

# Evolutionary modes in protein observable space: the case of thioredoxins

*Running Title: Protein evolutionary modes*

**Sara Del Galdo**, *Dept. of Chemical Science and Technology, University of Roma Tor*

*Vergata, via della Ricerca Scientifica, 00133 Rome, Italy*

**Josephine Alba**, *Dept. of Chemistry, Sapienza University of Rome, P.le A. Moro, 5, 00185,*

*Rome, Italy*

**Andrea Amadei\***, *Dept. of Chemical Science and Technology, University of Roma Tor*

*Vergata, via della Ricerca Scientifica, 00133 Rome, Italy*

**Marco D'Abramo\***, *Dept. of Chemistry, Sapienza University of Rome, P.le A. Moro, 5,*

*00185, Rome, Italy*

*Contact Information: Marco D'Abramo, email: [marco.dabramo@uniroma1.it](mailto:marco.dabramo@uniroma1.it)*

*Andrea Amadei, email: [andrea.amadei@uniroma2.it](mailto:andrea.amadei@uniroma2.it)*

## **Abstract**

1  
2 In this article we investigated the structural and dynamical evolutionary behaviour of a set of  
3  
4 10 thioredoxin proteins as formed by three extant forms and seven resurrected ones in  
5  
6 laboratory. Starting from the crystallographic structures, we performed all-atoms molecular  
7  
8 dynamics simulations and compare the trajectories in terms of structural and dynamical  
9  
10 properties. Interestingly, the structural properties related to the protein density (i.e. the number  
11  
12 of residues divided by the excluded molecular volume) well describe the protein evolutionary  
13  
14 behaviour. Our results also suggest that the changes in sequence as occurred during the  
15  
16 evolution have affected the protein essential motions, allowing us to discriminate between  
17  
18 ancient and extant proteins in terms of their dynamical behaviour. Such results are yet more  
19  
20 evident when the bacterial, archeal and eukaryotic thioredoxins are separately analysed.  
21  
22  
23  
24  
25  
26  
27  
28

## **Keywords**

29 Protein evolution, protein dynamics, thioredoxin, essential motions, molecular dynamics  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1  
2  
3  
4  
5 **Introduction**  
6

7 The knowledge of the causes which have driven the evolution of proteins – in terms of sequence  
8 and structure from their first appearance on the earth – represents a very intriguing challenge.  
9  
10 In fact, the possibility to describe the factors which have ultimately determined the protein  
11 molecular evolution could pave the way to rationalize the adaption of different forms of life to  
12 external changes as well as to a better understanding of the proteins as molecular machines.  
13  
14 However, dealing with protein evolution requires information on extinct molecules which have  
15 disappeared along the history. Starting from the seminal observation that “most or all  
16 apparently heterologous gene derive ultimately from a common gene ancestor” (Pauling et al.  
17 1963),  
18 different strategies have been proposed to reconstruct, with a certain degree of  
19 accuracy, putative sequences of proteins that no longer exist. These Ancestral Sequence  
20 Reconstruction (ASR) methods (Fitch 1971; Chang and Donoghue 2000; Hall 2006; Benner et  
21 al. 2007; Liberles 2008; Arenas et al. 2017) combine multiple alignments of extant protein  
22 sequences, phylogenetic analysis and the probability of amino acid substitution to infer a  
23 putative ancient protein sequence.  
24  
25

26 Differently from “horizontal” approaches, where the process determining the protein structure  
27 and function evolution is inferred from sequence comparison of extant proteins, thus losing the  
28 evolutionary aspects, phylogenetic approaches give the possibility to estimate protein  
29 sequences along the history. That is, a phylogenetic tree obtained from the analysis of extant  
30 protein family members belonging to the three different domains of life might be combined  
31 with a multiple sequence alignment and a substitution model of evolution to provide a statistical  
32 inference on the ancestral sequence at any internal node of the tree (Arenas 2015).  
33 By such  
34 an approach, many different ancestral proteins were resurrected (Chang et al. 2002; Benner et  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 al. 2007; Ortlund et al. 2007; Harms and Thornton 2010; Perez-Jimenez et al. 2011; Hart et al.  
2 2014) and some of them crystallized (Ortlund et al. 2007; Ingles-Prieto et al. 2013) From  
3  
4 an evolutionary point of view, such data constitute an invaluable source of knowledge which  
5  
6 can be used to shed light on the key factors regulating the evolution of the protein along the  
7  
8 history. It is worth to mention here that ASR methods are affected by errors because of the  
9  
10 limits of the substitution model used. Therefore, new strategies are continuously proposed to  
11  
12 improve the ASR accuracy and those considering structural constraints represent one of the  
13  
14 more promising approaches (Arenas et al. 2017) However, the evaluation of the ASR  
15  
16 accuracy in the case of the resurrected proteins studied in this work is beyond the scope of this  
17  
18 paper.  
19  
20  
21  
22  
23

24 To the best of our knowledge, two recent papers on thioredoxins represent the most complete  
25  
26 study on resurrected proteins, where a set of seven Precambrian thioredoxin enzymes were  
27  
28 resurrected, crystallized and tested to measure their enzymatic activity (Perez-Jimenez et al.  
29  
30 2011; Ingles-Prieto et al. 2013) These enzymes – ubiquitous in all living organisms  
31  
32 (Holmgren et al. 1975) – were probably present in primitive life forms, as suggested by the  
33  
34 archetypical active site (CXXC) and the conserved fold. It is worth noting that this structural  
35  
36 dataset is actually the largest present in literature (Perez-Jimenez et al. 2011; Ingles-Prieto et  
37  
38 al. 2013) In those papers (Perez-Jimenez et al. 2011; Ingles-Prieto et al. 2013), the melting  
39  
40 temperature ( $T_m$ ) was suggested as a possible evolutionary observable, which left genetic  
41  
42 footprints on ancestral organisms (Boussau et al. 2008)  
43  
44  
45  
46  
47

48 Considering the well-established link between protein function, structure and dynamics and the  
49  
50 availability of this set of thioredoxin protein structures mentioned above, several questions  
51  
52 immediately arise: how do the changes in protein sequences due to the selective pressure  
53  
54 influence the protein motions? To what extent? Do such differences depend on the organism  
55  
56 life domain?  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 Inspired by such questions, we report here our investigation on the structural and dynamical  
2 behaviour (as provided by molecular dynamics simulations) of this set of thioredoxin proteins  
3 (Fig. 1) with the aim to provide a link between protein evolution and protein structure and  
4 dynamics.  
5  
6  
7  
8  
9

## 10 **Materials and Methods**

11 The structures of the 10 thioredoxin (Trx) proteins were taken from the Protein Data Bank:  
12 1ERU, extant eukaryota (Human); 2TRX, extant bacteria (Ecoli); 2E0Q, extant archaea  
13 (Archea); 4BA7, last bacterial common ancestor (LBCA); 2YJ7, last common ancestor of the  
14 cyanobacterial, deinococcus and thermus groups (LPBCA); 2YN1, last common ancestor of  $\gamma$ -  
15 proteobacteria (LGPCA); 3ZIV, archaea/eukaryota common ancestor (AECA); 2YNX, last  
16 archaeal common ancestor (LACA); 2YOI, last eukaryotic common ancestor (LECA); 2YPM,  
17 last common ancestor of fungi and animals (LAFCA). The reconstruction of the ancestral Trx  
18 enzymes as well the associated divergence times (Fig. 1) have been taken from the literature  
19 (Ingles-Prieto et al. 2013)□. Briefly, in that paper the sequences of ancient thioredoxins have  
20 been estimated by phylogenetic tree encompassing more than 200 diverse Trxs sequences from  
21 the three domains of life. The associated phylogenetic tree has been used to estimate the  
22 divergence dates to nodes in the tree using multiple fossil calibrations (Yang et al. 2006;  
23 Rannala et al. 2007) on the hypothesis that root of the tree lies between bacteria and the  
24 common ancestor of archaea and eukaryotes.  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47

48 In this work, we used the reconstruction of the ancestral Trx enzymes previously published  
49 (Perez-Jimenez et al. 2011; Ingles-Prieto et al. 2013)□ (i.e., the inferred sequences of the  
50 ancient thioredoxins, the corresponding X-ray structures and the estimated divergence dates)  
51 to describe the behavior of different structural and dynamical properties of this protein set along  
52 the history as explained in the following subsections.  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 Note that all the definitions of ancestors depend on the accuracy of the reconstructed tree and  
2 the ancestral state reconstruction algorithms used.  
3  
4

### 5 6 7 *Molecular Dynamics simulations*

8  
9 The Molecular Dynamics (MD) simulations were performed using the Gromacs software  
10 package (Hess et al. 2008) and the amber99sb-ildn force field. The SPC model (Berendsen et  
11 al. 1993) was used to mimic the water and sodium ions were added to neutralize the total  
12 thioresoxin charge. All the ten thioresoxins were simulated with periodic boundary conditions  
13 in the isothermal-isochoric ensemble (NVT), using an integration step of 2 fs and keeping the  
14 temperature constant (298 K) by the velocity rescaling thermostat (Bussi et al. 2007). The  
15 bonds were constrained using the LINCS algorithm (Hess et al. 1997) and for short range  
16 interactions a cut-off radius of 1.1 nm was employed. To compute long range interactions the  
17 particle mesh Ewald method (Darden et al. 1997) was used with grid search and cut-off radii  
18 of 1.1 nm. We calibrated the density of the boxes containing the water-protein solutions in order  
19 to obtain an identical pressure, within the noise (~ 10 bar), to the one provided by an MD  
20 simulation of a pure SPC box with a density corresponding to the liquid water experimental  
21 density at 298 K (we used as reference density 33.32 molecules per nm<sup>3</sup>) according to the  
22 procedure described in our recent work (Del Galdo et al. 2015). We performed for all the  
23 systems a productive MD simulation lasting 100 ns. Essential dynamics analysis (Amadei et  
24 al. 1993) was applied to each single trajectory and to a combination of them in order to  
25 highlight the phylogenetic-based differences in protein essential motions. The overlap (s)  
26 between the covariance matrices (A, B) is defined by (Hess 2002):  
27  
28

$$29 \quad s(A, B) = 1 - d(A, B) / (\text{tr} A + \text{tr} B)^{1/2}$$

30 where *tr* is the trace of the covariance matrix and *d*(A,B) is the difference between the  
31 covariance matrices A and B as defined by Hess (2002).  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52

## *Structural analysis*

The protein mean excluded volume, i.e. the mean volume enclosed by the solvent-accessible surface, has been estimated by averaging (over the productive MD simulation) the protein excluded volume along MD and using a probe radius of 0.14 nm according to the method reported in literature (Eisenhaber et al. 1995)□. The protein partial molecular volume was computed by the method reported in our previous work (Del Galdo et al. 2015)□, which is based on the evaluation of the mean protein excluded volume, the mean volume of the protein hydration shell and the hydration shell SPC density increment with respect to the reference SPC density (bulk density).

The number of hydrogen bonds and the residues with secondary structures have been calculated by Gromacs tools. The solvent (polar, apolar and total) accessible and excluded surface areas were calculated by finding solvent-exposed vertices of intersecting atoms (Fraczkiewicz et al. 1998)□. The number of proline residues as well as the B-factors have been directly extracted by the corresponding pdb files.

## *Principal component analysis of the structural observables*

For these structure-related properties a principal component analysis has been performed. The melting temperature was added to this set, because it has been suggested as a possible evolutionary observable, as indicated by its increase along the history (Perez-Jimenez et al. 2011)□. Due to the different magnitude and physical meaning of such observables, the covariance matrix was built using the adimensional rescaled shifts with respect to their averages for all the observables. That is, for the melting temperature and density-related properties

$$\Delta T'_m = (T_m - \langle T_m \rangle) / \sigma_{T_m}$$

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

$$\Delta\rho' = (\rho - \langle\rho\rangle) / \sigma_\rho$$

$$\Delta v' = (v - \langle v\rangle) / \sigma_v$$

with the angle brackets indicating the mean over time and  $\sigma$  the square root of the variance.

## **Results**

To investigate the evolutionary behaviour of the 10 thioredoxin proteins, we analysed the corresponding molecular dynamics trajectories using both fluctuation-related properties (protein principal motions and entropy estimates) and structure-related observables (protein density, molecular volume, hydrogen bond contents, amount of secondary structures and solvent accessible surface area).

It is worth noting that although our set is composed by proteins belonging to the same class, these proteins span a quite large sequence identity interval (between 0.25 and 0.92, see Table 1).

### **Structural behaviour**

To address the evolutionary behaviour of thioredoxin observables and possibly uncover their correlation, we considered 12 different properties: the experimental melting temperature, the solvent excluded surface area, the polar, apolar and total solvent accessible surface area, the number of proline residues, the B-factors, the number of hydrogen bonds, the fraction of residues having secondary structure, the residue density within the protein excluded volume ( $\rho$ ) and the partial molecular volume ( $v$ ).

We excluded the relative density increment of the solvent density within the protein hydration shell with respect to the solvent bulk density, because it is nearly constant all along the evolutionary time (Fig. S1 in S.I.). This set of 11 observables was analyzed by means of



1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

principal component analysis, thus providing a set of eigenvectors and associated eigenvalues describing the evolutionary behaviour within such an 11-dimensional space.

The spectrum of eigenvalues (Fig. S2 in S.I.) obtained from the diagonalisation of the covariance matrix shows that the first eigenvector accounts for ~ 50 % of the total fluctuations, the second for ~ 30 % and the third for ~ 10 %. The first eigenvector has the (nearly identical) major component values for the  $T_m$ , the residue density ( $\rho$ ) and the partial molecular volume ( $v$ ). Due to the limited sample size formed by 10 proteins, which is further reduced in the case of separate analysis for the archeal (3 structures), eukaryotic (4 structures), and bacterial (4 structures) subsets, we decided to restrict our analysis on these three observables. Such a choice is largely justified by the principal component analysis performed on the global observable space defined by the 12 observables previously described, which shows that the major component of the first eigenvector (explaining about the 50% of the total fluctuations) is dominated by the  $T_m$ , the residue density ( $\rho$ ) and the partial molecular volume ( $v$ ).

Therefore, these three observables not only represent the best choice for the description of the evolutionary behaviour, but also include the melting temperature and the density related properties which are simple, physically-sound observables. We recently observed that the protein density well correlates with the protein optimal growth temperature (Amadei et al. 2017)□. We would like to stress here that there are not direct thermodynamic relations connecting the difference in the heat capacity and/or the melting temperature – experimentally found to be correlated with thioredoxin evolution – with the protein density and/or the partial molecular volume.

In figure 2, we show the evolutionary trend of these three observables for archeal, eukaryotic and bacterial thioredoxins.

1 From this figure, it is evident that the  $T_m$ , the residue density ( $\rho$ ) and the partial molecular  
2 volume ( $v$ ) significantly change along the evolutionary time (the sole exception being the  $T_m$   
3 of the archeal thioredoxins, which decreases of only 3 K along the history).  
4

5  
6  
7 In figure 3 we report the evolutionary trajectory for the sample of proteins considered within  
8 such a re-scaled 3-D observable space. Very interestingly, all the points are not spread over the  
9 planes, but they are rather well aligned along the diagonal, indicating a significant correlation  
10 among these properties.  
11

12  
13  
14 In order to identify a possible single *generalised* observable able to describe the evolutionary  
15 trend, we performed the diagonalization of the  $\Delta T'_m, \Delta \rho', \Delta v'$  covariance matrix, providing by  
16 means of its eigenvectors, the relevant modes and corresponding observables within such a  
17 space (Amadei et al. 1993).  
18

19  
20  
21 By using such eigenvectors ( $v_1, v_2, v_3$ ) as new basis set of the observable space, we can express  
22 the evolutionary trajectory in terms of the three corresponding *generalised* observables ( $g_1, g_2,$   
23  $g_3$ , each corresponding to a specific linear combination of the original properties and defined  
24 by the projection of the observable trajectory on the eigenvectors). For eukaryotes, archea and  
25 bacteria the largest covariance matrix eigenvalue (corresponding to the  $v_1$  eigenvector)  
26 provides 90%, 99% and 94% of the total square fluctuations, respectively, clearly indicating  
27 that the essential information on the evolutionary trend of thioredoxins can be obtained by the  
28  $g_1$  evolutionary behaviour (i.e. filtering out the small fluctuations along  $v_2$  and  $v_3$ ). It is worth  
29 to note that the  $v_1$  eigenvectors have nearly identical component absolute values ( $\sim 0.58$ ) with  
30 the  $\Delta \rho'$  component sign opposite to the others, thus indicating that the corresponding  
31 *generalized* observable is given by a virtually homogenous mixing of the original properties  
32 ( $\Delta T'_m, \Delta \rho', \Delta v'$ ), which are characterized by anti-correlation between  $\Delta \rho'$  vs  $\Delta T'_m, \Delta v'$ .  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56

57 As shown in figure 3 (lower right panel) the evolutionary trajectories along  $v_1$ , show a  
58 remarkable variation around 2 Gyears ago, well matching the archaean-proterozoic era  
59  
60  
61  
62  
63  
64  
65

1 transition believed to correspond to a significant decrease of the global earth temperature  
2 (Lowe et al. 2004).  
3

4 Therefore, our results suggest that global temperature decrease not only induced the  $T_m$   
5 decrease as reported in literature (Perez-Jimenez et al. 2011)□, but possibly provided a  
6 significant residue density increase (partial molecular volume decrease), at least for  
7 thioredoxins. Interestingly, the partial molecular volume decrease is essentially due to the  
8 protein volume decrease as the hydration shell water density is virtually constant for all the  
9 investigated thioredoxins (see Fig. S1 for the eukaryotic thioredoxins). Finally, it is worth to  
10 mention that our results pinpoint the evolutionary trend of these three observables as occurring  
11 in thioredoxins; in different proteins, as bacterial ribonuclease H1 (RNH) proteins where the  
12  $T_m$  is poorly correlated with evolution (Hart et al. 2014), it has been found that the low heat  
13 capacity of unfolding, due to the presence of residual structure in the unfolded state, is the  
14 major determinant for the RNHs difference in thermostability. Unfortunately, in that paper  
15 (Hart et al. 2014) □ only one structure of a resurrected RNH was made available, preventing  
16 the possibility to estimate the behavior of the observables used in the Trxs data set along the  
17 history.  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40

## 41 **Dynamical behavior**

### 42 *Principal motions*

43 The essential motions describing the overall thioredoxin dynamics were calculated by principal  
44 component analysis performed on a single trajectory obtained by concatenating the ten  
45 thioredoxin trajectories of the C-alpha atoms.  
46  
47  
48  
49  
50

51 The first eigenvector describes the 60% of the total variance, indicating that a remarkable  
52 amount of the protein motion of the thioredoxins is concentrated along this direction. The  
53 components of the first eigenvector (Fig. S3) show that the first essential motion is mainly  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 described by the N-terminal region, the  $\alpha$ -helix 1, 2 and 3 and the  $\beta$ -sheet 1, 4 and 5, (see Fig.  
2 1 for the secondary structure assignment).  
3  
4  
5  
6

7 The analysis of the sampled structures on the essential subspace as characterized by their  
8 projection on the first two eigenvectors shows three main regions, each corresponding to a  
9 different branch (Fig. S4). That is, the region of the subspace explored is characteristic of the  
10 life-domain. It is worth noting that the first eigenvector discriminates between archeal/bacterial  
11 vs eukaryotic thioredoxins, whereas the second eigenvector is able to discriminate between  
12 archeal and bacterial proteins, too.  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23

24 To shed light on this different behaviour among eukaryotic, archeal and bacterial branches, the  
25 essential motions as described by the single thioredoxin trajectories were compared by pairs.  
26  
27

28 The degree of similarity between the essential motions as well as the sequence identities are  
29 reported in the Table 1 for all the couples.  
30  
31  
32  
33

34 These values clearly show that the essential motions within the bacterial branch are well  
35 conserved, being the overlap between 0.65 (*E. coli* vs LGPCA) and 0.48 (*E. coli* and LPBCA).  
36  
37

38 The comparison between the overlap of the essential motions and the sequence identities also  
39 points out that consecutive (along the history) thioredoxins share high sequence identities and  
40 high motion overlaps only within the same branch. In fact, LACA and LECA although sharing  
41 a high sequence identity (0.56) shows a relatively low motion overlaps (0.209), thus indicating  
42 that evolutionary steps have affected the protein essential motions.  
43  
44  
45  
46  
47  
48  
49  
50

51 Interestingly, also the AECA and LACA show a quite large overlap of the essential motions  
52 when compared to the bacterial thioredoxins, whereas the three remaining thioredoxins  
53 belonging to the eukaryotic branch (LECA, LAFCA and Human) show a very limited overlap  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

with respect to the other thioredoxins, i.e. the maximum overlap is 0.49 between LAFCA and Human which share the 85% of sequence identity.

All these results indicate that the archeal thioredoxins are rather similar in terms of sequence identity (Fig. S5) and protein motions to their ancestor as well as to the bacterial thioredoxins. The eukaryotic branch thioredoxins, on the other hand, are more distant in both sequence identity and essential motions with respect to the other thioredoxins.

### *Entropy*

To obtain further insights in the thioredoxin protein behaviour, the estimation of their configurational fluctuation entropies has been performed using the covariance matrix of atomic coordinates as indicated by Schlitter (1993)□. In all the branches, the entropy shows a remarkable decrease along the history (Fig. 4), with a pronounced variation around 2 Gyars ago corresponding to the archaean-proterozoic era transition, believed to correspond to a significant decrease of the global earth temperature (Lowe et al. 2004)□. Interestingly, such a steep variation is also observed – in the same time interval – for the  $T_m$ , the partial molecular volume, the residue density and the associated *generalised* observable derived in the previous section.

From these results, it is tempting to suggest the hypothesis that the evolutionary behaviour of thioredoxins could be essentially entropically driven, with unfolded states of ancient and extant proteins similar in free energy and the ancient folded states (characterized by high melting temperatures) entropically stabilized by the residue density decrease.

1 After the Archean period, proteins were allowed to enhance local interactions – for example to  
2 optimize their biological functions – with a consequent entropy loss and melting temperature  
3 decrease, both permitted by the lower earth temperature.  
4  
5  
6  
7  
8

## 9 **Conclusions**

10 In this work we calculated several global protein observables to characterize the  
11 structural/molecular evolution of thioredoxins.  
12  
13  
14

15 In addition to the melting temperature, which was already found to change along the evolution,  
16 we calculated several other properties routinely used in large-scale classification of proteins as  
17 well as a new set of properties related to the protein density.  
18  
19  
20  
21  
22

23 Our data point out that the melting temperature, the protein density and the partial molecular  
24 volume are the major components of the main evolutionary motion in thioredoxins. Therefore,  
25 by means of principal component analysis, we studied the evolutionary behavior within such a  
26 3-dimensional space, which highlights an interesting trend along the history. In fact, the  
27 supposed global earth temperature decrease, dated in the archaean-proterozoic era transition,  
28 well matches the remarkable variation of the evolutionary trajectory along the main essential  
29 evolutionary mode.  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40

41 The analysis of the molecular dynamics trajectories points out that the dynamical behavior of  
42 the thioredoxins is – as expected – driven by the structural changes induced by the protein  
43 sequence variations (upon evolution). However, the evolutionary step corresponding to the  
44 archaean-proterozoic era transition determines remarkable changes of the protein entropy.  
45  
46  
47  
48  
49  
50

51 Our study shows that structural-molecular evolution of thioredoxin proteins can be well  
52 described by a set of generalized protein observables, and that, among several properties, those  
53 related to the protein density are some of the most representative. The possibility to extend our  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 study to different protein families – using experimental and/or homology modeling approaches  
2 – could provide additional information on the relationship between protein structure and  
3 protein evolution, thus representing a very interesting research theme to investigate in the next  
4 future.  
5  
6  
7  
8  
9

## 10 **Acknowledgments**

11 This work was partially funded by Sapienza, University of Rome (Progetto di Ateneo 2017).  
12  
13 The authors gratefully acknowledge NVIDIA and CINECA for computational support.  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23

## 24 **References**

- 25  
26 Amadei A, Del Galdo S, D’Abramo M (2017) Density discriminates between thermophilic and  
27 mesophilic proteins. *J Biomol Struct Dyn* 1–9. doi:10.1080/07391102.2017.1385537  
28  
29  
30  
31  
32 Amadei A, Linssen AB, Berendsen HJ (1993) Essential dynamics of proteins. *Proteins* 17:412–  
33 25. doi: 10.1002/prot.340170408  
34  
35  
36  
37 Arenas M (2015) Trends in substitution models of molecular evolution. *Front Genet.* doi:  
38 10.3389/fgene.2015.00319  
39  
40  
41  
42 Arenas M, Weber CC, Liberles DA, Bastolla U (2017) ProtASR: An Evolutionary Framework  
43 for Ancestral Protein Reconstruction with Selection on Folding Stability. *Syst Biol.* doi:  
44 10.1093/sysbio/syw121  
45  
46  
47  
48  
49  
50 Benner S, Sassi S, Gaucher E (2007) Molecular paleoscience: systems biology from the past.  
51 In: Toone E by EJ (ed) *Adv Enzymol Relat Areas of Molecular Biology, Volume 75: Protein*  
52 *Evolution.* pp 1–132  
53  
54  
55  
56  
57 Berendsen HJC, Grigera JR, Straatsma TP (1993) The missing term in effective pair potentials.  
58 *J Chem Phys* 98:10089–10092  
59  
60  
61  
62  
63  
64  
65

1  
2 Boussau B, Blanquart S, Necsulea A, et al (2008) Parallel adaptations to high temperatures in  
3 the Archaeon eon. *Nature*. doi: 10.1038/nature07393  
4

5  
6  
7 Bussi G, Donadio D, Parrinello M (2007) Canonical sampling through velocity rescaling. *J*  
8  
9 *Chem Phys* 126:014101. doi: 10.1063/1.2408420  
10

11  
12 Chang BS, Jonsson K, Kazmi MA, et al (2002) Recreating a functional ancestral archosaur  
13 visual pigment. *Mol Biol Evol* 19:1483–1489. doi: 10.1093/oxfordjournals.molbev.a004211  
14  
15

16  
17  
18 Chang BSW, Donoghue MJ (2000) Recreating ancestral proteins. *Trends Ecol. Evol.* 15:109–  
19  
20 114  
21

22  
23 Darden T, Tork D, Pedersen L (1997) Particle mesh Ewald: An N-log(N) method for Ewald  
24 sums in large systems. *J Comput Chem* 18:1463–1472  
25  
26

27  
28  
29 Del Galdo S, Marracino P, D’Abramo M, Amadei A (2015) In silico characterization of protein  
30 partial molecular volumes and hydration shells. *Phys Chem Chem Phys* 17:31270–31277. doi:  
31  
32 10.1039/C5CP05891K  
33  
34

35  
36 Eisenhaber F, Lijnzaad P, Argos P, et al (1995) The double cubic lattice method: Efficient  
37 approaches to numerical integration of surface area and volume and to dot surface contouring  
38 of molecular assemblies. *J Comput Chem* 16:273–284. doi: 10.1002/jcc.540160303  
39  
40  
41

42  
43 Fitch WM (1971) Toward defining the course of evolution: Minimum change for a specific tree  
44 topology. *Syst Biol* 20:406–416. doi: 10.1093/sysbio/20.4.406  
45  
46

47  
48  
49 Fraczkiewicz R, Braun W (1998) Exact and efficient analytical calculation of the accessible  
50 surface areas and their gradients for macromolecules. *J Comput Chem* 19:319–333. doi:  
51  
52 10.1002/(SICI)1096-987X(199802)19:3<319::AID-JCC6>3.0.CO;2-W  
53  
54

55  
56 Hall BG (2006) Simple and accurate estimation of ancestral protein sequences. *Proc Natl Acad*  
57  
58 *Sci U S A* 103:5431–6. doi: 10.1073/pnas.0508991103  
59  
60



1 Harms MJ, Thornton JW (2010) Analyzing protein structure and function using ancestral gene  
2 reconstruction. *Curr Opin Struct Biol* 20:360–6. doi: 10.1016/j.sbi.2010.03.005  
3  
4

5 Hart KM, Harms MJ, Schmidt BH, et al (2014) Thermodynamic system drift in protein  
6 evolution. *PLoS Biol* 12:e1001994. doi: 10.1371/journal.pbio.1001994  
7  
8  
9

10 Hess B (2002) Convergence of sampling in protein simulations. *Phys Rev E* 65:031910. doi:  
11 10.1103/PhysRevE.65.031910  
12  
13  
14

15 Hess B, Bekker H, Berendsen HJC, Fraaije J (1997) LINCS: A Linear Constant Solver for  
16 Molecular Simulations. *J Comput Chem* 18:1463–1472  
17  
18  
19  
20

21 Hess B, Kutzner C, van der Spoel D, Lindahl E (2008) GROMACS 4: Algorithms for Highly  
22 Efficient, Load-Balanced, and Scalable Molecular Simulation. *J Chem Theory Comput* 4:435–  
23 447. doi: 10.1021/ct700301q  
24  
25  
26  
27

28 Holmgren A, Soderberg BO, Eklund H, Branden CI (1975) Three-dimensional structure of  
29 Escherichia coli thioredoxin-S2 to 2.8 Å resolution. *Proc Natl Acad Sci* 72:2305–2309. doi:  
30 10.1073/pnas.72.6.2305  
31  
32  
33  
34

35 Ingles-Prieto A, Ibarra-Molero B, Delgado-Delgado A, et al (2013) Conservation of protein  
36 structure over four billion years. *Structure* 21:1690–1697. doi: 10.1016/j.str.2013.06.020  
37  
38  
39  
40

41 Liberles DA (2008) *Ancestral Sequence Reconstruction*. Oxford University Press  
42  
43  
44

45 Lowe DR, Tice MM (2004) Geologic evidence for Archean atmospheric and climatic  
46 evolution: Fluctuating levels of CO<sub>2</sub>, CH<sub>4</sub>, and O<sub>2</sub> with an overriding tectonic control.  
47 *Geology* 32:493–496. doi: 10.1130/G20342.1  
48  
49  
50  
51

52 Ortlund EA, Bridgham JT, Redinbo MR, Thornton JW (2007) Crystal Structure of an Ancient  
53 Protein: Evolution by Conformational Epistasis. *Science* (80- ) 317:1544–1548. doi:  
54 10.1126/science.1142819  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 Pauling L, Zuckerkandl E (1963) Chemical paleogenetics: Molecular 'restoration studies' of  
2 extinct forms of life. *Acta chem. scand* 17:S9–S16  
3  
4

5 Perez-Jimenez R, Inglés-Prieto A, Zhao Z-M, et al (2011) Single-molecule paleoenzymology  
6 probes the chemistry of resurrected enzymes. *Nat Struct Mol Biol* 18:592–596. doi:  
7 10.1038/nsmb.2020  
8  
9

10  
11  
12 Rannala B, Yang Z (2007) Inferring speciation times under an episodic molecular clock. *Syst*  
13 *Biol* 56:453–466. doi: 10.1080/10635150701420643  
14  
15

16  
17  
18 Schlitter J (1993) Estimation of absolute and relative entropies of macromolecules using the  
19 covariance matrix. *Chem Phys Lett* 215:617–621. doi: 10.1016/0009-2614(93)89366-P  
20  
21

22  
23 Yang Z, Rannala B (2006) Bayesian estimation of species divergence times under a molecular  
24 clock using multiple fossil calibrations with soft bounds. *Mol Biol Evol* 23:212–226. doi:  
25 10.1093/molbev/msj024  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

**Table 1.** The Matrix of the sequence identity (elements below the diagonal) and overlap between essential motions as obtained from the overlap of the covariance matrices (elements above the diagonal)

Protein		LBCA	LPBCA	LGPCA	Ecoli	AECA	LACA	LECA	LAFCA	Human	Archea
	pdb id.	<b>4ba7</b>	<b>2yj7</b>	<b>2yn1</b>	<b>2trx</b>	<b>3ziv</b>	<b>2ynx</b>	<b>2yoi</b>	<b>2ypm</b>	<b>1eru</b>	<b>2e0q</b>
LBCA	<b>4ba7</b>	1	0.653	0.618	0.546	0.682	0.406	0.238	0.264	0.169	0.401
LPBCA	<b>2yj7</b>	0.87	1	0.544	0.476	0.614	0.463	0.216	0.224	0.157	0.390
LGPCA	<b>2yn1</b>	0.67	0.68	1	0.655	0.476	0.283	0.206	0.223	0.159	0.299
Ecoli	<b>2trx</b>	0.56	0.57	0.83	1	0.449	0.261	0.201	0.258	0.154	0.243
AECA	<b>3ziv</b>	0.85	0.76	0.60	0.53	1	0.469	0.272	0.277	0.181	0.426
LACA	<b>2ynx</b>	0.77	0.72	0.56	0.51	0.92	1	0.209	0.208	0.148	0.228
LECA	<b>2yoi</b>	0.57	0.53	0.44	0.37	0.59	0.56	1	0.490	0.390	0.201
LAFCA	<b>2ypm</b>	0.51	0.48	0.42	0.36	0.54	0.51	0.85	1	0.323	0.187
Human	<b>1eru</b>	0.36	0.34	0.30	0.25	0.38	0.35	0.55	0.58	1	0.154
Archea	<b>2e0q</b>	0.51	0.46	0.41	0.40	0.61	0.53	0.48	0.46	0.38	1

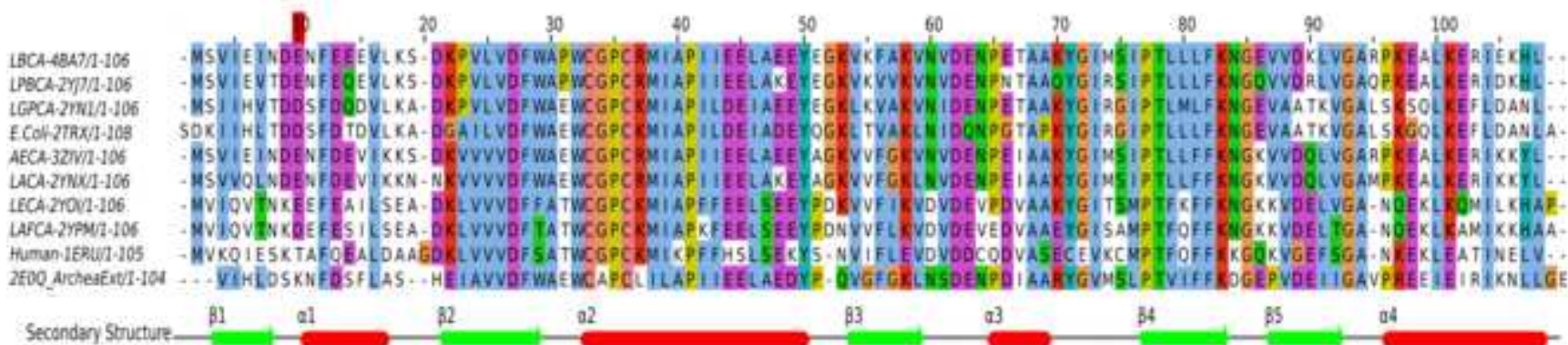
## Figure Captions

1  
2  
3  
4 **Fig. 1** Sequence alignment for the ten thioredoxins studied in this work and the associated  
5 secondary structure (upper panel). Thioredoxin protein names and associated geological time  
6 (lower left table). Lower right panel: ribbon representation of the bacterial extant thioredoxin  
7 (PDB id. 2trx)  
8  
9

10  
11  
12  
13 **Fig. 2** Time evolution of the melting temperature ( $T_m$  upper left panel), the partial molecular  
14 volume (lower left panel) and the residue density within the protein excluded volume ( $\rho$ , upper  
15 right panel). Black, red and green circles refer to proteins associated to the eukaryotic, bacterial  
16 and archeal branch, respectively. The dashed lines serve as a guide for the eye. Note that the  
17 melting temperatures of the two oldest resurrected thioredoxins belonging to the archeal-  
18 eukaryotic and bacterial branches coincide  
19  
20  
21  
22  
23

24  
25  
26 **Fig. 3** Plot of the adimensional rescaled shifts ( $\Delta T'_m, \Delta \rho', \Delta v'$ ) for the ten thioredoxin proteins  
27 (upper plots and lower left plot) and the projection of the protein observables ( $\Delta T'_m, \Delta \rho', \Delta v'$ )  
28 on the first eigenvector along the time (lower right plot; the dashed lines serve as a guide for  
29 the eye). Black, green and red circles refer to eukaryotic, archeal and bacteria thioredoxins,  
30 respectively  
31  
32  
33  
34  
35

36  
37 **Fig. 4** Changes in the thioredoxin entropy with respect to the corresponding common ancestors  
38 (AECA and LBCA). Black, green and red points indicate eukaryotic, archeal and bacterial  
39 thioredoxins, respectively. The dashed lines serve as a guide for the eye  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65



Name	Trx corresponding to	PDB id.	Time before present (billion years)
LBCA	last bacterial common ancestor	4ba7	-4.2
LPBCA	last common ancestor of the cyanobacterial, deinococcus and thermus groups	2yj7	-2.4
LGPCA	last common ancestor of g-proteobacteria	2yn1	-1.6
Ecoli	Ecoli	2trx	0
AECA	archaea/eukaryota common ancestor	3ziv	-4.2
LACA	last archaeal common ancestor	2ynx	-4.1
LECA	last eukaryotic common ancestor	2yoi	-1.6
LAFCA	last common ancestor of fungi and animals	2ypm	-1.4
Human	Human	1eru	0
Archea	Archea	2e0q	0

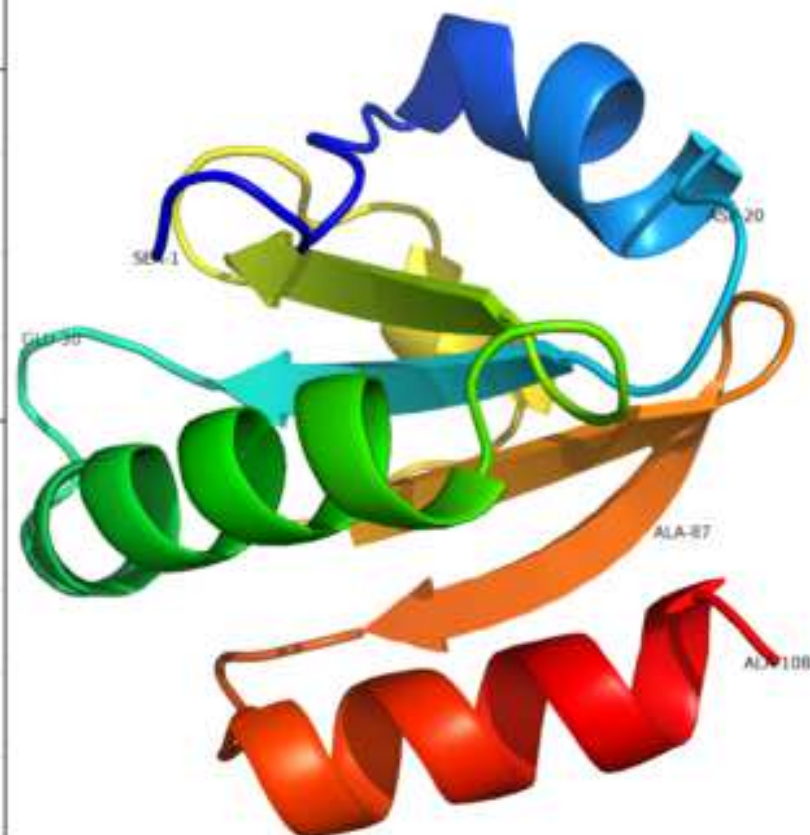


Figure 2

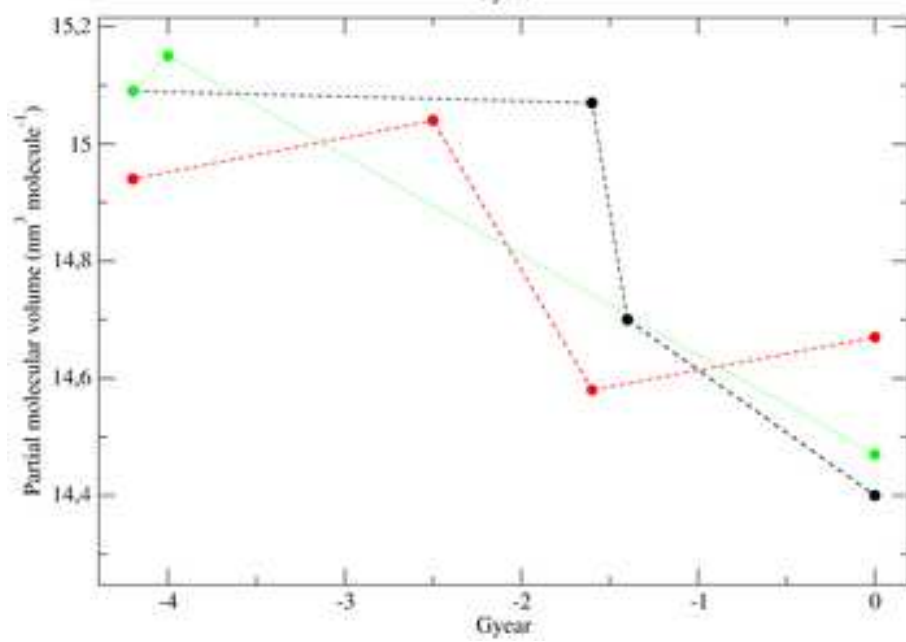
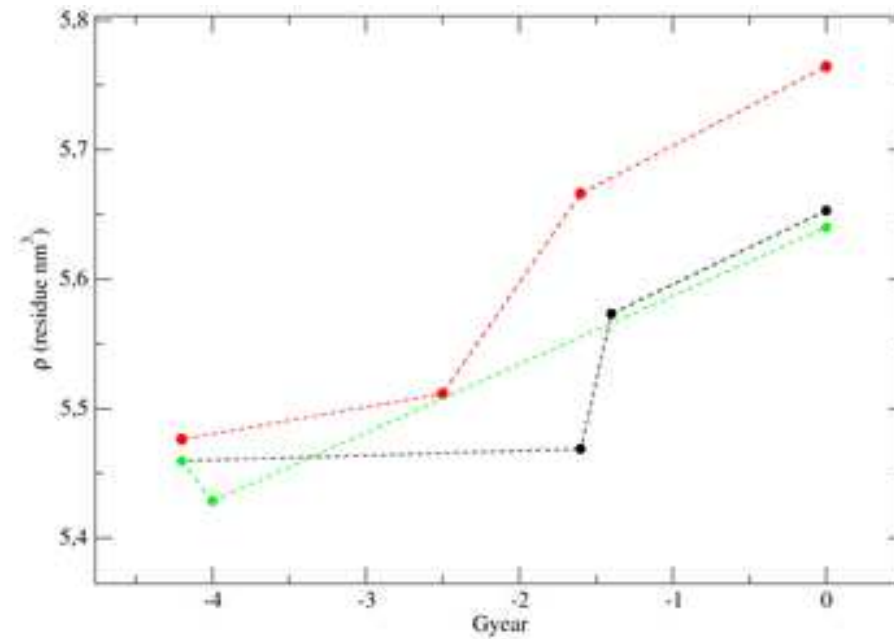
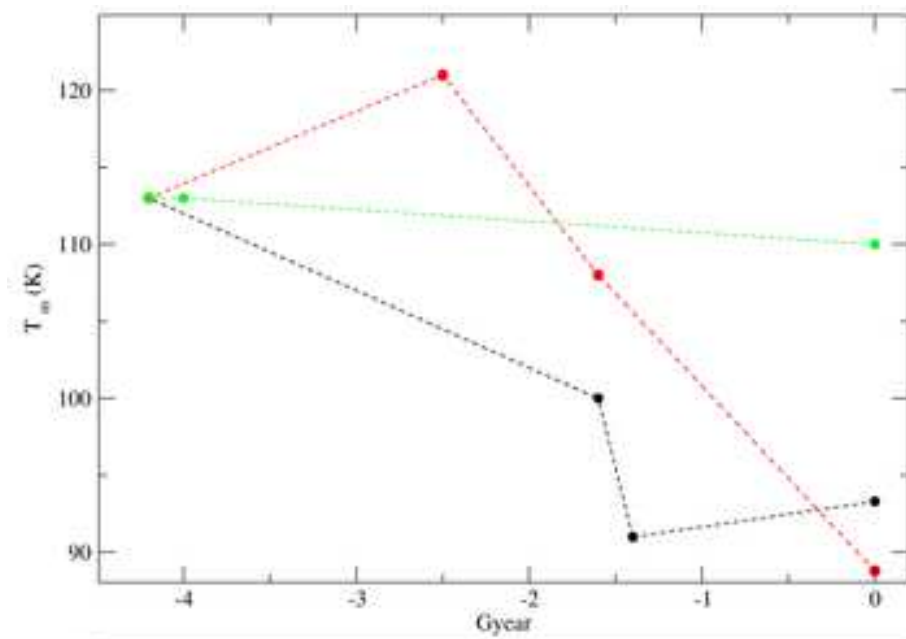
[Click here to access/download;Figure;Fig2.tif](#)

Figure 3

[Click here to access/download;Figure;Fig3.tif](#)

