

GERA: LABEL-EFFICIENT GEOMETRICALLY REGULARIZED ALIGNMENT

Anonymous authors

Paper under double-blind review

ABSTRACT

1 Pretrained unimodal encoders incorporate rich semantic information into embed-
 2 ding space structures. To be similarly informative, multi-modal encoders typically
 3 require massive amounts of paired data for alignment and training. We introduce a
 4 semi-supervised **Geometrically Regularized Alignment** (GeRA) method to align
 5 the embedding spaces of pretrained unimodal encoders in a label-efficient way.
 6 Our method leverages the manifold geometry of unpaired (unlabeled) data to im-
 7 prove alignment performance. To prevent distortions to local geometry during the
 8 alignment process —potentially disrupting semantic neighborhood structures and
 9 causing misalignment of unobserved pairs — we introduce a geometric loss term.
 10 This term is built upon a diffusion operator that captures the local manifold geom-
 11 etry of the unimodal pretrained encoders. GeRA is modality-agnostic and thus can
 12 be used to align pretrained encoders from any data modalities. We provide em-
 13 pirical evidence to the effectiveness of our method in the domains of speech-text
 14 and image-text alignment. Our experiments demonstrate significant improvement
 15 in alignment quality compared to a variety of leading baselines, especially with a
 16 small amount of paired data, using our proposed geometric regularization.

17 1 INTRODUCTION

18 Data comes in many modalities, including text, speech, images, and video. Unimodal encoders aim
 19 to extract the intrinsic features of data drawn from a single modality, representing it in an embedding
 20 space. The goal of multi-modal learning is to learn a *shared* representation space for encoders of
 21 different modalities. In this setting, objects captured in different modalities have common represen-
 22 tations in this shared space. This task is commonly referred to as *multi-modal alignment* (Baltrušaitis
 23 et al., 2018). Finding unified representations unlocks applications that require multiple modalities,
 24 like retrieving and generating descriptions of visual content.

25 In this paper, we consider multi-modal alignment using pretrained unimodal encoders. We are given
 26 paired and unpaired multi-modal data of potentially different dimensionalities and aim to learn an
 27 alignment transformation into a common embedding space. Although the domain of image and
 28 text alignment has been extensively explored thanks to large, publicly available image-text datasets
 29 (Schuhmann et al., 2021), one quickly runs into data availability problems when looking at new
 30 modalities. Indeed, for most modality pairs, such as speech and text or protein sequences and
 31 biomedical texts (Xu et al., 2023), there are far fewer paired data points than for images and text.

32 With the scenario above in mind, we present a robust and data-efficient alignment method that
 33 generalizes to new modalities, even under limited paired data availability. Our key idea is to preserve
 34 the local geometric structure learned by the pretrained encoders (Moschella et al., 2023; Antonello
 35 et al., 2021). These geometric structures, however, are not explicitly leveraged by existing alignment
 36 methods. Specifically, learning an alignment using only a contrastive objective, as explored by
 37 Radford et al. (2021) and others, seemingly does not maintain the manifold geometry (see Figure
 38 1) and requires substantial paired data for alignment. Conversely, the Procrustes method (Gower,
 39 1975) aligns the datasets through an isometric rotation transformation and hence fully preserves the
 40 geometric structure. However, Procrustes has low plasticity.

41 Our proposed **Geometrically Regularized Alignment** (GeRA) method leverages semantically rich
 42 manifold structures and preserves local geometry, while allowing enough flexibility to learn a mean-
 43 ingful alignment. We use a regularization loss which optimizes for local geometry preservation,

44 built on a diffusion operator to capture the local
 45 local geometry. We freeze the unimodal encoders
 46 during the alignment process, reducing computa-
 47 tional costs. Our approach falls into the regime of
 48 semi-supervised learning, as we can leverage the
 49 vast amount of unpaired (unlabeled) data with rel-
 50 atively few pairs to establish alignment. See Fig-
 51 ure 2 for an overview of our method.

52 **Contributions.** Our work advances the field of
 53 data-efficient multi-modal alignment by address-
 54 ing several limitations of existing methods. Our
 55 main contributions are three-fold:

- 56 • **Geometry-Preserving Alignment:** We intro-
 57 duce a semi-supervised alignment method that
 58 aligns multi-modal data distributions while pre-
 59 serving local geometry. It exhibits both global
 60 flexibility to align the paired points and local ge-
 61 ometric preservation to incorporate the rich se-
 62 mantic information of the manifold structure.
- 63 • **Efficiency:** A key advantage of the proposed
 64 method is its label efficiency, as it employs a
 65 semi-supervised approach to use unlabeled data.
 66 This enables the alignment to capture additional
 67 information from the pretrained unimodal en-
 68 coders in regions where there are no labeled
 69 pairs.
- 70 • **Modality-Agnostic Formulation:** GeRA is ag-
 71 nostic to the choice of encoders and modalities;
 72 it does not rely on domain-specific knowledge
 73 like augmentation. We experiment across mul-
 74 tiple encoders and data modalities to show that
 75 our method is effective across configurations. It
 76 can be efficiently applied whenever pretrained
 77 models are available.

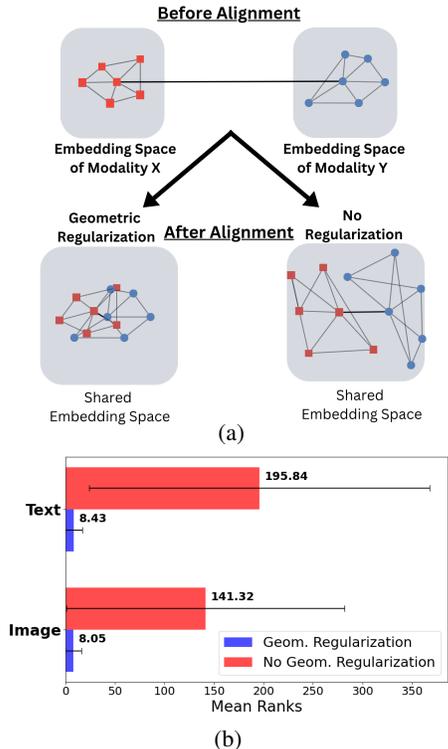


Figure 1: (a) Illustration of the effect of GeRA on alignment quality; GeRA preserves local neighborhoods, whereas non-regularized methods might distort them. Inter-modality black lines denote known pairs and gray lines denote neighbors. (b) Average ranking of the five nearest neighbors (before alignment) in the learned aligned spaces, using contrastive loss with and without our geometric regularization.

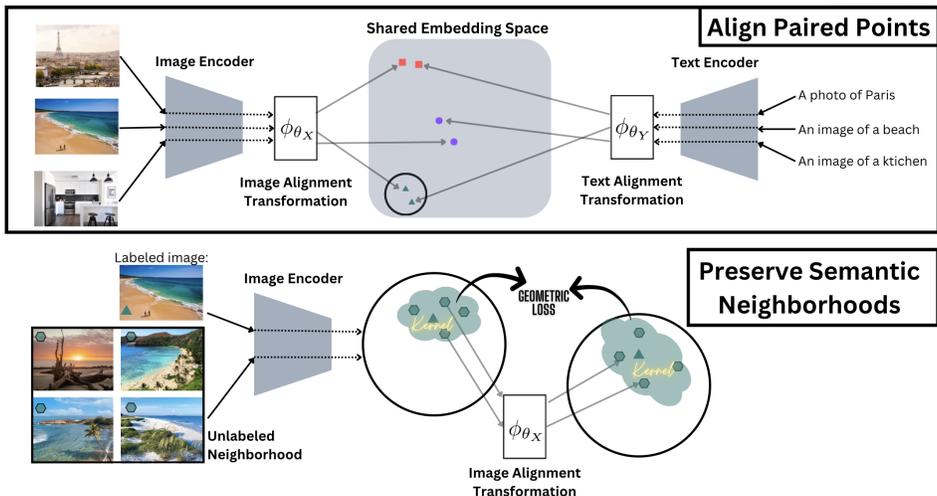


Figure 2: GeRA Training Approach: We optimize image and text alignment functions, focusing on achieving both global alignment of paired points and the preservation of local geometric structures.

78 2 RELATED WORK

79 Various multi-modal alignment methods have been introduced, each based on different assumptions
80 on data availability and computational needs; most have been applied to text and image modalities.

81 **Training Multi-Modal Encoders:** Radford et al. (2021); Chen et al. (2022); Jia et al. (2021) jointly
82 train image and text encoders from scratch, learning a shared representation for both modalities using
83 a contrastive objective (Wang & Isola, 2020). This approach outperforms many existing models
84 (Kolesnikov et al., 2020; Chen et al., 2020; He et al., 2016) in zero-shot classification on ImageNet
85 (Deng et al., 2009). These methods, however, demand large training datasets (Gadre et al., 2023;
86 Thomee et al., 2016; Sun et al., 2017) and consume significant computational resources. Zhai et al.
87 (2021) reduce computational costs by freezing the image encoder. Due to the large training datasets
88 (Sun et al., 2017; Thomee et al., 2016) this method remains computationally intensive.

89 **Relative and Anchor-based Encodings:** Moschella et al. (2023); Antonello et al. (2021) demon-
90 strate that high-quality encoders produce semantically rich and consistent manifold structures. This
91 observation suggests the concept of relative encodings, where a sample is encoded based on its
92 neighborhood. Such relative encodings have been shown to remain consistent across various en-
93 coders and modalities (Moschella et al., 2023). Building on this idea, Norelli et al. (2022) ensures
94 consistent encodings across different modalities using frozen pretrained encoders (Song et al., 2020;
95 Dosovitskiy et al., 2020), eliminating the need for a training phase. Their method achieves high
96 performance, coming close to models trained with substantially more data. However, a trade-off
97 arises: inference time increases with the number of anchor points (labeled data) used.

98 **Unsupervised Alignment Techniques:** There has been also research efforts towards unsupervised
99 alignment of embedding spaces without relying on paired modality data. The study Alvarez-Melis
100 & Jaakkola (2018) uses the Gromov-Wasserstein optimal transport objective (Nekrashevich et al.,
101 2023) to align word embeddings from various languages. Despite its advantage of not requiring
102 labeled data, the method poorly scales to the 4th power in terms of the number of embedding points.

103 **Manifold Geometry:** Early works in manifold learning, such as Locally Linear Embedding (LLE)
104 (Roweis & Saul, 2000), Isomap (Tenenbaum et al., 2000), and multi-dimensional scaling (Saeed
105 et al., 2018), capture geometric properties of data manifolds while mapping them into simpler
106 spaces. These methods leverage the rich local structure of datasets, constructing a lower-dimensional
107 embedding that retains the topological and geometric characteristics of local neighborhoods in high-
108 dimensional data space. In the context of semi-supervised learning, Sindhwani et al. (2005); Zhu
109 et al. (2003) propose frameworks for integrating geometry learned from both labeled and unlabeled
110 data into classification algorithms based on the graph Laplacian (Sindhwani et al., 2005), and based
111 on a Gaussian random field model (Zhu et al., 2003). These works, however, focus on a unimodal
112 setting and do not address semi-supervised alignment of data from multiple modalities.

113 3 PROBLEM FORMULATION

114 3.1 MULTI-MODAL ALIGNMENT

115 Consider two datasets $X \in \mathbb{R}^{N_X \times d_X}$ and $Y \in \mathbb{R}^{N_Y \times d_Y}$, originating from two distinct modalities.
116 Here, N_X and N_Y denote the number of samples in each dataset, and d_X and d_Y represent the
117 dimensions of X and Y , respectively, which may differ. These datasets denote points embedded
118 by pretrained unimodal encoders. We assume pairwise correspondence for only a small subset of
119 these points, denoted by $\{(x_{p_i}, y_{q_i})\}_{i=1}^M$, where $x_{p_i} \in X$, $y_{q_i} \in Y$, $p_i \in [1, N_X]$ and $q_i \in [1, N_Y]$
120 are some index permutations, and $M \ll N_X, N_Y$. Those pairings correspond to the same objects
121 captured by different modalities. All other points in X and Y are unpaired (unlabeled).

122 The task at hand is to align the data distributions into a common embedding space. While most meth-
123 ods focus only on aligning the paired data points, we propose to leverage unlabeled (unpaired) points
124 from each modality to preserve the rich geometric structure of their original embedding spaces.

125 To approach this problem, we define two trainable alignment functions, namely $\phi_{\theta_X} : \mathbb{R}^{d_X} \rightarrow \mathbb{R}^d$
126 and $\phi_{\theta_Y} : \mathbb{R}^{d_Y} \rightarrow \mathbb{R}^d$, where d is the dimension of the joint embedding space. These alignment
127 transformations are modeled as neural networks. This approach is *encoder and modality agnostic*,
128 requiring some pretrained unimodal encoder for each modality and a small paired dataset.

129 3.2 PRESERVING MANIFOLD GEOMETRY

130 Unimodal encoders, trained on large and often self-supervised datasets, learn to encode the data into
 131 a rich representation that accurately reflects the intrinsic structure of the data. When training the
 132 alignment using a smaller paired dataset, this manifold structure might be distorted, degrading the
 133 quality of the alignment. Existing methods are not required to preserve these structures, leaving a
 134 potential source of information unused. Specifically, learning an alignment using only a contrastive
 135 objective, as explored by (Radford et al., 2021), distorts the neighborhood geometry (see Figure 1b)
 136 and requires substantial paired data to learn an effective alignment.

137 We propose a geometrically regularized alignment that aligns the paired points while preserving
 138 local neighborhood structures, which is motivated by the relation between local neighborhoods and
 139 the (Riemannian) manifold geometry (e.g. in approximating geodesic distances) (Coifman & La-
 140 fon, 2006; Li & Dunson, 2019). This approach offers global flexibility for obtaining a meaningful
 141 alignment and local regularization to maintain the neighborhood structure. The geometry regulariza-
 142 tion pursues an intuitive goal of keeping similar objects close in the aligned space. This allows for
 143 more effective generalization of the learned alignment to nearby (unpaired) points, as evident from
 144 the improved performance by our approach (see Section 5.3.) For preserving local neighborhoods,
 145 unlabeled data can be leveraged, thus allowing a semi-supervised approach.

146 4 GEOMETRICALLY REGULARIZED ALIGNMENT METHOD

147 4.1 GERA LOSS FUNCTION

148 We introduce the GeRA loss, which optimizes for both aligning paired points and preserving the
 149 neighborhood structure of nearby unpaired points. This loss is semi-supervised, as it uses paired
 150 data for the alignment and captures the local geometry using both paired and unpaired data. The
 151 loss is defined as follows:

$$\begin{aligned} \mathcal{L}_{GeRA}(\theta_X, \theta_Y) = \mathbb{E}_{(X_B, Y_B) \sim P_{Pos}} & \left[\underbrace{\mathcal{L}_{Con}(X_B, Y_B; \theta_X, \theta_Y) + \mathcal{L}_{Con}(Y_B, X_B; \theta_Y, \theta_X)}_{\text{Alignment}} \right. \\ & \left. + \underbrace{\alpha \cdot (\mathcal{L}_{Geo}(X_B; \theta_X) + \mathcal{L}_{Geo}(Y_B; \theta_Y))}_{\text{Geometric Regularization}} \right] \end{aligned} \quad (1)$$

152 where θ_X and θ_Y parameterize the alignment transformations, ϕ_{θ_X} and ϕ_{θ_Y} , respectively, P_{Pos}
 153 represents the uniform distribution over all paired points from both modalities, and B represents the
 154 number of paired data points in a batch.

155 **Alignment:** We align the labeled points via a contrastive loss, denoted by $\mathcal{L}_{Con}(X_B, Y_B; \theta_X, \theta_Y)$
 156 as proposed by Radford et al. (2021). It minimizes the distance between positive pairs (paired
 157 points) while maximizing the distance of negative samples. We apply this loss to the alignment
 158 transformation outputs:

$$\mathcal{L}_{Con}(X_B, Y_B; \theta_X, \theta_Y) = -\frac{1}{2} \sum_{x \in X_B} \log \frac{\exp(\text{cossim}(\phi_{\theta_X}(x), \phi_{\theta_Y}(y))/t)}{\sum_{y \in Y_B} \exp(\text{cossim}(\phi_{\theta_X}(x), \phi_{\theta_Y}(y))/t)} \quad (2)$$

159 where t is a temperature hyperparameter.

160 **Geometric Regularization:** Our geometric loss term aims to preserve the local geometric structure:

$$\mathcal{L}_{Geo}(X_B; \theta_X) = \frac{1}{B} \sum_{x \in X_B} \mathbb{E}_{N_K(x) \sim S(x)} \left[\left\| \mathbf{W}_{N_K(x)} - \mathbf{W}_{\phi_{\theta_X}(N_K(x))} \right\|_F^2 \right]. \quad (3)$$

161 where $\mathbf{W}_{N_K(x)}$ and $\mathbf{W}_{\phi_{\theta_X}(N_K(x))}$ are some matrices encoding the neighborhood structure of sam-
 162 ple x , and $N_K(x)$ denotes a sampled set of K neighbors of x (according to the original embedded
 163 space). This loss operates only within a single modality and is independent of the other modality.
 164 For a given batch of unimodal samples X_B , we sample a set of K neighbors $N_K(x)$ (defined based
 165 on proximity in the original space) for each sample in the batch, drawn from a precomputed larger
 166 neighborhood distribution $S(x)$. We investigate various sampling methods:

- 167 • The “closest” method deterministically takes the K nearest neighbors.
- 168 • The “uniform” method samples K neighbors uniformly from the larger neighborhood.
- 169 • The “biased” method samples proximate neighbors with higher probability than distant neighbors.

170 The loss in equation 3, penalizes local distortion during the alignment, thus preserving the geometry.
 171 The choice of a suitable neighborhood encoding to capture local structure is a crucial consideration.
 172 We propose to use an approximation of the heat kernel, discussed in Section 4.2.1. Additionally, we
 173 report alternative choices and evaluate the differences in performance in Section 5.5.

174 4.2 KERNEL ENCODINGS

175 4.2.1 THE HEAT KERNEL

176 We next present the different choices of kernels to encode the neighborhood structure. Given a set
 177 of points, $\{x_i\}_{i=1}^N$, assumed to lie on some low dimensional manifold, \mathcal{X} , the diffusion operator
 178 (Coifman & Lafon, 2006), denoted by \mathbf{W}^{Heat} , is defined by:

$$\begin{aligned} \mathbf{K}^{\text{Heat}}(x_i, x_j) &= e^{-\|x_i - x_j\|_2^2 / 4\epsilon} \\ \mathbf{W}^{\text{Heat}}(x_i, x_j) &= \frac{\mathbf{K}^{\text{Heat}}(x_i, x_j)}{\sum_l \mathbf{K}^{\text{Heat}}(x_i, x_l)} \end{aligned} \quad (4)$$

179 This operator was shown to converge pointwise to the Neumann heat kernel of the underlying data
 180 manifold as ϵ approaches zero and the number of points tends to infinity. Below, we articulate some
 181 advantages of using \mathbf{W}^{Heat} in our formalism:

182 **Intrinsic.** The heat kernel and its approximation are intrinsic, meaning that they are independent of
 183 the choice of coordinates. As a result, they are invariant to isometric transformations.

Informative. The heat kernel captures essential intrinsic geometric information. For example, the
 geodesic distance g between two points x, y on a manifold can be recovered from the heat kernel via
 the limit (Varadhan, 1967):

$$g(x, y) = \lim_{t \rightarrow 0} \sqrt{-4t \log h_t(x, y)},$$

184 where $h_t(x, y)$ denotes the continuous heat kernel, which relates to \mathbf{W}^{Heat} by $h_t =$
 185 $\lim_{\epsilon \rightarrow 0, N \rightarrow \infty} (\mathbf{W}^{\text{Heat}})^{t/\epsilon}$ (under slightly different normalization) (Coifman & Lafon, 2006).

186 **Multi-Scale.** The locality of the heat kernel is sensitive to the time variable, t . In its discrete
 187 approximation, \mathbf{W}^{Heat} , the locality is governed by the kernel scale, ϵ , and the sample density of
 188 the point cloud. Through these parameters, the heat kernel and its approximation are capable of
 189 capturing multi-scale features. Specifically, a smaller ϵ in equation 4 results in a more local kernel.

190 4.2.2 ALTERNATIVE KERNEL ENCODINGS

191 The majority of our experiments use the diffusion operator to capture local neighborhood geometry,
 192 but other choices are possible. For example, the following kernels capture the pairwise L^2 distance
 193 and related values:

$$\mathbf{K}^{\text{Linear}}(x_i, x_j) = \|x_i - x_j\|_2 \quad \forall x_i, x_j \in X \quad (5)$$

$$\mathbf{K}^{\text{Squared}}(x_i, x_j) = \|x_i - x_j\|_2^2 \quad \forall x_i, x_j \in X \quad (6)$$

$$\mathbf{K}^{\text{Inverse}}(x_i, x_j) = \frac{1}{1 + \|x_i - x_j\|_2^2} \quad \forall x_i, x_j \in X \quad (7)$$

194 We normalize each kernel by the average column values, similarly to the diffusion operator, resulting
 195 in the neighborhood encoding \mathbf{W}_X^Z , where Z stands for “Linear”, “Squared” or “Inverse”:

$$\mathbf{W}^Z(x_i, x_j) = \frac{\mathbf{K}^Z(x_i, x_j)}{\sum_l \mathbf{K}^Z(x_i, x_l)} \quad (8)$$

196 In Section 5.5, we empirically demonstrate that the heat kernel yields the best performance, indicat-
 197 ing better preservation of local neighborhood information.

198 5 EXPERIMENTS

199 5.1 EXPERIMENTAL DETAILS

200 We conduct extensive experiments to show the performance of GeRA under limited paired data
201 availability with images and text. In addition, in Section 5.4 we present results with speech and text.

202 Our default experimental setup is adapted from the setup used in ASIF (Norelli et al., 2022), which
203 serves as a baseline.

204 **Dataset:** Our training dataset for the image and text experiments is the Conceptual 12M (CC12M)
205 dataset (Changpinyo et al., 2021). This dataset consists of 12 million paired entries of images and
206 their corresponding textual descriptions, spanning a broad spectrum of visual concepts.

207 **Unpaired points:** To preserve the local geometry of the pretrained unimodal models, we use un-
208 paired points from each modality to compute the geometric regularization in equation 1. For the
209 image and text experiments, we discard the pairing information of 6×10^6 data points from CC12M
210 and treat them as unpaired points used in the geometric regularization.

211 **Encoders:** For our first experiments we used the **Vision Transformer (ViT)** (Dosovitskiy et al.,
212 2020) and the **Masked and Permuted Network (MPNet)** (Song et al., 2020). The base model of
213 the ViT has 86 million parameters and the base model of MPNet has 109 million parameters.

214 **Zero-Shot Accuracy Metric:** We use zero-shot accuracy (Xia et al., 2023) as the metric to assess
215 the quality of our alignment method, measured on ImageNet (Deng et al., 2009). The ImageNet
216 dataset has 1,000 classes, each class is represented by 50 images in the evaluation split. As in
217 Radford et al. (2021), we encode the class names using various prompt templates and average them
218 in the shared embedding space. The images are directly mapped into the common embedding. We
219 calculate the proximity between the image embeddings and the class embedding vectors using cosine
220 similarity. Image classification is determined by computing the nearest class within the embedding
221 space. Clearly, as the alignment method improves, the zero-shot accuracy increases.

222 **Precision@ k Metric:** For evaluation beyond image-text alignment we use the Precision@ k metric,
223 applied to the test split of the same dataset used for training. We select 10,000 test pairs resulting in
224 10,000 classes, such that the samples from one modality form classes, and we attempt their retrieval
225 based on corresponding samples of the other modality, and vice versa. Our findings are reported in
226 terms of precision@1 and precision@5. The test samples remain consistent across all experiments.

227 5.2 BASELINES

228 We verify the effectiveness of our proposed method by comparing it to established baseline models.
229 First, we examine the Procrustes alignment method, which is designed to learn a rotation matrix
230 that aligns one embedding with another. Then, we assess the performance of our alignment trans-
231 formation functions when trained solely with the contrastive loss, without including our geometry-
232 preserving regularization method. Lastly, we provide a comparison with the ASIF method, as de-
233 tailed in the related work section.

234 5.3 IMAGE AND TEXT ALIGNMENT THROUGH GERA

235 5.3.1 BENCHMARKING ON IMAGENET AND CC12M

236 We test GeRA with a neighborhood size of $K = 150$, using the heat kernel approximation as the
237 neighborhood encoding scheme, and with the “biased” sampling method. We evaluate our perfor-
238 mance compared to the baseline methods on the default configuration as described in section 5.1.

239 **Results:** Figure 3 shows that GeRA consistently outperforms both Procrustes Alignment and the
240 unregularized alignment based on the contrastive loss. This validates GeRA’s design choice of
241 balancing local preservation of geometric structures with global flexibility in the alignment process.
242 GeRA demonstrates a significant improvement of almost 9% over the unregularized alignment. The
243 increase in performance is particularly notable in situations where data availability is highly limited,
244 where accuracy improves from 3% for the contrastive loss trained with 1000 samples to almost 9%.

245 Figure 3 depicts that GeRA is the best-performing model in the low-data regimes. As the volume
 246 of data increases, ASIF slightly outperforms GeRA. However, in Figure 4, GeRA exhibits better
 247 results when evaluated on precision@5. In Section 5.7, we further demonstrate the advantage of
 248 GeRA over ASIF, in terms of inference time.

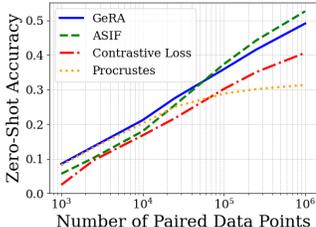


Figure 3: GeRA performance evaluated at zero-shot accuracy on ImageNet against the baselines.

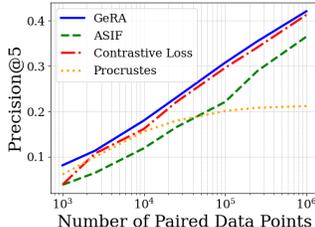


Figure 4: GeRA performance evaluated at precision@5 on in-distribution CC12M data against the baselines.

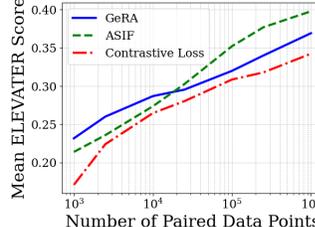


Figure 5: Comparison using mean scores of GeRA, non-regularized method, and ASIF across 20 ELEVATER datasets.

249 5.3.2 BENCHMARKING ON ELEVATER

250 To further validate the generality of our method, we expand our experiments beyond ImageNet and
 251 CC12M, incorporating multiple vision-language datasets into our evaluation pipeline. We employed
 252 the ELEVATER (Li et al., 2022) benchmark, which contains 20 image classification datasets. These
 253 datasets cover a broad spectrum of visual concepts, each presenting varying levels of difficulty.

254 **Results:** Figure 5 demonstrates that GeRA consistently outperforms the non-regularized method on
 255 the ELEVATER benchmark. In low-data regimes, the benefits of geometric regularization become
 256 especially clear. With 1,000 training pairs, the performance gain exceeds 5%. Even with 1 mil-
 257 lion training pairs, GeRA delivers an average performance improvement of 2.7%. Considering the
 258 diverse visual concepts, it becomes clear that GeRA has superior generalization capabilities.

259 Compared to ASIF, GeRA demonstrates superior performance in low-data regimes. When trained
 260 with 2,500 paired points, GeRA yields a mean score that is almost 3% higher. The advantage
 261 of GeRA diminishes as the number of paired points increases; ASIF surpasses GeRA when more
 262 paired points are available in training. Specifically, when trained with 1,000,000 paired points, ASIF
 263 achieves a mean score 3% higher than GeRA’s. However, in this regime, ASIF is more than 100×
 264 slower at inference time, as demonstrated in Figure 8.

265 5.4 SPEECH AND TEXT ALIGNMENT THROUGH GERA

266 To further validate adaptability and performance across diverse modalities, we consider the domain of speech-
 267 text alignment. We show that our method’s efficacy is not confined to a specific modality and that our hyperparam-
 268 eter choices, optimized for the image-text scenario, are not overfit to that context.
 269

272 **Encoder:** We use Whisper (Radford et al., 2023) as the
 273 speech encoder consisting of 74 million parameters. For
 274 text, we again use MPNet (see Section 5.1).

275 **Dataset:** Our training uses the LibriTTS dataset (Zen
 276 et al., 2019). This dataset is an assembly of text-speech
 277 pairs, aggregating to 585 hours of read English speech.
 278 Each entry corresponds to distinct sentences of speech
 279 and their textual counterparts. Entries with significant
 280 background noise are filtered out. The dataset includes
 281 205,044 pairs in totals, which is considerably smaller than
 282 the text-image alignment dataset of 12 million pairs. In

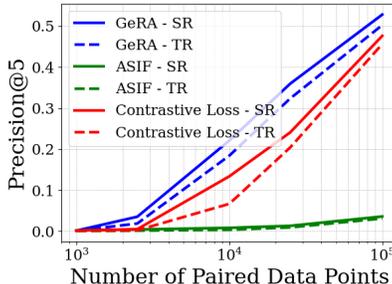


Figure 6: Performance of GeRA compared to ASIF and the pure contrastive learning evaluated at precision@5 for the speech and text alignment, using in-distribution LibriTTS data.

283 this experiment, we used up to 10^5 paired points in the contrastive loss, and additional 10^5 unpaired
 284 points in the geometric regularization.

285 **Results:** Figure 6 shows that GeRA significantly outperforms ASIF on speech–text alignment. The
 286 discrepancy increases as the number of paired training points increases. With 100,000 training pairs,
 287 GeRA achieves a precision@5 score for speech retrieval (SR) of over 51% while ASIF’s score is
 288 below 5%. These results show the generalizing capabilities of GeRA to the speech-text domain,
 289 while ASIF struggles with this modality. In addition, GeRA surpasses the model trained solely
 290 with the contrastive loss, which attains a precision@5 score of 48% for speech retrieval for 100,000
 291 paired training points.

292 5.5 INFLUENCE OF GEOMETRY PRESERVATION

293 We analyze the influence of geometric regularization on GeRA via ablation studies measuring the
 294 impact of various design choices. These choices include the size of the neighborhood kernel (number
 295 of neighbors), the kernel encodings, and the neighborhood sampling method.

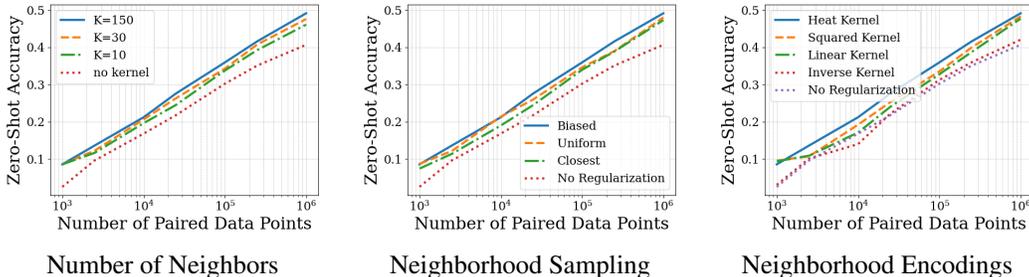


Figure 7: The impact of various design choices in GeRA, including neighborhood kernel size, geometry encoding scheme, and neighborhood sampling method.

296 5.5.1 RESULTS

297 **Number of Neighbors:** As the kernel matrix size increases, performance improves, but with di-
 298 minishing returns. Initial increases in neighborhood size yield substantial gains, but this increase
 299 plateaus, yielding a trade-off between accuracy and computational cost. We achieve a zero-shot
 300 accuracy of over 49% when trained on 1,000,000 paired points using a neighborhood size of 150.
 301 Comparing unregularized baseline model to GeRA yields a 9% increase in top-1 accuracy.

302 **Neighborhood Sampling:** Our experiments show that the sampling method affects accuracy. Bi-
 303 ased sampling, using primarily close but also including distant neighbors, proves most effective.
 304 The uniform distribution ranks second, including equally close and distant neighbors, whereas the
 305 least effective sampling method is the “closest” method, which only includes the nearest neighbors
 306 All methods surpass the neighborhood-free baseline.

307 **Neighborhood Encoding:** Regularization with any of our neighborhood encodings performs better
 308 than the contrastive loss alone. Among all, the heat kernel encoding consistently outperforms the
 309 other encodings by 2% on average over all training sizes, echoing the theoretical properties inspiring
 310 its choice. Overall, our choice of the heat kernel is confirmed to capture geometric information and
 311 demonstrates the benefit of geometric regularization in alignment tasks with limited paired data.

312 5.6 INFLUENCE OF PRETRAINED ENCODERS

313 GeRA is encoder-agnostic and hence not tied to a specific choice of encoders. We initially adopted
 314 the configurations that were previously tested with ASIF. Next, we discuss the generality of GeRA
 315 across different encoders, demonstrated empirically. See results in Figure 10 in the Appendix.

316 **Results:** Our method frequently surpasses ASIF by a significant margin with different encoders.
 317 This includes CLIP Encoders, the combinations of ViT-RoBERTa, ViT-BERT, and MAE-MPNet. In
 318 the setting recommended by ASIF, namely ViT-MPNet or using ViT-SentenceTransformer BERT,

319 our performance is either on par with or slightly below that of ASIF. Overall, our method offers
 320 more consistent and stable results compared to ASIF.

321 5.7 TRAINING AND INFERENCE TIME

322 In Figure 9, GeRA’s training time increases with the neighborhood size used for geometric regu-
 323 larization. Even with the highest number of paired training points (1 million pairs) and the largest
 324 kernel size ($K = 150$), however, training GeRA on an NVIDIA GeForce RTX 3090 only takes 20
 325 hours. Figure 8 shows that GeRA has consistent inference times, as the alignment transformation
 326 during inference is not affected by neighborhood size or number of training pairs. Conversely, ASIF
 327 has significant overhead, with inference times increasing linearly in the number of anchor points.
 328 For example, with 1 million training pairs, one ASIF inference takes over 2 hours for the retrieval
 329 task, while GeRA completes in under 16 seconds.

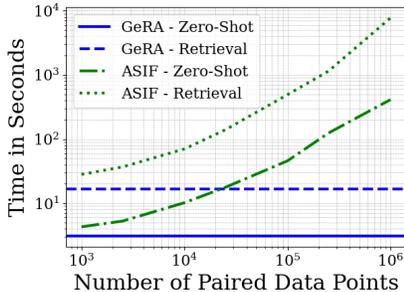


Figure 8: Comparison of inference times between GeRA and ASIF for the Zero-Shot and the Retrieval Evaluation.

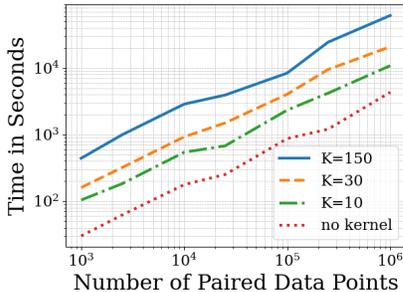


Figure 9: Time for training GeRA with different neighborhood sizes using an NVIDIA GeForce RTX 3090.

330 6 DISCUSSION

331 **Limitations:** Preserving local geometry requires taking neighborhood information into account,
 332 which leads to quickly increasing batch sizes, i.e., for batch size B and geometric regularization
 333 with K neighbors, the effective batch size becomes $B \cdot K$. This limits the ability to use larger batch
 334 sizes, which may slow down the convergence.

335 Moreover, GeRA depends on powerful pretrained models that define the geometry. In the absence
 336 of powerful pretrained models, the regularization’s effectiveness diminishes. Our experiments in
 337 section 5.6 show that selecting powerful encoders are necessary for both GeRA and ASIF. One po-
 338 tential solution in the absence of powerful pretrained encoders is to collect corrected neighborhood
 339 information for our loss term using human annotations or rules defined by a domain expert.

340 **Future Work:** Our work opens several interesting future work directions. In terms of the attrac-
 341 tive capability of GeRA to align domains with limited paired data supervision, there are several
 342 other modalities and downstream tasks that could be explored. Examples include aligning protein
 343 sequences and biomedical texts, which is needed for protein representation learning. Traditional
 344 unimodal approaches, which only focus on protein sequences, often miss functional aspects of pro-
 345 teins. Recent efforts incorporate text data on protein functions as an additional modality, enriching
 346 representations (Xu et al., 2023). However, the available datasets are relatively small, featuring
 347 only half a million paired data points. As a result, this domain is a key target for future work on
 348 GeRA. Additional future work directions for GeRA include exploring learnable parametric geomet-
 349 ric kernels (e.g. realized as self-attention blocks or small transformers), simultaneous co-training
 350 and multi-task training of both the encoders (on the unimodal data components) and the GeRA
 351 alignment module leading to dynamically changing manifolds landscape and potentially requiring
 352 exploring into momentum models for increased training stability, exploring multi-scale (coarsen-
 353 ing, multi-grid) manifold mapping methods to further enhance the preservation of the more global
 354 manifold structure after alignment, and many more.

355 REFERENCES

- 356 David Alvarez-Melis and Tommi S. Jaakkola. Gromov-wasserstein alignment of word embedding
357 spaces. *CoRR*, abs/1809.00013, 2018. URL <http://arxiv.org/abs/1809.00013>.
- 358 Richard Antonello, Javier S Turek, Vy Vo, and Alexander Huth. Low-dimensional structure in the
359 space of language representations is reflected in brain responses. *Advances in neural information*
360 *processing systems*, 34:8332–8344, 2021.
- 361 Tadas Baltrušaitis, Chaitanya Ahuja, and Louis-Philippe Morency. Multimodal machine learning:
362 A survey and taxonomy. *IEEE transactions on pattern analysis and machine intelligence*, 41(2):
363 423–443, 2018.
- 364 Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12m: Pushing
365 web-scale image-text pre-training to recognize long-tail visual concepts. In *Proceedings of the*
366 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3558–3568, 2021.
- 367 Ting Chen, Simon Kornblith, Kevin Swersky, Mohammad Norouzi, and Geoffrey E Hinton. Big
368 self-supervised models are strong semi-supervised learners. *Advances in neural information pro-*
369 *cessing systems*, 33:22243–22255, 2020.
- 370 Zhongzhi Chen, Guang Liu, Bo-Wen Zhang, Fulong Ye, Qinghong Yang, and Ledell Wu. Alt-
371 clip: Altering the language encoder in clip for extended language capabilities. *arXiv preprint*
372 *arXiv:2211.06679*, 2022.
- 373 Ronald R. Coifman and Stéphane Lafon. Diffusion maps. *Applied and Computational Har-*
374 *monic Analysis*, 21(1):5–30, 2006. ISSN 1063-5203. doi: [https://doi.org/10.1016/j.acha.](https://doi.org/10.1016/j.acha.2006.04.006)
375 [2006.04.006](https://doi.org/10.1016/j.acha.2006.04.006). URL [https://www.sciencedirect.com/science/article/pii/](https://www.sciencedirect.com/science/article/pii/S1063520306000546)
376 [S1063520306000546](https://www.sciencedirect.com/science/article/pii/S1063520306000546). Special Issue: Diffusion Maps and Wavelets.
- 377 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hier-
378 archical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*,
379 pp. 248–255, 2009. doi: 10.1109/CVPR.2009.5206848.
- 380 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
381 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszko-
382 reit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at
383 scale. *CoRR*, abs/2010.11929, 2020. URL <https://arxiv.org/abs/2010.11929>.
- 384 Samir Yitzhak Gadre, Gabriel Ilharco, Alex Fang, Jonathan Hayase, Georgios Smyrnis, Thao
385 Nguyen, Ryan Marten, Mitchell Wortsman, Dhruva Ghosh, Jieyu Zhang, et al. Datacomp: In
386 search of the next generation of multimodal datasets. *arXiv preprint arXiv:2304.14108*, 2023.
- 387 John C Gower. Generalized procrustes analysis. *Psychometrika*, 40:33–51, 1975.
- 388 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-
389 nition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.
390 770–778, 2016.
- 391 Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan
392 Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning
393 with noisy text supervision. In *International conference on machine learning*, pp. 4904–4916.
394 PMLR, 2021.
- 395 Jeff Johnson, Matthijs Douze, and Hervé Jégou. Billion-scale similarity search with gpus. *IEEE*
396 *Transactions on Big Data*, 7(3):535–547, 2019.
- 397 Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Joan Puigcerver, Jessica Yung, Sylvain Gelly,
398 and Neil Houlsby. Big transfer (bit): General visual representation learning. In *Computer Vision–*
399 *ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part*
400 *V 16*, pp. 491–507. Springer, 2020.

- 401 Chunyuan Li, Haotian Liu, Liunian Li, Pengchuan Zhang, Jyoti Aneja, Jianwei Yang, Ping Jin,
402 Houdong Hu, Zicheng Liu, Yong Jae Lee, et al. Elevater: A benchmark and toolkit for evaluating
403 language-augmented visual models. *Advances in Neural Information Processing Systems*, 35:
404 9287–9301, 2022.
- 405 Didong Li and David B Dunson. Geodesic distance estimation with spherelets. *arXiv preprint*
406 *arXiv:1907.00296*, 2019.
- 407 Luca Moschella, Valentino Maiorca, Marco Fumero, Antonio Norelli, Francesco Locatello, and
408 Emanuele Rodolà. Relative representations enable zero-shot latent space communication, 2023.
- 409 Maksim Nekrashevich, Alexander Korotin, and Evgeny Burnaev. Neural gromov-wasserstein opti-
410 mal transport. *arXiv preprint arXiv:2303.05978*, 2023.
- 411 Antonio Norelli, Marco Fumero, Valentino Maiorca, Luca Moschella, Emanuele Rodola, and
412 Francesco Locatello. Asif: Coupled data turns unimodal models to multimodal without train-
413 ing. *arXiv preprint arXiv:2210.01738*, 2022.
- 414 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agar-
415 wal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya
416 Sutskever. Learning transferable visual models from natural language supervision. *CoRR*,
417 abs/2103.00020, 2021. URL <https://arxiv.org/abs/2103.00020>.
- 418 Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever.
419 Robust speech recognition via large-scale weak supervision. In *International Conference on Ma-*
420 *chine Learning*, pp. 28492–28518. PMLR, 2023.
- 421 S. T. Roweis and L. K. Saul. Nonlinear dimensionality reduction by locally linear embedding. In
422 *Science*, pp. 2323–2326, 2000.
- 423 Nasir Saeed, Haewoon Nam, Mian Imtiaz Ul Haq, and Dost Bhatti Muhammad Saqib. A survey on
424 multidimensional scaling. *ACM Computing Surveys (CSUR)*, 51(3):1–25, 2018.
- 425 Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis,
426 Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of
427 clip-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114*, 2021.
- 428 Vikas Sindhwani, Partha Niyogi, and Mikhail Belkin. Beyond the point cloud: from transductive
429 to semi-supervised learning. In *Proceedings of the 22nd international conference on Machine*
430 *learning*, pp. 824–831, 2005.
- 431 Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. Mpnet: Masked and permuted pre-
432 training for language understanding. *Advances in Neural Information Processing Systems*, 33:
433 16857–16867, 2020.
- 434 Chen Sun, Abhinav Shrivastava, Saurabh Singh, and Abhinav Gupta. Revisiting unreasonable ef-
435 fectiveness of data in deep learning era. In *Proceedings of the IEEE international conference on*
436 *computer vision*, pp. 843–852, 2017.
- 437 Joshua B Tenenbaum, Vin de Silva, and John C Langford. A global geometric framework for
438 nonlinear dimensionality reduction. *science*, 290(5500):2319–2323, 2000.
- 439 Bart Thomee, David A Shamma, Gerald Friedland, Benjamin Elizalde, Karl Ni, Douglas Poland,
440 Damian Borth, and Li-Jia Li. Yfcc100m: The new data in multimedia research. *Communications*
441 *of the ACM*, 59(2):64–73, 2016.
- 442 S. R. S. Varadhan. *On the behavior of the fundamental solution of the heat equation with variable*
443 *coefficients*. Communications on Pure and Applied Mathematics, 1967.
- 444 Gary Wang, Kyle Kastner, Ankur Bapna, Zhehuai Chen, Andrew Rosenberg, Bhuvana Ramabhad-
445 ran, and Yu Zhang. Understanding shared speech-text representations. In *ICASSP 2023-2023*
446 *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 1–5.
447 IEEE, 2023.

- 448 Tongzhou Wang and Phillip Isola. Understanding contrastive representation learning through align-
449 ment and uniformity on the hypersphere. In *International Conference on Machine Learning*, pp.
450 9929–9939. PMLR, 2020.
- 451 Heming Xia, Qingxiu Dong, Lei Li, Jingjing Xu, Ziwei Qin, and Zhifang Sui. Imagenetvc:
452 Zero-shot visual commonsense evaluation on 1000 imagenet categories. *arXiv preprint*
453 *arXiv:2305.15028*, 2023.
- 454 Minghao Xu, Xinyu Yuan, Santiago Miret, and Jian Tang. Protst: Multi-modality learning of protein
455 sequences and biomedical texts. *arXiv preprint arXiv:2301.12040*, 2023.
- 456 Heiga Zen, Viet Dang, Rob Clark, Yu Zhang, Ron J Weiss, Ye Jia, Zhifeng Chen, and Yonghui Wu.
457 Libritts: A corpus derived from librispeech for text-to-speech. *arXiv preprint arXiv:1904.02882*,
458 2019.
- 459 Xiaohua Zhai, Xiao Wang, Basil Mustafa, Andreas Steiner, Daniel Keysers, Alexander Kolesnikov,
460 and Lucas Beyer. Lit: Zero-shot transfer with locked-image text tuning. *CoRR*, abs/2111.07991,
461 2021. URL <https://arxiv.org/abs/2111.07991>.
- 462 Xiaojin Zhu, Zoubin Ghahramani, and John D Lafferty. Semi-supervised learning using gaussian
463 fields and harmonic functions. In *Proceedings of the 20th International conference on Machine*
464 *learning (ICML-03)*, pp. 912–919, 2003.

465 A APPENDIX

466 A.1 ADDITIONAL DETAILS ON FIGURE 1(B)

467 The metrics reported in Figure 1(b) are computed based on the models trained in Section 5.3.1
 468 on the CC12M dataset (see results in Figure 3). The model labeled in Figure 1(b) as “No Geom.
 469 Regularization” is an alignment model trained with the contrastive loss only, and the modeled labeled
 470 as “Geom. Regularization” is an alignment model trained with the GeRA model using $K = 150$
 471 neighbors. Both models were trained using 10^6 paired points for the contrastive loss, and 6×10^6
 472 points for the geometric regularization.

473 In order to further demonstrate the neighborhood distortion effect, we evaluate the image-to-image
 474 kNN classification accuracy of ImageNet data using the pretrained image encoder, with and with-
 475 out the alignment layers, trained with the contrastive loss or with GeRA loss as described above.
 476 More concretely, we use the validation set of ImageNet, which consists of 50 samples per class. We
 477 randomly select 10 samples from each class as labeled training data. We embed this training data
 478 and the remaining images using the models. Each image was then assigned to a class based on the
 479 majority vote among its $k = 5$ nearest neighbors. We note that different values of k led to similar
 480 results. We then compare the performance of GeRA against the unregularized contrastive learning
 481 model and the original embedding space generated by ViT before applying our transformation, re-
 482 ported in Table 1. This experiment demonstrates that GeRA preserves the geometry of the image
 483 space obtained by the ViT model pre-trained on the image domain, while alignment with vanilla
 484 contrastive loss disturbs it.

Table 1: ImageNet kNN accuracy computed in the embedding spaces of ViT, ViT+alignment layers trained with contrastive loss only, and ViT+alignment layers trained with GeRA.

Method	kNN Classifier Accuracy ($k = 5$)
ViT only	0.76
No Geom. Regularization	0.67
Geom. Regularization	0.75

485 A.2 ADDITIONAL EXPERIMENTAL DETAILS AND HYPERPARAMETERS

486 A.2.1 ADDITIONAL EXPERIMENTAL DETAILS

487 **Use of Unpaired Points:** To give a bit more detail on the use of unpaired data in our experiments,
 488 in the image-text experiments, the dataset we used for training is CC12M which is a paired dataset.
 489 However, during training we only consider the pairings for a small number of samples and use
 490 the (fraction of) remaining samples as unpaired data to simulate a scenario where there are limited
 491 amounts of paired data and many unpaired data points. More concretely, we take M paired samples
 492 (used for contrastive loss) and include $N \gg M$ unpaired samples (distinct from the paired points
 493 used in the contrastive loss), where the unpaired points for each modality are chosen randomly and
 494 independently of the other modality. We leverage the unpaired data in the neighborhoods of each
 495 paired datapoint, and construct the kernels in the geometric regularization based on these neighbor-
 496 ing unpaired points. Note that for a pair (x, y) , the neighborhood for x is in general not the same as
 497 the neighborhood for y , i.e., the neighbors do not have to be pairs themselves.

498 **Pre-computing Nearest Neighbors for the Geometric Regularization:** To speed up training time,
 499 we pre-compute the neighborhood distributions, $S(x)$, from which $N_k(x)$ is sampled for the geo-
 500 metric regularization in equation 3. For each paired point in each modality, we collect 800 nearest
 501 neighbors for constructing $S(x)$. We perform this nearest neighbor search using Faiss (Johnson
 502 et al., 2019), which takes approximately 45 – 55 minutes to compute on an NVIDIA GeForce RTX
 503 3090, for 800 nearest neighbors of 6×10^6 samples in a 768 dimensional space, takes .

504 A.2.2 HYPERPARAMETERS

505 **Alpha (α):** This parameter balances the geometric regularization term with the contrastive objective,
 506 determining the relative importance of each in the loss function (See Equation 1).

Table 2: Summary of the explored hyperparameter spaces and the optimal values discovered for each method. The table outlines the range of values over which each hyperparameter was tuned, denoted in the ‘Range’ column. Subsequent columns present the hyperparameter values that yielded the best performance for each respective method, GeRA, Contrastive Loss, and ASIF (hyperparameters as stated in the paper (Norelli et al., 2022)), during experimentation. In instances where a hyperparameter is not applicable to a method, the cell is left blank.

Hyperparameter	Range	GeRA	Contrastive Loss	ASIF
Batch Size	500-4,000	2000	2000	–
Learning Rate	1e-5 – 5e-4	2e-4	2e-4	–
Dropout	0.0 – 0.5	0.3	0.3	–
Number of Hidden Layers	1 – 3	1	1	–
Hidden Dimension	768 – 16,000	8000	8000	–
Output Dimension	512 – 768	768	768	–
Number of Neighbors	5 – 150	150	–	–
Alpha	0.1 – 2.0	0.5	–	–
Epsilon (σ value)	0.1 – 3.0	0.8	–	–
Temperature	0.01 – 0.4	0.04	0.04	–
p (Exponentiation)	1 – 8	–	–	8
k (Sparsification)	50 – 1600	–	–	800

507 **Epsilon** (ϵ): This represents the kernel size, influencing the locality of the kernel matrix. A smaller
 508 epsilon value implies that the heat kernel captures more localized features, thereby considering
 509 neighbors in closer proximity (See Equation 4). In our experiments with the heat kernel we compute
 510 ϵ by: $\epsilon = \sigma \times \text{mean}(\{\|x_i - x_j\|_2^2\}_{i,j})$, i.e., a constant, σ , multiplied by the mean of the pairwise
 511 euclidean distances in the neighborhood. We found that $\sigma = 0.8$ performs best. This kernel normal-
 512 ization adapts to the scale and characteristics of each local neighborhood, and facilitates handling
 513 neighborhoods of different sizes and densities.

514 **Temperature** (t): Applied in the output layer, the temperature parameter modulates the sharpness
 515 of the distribution. A higher temperature results in a softer probability distribution over classes,
 516 whereas a lower temperature makes the distribution more concentrated (See Equation 2).

517 **Number of Neighbors** (K): This parameter specifies the number of neighbors included into the
 518 geometric regularization loss (see Equation 3). The larger the amount of neighbors, i.e., the larger
 519 the kernel matrix W , the better we can capture the local geometry and hence preserve it. Our
 520 ablation study in Figure 7 (left) demonstrates that the number of neighbors strongly correlates with
 521 the downstream alignment performance. However, the marginal increase in performance seems to
 522 diminish with larger neighborhoods, indicating already good performance using relatively small
 523 numbers of neighbors.

524 **Sampling Technique:** We aim to select samples that best represent the local geometry. Hence,
 525 selecting only the closest neighbors preserves locality best. However, to obtain better continuity
 526 of the embedding space, and increase the amount of information gathered from the neighbors in
 527 different epochs, we subsample the neighbors from a larger neighborhood distribution $S(x)$. We
 528 examined different ways of sampling the neighbors, including:

- 529 • ‘Uniform’ sampling, where K neighbors are uniformly sampled from the pre-computed neighbor-
 530 hood distribution $S(x)$ (including 800 points in our experiments). This approach includes closer
 531 and farther points with equal probabilities.
- 532 • ‘Closest’ sampling, where only the closest K neighbors are chosen from $S(x)$ for each paired
 533 point.
- 534 • ‘Biased’ sampling, where K neighbors are sampled from $S(x)$ with higher probabilities given to
 535 closer points.

536 Figure 7 (middle) demonstrates that sampling neighborhood points with a bias towards closest neigh-
 537 bors, i.e., higher probability of sampling closer neighbors while still including some information
 538 about further points, performs best overall.

539 A.2.3 PERFORMANCE DIFFERENCES OF GERA AND ASIF IN SPEECH-TEXT ALIGNMENT

540 Figure 6 shows a significant performance gap between GeRA and ASIF in speech-text alignment,
 541 in contrary to results for the image-text modality presented in Figure 3, for example. The study by
 542 Wang et al. (2023) found that training shared encoders for speech-text data produce more compact
 543 and overlapping representations, whereas the embedding spaces of uni-modal encoders yield distinct
 544 representations for speech and text.

545 Taking this observation into account, our results in Section 5.4 highlight the main advantage of our
 546 approach over ASIF. When the uni-modal embedding spaces are different, as suggested by Wang
 547 et al. (2023) for the speech and text modalities, we hypothesize that ASIF needs a lot more paired
 548 samples to properly align the spaces, since extrapolation from a limited number of pairs is likely to
 549 be inaccurate. Specifically, the performance obtained by ASIF in Wang et al. (2023), which is better
 550 than our reported ASIF performance in this setting, is using encoders that were trained on data that
 551 included paired points from the two modalities. In contrast, in our experiments, we use encoders
 552 that were trained on purely uni-modal data, and that may be the source of the performance gap.

553 Unlike ASIF, our method for alignment of the uni-modal models combines the strengths of con-
 554 trastive alignment with paired data, to match the uni-modal embedding spaces, and preserving the
 555 geometry of the respective spaces, thus it is more robust to the differences in uni-modal embedding
 556 spaces. For instance, in Figure 6, we see that vanilla contrastive loss performs quite well in aligning
 557 (purely) uni-modal models, significantly outperforming ASIF, while our method further improves
 558 the performance of the contrastive loss.

559 A.3 EVALUATING GERA ON DIFFERENT PRETRAINED ENCODERS

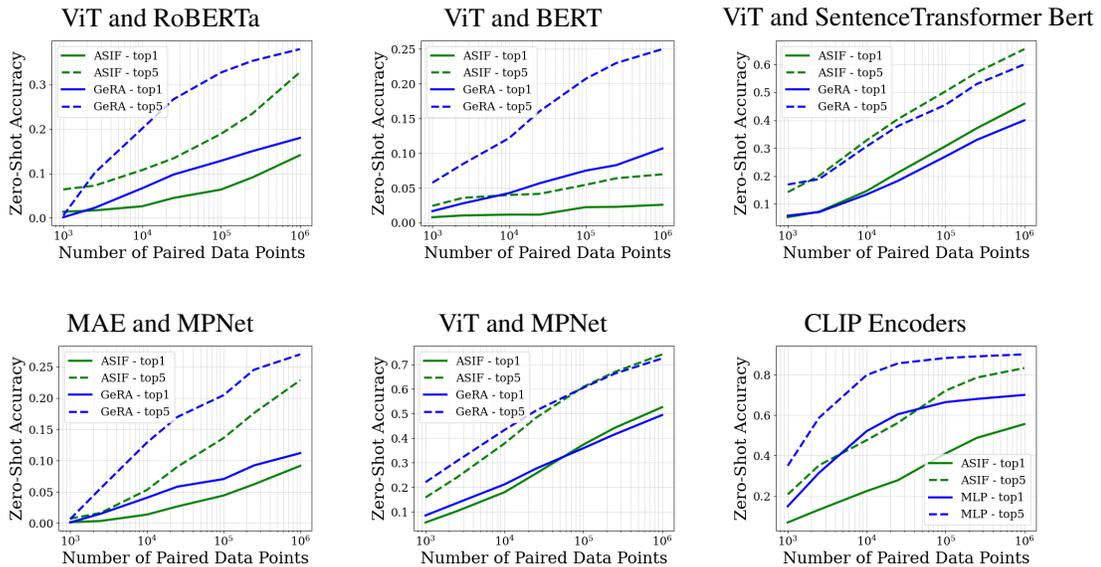


Figure 10: Performance comparison of GeRA and ASIF using various vision and language encoders.