

---

# Unveiling Transformer Perception by Exploring Input Manifolds

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 This paper introduces a general method for the exploration of equivalence classes in  
2 the input space of Transformer models. The proposed approach is based on sound  
3 mathematical theory which describes the internal layers of a Transformer architec-  
4 ture as sequential deformations of the input manifold. Using eigendecomposition  
5 of the pullback of the distance metric defined on the output space through the  
6 Jacobian of the model, we are able to reconstruct equivalence classes in the input  
7 space and navigate across them. We illustrate how this method can be used as a  
8 powerful tool for investigating how a Transformer sees the input space, facilitating  
9 local and task-agnostic explainability in Computer Vision and Natural Language  
10 Processing tasks.

## 11 1 Introduction

12 In this paper, we propose a method for exploring the input space of Transformer models by identifying  
13 equivalence classes with respect to their predictions. We define an *equivalence class* of a Transformer  
14 model as the set of vectors in the embedding space whose outcomes under the Transformer process  
15 are the same. The study of the input manifold on which the inverse image of models lies provides  
16 insights for both explainability and sensitivity analyses. Existing methods aiming at the exploration  
17 of the input space of Deep Neural Networks and Transformers either rely on perturbations of input  
18 data using heuristic or gradient-based criteria [16, 22, 17, 14], or they analyze specific properties of  
19 the embedding space [5].

20 Our approach is based on sound mathematical theory which describes the internal layers of a  
21 Transformer architecture as sequential deformations of the input manifold. Using eigendecomposition  
22 of the pullback of the distance metric defined on the output space through the Jacobian of the model,  
23 we are able to reconstruct equivalence classes in the input space and navigate across them. In the  
24 XAI scenario, our framework can facilitate local and task-agnostic explainability methods applicable  
25 to Computer Vision (CV) and Natural Language Processing (NLP) tasks, among others.

26 In Section 2, we summarise the preliminaries of the mathematical foundations of our approach.  
27 In Section 3, we present our method for the exploration of equivalence classes in the input of the  
28 Transformer models. In Section 4, we perform a preliminary investigation of some applicability  
29 options of our method on textual and visual data. In Section 5, we discuss the relevant literature about  
30 embedding space exploration and feature importance. Finally, in Section 6, we give our concluding  
31 remarks<sup>1</sup>.

---

<sup>1</sup>The code to reproduce our experiments can be found in the Supplementary Materials.

## 32 2 Preliminaries

33 In this Section, we provide the theoretical foundation of the proposed approach, namely the Geometric  
34 Deep Learning framework based on Riemannian Geometry [2].

35 A neural network is considered as a sequence of maps, the layers of the network, between manifolds,  
36 and the latter are the spaces where the input and the outputs of the layers belong to.

37 **Definition 1** (Neural Network). *A neural network is a sequence of  $C^1$  maps  $\Lambda_i$  between manifolds of  
38 the form:*

$$M_0 \xrightarrow{\Lambda_1} M_1 \xrightarrow{\Lambda_2} M_2 \xrightarrow{\Lambda_3} \dots \xrightarrow{\Lambda_{n-1}} M_{n-1} \xrightarrow{\Lambda_n} M_n \quad (1)$$

39 We call  $M_0$  the input manifold and  $M_n$  the output manifold. All the other manifolds of the sequence  
40 are called representation manifolds. The maps  $\Lambda_i$  are the layers of the neural network. We denote  
41 with  $\mathcal{N}_{(i)} = \Lambda_n \circ \dots \circ \Lambda_i : M_i \rightarrow M_n$  the mapping from the  $i$ -th representation layer to the output  
42 layer.

43 As an example, consider a shallow network with just one layer, the composition of a linear operator  
44  $A \cdot + b$  with a sigmoid function  $\sigma$ , where  $A \in \mathbb{R}^{m \times n}$  and  $b \in \mathbb{R}^m$ : then, the input manifold  $M_0$  and  
45 the output manifold  $M_1$  shall be  $\mathbb{R}^n$  and  $\mathbb{R}^m$ , respectively, and the map  $\Lambda_1(\cdot) = \sigma(A \cdot + b)$ . We  
46 generalize this observation into the following definition.

47 **Definition 2** (Smooth layer). *A map  $\Lambda_i : M_{i-1} \rightarrow M_i$  is called a smooth layer if it is the restriction  
48 to  $M_{i-1}$  of a function  $\bar{\Lambda}^{(i)}(x) : \mathbb{R}^{d_{i-1}} \rightarrow \mathbb{R}^{d_i}$  of the form*

$$\bar{\Lambda}_\alpha^{(i)}(x) = F_\alpha^{(i)} \left( \sum_\beta A_{\alpha\beta}^{(i)} x_\beta + b_\alpha^{(i)} \right) \quad (2)$$

49 for  $i = 1, \dots, n$ ,  $x \in \mathbb{R}^{d_i}$ ,  $b^{(i)} \in \mathbb{R}^{d_i}$  and  $A^{(i)} \in \mathbb{R}^{d_i \times d_{i-1}}$ , with  $F^{(i)} : \mathbb{R}^{d_i} \rightarrow \mathbb{R}^{d_i}$  a diffeomor-  
50 phism.

51 **Remark 1.** *Transformers implicitly apply for this framework, since their modules are smooth  
52 functions, such as fully connected layers, GeLU and sigmoid activations.*

53 Our aim is to transport the geometric information on the data lying in the output manifold to the  
54 input manifold: this allows us to obtain insight on how the network "sees" the input space, how it  
55 manipulates it for reaching its final conclusion. For fulfilling this objective, we need several tools  
56 from differential geometry. The first key ingredient is the notion of singular Riemannian metric,  
57 which has the intuitive meaning of a degenerate scalar product which changes point to point.

58 **Definition 3** (Singular Riemannian metric). *Let  $M = \mathbb{R}^n$  or an open subset of  $\mathbb{R}^n$ . A singular  
59 Riemannian metric  $g$  over  $M$  is a map  $g : M \rightarrow \text{Bil}(\mathbb{R}^n \times \mathbb{R}^n)$  that associates to each point  $p$  a  
60 positive semidefinite symmetric bilinear form  $g_p : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$  in a smooth way.*

61 Without loss of generality, we can assume the following hypotheses on the sequence (1): *i) The  
62 manifolds  $M_i$  are open and path-connected sets of dimension  $\dim M_i = d_i$ . ii) The maps  $\Lambda_i$  are  $C^1$   
63 submersions. iii)  $\Lambda_i(M_{i-1}) = M_i$  for every  $i = 1, \dots, n$ . iv) The manifold  $M_n$  is equipped with  
64 the structure of Riemannian manifold, with metric  $g^{(n)}$ . Definition 3 naturally leads to the definition  
65 of the pseudolength and of energy of a curve.*

66 **Definition 4** (Pseudolength and energy of a curve). *Let  $\gamma : [a, b] \rightarrow \mathbb{R}^n$  a curve defined on the  
67 interval  $[a, b] \subset \mathbb{R}$  and  $\|v\|_p = \sqrt{g_p(v, v)}$  the pseudo-norm induced by the pseudo-metric  $g_p$  at  
68 point  $p$ . Then the pseudolength of  $\gamma$  and its energy are defined as*

$$Pl(\gamma) = \int_a^b \|\dot{\gamma}(s)\|_{\gamma(s)} ds = \int_a^b \sqrt{g_{\gamma(s)}(\dot{\gamma}(s), \dot{\gamma}(s))} ds, \quad E(\gamma) = \int_a^b \|\dot{\gamma}(s)\|_{\gamma(s)}^2 ds \quad (3)$$

69 The notion of pseudolength leads naturally to define the distance between two points.

70 **Definition 5** (Pseudodistance). *Let  $x, y \in M = \mathbb{R}^n$ . The pseudodistance between  $x$  and  $y$  is then*

$$Pd(x, y) = \inf\{Pl(\gamma) \mid \gamma : [0, 1] \rightarrow M, \gamma \in C^1([0, 1]), \gamma(0) = x, \gamma(1) = y\}. \quad (4)$$

71 One can observe that endowing the space  $\mathbb{R}^n$  with a singular Riemannian metric leads to have  
 72 non trivial curves whose length is zero. A straightforward consequence is that there are distinct  
 73 points whose pseudodistance is therefore zero: a natural equivalence relation arises, *i.e.*  $x \sim y \Leftrightarrow$   
 74  $Pd(x, y) = 0$ , obtaining thus a metric space  $(\mathbb{R}^n / \sim, Pd)$ .

75 The second crucial tool is the notion of pullback of a function. Let  $f$  be a function from  $\mathbb{R}^p$  to  $\mathbb{R}^q$ ,  
 76 and fix the coordinate systems  $x = (x_1, \dots, x_p)$  and  $y = (y_1, \dots, y_q)$  on  $\mathbb{R}^p$  and on  $\mathbb{R}^q$ , respectively.  
 77 Moreover, we endow  $\mathbb{R}^q$  with the standard Euclidean metric  $g$ , whose associated matrix is the identity.  
 78 The space  $\mathbb{R}^p$  can be equipped with the pullback metric  $f^*g$  whose representation matrix reads as

$$(f^*g)_{ij} = \sum_{h,k=1}^q \left( \frac{\partial f_h}{\partial x_i} \right) g_{hk} \left( \frac{\partial f_k}{\partial x_j} \right). \quad (5)$$

79 The sequence (1) shows that a neural network can be considered simply as a function, a composition  
 80 of maps: hence, taking  $f = \Lambda_n \circ \Lambda_{n-1} \circ \dots \circ \Lambda_1$  and supposing that  $M_0 = \mathbb{R}^p$ ,  $M_n = \mathbb{R}^q$ , the  
 81 generalization of (5) applied to (1) provides with the pullback of a generic neural network.

82 Hereafter, we consider in (1) the case  $M_n = \mathbb{R}^q$ , equipped with the trivial metric  $g^{(n)} = I_q$ , *i.e.*,  
 83 the identity. Each manifold  $M_i$  of the sequence (1) is equipped with a Riemannian singular metric,  
 84 denoted with  $g^{(i)}$ , obtained via the pullback of  $\mathcal{N}_i$ . The pseudolength of a curve  $\gamma$  on the  $i$ -th  
 85 manifold, namely  $Pl_i(\gamma)$ , is computed via the relative metric  $g^{(i)}$  via (3).

## 86 2.1 General results

87 We depict hereafter the theoretical bases of our approach. We denote with  $\mathcal{N}_i$  the submap  $\Lambda_i \circ \dots \circ \Lambda_n :$   
 88  $M_i \rightarrow M_n$ , and with  $\mathcal{N} \equiv \mathcal{N}_0$  the map describing the action of the complete network. The starting  
 89 point is to consider the pair  $(M_i, Pd_i)$ : this is a pseudometric space, which can be turned into a  
 90 full-fledged metric space  $M_i / \sim_i$  by the metric identification  $x \sim_i y \Leftrightarrow Pd_i(x, y) = 0$ . The first  
 91 result states that the length of a curve on the  $i$ -th manifold is preserved among the mapping on the  
 92 subsequent manifolds.

93 **Proposition 1.** *Let  $\gamma : [0, 1] \rightarrow M_i$  be a piecewise  $\mathcal{C}^1$  curve. Let  $k \in \{i, i+1, \dots, n\}$  and consider*  
 94 *the curve  $\gamma_k = \Lambda_k \circ \dots \circ \Lambda_i \circ \gamma$  on  $M_k$ . Then  $Pl_i(\gamma) = Pl_k(\gamma_k)$ .*

95 In particular this is true when  $k = n$ , *i.e.*, the length of a curve is preserved in the last manifold. This  
 96 result leads naturally to claim that if two points are in the same class of equivalence, then they are  
 97 mapped into the same point under the action of the neural network.

98 **Proposition 2.** *If two points  $p, q \in M_i$  are in the same class of equivalence, then  $\mathcal{N}_i(p) = \mathcal{N}_i(q)$ .*

99 The next step is to prove that the sets  $M_i / \sim_i$  are actually smooth manifolds: to this aim, we introduce  
 100 another equivalence relation:  $x \sim_{\mathcal{N}_i} y$  if and only if there exists a piecewise  $\gamma : [0, 1] \rightarrow M_i$  such  
 101 that  $\gamma(0) = x, \gamma(1) = y$  and  $\mathcal{N}_i \circ \gamma(s) = \mathcal{N}_i(x) \forall s \in [0, 1]$ . The introduction of this equivalence  
 102 relation allows us to easily state the following proposition.

103 **Proposition 3.** *Let  $x, y \in M_i$ , then  $x \sim_i y$  if and only if  $x \sim_{\mathcal{N}_i} y$ .*

104 The following corollary contains the natural consequences of the previous result; the second point of  
 105 the claim below is the counterpart of Proposition 2.

106 **Corollary 1.** *Under the hypothesis of Proposition 3, one has that  $M_i / \sim_i = M_i / \sim_{\mathcal{N}_{i+1}}$ . Moreover,*  
 107 *if two points  $p, q \in M_i$  are connected by a  $\mathcal{C}^1$  curve  $\gamma : [0, 1] \rightarrow M_i$  satisfying  $\mathcal{N}_i(p) = \mathcal{N}_i \circ \gamma(s)$*   
 108 *for every  $s \in [0, 1]$ , then they lie in the same class of equivalence.*

109 Making use of the Godement's criterion, we are now able to prove that the set  $M_i / \sim_i$  is a smooth  
 110 manifold, together with its dimension.

111 **Proposition 4.**  $\frac{M_i}{\sim_i}$  *is a smooth manifold of dimension  $dim(\mathcal{N}(M_0))$ .*

112 This last achievement provides practical insights about the projection  $\pi_i$  on the quotient space, that  
 113 consists the building block of the algorithms used for recovering and exploring the equivalence  
 114 classes of a neural network.

115 **Proposition 5.**  $\pi_i : M_i \rightarrow M_i / \sim_i$  is a smooth fiber bundle, with  $Ker(d\pi_i) = \mathcal{V}M_i$ , which is  
 116 therefore an integrable distribution.  $\mathcal{V}M_i$  is the vertical bundle of  $M_i$ . Every class of equivalence  
 117  $[p]$  is a path-connected submanifold of  $M_i$  and coincide with the fiber of the bundle over the point  
 118  $p \in M_i$ .

### 119 3 Methodology

120 The results depicted in Section 2.1 provide powerful tools for investigating how a neural network  
 121 sees the input space starting from a point  $x$ . In particular we point out the following remarks: i) If  
 122 two points  $x, y$  belonging to the input manifold  $M_0$  are such that  $x \sim_0 y$ , then  $\mathcal{N}(x) = \mathcal{N}(y)$ ; ii)  
 123 given a point  $p \in M_n$ , the counterimage  $\mathcal{N}^{-1}(p)$  is a smooth manifold, whose connected components  
 124 are classes of equivalences in  $M_0$  with respect to  $\sim_0$ . A necessary condition for two points  $x, y \in M_0$   
 125 to be in the same class of equivalence is that  $\mathcal{N}(x) = \mathcal{N}(y)$ ; iii) any class of equivalence  $[x]$ ,  $x \in M_0$ ,  
 126 is a maximal integral submanifold of  $\mathcal{V}M_0$ . The above observations directly provide with a strategy  
 127 to build up the equivalence class of an input point  $x \in M_0$ . Proposition 5 tells us that  $\mathcal{V}M_0$  is an  
 128 integrable distribution, with dimension equal to the dimension of the kernel of  $g^{(0)}$ : we can hence find  
 129  $dim(Ker(g^{(0)}))$  vector fields which are a base for the tangent space of  $M_0$ . This means that we can  
 130 compute the eigenvalue decomposition of  $g_x^{(0)}$  and consider the  $L$  linearly independent eigenvectors,  
 131 namely  $\{v_l\}_{l=1, \dots, L}$ , associated to the null eigenvalue: these eigenvectors depend *smoothly* on the  
 132 point, a fact that is not trivial when the matrix associated to the metric depends on several parameters  
 133 [15]. We can build then all the null curves by randomly selecting one eigenvector  $\tilde{v} \in \{v_l\}$  and then  
 134 reconstruct the curve along the direction  $\tilde{v}$  from the starting point  $x$ . From a practical point of view,  
 135 one is led to solve the Cauchy problem, a first order differential equation, with  $\dot{\gamma} = \tilde{v}$  and initial  
 136 condition  $\gamma(0) = x$ .

#### 137 3.1 Input Space Exploration

138 This whole procedure is coded in the Singular Metric Equivalence Class (SiMEC) and the Singular  
 139 Metric Exploration (SiMExp) algorithms, whose general schemes are depicted in Algorithms 1 and 2.  
 140 SiMEC reconstructs the class of equivalence of the input via the exploration of the input space by  
 141 randomly selecting one of the eigenvectors related to the zero eigenvalue. On the opposite, in SiMExp,  
 142 in order to move from a class of equivalence to another we consider the eigenvectors relative to the  
 143 nonzero eigenvalues. This requires the slight difference in lines 5 to 7 between Algorithm 1 and  
 144 Algorithm 2.

145

---

**Algorithm 1** The Singular Metric Equivalence Class (SiMEC) algorithm.

---

- 1: Set the network  $\mathcal{N}$ ; choose the maximum number of iterations. Choose the input  $p_0$ .
  - 2: **for**  $k = 0, 1, \dots, K - 1$  **do**
  - 3:   Compute  $g_{\mathcal{N}(p_k)}^n$
  - 4:   Compute the pullback metric  $g_{p_k}^0$
  - 5:   Diagonalize  $g_{p_k}^0$  and find the eigenvectors  $\{v_l\}_l$  associated to the zero eigenvalue
  - 6:   Randomly select  $\tilde{v} \in \{v_l\}_l$
  - 7:    $\delta = 1/\sqrt{\max(\text{eigenvalues of } g_{p_k}^0)}$
  - 8:    $p_{k+1} \leftarrow p_k + \delta \tilde{v}$
  - 9: **end for**
  - 10: Optionally: store  $\{p_k\}_{k=0, \dots, K}$  for optimizing future computations
  - 11: Project  $p_k$  to the nearest feasible region
- 

---

**Algorithm 2** The Singular Metric Exploration (SiMExp) algorithm.

---

- 1: Set the network  $\mathcal{N}$ ; choose the maximum number of iterations. Choose the input  $p_0$ .
  - 2: **for**  $k = 0, 1, \dots, K - 1$  **do**
  - 3:   Compute  $g_{\mathcal{N}(p_k)}^n$
  - 4:   Compute the pullback metric  $g_{p_k}^0$
  - 5:   Diagonalize  $g_{p_k}^0$  and find the eigenvectors  $\{w_l\}_l$  associated to the non-zero eigenvalue
  - 6:   Randomly select  $\tilde{w} \in \{w_l\}_l$
  - 7:    $\delta = 2/\sqrt{\max(\text{eigenvalues of } g_{p_k}^0)}$
  - 8:    $p_{k+1} \leftarrow p_k + \delta \tilde{w}$
  - 9: **end for**
  - 10: Optionally: store  $\{p_k\}_{k=0, \dots, K}$  for optimizing future computations
  - 11: Project  $p_k$  to the nearest feasible region
- 

146 There are some remarks to point out. From a numerical point of view, the diagonalization of the  
 147 pullback may lead to have even negative eigenvalues: hence one may use the notion of energy of  
 148 a curve, related to the pseudolength. The update rule for the new point (line 8) amounts to solve  
 149 the differential problem via the Euler method: for a reliable solution, we suggest to choose a small  
 150 step-length  $\delta$ . On the other hand, if the value of  $\delta$  is too small more iterations are needed to move

151 away from the starting point sensibly. Therefore there is a trade-off between the reliability of the  
 152 solution and the exploration pace. The proof of the well-posedness theorem for Cauchy problems,  
 153 cf. [18, Theorem 2.1], yields some insights, suggesting to set  $\delta$  equal to the inverse of the Lipschitz  
 154 constant of the map  $\mathcal{N}$  – which in practice we can estimate with the inverse of the square root of the  
 155 largest eigenvalue  $\lambda_M$  of the pullback metric  $g_{p_k}^0$ . This is our default choice for Algorithm 1. We also  
 156 note that Algorithm 1 is more sensitive to the choice of the parameter  $\delta$  compared to Algorithm 2.  
 157 To build points in the same equivalence class Algorithm 1 needs to follow a null curve closely with  
 158 as little approximations as possible, namely with a small  $\delta$ . In contrast Algorithm 2, whose goal is  
 159 to change the equivalence class from one iteration to the next, does not have the same problem and  
 160 larger  $\delta$  are allowed. Our default choice is therefore to set  $\delta = 2\lambda_M^{-1/2}$  for Algorithm 2. As for the  
 161 computational complexity of the two algorithms, the most demanding step is the computation of the  
 162 eigenvalues and eigenvectors, which is  $O(n^3)$ , with  $n$  the dimension of the square matrix  $g_{p_k}^0$  [20].  
 163 Since all the other operations are either  $O(n)$  or  $O(n^2)$ , we conclude that the complexity of both  
 164 Algorithms 1 and 2 is  $O(n^3)$ .

### 165 3.2 Interpretability

166 Algorithms 1 and 2 allow for the exploration of the equivalence classes in the input space of a Trans-  
 167 former model. However, the points explored by these algorithms may not be directly interpretable  
 168 by a human perspective. For instance, an image or a piece of text may need to be decoded to be  
 169 “readable” by a human observer. Furthermore, we present an interpretation of the eigenvalues of the  
 170 pullback metric which allows us to define a feature importance metric. We present two interpretability  
 171 methods for Transformers based on input space exploration. Both methods are then demonstrated on  
 172 a Vision Transformer (ViT) trained for digit classification [8], and two BERT models, one trained for  
 173 hate speech classification and the other trained for MLM [7, 19].

174

---

**Algorithm 3** Feature Importance Analysis Using Pull-  
back Metric  $g_{x^e}^0$

---

- 1: **Inputs:**
  - 2: Transformer model  $T$  with: Tokenizer  $t_T$ , Em-  
bedding layer  $e_T$ , Intermediate layers  $l_T$
  - 3: Input data  $x$
  - 4: Tokenize input  $x$  to obtain tokens  $x^t = t_T(x)$
  - 5: Compute embeddings  $x^e = e_T(x^t)$
  - 6: Compute intermediate representations  $g_{l_T}^n(x^e)$
  - 7: Calculate the pullback metric  $g_{x^e}^0$
  - 8: Diagonalize  $g_{x^e}^0$  to extract eigenvalues
  - 9: Identify the maximum eigenvalue for each embed-  
ding, indicating its importance
  - 10: **Output:** Heatmap of embedding importance based  
on the eigenvalues
- 

---

**Algorithm 4** Exploration of Embedding Space in  
Transformers

---

- 1: **Inputs:**
  - 2: Transformer model  $T$  with: Tokenizer  $t_T$ , Em-  
bedding layer  $e_T$ , Intermediate layers  $l_T$
  - 3: Input data  $x$  (image or text)
  - 4: Retrieve segments  $x^t = t_T(x)$ .
  - 5: Choose segments  $P = \{p | p \in x^t\}$  for updates;  
keep others unchanged.
  - 6: Compute embeddings  $x^e = e_T(x^t)$ .
  - 7: Apply SiMEC or SiMExp on  $x^e$ , updating embed-  
dings for segments in  $P$ .
  - 8: **Outputs:** Modified input embedding, one for each  
SiMEC/SiMExp iteration.
- 

175 **Feature importance.** Consider a Transformer model  $T$  whose architecture includes a tokenizer  $t_T$   
 176 (or patcher for images) that segments the input so that each segment can be converted into a continuous  
 177 representation by an embedding layer  $e_T$ . This results in a matrix of dimensions  $n_s \times h$ , where  $n_s$   
 178 represents the number of segments, and  $h$  denotes the hidden size of the model’s embeddings. The  
 179 eigenvalues of the pullback metric can be used to deduce the importance of each embedding and, by  
 180 extension, the significance of the segments they represent, with respect to the final prediction. The  
 181 process for determining the importance of textual tokens or image patches is outlined in Algorithm 3.  
 182 The appearance of the resulting heatmaps varies according to the type of input used. An example  
 183 of experiments with ViT on the MNIST dataset [12] is shown in Figure 1 that depicts heatmaps for  
 184 two MNIST instances. Figure 2, on the left, illustrates two experiment using Algorithm 3 on both a  
 185 BERT model for hate speech detection and a BERT model for MLM.

186 **Interpretation of input space exploration.** Using SiMEC and SiMExp to explore the embedding  
 187 space reveals how Transformer models perceive equivalence among different data points. Specifically,  
 188 these methodologies facilitate the sequential acquisition of embedding matrices  $p_0 \dots p_K$  at each  
 189 iteration, as detailed in Algorithms 1 and 2. Algorithm 4 implements a practical application of the

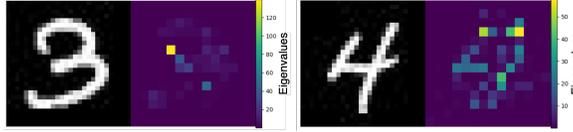


Figure 1: Example output from Algorithm 3 applied to digit classification. These two instances are predicted as 3 (left) and 4 (right). The brightness of the color indicates the eigenvalue’s magnitude. The brighter the color, the more sensitive the patch. This indicates that changes in the values of these sensitive patches are likely to have a greater impact on the prediction probabilities. Each patch in the heatmap corresponds to a  $2 \times 2$  square pixel.



Figure 2: Example outputs from Algorithm 3. The darker the color, the higher the token’s eigenvalue. Left: The sentence analysed is classified as “offensive” by the BERT for hate speech detection, with significant contributions from tokens [CLS], politicians, corrupt, and ##eit (part of the word *deceitful*). Right: Example instance processed by a BERT model for masked language modeling. [MASK] is predicted as “ham”, with the most influential tokens being pizza and cheese.

190 SiMEC/SiMExp approach with Transformer models. A key feature of this method is its ability  
 191 to selectively update specific tokens (for text inputs) or patches (for image inputs) during each  
 192 iteration. This selective updating allows us to explore targeted modifications that prompt the model  
 193 to either categorize different inputs as the same class or recognize them as distinct. Unlike traditional  
 194 approaches where modifications are predetermined, this method lets the model itself guide us to  
 195 understand which data points belong to specific equivalence classes. To interpret embeddings resulted  
 196 from the exploration process, they must be mapped back into a human-understandable form, such as  
 197 text or images. The interpretation of an embedding vector depends on the operations performed by  
 198 the Transformer’s embedding module  $e_T$ . If  $e_T$  consists only of invertible operations, it is feasible to  
 199 construct a layer that performs the inverse operation relative to  $e_T$ . The output can then be visualized  
 200 and directly interpreted by humans, allowing for a comparison with the original input to discern  
 201 how differences in embeddings reflect differences in their representations (e.g., text, images). If the  
 202 operations in  $e_T$  are non-invertible, a trained decoder is required to reconstruct an interpretable output  
 203 from each embedding matrix  $p_0 \dots p_K$ . When using a BERT model, it is feasible to utilize layers  
 204 that are specialized for the masked language modeling (MLM) task to map input embeddings back to  
 205 tokens. This approach is effective whether the BERT model in question is specifically designed for  
 206 MLM or for sentence classification. In the case of sentence classification models, it is necessary to  
 207 select a corresponding MLM BERT model that shares the same internal architecture, including the  
 208 number of layers and embedding size.

209 Algorithm 5 depicts the process of interpreting Algorithm 4 outputs for both ViT and BERT experi-  
 210 ments. After initializing the decoder according to the model type, the embeddings  $p_0 \dots p_K$  need to  
 211 be constrained to a feasible region. This region is defined by the distribution of embeddings derived  
 212 from the original input instances. Next, the embeddings are decoded, and the selected segments  
 213 for exploration are extracted. These segments are then used to replace the corresponding parts of  
 214 the original input instance. Figure 3 depicts an example outcome of Algorithm 5 applied on a ViT  
 215 exploration experiment. Given that the interpretation process includes both a capping step and a  
 216 decoding step (lines 10 and 11 of Algorithm 5), it’s important to note that there isn’t a direct 1:1  
 217 correspondence between each iteration’s update and the interpretation outcomes. Our primary focus  
 218 is on exploring the input embedding space, rather than the input image or input sentence spaces.  
 219 For further investigation, we provide a detailed discussion on considering interpretation outputs as  
 220 alternative prompts in Section 4.

---

**Algorithm 5** Interpretation for Exploration results for ViT and BERT models.

---

- 1: **Inputs:**
  - 2:     Transformer model  $T$  with: Tokenizer  $t_T$ , Embedding layer  $e_T$ , Intermediate layers  $l_T$
  - 3:     Modified embeddings  $p_0 \dots p_K$  resulted from Algorithm 4 applied on an input  $x$
  - 4:      $P = \{p|p \in x^t\}$  indices of updated segments
  - 5: **If  $T$  is ViT:**
  - 6:     Initialize decoder  $d$  with weights from  $e_T$ .
  - 7: **If  $T$  is BERT:**
  - 8:     Initialize decoder with intermediate and final layers of a BERT for MLM task.
  - 9:     Compute embeddings distributions for original input data
  - 10:     Use the original embeddings distributions to cap  $p_0 \dots p_K$
  - 11:     Decode modified embeddings  $p_0 \dots p_K$  using  $d$  to generate the corresponding images/sentences  $X' = x'_0 \dots x'_K$ .
  - 12: **For each  $x' \in X'$ :** replace segments relative to indices  $P$  in  $x$  with those in  $x'$ .
  - 13: **Outputs:**
  - 14:     Modified input images/sentences, one for each SiMEC/SiMExp iteration.
- 

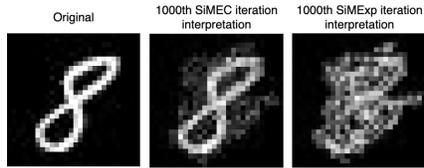


Figure 3: Example of SiMEC and SiMExp output interpretation for ViT digit classification. Left: Original MNIST image of an “8”. Center: Interpretation of a  $p_{1000}$  from a SiMEC experiment, where  $p_{1000}$  is predicted as “8”. Right: Interpretation of a  $p_{1000}$  from a SiMExp experiment, where  $p_{1000}$  is predicted as “4”. All patches are subject to SiMEC and SiMExp updates.

## 221 4 Experiments

222 Experiments are conducted on textual and visual data. We aim to perform a preliminary investigation  
223 of 3 features of our approach: (i) how the class probability changes on the decoded output of  
224 SiMEC/SiMExp, (ii) what is the trade-off between the quantity and the quality of the output, and (iii)  
225 how our method can be used to extract feature importance-based explanations.

226 In the textual case, we experiment with hate speech classification datasets: we use HateXplain<sup>2</sup> [13],  
227 which provides a ground truth for feature importance, plus a sample of 100 hate speech sentences  
228 generated by prompting ChatGPT<sup>3</sup>, which serve purposes (i) and (ii). In the visual case, we perform  
229 experiments on MNIST [12] dataset.

230 **Using interpretation outputs as alternative prompts** An interesting investigation is to determine  
231 if our interpretation algorithm (Algorithm 5) can generate alternative prompts that stay in the same  
232 equivalence class as the original input data or move to a different one, based on SiMEC and SiMExp  
233 explorations. We test how the probability assigned to the original equivalence class by the Transformer  
234 model changes as the SiMEC and SiMExp algorithms explore the input embedding manifold.

235 For BERT experiments we generate prompts to inspect the probability distribution over the vocabulary  
236 for tokens updated by Algorithms 1 and 2. We decode the updated  $p_0 \dots p_K$  using Algorithm 5,  
237 focusing on tokens updated through the iterations. For each of these decoded tokens, we extract  
238 the top-5 scores to obtain 5 alternative tokens to replace the original ones, creating 5 alternate  
239 prompts. We then extract the prediction  $i^* = \arg \max_i y_i$  for the original sentence, which represents  
240 the output whose equivalence class we aim to explore. Finally, we classify the new prompts,  
241 obtaining the corresponding predictions  $Y = \mathbf{y}^{(0)} \dots \mathbf{y}^{(K)}$ , where each  $\mathbf{y}^{(k)} \in \mathbb{R}^N$ ,  $N$  being  
242 the number of prediction classes. We visualize the prediction trend for the  $i^*$ th value in every  
243  $\mathbf{y}^{(0)} \dots \mathbf{y}^{(K)}$  categorizing the images into two subsets: those that lead to a change in prediction  
244  $Y_c = \{\mathbf{y}^{(k)} \in Y \mid \arg \max_i y_i^{(k)} \neq i^*\}$  and those that don't  $Y_s = \{\mathbf{y}_i \in Y \mid \arg \max_i y_i^{(k)} = i^*\}$ .

---

<sup>2</sup>MIT License

<sup>3</sup>Used prompts are included in the Supplementary Materials.

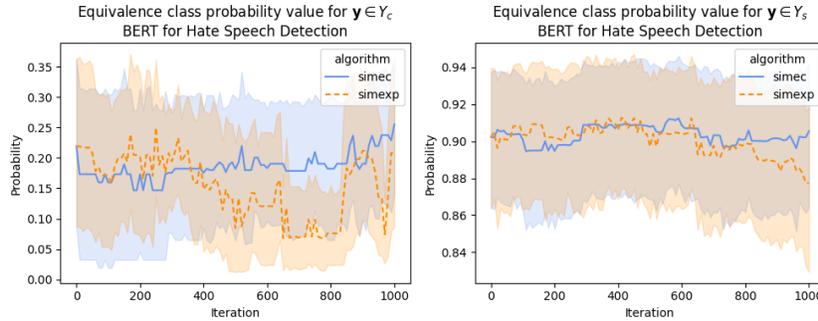


Figure 4: Analysis involving results SiMEC and SiMExp applied to BERT for hate speech detection. Left: Prediction values for  $i^*$  for each  $y \in Y_c$ . Right: Prediction values for  $y \in Y_s$ .

245 Sentence classification experiments<sup>4</sup> involved 1000 iterations from both SiMEC and SiMExp, applied  
 246 to a subset of 8 sentences from the ChatGPT hate speech dataset. The plot on the left side of Figure 4  
 247 illustrates that, as the original embeddings are increasingly modified, SiMExp tends to produce  
 248 alternatives with lower prediction values for  $i^*$  compared to SiMEC. Thus, even if predictions change  
 249 in SiMEC experiments, the equivalence class prediction value remains approximately constant and  
 250 higher than in SiMExp. Considering the plot on the right side of Figure 4, SiMExp identifies prompts  
 251 that lower the prediction value for  $i^*$ . ViT and MLM experiments are detailed in the Supplementary  
 252 Materials.

253 **Input space exploration** We measure the time required to explore the input space of a ViT with  
 254 the SiMEC algorithm and compare it with a perturbation-based method. The perturbation-based  
 255 method mimics a trial-and-error approach as it takes an input image and, at each iteration, perturbs  
 256 it by a semi-random vector  $\mathbf{v}_{t+1} = a_t \mathbf{v}_t + \eta \epsilon$ , where  $a_t = 1$  if  $y_t = y_{t-1}$ ,  $a_t = -1$  otherwise,  $\epsilon$  is  
 257 an orthogonal random vector from a standard normal distribution and  $\eta$  is the step length. With the  
 258 perturbation, we obtain a new image, then check whether the model yields the same label for the  
 259 new image. The perturbation vector is re-initialized at random from a normal distribution 20% of the  
 260 times to allow for exploration. We construct this method to have a direct comparison with ours in the  
 261 absence of a consolidated literature about the task.

262 We train a ViT model having 4 layers and 4 heads per layer on the MNIST dataset<sup>5</sup>. The SiMEC  
 263 algorithm is run for 1000 iterations, so that it can generate 1000 examples starting from a single  
 264 image. In a sample of 100 images, the average time is approximately 339 seconds.<sup>6</sup> In the same  
 265 time, the perturbation-based algorithm can produce up to 36000 images. However, we notice that  
 266 the perturbation-based algorithm ends up producing monochrome (pixel color has zero variance) or  
 267 totally noisy images, which provide little information about the behavior of the model. Excluding  
 268 only the images with low color variance ( $< 0.01$ ), we are left, on average, with 19 images (standard  
 269 deviation 13.9). SiMEC, in contrast, doesn't present this behavior, as all 1000 images have high  
 270 enough intensity variance and are thus useful for explainability purposes.

271 As BERT has many more parameters with respect to our ViT model, processing textual data takes  
 272 longer. Specifically, in a sample of 16 sentences, the average time needed to run 1000 iterations on a  
 273 sentence is 7089 seconds, taking into account both MLM and classification experiments.

274 **Feature importance-based explanations** We compare our method against Attention Rollout  
 275 (AR) [1] and the Relevancy method proposed by Chefer et al. [6]. In the textual case, we provide a  
 276 quantitative evaluation using the HateXplain dataset, which contains 20147 sentences (of which 1924  
 277 in the test set) annotated with *normal*, *offensive* and *hate speech* labels as well as the positions of  
 278 words that support the label decision. We then measure the cosine similarity between the importance  
 279 assigned by each method to each word in a sentence and the ground truth. Notice that, since the

<sup>4</sup>Model used: [huggingface.co/ctoramam/hate-speech-bert](https://huggingface.co/ctoramam/hate-speech-bert)

<sup>5</sup>Using Adam optimizer, the model achieved the highest validation accuracy (96.25%) in 20 epochs.

<sup>6</sup>All experiments are based on the current PyTorch implementation of the algorithms and run on a Ubuntu 20.04 machine endowed with one NVIDIA A100 GPU and CUDA 12.4.

280 dataset contains multiple annotations, the ground truth  $y$  for each word  $w$  is obtained as the average  
281 of the binary labels assigned by each annotator, and therefore  $y(w) \in [0; 1]$ . We also normalize all  
282 scores in  $[0; 1]$  so to have them on the same scale. The average similarity achieved by our method is  
283 **0.707** (standard deviation  $\sigma = 0.302$ ), against **0.7** ( $\sigma = 0.315$ ) for Relevancy and **0.583** ( $\sigma = 0.318$ ) for  
284 AR. This proves our method to be more effective in finding the most sensitive tokens for classification.  
285 We provide an example on image classification in the Supplementary Materials.

## 286 5 Related work

287 Our work relates to embedding space exploration literature, and has at least one collateral applications  
288 in the XAI domain, namely producing feature importance-based explanations.

289 **Embedding space exploration.** Works dealing with embedding space exploration mostly focus  
290 on the study of specific properties of the embedding space of Transformers, especially in NLP. For  
291 instance, Cai et al. [5] challenge the idea that the embedding space is inherently anisotropic [10]  
292 discovering local isotropy, and find low-dimensional manifold structures in the embedding space  
293 of GPT and BERT. Biś et al. [3] argue that the anisotropy of the embedding space derives from  
294 embeddings shifting in common directions during training. In the field of CV, Vilas et al. [21] map  
295 internal representations of a ViT onto the output class manifold, enabling the early identification of  
296 class-related patches and the computation of saliency maps on the input image for each layer and  
297 head. Applying Singular Value Decomposition to the Jacobian matrix of a ViT, Salman et al. [17]  
298 treat the input space as the union of two subspaces: one in which image embedding doesn't change,  
299 and another one for which it changes. Except for the last one, all the aforementioned approaches rely  
300 on data samples. By studying the inverse image of the model, instead, we can do away with data  
301 samples.

302 **Feature importance-based explanations.** Feature importance is a measure of the contribution of  
303 each data feature to a model prediction. In the context of Computer Vision and Natural Language  
304 Processing, it amounts to giving a weight to pixels (or patches of pixels) in an image and tokens  
305 in a piece of text, respectively. In recent years, much research has focused on Transformers in  
306 both CV and NLP. Most approaches are based on the attention mechanism of the Transformer  
307 architecture. Abnar and Zuidema [1] quantify the overall attention of the output on the input by  
308 computing a linear combination of layer attentions (Attention Rollout) or applying a maximum  
309 flow algorithm (Attention Flow). To overcome the limitations [4] of attention-based methods, Hao  
310 et al. [11] use the concept of *attribution*, which is obtained by multiplying attention matrices by  
311 the integrated gradient of the model with respect to them. Chefer et al. [6] propose the Relevancy  
312 metric to generalize attribution to bi-modal and encoder-decoder architectures. Other methods are  
313 perturbation-based, where perturbations of input data are used to record any change in the output and  
314 draw a saliency map on the input. In order to overcome the main issue with such methods, i.e. the  
315 generation of outlier inputs, Englebort et al. [9] apply perturbations after the position encoding of the  
316 patches. In contrast with these methods, ours does not need arbitrary perturbations of inputs, and  
317 considers all parameters of the model, not only the attention query and key matrices.

## 318 6 Conclusions

319 Our exploration of the Transformer architecture through a theoretical framework grounded in Rie-  
320 mannian Geometry led to the application of our two algorithms, SiMEC and SiMExp, for examining  
321 equivalence classes in the Transformers' input space. We demonstrated how the results of these explo-  
322 ration methods can be interpreted in a human-readable form and conducted preliminary investigations  
323 into their potential applications. Notably, our methods show promise for ranking feature importance  
324 and generating alternative prompts within the same or different equivalence classes.

325 Future research directions include expanding our experimental results and delving deeper into the po-  
326 tential of our framework for controlled input generation within an equivalence class. This application  
327 holds significant promise for enhancing the explainability of Transformer models' decisions and for  
328 addressing issues related to bias and hallucinations.

## References

- [1] S. Abnar and W. Zuidema. Quantifying attention flow in transformers. 2020. arXiv:2005.00928.
- [2] A. Benfenati and A. Marta. A singular Riemannian geometry approach to Deep Neural Networks I. Theoretical foundations. *Neural Networks*, 158:331–343, 2023.
- [3] D. Biś, M. Podkorytov, and X. Liu. Too much in common: Shifting of embeddings in transformer language models and its implications. In *Proceedings of the 2021 conference of the North American chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5117–5130, 2021.
- [4] G. Brunner, Y. Liu, D. Pascual, O. Richter, M. Ciaramita, and R. Wattenhofer. On identifiability in transformers. *International Conference on Learning Representations*, 2019.
- [5] X. Cai, J. Huang, Y. Bian, and K. Church. Isotropy in the contextual embedding space: Clusters and manifolds. In *International conference on learning representations*, 2020.
- [6] H. Chefer, S. Gur, and L. Wolf. Generic attention-model explainability for interpreting bi-modal and encoder-decoder transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 397–406, 2021.
- [7] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In J. Burstein, C. Doran, and T. Solorio, editors, *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL <https://aclanthology.org/N19-1423>.
- [8] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale, 2021. arXiv:2010.11929.
- [9] A. Englebert, S. Stassin, G. Nanfack, S. A. Mahmoudi, X. Siebert, O. Cornu, and C. De Vleeschouwer. Explaining through transformer input sampling. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 806–815, 2023.
- [10] J. Gao, D. He, X. Tan, T. Qin, L. Wang, and T.-Y. Liu. Representation degeneration problem in training natural language generation models. In *International conference on learning representations*, volume abs/1907.12009, 2019. URL <https://api.semanticscholar.org/CorpusID:59317065>.
- [11] Y. Hao, L. Dong, F. Wei, and K. Xu. Self-attention attribution: Interpreting information interactions inside transformer. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 12963–12971, 2021.
- [12] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. doi: 10.1109/5.726791.
- [13] B. Mathew, P. Saha, S. M. Yimam, C. Biemann, P. Goyal, and A. Mukherjee. Hatexplain: A benchmark dataset for explainable hate speech detection. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pages 14867–14875, 2021.
- [14] N. Papernot, P. McDaniel, S. Jha, M. Fredrikson, Z. B. Celik, and A. Swami. The Limitations of Deep Learning in Adversarial Settings. In *2016 IEEE European Symposium on Security and Privacy (EuroS&P)*, pages 372–387, Mar. 2016. doi: 10.1109/EuroSP.2016.36.
- [15] F. Rellich and J. Berkowitz. *Perturbation Theory of Eigenvalue Problems*. New York University. Institute of Mathematical Sciences. Gordon and Breach, 1969. ISBN 9780677006802.
- [16] M. T. Ribeiro, S. Singh, and C. Guestrin. "why should i trust you?" explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1135–1144, 2016.

- 376 [17] S. Salman, M. M. B. Shams, and X. Liu. Intriguing equivalence structures of the embedding  
377 space of vision transformers, 2024. arXiv:2401.15568.
- 378 [18] M. E. Taylor. *Partial differential equations. I: Basic theory*, volume 115 of *Appl. Math. Sci.*  
379 Cham: Springer, 3rd corrected and expanded edition edition, 2023. ISBN 978-3-031-33858-8;  
380 978-3-031-33861-8; 978-3-031-33859-5. doi: 10.1007/978-3-031-33859-5.
- 381 [19] C. Toraman, F. Şahinuç, and E. H. Yilmaz. Large-scale hate speech detection with cross-  
382 domain transfer. In *Proceedings of the Language Resources and Evaluation Conference*, pages  
383 2215–2225, Marseille, France, June 2022. European Language Resources Association. URL  
384 <https://aclanthology.org/2022.lrec-1.238>.
- 385 [20] L. N. Trefethen and D. I. Bau. *Numerical linear algebra. Twenty-fifth anniversary edition*,  
386 volume 181 of *Other Titles Appl. Math.* Philadelphia, PA: Society for Industrial and Applied  
387 Mathematics (SIAM), 2022. ISBN 978-1-61197-715-8.
- 388 [21] M. G. Vilas, T. Schaumlöffel, and G. Roig. Analyzing vision transformers for image classi-  
389 fication in class embedding space. *Advances in Neural Information Processing Systems*, 36,  
390 2024.
- 391 [22] M. Wu, H. Wu, and C. Barrett. Verix: Towards verified explainability of deep neural networks.  
392 *Advances in neural information processing systems*, 36, 2024.

## 393 **NeurIPS Paper Checklist**

### 394 **1. Claims**

395 Question: Do the main claims made in the abstract and introduction accurately reflect the  
396 paper's contributions and scope?

397 Answer: [\[Yes\]](#)

398 Justification: In the abstract and introduction we claim that we present a method for the  
399 exploration of equivalence classes in the input space of Transformer models, which is  
400 analyzed in depth in Section 3. The mathematical theory we refer to is deepened in Section  
401 2.

402 Guidelines:

- 403 • The answer NA means that the abstract and introduction do not include the claims  
404 made in the paper.
- 405 • The abstract and/or introduction should clearly state the claims made, including the  
406 contributions made in the paper and important assumptions and limitations. A No or  
407 NA answer to this question will not be perceived well by the reviewers.
- 408 • The claims made should match theoretical and experimental results, and reflect how  
409 much the results can be expected to generalize to other settings.
- 410 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
411 are not attained by the paper.

### 412 **2. Limitations**

413 Question: Does the paper discuss the limitations of the work performed by the authors?

414 Answer: [\[Yes\]](#)

415 Justification: We discuss the limitations and tradeoff given by numeric integration in  
416 Subsection 3.1 and theoretical assumptions are enumerated in Section 2. Computational  
417 efficiency of our algorithms is discussed in Subsection 3.1. We conducted experiments on 3  
418 datasets only, one of which of small dimensions since our main focus is on the mathematical  
419 theory grounding the application of the method to Transformers, as stated at the beginning  
420 of Section 4. Other limitations are mentioned throughout Section 4, including the fact that  
421 our investigations in the human-readable scenario are at a preliminary stage.

422 Guidelines:

- 423 • The answer NA means that the paper has no limitation while the answer No means that  
424 the paper has limitations, but those are not discussed in the paper.
- 425 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 426 • The paper should point out any strong assumptions and how robust the results are to  
427 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
428 model well-specification, asymptotic approximations only holding locally). The authors  
429 should reflect on how these assumptions might be violated in practice and what the  
430 implications would be.
- 431 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
432 only tested on a few datasets or with a few runs. In general, empirical results often  
433 depend on implicit assumptions, which should be articulated.
- 434 • The authors should reflect on the factors that influence the performance of the approach.  
435 For example, a facial recognition algorithm may perform poorly when image resolution  
436 is low or images are taken in low lighting. Or a speech-to-text system might not be  
437 used reliably to provide closed captions for online lectures because it fails to handle  
438 technical jargon.
- 439 • The authors should discuss the computational efficiency of the proposed algorithms  
440 and how they scale with dataset size.
- 441 • If applicable, the authors should discuss possible limitations of their approach to  
442 address problems of privacy and fairness.
- 443 • While the authors might fear that complete honesty about limitations might be used by  
444 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
445 limitations that aren't acknowledged in the paper. The authors should use their best

446 judgment and recognize that individual actions in favor of transparency play an impor-  
447 tant role in developing norms that preserve the integrity of the community. Reviewers  
448 will be specifically instructed to not penalize honesty concerning limitations.

### 449 3. Theory Assumptions and Proofs

450 Question: For each theoretical result, does the paper provide the full set of assumptions and  
451 a complete (and correct) proof?

452 Answer: [Yes]

453 Justification: The full proofs are part of two previously published papers which we cannot  
454 disclose for anonymity requirements. We replicate the relevant proofs in the supplementary  
455 material, part of which will be removed from the final version of the paper, referencing to  
456 the other papers.

457 Guidelines:

- 458 • The answer NA means that the paper does not include theoretical results.
- 459 • All the theorems, formulas, and proofs in the paper should be numbered and cross-  
460 referenced.
- 461 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 462 • The proofs can either appear in the main paper or the supplemental material, but if  
463 they appear in the supplemental material, the authors are encouraged to provide a short  
464 proof sketch to provide intuition.
- 465 • Inversely, any informal proof provided in the core of the paper should be complemented  
466 by formal proofs provided in appendix or supplemental material.
- 467 • Theorems and Lemmas that the proof relies upon should be properly referenced.

### 468 4. Experimental Result Reproducibility

469 Question: Does the paper fully disclose all the information needed to reproduce the main ex-  
470 perimental results of the paper to the extent that it affects the main claims and/or conclusions  
471 of the paper (regardless of whether the code and data are provided or not)?

472 Answer: [Yes]

473 Justification: Pseudo-code of the proposed algorithms is reported in Subsections 3.1 and 3.2  
474 so to make the algorithms reproducible, plus our implementation is made available in the  
475 supplementary material. Experiments, including the complete setting, and the respective  
476 baselines are described in Section 4.

477 Guidelines:

- 478 • The answer NA means that the paper does not include experiments.
- 479 • If the paper includes experiments, a No answer to this question will not be perceived  
480 well by the reviewers: Making the paper reproducible is important, regardless of  
481 whether the code and data are provided or not.
- 482 • If the contribution is a dataset and/or model, the authors should describe the steps taken  
483 to make their results reproducible or verifiable.
- 484 • Depending on the contribution, reproducibility can be accomplished in various ways.  
485 For example, if the contribution is a novel architecture, describing the architecture fully  
486 might suffice, or if the contribution is a specific model and empirical evaluation, it may  
487 be necessary to either make it possible for others to replicate the model with the same  
488 dataset, or provide access to the model. In general, releasing code and data is often  
489 one good way to accomplish this, but reproducibility can also be provided via detailed  
490 instructions for how to replicate the results, access to a hosted model (e.g., in the case  
491 of a large language model), releasing of a model checkpoint, or other means that are  
492 appropriate to the research performed.
- 493 • While NeurIPS does not require releasing code, the conference does require all submis-  
494 sions to provide some reasonable avenue for reproducibility, which may depend on the  
495 nature of the contribution. For example
  - 496 (a) If the contribution is primarily a new algorithm, the paper should make it clear how  
497 to reproduce that algorithm.
  - 498 (b) If the contribution is primarily a new model architecture, the paper should describe  
499 the architecture clearly and fully.

- 500 (c) If the contribution is a new model (e.g., a large language model), then there should  
501 either be a way to access this model for reproducing the results or a way to reproduce  
502 the model (e.g., with an open-source dataset or instructions for how to construct  
503 the dataset).
- 504 (d) We recognize that reproducibility may be tricky in some cases, in which case  
505 authors are welcome to describe the particular way they provide for reproducibility.  
506 In the case of closed-source models, it may be that access to the model is limited in  
507 some way (e.g., to registered users), but it should be possible for other researchers  
508 to have some path to reproducing or verifying the results.

## 509 5. Open access to data and code

510 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
511 tions to faithfully reproduce the main experimental results, as described in supplemental  
512 material?

513 Answer: [Yes]

514 Justification: All experiments are made reproducible through scripts provided as supplemen-  
515 tary material.

516 Guidelines:

- 517 • The answer NA means that paper does not include experiments requiring code.
- 518 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/  
519 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 520 • While we encourage the release of code and data, we understand that this might not be  
521 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not  
522 including code, unless this is central to the contribution (e.g., for a new open-source  
523 benchmark).
- 524 • The instructions should contain the exact command and environment needed to run to  
525 reproduce the results. See the NeurIPS code and data submission guidelines ([https://  
526 nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 527 • The authors should provide instructions on data access and preparation, including how  
528 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 529 • The authors should provide scripts to reproduce all experimental results for the new  
530 proposed method and baselines. If only a subset of experiments are reproducible, they  
531 should state which ones are omitted from the script and why.
- 532 • At submission time, to preserve anonymity, the authors should release anonymized  
533 versions (if applicable).
- 534 • Providing as much information as possible in supplemental material (appended to the  
535 paper) is recommended, but including URLs to data and code is permitted.

## 536 6. Experimental Setting/Details

537 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-  
538 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the  
539 results?

540 Answer: [Yes]

541 Justification: All details are provided in Section 4: details of the analyzed architectures,  
542 number of iterations of the SiMEC/SiMExp algorithms, technical infrastructure on which  
543 the experiments were performed, amount of data the experiments were performed on.

544 Guidelines:

- 545 • The answer NA means that the paper does not include experiments.
- 546 • The experimental setting should be presented in the core of the paper to a level of detail  
547 that is necessary to appreciate the results and make sense of them.
- 548 • The full details can be provided either with the code, in appendix, or as supplemental  
549 material.

## 550 7. Experiment Statistical Significance

551 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
552 information about the statistical significance of the experiments?

553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604

Answer: [Yes]

Justification: The standard deviation is reported for the experiment that supports the claims of the paper, i.e. the one on feature importance-based explanations, in Section 4. Standard deviation is also reported for the number of uninformative images produced by the perturbation-based baseline method in the Input space exploration experiment.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

## 8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: All the experiments were performed on the same infrastructure, which is reported in a footnote in Section 4. Time of execution is one of the key indicators reported for the Input space exploration experiments. More computing power would be required for experiments on bigger Transformer models.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

## 9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

Answer: [Yes]

Justification: We comply with the terms of use of the datasets employed in the experiments, and we deem our work has no potentially harmful effect on people safety, security, discrimination, surveillance, harassment, nor on human rights. Our proposal does not contribute to spread bias and unfairness towards certain groups of people nor to harm the environment.

605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

**10. Broader Impacts**

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [No]

Justification: Although the impacts of XAI on society is broad and deep, in this paper we focus only on the technical problem of exploring the equivalence classes in the input space of Transformers, which doesn't add any specific impact to the discussion about XAI in general.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

**11. Safeguards**

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: The paper poses no such risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
710

## 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: The datasets used in the paper are explicitly mentioned in the references, as required by the terms of use. Where applicable, the license is also reported.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

## 13. New Assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: The paper does not release new assets.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

## 14. Crowdsourcing and Research with Human Subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: Our work doesn't include crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

711 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human**  
712 **Subjects**

713 Question: Does the paper describe potential risks incurred by study participants, whether  
714 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
715 approvals (or an equivalent approval/review based on the requirements of your country or  
716 institution) were obtained?

717 Answer: [NA]

718 Justification: Our work does not involve crowdsourcing nor research with human subjects.

719 Guidelines:

- 720 • The answer NA means that the paper does not involve crowdsourcing nor research with  
721 human subjects.
- 722 • Depending on the country in which research is conducted, IRB approval (or equivalent)  
723 may be required for any human subjects research. If you obtained IRB approval, you  
724 should clearly state this in the paper.
- 725 • We recognize that the procedures for this may vary significantly between institutions  
726 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
727 guidelines for their institution.
- 728 • For initial submissions, do not include any information that would break anonymity (if  
729 applicable), such as the institution conducting the review.