# Unsupervised Data Pattern Discovery on the Cloud

Thomas Cecconello[1], Lucas Puerari[1] and Giuseppe Vizzari[1]

[1]*Department of Informatics, Systems and Communication (DISCo), University of Milano-Bicocca, Milan, Italy*

## Abstract

Scientific research implies the production of data describing phenomena still not studied and well understood. Sometimes the amount and rate of generation of produced data can be overwhelming, and anyway tools supporting a computer assisted analysis of scientific data can support systematic forms of data driven analysis. Machine learning can be an instrument in an overall flow including domain experts and computer scientists. Adopted machine learning approaches need to be unsupervised, employing just the input data as a teacher. We propose a two-step workflow: (i) achieving a compact representation of elements of the dataset by means of representation learning techniques, shifting the analysis from cumbersome representations to compact vectors in a latent space, and (ii) clustering points associated to instances to suggest patterns to the domain experts that will evaluate their potential meaning within the domain. The paper presents the rationale of the approach within a cloud based setting, and first experiments on an image dataset from the literature.

## Keywords

Representation learning, Pattern discovery, Clustering, Cloud

## 1. Introduction

Scientific research is a human activity that often implies the production of new or improved measurement tools, leading to the generation of new data, that needs to be analyzed, either to corroborate existing theories, or to support the generation of new ones. Next generation astronomical facilities, for instance, such as the Square Kilometre Array will generate an overwhelming volume of data at a rate that simply cannot be matched by our ability to make sense out of them.

Sometimes, moreover, it is impossible to use supervised techniques to support these researches: the studied phenomena are often object of intense study and classification schemes are still not agreed upon, they might be uncertain, or new data was acquired exactly with the goal of defining a classification scheme that was impossible to come up with using previously available data. Automated/semi-automated tools are needed to support this kind of research: machine learning can be an instrument in an overall workflow including domain experts and computer scientists. We emphasize that adopted machine learning approaches need to be essentially unsupervised, employing just the input data as a teacher: as we will discuss more thoroughly in Sect. 3.1, this kind of investigation has seen recently a certain interest and attention with particular reference to so-called self–supervised learning approaches. This term refers to a way

of framing an unsupervised learning problem so as to apply supervised learning algorithms to solve it; this typically implies finding a way to create a loss function not requiring user labelling of data (although it might consider forms of automatic labelling). The workflow proposed in this paper comprises two steps: (i) achieving a compact representation of the elements of the dataset by means of representation learning [1] techniques, shifting the following analysis from cumbersome representations to compact vectors in a latent space, and (ii) clustering points associated to instances of the starting dataset to suggest patterns to the domain experts that will evaluate their potential meaning in the studied domain. The steps could, of course, be iterated, with changes in the operating parameters and hyperparameters of the involved algorithms, and proper user interfaces are required to support these activities by end users that are not computer scientists. This kind of investigation certainly has significant relationships with techniques sometimes referred to as *deep clustering* [2], which also mostly (but not exclusively) consider unsupervised techniques employing deep learning approaches and techniques supporting analyses based on clustering (some of which propose a very similar workflow, mostly excluding forms of visual analysis, but clearly comprising forms of representation learning before actual clustering). Nonetheless, we want to emphasize that we need to be as informative as possible, taking a human-in-the-loop approach, and also consider that the human is not necessarily an expert in machine learning techniques. Our goal here is basically to support knowledge creation taking a data-driven perspective, either to suggest or support a form of validation of more theoretically guided approaches.

Another relevant aspect of our approach is cloud orientation. Due to the potentially significant computational requirements of the tasks included in the workflow, cloud architectures can represent extremely promising approaches, also supporting the integration with data repositories complying FAIR[1] principles within an Open Science perspective[2]. This is the context in which the NEANIAS project (Novel EOSC Services for Emerging Atmosphere, Underwater & Space Challenges), and the present work are set. The project is aimed at contributing to the European Open Science Cloud (EOSC)[3], also by means of the development and integration of innovative cross-cutting services for tackling operationally space-related studies. The following Section will elaborate the motivations and context in which this research is set. Sect. 3 will describe the overall workflow and the comprised steps, while Sect 4 will describe experimental results achieved with a preliminary version of the service implementing the workflow. Conclusions and future developments will end the paper.

## 2. Motivation and Context

The European Open Science Cloud (EOSC) is a long term initiative recognised and funded by the Council of the European Union having the ambition to European researchers, innovators, companies and citizens with a federated and open multi-disciplinary environment where they can publish, find and re-use data, tools and services for research, innovation and educational purposes. Until 2020, i.e. within the Horizon 2020 programme, EOSC federated existing research

---

[1]https://www.go-fair.org/fair-principles/
[2]https://www.oecd.org/science/inno/open-science.htm
[3]https://eosc-portal.eu/about/eosc

data infrastructures in Europe, and started the realization of a web of FAIR data and related services for science, making research data interoperable and machine actionable following the FAIR guiding principles [3]. Within this context, the NEANIAS project aims at contributing by developing, integrating, and disseminating to the relevant communities a set of innovative cross-cutting services, in particular for tackling operationally space-related studies. Some of these services employ AI services and, in particular, Machine Learning tools supporting data analysis.

Within the Work Package 4 of the project (Space Research Services), we face several situations that imply forms of analysis of astronomical images that can be framed as supervised machine learning approaches: for instance, images depicting radio maps of the galactic plane in different bands can be analyzed for the detection and classification of sources of signals [4]. In these cases, in general, there are available annotated datasets, and of course experts have accumulated sufficient knowledge supporting annotators in their operations. There are other situations, however, in which studied images are still object of intense study, maybe because the number of considered objects is still not really high (for instance, supernovae remnants [5]), or because knowledge required for annotating astronomical images is still being constructed (the interpretation of certain physical observations is still object of debate). Discussions with domain experts led to consider the possibility to support analyses of the available (and soon to be generated and distributed) datasets by means of unsupervised ML approaches, that could be used to generate visualizations of the overall dataset and comprising images within a representation of a latent space associated to the dataset itself, in the vein of [6] and not far from the extremely recent approach described in [7], or even to perform pattern discovery operations and clustering, suggesting potential classifications schemes. This kind of computer support to domain expert research activities could be useful at least in two cases within NEANIAS (in particular the above mentioned study of infrared and radio images of supernova remnants and also analysis of images associated to compact sources (clumps) that could lead to interesting innovative insights on star formation), but we are confident that this might be developed and deployed as a service of wider interest and applicability within the EOSC.

## 3. Proposed Workflow

The proposed workflow and the set of implied tasks for unsupervised data pattern discovery, in particular in image datasets, is summarized in Figure 1. For simplicity we excluded a data preparation step, that however is often extremely important and that is definitely not trivial within the context of astronomical and multi-spectral images (with serious methodological issues about how to manage normalization especially considering images acquired by means of different types of sensors). Nonetheless, for sake of space, we are going to focus on the steps bringing to the construction of a latent space based on the dataset, and in subsequent phases of analysis of this space, both through visualization techniques and clustering.

While the *representation learning* task must be carried out immediately after data preparation, so as to achieve a latent space from the points associated to images within the dataset, clustering and visual analysis can be considered as potentially parallel tasks: first of all, the user might be proposed immediately a visualization of the structure of the latent space, in which images that
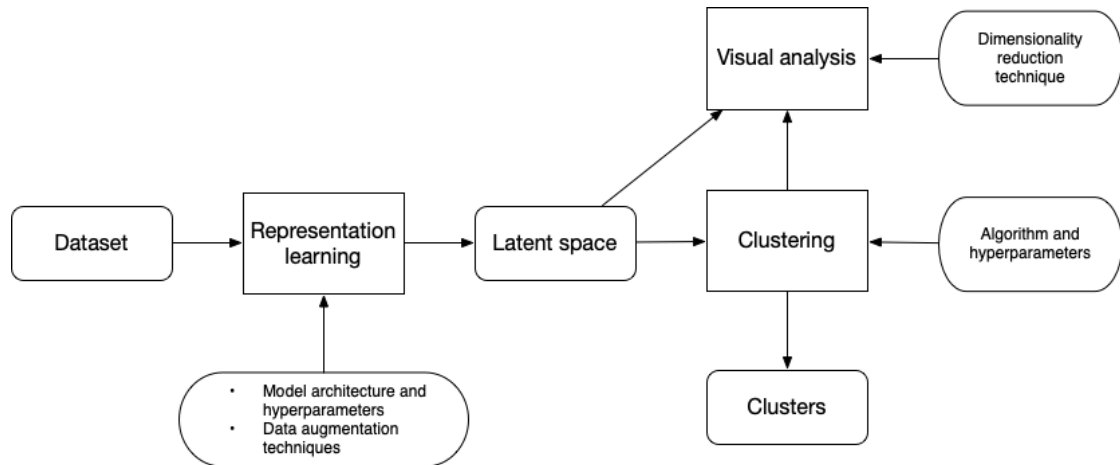
**Figure 1:** Proposed workflow for unsupervised data pattern discovery.

generated it are positioned. To do this, however, since the latent space is typically characterized by a relatively high dimensionality, a specific technique for dimensionality reduction must be adopted to achieve a 2D or 3D map of the latent space. Clustering points in this space does not strictly require a visualization, but, of course, being able to have an at-a-glance representation of clustered images, their neighborhood in the latent space, being able to navigate it and visually inspect the grouping suggested by the clustering algorithm execution can help the user evaluating the achieved results, and their plausibility within the domain of research. In the following we will articulate the techniques we selected for the different above mentioned tasks, that is, representation learning, dimensionality reduction and visualization, and clustering.

## 3.1. Representation learning

A growing number of unsupervised or self-supervised representation learning techniques developed to deal with pictorial data are being developed and experimented to avoid or simplify (partly automating) time consuming and potentially expensive image labelling tasks. Despite the recent interest, one of the earliest approaches for this kind of task is based on Auto Encoders (AE) [8]: within this approach, encoding and decoding neural networks are trained using a loss signal related to the difference between the input and reconstructed image. At the point of contact between encoder and decoder we have a *latent* vector representation of the image presented as input. From this venerable approach, that produced extremely interesting results for noise reduction and image compression limiting the loss of quality, several additional techniques have been generated, in particular Variational Auto Encoders (VAE) [9] and Generative Adversarial Networks (GAN) [10]. While basic AEs are very simple and easily implemented, the overall approach has the final goal of effectively reconstructing the inputs minimizing errors, but the resulting latent space might not help identifying meaningful groups of elements within the dataset that generated the latent space. For related reasons, also VAE and GAN are not generally suited to generate latent spaces simplifying an effective visual analysis and clustering of the starting

dataset[4].

More recent works in the self-supervised representation learning context were specifically devised for simplifying image classification tasks, significantly reducing the need of annotated data. In particular, some self-supervised techniques are based on a two step approach [12]: a first step not requiring supervision is used essentially to provide a first structure of the latent space that is subsequently tuned by using a relatively small labeled dataset. The adopted transfer learning approach implies that the labeled dataset covers relatively well the categories of subjects included in the initial unlabeled dataset. Relevant representatives of this approach are MOCO [13] and SimCLR [14]. Although the last approach is just partly applicable to our specific context, since we do not have any annotated image due to the fact that the classification scheme is unknown, the unsupervised part of SimCLR has been successfully used as a generator of latent space and encoder for a subsequent clustering approach, which we will describe later on. Moreover, there is a recent report of a successful adoption of this technique with a small initial dataset, with small images and small network architectures [15]: in our case, the low requirements at least on the number of images required for training is very important and it suggests this could be a promising representation learning approach.

## 3.2. Latent Space Visualization

In order to support an at-a-glance evaluation of both the structure of the latent space and the intuitive possibility to perform effectively a clustering of the images within the dataset, it can be useful to produce a reduction of the latent space to 2D or 3D spaces. This of course implies a reduction of dimensionality, since the dimension of the latent space is generally quite substantial. Commonly used algorithms for performing this reduction are Principal Component Analysis (PCA) [16], t-distributed stochastic neighbor embedding (TSNE) [17] and Uniform Manifold Approximation and Projection (UMAP) [18]. These techniques can have problems in simultaneously preserving the capability to represent local and global structures within the latent space, so more recent approaches (in particular TriMAP[19], PaCMAP[20], and denseMAP [21]) are working towards an improvement of this aspect.

Besides the techniques, it is relevant to mention the fact that relevant projects aimed at providing fully fledged tools for the visualization of latent spaces and datasets that were used to generate them can be found in the literature. In particular, Latent Space Cartography tool [6] is a general tool providing 2D maps in which images within the dataset are positioned trying to capture "semantic" dimensions (e.g. using a dataset of emojis, the tool can show the interpolation between to points like a happy face and a sad one: moving along the line the smile fades gradually). Moving to our context of application, a very recent work [7] visualizing a latent space of a multidimensional dataset of galaxies generated through a self-supervised approach (SimCLR) and adopting UMAP for sake of dimensionality reduction helps understanding the physics of a certain type of analysed galaxies. In particular, in this case there is no particular discovery, but a confirmation of the presence of two well-known categories of galaxies (rotating main-sequence disks and massive slow rotators) from a purely data driven perspective. This application suggests that the proposed workflow and overall approach is promising also for supporting the analysis of

---

[4]Although the are some GAN based approaches actually focused on this task, such as [11]

situations in which instead classifications are still not present.

## 3.3. Clustering Latent Spaces

We do not intend, nor we think it is reasonable within this paper, to provide a brief introduction to clustering here. In fact, we are not really facing a general clustering problem, since we are actually proposing to analyze a latent space whose structure and "semantics" is actually unknown. We will just briefly report a few relevant already tested or promising approaches that are currently being evaluated, also considering that the project is still ongoing and we are still focused on the representation learning and visualization steps.

Within Sect. 3.1 we already suggested that the SimCLR self-supervised approach has been used to generate a latent space later analysed by means of a clustering approach called SCAN [22]. Within this approach, authors propose training a network performing clustering (i.e. predicting the cluster assignment for a given individual): of course, the number of clusters in which the dataset must be subdivided must be known in advance, similarly as for K-Means. Unlike K-Means, authors make assumptions on the distribution of points in the different clusters (since they are considering image datasets they assume that it is very unlikely the situation in which one or very few classes dominate others), and this is reflected on the defined loss function. Whereas this assumption does not make sense in our case, this represents an interesting success case for clustering a latent space.

Moreover, within our situation, the number of clusters is actually unknown and also their shapes within the latent space. For this reason we think that density based approaches such as DBSCAN [23] could represent a reasonable starting point, although a distance metric should be identified and proper calibration should be carried out. While this aspect can be problematic in a situation in which complete automation is desirable, within our case the domain experts should be kept in the loop and therefore it does not really represent a insurmountable problem, provided that proper support to the selection of relevant hyperparameters is provided.

# 4. Experimental Framework and Results

Starting from the above discussion, we started developing and testing Latent Space Explorer, a cloud oriented framework supporting the workflow described in Figure 1, first of all to support the exploration of the practical issues related to latent space visual analysis, even before moving in the clustering step. In fact, the number of alternative options already in the representation learning and in the latent space dimensionality reduction steps call for an initial focusing on these aspects, also on already existing datasets before moving to the specific context of application at hand.

## 4.1. Architecture

A decade ago a manifesto named "twelve-factor app"[5] was published: it describes a guidelines for building software-as-a-service (SaaS) systems. Most of the proposed best-practices are actually
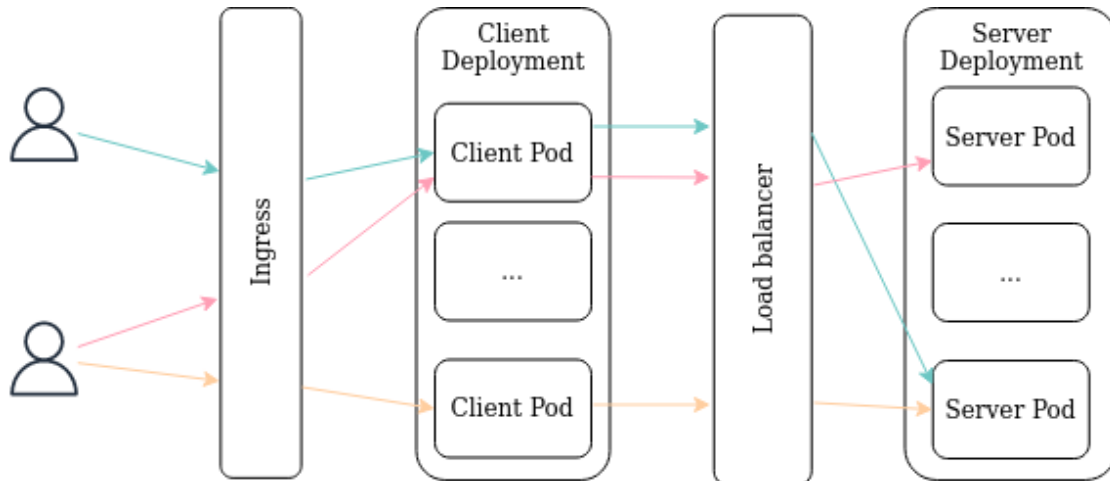
---

[5]https://12factor.net

**Figure 2:** Architecture view focused on *stateless* and *scalablility* features. Two users contact the same endpoint that redirects the requests to most available client resource. The pods does not store any information locally, so the next request does not need to be redirected at the same pod. The same happens for client server communication.

fulfilled by design through a proper adoption of Kubernetes[6] and by current CI/CD systems, while other ones still need attention in the design of an architecture.

Looking at Figure 2, two related features are highlighted: *stateless* processes and process *scalability*. If each pod receives complete information required to fulfill an incoming request without saving any information, then the requester does not need to be aware of which pod to contact, and therefore it could submit the request to any available pod. This feature, therefore, helps scaling pods as needed.

Another important feature for making cloud native applications more in tune with the "twelve-factor app" manifesto is related to the management of enabling services. Figure 3, on the top portion, represents an example of a a situation in which the actual system architecture depends on unique services for tasks such as cloud storage and scheduling. On the other hand, if for some reason a project needs to migrate to another storage service, this dependency leads to hard refactoring of the code to support the migration. The "twelve-factor app" design guide suggests to separate the management of supporting services into another component, in order to isolate code changes.

The approach adopted in the desing and implementation of the Latent Space Explorer tries to combine a modern cloud oriented approach with the two above principles. The first version of the service is publicly available for usage[7] and a small user guide has also been released[8]. The overall project is going to be released adopting an open source license (still to be defined), and all the relevant information will also be made available through the above links and through the NEANIAS project communication channels.

---

[6]Kubernetes is an open-source container-orchestration system for automating computer application deployment, scaling, and management - See: https://kubernetes.io/

[7]https://lse.neanias.eu/

[8]https://docs.neanias.eu/projects/s3-service/en/latest/services/latent_space_explorer.html
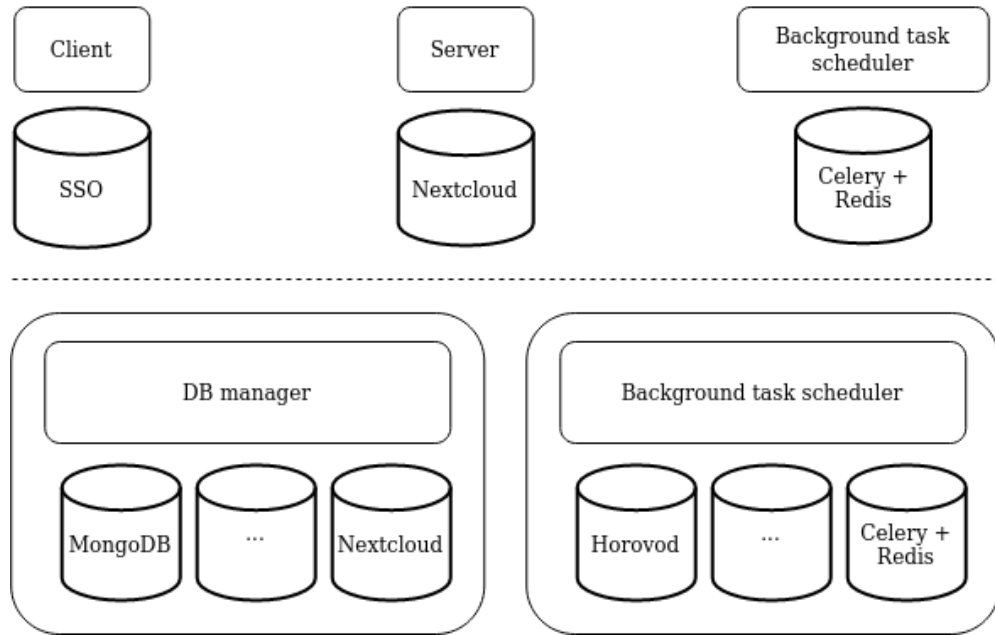
**Figure 3:** Architecture view focused on *backing services*

## 4.2. Workflow and Visualization

While the framework is still under active development and it has still not reached the desired level for the first release, it can already be used to support forms of visualizations of reductions (to 2D and 3D spaces) of a latent space resulting from the representation learning step of the proposed workflow. A sample screenshot of the working version of Latent Space Explorer is shown in Figure 4: the left column shows information about the specific experiment to perform and visualize, allowing the specification/selection of the dataset, associated latent space, dimensionality reduction technique, and even clustering algorithm (each with the associated hyperparameters when applicable). The central part is the navigable visualization of the latent space, including points associated to the images of the dataset. The right column shows details about the selected point, or statistics about the latent space and/or clustering.

Within this phase of the work, the main value of the potential usage of the tool is actually the visualization of the 2D or 3D reduction of a latent space, to have an at-a-glance idea of the structure achieved through the representation learning step. Figure 5 shows a visualization of the UMAP reduction the latent space associated to the Standford Dogs dataset [24] (a subset of Imagenet with more fine-grained labels that represent 120 dog breeds) considering the latent space associated to a resnet158 3x pre-trained on Imagenet with SimCLRV2. The color of the points is associated to the actual breed of the dog associated to the image, and it can be easily seen that visually peculiar dogs are positioned in areas of the space that are well separated from the others; visual features of dogs are of course most relevant (it seems clear, for instance, that fur length is short in the left hand side and long in the right hand side of the space). Nonetheless, contextual elements (e.g. snowy background or a beach, presence of a cage) are also plausibly relevant in
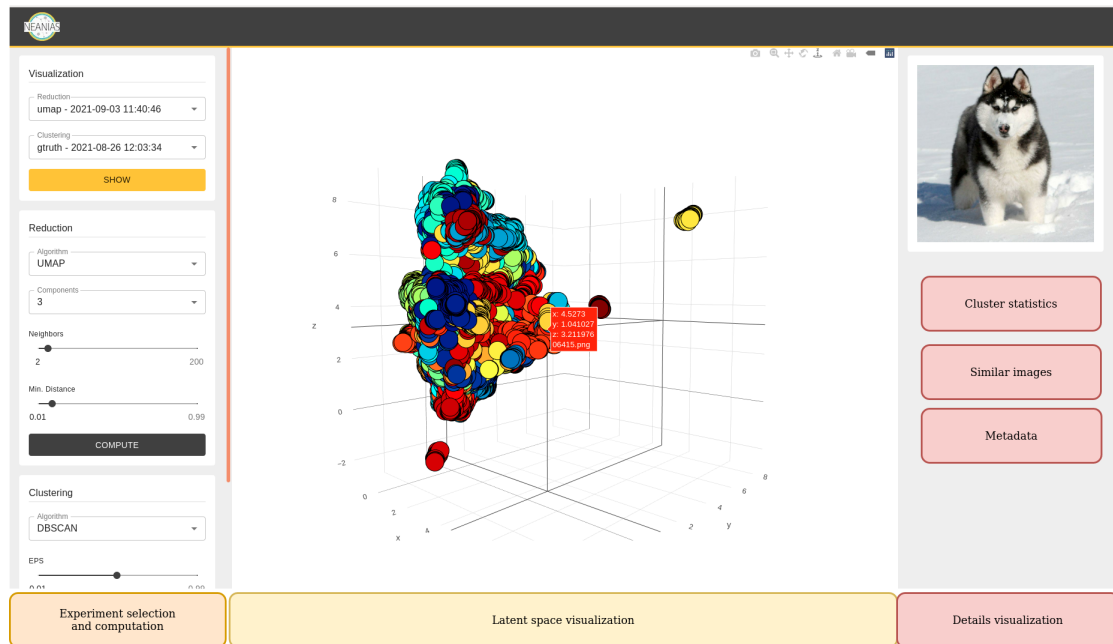
**Figure 4:** A screenshot of Latent Space Explorer.

determining the positioning of points in space. Although the visual inspection of this kind of latent space reveals potentially useful information to a viewer, clustering this kind of space with the aim of identifying dog breeds would hardly produce immediately usable results. Let us consider, however, different network architectures, to have an idea of the effect of changing the size of the network on the capability of better representing "semantic" aspects of the dataset. Figure 6 shows a comparison of different latent spaces associated to different network architectures, selected from the list of available pre-trained (on Imagenet) networks. The base model is a Resnet [25] with different hyperparameters: in particular, use of selective kernel (SK) [26], different depths [50,101,158], and width [1x, 3x]. In general we can see that latent spaces associated to deeper and wider networks are generally better at discriminating dog breeds, plausibly being able to grasp more information about the dataset and comprised images. Even the selective kernel achieves better results.

The visual analysis of the reduction of the latent space to a 2D/3D structure is therefore clearly important to have an idea of how the result of the representation learning step. It is very reasonable to inspect this kind of visualization, maybe trying different alternative representation learning approaches (or sets of values for the hyperparameters) before actually even trying to move to the clustering phase. Our situation, the analysis of astronomical images, is characterised by smaller datasets (although, as suggested in the introduction, initiatives like the Square Kilometre Array will plausibly change the situation significantly) of images that are much less cluttered than those present in Imagenet (and in its subset represented by the Stanford Dogs). On the one hand, the representation learning approaches will have to face less "distractions" (although the data preparation phase will necessary face previously mentioned challenges related to normalization
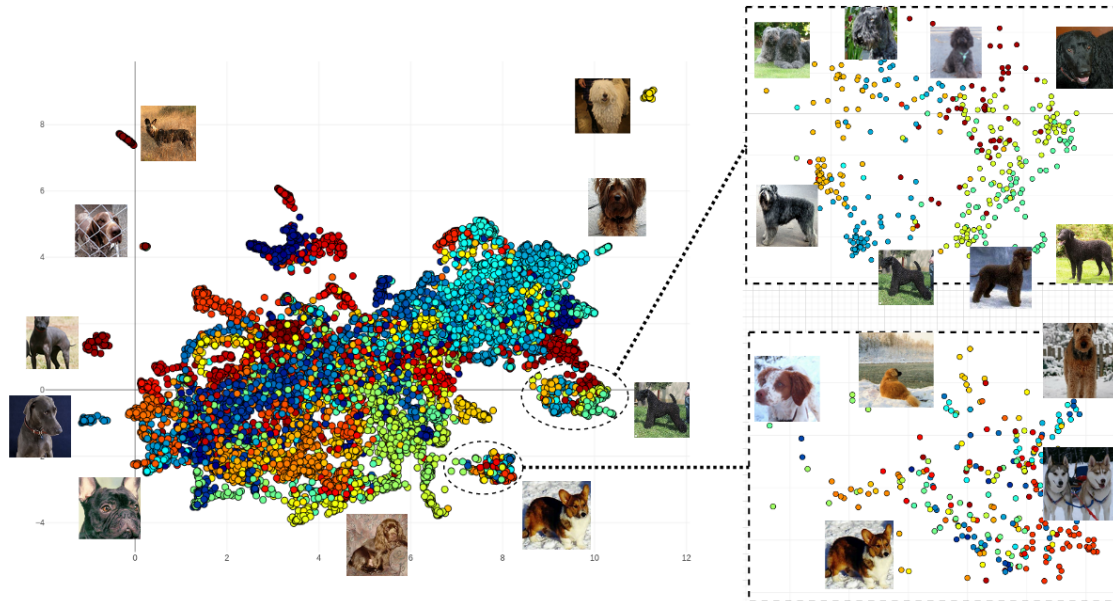
**Figure 5:** Visualization of the UMAP reduction ($mindistance = 0.1$, $n_neighbors = 15$) of the latent space associated to the Standford Dogs dataset using Resnet152 with layer widths 3x, pre-trained on Imagenet with SimCLRV2.

within a multi-spectral images context), on the other we will not be able to use large deep neural network architectures (the training would simply not converge with so little data available), and this could make it hard to process some fine details that might be extremely important. In order to produce interesting results we will need to consider recent developments within the self-supervised techniques [15] context that promise to be able to provide useful results even with small datasets (of potentially small images), with small networks.

## 5. Conclusions

The paper has presented the current status of development of a project aimed at supporting unsupervised data pattern discovery within datasets of images, particulrarly focused on astronomical images, but in general usable for performing workflows including representation learning, leading to the definition of a latent space associated to the dataset, and forms of analysis ranging from visualization of 2D/3D reductions of the latent space and even clustering. We have set the work within the context of the NEANIAS project and more generally within EOSC, and we have discussed the motivation, rationale of the approach, discussing the relevant state of the art. The current status of the project was presented, and its current potential usage to visually inspect latent spaces was discussed. We tried to clarify why the visual analysis of the latent space structure is important, not just because a human-in-the-loop perspective is crucial within a knowledge discovery and creation for scientific research scenario, but also since it is actually useful to have a preliminary idea of the potential outcomes of clustering operations on the latent space. Current and future works on the project are focused on (i) a preliminary evaluation of the potential
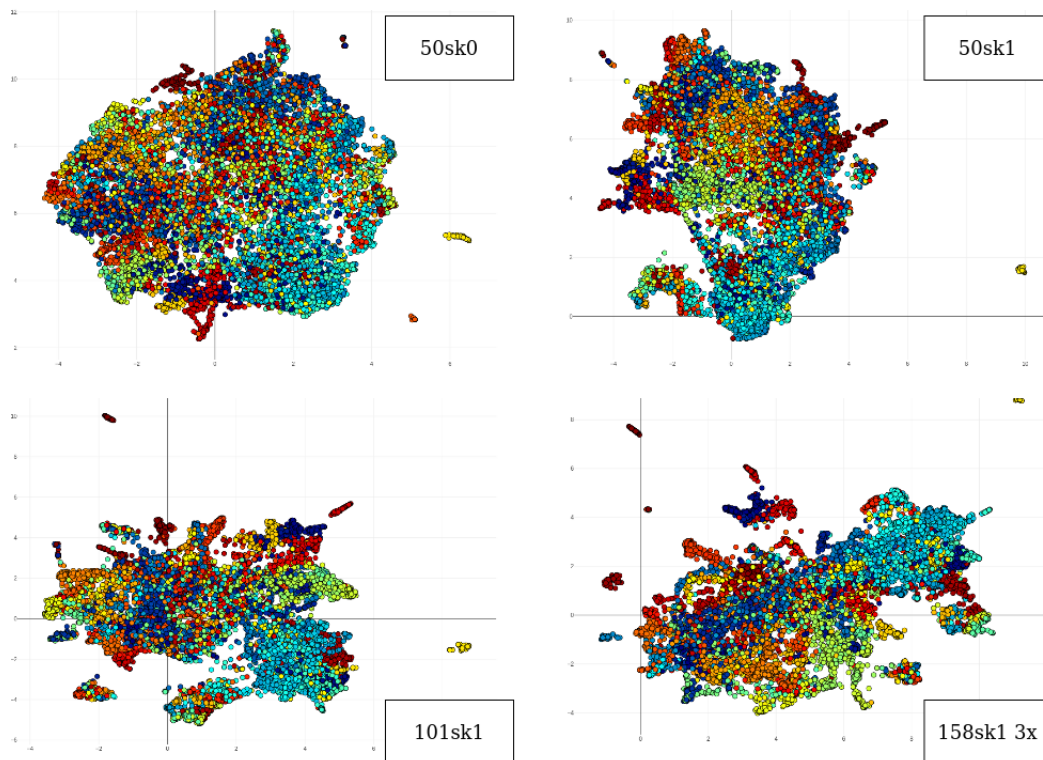
**Figure 6:** A comparison of different latent spaces representations including actual ground truth (based on the Stanford dogs dataset).

results employing target datasets related to infrared and radio images of supernova remnants and compact sources (clumps) that could lead to interesting innovative insights on star formation, (ii) the adoption of techniques for the representation learning step that are particularly suited to deal with small datasets. In the medium run, we also intend to release the produced framework as a service within the NEANIAS project ecosystem and also as an open source project.

## Acknowledgments

## References

[1] Y. Bengio, A. Courville, P. Vincent, Representation learning: A review and new perspectives, IEEE Transactions on Pattern Analysis and Machine Intelligence 35 (2013) 1798–1828.

[2] E. Min, X. Guo, Q. Liu, G. Zhang, J. Cui, J. Long, A survey of clustering with deep learning: From the perspective of network architecture, IEEE Access 6 (2018) 39501–39514.

[3] M. D. Wilkinson, M. Dumontier, I. J. Aalbersberg, G. Appleton, M. Axton, A. Baak, N. Blomberg, J.-W. Boiten, L. B. da Silva Santos, P. E. Bourne, J. Bouwman, A. J. Brookes, T. Clark, M. Crosas, I. Dillo, O. Dumon, S. Edmunds, C. T. Evelo, R. Finkers, A. Gonzalez-Beltran, A. J. G. Gray, P. Groth, C. Goble, J. S. Grethe, J. Heringa, P. A. C. 't Hoen, R. Hooft, T. Kuhn, R. Kok, J. Kok, S. J. Lusher, M. E. Martone, A. Mons, A. L. Packer, B. Persson, P. Rocca-Serra, M. Roos, R. van Schaik, S.-A. Sansone, E. Schultes, T. Sengstag, T. Slater, G. Strawn, M. A. Swertz, M. Thompson, J. van der Lei, E. van Mulligen, J. Velterop, A. Waagmeester, P. Wittenburg, K. Wolstencroft, J. Zhao, B. Mons, The FAIR guiding principles for scientific data management and stewardship, Scientific Data 3 (2016) 160018.

[4] S. Riggi, F. Vitello, U. Becciani, C. Buemi, F. Bufano, A. Calanducci, F. Cavallaro, A. Costa, A. Ingallinera, P. Leto, et al., Caesar source finder: Recent developments and testing, Publications of the Astronomical Society of Australia 36 (2019) e037.

[5] R. Z. E. Alsaberi, L. A. Barnes, M. D. Filipović, N. I. Maxted, H. Sano, G. Rowell, L. M. Bozzetto, S. Gurovich, D. Urošević, D. Onić, B. Q. For, P. Manojlović, G. Wong, T. J. Galvin, P. Kavanagh, N. O. Ralph, E. J. Crawford, M. Sasaki, F. Haberl, P. Maggi, N. F. H. Tothill, Y. Fukui, Radio emission from interstellar shocks: Young type ia supernova remnants and the case of n 103b in the large magellanic cloud, Astrophysics and Space Science 364 (2019) 204.

[6] Y. Liu, E. Jun, Q. Li, J. Heer, Latent space cartography: Visual analysis of vector space embeddings, Computer Graphics Forum 38 (2019) 67–78.

[7] R. Sarmiento, M. Huertas-Company, J. H. K. S. F. Sànchez, H. D. Sànchez, N. Drory, J. Falcón-Barroso, Capturing the physics of MaNGA galaxies with self-supervised machine learning, The Astrophysical Journal (to appear). URL: https://arxiv.org/abs/2104.08292.

[8] P. Baldi, K. Hornik, Neural networks and principal component analysis: Learning from examples without local minima, Neural Networks 2 (1989) 53–58.

[9] D. P. Kingma, M. Welling, Auto-encoding variational bayes, in: Y. Bengio, Y. LeCun (Eds.), 2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings, 2014.

[10] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. C. Courville, Y. Bengio, Generative adversarial nets, in: Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, K. Q. Weinberger (Eds.), Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada, 2014, pp. 2672–2680.

[11] J. Donahue, P. Krähenbühl, T. Darrell, Adversarial feature learning, in: 5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings, 2017.

[12] L. Jing, Y. Tian, Self-supervised visual feature learning with deep neural networks: A survey, IEEE Transactions on Pattern Analysis and Machine Intelligence 43 (2021) 4037–4058. ((2021)).

[13] K. He, H. Fan, Y. Wu, S. Xie, R. B. Girshick, Momentum contrast for unsupervised visual representation learning, in: 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020, Computer Vision

Foundation / IEEE, 2020, pp. 9726–9735.

[14] T. Chen, S. Kornblith, K. Swersky, M. Norouzi, G. E. Hinton, Big self-supervised models are strong semi-supervised learners, in: H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, H. Lin (Eds.), Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual, 2020.

[15] Y. Cao, J. Wu, Rethinking self-supervised learning: Small is beautiful, CoRR abs/2103.13559 (2021). URL: https://arxiv.org/abs/2103.13559.

[16] H. Hotelling, Analysis of a complex of statistical variables into principal components., Journal of educational psychology 24 (1933) 417.

[17] L. Van der Maaten, G. Hinton, Visualizing data using t-sne., Journal of machine learning research 9 (2008).

[18] L. McInnes, J. Healy, N. Saul, L. Grossberger, UMAP: Uniform manifold approximation and projection, The Journal of Open Source Software 3 (2018) 861.

[19] E. Amid, M. K. Warmuth, Trimap: Large-scale dimensionality reduction using triplets, arXiv preprint arXiv:1910.00204 (2019).URL: https://arxiv.org/abs/1910.00204.

[20] Y. Wang, H. Huang, C. Rudin, Y. Shaposhnik, Understanding how dimension reduction tools work: an empirical approach to deciphering t-SNE, UMAP, TriMAP, and PaCMAP for data visualization, arXiv preprint arXiv:2012.04456 (2020).URL: https://arxiv.org/abs/2012.04456.

[21] A. Narayan, B. Berger, H. Cho, Density-preserving data visualization unveils dynamic patterns of single-cell transcriptomic variability, bioRxiv (2020). URL: https://www.biorxiv.org/content/early/2020/05/14/2020.05.12.077776.

[22] W. V. Gansbeke, S. Vandenhende, S. Georgoulis, M. Proesmans, L. V. Gool, SCAN: learning to classify images without labels, in: A. Vedaldi, H. Bischof, T. Brox, J. Frahm (Eds.), Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part X, volume 12355 of *Lecture Notes in Computer Science*, Springer, 2020, pp. 268–285.

[23] M. Ester, H. Kriegel, J. Sander, X. Xu, A density-based algorithm for discovering clusters in large spatial databases with noise, in: E. Simoudis, J. Han, U. M. Fayyad (Eds.), Proceedings of the Second International Conference on Knowledge Discovery and Data Mining (KDD-96), Portland, Oregon, USA, AAAI Press, 1996, pp. 226–231.

[24] A. Khosla, N. Jayadevaprakash, B. Yao, F.-F. Li, Novel dataset for fine-grained image categorization: Stanford dogs, in: Proc. CVPR Workshop on Fine-Grained Visual Categorization (FGVC), volume 2, 2011.

[25] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 770–778.

[26] X. Li, W. Wang, X. Hu, J. Yang, Selective kernel networks, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2019, pp. 510–519.