

# **Structure and functional implications of WYL-domain-containing transcription factor PafBC involved in the mycobacterial DNA damage response**

Andreas U. Müller<sup>1</sup>, Marc Leibundgut<sup>1</sup>, Nenad Ban<sup>1</sup>, and Eilika Weber-Ban<sup>1\*</sup>

<sup>1</sup>ETH Zurich, Institute of Molecular Biology & Biophysics, CH-8093 Zurich, Switzerland

\*corresponding author

Keywords: PafBC, transcription factor, DNA-damage response pathway, WYL domain, Sm-fold, proteasome, mycobacteria

## Abstract

In mycobacteria, transcriptional activator PafBC is responsible for upregulating the majority of genes induced by DNA damage. Understanding the mechanism of PafBC activation is impeded by a lack of structural information on this transcription factor that contains a widespread, but poorly understood WYL domain frequently encountered in bacterial transcription factors. Here, we determined the crystal structure of *Arthrobacter aurescens* PafBC. The protein consists of two modules, each harboring an N-terminal helix-turn-helix DNA binding domain followed by a central WYL and a C-terminal extension (WCX) domain. The WYL domains exhibit Sm-folds, while the WCX domains adopt ferredoxin-like folds, both characteristic for RNA binding proteins. Our results suggest a mechanism of regulation in which WYL domain-containing transcription factors may be activated by binding RNA molecules. Using an *in vivo* mutational screen in *Mycobacterium smegmatis*, we identify potential co-activator binding sites on PafBC.

## Introduction

DNA damage represents a threat to the integrity of genetic information and is therefore counteracted in all organisms by an arsenal of DNA repair processes that are activated by specific DNA damage response pathways. Mycobacteria and many other actinobacteria employ two distinct yet interconnected pathways in order to upregulate the expression of specific sets of genes required for repair and survival of DNA damage.

The “SOS response”, the canonical pathway described in most bacterial species, relies on cleavage and removal of LexA, a transcriptional repressor of DNA repair genes (“SOS genes”) (reviewed in (Kreuzer, 2013; Maslowska et al., 2018)). Under normal conditions, LexA ensures low expression levels of the SOS genes by binding to a promoter motif called “SOS box” (Little et al., 1981). Single-stranded DNA (ssDNA) fragments accumulating under DNA damage conditions serve as DNA stress signal for the SOS response and are sensed by the ATPase RecA. RecA forms a filamentous complex with ssDNA that is able to induce autoproteolytic cleavage of the LexA repressor, leading to derepression of the SOS genes (Galletto et al., 2006; Little et al., 1980; Phizicky and Roberts, 1981). In *Mycobacterium tuberculosis* (Mtb), the LexA repressor controls about 20 genes (Davis et al., 2002a; Smollett et al., 2012).

In contrast, the second pathway regulates over 150 genes, including many of the LexA-controlled genes, like for example *recA*. This predominant pathway operates independently of LexA and RecA, as demonstrated by deletion of *recA* in Mtb, which leaves upregulation of most DNA repair genes intact (Davis et al., 2002b; Rand et al., 2003). Different from the regulatory principle of derepression, these genes are regulated by transcriptional activation by the heterodimeric protein complex PafBC (Fudrini Olivencia et al., 2017; Müller et al., 2018). The complex consists of the close sequence homologs PafB and PafC (proteasome accessory factors B and C) that are encoded together in an operon that is tightly associated with the bacterial proteasome gene locus, suggesting a functional connection. Indeed, many DNA repair proteins are removed by proteasomal degradation after the DNA damage has been repaired, thereby helping to shut down the stress response and preventing negative impact of DNA-modifying activities under normal conditions (Müller et al., 2018).

PafBC activates its target genes via a promoter motif called RecA-NDp (RecA-independent promoter), which was demonstrated by *in vivo* identification of PafBC binding sites using cell culture

cross-linking followed by immunoprecipitation of PafBC-DNA complexes (Müller et al., 2018). However, PafBC protein levels are not changing in response to DNA stress (Fudrini Olivencia et al., 2017). Furthermore, specific interaction between PafBC and the identified DNA-target regions could not be reconstituted *in vitro*. Taken together, these results suggest that an additional “response-producing” event must take place to initiate PafBC transcription activation.

In order to establish the mechanistic principles employed by PafBC to activate transcription at a molecular level, understanding of the structural framework is crucial. Based on sequence similarity, PafBC belongs to a family of bacterial regulators featuring a winged helix-turn-helix (HTH) domain at the N-terminus, followed by a C-terminal domain of unknown function named WYL domain after a consecutive W-Y-L sequence motif. It has been suggested that the WYL domain might play the role of a ligand-binding domain in the context of this class of transcription factors. A handful of other WYL domain-containing proteins were studied to date: (1) DriD, an SOS response-independent transcriptional activator of a cell division inhibitor protein in *Caulobacter crescentus* (Modell et al., 2014) (2) SII7009, SII7062, SII7078, transcriptional repressors of CRISPR/Cas system mature crRNA in *Synechocystis* 6803 (Hein et al., 2013), (3) PIF1 helicase from *Thermotoga elfii* (Andis et al., 2018) and (4) WYL domain-containing proteins stimulating RNA cleavage by Cas13d in *Eubacterium siraeum* and *Ruminococcus* sp. (Yan et al., 2018). However, structural information on WYL domain-containing transcriptional regulators is missing, and evidence as to how they exert their functions mechanistically has remained elusive.

In this study, we determine the crystal structure of PafBC from *Arthrobacter aurescens* in its non-activated, DNA-free state. The structure reveals that the WYL domain exhibits an Sm-fold, commonly encountered in RNA-binding proteins, and is followed by an additional C-terminal extension (WCX) domain featuring a ferredoxin-like fold. Based on the structure of the PafBC WYL-domain, we carry out a comprehensive computational analysis of WYL domain-containing proteins, and demonstrate that the WYL domain is a widespread feature of bacterial transcription factors present in almost all bacterial taxa. Our study shows that Sm-fold proteins are a much more frequent occurrence in bacteria than previously thought. Based on the high structural similarity of the WYL motif-containing domain to the bacterial RNA chaperone Hfq and the known binding sites of Hfq, we identify functionally essential residues in the WYL domain of PafBC, which are likely involved in binding of a response-producing ligand in this distinct class of transcriptional regulators.

## Results

### The crystal structure of PafBC in the non-activated state exhibits an asymmetric domain arrangement incompatible with DNA binding

In order to obtain information about the architecture of the PafBC class of transcriptional regulators, we set out to determine the crystal structure of PafBC. We carried out crystallization experiments using a range of PafBC orthologs from different actinobacterial organisms, also including PafBC proteins from organisms encoding a naturally fused PafBC complex (i.e. from *Kocuria rhizophila*, *Thermobifida fusca*, and *Arthrobacter aurescens*). Three-dimensional crystals suitable for data collection were ultimately obtained using an *A. aurescens* PafBC construct (<sup>Aau</sup>PafBC $\Delta$ NC), which was shortened by 17 amino acids at the N-terminus and 7 amino acids at the C-terminus based on

sequence alignment, since these residues are not conserved amongst the orthologs and not even present in most of them. Indeed, it is likely that the start site for the *A. aurescens* protein was misassigned, since a valine (encoded by GTG) is present at the position where the other PafBC proteins feature the conserved initiator methionine, and GTG is a frequent start codon in actinobacteria (Belinky et al., 2017; DeJesus et al., 2013) (Figure S1). Selenomethionine-labelled protein was used for crystallization and the structure was determined *de novo* by single-wavelength anomalous diffraction. <sup>Aau</sup>PafBC $\Delta$ NC crystallized in space group  $P2_12_12_1$  with two molecules in the asymmetric unit, and the structure was refined to 2.2 Å with an  $R_{\text{work}}/R_{\text{free}}$  of 20%/24% (Table 1). The structural model was built to near completion, encompassing 1279 residues out of 1328. The missing residues are all located in three poorly ordered loop regions. Although <sup>Aau</sup>PafBC is a natural fusion protein and thus consists of a single polypeptide chain, for simplicity and easier comparison to the majority of actinobacteria encoding separate PafB/PafC proteins we will refer to the PafB- and PafC-corresponding parts as PafB or PafC, respectively.

Although most transcriptional regulators with an HTH domain feature an internal symmetry axis when bound to DNA (Jones et al., 1999), PafBC in our structure adopts an asymmetric conformation (Figure 1, S2 and S3). The asymmetric arrangement of PafBC is not surprising, since PafBC is not in complex with its consensus DNA binding site and its domain arrangement reflects the non-activated form of PafBC. This is in agreement with the observation that specific interaction between PafBC and the identified DNA-target regions takes place *in vivo*, but could not be reconstituted *in vitro* (Fudrini Olivencia et al., 2017; Müller et al., 2018), further supporting the notion that PafBC is in a non-activated state in absence of the putative DNA-stress sensing ligand.

The PafBC structure features six distinct domains, three in each of the homologous PafB and PafC modules. Each module includes an N-terminal HTH domain followed by two other domains (Figure 1 and Figure S2). To distinguish between the homologous domains of each module, we refer to the domains belonging to the PafB module as helix-turn-helix (HTH-B), WYL (WYL-B) and C-terminal extension of the WYL (WCX-B) domains, while in PafC the corresponding domains are termed HTH-C, WYL-C and WCX-C, respectively. All individual domains are connected by long loops, which are particularly pronounced between the WYL and WCX domains, suggesting a high degree of flexibility and conformational adaptability of the PafBC complex.

Both HTH domains adopt a classical winged helix-turn-helix fold with a three-helix bundle consisting of helices H1, H2 and H3 followed by a two-stranded wing, a topology that is typical for many transcriptional regulators (Figure 2 and S4a). From structures of other winged HTH domains in complex with DNA, it is known that H3 usually mediates the specific DNA recognition by binding into the major groove, while the wing provides additional contacts in the minor groove (Clark et al., 1993; Rajagopalan et al., 2013; Wisedchaisri et al., 2007). The putative recognition helix of PafB features two highly conserved phenylalanines (F42 and F46) involved in forming the hydrophobic core of the three-helix bundle (Figure S4b). From the opposing surface-accessible side of the recognition helix, two strictly conserved arginine side chains project outwards (R44 and R48). These residues might be involved in sequence-specific DNA binding with guanine and cytosine bases. Notably, the putative recognition helix of PafC is shorter and comprises a different amino acid sequence that would suggest recognition of non-palindromic DNA sequence, which is in accordance with the non-palindromic nature of the PafBC binding motif (RecA-NDp) (Gamulin et al., 2004; Müller et al., 2018).

However, the most striking difference between the HTH domains in PafB and PafC is their location and accessibility in the solved complex structure.

The HTH domain of PafB (HTH-B) is accessible in the structure, since only helix H1 and a small part of H3 make hydrophobic contacts to the rest of the molecule (Figure S4b). In contrast, the HTH domain of PafC (HTH-C) interacts extensively with the other domains of the protein (Figure S2b and S4c-f). Importantly, the putative recognition helix (H3) appears wedged into the protein core, and following H3, a long loop extends around the back of the central helix contacting the central  $\alpha$ -helix of PafB ( $\alpha$ 4) and the WYL-B domain (Figure 2b and S4c). Superposition of the HTH-B and HTH-C domains reveals a good agreement of the folds except for a much longer  $\beta$ 1/ $\beta$ 2-loop (wing) in HTH-B, while HTH-C displays a very short connection between  $\beta$ 1 and  $\beta$ 2 and instead features a long loop between H3 and  $\beta$ 1 (Figure S4a). Sequence alignment-assisted comparison of the structural elements reveals that the  $\beta$ 2 strands of both HTH domains occupy the same position relative to the helices, while the  $\beta$ 1 position differs. This opens up the possibility that the wing of HTH-C has undergone a register shift to accommodate HTH-C in the protein core of the non-activated PafBC complex (Figure S5).

In all PafB and PafC homologs, the HTH domains are connected to the WYL domains via a sequence stretch of roughly 30 residues. This region forms a long single helix in the PafB module of our structure. The helix is located in the core of the PafBC structure, traversing the entire complex (Figure 1 and S3). In contrast, the equivalent region in PafC consists of three smaller helices forming a bundle between HTH-C and WYL-C ( $\alpha$ 4',  $\alpha$ 4'',  $\alpha$ 4''', Figure S2). The long central helix of PafB packs against helix H1 of the HTH domain in PafC (Figure S4c). The WYL domains of PafB and PafC are located at the perimeter on opposite sides of the complex.

### **The PafBC WYL domains exhibit an Sm-fold and are followed by a ferredoxin-like domain**

The structure shows that the domain previously referred to as the “WYL domain” of PafB or PafC in fact consists of two separate domains, with only the first domain featuring the characteristic WYL sequence motif. In the context of this study we refer to this domain as the WYL domain and to the second domain as WCX domain. The PafBC WYL domains are located at the periphery of the molecule on opposite sides, while the WCX domains come together to form a dimeric interaction module.

Notably, in this interaction module the WXC domains are arranged in a two-fold pseudo-symmetric manner, in spite of the overall asymmetric arrangement of the domains in the PafBC structure (Figure 3a and 3b). Each WCX domain harbors a four-stranded anti-parallel  $\beta$ -sheet of 4-1-3-2 topology framed by two short  $\alpha$ -helices, which contain a hydrophobic core (Figure 3b, 3c left and S6a). The main chain sharply bends at a highly conserved cis-proline into a C-terminal  $\alpha$ -helix, which crosses the C-terminal  $\alpha$ -helix of the other WCX domain (Figure S6b). Interaction of the WCX domains arises through interdigitation of two pairs of helices. The contacts are stabilized by salt-bridges, hydrogen bonds and a conserved hydrophobic island containing a pair of highly conserved leucines (Figure S6b-e). In fact, the majority of hydrogen bonding occurs at the ends of the crossed C-terminal  $\alpha$ -helices, while high conservation among the interacting residues seems to be restricted to the two leucines (Figure S6b). The PafBC interaction via the WCX domains represents a strong element in the PafBC non-covalent interaction and is probably maintained also in the active DNA-binding form.

We carried out a comparative analysis based on the protein fold of the WYL and WCX domains using the Dali fold recognition program to discover structural homologues (Holm and Laakso, 2016). The search using the isolated WCX domain yielded a very broad spectrum of hits ranging from spliceosomal proteins, elongation factors, oxidoreductases to proteases. After further manual assessment of the individual hits, we discovered that the WCX domain follows a typical ferredoxin-like fold ( $\beta\alpha\beta\beta\alpha\beta$ ) with an ancillary C-terminal  $\alpha$ -helix. Interestingly, a number of RNA-binding proteins such as human hnRNP A1 (Figure 3c, middle), CRISPR/Cas protein Cse3 (Figure 3c, right) and ribosomal proteins also contain the ferredoxin-like fold and these domains are directly involved in RNA interaction (Figure S7).

The fold of the WYL domain consists of a five-stranded anti-parallel  $\beta$ -sheet with a 5-1-2-3-4 topology preceded by an  $\alpha$ -helix (Figure 4a). The strands are strongly curved and the middle  $\beta$ 2-strand is almost twice the length of the other five, causing it to arch back over itself and resulting in a  $\beta$ -sandwich topology, where  $\beta$ -strands 5-1-2 make up one half and strands 2-3-4 the other. Middle strand  $\beta$ 2 is participating in both and connects the two halves. The eponymous WYL residues are located in  $\beta$ 3, with the highly conserved tyrosine pointing away from the hydrophobic core. Structure similarity searches using the isolated WYL domain on the Dali webserver (Holm and Laakso, 2016) returned PDB entries of proteins containing an Sm-fold, like the bacterial RNA chaperone Hfq (host factor for RNA bacteriophage  $\phi$ Q $\beta$  replication) and certain spliceosomal proteins. Proteins containing an Sm-fold are very abundant in eukaryotes, while only few examples (amongst them Hfq) have been described in bacteria. Closer comparison of the WYL domain with Hfq shows that the WYL domain Sm-fold features a slightly longer N-terminal helix, longer  $\beta$ 2, and  $\beta$ 3 strands as well as longer loops between strands  $\beta$ 1/ $\beta$ 2 and  $\beta$ 3/ $\beta$ 4 (Figure 4a and b). Many proteins of the Sm-like family were shown or predicted to bind RNA. The structural similarity of the WYL domain to Hfq and other Sm-like proteins therefore suggests that the WYL domains provide a binding site for an RNA molecule.

### **The WYL domain Sm-2 loop contains essential residues for PafBC function**

Previously, we showed that PafBC levels do not change under stress conditions and we could not detect any specific DNA binding activity towards the RecA-NDp motif *in vitro* (Fudrini Olivencia et al., 2017; Müller et al., 2018), suggesting that PafBC requires a co-activator for its activity, which is only present during stress conditions. The structural homology between the WYL domains of PafBC and Sm-folds involved in RNA binding indicates that the response-producing ligand might be an RNA molecule and the WYL domain could act as a ligand-sensing domain.

In order to deduce potential ligands and ligand binding locations from the homology between the PafBC WYL domains and the Sm-fold-containing bacterial RNA chaperone Hfq, we carried out a comparative analysis of potential binding regions. Hfq forms a homohexameric ring-shaped complex that was shown to bind RNA at three distinct sites (Figure 4b) (Updegrove et al., 2016): the proximal site exposing the  $\alpha$ -helices and binding sRNA and mRNA (shown with ligand in Figure 4b); the distal site binding A-rich oligonucleotides; and the rim (also called lateral) site literally represented by the rim of the Hfq ring. There is also increasing evidence that the C-terminus is functionally involved in RNA binding (Santiago-Frangos et al., 2016; Vecerek et al., 2008). Given the structural similarity of PafBC's WYL domains to Hfq, we compared sequence conservation patterns in both proteins. Hfq exhibits a highly conserved loop in its Sm-2 motif, containing residues that contact the RNA backbone at the proximal binding site (Figure 4b-d). Such a highly conserved loop is also



found at the corresponding locations in the PafBC WYL domains, where two arginine side chains point into the direction where the RNA ligand is positioned in Hfq (Figure 4d and e). A potential ligand may be bound at this location and transduce the signal of DNA damage to PafBC, which in turn becomes activated to carry out its role as transcriptional activator.

To test the hypothesis that the conserved sequence stretch between strands  $\beta 4$  and  $\beta 5$  in the PafBC WYL domains has functional significance, we complemented the *M. smegmatis*  $\Delta pafBC$  strain with PafBC mutants featuring amino acid substitutions at this location and assessed the viability of the mutants in presence of the DNA damaging agent mitomycin C (MMC). We also chose residues at other sites in the WYL domain based on sequence conservation and structural similarity to Hfq. Specifically, we selected the conserved tyrosine that is part of the WYL triplet, another conserved tyrosine (sometimes histidine) in strand  $\beta 1$ , and the conserved patch between  $\beta 4$  and  $\beta 5$  for mutation. The chosen residues were mutated to alanine and the mutations were introduced separately into PafB or PafC or into both proteins simultaneously. To reduce the permutation space, we decided to treat the two arginine residues in the  $\beta 4/\beta 5$  loop as functionally redundant (i.e. they were simultaneously substituted with alanine). Since the HTH domain of PafC would not be able to bind DNA in the observed conformation (Figure 1 and 2b), we also deleted the HTH domains individually to establish if they are required for PafBC function.

To assess the viability of the PafBC mutant strains, the cells were first grown for a defined period of time in presence of increasing concentrations of MMC. Subsequently, the dye resazurin was added, which is reduced by living (but not by dead) cells to resorufin, giving rise to a color change. Wild type *M. smegmatis* cells grow in presence of up to 100 ng/ml MMC, while the  $\Delta pafBC$  strain shows growth only up to about 8 ng/ml MMC, which is in agreement with the previously determined minimal inhibitory concentrations for these strains (Figure 5a) (Fudrini Olivencia et al., 2017).

Complementation of the *pafBC* knockout strain with wild type PafBC restores the viability to the level observed for the wild type. Deletion of either the HTH domain of PafB or PafC leads to the same reduced viability as observed for the  $\Delta pafBC$  strain (Figure 5a), indicating that both domains are required for the function of PafBC. It has to be noted that, in contrast to  $\Delta HTH-C$ , which expressed well, the expression of  $\Delta HTH-B$  was barely detectable and may therefore not be sufficient for complementation (Figure 5g and 5h). Nevertheless, the complementation experiment demonstrates that the second HTH domain (HTH-C), which in our structure is in an inaccessible conformation for DNA-binding, is required for a fully functional PafBC complex.

We then tested the alanine point mutants for complementation. For most alanine-substitution mutants a pattern could be observed (Figure 5b-f): If the mutation is present in only one of the WYL domains, the viability is only moderately affected, but in case both WYL domains carry the mutation, the effect seems to be additive and the viability of the cells is lowered to the level of the knockout strain. This is the case for the double arginine mutants (Figure 5c), the tyrosine of the WYL triplet (Figure 5d), and the phenylalanine of the  $\beta 4/\beta 5$  loop. Mutation of the  $\beta 1$  histidine/tyrosine leads to a comparable result, except that the decrease in viability is milder if one of the WYL domains carries the mutations, and the effect is less severe than in the knockout strain if

both WYL domains are mutated (Figure 5b). Furthermore, mutation of the serine/aspartate in the  $\beta 4/\beta 5$  loop did not affect viability (Figure 5e).

Our results demonstrate that the conserved residues in the  $\beta 4/\beta 5$  loop and  $\beta 1$  of the WYL domain are required for the function of PafBC, likely because they interact with a signal-transducing ligand. Furthermore, successive inactivation of the subunits has an additive effect, suggesting that there are two ligand binding sites present, one at each WYL domain. In summary, our observations show that both subunits of the PafBC complex are functional and both WYL domains are required for full viability.

### **The WYL domain is mainly associated with DNA-binding domains**

Based on the structural analysis of the WYL domains in the PafBC complex and the fact that this domain occurs also in other bacterial regulators, we carried out a thorough bioinformatics analysis of WYL domain-containing proteins to understand, in which functional context they occur and how widely distributed they are.

We computationally analyzed the co-occurrence of the WYL domain with other domains along with its taxonomic distribution using hidden Markov models (HMMs) (Eddy, 2011). HMMs are widely used for finding distant protein homologs and they provide the basis for one of the largest protein family databases, Pfam, which groups proteins containing the same domain into families. Our structural analysis has shown that the PafBC C-terminal part originally assigned as “WYL” domain as a whole, in fact consists of two domains, the actual WYL domain and a C-terminal extension (WCX) domain. Based on the domain boundaries of the WYL domains in our structure and sequence alignment with other PafBC orthologs, we generated a WYL domain HMM and used it to retrieve all WYL domain-containing proteins from the UniProt reference proteomes yielding 15’079 entries (Table S1). The resulting entries were distributed across 5’330 different species with only 81 sequences from 50 species among *Eukaryota*, *Archaea* or *Viruses*, which were mostly candidate species (Table S2). Thus, the WYL domain appears to be limited to bacteria and we restricted our subsequent analyses to bacterial sequences.

In order to identify domain families associated with WYL domain proteins, the retrieved WYL domain-containing bacterial sequences were annotated based on all Pfam HMMs and additional HMM profiles we generated for the WCX domain and the PafBC N-terminal winged HTH domain that was not recognized by any of the existing Pfam HMMs. Two observations can immediately be made from the final set of domain architecture classes (Figure 6a and S8b): First, the majority of classes, covering more than 90% of sequences, shows co-occurrence of the WYL domain with an HTH domain preceding it. Second, the WYL domain is primarily present together with a C-terminally located WCX domain, and only about 25% of sequences exhibit the WYL domain alone.

About two thirds of all sequences are found in class A, which also comprises all PafB and PafC sequences. The second largest group, class B (15% of all hits), contains proteins with only an HTH and a WYL domain, lacking the WCX domain. Notably, class C could also be viewed as a subgroup of class A, as it is made up of natural fusion proteins of actinobacterial PafB and PafC homologs. A significant number of sequences contain a Helicase C3 domain in combination with the WYL domain, but lacking the WCX domain.



The distribution of WYL domain-containing proteins among bacterial phyla reflects the distribution of these phyla in the reference proteomes, showing that WYL domains are ubiquitous among bacteria (Figure S8a). Interestingly, the gram-positive phyla of *Actinobacteria* and *Firmicutes* exhibit on average roughly 5.2 and 2.8 WYL domain-containing proteins per organism, respectively, while the gram-negative phyla of *Proteobacteria* and *Bacteroidetes* show only 2.0 and 2.1 average WYL domain-containing proteins per organism, respectively (Figure 6b). By analyzing the taxonomic distribution for each domain architecture class, we observed that the PafBC-like class A is most prominently found in *Actinobacteria*, *Firmicutes*, and *Bacteroidetes*, but much less abundant in *Proteobacteria* (Figure 6c). On the other hand, more than two thirds of the HTH-WYL architecture members (class B) are found in proteobacterial species (Figure 6d).

Taken together, our analysis shows that the majority of WYL domain-containing proteins are transcriptional regulators based on the presence of an HTH domain. It therefore seems likely that the mechanism of transcriptional regulation and signal relay employed by PafBC, although currently unknown, is a widespread principle found in almost all bacteria. Moreover, in some phyla multiple of these transcriptional regulators are present in one organism, suggesting that WYL domain-containing regulators may be involved in different pathways.

## Discussion

During the mycobacterial DNA damage response, the heterodimeric transcriptional regulator PafBC activates most of the genes required for an adequate response to DNA stress (Fudrini Olivencia et al., 2017; Müller et al., 2018). However, understanding and experimentation concerning the regulatory mechanism of PafBC was largely hampered by a lack of knowledge about its molecular structure. This limitation has manifested also in other studies concerning WYL domain-containing proteins (Andis et al., 2018; Hein et al., 2013; Modell et al., 2014; Yan et al., 2018). Our computational analysis showed that roughly 90% of all WYL domain-containing proteins possess an N-terminal HTH domain suggesting that these are transcriptional regulators (Figure 6a). Furthermore, our analysis revealed the WYL domain as a domain specific to bacteria that is present in nearly all bacterial phyla (Figure 6b and S8). Thus, it is conceivable that the regulatory mechanism employed by PafBC might represent a shared mode of action for all of these regulators. This is not only an exciting possible concept, but might also be helpful in gaining a full understanding of the nature of the regulation.

We obtained the crystal structure of PafBC in the absence of DNA, in a largely asymmetric domain arrangement that is likely characteristic for the non-activated state. The domains are connected through long loops, suggesting a great degree of flexibility for the entire protein and that the protein might undergo large domain movements upon DNA binding, where it could eventually adopt a more symmetric arrangement, as observed for the WCX domains (Figure 7). In such a state, the helices connecting HTH-C and WYL-C ( $\alpha 4'$ ,  $\alpha 4''$ ,  $\alpha 4'''$ ) could merge into a single helix and act as the counterpart to PafB helix  $\alpha 4$  in a coiled-coiled fashion at the center of the protein. Such an interaction could be mediated by the row of hydrophobic residues that are featured along the axes of helices  $\alpha 4'$ - $\alpha 4'''$ . Also, the HTH-C domain was observed in a state inaccessible for DNA binding (Figure 2b). Besides their role in protein-DNA interaction, winged HTH domains were found to mediate protein-protein interactions (Littlefield and Nelson, 1999; Woo et al., 2009; Zheng et al., 2002). Thus, the conformation of HTH-C may represent a state that is part of a regulatory mechanism, in which PafBC is prevented from efficient DNA binding under non-stress conditions. In

agreement with this notion, the PafBC mutant lacking HTH-C cannot complement the phenotype of  $\Delta pafBC$  observed under DNA stress (Figure 5a), suggesting that HTH-C must fulfill an essential function, i.e. DNA binding/recognition. Furthermore, the recognition helices in the HTH domains of PafB and PafC are different in length and also in amino acid composition (Figure S4a and S5), and likewise the PafBC binding motif (RecA-NDp) is non-palindromic (Müller et al., 2018). Together with the regulatory switch, this could then also explain why PafBC is a heterodimer.

Our results provide key insights into the WYL and WCX domains, revealing that they adopt folds similar to proteins associated with RNA binding (Figure 3, 4 and S7). Considering that the activation of PafBC upon DNA damage does not rely on changes in protein levels (Fudrini Olivencia et al., 2017), the possibility of another, stress-dependent factor required to elicit PafBC activity is likely. The results obtained in the complementation study with single amino acid substitutions in the WYL domains strongly suggest a binding interface for a signal-transducing ligand (Figure 5). Such a potential factor may thus well be an RNA molecule or, in a broader context, a nucleic-acid or nucleic acid derivative, relaying the signal of DNA damage to PafBC by recognition at the WYL and/or WCX domains. In fact, the WYL domain of PIF1 helicase from *Thermotoga elfii* was shown to bind single-stranded DNA, thereby stimulating helicase activity (Andis et al., 2018). Binding of single-stranded DNA may be conceivable for the WYL domain proteins of class J of our computational analysis, which are also associated with a helicase domain.

The Sm-fold of the WYL domain is characteristic of eukaryotic RNA binding proteins, the Sm proteins. Their ring-shaped assemblies are core components of the spliceosomal snRNPs (small nuclear ribonucleoproteins) (Bertram et al., 2017). Through analogy, the bacterial protein Hfq is considered the sole representative of the Sm-like/LSm family based on its hexameric assembly state and RNA chaperone function (Khusial et al., 2005). Interestingly, no Hfq homolog has been identified in actinobacteria to date using sequence searches (Chao and Vogel, 2010), but we found WYL domain-containing proteins to be significantly enriched in the actinobacterial phylum (Figure 6b). It is possible that some of these actinobacterial WYL domain-containing proteins carry out a similar function to Hfq.

The crystal structure of PafBC provides the framework for understanding the mechanism by which PafBC connects the signal of DNA stress with a transcriptional response through use of its WYL/WCX domains. These results will also help us to better understand WYL domain-containing proteins in general.

## Methods

### Protein expression and purification of *Arthrobacter aureescens* PafBC

Full-length <sup>Aau</sup>PafBC was amplified from genomic DNA of *Arthrobacter aureescens* strain 579 (DSM-20116) using the primers ACGCGCTTGCTGCTTTCC (forward) and CTAGCCAGCCTTGCTGCCG (reverse), which were designed based on the sequenced genome of strain TC1 (NC\_008711; locus tag AAur\_2182). The amplicon was cloned into a temporary vector and a truncated variant of <sup>Aau</sup>PafBC (<sup>Aau</sup>PafBC $\Delta$ NC) missing the first 17 amino acids at the N-terminus and the last 7 amino acids at the C-terminus was amplified from this vector using primers GCATCCCGCACCGAACG (forward) and TGAGTCGTACTGCACCAAAG (reverse). The amplicon of <sup>Aau</sup>PafBC $\Delta$ NC was cloned into an isopropyl- $\beta$ -D-thiogalactopyranosid (IPTG)-inducible expression vector with a cleavable His<sub>6</sub>-TEV tag at the N-terminus. Selenomethionine-labeled protein was expressed according to a procedure adapted from

(Doublié, 2007): *E. coli* Rosetta (DE3) cells harboring the expression vector were grown as shaking cultures at 37°C in M9 medium (M9 salts supplemented with 2 mM MgSO<sub>4</sub>, 0.1 mM CaCl<sub>2</sub>, 0.5% w/v glucose, 2 mg/l biotin, 2 mg/l thiamine, 0.03 mg/l FeSO<sub>4</sub>). At an OD<sub>600</sub> of 0.5, 100 mg/ml of phenylalanine, lysine, and threonine, 50 mg/ml of isoleucine, leucine, and valine, as well as 80 mg/ml of selenomethionine (Chemie Brunschwig) were added as solid powder to the cultures, which were further incubated for 30 min. Expression was then induced with 0.5 mM IPTG and cells were further incubated at 16°C overnight. Cells were harvested (F9S, 7'000 rpm, 10 min, 4°C) and pellets were resuspended in lysis buffer (50 mM HEPES-NaOH pH 7.8/4°C, 300 mM NaCl, 2 mM TCEP). The cell suspension was homogenized using a Heidolph DIA600 mixer and cells were lysed by high pressure shear force using a Microfluidizer M110-L device (Microfluidics; 5 passes, 11'000 psi chamber pressure). After removal of cell debris (SS34, 20'000 rpm, 4°C, 30 min), the cleared lysate was supplemented with 1 mM PMSF, 1x cOmplete EDTA-free protease inhibitors (Roche), 50 U/ml DNase I, 10 mM imidazole and incubated for 30 min on ice. The lysate was passed over a self-packed Ni<sup>2+</sup>-charged IMAC Sepharose 6 Fast Flow (GE Healthcare) column, and bound protein was eluted step-wise with lysis buffer containing 80 mM to 250 mM imidazole. After pooling protein-containing elution fractions, His-tagged TEV protease was added to a 1:30 molar ratio and the protein sample was dialyzed against 25 mM HEPES-NaOH pH 7.8/4°C, 150 mM NaCl, 2 mM DTT, 1 mM EDTA at 4°C overnight. TEV protease was removed by affinity chromatography and the protein sample was dialyzed against 25 mM HEPES-NaOH pH 7.8/4°C, 40 mM NaCl, 2 mM DTT, 1 mM EDTA at 4°C overnight. The protein was further loaded on a Source 30Q column and eluted with a 50 mM to 400 mM NaCl gradient in 25 mM HEPES-NaOH pH 7.8/4°C, 1 mM TCEP. Protein-containing elution fractions were pooled and concentrated using an Amicon Ultra 30K centrifugal filter (3'500 g, 4°C). The concentrated protein was run on a self-packed 100 ml Superose 12 prep grade column in crystallization buffer (10 mM HEPES-NaOH pH 7.8/4°C, 50 mM NaCl, 1 mM TCEP, 0.1 mM EDTA). Protein was concentrated as above to 22 mg/ml, aliquotted, frozen in liquid nitrogen and stored at -80°C until use.

### Crystallization of *Arthrobacter aureescens* PafBC

Crystals of <sup>Aau</sup>PafBCΔNC were grown in sitting drop vapor diffusion plates (Hampton Research) at 4°C by mixing 1-2 μl protein solution (5.6-7.5 mg/ml) and 1 μl reservoir solution. Crystals appeared after 2-3 days using reservoir solutions containing 100 mM Bis-Tris-propane pH 8.8 to 9.2 (20°C), 80-160 mM KSCN, 18-21% PEG-2000 (Sigma). PEG-2000 was added in 5% steps to reach a final concentration of 36% w/v in drops containing crystals before flash cooling the crystals in liquid nitrogen.

### Data collection, experimental phasing, structure determination, and refinement

Reflection image data was collected at the X06SA beamline of the Swiss Light Source (SLS, Paul-Scherrer-Institut, Villigen, Switzerland) at 100 K and 12'670 eV beam energy (0.978561 Å). Diffraction images were processed using XDS (Kabsch, 2010) and scaled using AIMLESS (Evans and Murshudov, 2013). Determination of heavy atom sites, initial phases, and crude main chain tracing were carried out using the SHELX programs (Sheldrick, 2010). The resulting experimental electron density map displayed easily discernible protein features. The initial model from SHELX was further extended with PHENIX AutoBuild (Terwilliger et al., 2008) and subsequent iterative model building and refinement was carried out using Coot (Emsley et al., 2010) and phenix.refine (Afonine et al., 2012), respectively.

## Structure visualization

Graphical representations of protein structures were prepared using UCSF Chimera v1.12 build 41623 (Pettersen et al., 2004). Amino acid residue conservation was calculated using AL2CO (independent counts, BLOSUM-62 matrix) (Pei and Grishin, 2001). For that, PafBC protein sequences of 23 different organisms spanning the entire actinobacterial phylum (UniProt accession numbers P9WIM1, I7G3U5, C0ZZU3, A1SK18, A7BCC5, A4X749, A0LU62, A6W976, Q8NQE2, Q9RJ64, H6RJ02, D2ATU2, C8XAP4, C7PVW0, A0A160VN40, Q0RLT0, A0A1D7W444, A0A1H2KTF9, B2GIN6, Q47P13, C5CBV3, A9WSH6, A1R6R2, P9WIL9, A0QZ41, C0ZZU2, A1SK19, A7BCC6, A4X750, A0LU63, A6W977, Q8NQE3, Q9RJ65, H6RJ01, D2ATU1, C8XAP3, C7PVW1, A0A161KHT8, Q0RLS9, A0A1D7W495, A0A1H2KTV3) were aligned against the *A. aureescens* strain 579 PafBC protein sequence. To be able to calculate conservation across the entire length of naturally fused PafBC proteins, sequences of separately encoded PafB/PafC proteins were concatenated prior to alignment.

## Mutational screening of *Mycobacterium smegmatis* PafBC

The coding sequence of *pafBC* was amplified from *Mycobacterium smegmatis* mc2-155 SMR5 (Sander et al., 1995) genomic DNA and cloned into a pMyNT-derived integrative plasmid (pMyNT template provided by A. Geerlof, EMBL Hamburg). As promoter, a DNA fragment containing 347 bp upstream of the *pafA* coding sequence was inserted in front of *pafBC* to generate the *pafBC* complementation plasmid. Variants were then generated by KLD site-directed mutagenesis or by Gibson assembly using the *pafBC* complementation plasmid as template. The various plasmids were then transformed into the *Mycobacterium smegmatis*  $\Delta$ *pafBC* strain (Fudrini Olivencia et al., 2017), and viability in presence of mitomycin C was assessed with the resazurin assay as described previously (Müller et al., 2018).

## Computational analysis of WYL domain-containing domains

All alignments were generated using ClustalO v1.2.4 with default settings (Sievers et al., 2011). Manual analyses of alignments were performed in Jalview v2.10.5 (Waterhouse et al., 2009). Steps involving hidden Markov-models (HMMs) were conducted using the programs “hmmbuild”, “hmmsearch”, and “hmmscan” from the software suite HMMER v3.2.1 (hmm.org) (Eddy, 2011). A local search database was created from all UniProt reference proteomes (September 2018 release, uniprot.org) (The UniProt Consortium, 2017).

For the initial analysis, PafBC homologs from other actinobacteria were identified by BLAST (blast.ncbi.nlm.nih.gov; restricted search to actinobacterial species, otherwise default settings) (Camacho et al., 2009) using *Mycobacterium smegmatis* PafB or PafC as input. The identity of the obtained PafBC homologs was cross-checked on the genome level for the operon organization of the genes and association with the Pup-proteasome system gene locus. In total, ten PafB/PafC sequence pairs were retrieved (UniProt accession numbers P9WIM1, P9WIL9, I7G3U5, A0QZ41, A7BCC5, A7BCC6, A4X749, A4X750, C0ZZU3, C0ZZU2, A1SK18, A1SK19, A0LU62, A0LU63, A6W976, A6W977, Q8NQE2, Q8NQE3, Q9RJ64, Q9RJ65). A global alignment of these sequences was used to generate an HMM with “hmmbuild” (default settings), which was subsequently used to search against the reference proteomes database with “hmmsearch” (command line options were -E 1 --domE 1 --incE 0.01 --incdomE 0.03). A list of domains contained in the retrieved sequences was obtained with the “hmmscan” module (command line options were -E 0.1 --domE 0.1 --incE 0.01 --incdomE 0.03) using all HMMs from the Pfam database release 32.0 (El-Gebali et al., 2019). Domain lengths were analyzed based on the envelope boundaries given by “hmmscan” for sequences with a domain score above 30.0 (independent domain e-value < 0.001).

Because the Pfam HMM profile of the WYL domain (Pfam-WYL) includes both WYL and WCX domain, we had to define a new WYL domain HMM in order to annotate it correctly for our analysis. To build the WYL HMM, 250 sequences were randomly sampled from sequences with a Pfam-WYL length greater than 127 residues and another 250 sequences were randomly sampled from sequences with a Pfam-WYL length less than 127 residues (other thresholds: domain score > 30.0, independent domain e-value < 0.001). The threshold was chosen based on the domain boundaries seen in the crystal structure of <sup>Aau</sup>PafBC and the length distribution of C-terminal region of the retrieved Pfam-WYL-containing proteins. Sequences with obvious defects in the WYL region were discarded manually resulting in 477 entries that were used for alignment. From this alignment, the WYL domain boundaries were established using the N-terminal boundary of the Pfam-WYL, while the C-terminal boundary was chosen by the C-termini of short Pfam-WYs and the crystal structure. The alignment was then trimmed to the WYL boundaries, sequences with a pairwise identity above 70% were clustered using CD-HIT v4.6.8 (command line options -n 4 -c 0.7) (Li and Godzik, 2006) and an HMM was generated as above. To mature the WYL HMM, an iterative approach similar to the Pfam HMM generation was chosen. Sequences were retrieved from the reference proteomes using the WYL HMM to generate a full alignment, which was then again trimmed to the WYL domain boundaries to make up a new seed alignment used to build a new HMM (gathering threshold: domain score > 27.0). The process was repeated for a total of three iterations.

To build the HMM for the unrecognized (winged) HTH domains in PafBC and many other WYL domain-containing proteins, 500 sequences were randomly sampled from the group of sequences with a median length of 326 amino acids and containing only a WYL domain (thresholds: domain score > 30.0, independent domain e-value < 0.001). The sequences were aligned and curated based on the presence of the conserved blocks representing the helices of the HTH fold. The sequences were further split into groups exhibiting a PafB-like or PafC-like wing sequence of 184 and 235 sequences, respectively. For each group, an HMM was built and searched against reference proteomes database as described above, but using an identity threshold of 90% and iterating only once in order to keep a separation to other HTH-type domains.

To build the HMM for the C-terminal extension found in many WYL domains (WCX), 500 sequences were randomly sampled from the group containing a Pfam-WYL with a length greater than 127 amino acids (thresholds: domain score > 30.0, independent domain e-value < 0.001). Sequences missing the C-terminal conserved block were removed, leaving 434 sequences for HMM generation. The HMM generated from the manually curated seed alignment was used without iteration.

The custom HMMs were then added to the local search database of the Pfam HMM database (see also above).

To establish the domain architecture classes of the WYL domain-containing proteins, the WYL HMM was used to search against the reference proteome database with “hmmsearch” (command line options as above). Obtained sequences with a domain score > 27.0 were analyzed for the presence of other domains using “hmmsearch” together with the local Pfam HMM database including the custom HMMs as described above. Protein sequences were then categorized according to their sequence length, their identified domains (on the level of Pfam clans; independent domain e-value < 0.01) and the domain length. Domain architecture classes with large regions apparently containing no domain were checked manually by alignment for presence of conserved features, by alignment to other classes and by alignment to the manually curated PafBC seed sequences (see above). They



were then grouped together with other classes, where appropriate. Due to low overall abundance and being mostly candidate species, all non-bacterial sequences were excluded from the main analysis.

## Acknowledgements

We thank Takashi Tomizaki and the staff of beamline X06SA at the Swiss Light Source (Villigen, Switzerland) for support with data collection; Beat Blattmann and Céline Stutz-Ducommun of the Protein Crystallization Center (University of Zurich) for support with the initial screens; Marcel Bolten for support with X-ray data analysis. The research was funded by the Swiss National Science Foundation (SNSF), grant no. 31003A-163314.

## Author Contributions

AUM, NB, and EWB designed research and analyzed data. AUM performed protein purification, crystallization, and *in vivo* experiments. AUM and ML analyzed crystallographic data. AUM, NB, and EWB wrote the manuscript. All authors contributed to editing of the manuscript.

## Competing Interests statement

The authors declare no conflict of interest.

## References

- Afonine, P.V., Grosse-Kunstleve, R.W., Echols, N., Headd, J.J., Moriarty, N.W., Mustyakimov, M., Terwilliger, T.C., Urzhumtsev, A., Zwart, P.H., and Adams, P.D. (2012). Towards automated crystallographic structure refinement with phenix.refine. *Acta Crystallographica Section D* *D68*, 352-367.
- Andis, N.M., Sausen, C.W., Alladin, A., and Bochman, M.L. (2018). The WYL Domain of the PIF1 Helicase from the Thermophilic Bacterium *Thermotoga elfii* is an Accessory Single-Stranded DNA Binding Module. *Biochemistry* *57*, 1108-1118.
- Belinky, F., Rogozin, I.B., and Koonin, E.V. (2017). Selection on start codons in prokaryotes and potential compensatory nucleotide substitutions. *Sci Rep* *7*, 12422.
- Bertram, K., Agafonov, D.E., Dybkov, O., Haselbach, D., Leelaram, M.N., Will, C.L., Urlaub, H., Kastner, B., Luhrmann, R., and Stark, H. (2017). Cryo-EM Structure of a Pre-catalytic Human Spliceosome Primed for Activation. *Cell* *170*, 701-713 e711.
- Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K., and Madden, T.L. (2009). BLAST+: architecture and applications. *BMC Bioinformatics* *10*.
- Chao, Y., and Vogel, J. (2010). The role of Hfq in bacterial pathogens. *Curr Opin Microbiol* *13*, 24-33.
- Clark, K.L., Halay, E.D., Lai, E., and Burley, S.K. (1993). Co-crystal structure of the HNF-3/fork head DNA-recognition motif resembles histone H5. *Nature* *364*, 412-420.
- Davis, E.O., Dullaghan, E.M., and Rand, L. (2002a). Definition of the Mycobacterial SOS Box and Use To Identify LexA-Regulated Genes in *Mycobacterium tuberculosis*. *Journal of Bacteriology* *184*, 3287-3295.
- Davis, E.O., Springer, B., Gopaul, K.K., Papavinasasundaram, K.G., Sander, P., and Böttger, E.C. (2002b). DNA damage induction of *recA* in *Mycobacterium tuberculosis* independently of RecA and LexA. *Molecular Microbiology* *46*, 791-800.
- DeJesus, M.A., Sacchetti, J.C., and Ioerger, T.R. (2013). Reannotation of translational start sites in the genome of *Mycobacterium tuberculosis*. *Tuberculosis (Edinb)* *93*, 18-25.
- Doublié, S. (2007). Production of Selenomethionyl Proteins in Prokaryotic and Eukaryotic Expression Systems. In *Macromolecular Crystallography Protocols Methods in Molecular Biology*, J. Walker, and S. Doublié, eds.
- Eddy, S.R. (2011). Accelerated Profile HMM Searches. *PLoS Computational Biology* *7*, e1002195.



- El-Gebali, S., Mistry, J., Bateman, A., Eddy, S.R., Luciani, A., Potter, S.C., Qureshi, M., Richardson, L.J., Salazar, G.A., Smart, A., *et al.* (2019). The Pfam protein families database in 2019. *Nucleic Acids Research* 47, D427-D432.
- Emsley, P., Lohkamp, B., Scott, W.G., and Cowtan, K. (2010). Features and Development of Coot. *Acta Crystallographica Section D* 66, 486-501.
- Evans, P.R., and Murshudov, G.N. (2013). How good are my data and what is the resolution? *Acta Crystallographica Section D* 69, 1204-1214.
- Fudrini Olivencia, B., Müller, A.U., Roschitzki, B., Burger, S., Weber-Ban, E., and Imkamp, F. (2017). *Mycobacterium smegmatis* PafBC is involved in regulation of DNA damage response. *Scientific Reports* 7.
- Galletto, R., Amitani, I., Baskin, R.J., and Kowalczykowski, S.C. (2006). Direct observation of individual RecA filaments assembling on single DNA molecules. *Nature* 443, 875-878.
- Gamulin, V., Cetkovic, H., and Ahel, I. (2004). Identification of a promoter motif regulating the major DNA damage response mechanism of *Mycobacterium tuberculosis*. *FEMS Microbiology Letters* 238, 57-63.
- Hein, S., Scholz, I., Voß, B., and Hess, W.R. (2013). Adaptation and modification of three CRISPR loci in two closely related cyanobacteria. *RNA Biology* 10, 852-864.
- Holm, L., and Laakso, L.M. (2016). Dali server update. *Nucleic Acids Research* 44, W351-W355.
- Jones, S., van Heyningen, P., Berman, H.M., and Thornton, J.M. (1999). Protein-DNA interactions: a structural analysis. *Journal of molecular biology* 287, 877-896.
- Kabsch, W. (2010). XDS. *Acta Crystallographica Section D* 66, 125-132.
- Khusial, P., Plaag, R., and Zieve, G.W. (2005). LSm proteins form heptameric rings that bind to RNA via repeating motifs. *Trends in Biochemical Sciences* 30, 522-528.
- Kreuzer, K.N. (2013). DNA damage responses in prokaryotes: regulating gene expression, modulating growth patterns, and manipulating replication forks. *Cold Spring Harb Perspect Biol* 5, a012674.
- Li, W., and Godzik, A. (2006). CD-HIT: a fast program for clustering and comparing large sets of protein or nucleotide sequences. *Bioinformatics* 22, 1658-1659.
- Little, J.W., Edmiston, S.H., Pacelli, L.Z., and Mount, D.W. (1980). Cleavage of the *Escherichia coli* *lexA* protein by the *recA* protease. *Proceedings of the National Academy of Sciences, USA* 77, 3225-3229.
- Little, J.W., Mount, D.W., and Yanischperron, C.R. (1981). Purified LexA Protein Is a Repressor of the RecA and LexA Genes. *PNAS* 78, 4199-4203.
- Littlefield, O., and Nelson, H.C. (1999). A new use for the 'wing' of the 'winged' helix-turn-helix motif in the HSF-DNA cocrystal. *Nat Struct Biol* 6, 464-470.
- Maslowska, K.H., Makiela-Dzbenka, K., and Fijałkowska, I.J. (2018). The SOS System: A Complex and Tightly Regulated Response to DNA Damage. *Environ Mol Mutagen*.
- Modell, J.W., Kambara, T.K., Perchuk, B.S., and Laub, M.T. (2014). A DNA Damage-Induced, SOS-Independent Checkpoint Regulates Cell Division in *Caulobacter crescentus*. *PLoS Biology* 12, e1001977.
- Müller, A.U., Imkamp, F., and Weber-Ban, E. (2018). The Mycobacterial LexA/RecA-Independent DNA Damage Response Is Controlled by PafBC and the Pup-Proteasome System. *Cell Reports* 23, 3551-3564.
- Pei, J., and Grishin, N.V. (2001). AL2CO: calculation of positional conservation in a protein sequence alignment. *Bioinformatics* 17, 700-712.
- Pettersen, E.F., Goddard, T.D., Huang, C.C., Couch, G.S., Greenblatt, D.M., Meng, E.C., and Ferrin, T.E. (2004). UCSF Chimera--a visualization system for exploratory research and analysis. *J Comput Chem* 25, 1605-1612.
- Phizicky, E.M., and Roberts, J.W. (1981). Induction of SOS functions: regulation of proteolytic activity of *E. coli* RecA protein by interaction with DNA and nucleoside triphosphate. *Cell* 25, 259-267.
- Rajagopalan, S., Teter, S.J., Zwart, P.H., Brennan, R.G., Phillips, K.J., and Kiley, P.J. (2013). Studies of IscR reveal a unique mechanism for metal-dependent regulation of DNA binding specificity. *Nat Struct Mol Biol* 20, 740-747.

- Rand, L., Hinds, J., Springer, B., Sander, P., Buxton, R.S., and Davis, E.O. (2003). The majority of inducible DNA repair genes in *Mycobacterium tuberculosis* are induced independently of RecA. *Molecular Microbiology* 50, 1031-1042.
- Sander, P., Meier, A., and Böttger, E.C. (1995). *rpsL*<sup>+</sup>: A dominant selectable marker for gene replacement in mycobacteria. *Molecular Microbiology* 16, 991-1000.
- Santiago-Frangos, A., Kavita, K., Schu, D.J., Gottesman, S., and Woodson, S.A. (2016). C-terminal domain of the RNA chaperone Hfq drives sRNA competition and release of target RNA. *Proc Natl Acad Sci U S A* 113, E6089-E6096.
- Sheldrick, G.M. (2010). Experimental phasing with SHELXC/D/E: combining chain tracing with density modification. *Acta Crystallographica Section D* D66, 479-485.
- Sievers, F., Wilm, A., Dineen, D., Gibson, T.J., Karplus, K., Li, W., Lopez, R., McWilliam, H., Remmert, M., Söding, J., *et al.* (2011). Fast, scalable generation of high-quality protein multiple sequence alignments using Clustal Omega. *Molecular Systems Biology* 7.
- Smollett, K.L., Smith, K.M., Kahramanoglou, C., Arnvig, K.B., Buxton, R.S., and Davis, E.O. (2012). Global Analysis of the Regulon of the Transcriptional Repressor LexA, a Key Component of SOS Response in *Mycobacterium tuberculosis*. *Journal of Biological Chemistry* 287, 22004-22014.
- Terwilliger, T.C., Grosse-Kunstleve, R.W., Afonine, P.V., Moriarty, N.W., Zwart, P.H., Hung, L.-W., Read, R.J., and Adams, P.D. (2008). Iterative model building, structure refinement and density modification with the PHENIX AutoBuild wizard. *Acta Crystallographica Section D* D64, 61-69.
- The UniProt Consortium (2017). UniProt: the universal protein knowledgebase. *Nucleic Acids Research* 45, D158-D169.
- Updegrave, T.B., Zhang, A., and Storz, G. (2016). Hfq: the flexible RNA matchmaker. *Curr Opin Microbiol* 30, 133-138.
- Vecerek, B., Rajkowitsch, L., Sonnleitner, E., Schroeder, R., and Blasi, U. (2008). The C-terminal domain of Escherichia coli Hfq is required for regulation. *Nucleic Acids Research* 36, 133-143.
- Waterhouse, A.M., Procter, J.B., Martin, D.M., Clamp, M., and Barton, G.J. (2009). Jalview Version 2— a multiple sequence alignment editor and analysis workbench *Bioinformatics* 25, 1189-1191.
- Wisedchaisri, G., Chou, C.J., Wu, M., Roach, C., Rice, A.E., Holmes, R.K., Beeson, C., and Hol, W.G. (2007). Crystal structures, metal activation, and DNA-binding properties of two-domain IdeR from *Mycobacterium tuberculosis*. *Biochemistry* 46, 436-447.
- Woo, J.S., Lim, J.H., Shin, H.C., Suh, M.K., Ku, B., Lee, K.H., Joo, K., Robinson, H., Lee, J., Park, S.Y., *et al.* (2009). Structural studies of a bacterial condensin complex reveal ATP-dependent disruption of intersubunit interactions. *Cell* 136, 85-96.
- Yan, W.X., Chong, S., Zhang, H., Makarova, K.S., Koonin, E.V., Cheng, D.R., and Scott, D.A. (2018). Cas13d Is a Compact RNA-Targeting Type VI CRISPR Effector Positively Modulated by a WYL-Domain-Containing Accessory Protein. *Molecular Cell* 70, 327-339.
- Zheng, N., Schulman, B.A., Song, L., Miller, J.J., Jeffrey, P.D., Wang, P., Chu, C., Koepp, D.M., Elledge, S.J., Pagano, M., *et al.* (2002). Structure of the Cul1-Rbx1-Skp1-F boxSkp2 SCF ubiquitin ligase complex. *Nature* 416, 703-709.

## Figure Legends

**Figure 1: Crystal structure of the PafBC complex from *Arthrobacter aurescens* (<sup>Aau</sup>PafBCΔNC).** The protein is the product of a natural gene fusion and encompasses a PafB and a PafC part (bottom scheme) connected via a 20 amino acids long linker (yellow). The winged helix-turn-helix domains (HTH) are colored in red, WYL domains are colored in blue and the C-terminal extension of the WYL domain (WCX) is colored in light blue. Other parts of the structure belonging to PafB are colored in dark gray, while the remaining PafC parts are colored in light grey. Dashed lines bridge gaps in the model. See also Figure S1, S2, and S3.

**Figure 2: The non-activated conformation of PafBC buries the HTH domain of PafC.** (a) The HTH domain of PafB (HTH-B, red) adopts a typical winged helix-turn-helix fold (wHTH). The recognition helices (H3) of wHTH domains typically insert into the major groove of the DNA, while the wings establish contact to the minor groove (small inset, middle). The recognition helix of HTH-B (orange) is exposed and available to accept a DNA ligand. The main chain of H1 and the wing coordinate a potassium ion (lilac sphere). (b) The HTH domain of PafC (HTH-C, red) on the other hand forms part of the protein core. H3 of HTH-C (orange) is shorter by two helical turns compared to H3 of HTH-B, and an unstructured loop reaches into the protein core. Additionally, the β-sheet of the wing extends the β-sheet of the WYL domain of PafB (WYL-B, blue). Dashed lines bridge gaps in the models. See also Figure S4.

**Figure 3: The C-terminal extension (WCX) domains of PafB and PafC contain a ferredoxin-like fold and display two-fold rotational pseudosymmetry.** (a) The WCX domains of PafB (WCX-B, lilac) and PafC (WCX-C, pink) contact each other in the crystal structure. Dashed lines bridge gaps in the model. (b) WCX-B and WCX-C exhibit a two-fold rotational symmetry axis (circle and dashed line). (c) The WCX domains (left; shown for PafC; residues 585-664) contain a ferredoxin-like fold with an additional C-terminal α-helix (α3). Very versatile and present in proteins with highly diverse functions, the ferredoxin-like fold is also found in many RNA-binding proteins such as human hnRNP A1 (also known as UP1; middle; brown; PDB 6DCL; residues 7-89 shown) or the C-terminal domain of CRISPR-Cas protein Cse3 (right; beige; PDB 2Y8W; residues 90-211 shown). RNA ligands are colored in orange.

**Figure 4: A conserved loop responsible for RNA binding in the Sm-fold protein Hfq is also conserved in the WYL domain.** (a) The structure of the WYL domain (shown for PafC; residues 485-561) is highly similar to the Sm-fold of the bacterial hexameric RNA chaperone Hfq. (b) Each of the Hfq subunits (beige; PDB 1KQ2; one subunit colored in brown) adopts the Sm-fold. The Hfq ring can bind RNA at three distinct sites (proximal, distal, and rim); here shown with an RNA ligand bound at the proximal site. The RNA ligand of Hfq is colored in orange and is not shown for the top view to visualize the fold. (c) Multiple sequence alignments of Hfq (top) or PafB and PafC protein sequences (bottom) highlight a patch of strongly conserved residues (red boxes) located in the β4/β5-loop of the Sm-fold. Secondary structure elements of *Staphylococcus aureus* (Sau) Hfq (beige; PDB 1KQ2; UniProt Q2FYZ1) and *Arthrobacter aurescens* (Aau) PafBC (blue) are shown below each alignment. Naturally fused PafBC proteins were separated into PafB and PafC parts before alignment (asterisks). Alignment is colored according to percent identity. Aae = *Aquifex aeolicus* (O66512), Lin = *Leptospira interrogans* (Q8F5Z7), Kve = *Koribacter versatilis* (Q1IIF9), Eco = *Escherichia coli* (P0A6X3), Tma = *Thermotoga maritima* (Q9WYZ6), Tfu = *Thermobifida fusca* (Q47P13), Sco = *Streptomyces coelicolor* (Q9RJ64, Q9RJ65), Mtb = *Mycobacterium tuberculosis* (P9WIM1, P9WIL9), Msm = *Mycobacterium smegmatis*

(I7G3U5, A0QZ41), Krh = *Kocuria rhizophila* (B2GIN6), Cgl = *Corynebacterium glutamicum* (Q8NQE2, Q8NQE3). UniProt sequence identifiers in brackets. (d) Two of the highly conserved residues of the  $\beta$ 4/ $\beta$ 5-loop are involved in substrate binding at the proximal face of Hfq (PDB 1KQ2, one subunit shown). Hydrogen bonds are colored in green. (e) The  $\beta$ 4/ $\beta$ 5-loop residues of <sup>Aau</sup>PafBC $\Delta$ NC present an interface for potential ligand binding in a similar manner as Hfq.

**Figure 5: Mutational complementation screen of residues in the WYL domain of PafBC in *Mycobacterium smegmatis*.** (a) The *M. smegmatis*  $\Delta$ pafBC strain ( $\Delta$ pafBC, dashed black line) exhibits a lower viability by an order of magnitude in presence of the DNA-damaging agent mitomycin C (MMC) than the wild-type strain (WT, solid black line), which is restored by expressing wild type PafBC (comBC, green). Expression of PafBC with either the HTH domain of PafB or PafC deleted ( $\Delta$ HTH-B, orange or  $\Delta$ HTH-C, blue) results in the same phenotype as observed for the knockout strain. (b-f) Complementation with alanine substitution mutants in either PafB (green) or PafC (orange) moderately affect viability, while mutations in both proteins (blue) additively reduce the viability, in three cases to knockout levels (c, d, f). Each data point represents the mean of three or more individual experiments. Error bars represent the standard deviation of the mean. (g-h) PafBC variants were expressed from an integrative plasmid in the *M. smegmatis*  $\Delta$ pafBC strain and the expression levels were compared to the knockout ( $\Delta$ pafBC) and wild-type (WT) strains carrying the empty plasmid. RpoB served as loading control. A representative immunoblot of four individual experiments is shown.

**Figure 6: Domain architectures and taxonomic distribution of WYL domain-containing proteins.** (a) Domain architecture classes of WYL domain-containing proteins reveal a tight association with an N-terminal HTH domain. Median and standard deviation of protein length are given next to the domain architecture sketch of each class followed by the number of sequences. Domain architecture sketches are drawn to scale based on the median values of protein length, domain boundaries, and domain length. The scale bar equals 100 amino acids (aa). Domains with dashed line borders were assigned manually. For clarity, architectures with less than 100 sequences are not shown. (b) The taxonomic distribution of all WYL domain proteins shows a prevalent occurrence in *Actinobacteria*. (c) Class A, featuring the WCX domain located C-terminally of the WYL domain, is mainly found in gram-positive bacteria, namely *Actinobacteria* and *Firmicutes*, while (d) Class B, exhibiting only the WYL domain, is mostly found in *Proteobacteria*. The segment radian represents the number of unique species, while the thickness of the segment represents the average number of sequences per species within that taxonomic group. The number of species is given below the class labels with the number of sequences in parentheses. For clarity, taxonomic groups smaller than 1.5% of the total number are not shown. See also Figure S1, S2 and S3.

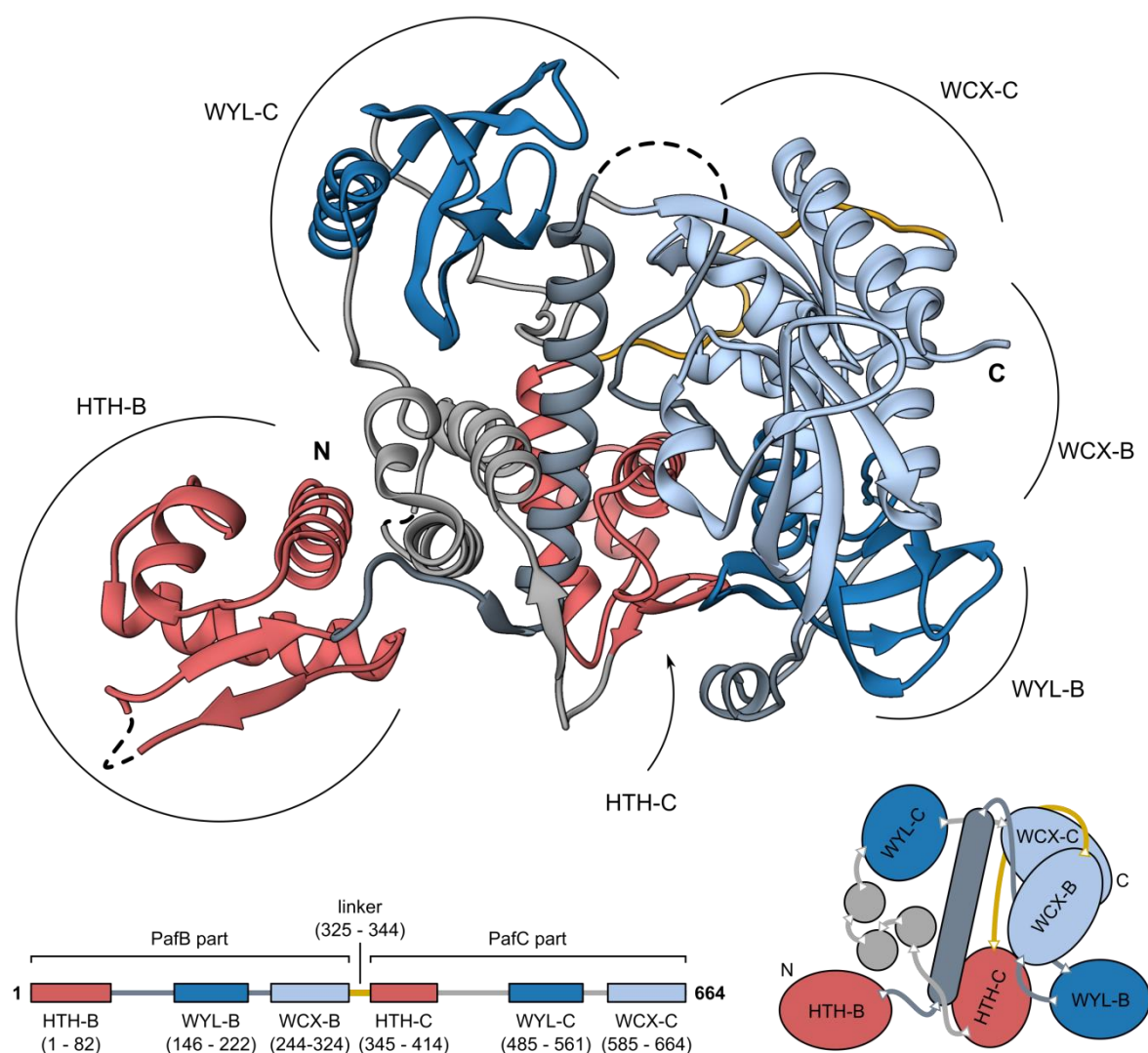
**Figure 7: Hypothetical model of DNA-binding by activated PafBC.** In the non-activated state, PafBC buries the recognition helix (orange) of PafC's HTH domain (HTH-C) and cannot bind to its cognate promoter motif. Upon activator binding, PafBC likely undergoes large structural rearrangements of its domains to release the HTH-C domain, allowing promoter recognition and transcriptional activation of DNA repair genes.

## Tables

**Table 1: Data collection and refinement statistics of the crystal structure of *Arthrobacter aurescens* PafBC (<sup>Aau</sup>PafBCΔNC).**

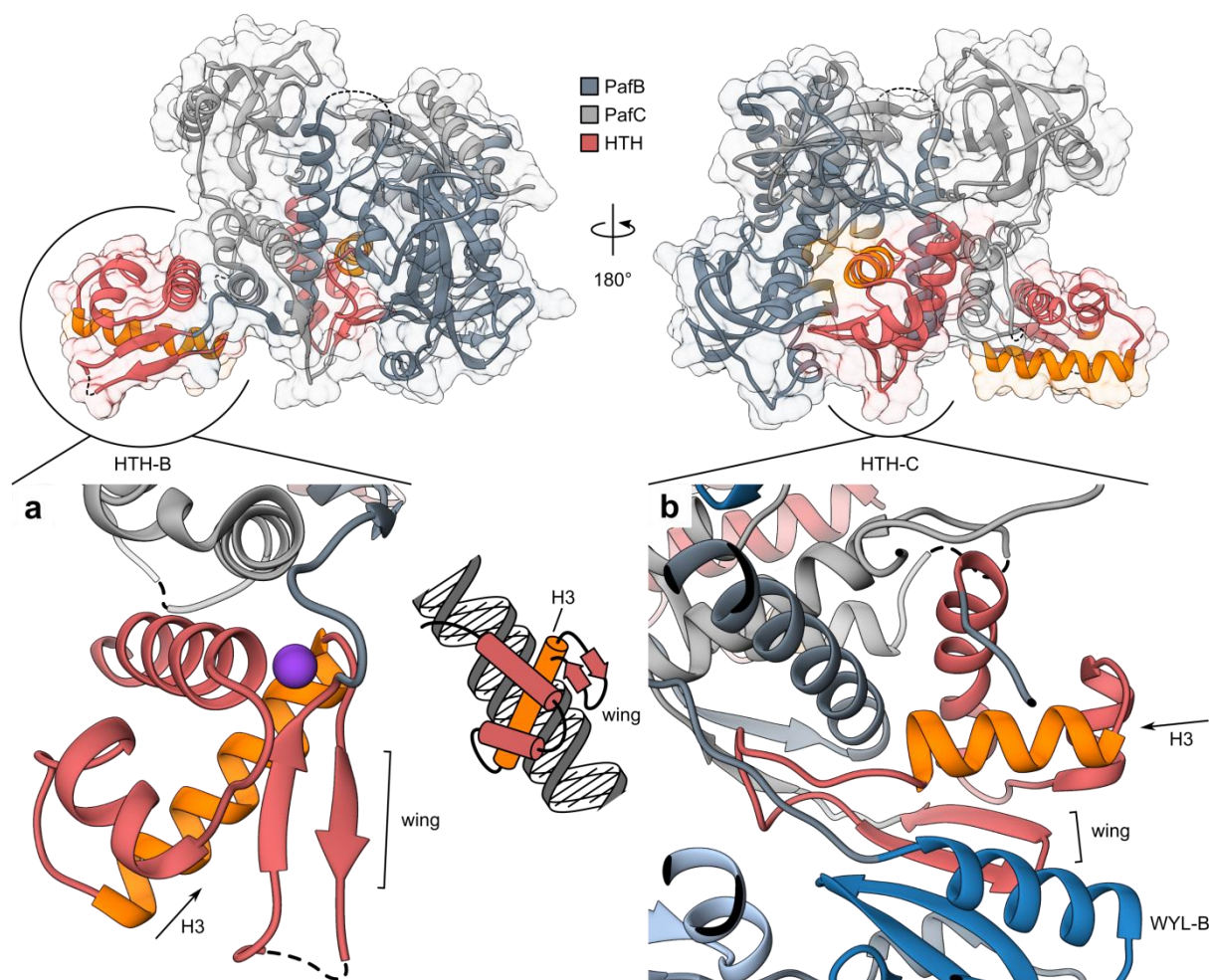
<sup>Aau</sup> PafBCΔNC (SeMet)	
<b>Data Collection</b>	
Wavelength	0.978561 Å (12'670 eV)
Resolution range	48.72 - 2.2 (2.279 - 2.2)
Space group	<i>P</i> 2 <sub>1</sub> 2 <sub>1</sub> 2 <sub>1</sub>
Unit cell dimensions (Å)	77.88, 119.03, 160.21
α, β, γ (°)	90, 90, 90
Total reflections	5'208'738 (520'414)
Unique reflections	76'285 (7'554)
Multiplicity	68.3 (68.9)
Completeness (%)	99.95 (99.96)
Mean I/sigma(I)	31.58 (2.76)
Wilson B-factor	48.04
R-merge	0.1226 (2.069)
R-meas	0.1235 (2.084)
R-pim	0.01484 (0.2494)
CC1/2	1 (0.871)
CC*	1 (0.965)
Matthews coefficient	2.49
Solvent fraction	0.505
Molecules/ASU	2
<b>Refinement</b>	
Reflections used in refinement	76'278 (7'555)
Reflections used for R-free	3'814 (378)
R-work	0.2023 (0.2896)
R-free	0.2336 (0.3364)
CC(work)	0.958 (0.831)
CC(free)	0.942 (0.784)
Number of TLS groups	10
<b>Model statistics</b>	
Number of non-hydrogen atoms	10'252
- in macromolecules	9'936
- in ligands	49
- in solvent	267
Protein residues	1'279
Heavy atom sites (Se)	14
RMS(bonds)	0.012
RMS(angles)	1.46
Ramachandran favored (%)	98.26
Ramachandran allowed (%)	1.66
Ramachandran outliers (%)	0.08
Rotamer outliers (%)	2.13
Clashscore	6.94
Average B-factor	74.09
- macromolecules	74.58
- ligands	88.01
- solvent	53.46



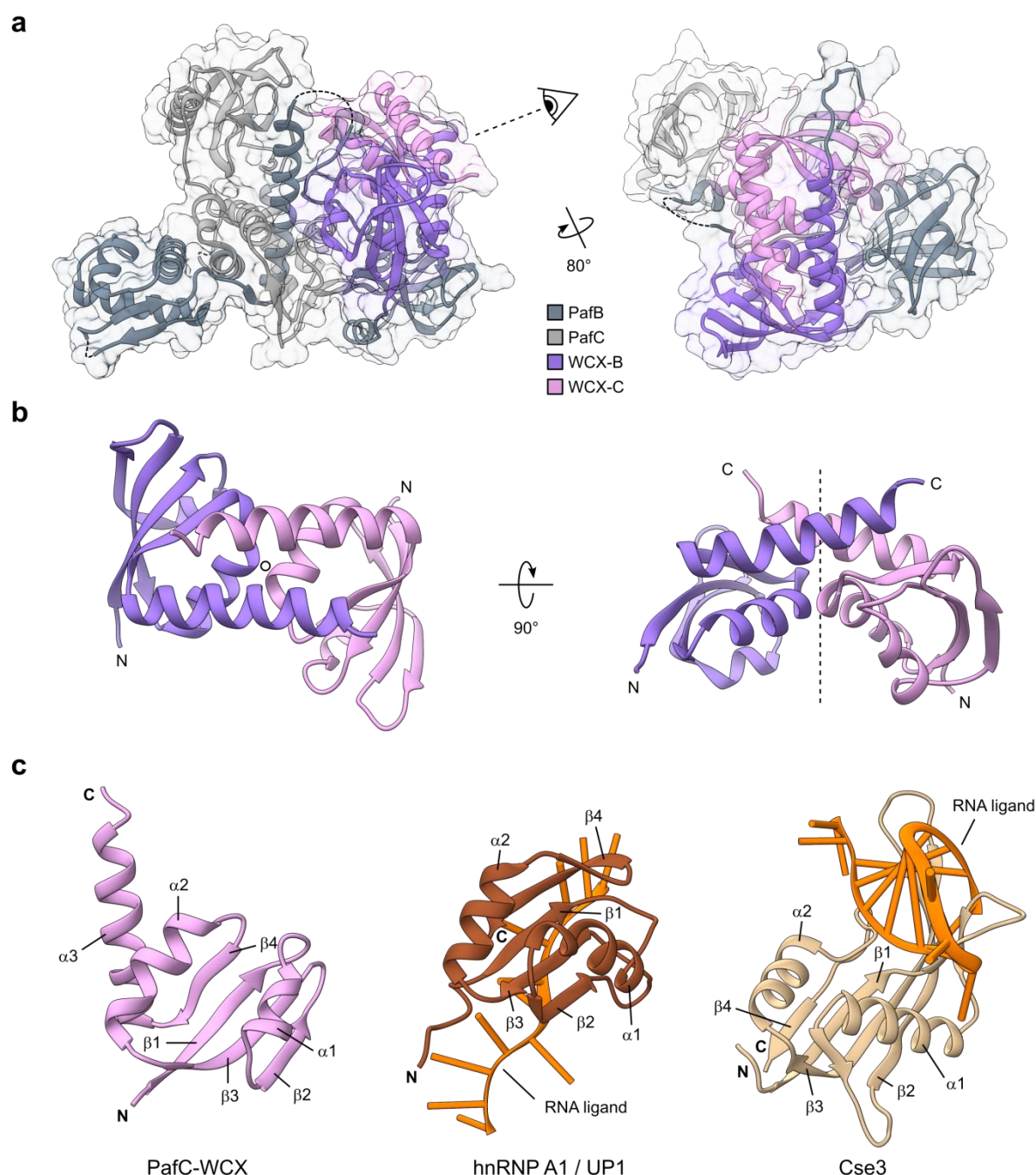


**Figure 1: Crystal structure of the PafBC complex from *Arthrobacter aurescens* ( $^{Au}$ PafBC $\Delta$ NC).** The protein is the product of a natural gene fusion and encompasses a PafB and a PafC part (bottom scheme) connected via a 20 amino acids long linker (yellow). The winged helix-turn-helix domains (HTH) are colored in red, WYL domains are colored in blue and the C-terminal extension of the WYL domain (WCX) is colored in light blue. Other parts of the structure belonging to PafB are colored in dark gray, while the remaining PafC parts are colored in light grey. Dashed lines bridge gaps in the model. See also Figure S1, S2, and S3.



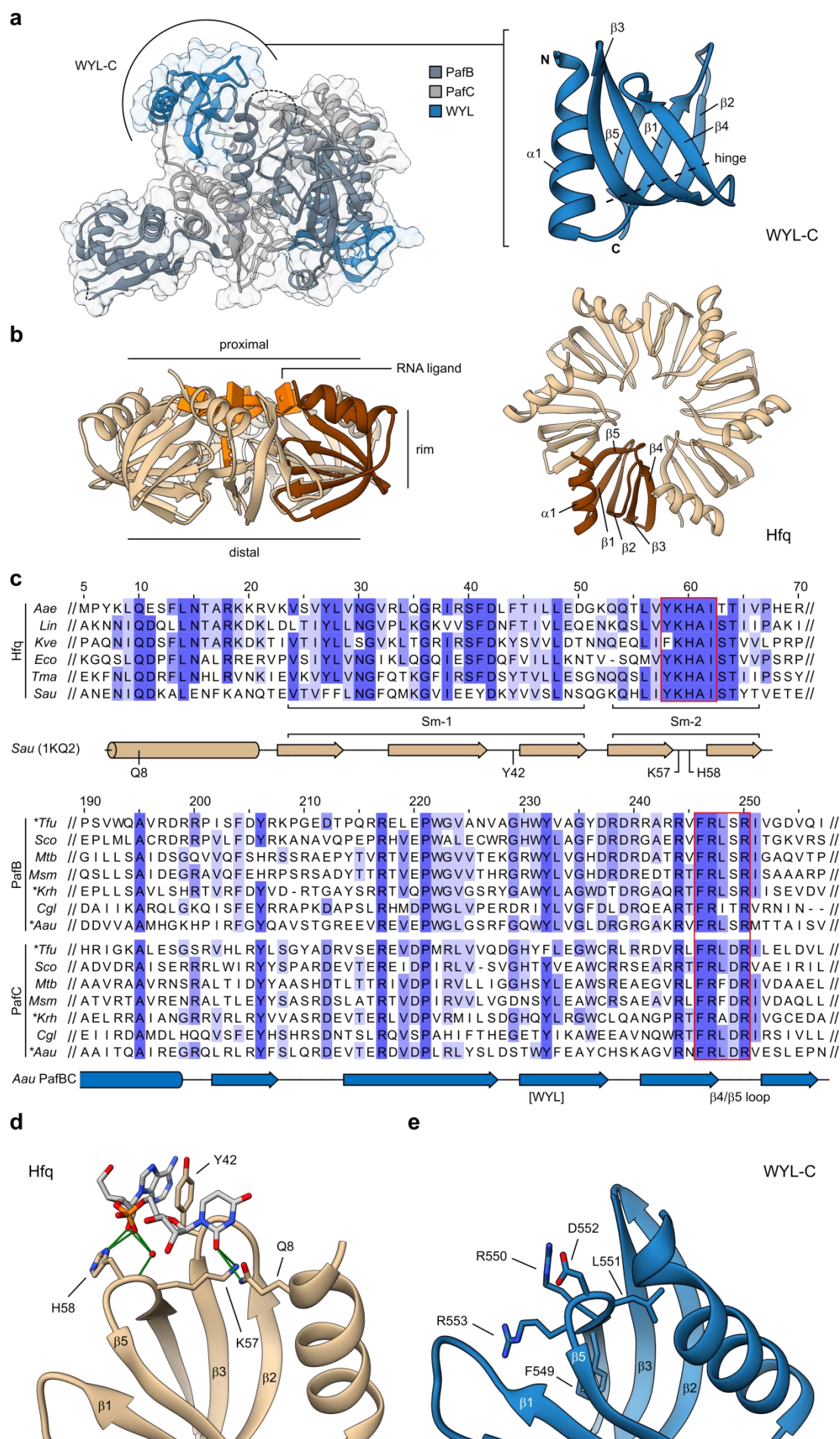


**Figure 2: The non-activated conformation of PafBC buries the HTH domain of PafC.** (a) The HTH domain of PafB (HTH-B, red) adopts a typical winged helix-turn-helix fold (wHTH). The recognition helices (H3) of wHTH domains typically insert into the major groove of the DNA, while the wings establish contact to the minor groove (small inset, middle). The recognition helix of HTH-B (orange) is exposed and available to accept a DNA ligand. The main chain of H1 and the wing coordinate a potassium ion (lilac sphere). (b) The HTH domain of PafC (HTH-C, red) on the other hand forms part of the protein core. H3 of HTH-C (orange) is shorter by two helical turns compared to H3 of HTH-B, and an unstructured loop reaches into the protein core. Additionally, the  $\beta$ -sheet of the wing extends the  $\beta$ -sheet of the WYL domain of PafB (WYL-B, blue). Dashed lines bridge gaps in the models. See also Figure S4.

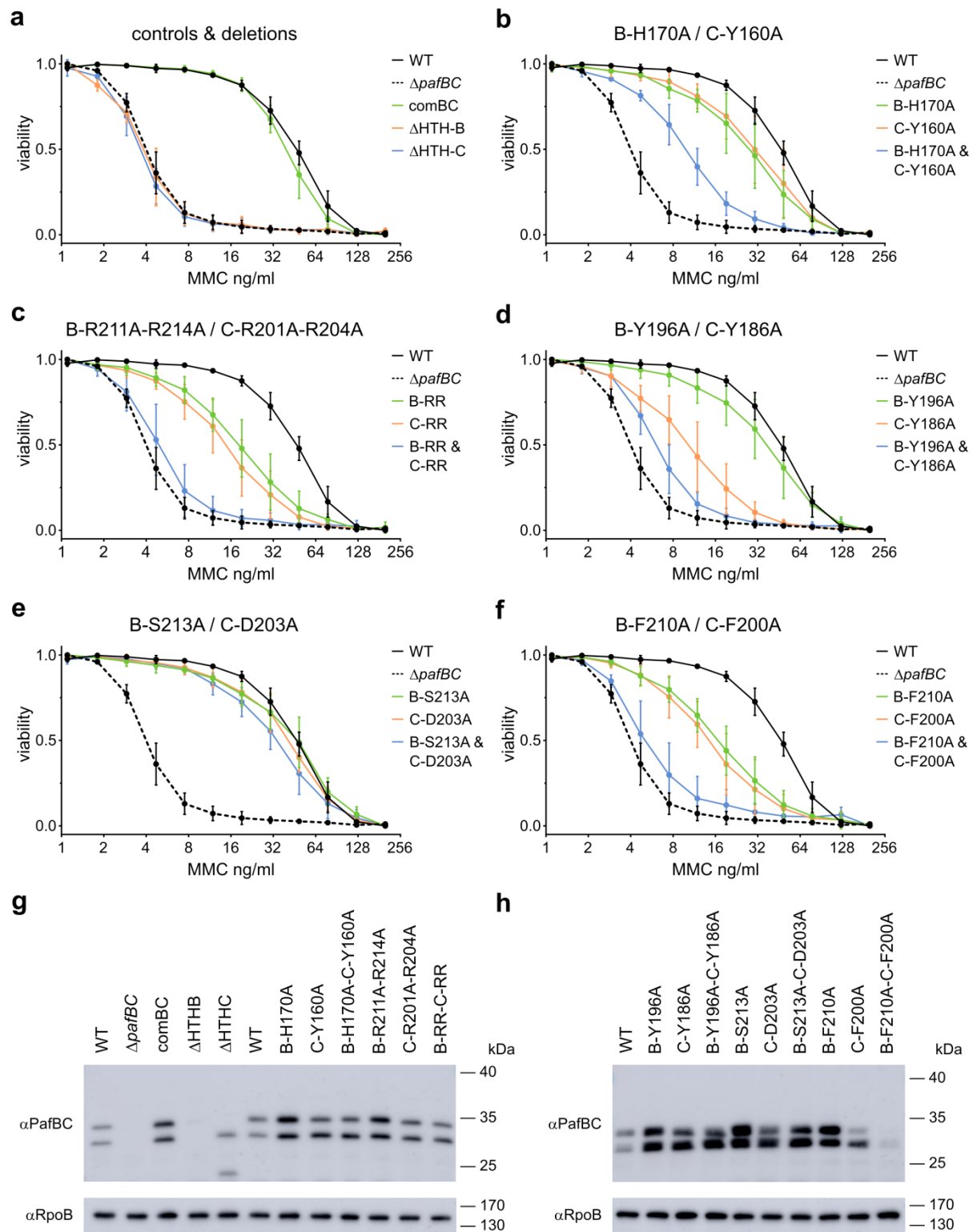


**Figure 3: The C-terminal extension (WCX) domains of PafB and PafC contain a ferredoxin-like fold and display two-fold rotational pseudosymmetry.** (a) The WCX domains of PafB (WCX-B, lilac) and PafC (WCX-C, pink) contact each other in the crystal structure. Dashed lines bridge gaps in the model. (b) WCX-B and WCX-C exhibit a two-fold rotational symmetry axis (circle and dashed line). (c) The WCX domains (left; shown for PafC; residues 585-664) contain a ferredoxin-like fold with an additional C-terminal  $\alpha$ -helix ( $\alpha 3$ ). Very versatile and present in proteins with highly diverse functions, the ferredoxin-like fold is also found in many RNA-binding proteins such as human hnRNP A1 (also known as UP1; middle; brown; PDB 6DCL; residues 7-89 shown) or the C-terminal domain of CRISPR-Cas protein Cse3 (right; beige; PDB 2Y8W; residues 90-211 shown). RNA ligands are colored in orange.





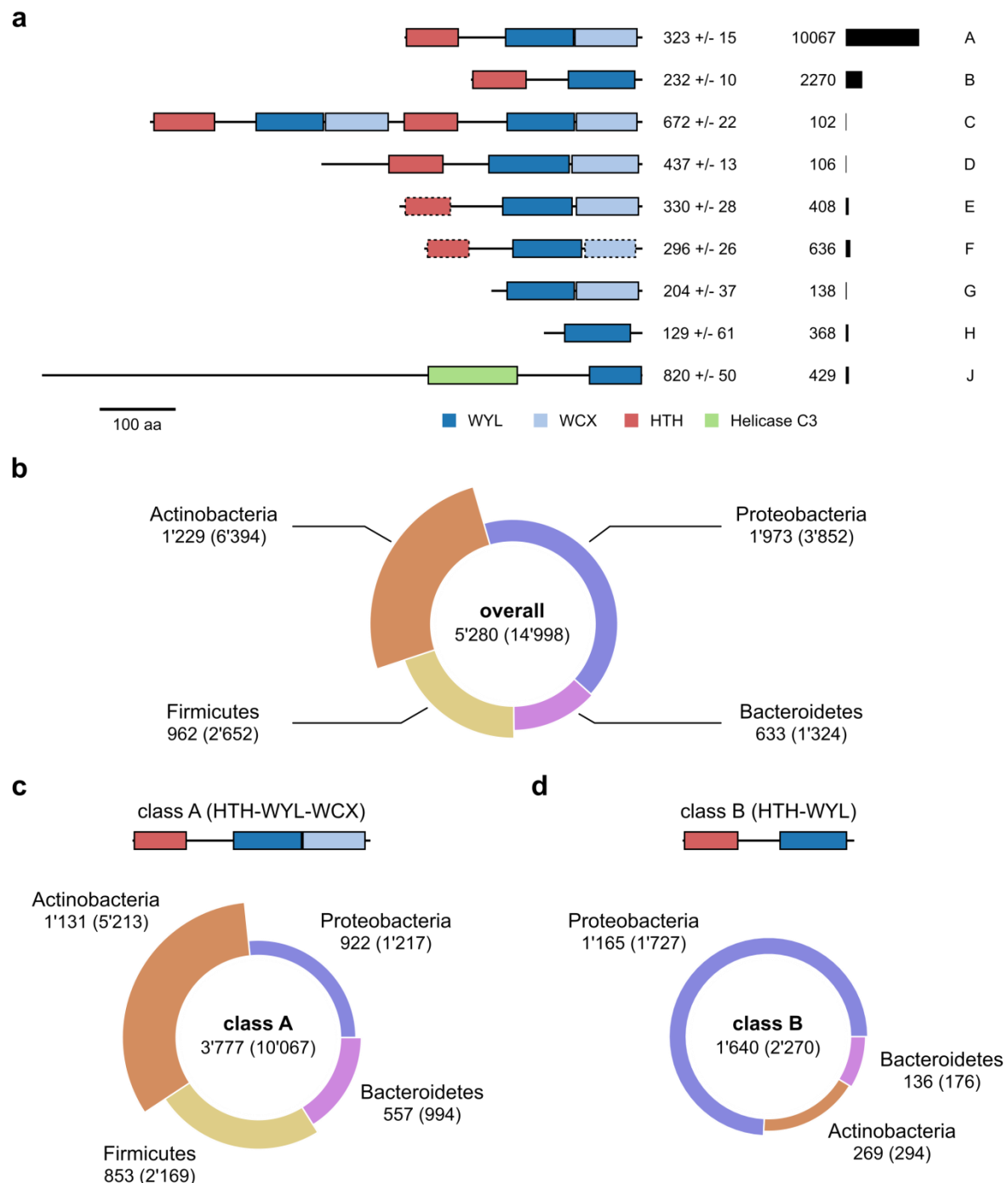
**Figure 4: A conserved loop responsible for RNA binding in the Sm-fold protein Hfq is also conserved in the WYL domain.** (a) The structure of the WYL domain (shown for PafC; residues 485-561) is highly similar to the Sm-fold of the bacterial hexameric RNA chaperone Hfq. (b) Each of the Hfq subunits (beige; PDB 1KQ2; one subunit colored in brown) adopts the Sm-fold. The Hfq ring can bind RNA at three distinct sites (proximal, distal, and rim); here shown with an RNA ligand bound at the proximal site. The RNA ligand of Hfq is colored in orange and is not shown for the top view to visualize the fold. (c) Multiple sequence alignments of Hfq (top) or PafB and PafC protein sequences (bottom) highlight a patch of strongly conserved residues (red boxes) located in the  $\beta 4/\beta 5$ -loop of the Sm-fold. Secondary structure elements of *Staphylococcus aureus* (Sau) Hfq (beige; PDB 1KQ2; UniProt Q2FYZ1) and *Arthrobacter aurescens* (Aau) PafBC (blue) are shown below each alignment. Naturally fused PafBC proteins were separated into PafB and PafC parts before alignment (asterisks). Alignment is colored according to percent identity. Aae = *Aquifex aeolicus* (O66512), Lin = *Leptospira interrogans* (Q8F5Z7), Kve = *Koribacter versatilis* (Q1IIF9), Eco = *Escherichia coli* (P0A6X3), Tma = *Thermotoga maritima* (Q9WYZ6), Tfu = *Thermobifida fusca* (Q47P13), Sco = *Streptomyces coelicolor* (Q9RJ64, Q9RJ65), Mtb = *Mycobacterium tuberculosis* (P9WIM1, P9WIL9), Msm = *Mycobacterium smegmatis* (I7G3U5, A0QZ41), Krh = *Kocuria rhizophila* (B2GIN6), Cgl = *Corynebacterium glutamicum* (Q8NQE2, Q8NQE3). UniProt sequence identifiers in brackets. (d) Two of the highly conserved residues of the  $\beta 4/\beta 5$ -loop are involved in substrate binding at the proximal face of Hfq (PDB 1KQ2, one subunit shown). Hydrogen bonds are colored in green. (e) The  $\beta 4/\beta 5$ -loop residues of <sup>Aau</sup>PafBC $\Delta$ NC present an interface for potential ligand binding in a similar manner as Hfq.



**Figure 5: Mutational complementation screen of residues in the WYL domain of PafBC in *Mycobacterium smegmatis*.** (a) The *M. smegmatis*  $\Delta pafBC$  strain ( $\Delta pafBC$ , dashed black line) exhibits a lower viability by an order of magnitude in presence of the DNA-damaging agent mitomycin C (MMC) than the wild-type strain (WT, solid black line), which is restored by expressing wild type PafBC (comBC, green). Expression of PafBC with either the HTH domain of PafB or PafC deleted ( $\Delta HTH-B$ , orange or  $\Delta HTH-C$ , blue) results in the same phenotype as observed for the knockout strain. (b-f) Complementation with alanine substitution mutants in either PafB (green) or PafC (orange) moderately affect viability, while mutations in both proteins (blue) additively reduce the

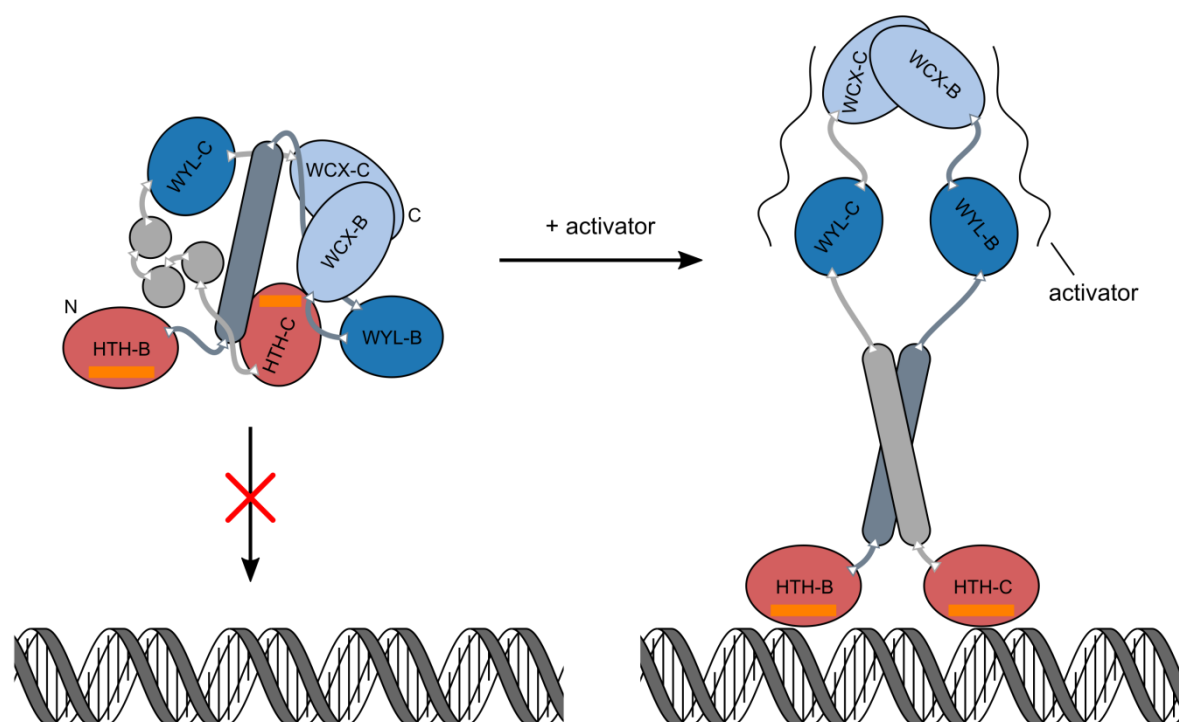
viability, in three cases to knockout levels (c, d, f). Each data point represents the mean of three or more individual experiments. Error bars represent the standard deviation of the mean. (g-h) PafBC variants were expressed from an integrative plasmid in the *M. smegmatis*  $\Delta pafBC$  strain and the expression levels were compared to the knockout ( $\Delta pafBC$ ) and wild-type (WT) strains carrying the empty plasmid. RpoB served as loading control. A representative immunoblot of four individual experiments is shown.



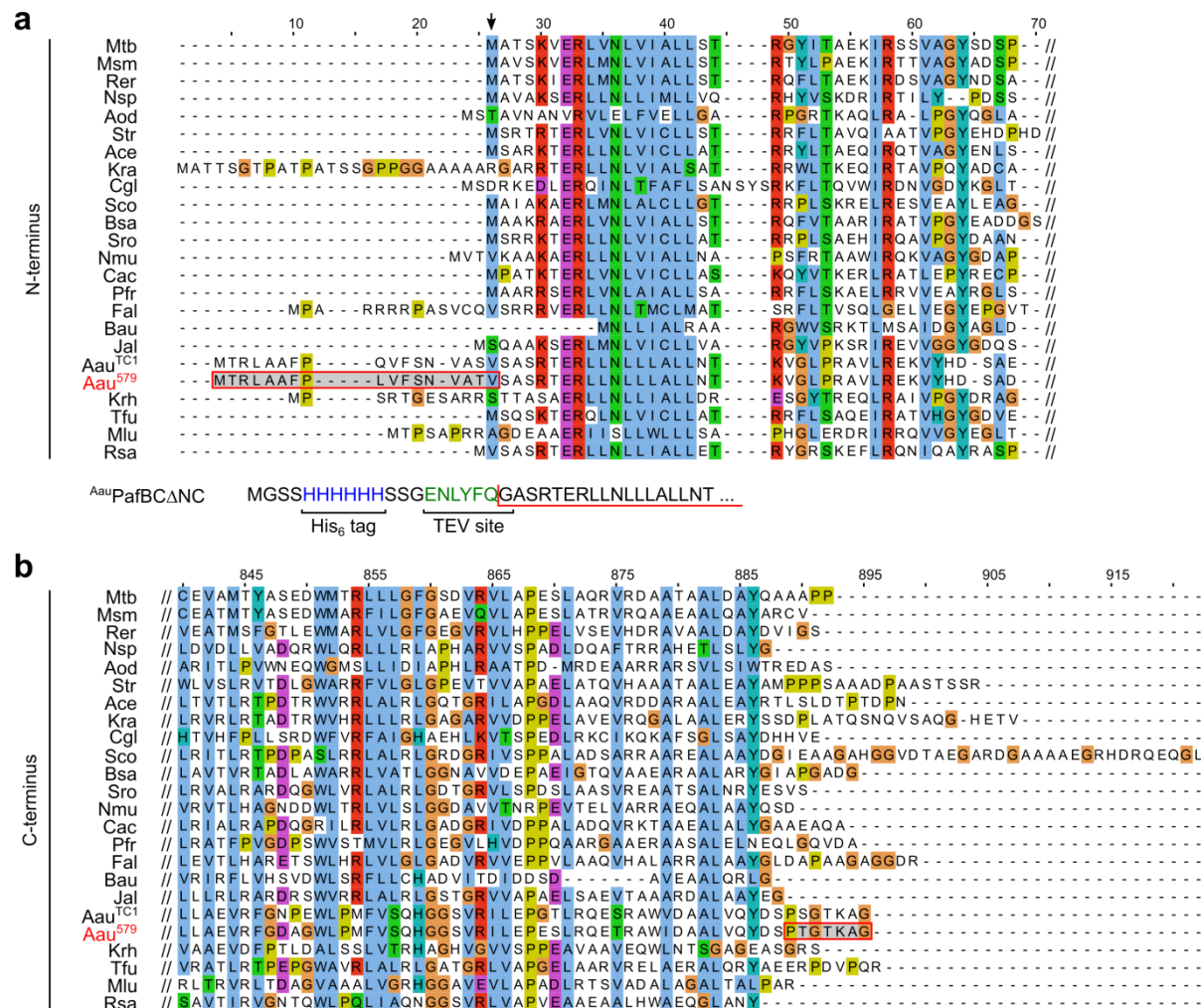


**Figure 6: Domain architectures and taxonomic distribution of WYL domain-containing proteins.** (a) Domain architecture classes of WYL domain-containing proteins reveal a tight association with an N-terminal HTH domain. Median and standard deviation of protein length are given next to the domain architecture sketch of each class followed by the number of sequences. Domain architecture sketches are drawn to scale based on the median values of protein length, domain boundaries, and domain length. The scale bar equals 100 amino acids (aa). Domains with dashed line borders were assigned manually. For clarity, architectures with less than 100 sequences are not shown. (b) The taxonomic distribution of all WYL domain proteins shows a prevalent occurrence in *Actinobacteria*. (c) Class A, featuring the WCX domain located C-terminally of the WYL domain, is mainly found in gram-positive bacteria, namely *Actinobacteria* and *Firmicutes*, while (d) Class B, exhibiting only the WYL domain, is mostly found in *Proteobacteria*. The segment radian represents the number of

unique species, while the thickness of the segment represents the average number of sequences per species within that taxonomic group. The number of species is given below the class labels with the number of sequences in parentheses. For clarity, taxonomic groups smaller than 1.5% of the total number are not shown. See also Figure S1, S2 and S3.



**Figure 7: Hypothetical model of DNA-binding by activated PafBC.** In the non-activated state, PafBC buries the recognition helix (orange) of PafC's HTH domain (HTH-C) and cannot bind to its cognate promoter motif. Upon activator binding, PafBC likely undergoes large structural rearrangements of its domains to release the HTH-C domain, allowing promoter recognition and transcriptional activation of DNA repair genes.

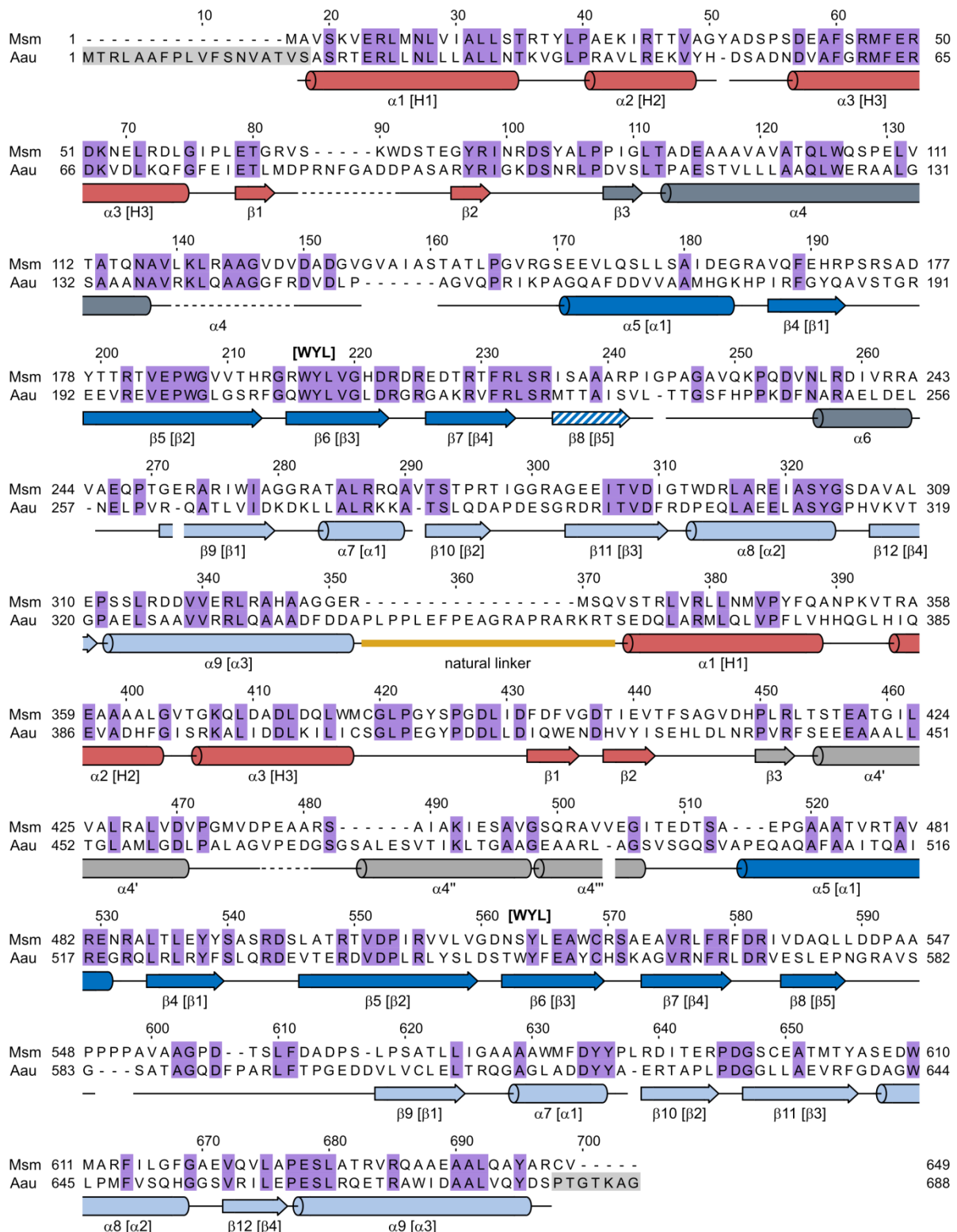


**Figure S1: Sequence alignment of the N- and C-terminal regions of PafBC orthologs.** (a) The N-terminus of <sup>Aau</sup>PafBC from strain 579 (red, Aau<sup>579</sup>) is likely a misannotation in the database, because most other PafB proteins contain a conserved initiator methionine at alignment position 26 (arrow at the top) and the valine in *A. aureus* may be due to usage of an alternative start codon, which is common among *Actinobacteria*. For crystallization experiments, the additional 17 residues of Aau<sup>579</sup> from the N-terminus were removed (including the theoretical initiator methionine, red box with gray background), and the first serine was replaced with glycine to generate a TEV protease cleavage site (see also experimental procedures). The N-terminal sequence of the construct expressed and used for crystallization (<sup>Aau</sup>PafBCΔNC) is shown below. The red line marks the N-terminus of the purified protein. (b) The C-terminal residues of Aau<sup>579</sup> are not conserved among other actinobacterial PafC proteins, and thus, 7 C-terminal residues are not contained in <sup>Aau</sup>PafBCΔNC, which was used for crystallization. Residues are colored according to the ClustalX color scheme.

Mtb = *Mycobacterium tuberculosis* (P9WIM1, P9WIL9), Msm = *Mycobacterium smegmatis* (I7G3U5, A0QZ41), Rer = *Rhodococcus erythropolis* (C0ZZU3, C0ZZU2), Nsp = *Nocardioideis* sp. (A1SK18, A1SK19), Aod = *Actinomyces odontolyticus* (A7BCC5, A7BCC6), Str = *Salinispora tropica* (A4X749, A4X750), Ace = *Acidothermus cellulolyticus* (A0LU62, A0LU63), Kra = *Kineococcus radiotolerans* (A6W976, A6W977), Cgl = *Corynebacterium glutamicum* (Q8NQE2, Q8NQE3), Sco = *Streptomyces coelicolor* (Q9RJ64, Q9RJ65), Bsa = *Blastococcus saxobidens* (H6RJ02, H6RJ01), Sro = *Streptosporangium roseum* (D2ATU2, D2ATU1), Nmu = *Nakamurella multipartita* (C8XAP4, C8XAP3), Cac = *Catenulisporea acidiphila* (C7PVW0, C7PVW1), Pfr = *Propionibacterium freudenreichii*

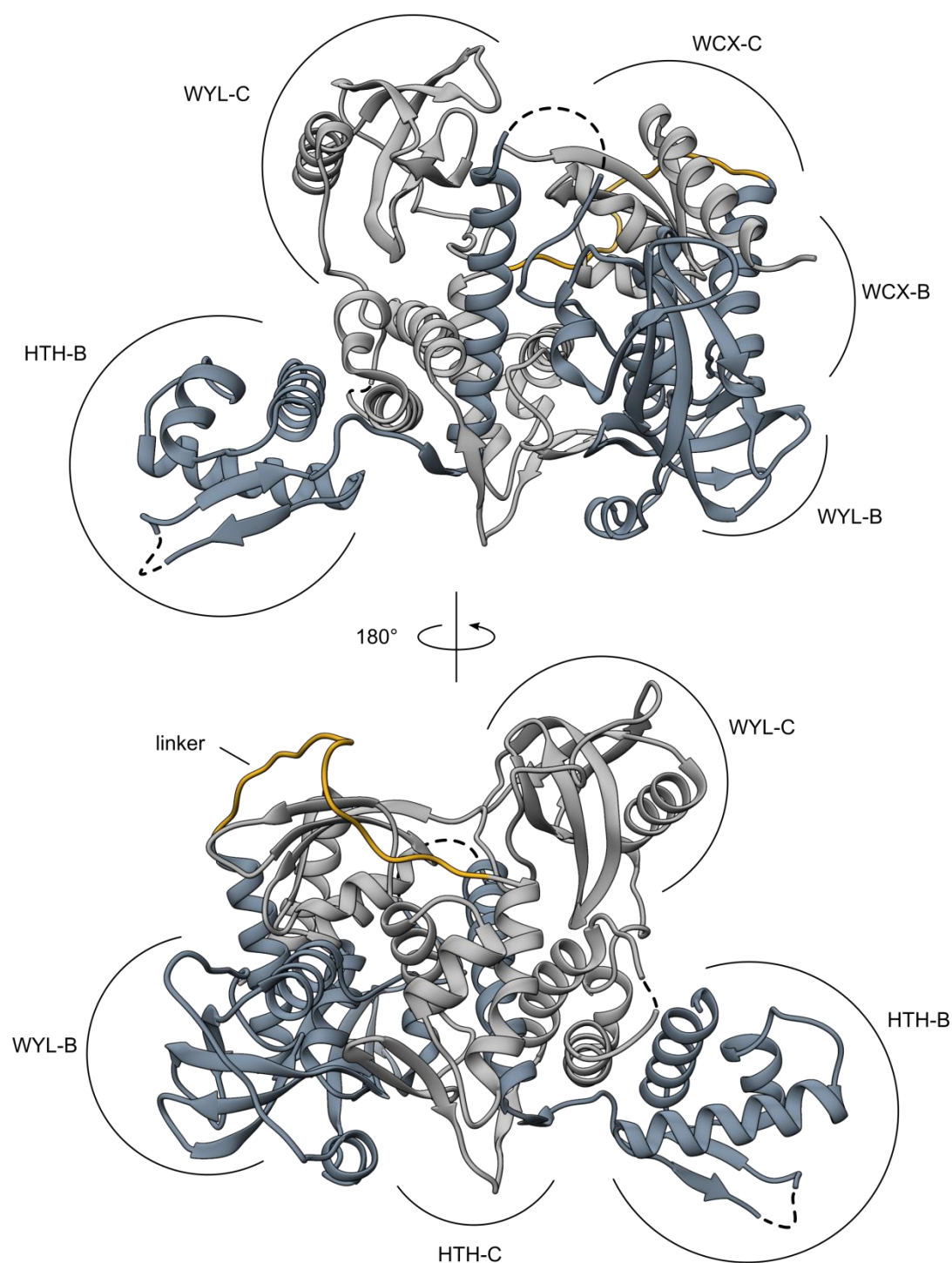
(A0A160VN40, A0A161KHT8), Fal = *Frankia alni* (Q0RLT0, Q0RLS9), Bau = *Brevibacterium aurantiacum* (A0A1D7W444, A0A1D7W495), Jal = *Jiangella alkaliphila* (A0A1H2KTF9, A0A1H2KTV3), Aau<sup>TC1</sup> = *Arthrobacter aurescens* strain TC1 (A1R6R2), Krh = *Kocuria rhizophila* (B2GIN6), Tfu = *Thermobifida fusca* (Q47P13), Mlu = *Micrococcus luteus* (C5CBV3), Rsa = *Renibacterium salmoninarum* (A9WSH6)

UniProt accession numbers are given in parentheses. Refers to Figure 1 and experimental procedures.

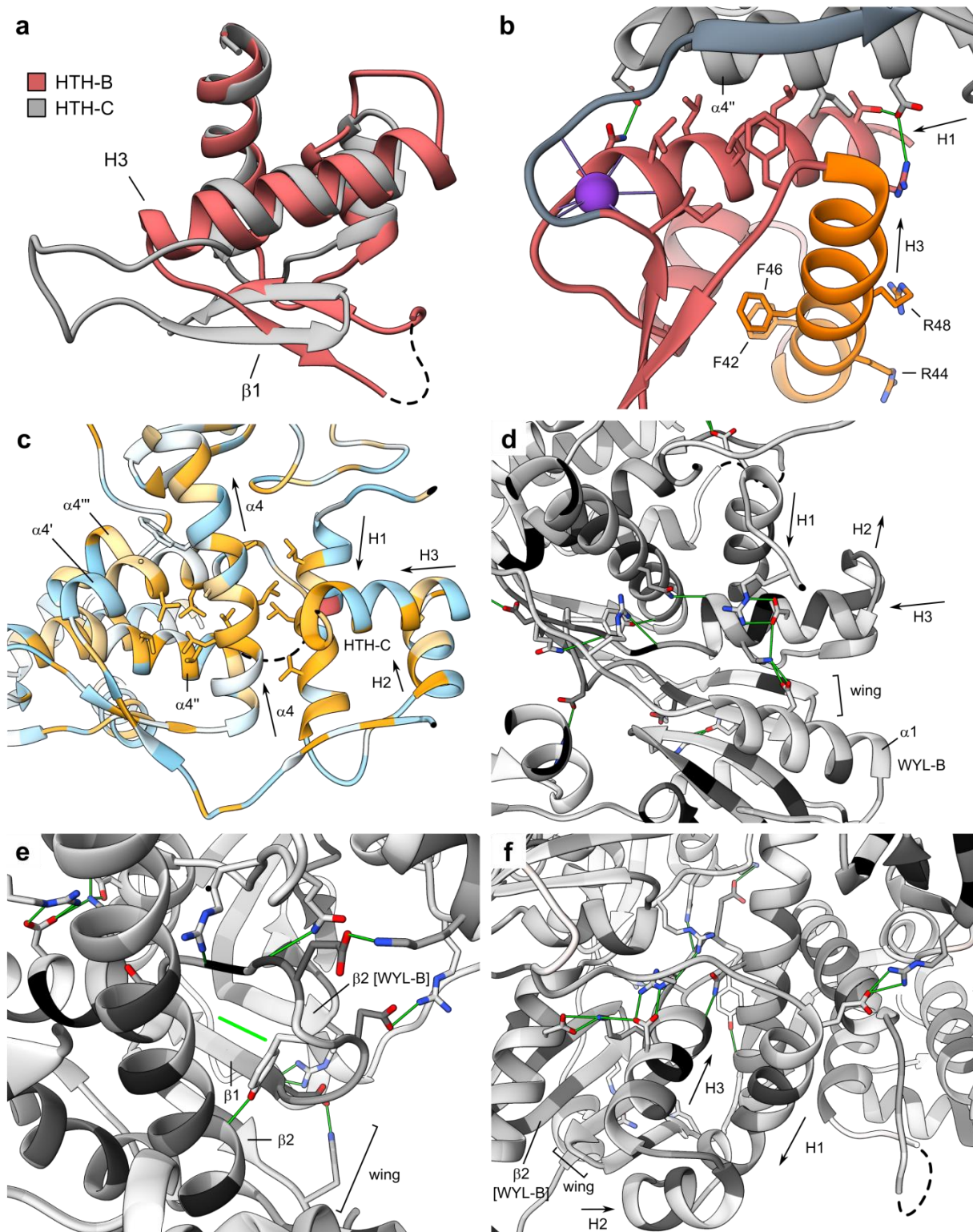




in both organisms are colored in violet. Gray residues were not part of the protein used for crystallization. Refers to Figure 1.



**Figure S3: Extended overview of the *AauPafBCΔNC* crystal structure.** The PafB part is colored in dark gray and the PafC part in light gray. The natural linker present in *AauPafBC* is colored in yellow. Refers to Figure 1.



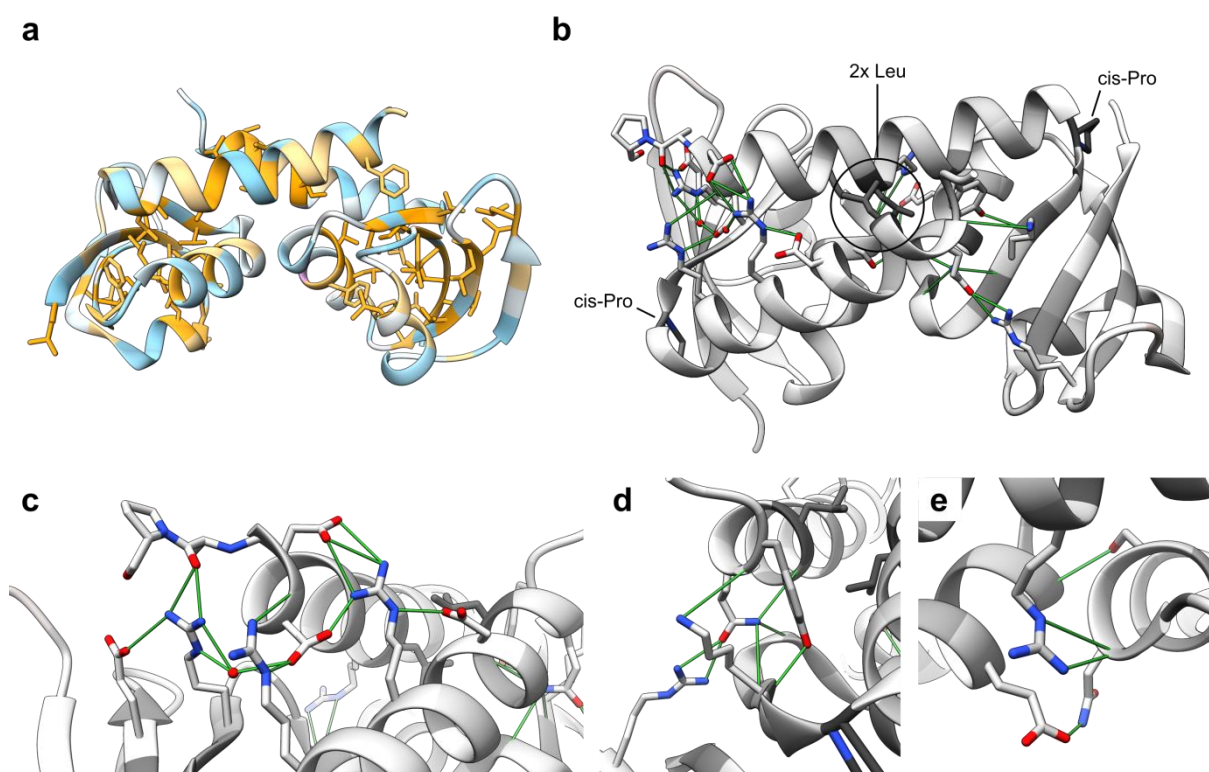
**Figure S4: HTH domains of *Aau*PafBCΔNC interact with other parts of the molecule to a very different degree.** (a) The HTH domain of PafC (HTH-C, gray) has a shorter H3 and a longer loop connecting to the β1 strand of the wing compared to the HTH domain of PafB (HTH-B, red), while the β1/β2 loop is much shorter in HTH-C. (b) H1 of HTH-B interacts with helix α4'' (part of the helix bundle in PafC) almost exclusively through hydrophobic interactions (residues in stick representation). The recognition helix of HTH-B (H3; orange) contains two highly conserved phenylalanines involved in forming the hydrophobic core of the domain and two highly conserved, exposed arginines that could make specific base contacts in the DNA-bound form of PafBC. A

potassium ion (lilac) is caged in between the main chain of H1 and the wing. (c) The central helix  $\alpha 4$  (part of PafB) makes hydrophobic interactions with H1 of PafC and  $\alpha 4''$  of PafC. Residues were colored based on their hydrophobicity according to the Kyte-Doolittle scale from orange (hydrophobic) to cyan (polar). Selenomethionines are colored in light red. Leucine, valine and phenylalanine are shown in stick representation. (d-f) Additionally, HTH-C interacts with the protein core through a variety of hydrogen bonds, mostly involving the main chain and non-conserved residues. The hydrogen bonds stabilizing the extended  $\beta$ -sheet formed by the HTH-C wing and WYL-B are collectively depicted as single green line in panel e. Conservation is colored from black (conserved) to white (no conservation) based on the BLOSUM-62 matrix. Dashed lines depict gaps in the model. Refers to Figure 2.

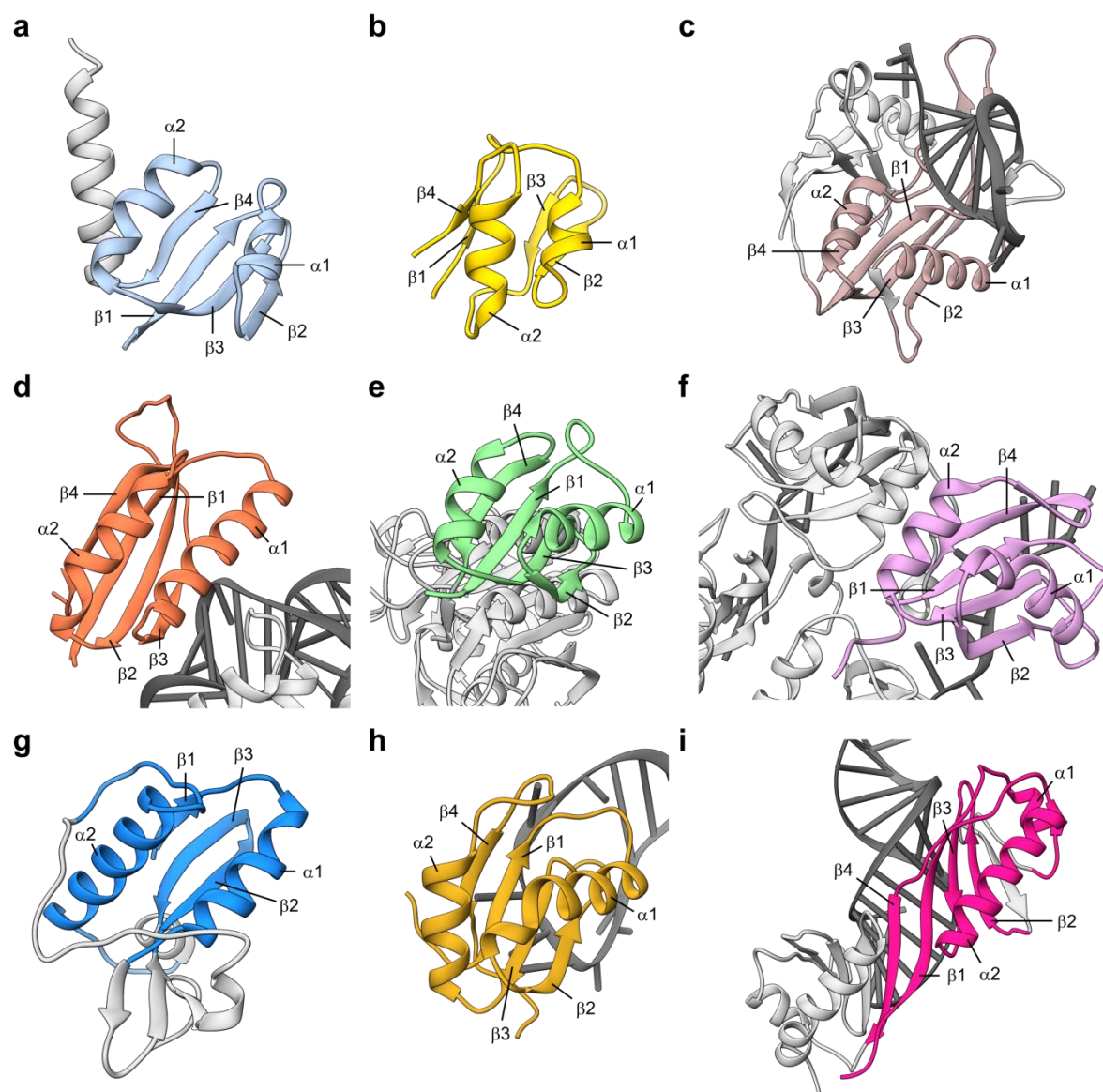




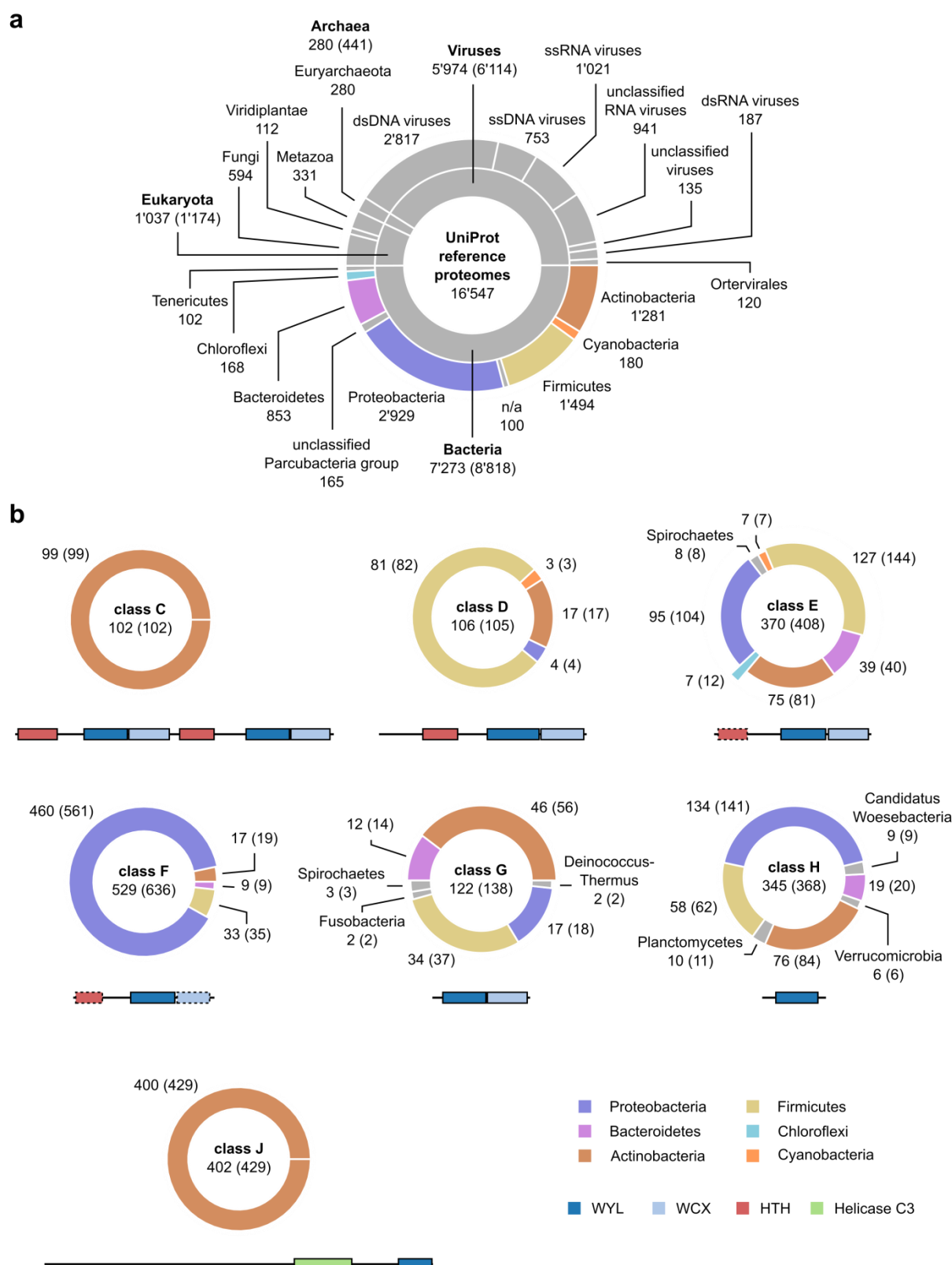




**Figure S6: The C-terminal extension (WCX) domains of PafBC form an interaction module.** (a) Hydrophobicity coloring reveals a hydrophobic core in the ferredoxin-like fold of each WCX domain. Hydrophobicity is colored from orange (hydrophobic) to cyan (polar) according to the Kyte-Doolittle scale. Leu, Val, Phe are shown in stick representation. (c-e) The WCX domains of PafBC contain a network of hydrogen bonds and salt bridges (green) around two highly conserved leucines located in a small hydrophobic island at the C-terminal helix. The residue conservation is colored from black (conserved) to white (no conservation) based on the BLOSUM-62 matrix. Refers to Figure 3.



**Figure S7: Structural comparison of selected proteins containing a ferredoxin-like fold.** (a) *Arthrobacter aureus* PafBC WCX-C domain (b) *Thermotoga maritima* ferredoxin (PDB 1VJW) (c) *Thermus thermophilus* Cse3 bound to RNA substrate (PDB 2Y8W) (d) Ribosomal protein S6 from *T. thermophilus* bound to rRNA (PDB 1G1X) (e) *Salmonella typhimurium* subtilisin (PDB 1SBP) (f) Human hnRNP A1 protein UP1 bound to an RNA ligand (PDB 6DCL) (g) T4 phage translational regulator protein RegA (PDB 1REG) (h) *Drosophila melanogaster* protein U1A/SNF bound to RNA (PDB 6F4H) (i) Yeast TATA-binding protein (PDB 1YTB). Only the ferredoxin-like fold is shown in color, while other parts are shown in light gray. Secondary structure elements in addition to the classical βαβαβ topology of the ferredoxin-like fold are colored in light gray as well. Nucleic acid ligands are colored in dark gray. Refers to Figure 3.



**Figure S8: Taxonomic distributions of the UniProt reference proteomes and domain architecture classes of proteins containing the WYL domain.** (a) Taxonomic distribution within UniProt reference proteomes. The six largest taxonomic groups of the bacteria superkingdom were color coded. Groups smaller than 5% of the total number of reference proteomes were omitted for better visualization. (b) Taxonomic distributions among domain architecture classes show that class C, D, F, and J are

largely specific to one of the bacterial superkingdoms. Groups smaller than 1.5% of the total number of reference proteomes were omitted for better visualization. Refers to Figure 6.