

---

# Unsupervised local structure learning in elemental zirconium using Topological Data Analysis

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We discuss in this extended abstract a new application in materials science where  
2 topological data analysis (TDA) has proven to be relevant. TDA is used to analyze  
3 local atomic structures during homogeneous nucleation phenomena of pure zirco-  
4 nium in the undercooled liquid namely below its melting temperature. Persistent  
5 diagrams (PDs) are considered as descriptors to represent configurations of the  
6 physical system generated by large-scale molecular dynamics (MD) simulations  
7 up to one million atoms, and are evaluated by comparison to classical physical de-  
8 scriptors. This work is a first step towards a deeper understanding of the nucleation  
9 behavior of the system, where topological machine learning is a challenging but  
10 promising perspective.

## 11 1 Context

12 **Physical motivation** As crystal nucleation of matter from the liquid state is essentially initiated  
13 by heterogeneous nucleation from impurities, surfaces, or near grain boundaries, understanding the  
14 homogeneous crystal nucleation, driven by complex phenomena difficult to capture experimentally,  
15 remains a challenging issue. In [3], the glass formation (GF) has been studied as well as the onset of  
16 crystallization features of pure zirconium as a function of the degree of undercooling by means of  
17 large-scale MD simulations up to one million atoms. The GF results from a geometrical frustration  
18 and a strong competition between the crystal and icosahedral structural orderings. Setting up a  
19 topological machine learning approach to understand how this structural competition impacts the  
20 crystallization kinetics is the main aim of this work.

21 **Data** Homogeneous nucleation is observed along an isotherm at  $T = 1250$  K at constant ambient  
22 pressure between the melting and the glass transition temperature ( $T = 890$  K), corresponding to the  
23 lowest nucleation time [3] (see Fig. 2 in Supplementary Material). Configurations are recorded and  
24 their inherent structures are produced by an energy minimization using a conjugate gradient. From  
25 those data, the goal is to characterize, in an unsupervised way, the set of local structures, namely a  
26 central atom with its first neighbours up to  $4.4 \text{ \AA}$  (corresponding to the first minimum of the pair  
27 correlation function and thus defining the coordination sphere, see Fig. 3 in Supplementary Material).

**Descriptors** To study these local structures we need descriptors which are invariant by translation and rotation. Usually in material science, the following classical descriptors are used: the Bond Angle Analysis (BAA) [1], the Bond-Orientational Order Analysis (BOOA) [13, 10] and the Common Neighbour Analysis [9] among others [2]. They characterize the radial and/or angular distribution of bonded pairs of atoms in the coordination sphere. In this work, we use for the first time TDA, and more precisely persistent homology, to study the topology of these local structures. The local environments of atoms are seen as a points cloud, and we draw PDs from them. PDs are usually

compared using bottleneck or Wasserstein distances, which are stable under small perturbations of the point clouds [7]. However, the induced metric spaces are not isometrically embeddable (even with a bi-lipschitz map) into a finite dimensional Hilbert space [5]. We propose to use the signature introduced in [6], classically used to study 3D shapes, in order to get a stable representation of our PDs into a Hilbert space. For each pair of points  $(x, y)$  in a PD  $D$ , we consider

$$m_D(x, y) = \min\{\|x - y\|_\infty, d_\Delta(x), d_\Delta(y)\},$$

28 where  $d_\Delta(\cdot)$  denotes the  $\ell^\infty$  distance to the diagonal. The signature of  $D$  corresponds then to the  
 29 vector with all those values sorted in decreasing order. Remark that the topological signature obtained  
 30 here and the usual physical descriptors are highly correlated (see Fig. 5 in Supplementary Material).

## 31 2 Method and results

32 The topological signature is computed for 20 000 local structures, subsampled from the data generated  
 33 in [3]. PDs have been denoised by removing points too close to the diagonal, using a threshold  
 34 learned by bootstrap. Then, the unsupervised learning is done using a Gaussian Mixture Model  
 35 (GMM), with full covariance matrices, fitted by an Expectation-Maximization algorithm (EM) on the  
 36 topological descriptors. The number of clusters is selected using the Integrated Completed Likelihood  
 37 criterion [4] (ICL)<sup>1</sup>. From our data, we get a signature vector with at most 17 coefficients for  $H_0$   
 38 and at most 18 coefficients for  $H_1$ . This vector is illustrated by drawing several PDs in Fig. 4 in  
 39 Supplementary Material. The GMM fits a model with 39 clusters (see Fig. 6 to see the ICL criterion  
 40 in Supplementary Material).

41 Since our analysis is performed on the structure obtained during  
 42 crystallization, we focus our attention on body-centered cubic  
 43 (BCC) structures which are described in CNA by the signature  
 44 [444]: 6 and [666]: 8. Among the 39 clusters, 6 of them have  
 45 a pure BCC structure from the CNA signature (see Fig. 8 in  
 46 Supplementary Material), thus would have been merged when  
 47 only considering CNA; but have different signatures in TDA  
 48 illustrated in Fig. 9 in Supplementary Material. It illustrates the  
 49 strength of TDA, which is able to detect perfect BCCs (where  
 50 coefficients should be zero, as all the points in  $H_0$  and  $H_1$  should  
 51 be equal), as well as their distorted counterparts due to temper-  
 52 ature effects in the simulation where coefficients are nonzero.  
 53 From Fig. 9, we remark that the cluster 16 seems to be the most  
 54 perfect, with the smallest coefficients.

55 By analyzing the physical system, we notice that several nuclei  
 56 are forming, showing the crystallization of the system. We rep-  
 57 resent particles with a complete BCC local structure in Fig. 1,  
 58 where particles in white are from cluster 16, and in red from  
 59 other pure BCC clusters. Particles from cluster 16 corresponds  
 60 to the center of the nuclei, and particles from clusters with higher  
 61 TDA signatures (clusters 1, 4, 14, 23, 37) surround them. In par-  
 62 ticular, this shows that TDA detects a gradual distortion of BCC local structures which are physically  
 63 relevant.

64 Clusters with BCCs but not pure are also interpretable from the CNA signature. We detect that they  
 65 are composed of structures close to BCC, with missing bonds or neighbours. They are also mainly  
 66 outside the main nuclei in the 3D representation. See Fig. 10 and Fig. 11 in Supplementary Material.

## 67 3 Conclusion and future work

68 This application shows a first attempt to understand the local structures which arise in elemental  
 69 zirconium during solidification, using TDA as a relevant descriptor. The TDA is found to represent  
 70 thin important patterns for local structure, even more than classical material science features, which  
 71 are discriminated using a GMM. This first step show good advantages of the descriptor, and we plan  
 72 to take benefit of it also for a study of the global structure.

<sup>1</sup>In practice, we use Python with Gudhi [11] and scikit-learn [12] libraries.

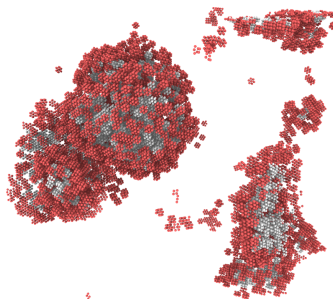


Figure 1: Representation of the system obtained with the OVITO software [14]. We plot only the pure BCC clusters for clarity. The perfect cluster is in the center of the nuclei, whereas higher values in TDA correspond to particles around the nuclei. Other particles (not BCC) are outside of the nuclei.

73 **References**

- 74 [1] G. J. Ackland and A. P. Jones, "Applications of local crystal structure measures in exper-  
75 iment and simulation," *Phys. Rev. B*, vol. 73, no. 5, p. 054104, 2006, doi: 10.1103/Phys-  
76 RevB.73.054104.
- 77 [2] A. P. Bartók, R. Kondor, G. Csányi. "On representing chemical environments", *Phys. Rev. B*,  
78 vol. 87, no. 18, 2013.
- 79 [3] S. Becker, E. Devijver, R. Molinier, and N. Jakse, "Glass-forming ability of elemental zirco-  
80 nium", *Phys. Rev. B*, vol. 102, 104205, 2020, doi: 10.1103/PhysRevB.102.104205
- 81 [4] C. Biernacki, G. Celeux and G. Govaert, "Assessing a mixture model for clustering with  
82 the integrated completed likelihood," in *IEEE Transactions on Pattern Analysis and Machine*  
83 *Intelligence*, vol. 22, no. 7, pp. 719-725, 2000, doi: 10.1109/34.865189.
- 84 [5] M. Carrière, U. Bauer. "On the Metric Distortion of Embedding Persistence Diagrams into  
85 Separable Hilbert Spaces", 35th International Symposium on Computational Geometry (SoCG  
86 2019), vol. 129, pp. 21:1–21:15, 2019, doi:10.4230/LIPIcs.SoCG.2019.21.
- 87 [6] M. Carrière, S.Y. Oudot and M. Ovsjanikov, "Stable Topological Signatures for Points on 3D  
88 Shapes", *Computer Graphics Forum*, 34: 1-12, 2015, doi:10.1111/cgf.12692
- 89 [7] D. Cohen-Steiner, H. Edelsbrunner, J. Harer. "Stability of Persistence Diagrams", *Discrete &*  
90 *Computational Geometry*, vol. 37, no. 1, pp. 103-120, 2007, doi:10.1007/s00454-006-1276-5.
- 91 [8] D. Faken and H. Jónsson, "Systematic analysis of local atomic structure combined with 3D  
92 computer graphics," *Computational Materials Science*, vol. 2, no. 2, pp. 279-286, 1994, doi:  
93 10.1016/0927-0256(94)90109-0.
- 94 [9] J. D. Honeycutt and H. C. Andersen, "Molecular dynamics study of melting and freezing of  
95 small Lennard-Jones clusters," *J. Phys. Chem.*, vol. 91, no. 19, pp. 4950-4963, 1987, doi:  
96 10.1021/j100303a014.
- 97 [10] W. Lechner and C. Dellago, "Accurate determination of crystal structures based on averaged  
98 local bond order parameters," *The Journal of Chemical Physics*, vol. 129, no. 11, p. 114707,  
99 2008, doi: 10.1063/1.2977970.
- 100 [11] C. Maria, J.-D. Boissonnat, M. Glisse, M. Yvinec. "The Gudhi Library: Simplicial Complexes  
101 and Persistent Homology", *Mathematical Software, ICMS*, 167-174, 2014.
- 102 [12] F. Pedregosa et al., "Scikit-learn: Machine Learning in Python", *JMLR* 12, pp. 2825-2830,  
103 2011.
- 104 [13] P. J. Steinhardt, D. R. Nelson, and M. Ronchetti, "Bond-orientational order in liquids and  
105 glasses," *Phys. Rev. B*, vol. 28, no. 2, pp. 784-805, 1983, doi: 10.1103/PhysRevB.28.784.
- 106 [14] A. Stukowski, "Visualization and analysis of atomistic simulation data with OVITO—the Open  
107 Visualization Tool," *Modelling Simul. Mater. Sci. Eng.*, vol. 18, no. 1, p. 015012, Jan. 2010,  
108 doi: 10.1088/0965-0393/18/1/015012.

## 109 A Data

### 110 A.1 Physical context

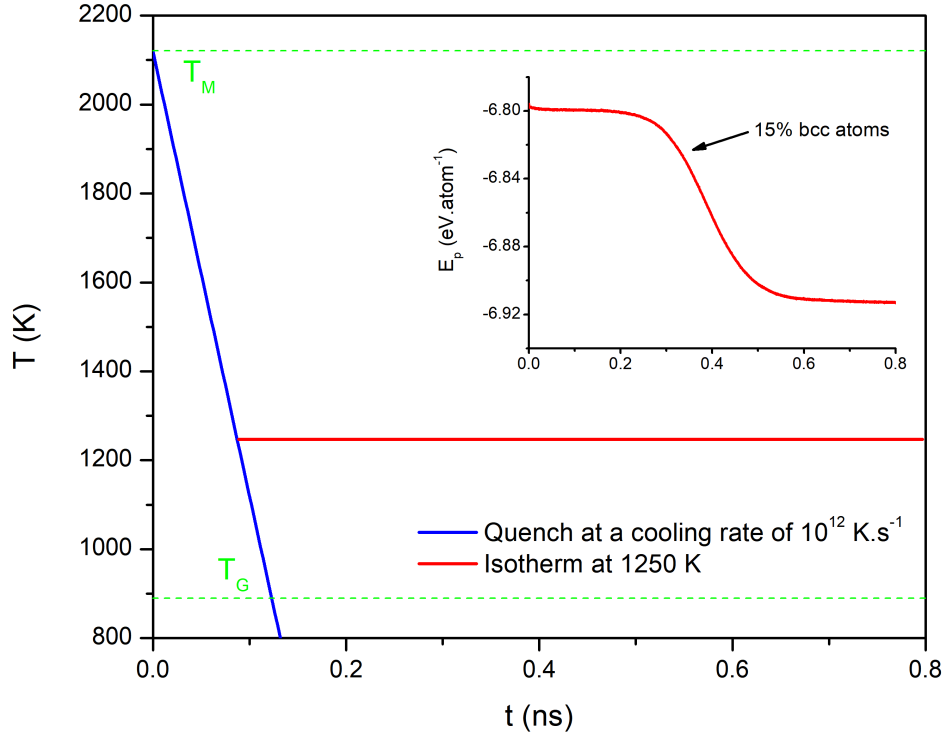


Figure 2: Time-temperature schematic representation of the quench and isotherm modelised by molecular dynamics simulations. The inset represents the evolution of the total potentiel energy along the isotherm at 1250 K.

111 To study the solidification process of a material, a rapid quench from the liquid to the glassy state  
112 must be performed. Then, along this quench, several metastable undercooled configurations of the  
113 simulation box, at specific times and temperatures, are extracted in order to do isotherms and let  
114 them evolve, through nucleation, to the solid stable state. For this study, we have extracted one  
115 configuration at the onset of crystallization, where 15 % of solid bcc atoms have been formed. The  
116 Fig. 2 summarize this steps, with  $T_M$  the thermodynamic melting temperature,  $T_G$  the glass transition  
117 temperature and in inset the evolution of the potential energy of the isotherm at 1250 K through the  
118 time of nucleation.

119 For an atom set at the origin, the radial distribution function  $r \mapsto g(r)$  gives the distribution of the  
120 other atoms in a sphere around it. A typical appearance of this function is shown for the liquid state of  
121 our material in Fig. 3. Graphically, one can observe a succession of peaks, each of which represents  
122 successive layers of neighbours around the central atom. In the case of the local structure, the first  
123 minimum of this function is set as a cut-off radius delimiting the end of the first neighbours shell.

### 124 A.2 Descriptors

125 We plot 3 persistence diagrams in Fig. 4 that reflect the diversity in our data. We consider an  
126 observation in cluster 16, the perfect one discussed in the main paper, an observation in cluster 14  
127 which is a pure BCC cluster with higher values in the TDA signature, and an observation in cluster 6

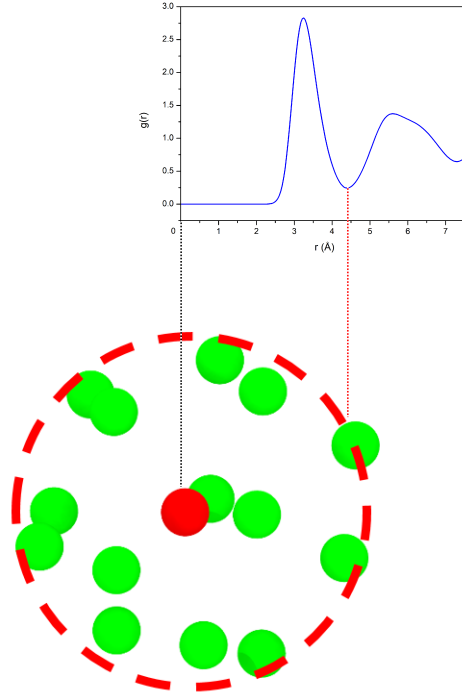


Figure 3: First neighbours shell of a central particle (in red) with a cut-off radius determined by the radial distribution function  $r \mapsto g(r)$  of undercooled Zirconium at 1250 K.

128 which is a cluster with BCC but also other local structures. More precisely, we select the observation  
 129 which is the closest to the center of each cluster.

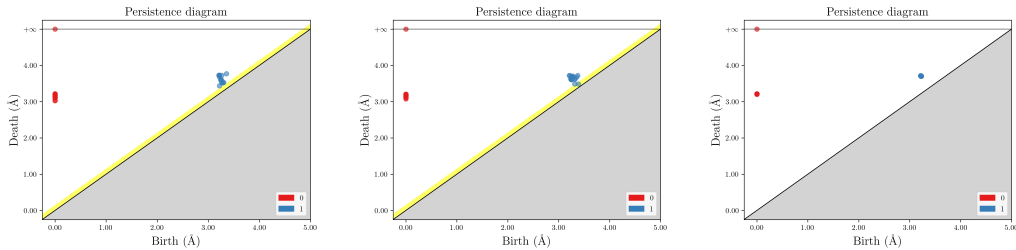


Figure 4: Persistence diagram on our data for observations in the clusters 6 (left), 14 (middle) and 16  
 (right). The yellow band is representing the noise detected by bootstrap.

129

130 To compare TDA descriptor with physical descriptors, we plot the empirical correlation matrix  
 131 between BAA, CNA, TDA and BOAA in Fig. 5. We remark that TDA is highly correlated with all  
 132 the other descriptors.

133 The correlation matrix associated to the TDA shows some specificity of the data: all local neighbours  
 134 contain at least 10 particles, so the first 10 components of  $H_0$  are highly correlated. There are  
 135 correlated with high values of  $H_1$ , because most of the local structures have few  $H_1$  components  
 136 (50% of the population have less than 5 components).

137 The correlation matrix associated to the CNA is sparse, because only few coefficients (the first ones)  
 138 are present in most of the structures. It also explains the low correlation between the last coefficients  
 139 of CNA with the other descriptors.

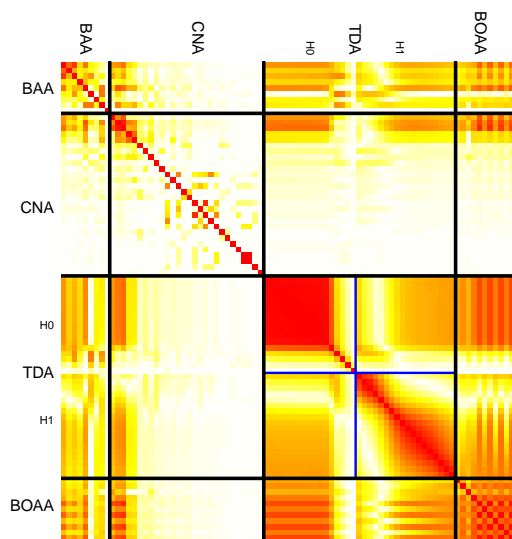


Figure 5: Empirical correlation matrix (absolute value) between several descriptor families: BAA, CNA, TDA and BOAA. Low values are white, high values are red. We distinguish between the descriptors by black lines, and between levels of homology with blue line ( $H_0$  and  $H_1$ ).

140 **B Method**

141 **B.1 Number of clusters**

142 We plot the ICL criterion in Fig. 6 to illustrate the model selection criterion. The criterion seems to  
143 be constant from a large enough dimension, but we smooth this curve (by a polynomial function) and  
144 compare several model selection criterion (BIC, ICL and the Slope Heuristics), which were focusing  
on the same value, 39 clusters.

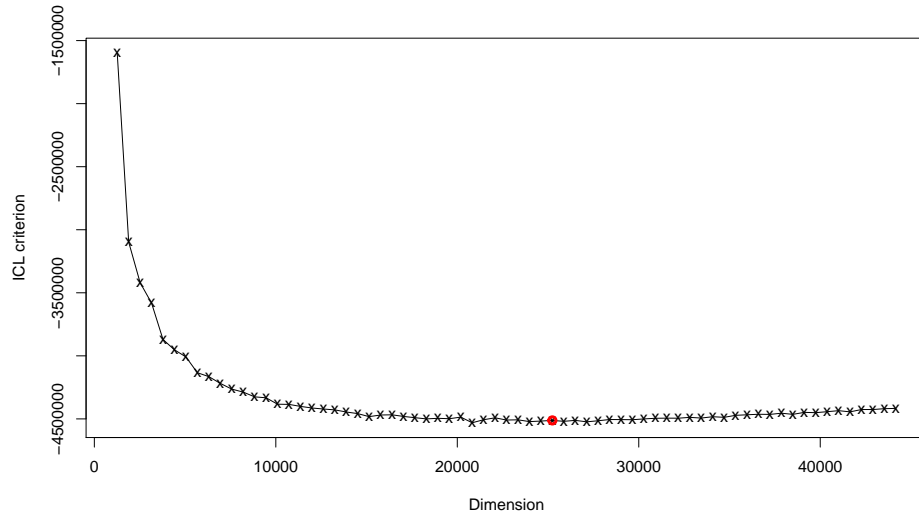


Figure 6: ICL criterion for each model in the collection with respect to the dimension of the model (number of parameters to be estimated). The red point corresponds to the selected model.

145

146 **B.2 Pure clusters of BCC**

147 In Fig. 7, we provide boxplots for CNA and TDA descriptor for all the observations. It illustrates  
148 the diversity of the structures we are looking at, despite of the focus we are doing on BCC. It also  
149 illustrates the analysis we did for the empirical correlation matrices, with range of values for each  
150 descriptor.

151 In Fig. 8 and 9, we provide boxplots for the 6 clusters with only BCC structures. Although there is  
152 no difference on the CNA descriptor (there are labeled BCC), the TDA signatures are very different  
153 from a cluster to another. We particularly remark that the cluster 16 has very small values. We  
154 have illustrated in Fig. 1 that they are corresponding to the structure which are inside the nuclei and  
155 surrounded by structure from other BCC clusters.

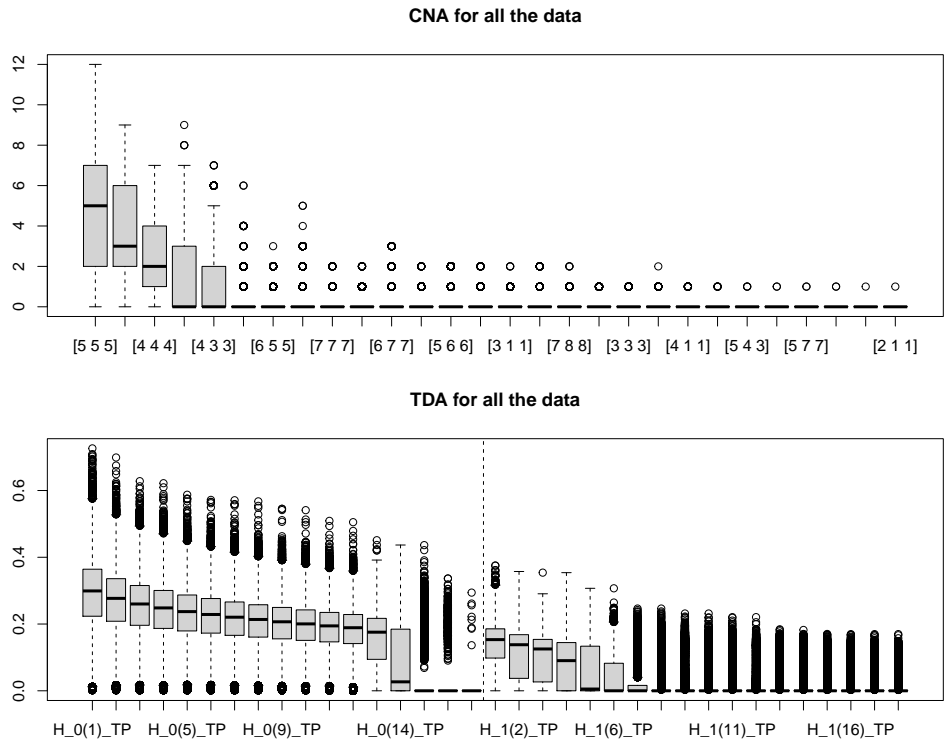


Figure 7: Boxplot of CNA and TDA descriptors for all the observations.

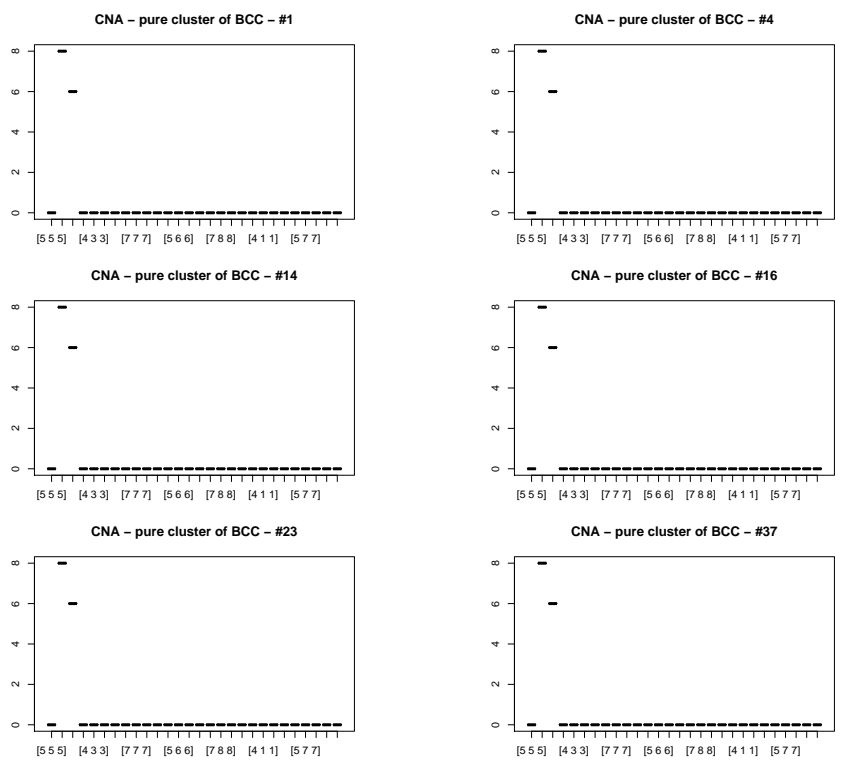


Figure 8: Boxplot for the CNA descriptor for the 6 pure clusters with BCCs.

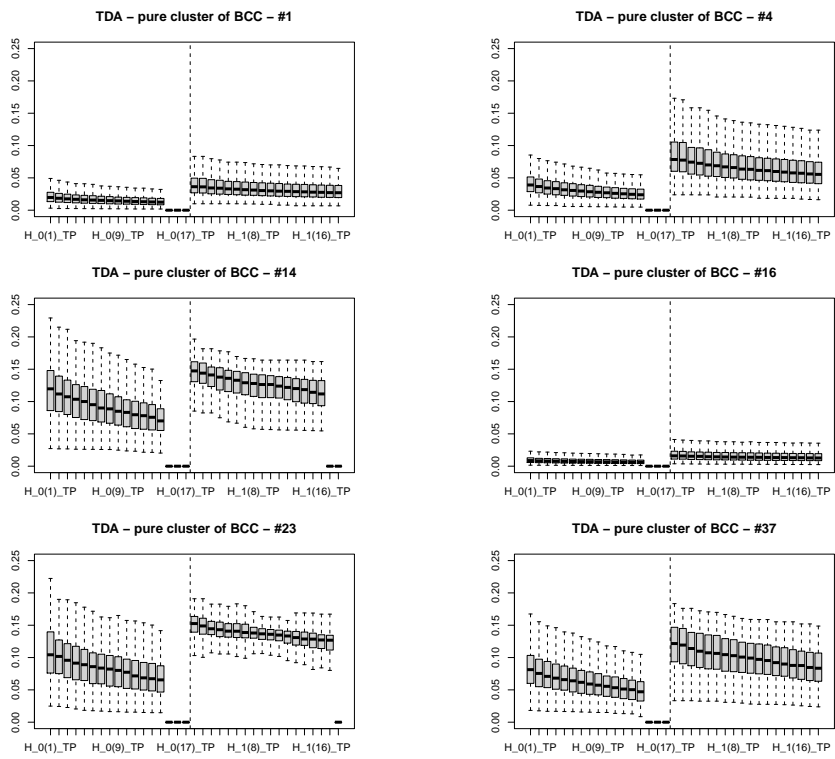


Figure 9: Boxplot for the TDA descriptor for the 6 pure clusters with BCCs.

156 **B.3 Clusters of BCC (not pure)**

157 In Fig. 10, we provide the CNA signature for observations in clusters that contain BCC structures  
 158 (but which are not pure). We remark that those signatures are close to the BCC ones ([444]: 6 and  
 159 [666]: 8), particularly for clusters 8 and 32. We discussed that in the main paper.

160 In Fig. 11, we provide the TDA signature for observations in the same clusters. We remark that the  
 161 discrimination is mainly over  $H_1$  coefficients (zero or nonzero).

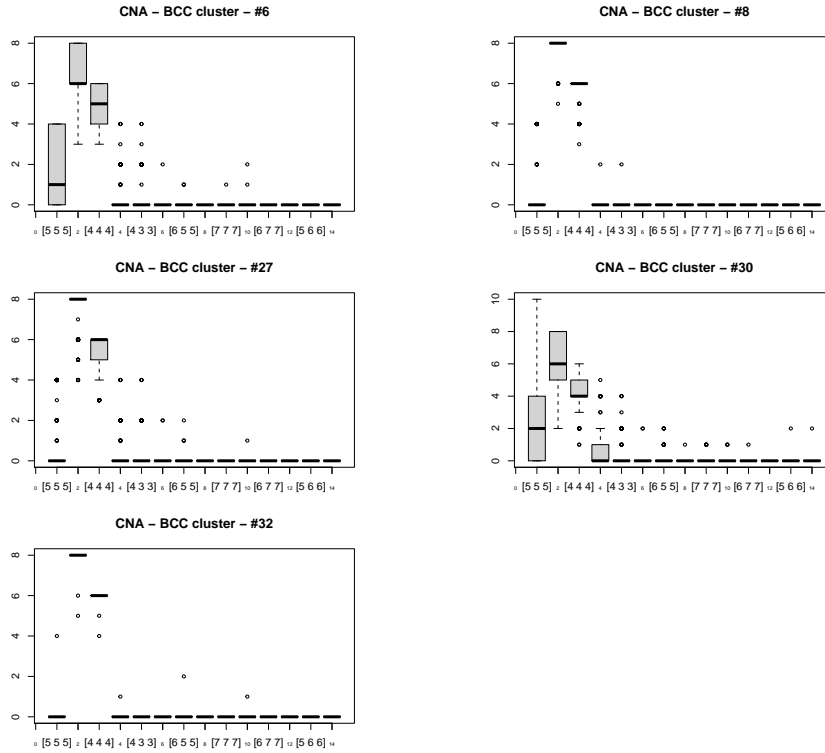


Figure 10: Boxplot for the CNA descriptor for the clusters with BCCs.

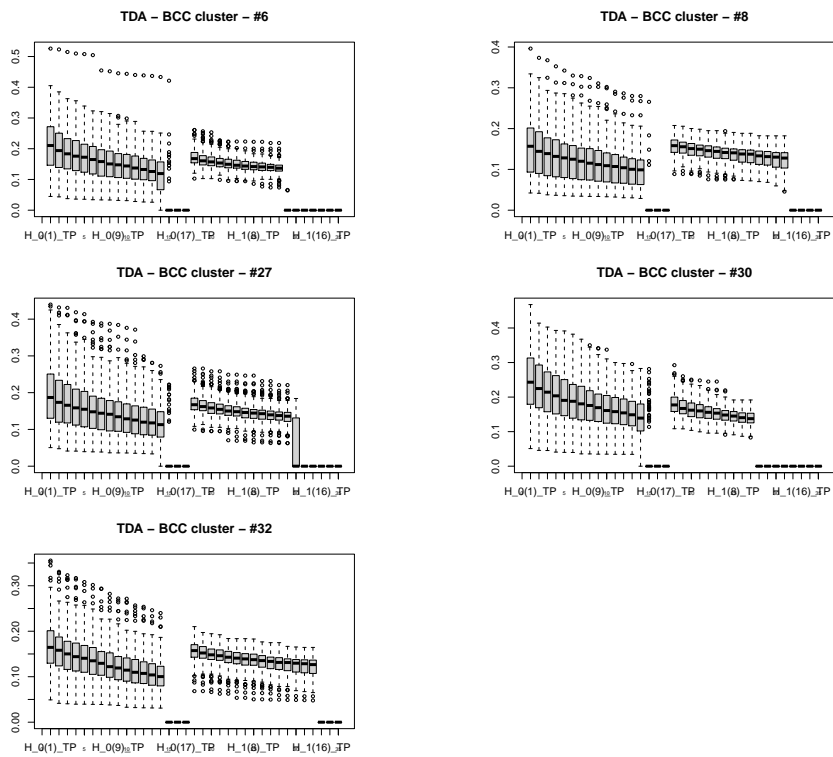


Figure 11: Boxplot for the TDA descriptor for the clusters with BCCs.