$$

As expected, the sequential composition of  $c_1$  and  $c_2$  is possible only when the outputs of  $c_2$  coincides with the inputs of  $c_1$ .

*Remark 1.* The reader may have noticed that different syntactic terms are rendered equal by the diagrammatic representation. For instance both  $c_1 \circ (c_2 \circ c_3)$  and  $(c_1 \circ c_2) \circ c_3$  are drawn as

$$\begin{array}{c} A_1 \\ \vdots \\ A_n \end{array} \left[ \begin{array}{c} B_1 \\ \vdots \\ B_m \end{array} \right] c_3 \left[ \begin{array}{c} C_1 \\ \vdots \\ C_o \end{array} \right] c_2 \left[ \begin{array}{c} D_1 \\ \vdots \\ D_p \end{array} \right] c_1$$

This is not an issue since the two terms represent the same GEO via the semantics that we illustrate here below, after a minimal background on categories.

*Categories.* Diagrams are arrows of the (strict) CD category freely generated by the monoidal signature  $\Gamma$ . The reader who is *not* an expert in category theory may safely ignore this fact and only know that a *category*  $\mathbf{C}$  consists of (1) a collection of objects denoted by  $Ob(\mathbf{C})$ ; (2) for all objects  $A, B \in Ob(\mathbf{C})$ , a collection of arrows  $f: A \rightarrow B$  with source object  $A$  and target object  $B$ ; (3) for all objects  $A$ , an identity arrow  $id_A: A \rightarrow A$  and (4) for all arrows  $f: A \rightarrow B$  and  $g: B \rightarrow C$ , a composite arrow  $g \circ f: A \rightarrow C$  satisfying

$$f \circ (g \circ h) = (f \circ g) \circ h \quad f \circ id_A = f = id_B \circ f$$

for all  $f: A \rightarrow B$ ,  $g: B \rightarrow C$  and  $h: D \rightarrow E$ .

Three categories will be particularly relevant for our work: the category  $\mathbf{Diag}_\Gamma$  having words in  $\mathcal{S}^*$  as objects and diagrams as arrows, the category  $\mathbf{GEO}$  having perception spaces as objects and GEOs as arrows and the category  $\mathbf{GENEO}$  having perception spaces as objects and GENEOS as arrows.

*Semantics.* As mentioned at the beginning of this section, our diagrammatic language allows one to express combinations of GEOs. Intuitively, the symbols in  $\Gamma$  are basic *building blocks* that can be composed in sequence and in parallel with the aid of some wiring technology. The building blocks have to be thought of as atomic GEOs, while diagrams as composite ones.

To formally provide semantics to diagrams in terms of GEOs, the key ingredient is an *interpretation*  $\mathcal{I}$  of the monoidal signature  $\Gamma$  within the (monoidal) category  $\mathbf{GEO}$ , shortly, a function assigning to each symbol  $g \in \Gamma$  a corresponding GEO. Then, by means of a universal property (or, depending on one's perspective, abstract mumbo jumbo), one obtains a function (actually a functor)  $\llbracket - \rrbracket_{\mathcal{I}}: \mathbf{Diag}_\Gamma \rightarrow \mathbf{GEO}$  assigning to each diagram the denoted GEOs (see Table 5 in the Appendix for a simple inductive definition).

Note that  $\llbracket - \rrbracket_{\mathcal{I}}$  may not be surjective, in the sense that not all GEOs are denoted by some diagrams: we call  $\mathcal{G}_{\mathcal{I}}^\Gamma$  the image of  $\mathbf{Diag}_\Gamma$  through  $\llbracket - \rrbracket_{\mathcal{I}}$ , i.e.,

$$\mathcal{G}_{\mathcal{I}}^\Gamma := \{(f, t) \mid \exists c \in \mathbf{Diag}_\Gamma \text{ s.t. } \llbracket c \rrbracket_{\mathcal{I}} = (f, t)\}.$$

Hereafter, we fix a monoidal signature  $\Gamma$  and an interpretation  $\mathcal{I}$  and we write  $\mathcal{G}_{\mathcal{I}}^\Gamma$  simply as  $\mathcal{G}$ . This represents the universe of GEOs that are interesting for the observer, which we are going to introduce in the next section.

### 3 Observers-Based Approximation and Complexity

This paper aims at developing an applicable mathematical theory of interpretable models, which is based on the following intuition: an agent  $\mathbb{A}$  can be interpreted via another agent  $\mathbb{B}$  from the perspective of an observer  $\mathbb{O}$  if: i)  $\mathbb{O}$  perceives  $\mathbb{B}$  as similar to  $\mathbb{A}$  and ii)  $\mathbb{O}$  perceives  $\mathbb{B}$  as less complex than  $\mathbb{A}$ . This perspective motivates us to build a framework allowing the modeling of distance measures for GEOs (Sect. 3.1) and their degree of complexity (opaqueness/not interpretability, Sect. 3.2), w.r.t. the specification of a certain observer.

**Definition 3.** An observer  $\mathbb{O}$  interested in  $\mathcal{G}$  is a couple  $(\mathbf{T}, \mathcal{C})$  where:

- $\mathbf{T}$  is a category of translations GENEOS, namely a category having as objects  $Ob(\mathbf{T})$  those perception spaces that are sources and targets of GEOs in  $\mathcal{G}$  and as arrows  $Hom(\mathbf{T})$  a selected set of GENEOS.
- $\mathcal{C}$  is a complexity assignment, namely a function  $\mathcal{C}: \Gamma \rightarrow \mathbb{R}^+$ .

The *translation GENEOS* in  $\mathbf{T}$  describe all the possible ways that the observer can “translate” data belonging to one perception space into data belonging to another perception space. Requiring these to be GENEOS, i.e., non-expansive, ensures that such translations performed by the observer cannot enlarge distances between data. For example, the observer may admit only isometries as morphisms in  $\mathbf{T}$ , or the observer may not admit any translation at all, meaning that  $\mathbf{T}$  only contains identities (note that this is the smallest possible  $\mathbf{T}$ ).

The *complexity assignment*  $\mathcal{C}: \Gamma \rightarrow \mathbb{R}^+$  maps any building block  $g$  from  $\Gamma$  into a positive real number, a quantity that represent how *complex* is perceived  $g$  by the observer. Here complexity does not refer to the usual computational complexity but rather to the degree of *stress* that the observer perceives in dealing with  $g$ . Note that such assignment is completely arbitrary and thus, different observers may assign different complexities to the same building block. Any observer can specify what are the types of functions that they deem interpretable and/or more informative, from their perspective, for a given problem.

### 3.1 Surrogate Distance of GEOs

To formalize the notion of a surrogate model for an observer  $\mathbb{O}$ , we introduce a new hemi-metric  $h_{\mathbb{O}}$ , which we call the *surrogate distance* of a GEO for another GEO. To proceed, it is fundamental the notion of *crossed translation pair*.

**Definition 4.** Let  $(f_{\alpha}, t_{\alpha}): (X_{\alpha}, G_{\alpha}) \rightarrow (Y_{\alpha}, K_{\alpha})$  and  $(f_{\beta}, t_{\beta}): (X_{\beta}, G_{\beta}) \rightarrow (Y_{\beta}, K_{\beta})$  be two GEOs in  $\mathcal{G}$ . A *crossed pair of translation*  $\pi$  from  $(f_{\alpha}, t_{\alpha})$  to  $(f_{\beta}, t_{\beta})$ , written  $\pi: (f_{\alpha}, t_{\alpha}) \rightleftharpoons_{\mathbf{T}} (f_{\beta}, t_{\beta})$ , is a couple  $\left( (l_{\alpha, \beta}, p_{\alpha, \beta}), (m_{\beta, \alpha}, q_{\beta, \alpha}) \right)$  where

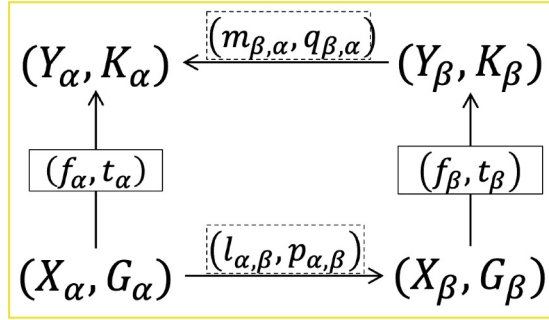
- $(l_{\alpha, \beta}, p_{\alpha, \beta}): (X_{\alpha}, G_{\alpha}) \rightarrow (X_{\beta}, G_{\beta})$  is a GENEOS in  $\mathbf{T}$ ,
- $(m_{\beta, \alpha}, q_{\beta, \alpha}): (Y_{\beta}, K_{\beta}) \rightarrow (Y_{\alpha}, K_{\alpha})$  is a GENEOS in  $\mathbf{T}$ .

Figure 2 provides an intuitive visualization of a crossed pair of translation GENEOS. Note that the two GENEOS have opposite directions.

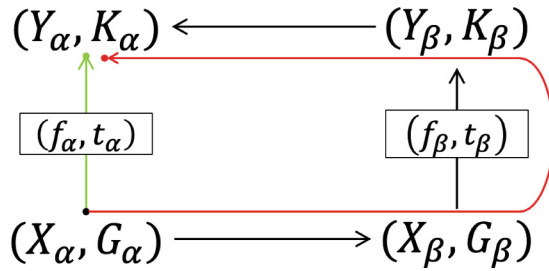
Next, we define the cost of a crossed translation pair.

**Definition 5.** Let  $\pi = \left( (l_{\alpha, \beta}, p_{\alpha, \beta}), (m_{\beta, \alpha}, q_{\beta, \alpha}) \right)$  be a *crossed translation pair* from  $(f_{\alpha}, t_{\alpha}): (X_{\alpha}, G_{\alpha}) \rightarrow (Y_{\alpha}, K_{\alpha})$  to  $(f_{\beta}, t_{\beta}): (X_{\beta}, G_{\beta}) \rightarrow (Y_{\beta}, K_{\beta})$ . The *functional cost* of  $\pi$ , written  $\text{cost}(\pi)$ , is defined as follows.

$$\text{cost}(\pi) := \frac{1}{|X_{\alpha}|} \sum_{x \in X_{\alpha}} d_{\text{dt}} \left( (m_{\beta, \alpha} \circ f_{\beta} \circ l_{\alpha, \beta})(x), f_{\alpha}(x) \right) \quad (1)$$



**Fig. 2.** Example of a crossed translation pair  $\pi: (f_\alpha, t_\alpha) \rightleftharpoons_{\mathbf{T}} (f_\beta, t_\beta)$ . We distinguish by solid and dashed blocks the GEs in  $\mathcal{G}$  from the GENEOS in  $\text{Hom}(\mathbf{T})$ .



**Fig. 3.** The surrogate distance measures how far the diagram is to commute.

*Remark 2.* Note that in Eq. (1),  $|X_\alpha|$  denotes the cardinality of the set  $X_\alpha$ . Whenever such set is infinite the cost is not defined. Although this never happens in practical cases, one can easily generalize (1) to deal with infinite sets by enriching  $X_\alpha$  with a Borel probability measure: see (4) in the Appendix.

Intuitively, the value  $\text{cost}(\pi)$  measures the distance of the two paths in the diagram in Fig. 3. With this, one can easily define a distance between GEs.

**Definition 6.** Let  $(f_\alpha, t_\alpha)$  and  $(f_\beta, t_\beta)$  be two GEs in  $\mathcal{G}$ . The surrogate distance of  $(f_\beta, t_\beta)$  from  $(f_\alpha, t_\alpha)$ , written  $h_\mathbb{O}((f_\alpha, t_\alpha), (f_\beta, t_\beta))$ , is defined as

$$\inf\{\text{cost}(\pi) \mid \pi: (f_\alpha, t_\alpha) \rightleftharpoons_{\mathbf{T}} (f_\beta, t_\beta)\} \tag{2}$$

We emphasize that all considered GENEOS to define crossed pairs of translations must be in  $\mathbf{T}$ . The possibility of choosing  $\mathbf{T}$  in different ways reflects the various approaches an observer can use to judge the similarity between data.

*Example 4.* Consider the smallest possible  $\mathbf{T}$  (that is, no arrows between different spaces and only the identity between equal spaces) representing an observer who cannot translate the data. In this case,  $h_\mathbb{O}((f_\alpha, t_\alpha), (f_\beta, t_\beta)) = \infty$

whenever  $(f_\alpha, t_\alpha)$  and  $(f_\beta, t_\beta)$  act on different perception spaces, since there is no translation pair  $\pi: (f_\alpha, t_\alpha) \Leftarrow_{\mathbf{T}} (f_\beta, t_\beta)$ . Whenever the perception spaces are the same, there is only one translation pair, formed by two identity GENEOS. Thus the surrogate distance of  $(f_\beta, t_\beta)$  from  $(f_\alpha, t_\alpha)$  collapses to the cost of such translation pair, that is,

$$\frac{1}{|X_\alpha|} \sum_{x \in X_\alpha} d_{\text{dt}}(f_\beta(x), f_\alpha(x))$$

Note that whenever  $d_{\text{dt}}$  assigns 0 to equal elements and 1 to different ones, this coincides with the standard notion of *fidelity* [32].

**Theorem 1.** *The function  $h_{\mathbb{O}}$  is a hemi-metric on  $\mathcal{G}$ .*

Notice that while  $h_{\mathbb{O}}$  is a hemi-metric, one can easily get a pseudo-metric by making it symmetric:  $d_{\mathbb{O}} := \max(h_{\mathbb{O}}((f_\alpha, t_\alpha), (f_\beta, t_\beta)), h_{\mathbb{O}}((f_\beta, t_\beta), (f_\alpha, t_\alpha)))$ . We choose to stay with the non-symmetric distance  $h_{\mathbb{O}}$  since it should measure how far the observer  $\mathbb{O}$  perceives the surrogate  $(f_\beta, t_\beta)$  from the GEO to interpret  $(f_\alpha, t_\alpha)$ . We believe that for this kind of measurement, it is more natural to drop symmetry, like, for instance, in the case of fidelity (Example 4).

### 3.2 Measures of Complexity

In Sect. 2.2 we have introduced string diagrams allowing for combining several building blocks taken from a given set of symbols  $\Gamma$  and we have illustrated how the semantics assigns to each diagram a GEO. Here, we establish a way to measure the *comfort* that an observer  $\mathbb{B}$  has in dealing with a certain diagram. We call such measure the *complexity* of a diagram relative to  $\mathbb{O}$ .

To give a complexity to each diagram, we exploit the complexity assignment  $\mathcal{C}: \Gamma \rightarrow \mathbb{R}^+$  of the observer  $\mathbb{O}$  that provides a complexity to each building block.

**Definition 7.** *Let  $c$  be a diagram in  $\text{Diag}_\Gamma$ . The complexity of a diagram  $c$  (relative to the observer  $\mathbb{O}$ ), written  $,c,$  is inductively as follows:*

$$\begin{aligned}
 & , \overbrace{\text{---} g \text{---}}^A := \mathcal{C}(g) & , \overbrace{\text{---} \text{---}}^A := 0 & , c_1 \otimes c_2 := , c_1 + , c_2 \\
 & , \text{---} \bullet \text{---}^A := 0 & , \text{---} \text{---}^A := 0 & , c_1 \circ c_2 := , c_1 + , c_2 \\
 & , \text{---} \text{---}^A := 0 & , \text{---} \bullet := 0
 \end{aligned}$$

Shortly, the complexity of a diagram  $c$  is the sum of all the complexities of the basic blocks occurring in  $c$ .

*Example 5 (Number of Parameters).* The set of basic blocks  $\Gamma$  may contain several generators that depend on one or more parameters whose value is usually learned during the training process. A common way to measure the complexity of a model is simply by counting the number of its parameter. This can be easily

accommodated in our theory by fixing the function  $\mathcal{C}: \Gamma \rightarrow \mathbb{R}^+$  to be the one mapping each generator  $g \in \Gamma$  into its number of parameters. It is thus trivial to see that for all circuit  $c, \mathcal{C}(c)$  is exactly the total number of parameters of  $c$ .

*Example 6 (Number of Nonlinearities).* Let us assume that  $\Gamma$  contains as building blocks the functions computing the linear combinations of  $n$  given inputs, for every  $n \in \mathbb{N}$  and for each tuple of real valued coefficients. Moreover,  $\Gamma$  contains as building blocks some classic activation functions in machine learning, such as the Sigmoid and the ReLu activation function. For instance, in our theory an observer may define the complexity  $\mathcal{C}: \Gamma \rightarrow \mathbb{R}^+$  to assign to each linear function the complexity of 0 and to each nonlinear function the complexity of 1. Then the complexity of each circuit  $c, \mathcal{C}(c)$  is exactly the number of nonlinear functions applied in the circuit, e.g. the number of neurons in a multi-layer perceptrons with ReLu activation functions in the hidden layers and Sigmoid activation function in the output layer.

We notice that we defined the complexity function on syntactic diagrams and not on semantic objects. Indeed, an operator, like e.g. a GEO, can be realized by possibly several different diagrams, however the complexity of the different diagrams should be different. To understand this choice, imagine one has to define the complexity of a function that, given a certain array of integers, returns the array in ascending order. Clearly the complexity of this function should depend on the specific algorithm that is used to produce the output given a certain input, and not on the function itself.

## 4 Learning and Explaining via GE(NE)Os Diagrams

Section 3 introduces the basic definitions that can be operatively used to instantiate our framework. Indeed, Eq. (2) defines a hemi-metric that can be used as a loss function to train a surrogate GEO to approximate another GEO, whereas Definition 7 establishes a way to measure their interpretability in terms of elementary blocks. This section first shows how the learning of surrogate models is defined (Sect. 4.1), and then how we can easily extract explanations from the learned surrogate models (Sect. 4.2). For the following we assume to have fixed an observer  $\mathbb{O} = (\mathbf{T}, \mathcal{C})$  interested in a set of GEOs  $\mathcal{G}$ .

### 4.1 Learning via GENEOS' Diagrams

Given two GEOs  $\alpha, \beta \in \mathcal{G}$ , with  $\alpha = (f_\alpha, t_\alpha) : (X_\alpha, G_\alpha) \rightarrow (Y_\alpha, K_\alpha)$  and  $\beta = (f_\beta, t_\beta) : (X_\beta, G_\beta) \rightarrow (Y_\beta, K_\beta)$ , and the category  $\mathbf{T}$  of translation GENEOS, the hemi-metric  $h_{\mathbb{O}}$  as defined in Eq. (1) expresses the cost of approximating  $\alpha$  with  $\beta$  via the available translation pairs, as illustrated in Fig. 3. In order to apply our framework to the problem of learning interpretable surrogate functions of a certain model on a certain dataset, from now on we assume that  $\alpha$  is given,  $\beta$  is learnable by depending on a set of parameters  $\theta \in \mathbb{R}^n$ , and  $X_{dt}$  denotes the training set collecting the available input data. Therefore,

learning  $f_\beta$  can be cast as the problem of finding the parameters  $\theta$ , such that  $h_{\mathbb{O}}(\alpha, \beta)$  is minimized on  $X_{dt}$ , i.e. that provide the lowest  $cost(\pi)$  amongst the  $\pi = \left( (l_{\alpha, \beta}^\pi, p_{\alpha, \beta}^\pi), (m_{\beta, \alpha}^\pi, q_{\beta, \alpha}^\pi) \right) : \alpha \rightleftharpoons_{\mathbf{T}} \beta$ :

$$\theta^* = \arg \min_{\theta} \left( \inf_{\frac{1}{|X_{dt}|}} \sum_{x \in X_{dt}} d_{dt} \left( m_{\beta, \alpha}^\pi (f_\beta(l_{\alpha, \beta}^\pi(x); \theta)), f_\alpha(x) \right) \right). \quad (3)$$

From our definition, the two perception spaces may be different. However, most frequently when learning surrogate functions, we have  $W_\alpha = W_\beta = W$ , for  $W \in \{X, Y, \mathbf{G}, \mathbf{K}\}$ , and there is only the translation pair  $\pi = ((id_X, id_G)(id_Y, id_K))$ . Thus, Eq. (3) simplifies in  $\arg \min_{\theta} \frac{1}{|X_{dt}|} \sum_{x \in X_{dt}} d_{dt} (f_\beta(x; \theta), f_\alpha(x))$ , which corresponds to the fidelity measure between  $f_\alpha$  and  $f_\beta$ , commonly used in XAI.

*Example 7 (Classifier Explanations).* Consider a classifier  $f_\alpha$  equivariant w.r.t. the groups  $\mathbf{G}_\alpha$  and  $\mathbf{K}_\alpha = \mathbf{1}$ , being  $\mathbf{1}$  the trivial group. As an example, Fig. 4 illustrates two different GEOs  $f_\beta$  and  $f_\gamma$  that can be used to explain  $f_\alpha$ . Notice that if the observer  $\mathbb{O}$  has no access to  $f_\alpha$ , i.e.  $\mathbb{O}$  does not know how  $f_\alpha$  is built (i.e.  $f_\alpha$  is a black-box for  $\mathbb{O}$ ), then  $f_\alpha$  should be an atomic block in  $\Gamma$ . In this case, the observer  $\mathbb{O}$  assigns to  $f_\alpha$  the complexity  $\mathcal{C}(f_\alpha) = \infty$ .

*Example 8 (Supervised Learning).* Wether  $f_\alpha$  denotes the function associating to each training input its label (i.e. the *supervisor*), then  $f_\beta$  and  $f_\gamma$  from Fig. 4 are simply two models trained via supervised learning, and their distance to  $f_\alpha$  is the accuracy (that can be thought of as the fidelity w.r.t. the ground-truth).

$f_\beta$  and  $f_\gamma$  differ in Example 7 only from the fact that  $f_\beta$  is equivariant on the same group  $\mathbf{G}_\alpha$  than  $f_\alpha$ , whereas  $f_\gamma$  might not. In fact, in case  $f_\gamma$  is not equivariant on  $\mathbf{G}_\alpha$  we may prove that  $f_\gamma$  will be surely a non-optimal approximation.

**Proposition 1.** *Let  $\mathbf{T}$ ,  $(f_\alpha, t_\alpha)$ ,  $(f_\beta, t_\beta)$  as in Example 4 and let  $NE$  be the set  $\{(g, x) \in G_\alpha \times X \mid f_\beta(x) \neq f_\beta(g * x)\}$ , i.e., the set containing all those couples falsifying equivariance of  $f_\beta$  w.r.t.  $G_\alpha$ . Then*

$$h_{\mathbb{O}}((f_\alpha, t_\alpha), (f_\beta, t_\beta)) \geq \frac{|NE|}{2 \cdot |G_\alpha|}$$

*Remark 3.* As stated in the introduction, single-hidden-layer neural networks are universal approximators but may require a large number of hidden neurons, increasing complexity. If we cap the model’s complexity, a neural network may not always approximate a given model accurately. Proposition 1 further establishes a fidelity lower bound based on non-equivariant datapoints.

## 4.2 Suitable Surrogate GEOs

We say that a GEO  $(f_\alpha, t_\alpha)$  is *explained* by another GEO  $(f_\beta, t_\beta)$  at the level  $\varepsilon$  for an observer  $\mathbb{O} = (\mathbf{T}, \mathcal{C})$  if:

$$1. h_{\mathbb{O}}((f_\alpha, t_\alpha), (f_\beta, t_\beta)) \leq \varepsilon; \quad 2., (f_\beta, t_\beta) \leq (f_\alpha, t_\alpha).$$

The second condition means that the complexity of the surrogate explaining model  $(f_\beta, t_\beta)$  should be lower than the complexity of the given model  $(f_\alpha, t_\alpha)$ . While not guaranteed, this requisite can be ensured by designing  $f_\beta$  with a suitable strategy. Recall that a model’s complexity is defined by atomic building blocks in  $\Gamma$ , which are combined to form the model. Using the simplest possible blocks helps limit complexity, though their selection depends on the observer’s knowledge and interpretability. Moreover, different studies [6] have shown how a proper domain-informed selection of GE(NE)Os, may strongly decrease the number of parameters necessary to solve a certain task w.r.t. a standard neural networks (as also shown by our experiments cf. Table 1).

*Example 9.* Given a set of GEOs  $(f_i, t_i) \in \Gamma$ , with complexity  $k_i = \mathcal{C}((f_i, t_i))$ , we can define  $f_\beta$  as a linear combination of  $(f_1, t_1), \dots, (f_n, t_n)$ . According to Definition 7, the complexity  $f_\beta$  would be  $k_1 + \dots + k_n$ , plus eventually the complexities of the scalar multiplications.

## 5 Experiments

In order to validate experimentally our theory, we build a classification task on MNIST dataset and rely on our framework to appropriately define an interpretable surrogate model. With our experiments we aim to answer two main research questions: whether personalized complexity measures are able to properly formalize an observer subjectivity, and if knowledge of the domain and of the complexity measured by an observer can lead to ad-hoc surrogate models with a better trade-off between complexity and accuracy. Thus for all the reported results, we assume to have fixed one (or more) given observers.<sup>1</sup>

### 5.1 Data

MNIST contains 70,000 grayscale (values from 0 to 255) images of handwritten digits (0–9), each image being  $28 \times 28$  pixels. We linearly rescale the images so that the values lay in  $[0, 1]$ . The images rescaled belong to  $\{0, \frac{1}{255}, \dots, 1\}^{28 \times 28}$ . We split our dataset into three stratified random disjoint subsets: training, validation, and test set, of 60%, 20% and 20% of images respectively.

### 5.2 Models

As opaque model, we employ a standard CNN, with the Tiny-Vgg architecture, that is composed by two convolutional layers as tail and a linear classifier head. To realize our GEOs surrogate approximation, we use two different architectures. From the MNIST training set, we extract randomly a set of patterns  $p_i$ . These patterns are square cutouts of train images, with height ( $H$ ) and width ( $W$ )

---

<sup>1</sup> Our code is available at <https://github.com/jacopojoy98/GENEO>.

of choice and with a center point chosen with probability proportional to the intensity of the image  $x$ :

$$p_i = x|_{Q_i}, \quad (c_{x_i}, c_{y_i}) \sim x$$

$$Q_i = \left\{c_{x_i} - \frac{W}{2}, \dots, c_{x_i} + \frac{W}{2}\right\} \times \left\{c_{y_i} - \frac{H}{2}, \dots, c_{y_i} + \frac{H}{2}\right\}$$

For each image  $x$  we identify the presence of a pattern  $p_i$  in position  $(i, j)$  with the following function:

$$f(x)_{p_i} : \left\{0, \frac{1}{255}, \dots, 1\right\}^{28 \times 28} \rightarrow \left\{0, \frac{1}{199920}, \dots, 1\right\}^{28 \times 28}$$

$$f_{p_i}(x)_{n,m} = 1 - \frac{\sum_{(i,j) \in Q_i} |x((i, j) + (n, m)) - p_i((i, j))|}{\text{vol } Q_i}$$

The choice of these specific patterns can be motivated by a domain knowledge or by the preferences that an observer can inject through a thoughtful design of their GEOs' building blocks for the classification task.

The first GEO then performs a Image-Wide-Maxpool to create a flat vector with as many entries as are the patterns, and whose  $i^{\text{th}}$  entry indicates the intensity with which the pattern was identified within the image

$$L_i = \max_{n,m}(f_{p_i}(x)_{n,m})$$

These intensities are then linearly combined with an activation function to identify the correct digit

$$\text{OUT}^k = \sigma \left( \sum_j \gamma_j^k L_j + b^k \right)$$

The second GEO instead, after the identification of patterns, selects for each pattern the position with the maximum activation through the Channel-Wise-Max (*CWM*)

$$CWM(f_{p_i}(x))_{n,m} = \begin{cases} s & \text{if } s = \max(f_{p_i}(x)) \\ 0 & \text{otherwise} \end{cases}$$

These matrices of activations are then linearly combined with a downstream nonlinear activation function

$$L_{n,m} = \sigma \left( \left( \sum_i w_i \cdot CWM(f_{p_i}(x))_{n,m} \right) + b_i \right)$$

The entries of this matrix are then linearly combined with a final sigmoidal activation function to produce the output of the model

$$\text{OUT}^k = \sigma \left( \left( \sum_{ij} w_{ji}^k \cdot L_{ji} \right) + b^k \right)$$

**Table 1.** The different models utilized with the relative hyperparameters, chosen on the validation set.

Model	Params	Epochs	LR	Model	Params	Epochs	LR	PATTERNS
CNN	228010	3	$3e-3$	GEO <sub>1</sub>	5010	296	$3e-3$	500
MLP	31810	57	$2e-4$	GEO <sub>1</sub>	3510	148	$7e-3$	350
MLP	15910	57	$1e-4$	GEO <sub>1</sub>	1710	456	$2e-2$	170
MLP	7850	5	$2e-3$	GEO <sub>1</sub>	1510	564	$1e-2$	150
MLP	5575	58	$2e-4$	GEO <sub>1</sub>	1210	496	$2e-2$	120
MLP	3985	58	$2e-4$	GEO <sub>1</sub>	990	198	$5e-2$	98
MLP	3190	9	$2e-3$	GEO <sub>2</sub>	8101	39	$1e-3$	250
				GEO <sub>2</sub>	8051	496	$1e-3$	200
				GEO <sub>2</sub>	8001	483	$1e-3$	150
				GEO <sub>2</sub>	7951	335	$1e-3$	100
				GEO <sub>2</sub>	7901	451	$1e-3$	50

To compare results, we chose a series of simple Multi-Layer Perceptrons, trained directly on the MNIST dataset. In particular, we used MLPs with the following configurations: with no hidden layers, with one hidden layer of dimension 5, 7, 20 and 40. The two models with hidden layers of dimension 5 and 7 are chosen to create MLPs with number of parameters similar to our GEOs. In Table 1 we report the most relevant characteristics of all the models we compare in our experiments.

### 5.3 Experiment Setup

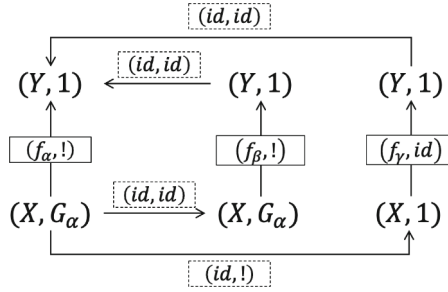
We performed the experiments training all models over the ground truth.

We employed early stopping on the validation set to determine the optimal number of training epochs. The accuracy was then evaluated on a separate test set. We also trained a portion of our models on a rescaled version of MNIST for which every separate group of  $2 \times 2$  points was substituted with the max of the four pixels, effectively reshaping the images to  $14 \times 14$  and allowing us to compare also models which start from different perception spaces.

### 5.4 Results

We first follow our theoretical framework to define the translation diagram of our experimental setup. Indeed, we are in a classical classification scenario, that can be easily represented by the graph in Fig. 4.

We start from the basic perception space  $(X, G_\alpha)$  that is, our image dataset  $X$  and the group of admissible transformations  $G_\alpha$ . Here, we have translations as admissible group actions in  $G_\alpha$  and  $f_\alpha$  is the opaque CNN. Our first GEO model  $f_\beta$  operates on the same perception space  $(X, G_\alpha)$ , as it works on the torus of the images, preserving translations. Therefore, the translation GENEIO is composed given by the couple  $(id_X, id_{G_\alpha})$ . Both the second GEO and the MLP, represented by  $f_\gamma$  instead do not preserve any transformation in the group. Therefore the perception space becomes  $(X, \mathbf{1})$  where  $\mathbf{1}$  denotes the trivial group. Being  $!$  the annihilator homomorphism from any group to the trivial group, the translation GENEIO for this GEO is given by the couple  $(id_X, !)$ . All models have  $(Y, \mathbf{1})$  as their output, since they all work on the space of output classes.



**Fig. 4.** Diagrams of two GEOs explaining a given GEO, where ! represents the annihilator homomorphism from any group to the trivial group 1.

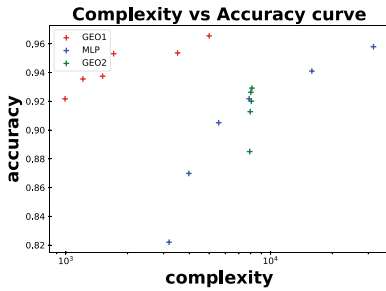
**Table 2.** Models with relative complexities, accuracies and fidelities w.r.t CNN.

Model	$C_1$	$C_2$	Acc	Fid	Model	$C_1$	$C_2$	Acc	Fid
					GEO <sub>1</sub>	5010	510	96.6%	92.5%
					GEO <sub>1</sub>	3510	360	95.4%	91.9%
CNN	228010	37578	97.8%		GEO <sub>1</sub>	1710	180	95.3%	92.4%
MLP	31810	50	96.3%	93.6%	GEO <sub>1</sub>	1510	160	93.7%	90.7%
MLP	15910	30	94.1%	93.5%	GEO <sub>1</sub>	1210	130	93.4%	91.5%
MLP	7850	10	91.8%	90.9%	GEO <sub>1</sub>	990	100	92.2%	89.3%
MLP	5575	17	90.3%	89.6%	GEO <sub>2</sub>	8101	511	92.9%	92.5%
MLP	3985	15	85.4%	86.1%	GEO <sub>2</sub>	8051	411	92.0%	91.8%
MLP	3190	14	85.1%	80.3%	GEO <sub>2</sub>	8001	311	92.6%	91.6%
					GEO <sub>2</sub>	7951	211	91.3%	91.1%
					GEO <sub>2</sub>	7901	111	88.5%	91.4%

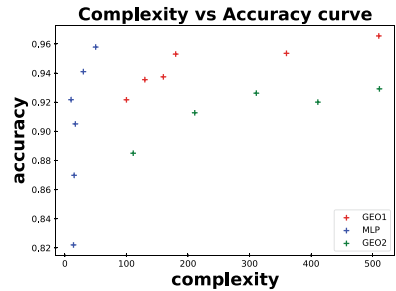
**Table 3.** The output of some of the models trained on a rescaled version of the starting perception space. The hyperparameters have been kept the same as the non rescaled experiments

Model	$C_1$	$C_2$	Acc	Model	$C_1$	$C_2$	Acc
				GEO <sub>1</sub>	5010	510	95.9%
MLP	8290	50	96.3%	GEO <sub>1</sub>	3510	360	95.5%
MLP	1970	10	91.8%	GEO <sub>1</sub>	1710	180	93.6%
MLP	1459	17	90.3%	GEO <sub>1</sub>	990	100	91.1%
MLP	1045	15	86.3%	GEO <sub>2</sub>	2221	511	93.1%
				GEO <sub>2</sub>	2121	411	92.5%
				GEO <sub>2</sub>	2021	111	89.5%

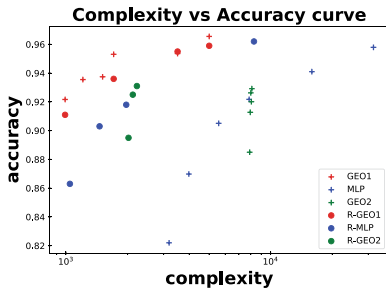
To show how the subjectivity of an observer may influence the results in practice, we measure complexity using two measures: Firstly we assign complexity 1 to each parameter of the model and we sum over all the parameters. Then we assign complexity 1 to all the non-linearities of the model, summing over all the non linearities. We report the performances obtained by the different models in Table 2 and we also compare the results with a different perception space in Table 3 where we present the results for resized images.



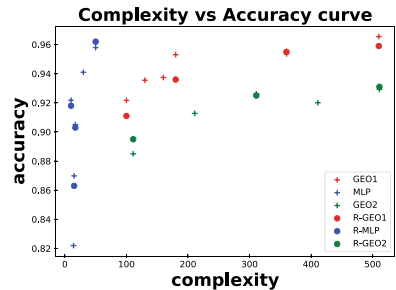
(a) The proposed GEOs can outperform, at similar complexity, model-agnostic MLPs. Notice that the translational equivariant GEO is able to perform much better at minor complexity.



(b) A second observer, who ascribes complexity only to the non-linearities, can have a different complexity vs accuracy curve.



(c) Changing the starting perception space does not affect significantly the performances, whereas the first observer sees GEO<sub>1</sub> with unchanging complexities and GEO<sub>2</sub> and the MLP with much smaller complexity.



(d) The different observer is not affected by any change in the measured complexity.

**Fig. 5.** Accuracy vs complexity comparisons.

The results show that the models built via thoughtful GEOs’ building blocks can approximate quite well the original task, providing models that are less complex for both the measure specified by the observer. The complexity vs accuracy curves representing the experiments are shown in Fig. 5.

## 6 Related Work

Explainable AI has become a fundamental field in AI that covers methodologies designed to provide understandable explanations of the inner workings of a ML model to a human being [27]. Roughly, XAI methods can be categorized into post-hoc methods, i.e., methods aiming to explain another trained opaque ML models, and interpretable-by-design methods, i.e., ML models that

provide explanations to the users inherently, by virtue of their intrinsic transparency [9, 42]. One of the most well-known techniques for post-hoc explanations is to train a surrogate interpretable model to reproduce the same output as an opaque model [15, 24, 28]. In this regard, our paper provides a solid mathematical framework that subsumes both these two paradigms in the same theory.

A key point in XAI is the way the quality of the provided explanations can be measured. For instance, explanations and interpretability can be evaluated qualitatively (user studies) or quantitatively (direct model metrics) [2, 30, 32, 43]. Qualitative measures include user performance, engagement, and explanation clarity [4, 16, 34, 37]. Quantitative measures include explanation completeness [40], fidelity [22], classification accuracy [23], and faithfulness [31]. Complexity measure of explanations is often used for logic-based explainers [13], but it is generally limited to be a count on the number of propositional variables in a formula. While this can easily be accommodated in our framework, up to the author knowledge, no other methods consider complexity measures from the perspective of an observer, offering flexibility in choosing suitable metrics for the task and models.

While there is a large agreement on the needs for XAI models, there are very few works that try to provide a formal mathematical theory of explanations and/or interpretability for ML models. For instance, in [39] the authors propose a new class of “compositionally-interpretable” models, which extend beyond intrinsically interpretable models to include causal models, conceptual space models, and more, by using category theory. [21] proposes a framework based on Category Theory and Institution Theory to define explanations and (explainable) learning agents mathematically. However, these works do not provide a practical measure for the interpretability of the models, completely omit the formalization of an observer, and do not take into account the notion of group equivariant operators. Another seminal work is [25, 26], which provides a more general foundation framework based on properties and desiderata for interpretable ML. However, it does not make any specific mention to a proper mathematical framework.

Finally, our framework is based on the theory of GE(NE)Os, which has been already used to bridge Topological Data Analysis (TDA) and ML. For instance, GENEOS originates from persistent homology with G-invariant non-expansive operators and have been successfully applied for 1D-signal comparisons and image recognition based on topological features [18]. Moreover, GENEOS have been applied to protein pocket detection [6, 8] and graph comparison [7]. While as observed in [8] GENEOS are more inherently interpretable due to a limited dependency on parameters, the theory we present in this paper significantly extend the previous applications, by aiming at the formalization of a more sound XAI theory evaluable quantitatively and based on observers’ preferences.

## 7 Conclusions and Future Work

This work explores the theoretical properties of GE(NE)Os to build a theoretical framework to build surrogate interpretable models, and measure in a rigorous

way the trade-off between complexity and performance. By formally proving the properties of our framework and with the experiments that we provide, we lay the groundwork for future research and opening avenues for practical applications in analyzing and interpreting complex data transformations. Our proposal highlights how it is possible to frame the theory of interpretable models through GE(NE)Os and opens new interesting research directions for Explainable AI. One such direction will be to formally describe existing machine learning models in terms of GE(NE)Os, to study the best interpretable approximations for typical tasks. Moreover, an interesting possible research could be to realize interpretable latent space compression through the use of GE(NE)Os.

**Acknowledgments.** This work has been partially supported by the Partnership Extended PE00000013 - “FAIR - Future Artificial Intelligence Research” - Spoke 1 “Human-centered AI” and ERC-2018-ADG G.A. 834756 “XAI: Science and technology for the eXplanation of AI decision making”. This research was partly funded by the Advanced Research + Invention Agency (ARIA) Safeguarded AI Programme and carried out within the National Centre on HPC, Big Data and Quantum Computing - SPOKE 10 (Quantum Computing) and by the European Union Next-GenerationEU - National Recovery and Resilience Plan (NRRP) M.4 C.2, I.N.1.4 CUP N. I53C22000690001. Bonchi is supported by the Ministero dell’Università e della Ricerca of Italy grant PRIN 2022 PNRR No. P2022HXNSC - RAP (Resource Awareness in Programming). P.F. conducted a portion of his research within the framework of the CNT WiLab National Laboratory and the WiLab-Huawei Joint Innovation Center. His work received partial support from INdAM-GNSAGA, the COST Action CaLISTA, and the HORIZON Research and Innovation Action PANDORA. This work was also funded by the European Union under Grant Agreement no. 101120763 - TANGO. Views and opinions expressed are however those of the author(s) only and do not necessarily reflect those of the European Union or the European Health and Digital Executive Agency (HaDEA). Neither the European Union nor the granting authority can be held responsible for them.

**Disclosure of Interests.** The authors have no competing interests.

## Appendix

**Proposition 2.** *Let  $(X, d_X, \mathbf{G}, *)$  be a perception pair. The followings hold.*

- (a)  $(\mathbf{G}, \circ)$  is a topological group.
- (b) The action of  $\mathbf{G}$  on  $X$  is continuous.

*Proof.* To prove (a) it is sufficient to prove that the maps  $(g', g'') \mapsto g' \circ g''$  and  $g \mapsto g^{-1}$  are continuous. First of all, we have to prove that if a sequence  $(g'_i)$  converges to  $g'$  and a sequence  $(g''_i)$  converges to  $g''$  in  $G$ , then the sequence

$(g'_i \circ g''_i)$  converges to  $g' \circ g''$  in  $G$ . We observe that, for every  $x \in X$ ,

$$\begin{aligned} d_X((g'_i \circ g''_i) * x, (g' \circ g'') * x) &= d_X(g'_i * (g''_i * x), g' * (g'' * x)) \\ &\leq d_X(g'_i * (g''_i * x), g'_i * (g'' * x)) \\ &\quad + d_X(g'_i * (g'' * x), g' * (g'' * x)) \\ &= d_X(g''_i * x, g'' * x) \\ &\quad + d_X(g'_i * (g'' * x), g' * (g'' * x)) \\ &\leq d_G(g''_i, g'') + d_G(g'_i, g'). \end{aligned}$$

Thus,  $d_G(g'_i \circ g''_i, g' \circ g'') \leq d_G(g''_i, g'') + d_G(g'_i, g')$ . This proves the first property. Then, we have to prove that if a sequence  $(g_i)$  converges to  $g$  in  $G$ , then the sequence  $(g_i^{-1})$  converges to  $g^{-1}$  in  $G$ . We have that

$$\begin{aligned} d_X(g_i^{-1} * x, g^{-1} * x) &= d_X(g_i * (g_i^{-1} * x), g_i * (g^{-1} * x)) \\ &= d_X((g_i \circ g_i^{-1}) * x, (g_i \circ g^{-1}) * x) \\ &= d_X(x, (g_i \circ g^{-1}) * x) \\ &= d_X((g \circ g^{-1}) * x, (g_i \circ g^{-1}) * x) \\ &= d_X(g * (g^{-1} * x), g_i * (g^{-1} * x)) \\ &\leq d_G(g, g_i). \end{aligned}$$

Therefore,  $d_G(g_i^{-1}, g^{-1}) \leq d_G(g, g_i)$ . This proves our second property.

Now we prove (b). We have to prove that if a sequence  $(x_i)$  converges to  $x$  in  $X$  and a sequence  $(g_i)$  converges to  $g$  in  $G$ , then the sequence  $(g_i * x_i)$  converges to  $g * x$  in  $X$ . Since  $\lim_{i \rightarrow \infty} x_i = x$  and  $\lim_{i \rightarrow \infty} g_i = g$ , then  $\lim_{i \rightarrow \infty} d_X(x_i, x) = 0$  and  $\lim_{i \rightarrow \infty} d_X(g_i * x, g * x) = 0$ . We have that, for every  $x \in X$ ,

$$\begin{aligned} d_X(g_i * x_i, g * x) &\leq d_X(g_i * x_i, g_i * x) + d_X(g_i * x, g * x) \\ &= d_X(x_i, x) + d_X(g_i * x, g * x) \\ &\leq d_X(x_i, x) + d_G(g_i, g). \end{aligned}$$

*Semantics of Diagrams.* It is convenient to first fix some notation.

*Remark 4 (Notation).* Given two sets  $X$  and  $Y$ , we write  $X \times Y$  for their Cartesian product and  $\sigma_{X,Y} : X \times Y \rightarrow Y \times X$  for the symmetry function mapping  $(x, y) \in X \times Y$  into  $(y, x) \in Y \times X$ ; given two functions  $f_1 : X_1 \rightarrow Y_1$  and  $f_2 : X_2 \rightarrow Y_2$ , we write  $f_1 \times f_2 : X_1 \times X_2 \rightarrow Y_1 \times Y_2$  for the function mapping  $(x_1, x_2) \in X_1 \times X_2$  into  $(f(x_1), f(x_2)) \in Y_1 \times Y_2$ ; Given  $f : X \rightarrow Y$  and  $g : Y \rightarrow Z$ , we write  $g \circ f : X \rightarrow Z$  for their composition. For an arbitrary set  $X$ , we write  $\text{id}_X : X \rightarrow X$  for the identity function, and  $\Delta_X : X \rightarrow X \times X$  for the copier function mapping  $x \in X$  into  $(x, x) \in X \times X$ ; We write  $1$  for a singleton set that we fix to be  $\{\star\}$  and  $!_X : X \rightarrow 1$  for the function mapping any  $x \in X$  into  $\star$ .

Given two perception spaces  $(X, G)$  and  $(Y, K)$ , their direct product written  $(X, G) \otimes (Y, K)$  is the perception space  $(X \times Y, G \times K)$ , where the distance on  $X \times Y$  is defined as  $d_{X \times Y}((x_1, y_1), (x_2, y_2)) := \max\{d_X(x_1, x_2), d_Y(y_1, y_2)\}$  while the group action is defined pointwise, that is  $(g, k) * (x, y) = (g * x, k * y)$ . We write  $\sigma_{(X,G),(Y,K)} : (X, G) \otimes (Y, K) \rightarrow (Y, K) \otimes (X, G)$  as  $(\sigma_{X,Y}, \sigma_{G,K})$ .

**Table 4.** The CD category of GEOs. Above  $(f, t): (X, G) \rightarrow (Y, K)$ ,  $(f', t'): (Y, K) \rightarrow (Z, L)$  and  $(f_1, t_1): (X_1, G_1) \rightarrow (Y_1, K_1)$ ,  $(f_2, t_2): (X_2, G_2) \rightarrow (Y_2, K_2)$  are GEOs. The notation on the right hand side is in Remark 4.

$$\begin{aligned}
 \text{id}_{X,G} &:= (\text{id}_X, \text{id}_G): (X, G) \rightarrow (X, G) \\
 \Delta_{X,G} &:= (\Delta_X, \Delta_G): (X, G) \rightarrow (X, G) \otimes (X, G) \\
 !_{X,G} &:= (!_X, !_G): (X, G) \rightarrow (1, 1) \\
 \sigma_{(X,G),(Y,K)} &:= (\sigma_{X,Y}, \sigma_{G,K}): (X, G) \otimes (Y, K) \rightarrow (Y, K) \otimes (X, G) \\
 (f, t) \circ (f', t') &:= (f' \circ f, t' \circ t): (X, G) \rightarrow (Z, L) \\
 (f_1, t_1) \otimes (f_2, t_2) &:= (f_1 \times f_2, t_1 \times t_2): (X_1, G_1) \otimes (X_2, G_2) \rightarrow (Y_1, K_1) \otimes (Y_2, K_2)
 \end{aligned}$$

**Table 5.** The semantics  $\llbracket - \rrbracket_{\mathcal{I}}: \mathbf{Diag} \rightarrow \mathbf{GEO}$  for an interpretation  $\mathcal{I}$ . Operations and constants occurring on the right hand side of the above equations are those in (4). Above  $\mathcal{I}_{\mathcal{S}}$  is a function mapping each  $A \in \mathcal{S}$  in a perception space such that, for all  $g \in \Gamma$  with arity  $A_1 \otimes \dots \otimes A_n$  and coarity  $B_1 \otimes \dots \otimes B_m$ , the source of  $\mathcal{I}(g)$  is  $\otimes_{i=1}^n \mathcal{I}_{\mathcal{S}}(A_i)$  and its target is  $\otimes_{j=1}^m \mathcal{I}_{\mathcal{S}}(B_j)$ .

$$\begin{aligned}
 \llbracket \begin{array}{c} \text{---} \\ \vdots \\ \text{---} \end{array} \begin{array}{c} \text{---} \\ \vdots \\ \text{---} \end{array} \rrbracket_{\mathcal{I}} &:= \mathcal{I}(g) & \llbracket \text{---} \rrbracket_{\mathcal{I}} &:= \text{id}_{1,1} & \llbracket c_1 \otimes c_2 \rrbracket_{\mathcal{I}} &:= \llbracket c_1 \rrbracket_{\mathcal{I}} \otimes \llbracket c_2 \rrbracket_{\mathcal{I}} \\
 \llbracket \begin{array}{c} \text{---} \\ \text{---} \end{array} \begin{array}{c} \text{---} \\ \text{---} \end{array} \rrbracket_{\mathcal{I}} &:= \Delta_{\mathcal{I}_{\mathcal{S}}(A)} & \llbracket \begin{array}{c} \text{---} \\ \text{---} \end{array} \rrbracket_{\mathcal{I}} &:= \sigma_{\mathcal{I}_{\mathcal{S}}(A), \mathcal{I}_{\mathcal{S}}(B)} & \llbracket c_1 \circ c_2 \rrbracket_{\mathcal{I}} &:= \llbracket c_1 \rrbracket_{\mathcal{I}} \circ \llbracket c_2 \rrbracket_{\mathcal{I}} \\
 \llbracket \text{---} \text{---} \rrbracket_{\mathcal{I}} &:= \text{id}_{\mathcal{I}_{\mathcal{S}}(A)} & \llbracket \text{---} \bullet \rrbracket_{\mathcal{I}} &:= !_{\mathcal{I}_{\mathcal{S}}(A)}
 \end{aligned}$$

With this notation one can extend the above structures of sets and functions to perception spaces and GEOs as illustrated in Table 4. By simply checking that the definitions in Table 4 provide GEOs, one can prove the following result.

**Lemma 1.** *GEO is a CD category in the sense of [12].*

From this fact, and the observation that  $\mathbf{Diag}_{\Gamma}$  is the (strict) CD category freely generated from the monoidal signature  $\Gamma$ , one obtains that, for each interpretation  $\mathcal{I}$ , there exists a unique CD functor  $\llbracket - \rrbracket_{\mathcal{I}}: \mathbf{Diag} \rightarrow \mathbf{GEO}$  extending  $\mathcal{I}$ . Its inductive definition is illustrated in Table 5

*Cost of Translation Pairs for Infinite Perception Spaces.* Here we explain how the cost of translation pairs defined in (1) can be defined for arbitrary sets  $X_{\alpha}$ .

To proceed, we need to equip each metric space  $X_{\alpha}$  with a Borel probability measure  $\mu_{\alpha}$ , in the spirit of [11]. In simple terms, the measure  $\mu_{\alpha}$  represents the probability of each data point in  $X_{\alpha}$  appearing in our experiments. We will assume that all GENEOS in  $\mathbf{T}$  are not just distance-decreasing (i.e., non-expansive) but also *measure-decreasing*, i.e., if  $(l_{\alpha,\beta}, p_{\alpha,\beta}): (X_{\alpha}, G_{\alpha}) \rightarrow (X_{\beta}, G_{\beta})$  belongs to  $\mathbf{T}$  and the set  $A \subseteq X_{\alpha}$  is measurable for  $\mu_{\alpha}$ , then  $l_{\alpha,\beta}(A)$  is measurable for  $\mu_{\beta}$ , and  $\mu_{\beta}(l_{\alpha,\beta}(A)) \leq \mu_{\alpha}(A)$ . Moreover, we assume that the function  $f_{\alpha,\beta}: X_{\alpha} \rightarrow \mathbb{R}$ , defined for every  $x \in X_{\alpha}$  as  $f_{\alpha,\beta}(x) := d_{\text{dt}}((m_{\beta,\alpha} \circ f_{\beta} \circ l_{\alpha,\beta})(x), f_{\alpha}(x))$ , is integrable with respect to  $\mu_{\alpha}$ .

**Definition 8.** Let  $\pi = \left( (l_{\alpha,\beta}, p_{\alpha,\beta}), (m_{\beta,\alpha}, q_{\beta,\alpha}) \right)$  be a crossed translation pair from  $(f_\alpha, t_\alpha): (X_\alpha, G_\alpha) \rightarrow (Y_\alpha, K_\alpha)$  to  $(f_\beta, t_\beta): (X_\beta, G_\beta) \rightarrow (Y_\beta, K_\beta)$ . The functional cost of  $\pi$ , written  $\text{cost}(\pi)$ , is defined as follows.

$$\text{cost}(\pi) = \int_{X_\alpha} d_{\text{dt}} \left( (m_{\beta,\alpha} \circ f_\beta \circ l_{\alpha,\beta})(x), f_\alpha(x) \right) d\mu_\alpha \tag{4}$$

*Proof of Theorem 1.* For sake of generality, we illustrate the proof for the case where  $\text{cost}(\pi)$  is defined as in (4). The case of  $\text{cost}(\pi)$  as in (1) follows by fixing

$\mu_\alpha$  as uniform Borel measure. Let us prove that  $h_\mathbb{O}$  enjoys the triangle inequality, i.e.,  $h_\mathbb{O}(\alpha, \gamma) \leq h_\mathbb{O}(\alpha, \beta) + h_\mathbb{O}(\beta, \gamma)$ , where  $\alpha, \beta$  and  $\gamma$  are three

$$\begin{aligned} \alpha &:= (f_\alpha, t_\alpha): (X_\alpha, G_\alpha) \rightarrow (Y_\alpha, K_\alpha) \\ \beta &:= (f_\beta, t_\beta): (X_\beta, G_\beta) \rightarrow (Y_\beta, K_\beta) \\ \gamma &:= (f_\gamma, t_\gamma): (X_\gamma, G_\gamma) \rightarrow (Y_\gamma, K_\gamma) \end{aligned}$$

GEOs in  $\mathcal{G}$  illustrated on the right. We consider three translation pairs:

$$\begin{aligned} \pi_1 &:= \left( (l_{\alpha,\beta}, p_{\alpha,\beta}), (m_{\beta,\alpha}, q_{\beta,\alpha}) \right): \alpha \rightleftharpoons_{\mathbf{T}} \beta \\ \pi_2 &:= \left( (l_{\beta,\gamma}, p_{\beta,\gamma}), (m_{\gamma,\beta}, q_{\gamma,\beta}) \right): \beta \rightleftharpoons_{\mathbf{T}} \gamma \\ \pi_3 &:= \pi_2 \circ \pi_1 = \left( (l_{\beta,\gamma} \circ l_{\alpha,\beta}, p_{\beta,\gamma} \circ p_{\alpha,\beta}), (m_{\beta,\alpha} \circ m_{\gamma,\beta}, q_{\beta,\alpha} \circ q_{\gamma,\beta}) \right): \beta \rightleftharpoons_{\mathbf{T}} \gamma \end{aligned}$$

Please note that if no crossed pair like  $\pi_1$  or  $\pi_2$  exists, then  $h_\mathbb{O}(\alpha, \beta) + h_\mathbb{O}(\beta, \gamma) = \infty$ , and hence the triangle inequality trivially holds. By definition their costs are

$$\begin{aligned} \text{cost}(\pi_1) &= \int_{X_\alpha} d_{\text{dt}} \left( (m_{\beta,\alpha} \circ f_\beta \circ l_{\alpha,\beta})(x), f_\alpha(x) \right) d\mu_\alpha \\ \text{cost}(\pi_2) &= \int_{X_\beta} d_{\text{dt}} \left( (m_{\gamma,\beta} \circ f_\gamma \circ l_{\beta,\gamma})(y), f_\beta(y) \right) d\mu_\beta \\ \text{cost}(\pi_3) &= \int_{X_\alpha} d_{\text{dt}} \left( (m_{\beta,\alpha} \circ m_{\gamma,\beta} \circ f_\gamma \circ l_{\beta,\gamma} \circ l_{\alpha,\beta})(x), f_\alpha(x) \right) d\mu_\alpha \end{aligned}$$

Since  $(m_{\beta,\alpha}, q_{\beta,\alpha})$  is a GENEO, we have that for every  $y \in X_\beta$ ,

$$d_{\text{dt}} \left( (m_{\gamma,\beta} \circ f_\gamma \circ l_{\beta,\gamma})(y), f_\beta(y) \right) \geq d_{\text{dt}} \left( (m_{\beta,\alpha} \circ m_{\gamma,\beta} \circ f_\gamma \circ l_{\beta,\gamma})(y), (m_{\beta,\alpha} \circ f_\beta)(y) \right)$$

and hence, setting  $y := l_{\alpha,\beta}(x)$  and recalling that  $l_{\alpha,\beta}$  is measure-decreasing,

$$\begin{aligned} &\int_{X_\beta} d_{\text{dt}} \left( (m_{\gamma,\beta} \circ f_\gamma \circ l_{\beta,\gamma})(y), f_\beta(y) \right) d\mu_\beta \\ &\geq \int_{X_\alpha} d_{\text{dt}} \left( (m_{\beta,\alpha} \circ m_{\gamma,\beta} \circ f_\gamma \circ l_{\beta,\gamma} \circ l_{\alpha,\beta})(x), (m_{\beta,\alpha} \circ f_\beta)(y) \circ l_{\alpha,\beta}(x) \right) d\mu_\alpha. \end{aligned}$$

Therefore, we have that  $\text{cost}(\pi_1) + \text{cost}(\pi_2) =$

$$\begin{aligned} &= \int_{X_\alpha} d_{\text{dt}}\left((m_{\beta,\alpha} \circ f_\beta \circ l_{\alpha,\beta})(x), f_\alpha(x)\right) d\mu_\alpha + \int_{X_\beta} d_{\text{dt}}\left((m_{\gamma,\beta} \circ f_\gamma \circ l_{\beta,\gamma})(y), f_\beta(y)\right) d\mu_\beta \\ &\geq \int_{X_\alpha} d_{\text{dt}}\left((m_{\beta,\alpha} \circ f_\beta \circ l_{\alpha,\beta})(x), f_\alpha(x)\right) d\mu_\alpha \\ &\quad + \int_{X_\alpha} d_{\text{dt}}\left((m_{\beta,\alpha} \circ m_{\gamma,\beta} \circ f_\gamma \circ l_{\beta,\gamma} \circ l_{\alpha,\beta})(x), (m_{\beta,\alpha} \circ f_\beta \circ l_{\alpha,\beta})(x)\right) d\mu_\alpha \\ &\geq \int_{X_\alpha} d_{\text{dt}}\left((m_{\beta,\alpha} \circ m_{\gamma,\beta} \circ f_\gamma \circ l_{\beta,\gamma} \circ l_{\alpha,\beta})(x), f_\alpha(x)\right) d\mu_\alpha = \text{cost}(\pi_2 \circ \pi_1) \end{aligned}$$

where the second to last inequality follows from the triangle inequality for  $d_{\text{dt}}$ . Therefore,  $\text{cost}(\pi_1) + \text{cost}(\pi_2) \geq \text{cost}(\pi_2 \circ \pi_1)$ . It follows that

$$\begin{aligned} &\inf\{\text{cost}(\pi') \mid \pi' : \alpha \rightleftharpoons_{\mathbf{T}} \beta\} + \inf\{\text{cost}(\pi'') \mid \pi'' : \beta \rightleftharpoons_{\mathbf{T}} \gamma\} \\ &= \inf\{\text{cost}(\pi') + \text{cost}(\pi'') \mid \pi' : \alpha \rightleftharpoons_{\mathbf{T}} \beta, \pi'' : \beta \rightleftharpoons_{\mathbf{T}} \gamma\} \\ &\geq \inf\{\text{cost}(\pi'' \circ \pi') \mid \pi' : \alpha \rightleftharpoons_{\mathbf{T}} \beta, \pi'' : \beta \rightleftharpoons_{\mathbf{T}} \gamma\} \\ &\geq \inf\{\text{cost}(\pi) \mid \pi : \alpha \rightleftharpoons_{\mathbf{T}} \gamma\} \end{aligned}$$

and thus  $h_{\mathbb{O}}(\alpha, \beta) + h_{\mathbb{O}}(\beta, \gamma) \geq h_{\mathbb{O}}(\alpha, \gamma)$ . In other words, (T) holds.

To prove (R) i.e., that for all GEOs  $(f_\alpha, t_\alpha) : (X_\alpha, G_\alpha) \rightarrow (Y_\alpha, K_\alpha)$ , it holds that  $h_{\mathbb{O}}\left((f_\alpha, t_\alpha), (f_\alpha, t_\alpha)\right) = 0$ , observe that, since  $\mathbf{T}$  is a category there exists the crossed pair of translation  $\iota := \left((\text{id}_{X_\alpha}, \text{id}_{G_\alpha}), (\text{id}_{Y_\alpha}, \text{id}_{K_\alpha})\right)$  given by the identity morphisms. One can easily check that  $\text{cost}(\iota) = 0$  and thus

$$\inf\{\text{cost}(\pi) \mid \pi : (f_\alpha, t_\alpha) \rightleftharpoons_{\mathbf{T}} (f_\alpha, t_\alpha)\} = 0.$$

*Proof of Proposition 1.* Fix  $A := \{(g, x) \mid f_\alpha(x) = f_\beta(x)\}$ ,  $B := \{(g, x) \mid f_\alpha(g * x) = f_\beta(g * x)\}$  and  $C := \{(g, x) \mid f_\beta(x) = f_\beta(g * x)\}$  and observe that  $A \cap B \subseteq C$ . Thus, by denoting with  $\overline{X}$ , the complement of a set  $X$ , it holds that  $\overline{A} \cup \overline{B} \supseteq \overline{C}$  and thus

$$|\overline{A}| + |\overline{B}| \geq |\overline{C}|. \tag{5}$$

We now use the hypothesis that  $G_\alpha$  is a group, to show the bijection of  $\overline{A}$  and  $\overline{B}$ : define  $\iota : \overline{B} \rightarrow \overline{A}$  as  $\iota(g, x) := (g, g * x)$  and  $\kappa : \overline{A} \rightarrow \overline{B}$  as  $\kappa(g, x) := (g, g^{-1} * x)$ . Observe that the functions are well defined and that they are inverse to each other. Thus  $|\overline{A}| = |\overline{B}|$  that, thanks to (5) gives us

$$2 \cdot |\overline{A}| \geq |\overline{C}|.$$

To conclude observe that  $\overline{C}$  is NE and that  $|\overline{A}|$  is  $|G_\alpha| \cdot h_{\mathbb{O}}((f_\alpha, t_\alpha), (f_\beta, t_\beta))$ .

## References

1. Ahmad, F., Ferri, M., Frosini, P.: Generalized permutants and graph GENEOS. *Mach. Learn. Knowl. Extract.* **5**(4), 1905–1920 (2023). <https://doi.org/10.3390/make5040092>
2. Alangari, N., Menai, M.E.B., Mathkour, H., Almosallam, I.: Exploring evaluation methods for interpretable machine learning: a survey. *Information* **14**(8), 469 (2023)
3. Anselmi, F., Rosasco, L., Poggio, T.: On invariance and selectivity in representation learning. *Inf. Infer. J. IMA* **5**(2), 134–158 (2016). <https://doi.org/10.1093/imaiai/iaw009>
4. Arora, S., Pruthi, D., Sadeh, N.M., Cohen, W.W., Lipton, Z.C., Neubig, G.: Explain, edit, and understand: rethinking user study design for evaluating model explanations. In: *AAAI*, pp. 5277–5285. AAAI Press (2022)
5. Bengio, Y., Courville, A., Vincent, P.: Representation learning: a review and new perspectives. *IEEE Trans. Pattern Anal. Mach. Intell.* **35**(8), 1798–1828 (2013)
6. Bergomi, M.G., Frosini, P., Giorgi, D., Quercioli, N.: Towards a topological-geometrical theory of group equivariant non-expansive operators for data analysis and machine learning. *Nat. Mach. Intell.* **1**(9), 423–433 (2019). <https://doi.org/10.1038/s42256-019-0087-3>
7. Bocchi, G., Ferri, M., Frosini, P.: A novel approach to graph distinction through geneos and permutants. *Sci. Rep.* **15**(1), 6259 (2025). <https://doi.org/10.1038/s41598-025-90152-7>
8. Bocchi, G., et al.: A geometric XAI approach to protein pocket detection. In: *xAI-2024 Late-breaking Work, Demos and Doctoral Consortium Joint Proceedings - The 2nd World Conference on eXplainable Artificial Intelligence*, vol. 3793, pp. 217–224 (2024). [https://ceur-ws.org/Vol-3793/paper\\_28.pdf](https://ceur-ws.org/Vol-3793/paper_28.pdf)
9. Bodria, F., Giannotti, F., Guidotti, R., Naretto, F., Pedreschi, D., Rinzivillo, S.: Benchmarking and survey of explanation methods for black box models. *Data Min. Knowl. Disc.* **37**(5), 1719–1778 (2023)
10. Camporesi, F., Frosini, P., Quercioli, N.: On a new method to build group equivariant operators by means of permutants. In: Holzinger, A., Kieseberg, P., Tjoa, A.M., Weippl, E. (eds.) *CD-MAKE 2018. LNCS*, vol. 11015, pp. 265–272. Springer, Cham (2018). [https://doi.org/10.1007/978-3-319-99740-7\\_18](https://doi.org/10.1007/978-3-319-99740-7_18)
11. Cascarano, P., Frosini, P., Quercioli, N., Saki, A.: On the geometric and riemannian structure of the spaces of group equivariant non-expansive operators (2023). <https://arxiv.org/abs/2103.02543>
12. Cho, K., Jacobs, B.: Disintegration and bayesian inversion via string diagrams. *Math. Struct. Comput. Sci.* **29**(7), 938–971 (2019)
13. Ciravegna, G., et al.: Logic explained networks. *Artif. Intell.* **314**, 103822 (2023)
14. Cohen, T., Welling, M.: Group equivariant convolutional networks. In: *International Conference on Machine Learning*, pp. 2990–2999 (2016)
15. Collaris, D., Gajane, P., Jorritsma, J., van Wijk, J.J., Pechenizkiy, M.: LEMON: alternative sampling for more faithful explanation through local surrogate models. In: *IDA. Lecture Notes in Computer Science*, vol. 13876, pp. 77–90. Springer, Heidelberg (2023). [https://doi.org/10.1007/978-3-031-30047-9\\_7](https://doi.org/10.1007/978-3-031-30047-9_7)
16. Colley, A., Kalving, M., Häkkinä, J., Väänänen, K.: Exploring tangible explainable AI (tangxai): a user study of two XAI approaches. In: *OZCHI*, pp. 679–683. ACM (2023)

17. Cybenko, G.: Approximation by superpositions of a sigmoidal function. *Math. Control Signals Syst.* **2**(4), 303–314 (1989)
18. Frosini, P., Jabłoński, G.: Combining persistent homology and invariance groups for shape comparison. *Disc. Comput. Geom.* **55**(2), 373–409 (2016). <https://doi.org/10.1007/s00454-016-9761-y>
19. Gerken, J.E., et al.: Geometric deep learning and equivariant neural networks. *Artif. Intell. Rev.* (2023)
20. Gerken, J.E., et al.: Geometric deep learning and equivariant neural networks. *Artif. Intell. Rev.* **56**(12), 14605–14662 (2023)
21. Giannini, F., Fioravanti, S., Barbiero, P., Tonda, A., Liò, P., Di Lavore, E.: Categorical foundation of explainable AI: a unifying theory. In: *World Conference on Explainable Artificial Intelligence*, pp. 185–206. Springer, Heidelberg (2024). [https://doi.org/10.1007/978-3-031-63800-8\\_10](https://doi.org/10.1007/978-3-031-63800-8_10)
22. Guidotti, R., Monreale, A., Ruggieri, S., Turini, F., Giannotti, F., Pedreschi, D.: A survey of methods for explaining black box models. *ACM Comput. Surv.* **51**(5), 93:1–93:42 (2019)
23. Harder, F., Bauer, M., Park, M.: Interpretable and differentially private predictions. In: *AAAI*, pp. 4083–4090. AAAI Press (2020)
24. Heidari, F., Taslakian, P., Rabusseau, G.: Explaining graph neural networks using interpretable local surrogates. In: *TAG-ML. Proceedings of Machine Learning Research*, vol. 221, pp. 146–155. PMLR (2023)
25. Hoffman, R.R., Klein, G.: Explaining explanation, part 1: theoretical foundations. *IEEE Intell. Syst.* **32**(3), 68–73 (2017)
26. Hoffman, R.R., Mueller, S.T., Klein, G.: Explaining explanation, part 2: empirical foundations. *IEEE Intell. Syst.* **32**(4), 78–86 (2017)
27. Kay, J.: Foundations for human-AI teaming for self-regulated learning with explainable AI (XAI). *Comput. Hum. Behav.* **147**, 107848 (2023)
28. Lualdi, P., Sturm, R., Siefkes, T.: Exploration-oriented sampling strategies for global surrogate modeling: a comparison between one-stage and adaptive methods. *J. Comput. Sci.* **60**, 101603 (2022)
29. Marconato, E., Passerini, A., Teso, S.: Interpretability is in the mind of the beholder: a causal framework for human-interpretable representation learning. *Entropy* **25**(12), 1574 (2023)
30. Mirzaei, S., Mao, H., Al-Nima, R.R.O., Woo, W.L.: Explainable AI evaluation: a top-down approach for selecting optimal explanations for black box models. *Inf.* **15**(1), 4 (2024)
31. Murdoch, W.J., Singh, C., Kumbier, K., Abbasi-Asl, R., Yu, B.: Definitions, methods, and applications in interpretable machine learning. *Proc. Natl. Acad. Sci.* 22071–22080 (2019). <https://doi.org/10.1073/pnas.1900654116>
32. Nauta, M., et al.: From anecdotal evidence to quantitative evaluation methods: a systematic review on evaluating explainable AI. *ACM Comput. Surv.* **55**(13s), 295:1–295:42 (2023)
33. Palacio, S., Lucieri, A., Munir, M., Ahmed, S., Hees, J., Dengel, A.: Xai handbook: towards a unified framework for explainable ai. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 3766–3775 (2021)
34. Panigutti, C., et al.: Co-design of human-centered, explainable AI for clinical decision support. *ACM Trans. Interact. Intell. Syst.* **13**(4), 21:1–21:35 (2023)
35. Ruhe, D., Brandstetter, J., Forré, P.: Clifford group equivariant neural networks. *Adv. Neural. Inf. Process. Syst.* **36**, 62922–62990 (2023)

36. Satorras, V.G., Hoogeboom, E., Welling, M.: E (n) equivariant graph neural networks. In: International Conference on Machine Learning, pp. 9323–9332. PMLR (2021)
37. Schulze-Weddige, S., Zylowski, T.: User study on the effects explainable AI visualizations on non-experts. In: ArtsIT. Lecture Notes of the Institute for Computer Sciences, Social Informatics and Telecommunications Engineering, vol. 422, pp. 457–467. Springer, Heidelberg (2021). [https://doi.org/10.1007/978-3-030-95531-1\\_31](https://doi.org/10.1007/978-3-030-95531-1_31)
38. Selinger, P.: A survey of graphical languages for monoidal categories. In: New Structures for Physics, pp. 289–355. Springer, Heidelberg (2010). [https://doi.org/10.1007/978-3-642-12821-9\\_4](https://doi.org/10.1007/978-3-642-12821-9_4)
39. Tull, S., Lorenz, R., Clark, S., Khan, I., Coecke, B.: Towards compositional interpretability for xai. arXiv preprint [arXiv:2406.17583](https://arxiv.org/abs/2406.17583) (2024)
40. Wagner, J., Köhler, J.M., Gindele, T., Hetzel, L., Wiedemer, J.T., Behnke, S.: Interpretable and fine-grained visual explanations for convolutional neural networks. In: CVPR, pp. 9097–9107. Computer Vision Foundation/IEEE (2019)
41. Worrall, D.E., Garbin, S.J., Turmukhambetov, D., Brostow, G.J.: Harmonic networks: deep translation and rotation equivariance. In: Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR), vol. 2, pp. 7168–7177 (2017)
42. Yang, W., et al.: Survey on explainable AI: from approaches, limitations and applications aspects. *Hum. Centric Intell. Syst.* **3**(3), 161–188 (2023)
43. Zhukov, A., Benois-Pineau, J., Giot, R.: Evaluation of explanation methods of AI - cnns in image classification tasks with reference-based and no-reference metrics. *Adv. Artif. Intell. Mach. Learn.* **3**(1), 620–646 (2023)

**Open Access** This chapter is licensed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license and indicate if changes were made.

The images or other third party material in this chapter are included in the chapter's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the chapter's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder.

