

# Combining Ancestral Reconstruction with Folding-Landscape Simulations to Engineer Heterologous Protein Expression

Gloria Gamiz-Arco<sup>1†</sup>, Valeria A. Risso<sup>1†</sup>, Eric A. Gaucher<sup>2</sup>, Jose A. Gavira<sup>3</sup>, Athi N. Naganathan<sup>4\*</sup>, Beatriz Ibarra-Molero<sup>1\*</sup> and Jose M. Sanchez-Ruiz<sup>1\*</sup>

**1 - Departamento de Química Física, Facultad de Ciencias, Unidad de Excelencia de Química Aplicada a Biomedicina y Medioambiente (UEQ), Universidad de Granada, 18071 Granada, Spain**

**2 - Department of Biology, Georgia State University, Atlanta, GA 30303, USA**

**3 - Laboratorio de Estudios Cristalográficos, Instituto Andaluz de Ciencias de la Tierra, CSIC, Unidad de Excelencia de Química Aplicada a Biomedicina y Medioambiente (UEQ), Universidad de Granada, Avenida de las Palmeras 4, Armilla, Granada 18100, Spain**

**4 - Department of Biotechnology, Bhupat and Jyoti Mehta School of Biosciences, Indian Institute of Technology Madras, Chennai 600036, India**

**Correspondence to Athi N. Naganathan, Beatriz Ibarra-Molero and Jose M. Sanchez-Ruiz:** [athi@iitm.ac.in](mailto:athi@iitm.ac.in) (A.N. Naganathan), [beatriz@ugr.es](mailto:beatriz@ugr.es) (B. Ibarra-Molero), [sanchezr@ugr.es](mailto:sanchezr@ugr.es) (J.M. Sanchez-Ruiz)

[@jmsanchezruiz](https://doi.org/10.1016/j.jmb.2021.167321) (J.M. Sanchez-Ruiz), [@AthiNaganathan](https://twitter.com/AthiNaganathan) (A.N. Naganathan), [@Gavirius](https://twitter.com/Gavirius) (J.A. Gavira)

<https://doi.org/10.1016/j.jmb.2021.167321>

**Edited by Daniel Otzen**

## Abstract

Obligate symbionts typically exhibit high evolutionary rates. Consequently, their proteins may differ considerably from their modern and ancestral homologs in terms of both sequence and properties, thus providing excellent models to study protein evolution. Also, obligate symbionts are challenging to culture in the lab and proteins from uncultured organisms must be produced in heterologous hosts using recombinant DNA technology. Obligate symbionts thus replicate a fundamental scenario of metagenomics studies aimed at the functional characterization and biotechnological exploitation of proteins from the bacteria in soil. Here, we use the thioredoxin from *Candidatus Photodesmus katoptron*, an uncultured symbiont of flashlight fish, to explore evolutionary and engineering aspects of protein folding in heterologous hosts. The symbiont protein is a standard thioredoxin in terms of 3D-structure, stability and redox activity. However, its folding outside the original host is severely impaired, as shown by a very slow refolding *in vitro* and an inefficient expression in *E. coli* that leads mostly to insoluble protein. By contrast, resurrected Precambrian thioredoxins express efficiently in *E. coli*, plausibly reflecting an ancient adaptation to unassisted folding. We have used a statistical-mechanical model of the folding landscape to guide back-to-ancestor engineering of the symbiont protein. Remarkably, we find that the efficiency of heterologous expression correlates with the *in vitro* (*i.e.*, unassisted) folding rate and that the ancestral expression efficiency can be achieved with only 1–2 back-to-ancestor replacements. These results demonstrate a minimal-perturbation, sequence-engineering approach to rescue inefficient heterologous expression which may potentially be useful in metagenomics efforts targeting recent adaptations.

© 2021 The Author(s). Published by Elsevier Ltd. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

## Introduction

Many proteins of industrial and therapeutic interest are produced in heterologous hosts using recombinant DNA technology.<sup>1–3</sup> Moreover, heterologous expression is unavoidable in metagenomics studies aimed at the functional characterization and biotechnological exploitation of proteins from organisms that cannot be cultured in the lab,<sup>4,5</sup> i.e., the majority of the bacteria in soil.<sup>6</sup> Unfortunately, over-expression of a foreign gene poses a serious challenge to an organism and may lead to non-functional species, such as misfolded protein, proteolyzed protein or, more typically, insoluble protein aggregates.<sup>7</sup> These problems may be alleviated by using engineered hosts that have been modified for instance to minimize protease activity or to over-express molecular chaperones that assist correct folding.<sup>2,8</sup> Despite advances in host engineering, heterologous expression of functional proteins remains a major biotechnological bottleneck. For instance, about half of the proteins targeted in structural genomics initiatives could not be purified.<sup>9</sup>

Ancestral sequence reconstruction (ASR) uses phylogenetic analyses and sequences of modern protein homologs to compute statistically plausible approximations to the corresponding ancestral sequences.<sup>10,11</sup> During the last ~25 years, proteins encoded by reconstructed sequences (“resurrected” ancestral proteins) have been widely used as tools to address important problems in molecular evolution.<sup>12–16</sup> They have also been found to provide new possibilities for protein biomedical applications and protein engineering.<sup>17–21</sup> Ancestral proteins may considerably differ from their modern counterparts in terms of sequence and their experimental preparation necessarily involves heterologous expression, as the ancient original hosts are not available. Therefore, the fact that many resurrected ancestral proteins have been purified and studied experimentally emerges as remarkable in itself, even after acknowledging the obvious publication bias in favour of positive experimental outcomes.

Moreover, a substantial number of studies have actually reported improved heterologous expression of ancestral proteins as compared with their modern counterparts. Examples include: phosphate-binding protein,<sup>22</sup> periplasmic binding protein,<sup>23</sup> serum paraoxonase,<sup>24</sup> coagulation factor VIII,<sup>25</sup> titin,<sup>26</sup> haloalkane dehalogenases,<sup>27</sup> cytidine and adenine base editors,<sup>28</sup> diterpene cyclase,<sup>29</sup> rubisco,<sup>30</sup> endoglucanases,<sup>31</sup> L-amino acid oxidases,<sup>32</sup> laccases,<sup>33</sup> front-end  $\Delta 6$ -desaturases<sup>34</sup> and fatty acid photo-decarboxylases.<sup>35</sup> On a related note, recent studies on ancestral proteins have noted an improved capability to yield crystals suitable for X-ray structural determination.<sup>36,37</sup>

Regardless of the specific mechanisms responsible for efficient ancestral folding in

modern organisms, it is clear that ancestral reconstruction may provide a basis for the sequence engineering of efficient heterologous expression. Yet, reconstructed ancestral sequences typically display extensive differences with respect to the corresponding modern sequences while, in many cases, researchers will be interested in minimally modifying the targeted modern sequence, in such a way that the properties of the encoded protein are barely altered. Minimal sequence perturbation would be particularly desirable when targeting recent adaptations, as, for instance, in metagenomics efforts at contaminated sites aimed at obtaining pollutant-degrading enzymes.<sup>38</sup> Overall, it is of interest to determine whether rescue of inefficient heterologous folding can be engineered on the basis of a few selected back-to-ancestor mutations.

Obligate symbionts typically display high evolutionary rates<sup>39,40</sup> and their proteins may be expected to differ considerably from their modern and ancestral homologs in terms of both sequence and properties. Consequently, the outcome and implications of evolutionary processes may become apparent even in small symbiont proteins that are amenable to detailed biomolecular characterization. Thioredoxins are small proteins (about 110 amino acid residues) that function as general redox catalysts in all known cells.<sup>41</sup> Here we use the thioredoxin from *Candidatus Photodesmus katoptron*, an uncultured symbiont of flashlight fish, as a simple model to explore evolutionary and engineering aspects of protein folding in heterologous hosts.

Flashlight fish (*Anomalopidae*) use light from sub-ocular bioluminescent organs to communicate, hunt prey and disorient predators.<sup>42</sup> Light is produced by luminous bacteria of the *Vibrionaceae* family of *Proteobacteria*. These bacteria are obligate symbionts and have not been cultured in the lab, thus replicating a fundamental scenario of metagenomics studies. Still, their genomes can be sequenced, since the sub-ocular organs of *Anomalopidae* fish harbour large numbers of *Vibrionaceae* in the absence of other bacteria.<sup>43</sup> *Candidatus Photodesmus katoptron*, the luminous bacterium of the *Anomalops katoptron* fish, shows extensive genome reduction and it is highly evolutionary divergent.<sup>44–46</sup> As expected from the high evolutionary rate of its original host, the sequence of the thioredoxin from *Candidatus Photodesmus katoptron* (CPk thioredoxin from now on) differs substantially from all known sequences of thioredoxins from other species. Despite this observation, it is similar to other modern thioredoxins (in particular, to its *E. coli* homolog) in terms of 3D-structure, stability and redox activity. However, as described below, its folding outside the original host appears severely impaired.

CPk thioredoxin displays a very slow refolding *in vitro*, reaching the native state in the time scale of hours<sup>47</sup> and its expression in *E. coli* at 37 °C leads mostly to insoluble protein. By contrast, resurrected

Precambrian thioredoxins have been extensively studied<sup>47–51</sup> and have been found to fold fast *in vitro* and efficiently in *E. coli* despite their huge sequence differences with modern thioredoxins in general and *E. coli* thioredoxin in particular.

Albeit inefficient, the heterologous folding of *CPk* thioredoxin in *E. coli* is not fully impaired and leads to about 20% of soluble protein at 37 °C, a yield that can be increased by carrying out the expression at lower temperatures to decrease protein aggregation.<sup>52</sup> This is a crucial feature that allows us to interrogate the folding properties of *CPk* thioredoxin variants. That is, the effect of sequence modifications on heterologous folding efficiency at 37 °C can be determined and correlated with the biomolecular properties of the corresponding variants of *CPk* thioredoxin, since these variants can actually be prepared in the lab.

We have used computational modelling of the folding landscape to guide back-to-the-ancestor engineering of *CPk* thioredoxin. Specifically, we have used a recently-developed,<sup>53</sup> block version of the Wako-Saitô-Muñoz-Eaton statistical-mechanical model of the folding landscape<sup>54–57</sup> to determine regions of the symbiont thioredoxin that are likely to be unfolded in aggregation-prone intermediate states<sup>58</sup> and we have performed back-to-ancestor sequence-engineering targeted to those regions. The ancestral protein we have used as reference is LPBCA thioredoxin, a putative Precambrian thioredoxin that we have previously characterized in detail<sup>47–51</sup> and that folds efficiently in *E. coli*, despite having only 58% sequence identity with *E. coli* thioredoxin.

Our current study includes several modern/ancestral chimeras, as well as a many single-mutant, back-to-ancestor variants, and allows us to generate several conclusions of general interest:

Folding in the heterologous host is likely akin to unassisted folding. This is supported by (i) the success of the approach used, which involves computational modelling of the unassisted folding landscape, (ii) the fact that the efficiency of heterologous expression correlates with the *in vitro* folding rate, (iii) the very limited rescue of inefficient heterologous expression by chaperone over-expression. Consequently, it appears plausible that ancestral folding efficiency reflects an adaptation to ancient unassisted folding.<sup>47</sup>

Stabilization does improve heterologous folding efficiency, but this cannot be explained by global stabilization alone. Rather, it is linked to specific stabilizing mutations at crucial positions in late-folding regions.

Although the sequences of the ancestral LPBCA thioredoxin and the modern *CPk* thioredoxin differ at 60 positions, the ancestral folding efficiency can be re-enacted in the symbiont thioredoxin with only 1–2 back-to-ancestor mutations. This result provides proof of concept for a minimal-perturbation, sequence-engineering approach to

rescue inefficient heterologous folding with potential application in metagenomics.

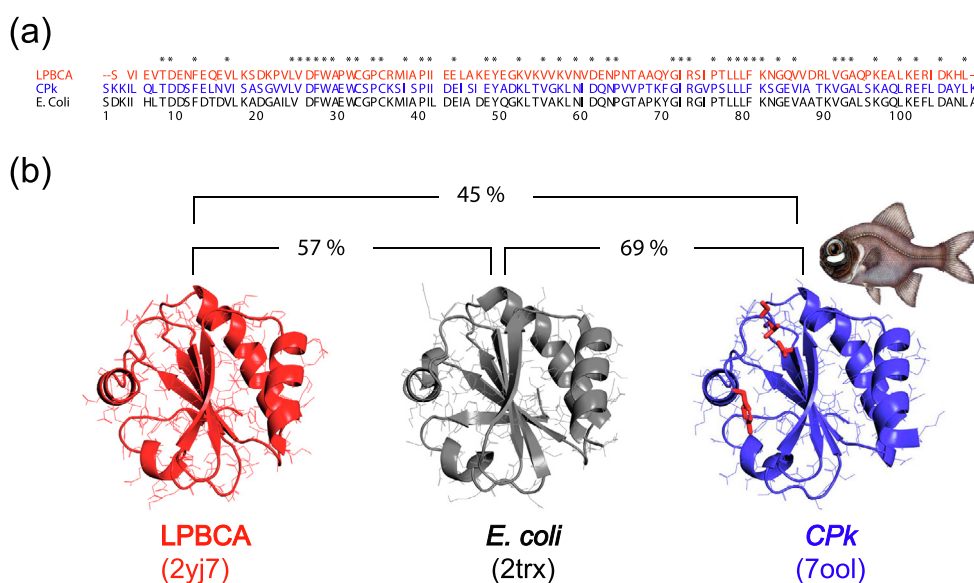
## Results and discussion

**Sequences and biomolecular properties of the modern and ancestral thioredoxins studied in this work.** Our current study utilizes the thioredoxin from the symbiont *Candidatus Photodesmus katoptron* (*CPk* thioredoxin), a modern *E. coli* homolog and a resurrected ancestral thioredoxin corresponding to the last common ancestor of the cyanobacterial, *Deinococcus* and *Thermus* groups, a Precambrian phylogenetic node dated at ~ 2.5 billion years ago. This LPBCA thioredoxin, as well as other resurrected Precambrian thioredoxins, has been previously characterized in detail.<sup>47–51</sup>

*CPk* thioredoxin is highly divergent at the sequence level. A BLAST search in the non-redundant protein sequences (nr) database using as query the sequence of *CPk* thioredoxin yields a *vibrio* protein with only 75% identity as the closest hit. Sequence identity of this thioredoxin from *Candidatus Photodesmus katoptron* (belonging to the *Vibrionaceae* family of *Proteobacteria*) with the thioredoxin from *E. coli* (belonging to the *Enterobacteriaceae* family of *Proteobacteria*) is even lower: 69%. The ancestral LPBCA thioredoxin displays even lower sequence identity to both modern proteins: 57% and 45% with the thioredoxins from *E. coli* and *Candidatus Photodesmus katoptron*, respectively.

Despite the extensive sequence differences (Figure 1a), the three proteins share the thioredoxin fold (Figure 1b). The 3D-structures of *E. coli* thioredoxin and LPBCA thioredoxin have been previously reported.<sup>49,59</sup> The determination of the X-ray structure for *CPk* thioredoxin has been addressed in this work. However, despite numerous attempts using different crystallization conditions and approaches, we failed to obtain crystals of diffraction quality for the wild-type *CPk* thioredoxin (see Methods for details). The structure shown in Figure 1b for *CPk* thioredoxin actually corresponds to an engineered version of the protein in which a short loop (70–77) of the symbiont protein has been replaced by the corresponding loop in the ancestral LPBCA thioredoxin (see details below), a replacement which involves only 4 mutational changes. This variant of *CPk* thioredoxin did produce crystals suitable for diffraction and led to its structural model at 2.85 Å resolution.

The agreement between the three structures shown in Figure 1b is consistent with previous structural work that supported conservation of thioredoxin structure over the span of life on Earth.<sup>49</sup> *In vitro* redox activity, as determined by the insulin aggregation assay and by the assay with thioredoxin reductase coupled to DTNB, is also similar for the three thioredoxins (Figure S1). The two



**Figure 1.** Sequences and structures of modern and ancestral thioredoxins. **a** Alignment of sequences from modern thioredoxins from *E. coli* and *Candidatus Photodesmus katopton* (CPK), and a resurrected ancestral thioredoxin corresponding to the last common ancestor of the cyanobacterial, *Deinococcus* and *Thermus* groups (LPBCA thioredoxin). Positions with identical residues in the three sequences are labelled with asterisks. **b** 3D-structures for the three thioredoxins studied. The experimental structure of CPK thioredoxin actually corresponds to a variant with four back-to-ancestor replacements (highlighted in red; see main text for details). The PDB identifiers are shown, as well as the sequence identity percentages between the three proteins. *Candidatus Photodesmus katopton* is a symbiont of flashlight fish. An illustration of a flashlight fish is shown here, alongside the CPK thioredoxin structure, as well as in the graphical abstract. Illustration used by courtesy of Encyclopædia Britannica, Inc., copyright 2011; used with permission.

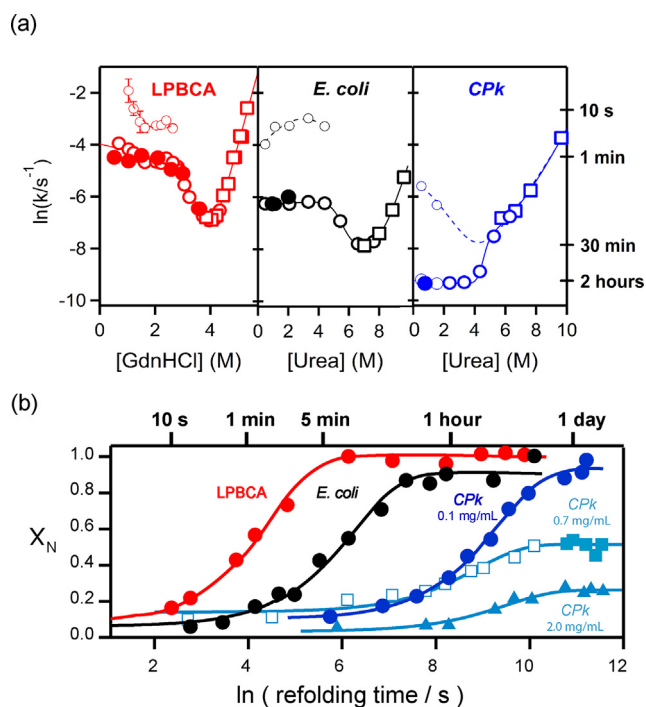
modern thioredoxins display similar high stability as shown by midpoint of about 8 M for urea denaturation experiments<sup>47</sup> and by denaturation temperatures above 80 °C (see details below). The ancestral LPBCA thioredoxin is a hyperstable protein that cannot be denatured by urea at room temperature and that has a denaturation temperature of about 123 °C.<sup>48,60</sup>

***In vitro* folding behaviour of the modern and ancestral thioredoxins studied in this work.** *In vitro* folding of thioredoxins has been known for many years to be a kinetically complex process involving intermediate states and parallel channels to arrive at the native state.<sup>61,62</sup> Such complexities reflect ruggedness of the folding landscape and are revealed by multi-exponential folding kinetics and rollovers in the folding branches of Chevron plots (*i.e.*, plots of folding-unfolding rate constant versus denaturant concentration). Figure 2a shows Chevron plots for the three thioredoxins studied here.<sup>47</sup> To identify the kinetic phase that leads to the native state, *i.e.*, the kinetic phase that defines the time-scale of refolding, we used double-jump unfolding assays,<sup>47,63,64</sup> which allow for a direct determination of the amount of native protein during the folding experiments. These assays (see Methods for details) are a specific instance of the well-known “jump assays” that were developed by pioneers of the *in vitro* protein folding field to resolve

the kinetic complexities of folding processes.<sup>65,66</sup> While the interpretation of the time dependence of a protein physical property may not be straightforward, double-jump unfolding assays lead to a profile for the fraction at native state *versus* time and reveal immediately the time scale in which the native state is reached upon refolding *in vitro*. Such profiles are given in Figs 2a and b for the three thioredoxins studied. It is clear that folding *in vitro* of CPK thioredoxin is substantially slower than the folding of *E. coli* thioredoxin and LPBCA thioredoxin, as discussed above. Furthermore, the *in vitro* folding of CPK thioredoxin is also inefficient, as shown by the fact that substantial amounts of protein fail to reach the native state (Figure 2b). The effect is more pronounced the higher the total concentration of protein, suggesting that *in vitro* folding inefficiency is linked to protein aggregation (which is, in fact, observed visually). The rate constants derived using double-jump unfolding assays (closed symbols in the plots of Figure 2a) indicate that the native state is reached mostly in the slow kinetic phase detected by fluorescence.

While the focus of this work lies on folding rates and folding efficiency, it is also of interest to comment briefly on the unfolding rates since these are related with kinetic stability, an important protein property. *E. coli* thioredoxin is known to be a kinetically stable protein. This is clearly shown





**Figure 2.** *In vitro* folding of modern and ancestral thioredoxins. **a** Chevron plots of folding-unfolding rate constant at pH 7 and 25 °C versus denaturant concentration for the three thioredoxins studied. Urea is used as denaturant for *E. coli* and *CPk* thioredoxins. LPBCA thioredoxin, however, is highly stable and cannot be denatured by urea at 25 °C. The chevron plot using the stronger denaturant guanidinium hydrochloride is shown for this protein. Still, the folding rates for LPBCA thioredoxin at low denaturant concentration obtained with urea and guanidine are in good agreement.<sup>47</sup> Circles and squares refer to experiments performed in the folding and unfolding directions, respectively. Error bars are standard errors derived from fits to the experimental profiles and are not shown when they are smaller than the size of the data point. Data are taken from Gamiz-Arco et al.<sup>47</sup> and were derived from fluorescence kinetic profiles. These profiles are often multiphasic in the folding direction, with the slow folding phase leading to the native state, as shown by the agreement with the folding rates (closed symbols) derived using double-jump unfolding assays. Lines shown are meant to guide eye. **b** Profiles of fraction of native state vs. time obtained by using double-jump unfolding assays. The refolding time is the time the protein is allowed to refold after the first jump, i.e. the time elapsed between the start of the folding process and its interruption (see Methods for details). Experiments were performed at pH 7, 25 °C in the presence of 1 M urea (see Methods for details). Protein concentration was 0.1 mg/mL, except for the profiles for *CPk* thioredoxin at 0.7 mg/mL and 2 mg/mL. The open symbols in the profile for *CPk* thioredoxin at 0.7 mg/mL are taken from Gamiz-Arco et al.<sup>47</sup>. The lines represent the best fits of a single exponential. Note that an exponential has a sigmoidal shape in a plot versus  $\ln t$ . In both **a** and **b**, typical values of the lifetime (calculated as the inverse of the first-order rate constant) are indicated to highlight the large differences in folding time-scale between the three proteins.

by the extrapolation to zero denaturant concentration of the unfolding branch in the Chevron plot which leads to a very low unfolding rate constant and an unfolding time scale on the order of a few months.<sup>67</sup> The unfolding of LPBCA thioredoxin is about three orders of magnitude slower than the unfolding of *E. coli* thioredoxin,<sup>51</sup> indicating a much-enhanced kinetic stabilization for the ancestral protein. By contrast, the unfolding of the symbiont *CPk* thioredoxin is somewhat faster than the unfolding of *E. coli* thioredoxin (Figure 2a). Still, an extrapolation to zero denaturant concentration of the unfolding branch of the Chevron plot yields an unfolding time scale on the order of days. Therefore, although the symbiont thioredoxin is less

kinetically stable than *E. coli* thioredoxin, it seems to retain some substantial level of kinetic stability. The differences in kinetic stability between the three thioredoxins may reflect, at least in part, the different living temperatures of the host organisms, as we have previously discussed.<sup>68</sup>

**Efficiency of the heterologous expression in *E. coli* of modern and ancestral thioredoxins.** We determined the efficiency of expression in *E. coli* at 37 °C as the ratio of soluble protein to total protein determined after overexpression for 3 h, as recommended by standard protocols (see Methods for details).

As was to be expected, expression of *E. coli* thioredoxin in *E. coli* is highly efficient, leading to

essentially 100% soluble protein. Remarkably, the expression of the ancestral LPBCA thioredoxin is also highly efficient, despite the extensive sequence differences with *E. coli* thioredoxin (only 57% sequence identity). On the other hand, expression of *CPk* thioredoxin in *E. coli* only leads to about 20% soluble protein (Figure 3).

Chaperone over-expression is a common host-engineering strategy to improve heterologous protein expression.<sup>2,8</sup> We attempted to rescue the inefficient folding of the *CPk* thioredoxin in *E. coli* by complementation with plasmids containing genes of the following *E. coli* chaperones: trigger factor, groES, groEL, dnaK, dnaJ. We used five separate combinations of these chaperones, as shown in Figure 3a. Only very moderate enhancements in folding efficiency were observed (Figure 3).

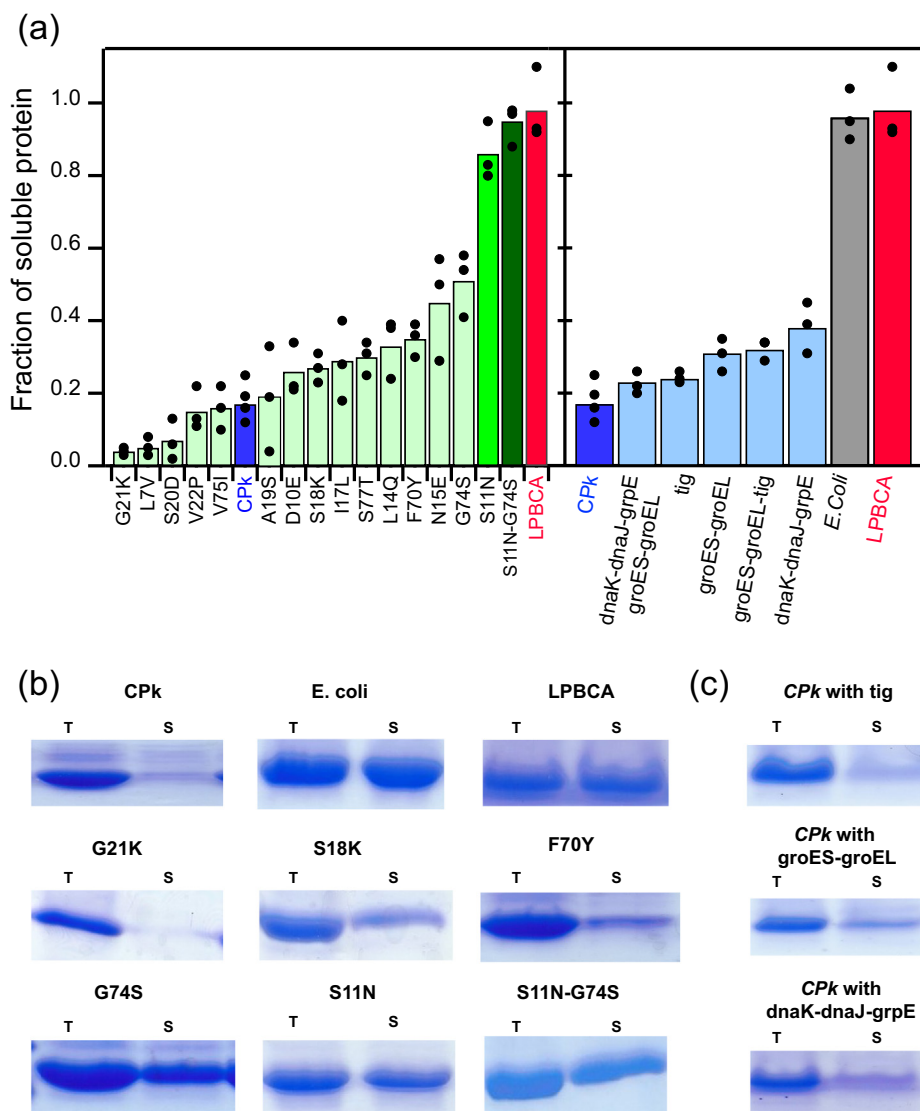
**A plausible evolutionary narrative.** The experiments described above (Figures 2 and 3) reveal a striking disparity between the folding behaviour of the *CPk* thioredoxin and LPBCA thioredoxin outside their original hosts. The ancestral protein folds fast *in vitro* and its expression in *E. coli* is efficient. By contrast, expression of the symbiont thioredoxin in *E. coli* produces a substantial amount of insoluble protein and refolding experiments show that it reaches the native state *in vitro* in the time scale of hours. This slow and inefficient folding can hardly be assumed to correspond to the situation *in vivo* in the original host. Note that the synthesis of a ~ 100 residue protein by bacterial ribosomes takes about 5 seconds.<sup>58</sup> Obviously, it is difficult to understand that a small protein, which can in principle fold fast, has been selected during evolution to fold in its original host in a time scale ~ 3 orders of magnitude above the time required for synthesis in the ribosome. A much more likely scenario is that folding of *CPk* thioredoxin in its original host is fast, plausibly allowing for co-translational folding, and efficient. This fast/efficient *in vivo* folding would obviously be the result of the interaction of the protein with the cellular folding assistance machinery, mainly the ribosome itself<sup>59,70</sup> and the ribosome-binding chaperones (the trigger factor), which would guide and assist co-translational folding, although a role for downstream chaperones and the specific environment in the symbiont (pH, redox, etc.) cannot be ruled out.

In fact, it is not unreasonable that the folding of thioredoxins in at least some modern organisms be assisted. A commonly accepted view is that folding is often inefficient and, therefore, it may require assistance, for proteins larger than 100 residues.<sup>58</sup> According to this, thioredoxins, with about 110 residues, would be a borderline case. However, thioredoxins have a serious folding problem that was already noted by Fred Richards many years ago<sup>61</sup>: the proline residue at position 76, which is essential for activity and strictly conserved,

is in *cis* conformation, which generates a well-known kind of folding kinetic bottleneck. In fact, folding of *E. coli* thioredoxin in the test tube is known to be a complex process.<sup>62</sup> Overall, it appears plausible that thioredoxin folding *in vivo* may benefit from assistance and there can be little doubt that this is in fact the case with the symbiont thioredoxin we study in this work, since its unassisted folding outside its original host is seriously impaired. Such assistance, however, would not be available for *CPk* thioredoxin in the heterologous *E. coli* host, since co-evolution in the original host would lead to its adaptation to the folding assistance machinery of the symbiont. Note that, not only the symbiont thioredoxin, but also most of the symbiont chaperones and ribosomal proteins are highly divergent at the sequence level (see Tables S2 and S3). Therefore, co-evolution of interacting symbiont proteins is a likely scenario.

As we have recently noted,<sup>47</sup> since evolution has no foresight, folding assistance cannot have arisen before protein folding itself and, therefore, the most ancient proteins could plausibly fold with little or no assistance. That is, ancestral folding efficiency plausibly relied on fast unassisted folding that limits the transient population of aggregation-prone partially-unfolded states. Consequently, efficient folding of resurrected ancestral proteins in modern organisms may plausibly reflect an ancient adaptation to unassisted folding. Thioredoxins could in fact provide a clear example of this ancient unassisted-folding scenario, because there are thioredoxins in the three domains of life and their emergence can, therefore, be traced back to the last universal common ancestor, LUCA. Certainly, it is not known when efficient folding assistance emerged. Still, it is clear that efficient unassisted folding would no longer be a useful feature after the evolutionary emergence of cellular folding-assistance, which would thus allow the evolutionary acceptance of mutations that impair the ancestral feature. Such degradation of unassisted folding would be of no consequence for folding in the original host, but will lead to inefficient expression in a heterologous host, where folding assistance would not be available due to co-evolution in the original host. Furthermore, it is not clear how long it took for the ancient trait (efficient unassisted folding) to degrade after cellular folding assistance was available. One possibility is that substantial degradation of ancient unassisted-folding of small proteins, such as thioredoxin, may have only occurred to a substantial extent in organisms with a high evolutionary rate, such as the obligate symbiont we consider in this work.

To summarize, it appears plausible that the folding of the most ancient thioredoxins was unassisted, that the folding of *CPk* thioredoxin in its original host is assisted and that, as a result of co-evolution in the original host, efficient assistance is not available for the symbiont thioredoxin in the heterologous *E. coli* host.

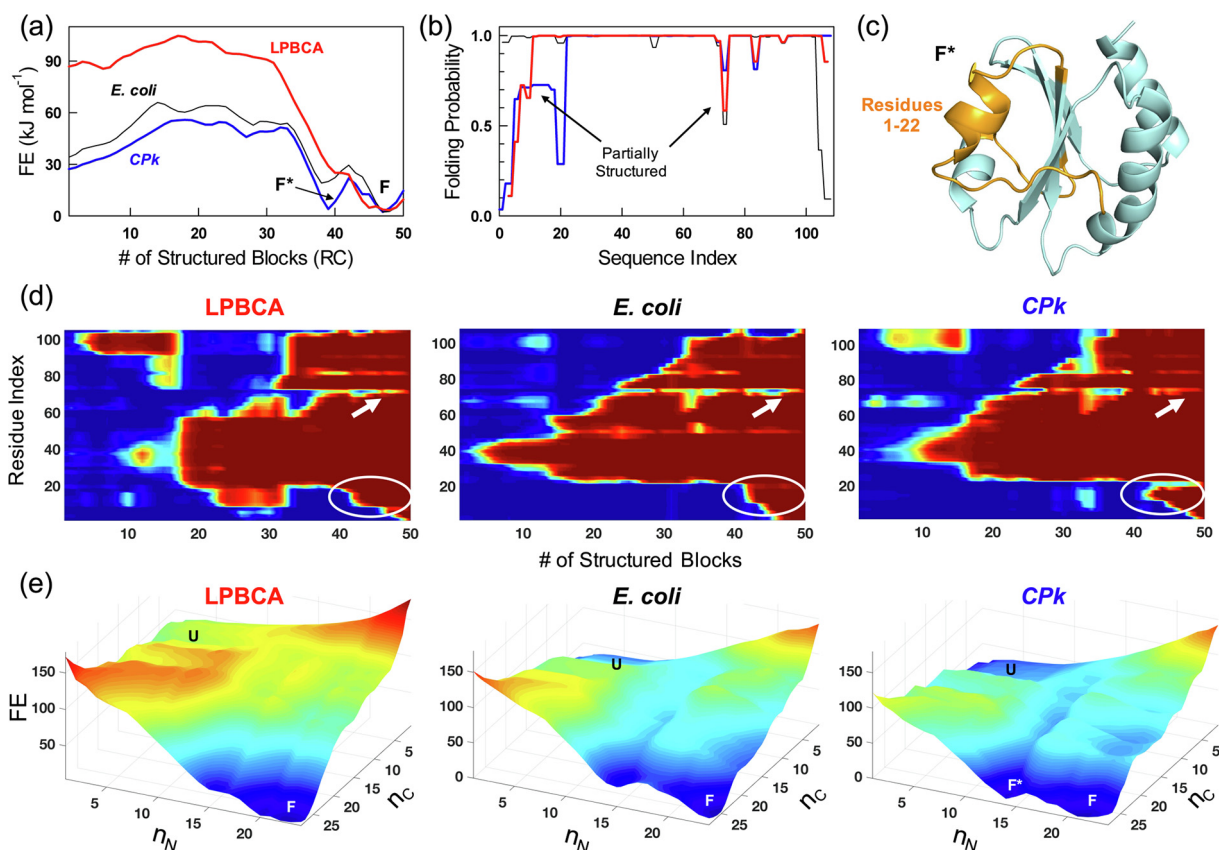


**Figure 3.** Expression in *E. coli* of modern and ancestral thioredoxins. **a** Fraction of soluble protein obtained upon expression of thioredoxin proteins in *E. coli* at 37 °C. The box at the left displays data for LPBCA thioredoxin, CPk thioredoxin, variants of the latter with single mutations and one double-mutant variant (S11N-G74S). The box at the right displays data for *E. coli* thioredoxin, LPBCA thioredoxin, CPk thioredoxin and data for the latter with over-expression of various chaperone teams. Bars represent the average of several independent determinations and the individual values are also shown. **b** Illustrative examples of the experimental determination of the fraction of soluble protein for LPBCA thioredoxin, *E. coli* thioredoxin, CPk thioredoxin and variants of the latter. **c** Representative examples of the experimental determination of the fraction of soluble protein for CPk thioredoxin with over-expression of chaperone teams. In both **b** and **c**, T and S represent “total” and “soluble”, respectively. For illustration, only sections of the SDS-PAGE gels with the thioredoxin bands are shown. Complete gels are shown Figures S10-S13.

Evolutionary narratives are necessarily speculative. Still, the narrative we propose has the merit of immediately suggesting an approach to the sequence-engineering of heterologous expression. That is, suitable back-to-ancestor mutations could lead to a more efficient heterologous expression. Furthermore, since folding in the heterologous host is, at least to some extent, unassisted, computational modelling of the unassisted folding-landscape may be used to guide back-to-ancestor

engineering for efficient heterologous expression. Computational modelling of the folding landscape for thioredoxins is described in the next section.

**Computational modelling of the folding landscape for modern and ancestral thioredoxins.** Here we use a recently developed version<sup>53</sup> of the Wako-Saitô-Muñoz-Eaton (WSME) statistical model of protein folding<sup>54–57</sup> to assess the main features of the folding landscape of modern and ancestral thioredoxins (Figure 4).



**Figure 4.** Statistical mechanical modelling of the folding landscape for modern and ancestral thioredoxins. A block version of the Wako-Saitô-Muñoz-Eaton model was used in all the calculations shown here. **a** Profiles of free energy versus number of structured blocks at 37 °C for the modern thioredoxins from *E. coli* and *Candidatus Photodesmus katoptron* (*CPk*), and the ancestral LPBCA thioredoxin. Note that a partially-unfolded intermediate ( $F^*$  arrow), clearly differentiated from the fully folded protein ( $F$ ), is distinctly observable only in *CPk* thioredoxin. **b** Residue folding probabilities as a function of sequence index at 37 °C following the colour code in panel **a**. **c** The predicted structure of  $F^*$  with the partially structured residues 1–22 highlighted in orange. **d** Folding probability, coloured in the spectral scale from blue (0) to red (1), as a function of a plausible reaction coordinate, the number of structured blocks, for the three thioredoxin studied. The N-terminal region that folds the last is highlighted by white ovals while the 70–77 region is highlighted by an arrow. **e** Free energy landscapes (z-axis in  $\text{kJ}\cdot\text{mol}^{-1}$ ) as function of  $n_N$  and  $n_C$ , the number of structured blocks in the N- and C-terminal half of the protein, respectively, for the three thioredoxins studied here.

As mentioned above, thioredoxin folding has been known from many years to be a complex process that involves intermediate states and parallel channels to arrive to the native state,<sup>61,62</sup> reflecting a rugged folding landscape. We do not aim here at reproducing these kinetic complexities in detail, but mostly at identifying regions of the thioredoxin molecule that are likely to be unfolded in intermediate states of the folding landscape. The rationale behind this approach is that such unfolded or partially-structured regions may be involved in aggregation and other undesirable interactions<sup>58</sup> and are, therefore, obvious targets for engineering efforts aimed at rescuing inefficient folding. This notion is consistent with early studies on the

directed evolution of proteins for enhanced heterologous expression which demonstrated improvements in the stability and solubility of intermediates.<sup>71</sup>

We employ the block version of WSME model (the bWSME model), which considers 2–3 consecutive residues as blocks to reduce the protein phase space and thus rapidly calculate free energy profiles (see Methods). In the block description, the model still considers > 490,000 microstates compared to the residue-level version that would involve considering the contribution to the partition function from > 9.9 million microstates. We first reproduce the apparent experimental equilibrium stability differences to calibrate the model. The resulting one-dimensional



folding free energy profiles as a function of the number of structured blocks, the reaction coordinate (RC), are quite similar for the three proteins, but with one major difference – *CPk* thioredoxin populates a partially structured intermediate ( $F^*$ ) on the folding side of the main barrier with a population of 27% (Figure 4a). However, neither of LPBCA or *E. coli* thioredoxins populate  $F^*$  significantly (<0.1%). It is important to note that this difference is not a consequence of the larger stability of *CPk* thioredoxin as *E. coli* thioredoxin does not populate  $F^*$  despite exhibiting similar stability (black in Figure 4a). Moreover, WSME model calculations reveal that the shape of the free-energy profiles is conserved under iso-stability conditions (of 25 kJ mol<sup>-1</sup>; Figure S2). These observations indicate that the intermediate  $F^*$  is intrinsic to the *CPk* thioredoxin conformational landscape and is independent of the overall thermodynamic stability.

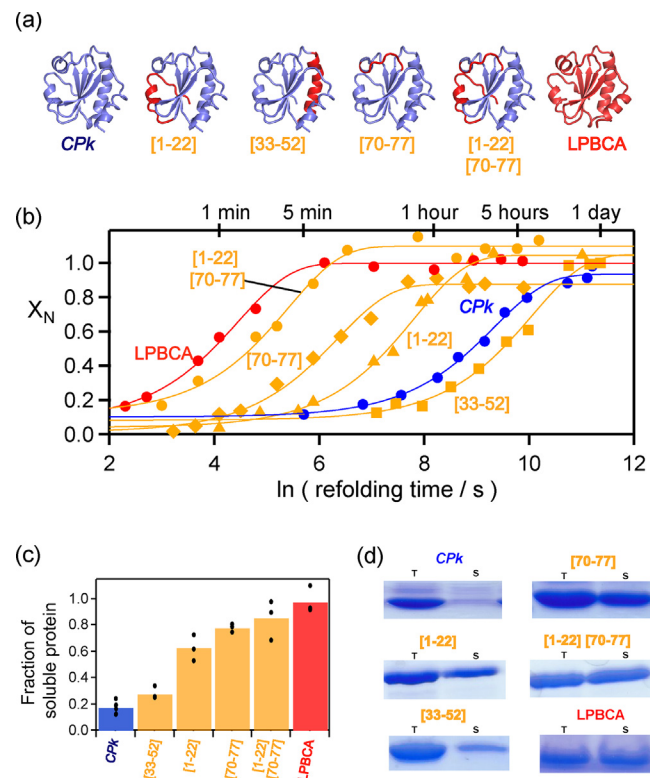
To identify the regions of *CPk* thioredoxin that are partially structured, we computed residue folding probabilities that quantify the extent to which every residue is structured at a given thermodynamic condition. At 37 °C, the N-terminal region of *CPk* thioredoxin (residues 1–22) is partially structured when compared to its ancestral counterpart (Figure 4b). Differences in folding probabilities are also evident in two other regions: residues 70–77 that harbours a critical cis-proline<sup>47,61</sup> and to a lesser extent in the residue stretch 83–84. We further calculated the probability of every region of the protein to be structured as a function of the reaction coordinate (Figure 4d), and find that the N-terminal region of the protein folds the last, thus revealing the identity of  $F^*$  (Fig 4c and white ovals in Figure 4d). It can also be seen that the residue-stretch 70–77 exhibits equilibrium fluctuations during the folding for both the proteins (arrows in Figure 4d). The two-dimensional free energy landscape (Figure 4e), constructed by accumulating partial partition functions involving combinations of a given number of residues structured at N- and C-terminal halves of the protein, highlights that  $F^*$  is the most populated state in *CPk* thioredoxin apart from multiple partially structured states that likely contribute to the slow folding of *CPk* thioredoxin (Figure 4e) compared to LPBCA thioredoxin.

There is a remarkable congruence between the computational predictions and the experimental folding kinetic data for *CPk* thioredoxin. Thus, a rate for reaching the native state much slower than expected from the typical shape of a chevron plot (compare continuous and discontinuous lines in the plot for *CPk* thioredoxin in Figure 2a) is the pattern expected from the accumulation of an intermediate. Furthermore, the structure predicted for  $F^*$  (Figure 4c) is supported by results obtained with modern/ancestral chimeras described in the next section.

**Efficiency of heterologous expression in *E. coli* of modern/ancestral thioredoxin chimeras.** On the basis of the folding-landscape computations described in the preceding section, we selected two regions as targets for the engineering of heterologous folding: the N-terminal 1–22 fragment, which includes a short  $\alpha$ -helix and large stretches of non-regular structure, and the 70–77 loop which includes a cis-proline residue that has been known for many years to be critical for thioredoxin folding.<sup>47,61</sup> These two regions are predicted to fold comparatively late and to remain unfolded during most of the protein folding process (Figure 4d). Plausibly, therefore, they may be involved in intermolecular processes that lead to insoluble protein. Since the heterologous expression of the ancestral LPBCA thioredoxin is efficient, we have prepared chimeras in which these regions are replaced by the corresponding ancestral sequences (see Figure 5a). These replacements involve 17 mutational changes in the case of the 1–22 fragment and 4 mutational changes in the case of the 70–77 loop. We designate the chimeras as *CPk*-[1-22] thioredoxin and *CPk*-[70-77] thioredoxin. We have also studied the “double chimera” in which both regions are simultaneously replaced by the corresponding ancestral sequences: *CPk*-[1-22]-[70-77] thioredoxin. The three chimeras show substantially improved heterologous expression in *E. coli* with *CPk*-[1-22]-[70-77] thioredoxin approaching 100% of soluble protein (Figure 5c and d).

In addition to the two regions referred to in the preceding paragraph, we have also selected for experimental analysis the 33–52 region which matches the longest  $\alpha$ -helix in the thioredoxin molecule. This region is selected to provide an obvious control experiment, since it is predicted to fold early (Figure 4d) and, according to our working hypothesis, we do not expect its replacement with the corresponding ancestral sequence to improve the efficiency of heterologous folding. The experimental results on *CPk*-[33-52] thioredoxin conform to this expectation (Figure 5c and d).

**Efficiency of heterologous expression of single-mutant variants of the symbiont thioredoxin.** As described in the preceding section, efficient heterologous expression of the symbiont thioredoxin is achieved through replacement of the 1–22 and 70–77 regions with the corresponding ancestral sequences in LPBCA thioredoxin. To explore the individual mutational contributions to the rescue, we have prepared 15 back-to-ancestor, single-mutant variants. These include the four back-to-ancestor mutations in the 70–77 region and 11 back-to-ancestor mutations in the 1–22 region. We have excluded from this analysis the initial 1–6 N-terminal segment, since it appears to be unstructured in the native structures. For all the 15 single-mutant variants



**Figure 5.** *In vitro* folding and expression in *E. coli* of modern/ancestral chimeras. **a** Definition and structural description of the studied modern/ancestral chimeras. The thioredoxin backbone is coloured to indicate the origin of the sequence: modern *CPk* thioredoxin (blue) or ancestral LPBCA thioredoxin (red). **b** Profiles of fraction of native state *versus* time for the *in vitro* folding of *CPk* thioredoxin, LPBCA thioredoxin and the several modern/ancestral chimeras defined in **a**. The values of the fraction of native state are derived from double-jump unfolding assays (see Methods for details). The plot is labelled with characteristic time values to highlight the wide range of folding times for the proteins studied. The lines represent the best fits of a single exponential. **c** Fraction of soluble protein obtained upon expression in *E. coli* at 37 °C of *CPk* thioredoxin, LPBCA thioredoxin and the four chimeras defined in **a**. Bars represent the average of several independent determinations and the individual values are also shown. **d** Representative examples of the experimental determination of the fraction of soluble protein for LPBCA thioredoxin, *CPk* thioredoxin and the modern/ancestral chimeras defined in **a**. Complete gels are shown Figures S10–S13.

we have determined the heterologous expression efficiency in *E. coli* (Figure 3a). The most remarkable result of these studies is that a single mutation at the 1–22 segment, S11N, rescues most of the inefficient heterologous folding. Combining this mutation with the best-rescuing mutation in the 70–77 loop leads to a double mutant variant of the symbiont thioredoxin, S11N/G74S, that approaches 100% soluble protein.

**Inefficient heterologous expression of CPk thioredoxin is rescued by back-to-the-most-ancient-ancestor mutations.** It is important to note, first of all, that the inefficient heterologous folding of the symbiont *CPk* thioredoxin in *E. coli* is rescued by mutations that are back-to-ancestor, but not back-to-*E. coli*. The reason is obviously that the residues at position 11 and 74 in *E. coli* thioredoxin and *CPk* thioredoxin are identical: S and G, respectively<sup>48</sup> (see also Figure 1 in Ingles-Prieto et al.<sup>49</sup>). Furthermore, the N at position 11 and S at position 74 are very likely the residues pre-

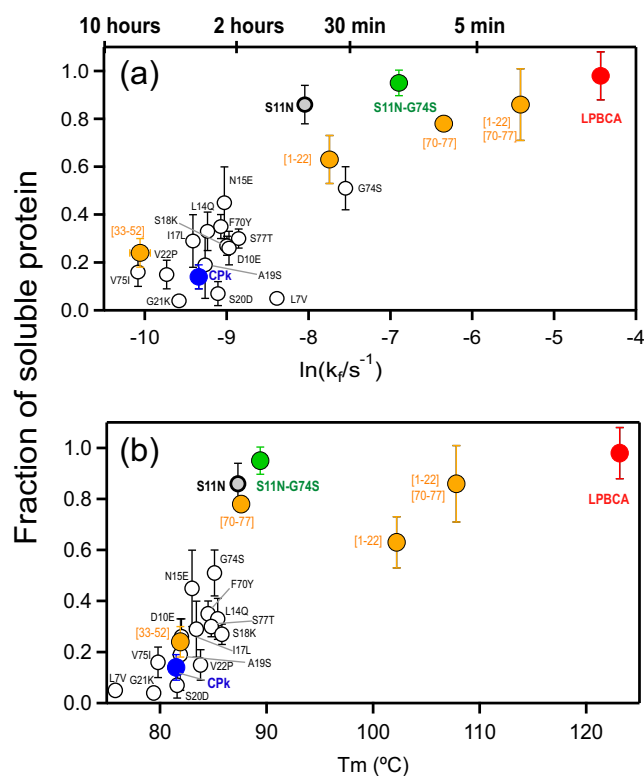
sent in the most ancient thioredoxins, as indicated by the posterior probabilities in the thioredoxins corresponding to the last common ancestor of bacteria (LBCA) and the last archaeal-eukaryotic common ancestor (AECA). The values are 0.999 (N at position 11, LBCA), 1.000 (S at position 74, LBCA), 0.991 (N at position 11, AECA) and 0.990 (S at position 74, AECA). As we have previously shown, sites with such high probability are very rarely, if ever, incorrectly predicted.<sup>72</sup> Our reason for providing the posterior probabilities at LBCA and AECA is that these nodes are immediately below the last universal common ancestor (LUCA) in the phylogeny used for ancestral reconstruction. The sequences of proteins in LUCA cannot be reconstructed by standard methodologies because an outgroup to determine LUCA in the phylogenetic tree is not available. Still, since N and S are with very high probability the residues at positions 11 and 74 in the nodes immediately below LUCA, it can be reasonably inferred that N and S were also the residues

at positions 11 and 74 in LUCA thioredoxin. Overall, even if ancestral reconstruction is unavoidable uncertain to some extent, there is little doubt about the identity of the amino acids in the most ancient thioredoxins at the crucial positions 11 and 74.

**Correlation between heterologous folding efficiency and the *in vitro* folding rate.** The modern-ancestral chimeras, the single-mutant variants and the double S11N/G74S variant studied here are all active and show levels of *in vitro* redox activity similar to the modern *CPk* thioredoxin and *E. coli* thioredoxin, as well as the ancestral LPBCA thioredoxin (Figure S1). They differ substantially, however, in terms of *in vitro* refolding rate. We have used double-jump unfolding assays to determine the rate constant for the last stage of *in vitro* refolding process, *i.e.*, the stage that leads to the native protein and defines the time scale of folding. These experiments reveal a very large (~400-fold) range of folding rates (Figures 2b, 5b and 6) with the time scale in which these proteins reach the native state *in vitro* varying between a few minutes and many hours. Remarkably, there is a good correlation between the efficiency of heterologous

folding and the folding rate *in vitro* (Figure 6a), supporting that efficient heterologous folding is achieved through the reduction of the time partially-unfolded states are significantly populated during the folding process.

**Relation between heterologous folding efficiency and protein stability.** We have used the denaturation temperature, as determined by differential scanning calorimetry, as a simple metric for the stability of the modern/ancestral chimeras, the single-mutant variants and the double S11N/G74S variant. For some selected variants, we have also performed urea denaturation studies. All the variants studied incorporate back-to-ancestor modifications. The ancestral sequence used as reference is that of LPBCA thioredoxin, a hyperstable protein with a very high denaturation temperature.<sup>48,60</sup> As anticipated, therefore, the back-to-ancestor modifications produce stability enhancements with respect to the *CPk* thioredoxin background in most cases. Furthermore, there appears to be a reasonable correlation between heterologous folding efficiency and stability, as described by the denaturation temperature values (Figure 6b).



**Figure 6.** Correlations of the efficiency of heterologous expression with *in vitro* folding rate and protein stability. **a** Plot of fraction of soluble protein obtained in the expression in *E. coli* versus the logarithm of the *in vitro* folding rate constant including *CPk* thioredoxin, LPBCA thioredoxin, several variants of *CPk* thioredoxin and several modern/ancestral chimeras (Figure 5). Typical values of the lifetime are indicated to highlight the wide range of folding time-scales. **b** Plot of fraction of soluble protein obtained in expression in *E. coli* versus denaturation temperature values derived from differential scanning calorimetry experiments. The proteins included here are the same as those included in the plot of panel **a**.

It is important to note, however, that global stability alone cannot explain the rescue of inefficient heterologous folding by back-to-ancestor modifications. This is more clearly seen when comparing the data of *CPk*-[1–22] thioredoxin with those for the variant of *CPk* thioredoxin with the single S11N mutation (Figure 7). Both scanning calorimetry data and chemical denaturation profiles indicate that replacement of the 1–22 segment in *CPk* thioredoxin with the corresponding LPBCA ancestral sequence brings about a very large stabilization, while the single S11N mutation produces a much more moderate stability enhancement (Figure 7a and b). Yet, heterologous folding of the single S11N variant is even more efficient than that of *CPk*-[1–22] thioredoxin (Figure 7c). Obviously, much of the stabilization brought about by the replacement of the 1–22 segment in *CPk* thioredoxin with its ancestral counterpart has little effect on heterologous folding efficiency.

The pattern described above for the heterologous folding efficiency is replicated by the *in vitro* folding rates. That is, the folding rate for the S11N variant of *CPk* thioredoxin is similar to the folding rate of the *CPk*-[1–22] thioredoxin (Figure 6a), despite of the much higher stability of the chimera as probed by the denaturation temperature value (Figure 6b). This result is not surprising, since stability enhancement does not necessarily lead to a faster folding rate. In fact, no significant correlation between stability and folding rate was found in an experimental analysis of a set 13 modern thioredoxins from different species.<sup>47</sup> As another example, the design of a completely symmetric  $\beta$ -trefoil led to a protein displaying both very high stability and very slow folding.<sup>73,74</sup> More generally, the uncoupling between thermodynamics and kinetics linked to the presence of populated intermediates is a known phenomenon, already noted in the seminal work of Agard and coworkers on  $\alpha$ -lytic proteases.<sup>75</sup>

Overall, our results are consistent with the notion that protein stabilization may improve heterologous expression,<sup>76</sup> but support that the rescuing effect of stabilization is linked to specific mutations in regions of the protein that are likely unfolded in aggregation-prone intermediate states and that have a strong effect on the unassisted folding rate.

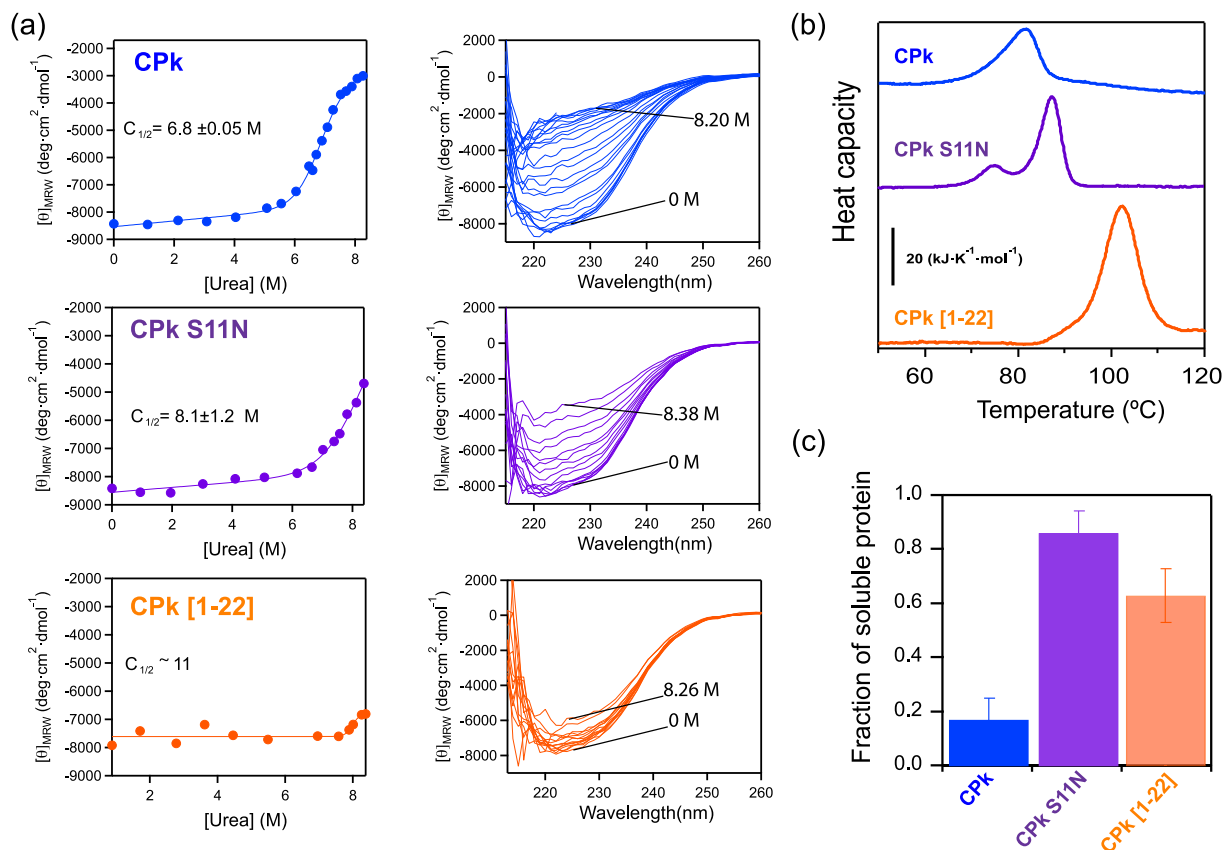
Finally, a somewhat intriguing result of our extensive experimental study on the stability of *CPk* thioredoxin variants deserves some attention. While most of the variants show the expected stability enhancement, this is not the case with *CPk*-[33–52] thioredoxin. This modern/ancestral chimera involves 10 back-to-ancestor amino acid replacements in the longest  $\alpha$ -helix of the thioredoxin structure, which is expected to fold early according to our statistical-mechanical calculations (Figure 4). *CPk*-[33–51] thioredoxin

displays a stability similar to that of the *CPk* thioredoxin background as shown by both the denaturation temperature values and the chemical denaturation profiles (Figure S3). One interesting possibility is that the ancestral stabilization is an adaptation to the need of efficient unassisted folding and high kinetic stability in an ancient environment. Consequently, it is not implemented in regions of the protein that are folded in aggregation-prone intermediates of the folding landscape and in the transition state that determines kinetic stability. Another possibility is that early folding of the long helix is crucial, mutations that impair its stability are not therefore accepted during evolution and, consequently, modern thioredoxins preserve the ancestral stability of the long helix. These interpretations are obviously speculative at this stage and will be explored in future work.

**Structural basis of the rescue of inefficient heterologous folding.** Inefficient heterologous expression of *CPk* thioredoxin in *E. coli* is rescued to a substantial extent by a single S11N mutation. Position 11 is part of a type IV turn involving residues 8–11 (Figure 8a). Turns are known to be crucial for protein folding in general, as they allow the polypeptide chain to fold onto itself and generate interactions that pertain to the native structure.<sup>77</sup> Presence of an asparagine residue at position 11 promotes the turn conformation because this residue can form stabilizing hydrogen bonds with the threonine at position 8 (Figure 8b). On the other hand, a serine at position 11 is not predicted to form stabilizing hydrogen bonds with the threonine at position 8 (Figure 8b) and consequently facilitates alternative conformations for the 8–11 segment in the high energy regions of the folding landscape. That is, the back to the ancestor mutation at position 11 should favour the local native conformation with respect local unfolded conformations, thus decreasing the time during which the corresponding partially unfolded states are significantly populated in the course of folding.

A similar explanation can be deduced for the effect of the G74S mutation included in the S11N/G74S variant that approaches 100% soluble protein in heterologous expression. As we have previously noted,<sup>47</sup> effects on folding of the G/S exchange at position 74 in thioredoxins are very likely related the fact that glycine has no side chain and places little restriction in local backbone conformation. The flexible link generated by the presence of a glycine residue will allow many different conformations in the high energy region of the folding landscape. This is particularly relevant for the 70–79 loop, since it also includes the proline residue at position 76, which is in the rare *cis*-conformer in the native structure.<sup>61</sup> Presence of a glycine residue at position 74 thus enables many conformations for the 70–79 loop that are not consistent with the native *cis* conformation for Pro76.





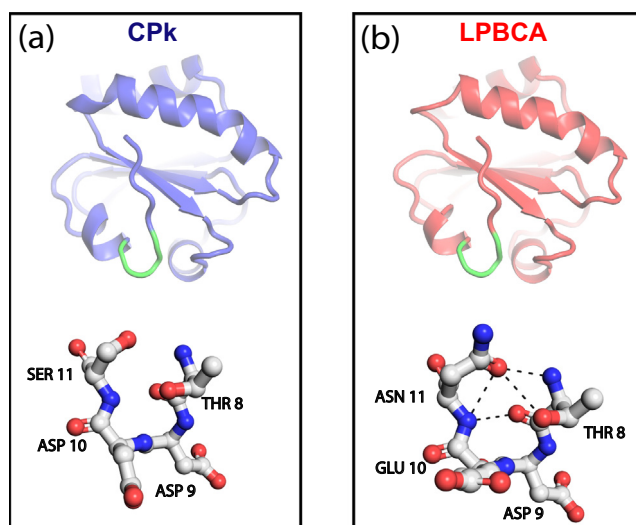
**Figure 7.** Relation between stabilization and rescue of inefficient heterologous expression. **a** and **b** Stability of the symbiotic *CPk* thioredoxin, its S11N variant and the chimera in which the 1–22 segment has been replaced by the corresponding ancestral LPBCA sequence (*i.e.*, *CPk*[1–22] thioredoxin). In the experiments shown in **a** stability is probed by urea-induced denaturation followed by circular dichroism. Both, the original spectra at several urea concentrations and the profiles of ellipticity at 222 nm *versus* denaturant concentration are shown. The continuous lines represent the best fits of a two-state model (see Methods) and mid-point urea concentrations derived from the fits are shown. In **b** stability is studied by differential scanning calorimetry (see Methods and Figure S15 for details). While the mutation S11N has a significant but moderate stabilizing effect, the chimera is highly stable as revealed by a high denaturation temperature and resistance to denaturation at 25 °C by high urea concentrations. **c** Fraction of soluble protein obtained upon expression in *E. coli* at 37 °C of *CPk* thioredoxin, its S11N variant and the *CPk*[1–22] thioredoxin chimera. Despite its much-enhanced stability, the chimera is less efficient at rescuing heterologous expression than the S11N variant.

## Concluding remarks

Our results support that the folding of proteins in heterologous hosts may be akin to some extent to unassisted folding. For the specific protein system studied here, this is supported by (i) the success of the approach used to rescue inefficient heterologous expression, which involved computational modelling of the unassisted folding landscape, (ii) the fact that the efficiency of heterologous expression correlates with the *in vitro* folding rate, (iii) the very moderate rescue of inefficient heterologous expression by chaperone over-expression.

Unassisted folding in heterologous hosts may conceivably result from the overexpression of the

protein exceeding the capacity of the folding-assistance machinery. This is a reasonable scenario, in particular since co-evolution may have led to the adaptation of the protein to the assistance machinery of its original host. Consequently, the fact that resurrected ancestral proteins often show improved heterologous expression as compared with their modern counterparts<sup>22–35</sup> plausibly reflects an ancient adaptation to unassisted folding. Efficient unassisted folding would no longer be a useful feature after the evolutionary emergence of cellular folding-assistance, thus allowing the evolutionary acceptance of mutations that impair the ancestral feature. Reversal of such mutations could then lead to a more efficient heterologous expression.



**Figure 8.** 3D-structures of the modern *CPk* thioredoxin and the ancestral LPBCA thioredoxin highlighting a type IV turn including position 11. Highlights of the turn show the hydrogen bonds involving the residue in position 11 with other residues in the turn as predicted by WHAT IF<sup>94</sup> with hydrogen-bond network optimization.<sup>95</sup> Hydrogen bonds are only predicted for the ancestral residue, which is therefore expected to stabilize the native turn conformation and disfavour non-native conformations.

Our results support that a few selected back-to-ancestor mutations can re-enact the folding efficiency of the resurrected ancestral proteins and point, therefore, to a minimal-perturbation, sequence-engineering approach to resolve inefficient heterologous expression. Some details of the practical application of the approach are noted below.

Prediction of back-to-ancestor mutations should be feasible for most protein systems, given the availability of large sequence databases and various software packages for the several steps of ancestral sequence reconstruction (for a recent account, see ref. <sup>21</sup>). Certainly, experimental screening for expression of all possible variants with single back-to-ancestor mutations may not be practical, in particular for large proteins. However, our results support that a limited screening of modern/ancestral chimeras may provide variants with enhanced expression. Furthermore, screening of individual mutations can be focused to protein regions that are expected to be unfolded in aggregation-prone intermediates populated during the folding process. Our results support that such regions can be predicted as the late-folding regions in folding landscape computations. Our version of the Wako-Saitô-Muñoz-Eaton statistical-mechanical model requires only a homology structure model as starting point and, most importantly, employs a block description to drastically reduce the number of microstates in the computation, thus allowing for a fast prediction even for large proteins. It is also worth noting that, for the protein system studied here, convincing molecular explanations

can be put forward for the effect of the S11N and G74S mutations. This suggests the additional possibility of using rational design to determine the specific back-to-ancestor mutations that rescue inefficient heterologous folding.

Overall, our results open up the possibility of rescuing inefficient heterologous expression linked to low solubility by introducing a comparatively small number of back-to-ancestor mutations targeted to specific protein regions. In this way, the sequence would be minimally altered, and, likely, the properties of the encoded protein would be barely altered. This approach could be particularly useful in metagenomics studies that depend on sequence-based screening. Such studies identify enzymes having potential new activities of interest (pollutant or plastic degradation, for instance) based on the sequence similarity with known enzymes. In this scenario, achieving efficient expression with minimal sequence alteration will effectively contribute to the characterization of the new activity in the laboratory.

## Methods

**Expression and purification of thioredoxin variants.** We followed procedures we have previously described in several publications<sup>47–49,60</sup> with small modifications. Briefly, genes encoding *Candidatus Photodesmus katoptron* thioredoxin (*CPk*) and the *CPk* chimeras (*CPk*[1–22] and *CPk*[33–52]) were synthesized with a His-tag at

the C-terminal and codon optimized for expression in *E. coli* cells. Mutations required for the single-mutant *CPk* variants (L7V, D10E, S11N, L14Q, N15E, I17L, S18K, A19S, S20D, G21K, V22P, F70Y, G74S, V75I, S77T), the double-mutant S11N/G74S variant and the chimera involving the loop<sup>70–77</sup> were introduced using the QuikChange Lighting Site-Directed Mutagenesis kit (Agilent Technologies) and the sequences were confirmed by DNA sequencing. Genes were cloned into pET24b(+) plasmid (GenScript Biotech) and transformed into *E. coli* BL21 (DE3) cells (Agilent). Protein expression was induced by 1 mM IPTG and cells were incubated overnight at 25 °C in LB medium and. Cell pellets were sonicated and His-tagged proteins were purified using affinity chromatography (HisGraviTrap column from GE Healthcare).

*E. coli* thioredoxin, LPBCA thioredoxin and the G74S variant of the later used in the experiments reported in this work were prepared following procedures similar to those described above, except that His-tags were not used. Therefore, these proteins were purified<sup>49</sup> by ion-exchange chromatography (Fractogel EMD DEAE column) followed by gel filtration chromatography on HiLoad Superdex 75 column. Our previous studies<sup>47</sup> indicate that the presence of a His-tag has a very small effect on the folding kinetic features of thioredoxins. Also, the presence of a His-tag does not have a significant effect on the efficiency of heterologous expression for *CPk* thioredoxin. The purification procedure based on ion-exchange chromatography and gel filtration was also used to prepare the non-His-tagged *CPk* thioredoxin used for crystallization (see further below).

Folding-unfolding experiments reported in this work were performed with thioredoxin solutions in 50 mM Hepes, pH 7. These solutions were prepared either by dialysis against the buffer at 4 °C or by passage through PD-10 desalting columns (GE Healthcare). Protein concentrations were measured spectrophotometrically using known values for the extinction coefficient. Guanidine and urea solutions in 50 mM HEPES, pH 7 were initially prepared by weight, but their concentrations were subsequently determined from refraction index measurements<sup>78,79</sup> using an Atago R500 hand refractometer. Urea solutions were purified by passage of the stock solution through an AG501-X8(D) ion-exchange resin (Bio-rad) before use.<sup>79</sup>

**Double-jump unfolding assays.** Double-jump unfolding assays, provide an estimate of the amount of native state in a protein solution.<sup>47,63,64</sup> They are based on the fact that the unfolding of the native state is much slower than the unfolding of intermediate (non-native and partially-unfolded) states. Therefore, the amount of native state in a protein solution can be estimated from the amplitude of the native-state unfolding kinetics observed upon transferring an aliquot to denaturing condi-

tions, since intermediate states will unfold in a much shorter time-scale. Double-jump unfolding assays can be used to follow *in vitro* protein refolding kinetics (by performing the assays at several times during refolding), which provides an immediate assessment of the folding time-scale, *i.e.*, the time scale in which the folding polypeptide chain reaches the native state. This is particularly convenient when folding is a complex process involving multi-exponential kinetics and parallel kinetic channels, as is the case with thioredoxins.<sup>61,62</sup> We have recently described and discussed in some detail the use of double-jump unfolding assays to follow thioredoxin refolding kinetics.<sup>47</sup> Briefly, folding is initiated by transferring protein unfolded in a solution with a high denaturant concentration to a solution with a low denaturant concentration (first jump). At different times ( $t_1$ ) after the first jump, aliquots are transferred to denaturing conditions (second jump) and the unfolding kinetics are followed by measuring fluorescence as a function of time ( $t_2$ ). Plots of fluorescence intensity *versus*  $t_2$  time after the second jump are well described by single exponentials, with lifetimes corresponding to the unfolding rate constant at the high denaturant concentration used (Figure S4). The amplitude of the exponential, however, varies reflecting the amount of native protein at the time ( $t_1$ ) at which refolding was interrupted (representative examples are given in Figures S5–S8). After normalization with a suitable control amplitude, the amplitudes lead to a profile of fraction of native protein *versus* refolding time ( $t_1$ ) (Figs. 2, 5 and S9). Specific details of the procedure used are discussed below.

In most cases, proteins were denatured in urea concentration within the range 7.5–9 M and, after fluorescence determinations had indicated that the unfolding process was essentially complete, the folding kinetics was initiated by dilution into typically 1 M urea, although some experiments at other final urea concentrations were performed, as shown in panel a of Figure 2). In most cases, the protein concentration in the folding kinetics experiments (*i.e.* after the first jump) was on the order 0.1 mg/mL, although additional experiments at higher protein concentrations were carried out with *CPk* thioredoxin (see Figure 2b). Typically, the second jumps involved a 1:15 dilution and the unfolding kinetics were determined by following the protein fluorescence at 350 nm as a function of time. The exact composition of the denaturing solution is immaterial for the result of the experiment, as long as the same composition is used in all the unfolding profiles corresponding to given folding experiment. Both, high urea concentrations (within the range 8–9.5 M) and high guanidine concentrations (within the range 3–5 M) were used. Of course, it is important that the denaturant concentration used does indeed unfolds the protein variant studied. In order to select denaturation concentrations that fulfill this

criterion, we determined the unfolding branches of Chevron plots for all the thioredoxins studied here (Figure S4). A few of the thioredoxins studied here are highly stable and they cannot be denatured by urea at 25 °C, not even using the highest urea concentrations experimentally available. In these cases, the initial unfolding step was performed in concentrated guanidine (within the range 3–4.5 M) and the dilution into native conditions was designed to ensure a low guanidine concentration in the folding kinetics experiment (within the range 0.1–0.3 M). Figures S5–S8 show several representative examples of the experiments we have described. Folding kinetic profiles for all the proteins studied here are shown in Figures 2b, 5b and S9.

In all cases, we carried out control experiments in which the native protein (at the same concentration used in the folding kinetic determinations) was transferred to the denaturing solution and the unfolding kinetics was followed by fluorescence at 350 nm. The amount of native protein at each time is then calculated as the ratio of the amplitude of the unfolding kinetics determined from the aliquot extracted at that time to the amplitude of the unfolding kinetics for the control. The resulting profiles of fraction of native state ( $X_N$ ) versus time conformed to a single exponential. Note, however, that the kinetic profiles shown in Figures 2b, 5b and S9 use logarithm of time in the x-axis to highlight differences in folding time scale and that a single  $X_N$  vs. t exponential appears as sigmoidal in a plot of  $X_N$  vs. Int.

Data of fraction of native protein vs. time were fitted with the following equation:  $X_N = X_\infty + (X_0 - X_\infty)\exp(-kt)$ , where  $k$  is the first-order rate constant, and  $X_0$  and  $X_\infty$  are short-time and long-time limiting values of the fraction of native protein. Note that  $X_0$  and  $X_\infty$  do not necessarily equal zero and unity, respectively. Small differences with the control may certainly cause  $X_\infty$  to depart somewhat from unity. More importantly, thioredoxin folding is a complex process involving parallel kinetic channels and intermediate states.<sup>61,62</sup> A value of  $X_0$  significantly higher than zero might reflect that a fraction of the molecules reaches the native state in a shorter time scale than that probed by our experiments (*i.e.*, a fast folding kinetic channel). Likewise, a value of  $X_\infty$  significantly smaller than unity that a fraction of the molecules reaches the native state in a longer time scale than that probed in our experiments (*i.e.*, a slow folding kinetic channel). In practice, however, the values of  $X_0$  and  $X_\infty$  determined from the fitting of the equation to our experimental folding profiles are reasonably close to 0 and 1 in essentially all cases. This implies that our experiments do identify the kinetic phase leading to the native state in the major folding channel. Therefore, we used as a metric of the folding time-scale the life-

time calculated from the rate constant value derived from the fittings (*i.e.*,  $1/k$ ).

**Protein solubility measurements.** Solubility of overexpressed thioredoxins variants in *E. coli* BL21(DE3) strain was checked based on SDS-PAGE, following standard protocols. Briefly, at least 3 independent clones of each thioredoxin variant were grown up to an optical density of 0.6 and induced with 1 mM IPTG for 3 hours at 37 °C. A 90 mL aliquot of the final culture was centrifuged at 4000 rpm, 10 min at 4 °C and the collected pellet was re-suspended in 6 mL of lysis buffer containing 20 mM Tris, pH 7.5, 50 mM NaCl and a protease inhibitor tablet (Roche cOmplete™). After sonication, two aliquots were taken. One aliquot was subjected to SDS-PAGE to estimate the total amount of protein. Other aliquot was centrifuged (15000 rpm for 10 min at 4 °C) and the supernatant was subjected to SDS-PAGE to provide the amount of soluble protein. The SDS-PAGE gels obtained in this work are shown in Figures S10–S13.

Quantification of total and soluble thioredoxin fractions was carried out by SDS-PAGE on 15% Tris-glycine SDS-polyacrylamide gels and using ImageJ software (<https://imagej.nih.gov/ij/>) for image analysis of the thioredoxin bands stained by Coomassie dye. Illustrative densitometry profiles are given in Figure S14. At least, three independent measurements were performed for each protein variant. The average value, the standard deviation and the individual values are given in Figures 3 and 5.

In addition, we attempted to rescue the inefficient heterologous folding of *CPk* thioredoxin by co-overexpressing *E. coli* chaperones. Five plasmids designed to express the following “chaperone teams” were purchased from TAKARA Bio Inc: pG-KJE8 (expressing dnaK-dnaJ-grpE-groES-groEL), pGro7 (expressing groES-groEL), pKJE7 (expressing dnaK-dnaJ-grpE), pG-Tf2 (expressing groES-groEL-tig) and pTf16 (expressing the trigger factor). Chaperone plasmids were transformed into BL21(DE3) chemical competent cells containing the different thioredoxin plasmids.

**Activity measurements.** Activity of thioredoxin proteins was measured using the insulin turbidimetric assay<sup>80</sup> (Holmgren, 1979) as we have previously described.<sup>47</sup> Briefly, in this assay, disulfides reduction by dithiothreitol (DTT) catalysed by thioredoxin causes insulin aggregation, which is followed spectrophotometrically at 650 nm. The reaction mixture contains 0.1 M phosphate buffer pH 6.5, 2 mM EDTA, 0.5 mg/mL of bovine pancreatic insulin and a final thioredoxin concentration of 1.5 μM. The reaction is initiated by addition of DTT to a 1 mM final concentration. Activity values for each variant reported were obtained from the maximum value of plots of  $dA_{650nm}/dt$  versus time. A total of 3 independent measurements were carried out for each thioredoxin variant. The resulting



average values and the corresponding standard deviations are reported in [Figure S1](#).

In addition, for some variants, thioredoxin activity was also assayed with thioredoxin reductase coupled to the reduction of DTNB<sup>81</sup> as we have previously described.<sup>47</sup> Final conditions in the cuvette were: 0.05 M Tris-HCl, 2 mM EDTA pH 8, 0.05 mg/mL BSA, 0.5 mM DTNB, 0.25 mM NADPH and 0.15  $\mu$ M *CPk* variants. Reaction was started by addition of thioredoxin reductase to a final concentration of 0.02  $\mu$ M and monitored spectrophotometrically. A total of three independent measurements were performed for each thioredoxin variant. The resulting average values and the corresponding standard deviations are reported in [Figure S1](#).

**Urea-induced equilibrium denaturation monitored by CD and fluorescence measurements.** The urea-induced equilibrium denaturation of *CPk* thioredoxin, its S11N single mutant and *CPk*[1–22] chimera was studied by using far-UV circular dichroism measurements at 25 °C. Protein concentration was  $\sim$  0.8 mg/mL in a 1 mm cuvette. The urea dependence of ellipticity at 222 nm could be adequately fitted by a two-state model that assumes a linear dependence of the unfolding free energy with denaturant concentration within the transition region and linear pre- and post-transition baselines, as previously described.<sup>47</sup> Values of the midpoint urea concentration ( $C_{1/2}$ ) and the slope of the urea-dependence of the unfolding free energy ( $m$ ) derived from these fits are given in [Figure 7](#).

**Denaturation Temperature measured by Differential Scanning Calorimetry.** The denaturation temperatures of wild *CPk* thioredoxin and single- and double-mutant variants (L7V, D10E, S11N, L14Q, N15E, I17L, S18K, A19S, S20D, G21K, V22P, F70Y, G74S, V75I, S77T, S11N/G74S) and modern/ancestral chimeras (*CPk* [1–22], *CPk* [70–77], *CPk* [1–22]/[70–77], *CPk* [33–52]) were determined by differential scanning calorimetry as the temperature for maximum of the calorimetric transition. Experimental DSC thermograms are shown in [Figure S15](#). Note that, in a few cases, two transitions were reproducibly seen in the thermograms, perhaps reflecting a decreased unfolding cooperativity upon mutation. In these cases, the maximum of the major transition was used as a metric of thermal stability. The experiments were performed with a MicroCal Auto-PEAQ DSC calorimeter (Malvern), at pH 7 in 50 mM HEPES buffer. Typically, protein concentration was within the 0.4–0.6 mg/mL range and scan rate was 240 K/h. Standard protocols well established in our laboratory for thioredoxins were followed.<sup>48,82</sup> In all cases, protein solutions for the calorimetric experiments were prepared by exhaustive dialysis and the buffer from the last dialysis step was used in the reference cell of the calorimeter. Calorimetric cells were kept under excess pressure to prevent degassing during the

scan. Several buffer–buffer baselines were recorded to ensure proper calorimeter equilibration prior to the protein run.

**Statistical-mechanical modeling of the folding landscape.** The WSME model approach, its variants and parameterization are described elsewhere in detail.<sup>57</sup> Briefly, the model is G $\bar{o}$ -like in its energetics requiring a starting structure and coarse-grains the polymer at the residue level. Thus, every residue is assumed to sample two sets of conformations – folded (represented as binary variable 1) and unfolded (0) – contributing to a set of  $2^N$  microstates for a  $N$ -residue protein. In the variant of the model used in the current work, we make two approximations to the original version.<sup>54,56</sup> First, we consider microstates restricted to only single-stretches of folded residues (single sequence approximation, SSA), two-stretches of folded residues (double sequence approximation, DSA) with no interaction across the island, and DSA with interaction across the structured islands if they interact in the native structure and even if the intervening residues are unfolded.<sup>83</sup> Second, we further reduce the accessible phase space by considering 2–3 consecutive residues to fold and unfold together (*i.e.* as a single block or unit). This model has been validated against the residue-level approximations and enables rapid predictions.<sup>53</sup> In addition, we include contributions from van der Waals interactions, all-to-all electrostatics,<sup>57</sup> simplified solvation terms and sequence and structure specific conformational entropy. At the time the simulations were performed, the 3D structure of *CPk* thioredoxin was not available. Therefore, homology modeling was employed to generate a model using the Robetta server<sup>84</sup> and the *E. coli* thioredoxin as a template. Still, the alpha carbon RMSD between the homology model and the experimental structure is 0.524 Å.

All simulations were performed at 37 °C, pH 7.0 and 0.1 M ionic strength conditions with the following parameters: atomic-level interaction energy ( $\xi$ ) of  $-70$  J mol $^{-1}$  for every heavy atom van der Waals interaction identified with a 6 Å spherical cut-off and excluding the nearest neighbours for LPBCA and *CPk* thioredoxin ( $-76$  J mol $^{-1}$  for *E. coli* thioredoxin), change in heat capacity ( $\Delta C_p^{cont}$ ) of  $-0.36$  J mol $^{-1}$  K $^{-1}$  per native contact, entropic penalty for fixing non-proline, non-glycine and residues in well-determined secondary structures ( $\Delta S^{conf}$ ) as  $-16.5$  J mol $^{-1}$  K $^{-1}$  per residue, entropic penalty for glycine and coil residues as  $-22.56$  J mol $^{-1}$  K $^{-1}$  per residue (accounting for the excess disorder in these regions,<sup>85</sup> and 0 J mol $^{-1}$  K $^{-1}$  per residue for proline residues given the limited backbone flexibility of prolines. Free energy profiles and surfaces as a function of the reaction coordinate, the number of structured blocks, are generated by accumulating partial partition functions corresponding to specific number of folded units.

Residue folding probabilities are calculated by summing up the probabilities of states in which the residue of interest is structured.

**Crystallization and structural determination.** Freshly purified CPk[70–77] thioredoxin was concentrated to 25 mg/ml prior setting the crystallization screening. Initial crystallization trials were carried out using the hanging-drop vapor diffusion method. Drops were prepared by mixing 1  $\mu$ L of protein solution with the reservoir in a 1:1 ratio, and equilibrated against 500  $\mu$ L of each precipitant cocktail of the HR-I & PEG/Ion™ crystallization screening (Hampton Research). Crystallization trials were kept at 293 K in an incubator. After one-week crystalline material was observed in conditions #33, #34, #43 and #11 of the HR-I kit and #13 of the PEG/Ion screening kit. Crystals were fished from the drop and transferred to the cryo-protectant solution prepared with the mother liquid supplemented with 15% (v/v) glycerol and subsequently flash-cooled in liquid nitrogen and stored until data collection.

Crystals were diffracted at ID23-1 beam-line of the European Synchrotron Radiation Facility (ESRF), Grenoble, France. The best diffracting crystals were obtained in condition #43 of the HR-I. The diffraction data were indexed and integrated using XDS<sup>86</sup> and scaled with SCALA from the CCP4 suite.<sup>87</sup> Thioredoxin crystals belonged to the  $I2_13$  space group with only two monomers in the asymmetric unit and therefore with a usually high water content, almost 74%, as determined from Matthews' coefficient, 4.71.<sup>88</sup> The molecular replacement solution was found using Molrep<sup>89</sup> and the coordinates of the PDB ID. 2TRX, chain A, locating the two monomers in the asymmetric unit. Refinement was done with phenix.refine<sup>90</sup> including manual building and water inspection with Coot<sup>91</sup> and using the Titration-Libration-Screw (TLS)<sup>92</sup> grouping since the initial steps. Model quality was checked using MolProbity<sup>93</sup> implemented within the Phenix suite.<sup>90</sup> Refinement statistics and quality indicators of the final model are summarized in Table S1. Coordinates and structure factors have been deposited at the PDB with accession code 7OOL.

Hydrogen bond analysis was done using WHAT IF<sup>94</sup> with hydrogen-bond network optimization.<sup>95</sup>

097142-B-100 (J.M.S.-R.) and BIO2016-74875-P (J.A.G.) and the Science, Engineering and Research Board (SERB, India) Grant MTR/2019/000392 (A.N.N.). We are grateful to the European Synchrotron Radiation Facility (ESRF), Grenoble, France, for the provision of time and the staff at ID23-1 beamline for assistance during data collection.

## Author contributions

J.M.S.R. designed the research. G.G.-A. purified the modern/ancestral chimeras and the thioredoxin variants; she also performed and analysed the experiments aimed at determining their folding kinetics and biomolecular properties. V.A.R. performed experiments addressed at determining the efficiency of heterologous expression and provided essential input for the molecular interpretation of mutational effects on expression efficiency. E.A.G. provided essential input for the evolutionary interpretation of the data. J.A.G. determined the X-ray structure of the symbiont protein and provided essential input regarding its interpretation and implications. A.N.N. performed the computational simulations of the folding landscape for thioredoxins and provided essential input regarding their engineering implications. B.I. M. and J.M.S.-R. directed the project. J.M.S.-R. wrote the first draft of the manuscript to which V. A.R., J.A.G., A.N.N. and B.I.M. added crucial paragraphs and sections. All authors discussed the manuscript, suggested modifications and improvements, and contributed to the final version.

## Declaration of Competing Interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## Appendix A. Supplementary material

Supplementary data to this article can be found online at <https://doi.org/10.1016/j.jmb.2021.167321>.

Received 17 August 2021;

Accepted 17 October 2021;

Available online 21 October 2021

## Acknowledgements

This work was supported by Human Frontier Science Program Grant RGP0041/2017 (J.M.S.-R. and E.A.G.), National Science Foundation Award #2032315 (E.A.G.), National Institutes of Health Award #R01AR069137 (E.A.G.), Department of Defense MURI Award #W911NF-16-1-0372 (E.A.G.), Spanish Ministry of Science and Innovation/FEDER Funds Grants RTI-2018-

### Keywords:

ancestral sequence reconstruction;  
computational modelling of protein folding landscapes;  
heterologous protein expression;  
proteins from uncultured organisms;  
obligate symbionts

† These authors contributed equally.

## References

- Walsh, G., (2010). Post-translational modifications of protein biopharmaceuticals. *Drug Discovery Today* **15**, 773–780.
- Baeshen, M.N., Al-Hejin, A.M., Bora, R.S., Ahmed, M.M. M., Ramadan, H.A.I., Saini, K.S., Baeshen, N.A., Redwan, E.M., (2015). Production of biopharmaceuticals in *E. coli*: current scenario and future perspectives. *J. Microbiol. Biotechnol.* **25**, 953–962.
- Tripathi, N.K., Shrivastava, A., (2019). Recent developments in bioprocessing of recombinant proteins: expression hosts and process development. *Front. Bioeng. Biotechnol.* **7**, 420.
- Uchiyama, T., Miyazaki, K., (2009). Functional metagenomics for enzyme discovery: challenges to efficient screening. *Curr. Opin. Biotechnol.* **20**, 616–622.
- Calderon, D., Peña, L., Suarez, A., Villamil, C., Ramirez-Rojas, A., Anzola, J.M., Garcia-Betancur, J.C., Cepeda, M. L., Uribe, D., Del Portillo, P., Mongui, A., (2019). Recovery and functional validation of hidden soil enzymes in metagenomic libraries. *MicrobiolOpen* **8**, e572
- Daniel, R., (2005). The metagenomics of soil. *Nat. Rev. Microbiol.* **3**, 470–478.
- Baneyx, F., Mujacic, M., (2004). Recombinant protein folding and misfolding in *Escherichia coli*. *Nat. Biotechnol.* **22**, 1399–1408.
- Selas Castiñeiras, T., Williams, S.G., Hitchcock, A.G., Smith, D.S., (2018). *E. coli* strain engineering for the production of advanced biopharmaceutical products. *FEMS Microbiol. Lett.* **365**, 15.
- I. Acebrón, L. Plaza-Vinuesa, B. de las Rivas, R. Muñoz, J. Cumella, F. Sánchez-Sancho, J.M. Mancheño, Structural basis of the substrate specificity and instability in solution of a glycosidase from *Lactobacillus plantarum*. *Biochim. Biophys. Acta Proteins Proteom.* **1865** (2017) 1227–1236.
- Pauling, L., Zuckerkandl, E., (1963). Chemical paleogenetics. Molecular “restoration studies” of extinct forms of life. *Acta Chem. Scan.* **17S**, 9–16.
- Liberles, D.A., (2007). Ancestral Sequence Reconstruction. Oxford University Press, Oxford.
- Benner, S.A., Sassi, S.O., Gaucher, E.A., (2007). Molecular paleoscience: systems biology from the past. *Adv. Enzymol. Relat. Areas Mol. Biol.* **75**, 1–132.
- Hochberg, G.K.A., Thornton, J.W., (2017). Reconstructing ancient proteins to understand the causes of structure and function. *Annu. Rev. Biophys.* **46**, 247–269.
- Gumulya, Y., Gillam, E.M., (2017). Exploring the past and the future of protein evolution with ancestral sequence reconstruction: the “retro” approach to protein engineering. *Biochem. J.* **474**, 1–19.
- S.D. Copley, Setting the stage for evolution of a new enzyme. *Curr. Opin. Struct. Biol.* **69** (2021) 41–49.
- Selberg, A.G.A., Gaucher, E.A., Liberles, D.A., (2021). Ancestral sequence reconstruction: from chemical paleogenetics to maximum likelihood algorithms and beyond. *J. Mol. Evol.* **89**, 157–164.
- Cole, M.F., Gaucher, E.A., (2011). Exploring models of molecular evolution to efficiently direct protein engineering. *J. Mol. Evol.* **72**, 193–203.
- Risso, V.A., Gavira, J.A., Mejia-Carmona, D.F., Gaucher, E.A., Sanchez-Ruiz, J.M., (2013). Hyperstability and substrate promiscuity in laboratory resurrections of Precambrian  $\beta$ -lactamases. *J. Am. Chem. Soc.* **135**, 2899–2902.
- Risso, V.A., Sanchez-Ruiz, J.M., Ozkan, S.B., (2018). Biotechnological and protein engineering implications of ancestral protein resurrection. *Curr. Opin. Struct. Biol.* **51**, 106–115.
- Trudeau, D.L., Tawfik, D.S., (2019). Protein engineers turned evolutionist – the quest for the optimal starting point. *Curr. Opin. Struct. Biol.* **60**, 46–52.
- M.A. Spence, J.A. Kaczmarek, J.W. Saunders, C.J. Jackson, Ancestral sequence reconstruction for protein engineers. *Curr. Opin. Struct. Biol.* **69** (2021) 131–141.
- Gonzalez, D., Hiblot, J., Darbinian, N., Miller, J.C., Gotthard, G., Shohreh, A., Chabriere, E., Elias, M., (2014). Ancestral mutations as a tool for solubilizing proteins: the case of a hydrophobic phosphate-binding protein. *FEBS Open Bio* **4**, 121–127.
- Withfield, J.H., Zhang, W.H., Herde, M.K., Clifton, B.E., Radziejewski, J., Janovjak, H., Henneberger, C., Jackson, C.J., (2015). Construction of a robust and sensitive arginine biosensor through ancestral protein reconstruction. *Protein Sci.* **24**, 1412–1422.
- Trudeau, D.L., Kaltenbach, M., Tawfik, D.S., (2016). On the potential origins of the high stability of reconstructed ancestral proteins. *Mol. Biol. Evol.* **33**, 2633–2641.
- Zakas, P.M., Brown, H.C., Knight, K., Meeks, S.L., Gaucher, E.A., Doering, C.B., (2017). Enhancing the pharmaceutical properties of protein drugs by ancestral sequence reconstruction. *Nat. Biotechnol.* **35**, 35–37.
- Manteca, A., Schönfelder, J., Alonso-Caballero, A., Fertin, M.J., Barrietabeña, N., Faria, B.F., Herrero-Galán, E., Alegre-Cebollada, J., De Sancho, D., Perez-Jimenez, R., (2017). Mechanochemical evolution of the giant muscle protein titin as inferred from resurrected proteins. *Nat. Struct. Mol. Biol.* **24**, 652–657.
- Babkova, P., Sebestova, E., Brezovsky, J., Chaloupkova, R., Damborsky, J., (2017). Ancestral haloalkane dehalogenases show robustness and unique substrate specificity. *ChemBioChem* **18**, 1448–1456.
- Koblan, L.W., Doman, J.L., Wilson, C., Levy, J.M., Tay, T., Newby, G.A., Maianti, J.P., Raguram, A., Liu, D.R., (2018). Improving cytidine and adenine base editors by expression optimization and ancestral reconstruction. *Nat. Biotechnol.* **36**, 843–846.
- Hendrikse, N.M., Charpentier, G., Nordling, E., Syrén, P.-O., (2018). Ancestral diterpene cyclases show increased thermostability and substrate acceptance. *FEBS J.* **285**, 4660–4673.
- Gomez-Fernandez, B.J., Garcia-Ruiz, E., Martin-Diaz, J., Gomez de Santos, P., Santos-Moriano, P., Plou, F.J., Ballesteros, A., Garcia, M., Rodriguez, M., Risso, V.A., Sanchez-Ruiz, J.M., Whitney, S.M., Alcalde, M., (2018). Directed *in vitro* evolution of Precambrian and extant Rubiscos. *Sci. Rep.* **8**, 5532.
- Barrietabeña, N., Alonso-Lerma, B., Galera-Prat, A., Joudeh, N., Barandiaran, L., Aldazabal, L., Arbulu, M., Alvalde, M., De Sancho, D., Gavira, J.A., Carrion-Vazquez, M., Perez-Jimenez, R., (2019). Resurrection of efficient Precambrian endoglucanases for lignocellulosic biomass hydrolysis. *Commun. Chem.* **2**, 76.
- Nakano, S., Minamino, Y., Hasebe, F., Ito, S., (2019). Deracemization and stereoinversion to aromatic D-amino acid derivatives with ancestral L-amino acid oxidase. *ACS Catal.* **9**, 10152–10158.



33. Gomez-Fernandez, B.J., Risco, V.A., Rueda, A., Sanchez-Ruiz, J.M., Alcalde, M., (2020). Ancestral resurrection and directed evolution of fungal mesozoic laccases. *Appl. Environ. Microbiol.* **86**, e0078–e120.
34. Li, D., Damry, A.M., Petrie, J.R., Vanhercke, T., Singh, S. P., Jackson, C.J., (2020). Consensus mutagenesis and ancestral reconstruction provide insight into the substrate specificity and evolution of the from-end  $\Delta 6$ -desaturase family. *Biochemistry* **59**, 1398–1409.
35. Sun, Y., Calderini, E., Kourist, R., (2021). A reconstructed common ancestor of the fatty acid photo-decarboxylase clade shows photo-decarboxylation activity and increased thermostability. *ChemBioChem* **22**, 1833–1840.
36. Nicoll, C.R., Bailleul, G., Fiorentini, F., Mascotti, M.L., Fraaije, M.W., Matevi, A., (2020). Ancestral sequence reconstruction unveils the structural basis of function in mammalian FMOs. *Nat. Struct. Mol. Biol.* **27**, 14–24.
37. Schriever, K., Saez-Mendez, P., Rudraraju, R.S., Hendrikse, N.M., Hudson, E.P., Biundo, A., Schnell, R., Syrén, P.O., (2021). Engineering of ancestors as a tool to elucidate structure, mechanism, and specificity of extant terpene cyclase. *J. Am. Chem. Soc.* **143**, 3794–3807.
38. Ufarté, L., Laville, É., Duquesne, S., Potocki-Veronese, G., (2015). Metagenomics for the discovery of pollutant degrading enzymes. *Biotechnol. Adv.* **33**, 1845–1854.
39. Moran, N.A., (1996). Accelerated evolution and Muller's ratchet in endosymbiotic bacteria. *Proc. Natl. Acad. Sci. U. S.A.* **93**, 2873–2878.
40. Woolfit, M., Bromhan, L., (2003). Increased rates of sequence evolution in endosymbiotic bacteria and fungi with small effective population sizes. *Mol. Biol. Evol.* **20**, 1545–1555.
41. Holmgren, A., (1985). Thioredoxin. *Annu. Rev. Biochem.* **54**, 237–271.
42. Morin, J.G., Harrington, A., Neelson, K., Krieger, N., Baldwin, T.O., Hastings, J.W., (1975). Light for all reasons: versatility in the behavioural repertoire of the flashlight fish. *Science* **190**, 74–76.
43. Hendry, T.A., de Wet, J.R., Dougan, K.E., Dunlap, P.V., (2016). Genome evolution of the obligate but environmentally active luminous symbionts of flashlight fish. *Genome Biol. Evol.* **8**, 2203–2213.
44. Hendry, T.A., Dunlap, P.V., (2011). The uncultured luminous symbiont of *Anomalops katoptron* (Beryciformes: Anomalopidae) represents a new bacterial genus. *Mol. Phylogenet. Evol.* **61**, 834–843.
45. Hendry, T.A., Dunlap, P.V., (2014). Phylogenetic divergence between the obligate luminous symbionts of flashlight fish demonstrates specificity of bacteria to host genera. *Environ. Microbiol. Rep.* **6**, 331–338.
46. Hendry, T.A., de Wet, J.R., Dunlap, P.V., (2014). Genomic signatures of obligate host dependence in the luminous bacterial symbiont of a vertebrate. *Environ. Microbiol.* **16**, 2611–2622.
47. G. Gamiz-Arco, V.A. Risco, A.M. Candel, A. Inglés-Prieto, M.L. Romero-Romero, E.A. Gaucher, J.A. Gavira, B. Ibarra-Molero, J.M. Sanchez-Ruiz, Non-conservation of folding rates in the thioredoxin family reveals degradation of ancestral unassisted folding. *Biochem. J.* **476** (2019) 3631–3647.
48. Perez-Jimenez, R., Ingles-Prieto, A., Zhao, Z.M., Sanchez-Romero, I., Alegre-Cebollada, J., Kosuri, P., Garcia-Manyes, S., Kappock, T.J., Tanokura, M., Holmgren, A., Sanchez-Ruiz, J.M., Gaucher, E.A., Fernandez, J.M., (2011). Single-molecule paleoenzymology probes the chemistry of resurrected enzymes. *Nat. Struct. Mol. Biol.* **18**, 592–596.
49. Ingles-Prieto, A., Ibarra-Molero, B., Delgado-Delgado, A., Perez-Jimenez, R., Fernandez, J.M., Gaucher, E.A., Sanchez-Ruiz, J.M., Gavira, J.A., (2013). Conservation of protein structure over four billion years. *Structure* **21**, 1690–1697.
50. V.A. Risco, F. Manssour-Triedo, A. Delgado-Delgado, R. Arco, Barroso-delJesus, A. Ingles-Prieto, R. Godoy-Ruiz, J.A. Gavira, E.A. Gaucher, B. Ibarra-Molero, J.M. Sanchez-Ruiz, Mutational studies on resurrected ancestral proteins reveal conservation of site-specific amino acid preferences throughout evolutionary history. *Mol. Biol. Evol.* **32** (2015) 440–455
51. Candel, A.M., Romero-Romero, M.L., Gamiz-Arco, G., Ibarra-Molero, B., Sanchez-Ruiz, J.M., (2017). Fast folding and slow unfolding of a resurrected Precambrian protein. *Proc. Natl. Acad. Sci. U.S.A.* **114**, E4122–E4123.
52. Rosano, G.L., Ceccarelli, E., (2014). Recombinant protein expression in *Escherichia coli*: advances and challenges. *Front. Microbiol.* **5**, 172.
53. Gopi, S., Aranganathan, A., Naganathan, A.N., (2019). Thermodynamics and folding landscapes of large proteins from a statistical mechanical model. *Curr. Res. Struct. Biol.* **1**, 6–12.
54. Wako, H., Saitô, N., (1978). Statistical mechanical theory of the protein conformation. 2. Folding pathway for protein. *J. Phys. Soc. Jpn.* **44**, 1939–1945.
55. Muñoz, V., Thompson, P.A., Hofrichter, J., Eaton, W.A., (1997). Folding dynamics and mechanism of  $\beta$ -hairpin formation. *Nature* **390**, 196–199.
56. Muñoz, V., Eaton, W.A., (1999). A simple model for calculating the kinetics of protein folding from three-dimensional structures. *Proc. Natl. Acad. Sci. U.S.A.* **96**, 11311–11316.
57. Naganathan, A.N., (2012). Predictions from an Ising-like statistical mechanical model on the dynamic and thermodynamic effects of protein surface electrostatics. *J. Chem. Theory Comput.* **8**, 4646–4656.
58. Balchin, D., Hayer-Hartl, M., Hartl, F.U., (2016). In vivo aspects of protein folding and quality control. *Science* **353**, aac4354.
59. Katti, S.K., LeMaster, D.M., Eklund, H., (1990). Crystal structure of thioredoxin from *Escherichia coli* at 1.68 Å resolution. *J. Mol. Biol.* **212**, 167–184.
60. Romero-Romero, M.L., Risco, V.A., Martinez-Rodriguez, S., Ibarra-Molero, B., Sanchez-Ruiz, J.M., (2016). Engineering ancestral protein hyperstability. *Biochem. J.* **473**, 3611–3620.
61. Kelley, R.F., Richards, F.M., (1987). Replacement of proline-76 with alanine eliminates the slowest kinetic phase in thioredoxin folding. *Biochemistry* **26**, 6765–6774.
62. Georgescu, R.E., Li, J.H., Goldberg, M.E., Tasayco, M.L., Chaffotte, A.F., (1998). Proline isomerization-independent accumulation of an early intermediate and heterogeneity of the folding pathway of mixed  $\alpha/\beta$  protein, *Escherichia coli* thioredoxin. *Biochemistry* **37**, 10286–10297.
63. Mücke, M., Schmidt, F.X., (1994). A kinetic method to evaluate the two-state character of solvent-induced protein denaturation. *Biochemistry* **33**, 12930–12935.
64. Ibarra-Molero, B., Sanchez-Ruiz, J.M., (1997). Are there equilibrium intermediates in the urea-induced unfolding of hen-egg-white lysozyme. *Biochemistry* **36**, 9616–9624.



65. Brandts, J.F., Halvorson, H.R., Brennan, M., (1975). Consideration of the possibility that the slow step on protein denaturation reactions is due to cis-trans isomerism of proline residues. *Biochemistry* **14**, 4953–4963.
66. Schmid, F.X., Baldwin, R.L., (1978). Acid catalysis of the formation of the slow-folding species of RNase A: evidence that the reaction is proline isomerization. *Proc. Natl. Acad. Sci. U.S.A.* **75**, 4764–4768.
67. Godoy-Ruiz, R., Ariza, F., Rodriguez-Larrea, D., Perez-Jimenez, R., Ibarra-Molero, B., Sanchez-Ruiz, J.M., (2006). Natural selection for kinetic stability is a likely origin of correlations between mutational effects on protein energetics and frequencies of amino acid occurrences in sequence alignments. *J. Mol. Biol.* **362**, 966–978.
68. Romero-Romero, M.L., Rizzo, V.A., Martinez-Rodriguez, S., Gaucher, E.A., Ibarra-Molero, B., Sanchez-Ruiz, J.M., (2016). Selection for protein kinetic stability connects denaturation temperatures to organismal temperatures and provides clues to Archaeal life. *PLoS ONE* **11**, e0156657.
69. Kaiser, C.M., Goldman, D.H., Chodera, J.D., Tinoco, I., Bustamante, C., (2011). The ribosome modulates nascent protein folding. *Science* **334**, 1723–1727.
70. Samelson, A.J., Bolin, E., Costello, S.M., Sharma, A.K., O'Brien, E.P., Marqusee, S., (2018). Kinetic and structural comparison of a protein's cotranslational folding and refolding pathways. *Sci. Adv.* **4**, eaas9098.
71. Roodveldt, C., Aharoni, A., Tawfik, D.S., (2005). Directed evolution of proteins for heterologous expression and stability. *Curr. Opin. Struct. Biol.* **15**, 50–56.
72. Randall, R.N., Radford, C.E., Roof, K.A., Natarajan, D.K., Gaucher, E.A., (2016). An experimental phylogeny to benchmark ancestral sequence reconstruction. *Nat. Commun.* **7**, 12847.
73. Broom, A., Doxey, A.C., Lobsanov, Y.D., Berthin, L.G., Rose, D.R., Howell, P.L., McConkey, B.J., Meiering, E.M., (2012). Modular evolution and the origins of symmetry: reconstruction of a three-fold symmetric globular protein. *Structure* **20**, 161–171.
74. Broom, A., Ma, M., Xia, K., Rafalia, H., Trainor, K., Colón, W., Gosavi, S., Meiering, E.M., (2015). Designed protein reveals structural determinants of extreme kinetic stability. *Proc. Natl. Acad. Sci. U.S.A.* **112**, 14605–14610.
75. Jaswal, S.S., Sohl, J.L., Davis, J.H., Agard, D.A., (2012). Energetic landscape of  $\alpha$ -lytic protease optimizes longevity through kinetic stability. *Nature* **415**, 343–346.
76. Goldenzweig, A., Goldsmith, M., Hill, S.E., Gertman, O., Laurino, P., Ashani, Y., Dym, O., Unger, T., Albeck, S., Prilusky, J., Lieberman, R.L., Aharoni, A., Silman, I., Sussman, J.L., Tawfik, D.S., Fleishman, S.J., (2016). Automated structure- and sequence-based design of proteins for high bacterial expression and stability. *Mol. Cell* **63**, 337–346.
77. Marcelino, A.M., Gierasch, L.M., (2008). Roles of  $\beta$ -turns in protein folding: from peptide models to protein engineering. *Biopolymers* **89**, 380–391.
78. Garcia-Mira, M.M., Sanchez-Ruiz, J.M., (2001). pH corrections and protein ionization in water/guanidinium chloride. *Biophys. J.* **81**, 3489–3502.
79. Acevedo, O., Guzman-Casado, M., Garcia-Mira, M.M., Ibarra-Molero, B., Sanchez-Ruiz, J.M., (2002). pH corrections in chemical denaturant solutions. *Anal. Biochem.* **306**, 158–161.
80. Holmgren, A., (1979). Thioredoxin catalyzes the reduction of insulin disulfides by dithiothreitol and dihydrolipoamide. *J. Biol. Chem.* **254**, 9627–9632.
81. Slaby, I., Holmgren, A., (1975). Reconstitution of *E. coli* thioredoxin from complementing peptide fragments obtained by cleavage at methionine-37 or arginine-73. *J. Biol. Chem.* **250**, 1340–1347.
82. Georgescu, R.E., Garcia-Mira, M.M., Tasayco, M.L., Sanchez-Ruiz, J.M., (2001). Heat capacity analysis of oxidized *Escherichia coli* thioredoxin fragments (1–73, 74–108) and their noncovalent complex. Evidence for the burial of apolar surface in protein unfolded states. *Eur. J. Biochem.* **268**, 1477–1485.
83. Kubelka, J., Henry, E.R., Cellmer, T., Hofrichter, J., Eaton, W.A., (2008). Chemical, physical, and theoretical kinetics of an ultrafast folding protein. *Proc. Natl. Acad. Sci. U.S.A.* **105**, 18655–18662.
84. Song, Y., DiMaio, F., Wang, R.Y.-R., Kim, D., Miles, C., Brunette, T., Thompson, J., Baker, D., (2013). High-resolution comparative modelling with RosettaCM. *Structure* **21**, 1735–1742.
85. Rajasekaran, N., Gopi, S., Narayan, A., Naganathan, A.N., (2016). Quantifying protein disorder through measures of excess conformational entropy. *J. Phys. Chem. B* **120**, 4341–4350.
86. W. Kabsch, XDS. *Acta Crystallogr. Sect. D Biol. Crystallogr.* **66** (2010) 125–132.
87. Winn, M.D., Ballard, C.C., Cowtan, K.D., Dodson, E.J., Emsley, P., Evans, P.R., Keegan, R.M., Krissinel, E.B., Leslie, A.G.W., McCoy, A., McNicholas, S.J., Murshudov, G.N., Pannu, N.S., Potterton, E.A., Powell, H.R., Read, R. J., Vagin, A., Wilson, K.S., (2011). Overview of the CCP4 suite and current developments. *Acta Crystallogr. Sect. D Biol. Crystallogr.* **67**, 235–242.
88. Kantardjieff, K.A., Rupp, B., (2003). Matthews coefficient probabilities: Improved estimates for unit cell contents of proteins, DNA, and protein-nucleic acid complex crystals. *Protein Sci.* **12**, 1865–1871.
89. Vagin, A., Teplyakov, A., (2010). Molecular replacement with MOLREP. *Acta Crystallogr. Sect. D Biol. Crystallogr.* **66**, 22–25.
90. Adams, P.D., Afonine, P.V., Bunkóczi, G., Chen, V.B., Davis, I.W., Echols, N., Headd, J.J., Hung, L.W., Kapral, G. J., Grosse-Kunstleve, R.W., McCoy, A.J., Moriarty, N.W., Oeffner, R., Read, R.J., Richardson, D.C., Richardson, J. S., Terwilliger, T.C., Zwart, P.H., (2010). PHENIX: A comprehensive Python-based system for macromolecular structure solution. *Acta Crystallogr. Sect. D Biol. Crystallogr.* **66**, 213–221.
91. Emsley, P., Lohkamp, B., Scott, W.G., Cowtan, K., (2010). Features and development of Coot. *Acta Crystallogr. Sect. D Biol. Crystallogr.* **66**, 486–501.
92. Painter, J., Merritt, E.A., (2006). TLSMD web server for the generation of multi-group TLS models. *J. Appl. Crystallogr.* **39**, 109–111.

93. Chen, V.B., Arendall, W.B., Headd, J.J., Keedy, D.A., Immormino, R.M., Kapral, G.J., Murray, L.W., Richardson, J.S., Richardson, D.C., (2010). MolProbity: All-atom structure validation for macromolecular crystallography. *Acta Crystallogr. Sect. D Biol. Crystallogr.* **66**, 12–21.
94. Vriend, G., (1990). WHAT IF: A molecular modelling and drug design program. *J. Mol. Graph.* **8**, 52–56.
95. Hooft, R.W.W., Sander, C., Vriend, G., (1996). Positioning hydrogen atoms by optimizing hydrogen-bond networks in protein structures. *Proteins* **26**, 363–376.