# The Intrinsic Manifolds of Radiological Images and their Role in Deep Learning

Nicholas Konz[1*(✉)][0000−0003−0230−1598], Hanxue Gu[1][0000−0003−2622−753X], Haoyu Dong[2][0000−0002−5132−0341], and Maciej A. Mazurowski[1,2,3,4(✉)][0000−0003−4202−8602]

[1] Department of Electrical and Computer Engineering, Duke University, NC, USA
*corresponding author nicholas.konz@duke.edu
[2] Department of Radiology, Duke University, NC, USA
[3] Department of Computer Science, Duke University, NC, USA
[4] Department of Biostatistics & Bioinformatics, Duke University, NC, USA
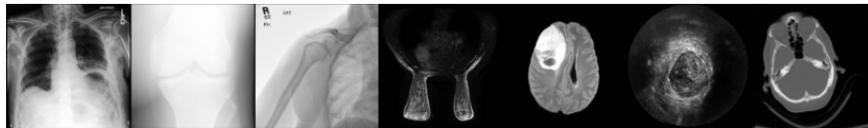maciej.mazurowski@duke.edu

**Abstract.** The manifold hypothesis is a core mechanism behind the success of deep learning, so understanding the intrinsic manifold structure of image data is central to studying how neural networks learn from the data. Intrinsic dataset manifolds and their relationship to learning difficulty have recently begun to be studied for the common domain of natural images, but little such research has been attempted for radiological images. We address this here. First, we compare the intrinsic manifold dimensionality of radiological and natural images. We also investigate the relationship between intrinsic dimensionality and generalization ability over a wide range of datasets. Our analysis shows that natural image datasets generally have a higher number of intrinsic dimensions than radiological images. However, the relationship between generalization ability and intrinsic dimensionality is much stronger for medical images, which could be explained as radiological images having intrinsic features that are more difficult to learn. These results give a more principled underpinning for the intuition that radiological images can be more challenging to apply deep learning to than natural image datasets common to machine learning research. We believe rather than directly applying models developed for natural images to the radiological imaging domain, more care should be taken to developing architectures and algorithms that are more tailored to the specific characteristics of this domain. The research shown in our paper, demonstrating these characteristics and the differences from natural images, is an important first step in this direction.

**Keywords:** Radiology · Generalization · Dimension · Manifold

## 1 Introduction

Although using deep learning-based methods to solve medical imaging tasks has become common practice, there lacks a strong theoretical understanding and

analysis of the effectiveness of such methods. This could be a potential problem for future algorithm development, as most successful methods for medical images are adapted from techniques solving tasks using natural image datasets [6]. Due to the apparent differences in relevant semantics between natural and medical domains [20], it is not clear what design choices are necessary when adapting these networks to medical images. This difference in domain is especially true when considering radiological images. Our goal is to provide a better, quantified footing for developing such radiology-specialized methods, by (1) analyzing the underlying structure of common radiological image datasets and determining how it relates to learning dynamics and generalization ability and (2) comparing these characteristics to common natural image datasets.



**Fig. 1. Sample images from each dataset studied.** From the left: CheXpert [15], OAI [23], MURA [25], DBC [26], BraTS 2018 [19], Prostate-MRI [29] and RSNA-IH-CT [10] (see Sec. 3).
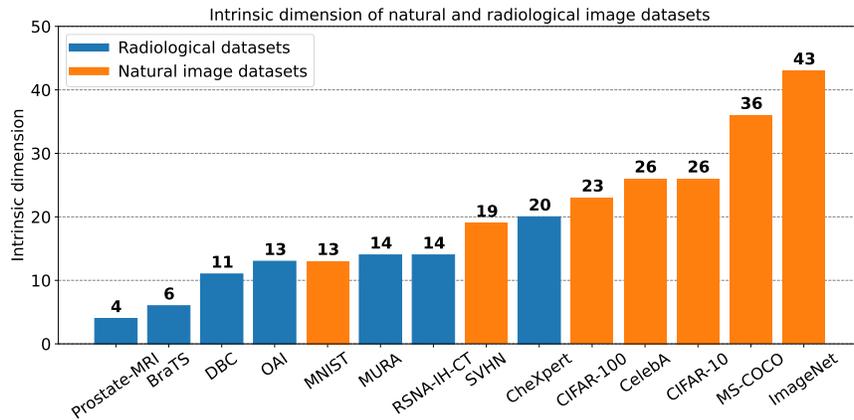
The Manifold Hypothesis [9,30,5] states that high-dimensional data, such as images, can be well described by a much smaller number of features/degrees of freedom than the number of pixels in an image; this number is the *intrinsic dimension* (ID) of the dataset. This is central to deep computer vision because these abstract visual features can be learned from data, allowing inference in a tractable, lower-dimensional space. It is therefore important to study the relationship of the ID of datasets with the learning process of deep models. This was recently explored for standard natural image datasets [24], but a similarly comprehensive study has yet to be conducted for radiological image datasets, which is important because of the apparent differences between these two domains.

**Contributions.** Our contributions are summarized as follows.

1. We investigate the intrinsic manifold structure of common radiology datasets (Fig. 1), and find that their IDs are indeed much lower than the number of pixels, and also generally lower than for natural image datasets (Fig. 2).
2. We also find that classification is generally harder with radiology datasets than natural images for moderate-to-low training set sizes.
3. We show that classification performance is negatively linear to dataset ID within both data domains, invariant to training set size. However, the absolute value of the slope of this relationship is much higher for radiological data than for natural image data.

4. We test these linearity findings on a wide range of common classification models, and find that performance for radiological images is almost independent of the choice of model, relying instead on the ID of the dataset.

Our results show that while the ID of a dataset affects the difficulty of learning from this dataset, what also matters is the complexity of the intrinsic features themselves, which we find to be indicated by the sharpness of the relationship between generalization ability and ID, that is more severe for radiological images. These findings give experimental evidence for the differences in intrinsic dataset structure and learning difficulty between the two domains, which we believe is the first step towards a more principled foundation for developing deep methods specially designed for radiology.



**Fig. 2. Intrinsic dimension of radiological (blue) and natural image (orange) datasets**, the latter from [24]. Figure recommended to be viewed in color.

**Related Work.** Intrinsic dimension (ID) estimation methods ([3]), have only recently been applied to datasets used in modern computer vision, beginning with [24], which explored the ID of common natural image datasets and how it relates to learning ability and generalization. There have been a few studies of the ID of medical datasets, *e.g.*, [7], but these are targeted at an individual modality or dataset. The most common ID estimator obtains a maximum-likelihood (MLE) solution for the ID by modeling the dataset as being sampled from a locally uniform Poisson process on the intrinsic data manifold [17]. Other estimators exist ([11,8]), but these are unreliable estimators for images [24], so we use the MLE estimator in this work. Note that we do not estimate dataset ID with some learned latent space dimension found by training an autoencoder-type model or similar; relatedly, we also do not study the ID of the learned feature structure,

or parameters of a trained feature extraction model, *e.g.*, [1,4]. In contrast, we study the intrinsic, model-independent structure of the dataset itself.

## 2   Methods: Intrinsic Dimension Estimation

Consider some dataset $\mathcal{D} \subset \mathbb{R}^d$ of $N$ images, where $d$ is the *extrinsic dimension*, *i.e.*, the number of pixels in an image. The Manifold Hypothesis [9] assumes that the datapoints $x \in \mathcal{D}$ lie near some low-dimensional manifold $\mathcal{M} \subseteq \mathbb{R}^d$ that can be described by an intrinsic dimension $\dim(\mathcal{M}) = m \ll d$.

In order to obtain an estimate of $m$, [17] models the data as being sampled from a Poisson process within some $m$-dimensional sphere at each datapoint $x$; the density of points is assumed to be approximately constant within the radius of the sphere. Rather than specifying the radius as a hyperparameter, the authors set this radius to be the distance of the $k^{th}$ nearest neighboring datapoint to $x$ (instead specifying $k$). By maximizing the likelihood of the data given the parameters of the Poisson model, we then obtain an ID for $x$ of

$$\hat{m}_k(x) = \left[ \frac{1}{k-1} \sum_{j=1}^{k-1} \log \frac{T_k(x)}{T_j(x)} \right]^{-1}, \tag{1}$$

where $T_j(x)$ is the $\ell_2$ distance from $x$ to its $j^{th}$ nearest neighbor. The authors of [18] then showed that MLE can be used again to obtain an estimate for the global dataset ID, as

$$\hat{m}_k = \left[ \frac{1}{N} \sum_{i=1}^{n} \hat{m}_k \left( x_i \right)^{-1} \right]^{-1} = \left[ \frac{1}{N(k-1)} \sum_{i=1}^{N} \sum_{j=1}^{k-1} \log \frac{T_k \left( x_i \right)}{T_j \left( x_i \right)} \right]^{-1}, \tag{2}$$

which is the estimator that we use for this work. Note here that $k$ is a hyperparameter of the estimator; just as in [24] we set $k = 20$, a moderate value as recommended by [17].

## 3   Datasets and Tasks

In this work, we use the common task of binary classification to analyze the effect of dataset intrinsic dimension on the learning of radiological images. We chose radiology datasets that are varied in modality and well-representative of the domain, while being large enough for a broad study of training set sizes for at least one realistic classification task. We explore using alternate tasks for the same datasets in Section 4.2.

The datasets are as follows. **(1)** CheXpert [15], where we detect pleural effusion in chest X-ray images. Next is **(2)** the Knee Osteoarthritis Initiative (OAI) [23], where we select the OAI-released screening packages 0.C.2 and 0.E.1 containing knee X-ray images. Following [2,31], we build a binary osteoarthritis

(OA) detection dataset by combining Kellgren-Lawrence (KL) scores of $\{0, 1\}$ as OA-negative and combining scores of $\{2, 3, 4\}$ as positive. **(3)** is MURA [25], where we detect abnormalities in musculoskeletal X-ray images. Next, **(4)** is the Duke Breast Cancer MRI (DBC) dataset [26], where we detect cancer in fat-saturated breast MRI volume slices; we take slices containing a tumor bounding box to be positive, and all other slices at least five slices away from the positives to be negative. We follow this same slice-labeling procedure for dataset **(5)**, BraTS 2018 [19] where we detect gliomas in T2 FLAIR brain MRI slices. Dataset **(6)** is Prostate-MRI-Biopsy [29], where we take slices from the middle 50% of each MRI volume, and label each slice according to the volume's assigned prostate cancer risk score; scores of $\{0, 1\}$ are negative, and scores of $\geq 2$, which correlates with risk of cancer [22], are positive. Our final dataset **(7)** is the RSNA 2019 Intracranial Hemorrhage Brain CT Challenge (RSNA-IH-CT) [10], where we detect any type of hemorrhage in Brain CT scans. Sample images are shown in Fig. 1.

## 4   Experiments and Results

### 4.1   The Intrinsic Dimension of Radiology Datasets

We first estimate the intrinsic dimension (ID) of the considered radiology datasets; the results are shown in Figure 2, alongside the natural image dataset results of [24]. For each dataset, we estimate the MLE ID (Equation (2)) on a sample of 7500 images such that there is an exact 50/50 split of negative and positive cases in the sample, to minimize estimator bias (although we found that when using fewer images, the estimates were little affected). Like natural images, we found the ID of radiological images to be many times smaller than the extrinsic dimension; however, radiological image datasets tend to have lower ID than natural images datasets. Intuitively, we found that modifying the dataset extrinsic dimension (resizing the images) had little effect on the ID.

### 4.2   Generalization Ability, Learning Difficulty and Intrinsic Dimension
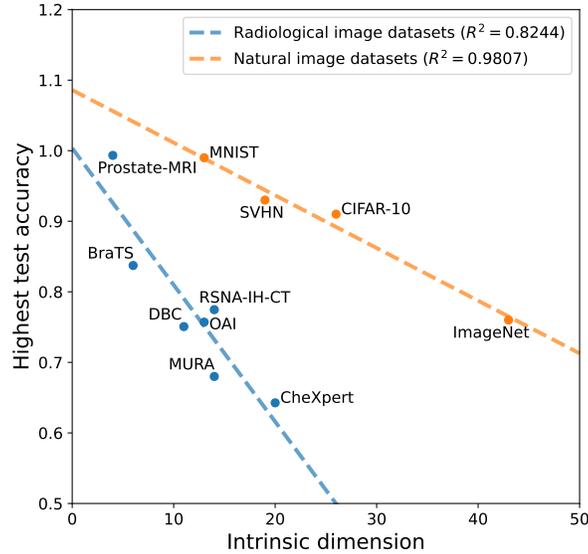
We now wish to determine what role the intrinsic dimension (ID) of a radiology dataset has in the degree of difficulty for a deep model to learn from it, and how this compares to the natural image domain. As in [24], we use the test set accuracy obtained when the model has maximally fit to the training set as a proxy for the generalization ability (GA) of the model on the dataset. We train and test with each dataset separately on it's respective binary classification task (Sec. 3), for a range of models and training set sizes $N_{\text{train}}$, from 500 to 2000. We train for 100 epochs with the same hyperparameters for all experiments; further details are provided in the supplementary materials. For each experiment with a studied dataset, we sample 2750 images from the given dataset such that there is an exact 50/50 split of negative and positive images. From this, we sample

750 test images, and from the remaining 2000 images, we sample some $N_{\text{train}}$ training examples.

Beginning with $N_{\text{train}} = 2000$ on ResNet-18, we plot the GA with respect to the dataset ID in Fig. 3, alongside the corresponding results for natural image datasets where binary classification was explored from [24]. Intriguingly, even across the range of tasks and datasets, we see that the relationship of GA with ID is approximately linear within each domain. Indeed, when fitting a simple ordinary least squares linear model

$$\text{GA} = a_{\text{GA,ID}} \times \text{ID} + b \tag{3}$$

to each of the two domains, we obtain a Pearson linear correlation coefficient of $R^2 = 0.824$ and slope of $a_{\text{GA,ID}} = -0.019$ for radiological images, and $R^2 = 0.981$ with $a_{\text{GA,ID}} = -0.0075$ for natural images.
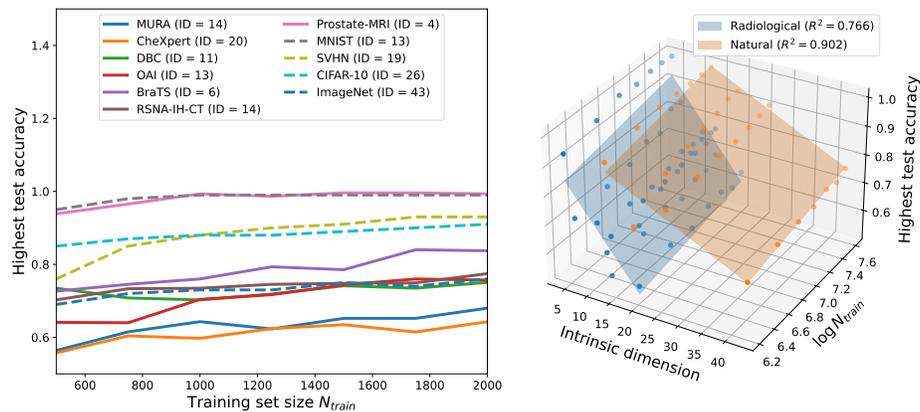


**Fig. 3. Linearity of model generalization ability with respect to dataset intrinsic dimension**, for radiological (blue) and natural (orange) image datasets, for $N_{\text{train}} = 2000$ on ResNet-18. Figure recommended to be viewed in color.

When repeating the same experiments over the aforementioned range of $N_{\text{train}}$ and ResNet models, we find that averaged over $N_{\text{train}}$, $R^2 = 0.76 \pm 0.05$, $a_{\text{GA,ID}} = -0.0199 \pm 0.0009$ and $R^2 = 0.91 \pm 0.12$, $a_{\text{GA,ID}} = -0.0077 \pm 0.0004$ ($\pm$ standard deviation) for the radiological and natural domains, respectively. These low deviations imply that even across a range of training sizes, both domains have a high, and mostly constant, correlation between GA and ID. Similarly,

both domains have approximately constant slopes of GA vs. ID with respect to $N_{\text{train}}$, but between domains, the slopes differ noticeably.

On the left of Fig. 4, we show how GA varies with respect to $N_{\text{train}}$ for datasets of both domains. We see that datasets with higher ID pose a more difficult classification task within both domains, *i.e.*, more training samples are required to achieve some test accuracy. However, between these two domains, radiological images are generally more difficult to generalize to than natural images, for these moderate-to-low training set sizes (that are typical for radiology tasks [28]). For example, OAI (ID = 13), MURA (ID = 14) and CheXpert (ID = 20), all prove to be more difficult than ImageNet, even though it has more than double the number of intrinsic dimensions (43).

This implies that the intrinsic dataset features described by these dimensions (and the correlations between them) can vary in learning difficulty between domains, indicated by the aforementioned sharper slope $a_{\text{GA,ID}}$. Our results show that radiological images generally have more difficult intrinsic features to learn than natural images, even if the number of these intrinsic feature dimensions is higher for the latter.



**Fig. 4. Left:** model generalization ability (GA) vs. training set size $N_{\text{train}}$ for various radiological (solid line) and natural (dashed line) image datasets, using ResNet-18. **Right:** linearity of GA with respect to $\log N_{\text{train}}$ and dataset intrisic dimension. Accompanies Fig. 3; recommended to be viewed in color.

We can also explore the dependence of generalization ability (GA) on both training set size $N_{\text{train}}$ and ID, shown on the right of Fig. 4. [21] found that learning requires a training sample size that grows exponentially with the data manifold's ID; this implies that GA should scale with $\log N_{\text{train}}$. Indeed we see that GA is approximately linear with respect to ID and $\log N_{\text{train}}$: for each domain of radiological and natural, we find multiple linear correlation coefficients $R^2$ of 0.766 and 0.902 between these variables, respectively. Given some new

radiology dataset, we could potentially use this to estimate the minimum $N_{\text{train}}$ needed to obtain a desired GA/test accuracy for the dataset, which would save overall training time. However, extensive experiments would need to be conducted before this could be widely applicable, due to confounding factors such as if the chosen model has a high enough capacity to fully fit to the training set; as such, we leave such an investigation for future work.

**Dependence on Model Choice.** As mentioned, we repeated the same experiments with a number of additional models for the radiology domain, to see if these linear relationships change with different model choices. In addition to ResNet-18, we tested on ResNet-34 and -50 ([12]), VGG-13, -16 and -19 ([27]), Squeezenet 1.1 ([14]), and DenseNet-121 and -169 ([13]). Averaged over $N_{\text{train}}$ and all models, we obtained $R^2 = 0.699 \pm 0.080$ and $a_{\text{GA,ID}} = -0.019 \pm 0.001$ (individual model results are provided in the supplementary materials). By the low deviations of both $R^2$ and the actual regressed slope $a_{\text{GA,ID}}$ of GA vs. ID, we see that the same linear relationship between GA and ID also exists for these models. We therefore infer that this relationship between classification GA and ID is largely independent of model size/choice, assuming that the model has a high enough capacity to fully fit to the training set.

**Dependence on Task Choice.** Logically, there should be other factors affecting a model's GA beyond the dataset's ID; *e.g.*, for some fixed dataset, harder tasks should be more difficult to generalize to. This section aims to determine how changing the chosen tasks for each dataset affects the preceding results. We will consider realistic binary classification tasks that have enough examples to follow the dataset generation procedure of Sec. 4.2. The datasets that support this are CheXpert—detect edema instead of pleural effusion, RSNA-IH-CT—detect subarachnoid hemorrhage, rather than any hemorrhage, and Prostate-MRI—detect severe cancer (score $> 2$), rather than any cancer. For robustness we will experiment on all three aforementioned ResNet models.

From switching all three datasets to their modified tasks, the linear fit parameters (averaged over $N_{\text{train}}$) changed as $R^2 = 0.76 \pm 0.05 \Rightarrow 0.78 \pm 0.15$ and $a_{\text{GA,ID}} = -0.0199 \pm 0.0009 \Rightarrow -0.012 \pm 0.002$ for ResNet-18; $R^2 = 0.77 \pm 0.07 \Rightarrow 0.76 \pm 0.21$ and $a_{\text{GA,ID}} = -0.019 \pm 0.001 \Rightarrow -0.011 \pm 0.003$ for ResNet-34; and $R^2 = 0.78 \pm 0.07 \Rightarrow 0.82 \pm 0.15$ and $a_{\text{GA,ID}} = -0.021 \pm 0.001 \Rightarrow -0.013 \pm 0.004$ for ResNet-50. We therefore conclude that the choice of task has some, but a small effect on the significance of the linear relationship of GA with ID ($R^2$), but can affect the parameters (slope $a_{\text{GA,ID}}$) of the relationship. Certainly, the choice of task represents a very large space of possibilities, so we leave more comprehensive investigations for future work.

## 5   Conclusion

Our results provide empirical evidence for the practical differences between the two domains of radiological and natural images, in both the intrinsic structure

of datasets, and the difficulty of learning from them. We found that radiological images generally have a lower intrinsic dimension (ID) than natural images (Fig. 2), but at the same time, they are generally harder to learn from (Fig. 4, left). This indicates that the intrinsic features of radiological datasets are more complex than those in natural image data. Therefore, assumptions about natural images and models designed for natural images should not necessarily be extended to radiological datasets without consideration for these differences. Further study of the differences between these two domains and the conceptual reasons for why they arise could lead to helpful guidelines for deep learning with radiology. We believe that the results in this work are an important step in this direction, and they lay the foundation for further research on this topic.

## References

1. Alessio Ansuini, Alessandro Laio, Jakob H Macke, and Davide Zoccolan. Intrinsic dimension of data representations in deep neural networks. *Advances in Neural Information Processing Systems*, 32, 2019.
2. Joseph Antony, Kevin McGuinness, Kieran Moran, and Noel E O'Connor. Automatic detection of knee joints and quantification of knee osteoarthritis severity using convolutional neural networks, 2017.
3. Jonathan Bac, Evgeny M Mirkes, Alexander N Gorban, Ivan Tyukin, and Andrei Zinovyev. Scikit-dimension: a python package for intrinsic dimension estimation. *Entropy*, 23(10):1368, 2021.
4. Tolga Birdal, Aaron Lou, Leonidas J Guibas, and Umut Simsekli. Intrinsic dimension, persistent homology and generalization in neural networks. *Advances in Neural Information Processing Systems*, 34, 2021.
5. Matthew Brand. Charting a manifold. *Advances in neural information processing systems*, 15, 2002.
6. Patryk Chrabaszcz, Ilya Loshchilov, and Frank Hutter. A downsampled variant of imagenet as an alternative to the cifar datasets. *arXiv preprint arXiv:1707.08819*, 2017.
7. Dietmar Cordes and Rajesh R Nandy. Estimation of the intrinsic dimensionality of fmri data. *Neuroimage*, 29(1):145–154, 2006.
8. Elena Facco, Maria d'Errico, Alex Rodriguez, and Alessandro Laio. Estimating the intrinsic dimension of datasets by a minimal neighborhood information. *Scientific reports*, 7(1):1–8, 2017.
9. Charles Fefferman, Sanjoy Mitter, and Hariharan Narayanan. Testing the manifold hypothesis. *Journal of the American Mathematical Society*, 29(4):983–1049, 2016.
10. Adam E Flanders, Luciano M Prevedello, George Shih, Safwan S Halabi, Jayashree Kalpathy-Cramer, Robyn Ball, John T Mongan, Anouk Stein, Felipe C Kitamura, Matthew P Lungren, et al. Construction of a machine learning dataset through collaboration: the rsna 2019 brain ct hemorrhage challenge. *Radiology: Artificial Intelligence*, 2(3):e190211, 2020.
11. Marina Gomtsyan, Nikita Mokrov, Maxim Panov, and Yury Yanovich. Geometry-aware maximum likelihood estimation of intrinsic dimension. In *Asian Conference on Machine Learning*, pages 1126–1141. PMLR, 2019.
12. Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

13. Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4700–4708, 2017.

14. Forrest N Iandola, Song Han, Matthew W Moskewicz, Khalid Ashraf, William J Dally, and Kurt Keutzer. Squeezenet: Alexnet-level accuracy with 50x fewer parameters and< 0.5 mb model size. *arXiv preprint arXiv:1602.07360*, 2016.

15. Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silviana Ciurea-Ilcus, Chris Chute, Henrik Marklund, Behzad Haghgoo, Robyn Ball, Katie Shpanskaya, et al. Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 590–597, 2019.

16. Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

17. Elizaveta Levina and Peter Bickel. Maximum likelihood estimation of intrinsic dimension. *Advances in neural information processing systems*, 17, 2004.

18. David JC MacKay and Zoubin Ghahramani. Comments on'maximum likelihood estimation of intrinsic dimension'by e. levina and p. bickel (2004). *The Inference Group Website, Cavendish Laboratory, Cambridge University*, 2005.

19. Bjoern H Menze, Andras Jakab, Stefan Bauer, Jayashree Kalpathy-Cramer, Keyvan Farahani, Justin Kirby, Yuliya Burren, Nicole Porz, Johannes Slotboom, Roland Wiest, et al. The multimodal brain tumor image segmentation benchmark (brats). *IEEE transactions on medical imaging*, 34(10):1993–2024, 2014.

20. Lia Morra, Luca Piano, Fabrizio Lamberti, and Tatiana Tommasi. Bridging the gap between natural and medical images through deep colorization. In *2020 25th International Conference on Pattern Recognition (ICPR)*, pages 835–842. IEEE, 2021.

21. Hariharan Narayanan and Sanjoy Mitter. Sample complexity of testing the manifold hypothesis. *Advances in neural information processing systems*, 23, 2010.

22. S. Natarajan, A. Priester, D. Margolis, J. Huang, and L. Marks. Prostate mri and ultrasound with pathology and coordinates of tracked biopsy (prostate-mri-us-biopsy) [data set]. DOI: 10.7937/TCIA.2020.A61IOC1A, 2020. Accessed: 2022-02-21.

23. M. Nevitt, D. Felson, and Gayle Lester. The osteoarthritis initiative. *Protocol for the cohort study*, 1, 2006.

24. Phillip Pope, Chen Zhu, Ahmed Abdelkader, Micah Goldblum, and Tom Goldstein. The intrinsic dimension of images and its impact on learning. *arXiv preprint arXiv:2104.08894*, 2021.

25. Pranav Rajpurkar, Jeremy Irvin, Aarti Bagul, Daisy Ding, Tony Duan, Hershel Mehta, Brandon Yang, Kaylie Zhu, Dillon Laird, Robyn L Ball, et al. Mura: Large dataset for abnormality detection in musculoskeletal radiographs. *arXiv preprint arXiv:1712.06957*, 2017.

26. Ashirbani Saha, Michael R Harowicz, Lars J Grimm, Connie E Kim, Sujata V Ghate, Ruth Walsh, and Maciej A Mazurowski. A machine learning approach to radiogenomics of breast cancer: a study of 922 subjects and 529 dce-mri features. *British journal of cancer*, 119(4):508–516, 2018.

27. Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.

28. Shelly Soffer, Avi Ben-Cohen, Orit Shimon, Michal Marianne Amitai, Hayit Greenspan, and Eyal Klang. Convolutional neural networks for radiologic images: a radiologist's guide. *Radiology*, 290(3):590–606, 2019.

29. Geoffrey A Sonn, Shyam Natarajan, Daniel JA Margolis, Malu MacAiran, Patricia Lieu, Jiaoti Huang, Frederick J Dorey, and Leonard S Marks. Targeted biopsy in the detection of prostate cancer using an office based magnetic resonance ultrasound fusion device. *The Journal of urology*, 189(1):86–92, 2013.
30. Joshua B Tenenbaum, Vin de Silva, and John C Langford. A global geometric framework for nonlinear dimensionality reduction. *science*, 290(5500):2319–2323, 2000.
31. Aleksei Tiulpin, Jérôme Thevenot, Esa Rahtu, Petri Lehenkari, and Simo Saarakkala. Automatic Knee Osteoarthritis Diagnosis from Plain Radiographs: A Deep Learning-Based Approach. *Scientific Reports*, 8(1):1727, 2018.

# Supplementary Materials for "The Intrinsic Manifolds of Radiological Images and their Role in Deep Learning"

## A   Experimental Settings (for reproducibility)

**Table 1.** Training/Inference Parameters and Settings

| Name | Value |
|------|-------|
| Loss Function | Binary Cross-Entropy |
| Optimizer | Adam [16] |
| Learning rate | 0.001 |
| Weight decay | 0.0001 |
| Image normalization | $[0, 255]$ |
| Image resolution | $224 \times 224$ |
| Training image augmentations | None |
| GPUs | $4\times$ NVIDIA GeForce RTX 3070 (8 GB) |
| Experiment Software | PyTorch 1.8.1, Numpy 1.22.1, CUDA 11.4 |
| Visualization Software | MatPlotLib 3.4.2 |
| Typical model training time (on $N_{\text{train}} = 2000$) | 10-30 min. |

**Table 2.** Model-Specific Details and Results (averaged over $N_{\text{train}}$ with std. dev.)

| Model | Training batch size | $R^2$ (GA vs. ID) | slope $a_{\text{GA,ID}}$ (GA vs. ID) |
|-------|---------------------|-------------------|--------------------------------------|
| ResNet-18 [12] | 200 | $0.756 \pm 0.052$ | $-0.0199 \pm 0.0009$ |
| ResNet-34 [12] | 128 | $0.772 \pm 0.071$ | $-0.0193 \pm 0.0012$ |
| ResNet-50 [12] | 64 | $0.781 \pm 0.066$ | $-0.0207 \pm 0.0010$ |
| VGG-13 [27] | 32 | $0.646 \pm 0.048$ | $-0.0194 \pm 0.0009$ |
| VGG-16 [27] | 32 | $0.623 \pm 0.066$ | $-0.0184 \pm 0.0008$ |
| VGG-19 [27] | 32 | $0.597 \pm 0.100$ | $-0.0168 \pm 0.0031$ |
| Squeezenet 1.1 [14] | 32 | $0.580 \pm 0.073$ | $-0.0173 \pm 0.0011$ |
| DenseNet-121 [13] | 32 | $0.770 \pm 0.073$ | $-0.0190 \pm 0.0009$ |
| DenseNet-169 [13] | 32 | $0.765 \pm 0.061$ | $-0.0189 \pm 0.0008$ |