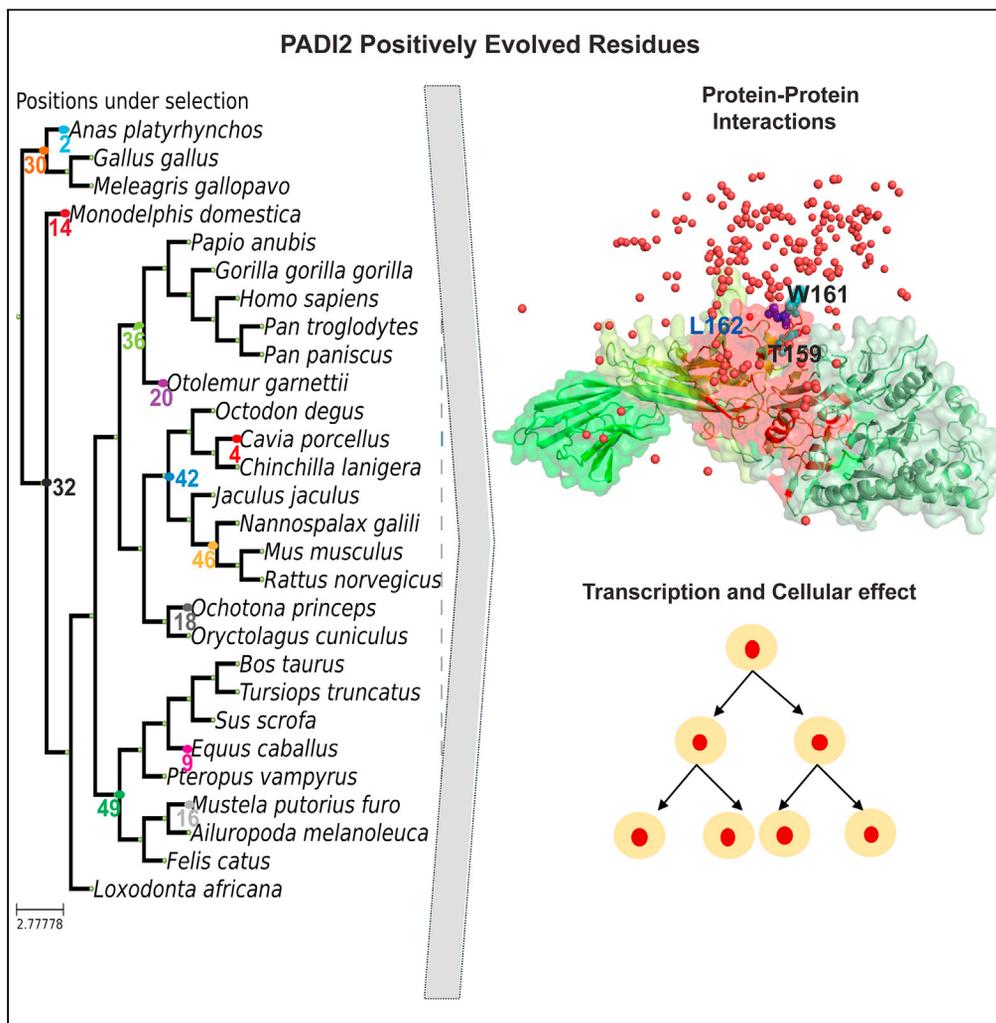


Article

# Evolutionary analysis reveals the role of a non-catalytic domain of peptidyl arginine deiminase 2 in transcriptional regulation



José Luis Villanueva-Cañas, Narcis Fernandez-Fuentes, Dominik Saul, ..., Cedric Notredame, Miguel Beato, Priyanka Sharma

priyanka.sharma@ipbs.fr

**Highlights**

Positively evolved residues are mainly in the non-catalytic domain of the PADI2

L162 of PADI2 is a positively evolved residue at the structurally exposed loop

PADI2-L162 mediates PADI2's effective interaction with the P-TEFb complex

PADI2-L162 dictates cell proliferation and c-MYC transcription

Villanueva-Cañas et al., iScience 27, 109584 April 19, 2024 Crown Copyright © 2024 Published by Elsevier Inc. <https://doi.org/10.1016/j.isci.2024.109584>



## Article

## Evolutionary analysis reveals the role of a non-catalytic domain of peptidyl arginine deiminase 2 in transcriptional regulation

José Luis Villanueva-Cañas,<sup>1,11,12</sup> Narcis Fernandez-Fuentes,<sup>2,12</sup> Dominik Saul,<sup>3,4</sup> Robyn Laura Kosinsky,<sup>5</sup> Catherine Teyssier,<sup>6</sup> Malgorzata Ewa Rogalska,<sup>1</sup> Ferran Pegenaute Pérez,<sup>7,8</sup> Baldomero Oliva,<sup>8,9</sup> Cedric Notredame,<sup>1,8</sup> Miguel Beato,<sup>1,8</sup> and Priyanka Sharma<sup>10,13,\*</sup>

## SUMMARY

**Peptidyl arginine deiminases (PADIs) catalyze protein citrullination, a post-translational conversion of arginine to citrulline. The most widely expressed member of this family, PADI2, regulates cellular processes that impact several diseases. We hypothesized that we could gain new insights into PADI2 function through a systematic evolutionary and structural analysis. Here, we identify 20 positively selected PADI2 residues, 16 of which are structurally exposed and maintain PADI2 interactions with cognate proteins. Many of these selected residues reside in non-catalytic regions of PADI2. We validate the importance of a prominent loop in the middle domain that encompasses PADI2 L162, a residue under positive selection. This site is essential for interaction with the transcription elongation factor (P-TEFb) and mediates the active transcription of the oncogenes *c-MYC*, and *CCNB1*, as well as impacting cellular proliferation. These insights could be key to understanding and addressing the role of the PADI2 *c-MYC* axis in cancer progression.**

## INTRODUCTION

Members of the peptidyl arginine deiminase (PADI) family catalyze the post-translational calcium-dependent citrullination of arginines, converting them into the non-coded amino acid citrulline.<sup>1–4</sup> Arginine citrullination is a widespread post-translational modification (PTM) that increases the hydrophobicity of proteins and can contribute to fine-tuning physiological processes. For instance, the citrullination of core histones weakens histone–nucleic acid interactions, impacting both chromatin organization and transcription.<sup>5</sup> Citrullination can also affect protein folding and consequently protein function.<sup>6–9</sup> The five members of the PADI family have distinct tissue-specific expression patterns: PADI1 is present only in the epidermis and uterus; PADI3, in the epidermis and hair follicles; PADI4, in immune cells, brain, uterus, and bone marrow; PADI6, in ovarian egg cells, embryonic tissues, and testicles; and PADI2, in brain, uterus, spleen, breast, pancreas, skin, and skeletal muscles.<sup>2,3</sup> Notably, PADI2 has been associated with multiple pathological states, including autoimmune disorders<sup>10</sup> and neurological diseases,<sup>11,12</sup> as well as with several cancers (e.g., breast, cervical, liver, lung, ovarian, thyroid, gastric, and prostate cancer).<sup>13–22</sup> The link between PADI2 and various diseases underscores an unmet need to elucidate its molecular mechanisms of action and to understand its pathophysiological function.

PADI family members have distinct substrates and target diverse arginine residues in different proteins.<sup>23</sup> All family members are Ca<sup>2+</sup>-dependent, but PADI2 also relies on ordered calcium binding to its active site for substrate binding and catalysis, as demonstrated in the deimination of arginine 26 in histone H3 (H3R26).<sup>24–27</sup> Importantly, H3R26 citrullination is critical for interaction with the chromatin remodeller SMARCD1 (SWI/SNF-related matrix-associated actin-dependent regulator of chromatin subfamily A containing DEAD/H box 1) and hence

<sup>1</sup>Centre for Genomic Regulation (CRG), The Barcelona Institute of Science and Technology (BIST), Dr. Aiguader 88, 08003 Barcelona, Spain

<sup>2</sup>Institute of Biological, Environmental and Rural Sciences, Aberystwyth University, Aberystwyth, Ceredigion, United Kingdom

<sup>3</sup>Division of Endocrinology, Mayo Clinic, Rochester, MN 55905, USA; Robert and Arlene Kogod Center on Aging, Mayo Clinic, Rochester, MN 55905, USA

<sup>4</sup>Department of Trauma and Reconstructive Surgery, BG Clinic, University of Tübingen, Tübingen, Germany

<sup>5</sup>Robert Bosch Center for Tumor Diseases, Stuttgart, Germany

<sup>6</sup>Institut de Recherche en Cancérologie de Montpellier (IRCM), INSERM U1194, Université de Montpellier, Institut Du Cancer de Montpellier (ICM), F-34298 Montpellier, France

<sup>7</sup>Live-Cell Structural Biology Laboratory, Department of Medicine and Life Sciences, E-08005 Barcelona, Spain

<sup>8</sup>Universitat Pompeu Fabra (UPF), Barcelona, Spain

<sup>9</sup>Structural Bioinformatics Laboratory (GRIB-IMIM), Department of Medicine and Life Sciences, E-08003 Barcelona, Spain

<sup>10</sup>Institut de Pharmacologie et de Biologie Structurale, IPBS, Université de Toulouse, CNRS, UPS, Toulouse, France

<sup>11</sup>Present address: Molecular Biology CORE (CDB), Hospital Clínic de Barcelona, Barcelona, Spain

<sup>12</sup>These authors contributed equally

<sup>13</sup>Lead contact

\*Correspondence: priyanka.sharma@ipbs.fr

<https://doi.org/10.1016/j.isci.2024.109584>



links PADI2 to gene activation in the context of naive pluripotency.<sup>28</sup> PADI2 also contributes to chromatin modification that promotes the differentiation of oligodendrocyte precursors and efficient myelination, processes required for motor and cognitive functions.<sup>29</sup> Recently, PADI2 was found to citrullinate MEK1, thereby promoting signaling by extracellular signal-regulated protein kinase 1/2 (ERK1/2) signaling in endometrium cancer.<sup>22</sup>

Recently, we identified another unique function of PADI2: the deimination of arginine 1810 (R1810) on repeat 31 of the C-terminal domain (CTD) of the large subunit of RNA polymerase 2 (RNAPII). This modification potentiates the interaction between the RNAPII-CTD and the positive transcription elongation factor b complex (P-TEFb) and consequently facilitates gene expression that is essential for cell identity.<sup>30</sup> Indeed, we reported that the citrullination of R1810 in the CTD of RNAPII facilitates its interaction with the P-TEFb kinase complex and promotes the recruitment of CDK9 to transcription start sites (TSS).<sup>30</sup> This contributes to overcoming the RNAPII pausing barrier, highlighting the functional connection between PADI2 and the P-TEFb kinase complex. These studies support the notion that PADI2 can selectively citrullinate arginine residues in specific proteins, which may explain its relevance to several pathophysiological conditions.

The *PADI* genes are ubiquitous in vertebrates but are absent from yeast, worms, and flies. A recent study focusing on the comprehensive identification of *PADI* homologs unveiled the evolutionary trajectory of *PADIs* within the animal lineage.<sup>31</sup> *PADIs* appear to have been introduced from cyanobacteria into animals by horizontal gene transfer (HGT),<sup>31</sup> supporting the previous hypothesis of HGT as a mechanism for introducing new genetic material into vertebrate genomes.<sup>32–35</sup> Enzymes mediating citrullination in human parasites and microbes are highly divergent in sequence and have different substrate specificities. These include pPAD, which is an extended agmatine deiminase found in *Porphyromonas gingivalis*, and *Giardia lamblia* ADI, an extended form of the free L-arginine deiminase gADI found in this human parasite.<sup>36,37</sup> Previous small-scale phylogenetic analyses of the *PADI* family have shown that PADI2 is the most conserved family member. Mammalian *PADIs* comprise three structural domains; the N-terminal domain (NTD) (PADI\_N, Pfam annotation: PF08526), the middle domain (PADI\_M, Pfam annotation: PF08527), and the catalytic C-terminal domain (PADI\_C, Pfam annotation: PF03068).<sup>38–41</sup> Of note, PADI1 and PADI3, and PADI4 and PADI6, are more closely related in evolutionary terms than PADI2 and any other family member.<sup>6</sup> This finding suggests that PADI2 is most closely related to the ancestral *PADI* enzyme and that the other *PADI* enzymes arose through gene duplication more recently in evolution.

The enzymes responsible for essential PTMs (including phosphorylation, acetylation, and glycosylation) are found across all domains of life, suggesting that they were present in the Last Universal Common Ancestor (LUCA).<sup>42,43</sup> Similarly, the PADI\_C domain was present in the LUCA, indicating an ancient origin of citrullination. The PADI\_M domain, which encompasses the sequential calcium-binding sites and maintains allosteric communication with PADI\_C,<sup>25</sup> is present in cyanobacteria. Indeed, recent work supports the degeneration of the PADI\_N domain in cyanobacteria and demonstrates the existence of catalytically active *PADI* proteins in cyanobacteria.<sup>31</sup> The relatively recent appearance of the NTD during the *PADI* family evolution suggests it has a functional relevance that is unique to higher organisms.

Previously, we found that PADI2 interacts with P-TEFb to maintain the expression of actively transcribed genes.<sup>30,44</sup> The P-TEFb complex comprises cyclin-dependent kinase 9 (CDK9) and its regulatory partner cyclin T1 (CCNT1).<sup>45</sup> Activation of the P-TEFb complex is required early in transcription to overcome RNAPII pausing and promote the productive phase of transcription elongation.<sup>45–48</sup> Lack of recruitment of the P-TEFb complex and associated components has been linked to several disease conditions,<sup>49–51</sup> highlighting the importance of this process in healthy physiology. While PADI2 could be a therapeutic target that specifically prevents p-TEFb-mediated expression of genes that promote disease states, more needs to be understood.

Here, we analyze the recent evolution of the *PADI* protein family. We performed a comparative genomics analysis of PADI2 and identified 20 putative amino acid substitutions, in different species, that might be important for PADI2 structure and functionality in mammals. The majority of the selected amino acids are exposed in the three-dimensional (3D) structure and belong to the non-catalytic NTD or to the middle domain of PADI2. These positions suggest that they can participate in protein-protein interactions. We investigated the functional relevance of positively selected exposed amino acids in the NTD and middle domain of PADI2 by analyzing their impact on PADI2 interactions with the P-TEFb complex. We found that leucine 162 (L162), in the exposed loop of the middle domain of PADI2, underwent positive selection during the initial divergence of primates. Indeed, this residue contributes to interactions between PADI2 and the P-TEFb complex. These results establish a link between the evolutionary selection of key residues in PADI2 and their role in regulating the protein-protein interactions required for functional transcriptional regulation.

## RESULTS

### Evolutionary analysis of the peptidyl arginine deiminase family

To obtain a global picture of evolutionary relations within the *PADI* family, we first searched for all *PADI* family members with sequence homology to human PADI2 in the OMA (Orthologous MAtRix) database<sup>52</sup> and collected all the available sequences in mammals (see [STAR Methods](#)). After close inspection, we discarded partial sequences and low-quality genomes, leaving 185 *PADI* orthologs that represent 25 mammalian species and three bird species ([Figure S1](#)). We used these sequences to build a multiple sequence alignment (MSA, [File S1](#)) and reconstructed a phylogenetic tree using the software RaxML, rooting the tree using the bird sequences as an outgroup ([Figure S2](#), see [STAR Methods](#)).

The tree we built is consistent with five *PADI* genes that resulted from the duplication of a common ancestral sequence. Indeed, each of the five well-defined clusters in the reconstructed tree corresponds to the five different *PADI* genes (*PADI1*, green; *PADI2*, purple; *PADI3*, blue; *PADI4*, red; and *PADI6*, orange), with known species relationships being broadly recapitulated within each of these clusters ([Figure S2](#)).

Robust bootstrap supports (82–100, [Figure S3](#)) for this tree strongly support the hypothesis that the genes within each cluster are orthologous. We observed clustering between *PADI1* and *PADI3* and between *PADI4* and *PADI6*, with *PADI2* being the most distant ortholog among family members. Notably, *PADI2* also has shorter branch lengths in comparison to other family members, indicating a higher level of conservation among family members in this cluster. *PADI6* family members have the longest branch lengths, indicating this cluster is the most divergent and suggesting these genes arose from an old duplication event. Interestingly, there are only three *PADI* genes in each bird species, and only one set of clusters within the *PADI2* family, suggesting that one or more duplication events took place after the divergence from birds. *PADI2* is also the closest subfamily to the remaining *PADI* outgroup species clusters, indicative of an ancestral position in this family.

Two bird genes (i.e., ANAPL06996 and ANAPL06995 in duck) in the unpainted cluster ([Figure S2](#)) did not exhibit clear orthology to specific *PADI* genes in either Ensembl or the OMA, though some were annotated as *PADI1* orthologs. Our data suggest that there was a bird-specific duplication, possibly from *PADI1*, which is the closest gene in the synteny analysis ([Figure S4A](#)). Our species comparison also revealed an inversion of the entire *PADI* locus that appears to have occurred in the common ancestor of humans and gibbons (Hominoidea; [Figure S4B](#)). Most mammalian species have the same gene order as that found in mice ([Figure S4C](#)). Overall, our analyses revealed key evolutionary changes at the *PADI* locus, including gene re-arrangements and duplications within mammals.

### Peptidyl arginine deiminase 2 is highly conserved across species

Since *PADI2* is likely to be the ancestral copy, we decided to focus on the most recent changes occurring in its sequence. By focusing on a specific phylogenetic selection of species, we aimed to identify changes that have occurred since the common amniote ancestor that occur with good resolution in different mammals. Our dataset contains *PADI2* representatives in all three major extant groups in class Mammalia and the principal orders or infraclasses: primates, rodents, carnivores, cetartiodactyla, lagomorphs, proboscidea, and marsupials. We built a multiple sequence alignment with all 28 *PADI2* sequences (see [STAR Methods, File S1](#)) and computed the ratio of substitution rates at non-synonymous sites and synonymous sites (dN/dS; also termed Ka/Ks) to infer the direction and magnitude of natural selection acting on protein-coding genes.<sup>53</sup> The difficulty of estimating the dN/dS ratio grows with the phylogenetic distance between sequences. Maximum likelihood methods were used to correct for multiple substitutions in the same site, and we estimated important parameters, such as the divergence between sequences or the transition/transversion ratio by deducing the most likely values to produce the input data.<sup>53</sup>

We calculated the dN/dS ratio for all internal branches and leaves in the tree using CodeML ([Figure 1](#)). We observed low dN/dS values whereby most of the branches had a value of <0.2. This was consistent with negative selection; in other words, changes in this genetic sequence were actively selected against. This analysis supports the idea that *PADI2* is a highly conserved protein, suggesting that it performs critical functions in the organism ([Figure 1](#)). However, we observed a relatively high dN/dS in the common ancestor of all primates, possibly indicating the relaxation of selection or the fixation of a few beneficial mutations due to positive selection at this point in evolution. The dN/dS ratio requires a high number of changes to detect a selection effect and our data indicate that dN/dS values are very low in *PADI2*. Another limitation of this method is that it does not distinguish amino acid substitutions that are chemically similar from others with very different properties. However, this analysis suggests that intriguing selection processes operated during *PADI2* evolution.

### Detection of positively selected residues in peptidyl arginine deiminase 2

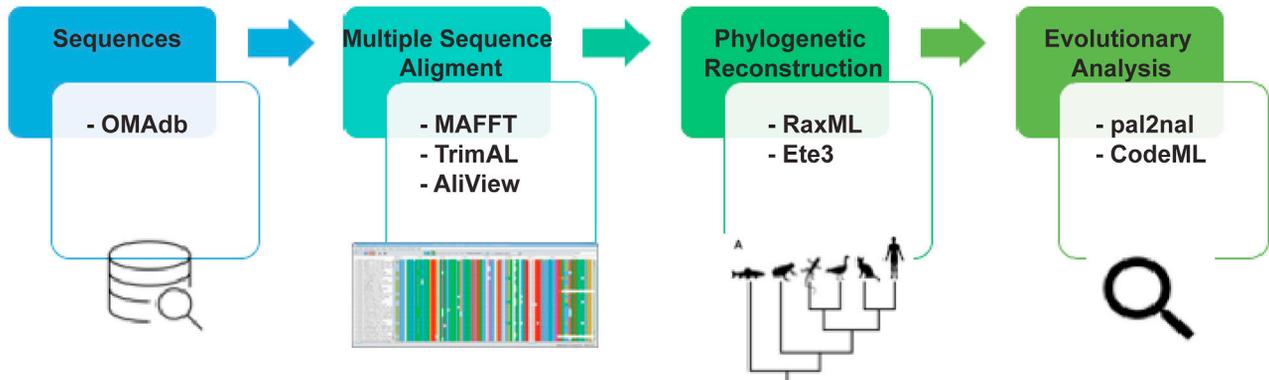
We used an orthogonal approach, the branch site positive selection test, to surpass the limitations of dN/dS analysis and identify specific residues that have undergone selection<sup>54,55</sup> This test can discriminate between the relaxation of selection and positive selection and can detect individual residues that are under positive selection in a particular lineage. This test compares a null model in which codon-based dN/dS for all branches can only be  $\leq 1$ ; in the alternative model, the labeled foreground branch may include codons evolving at dN/dS > 1.<sup>54</sup> This test can discriminate between the relaxation of selection or positive selection and can detect individual residues that are under positive selection in a particular lineage. After running the branch-site positive selection test on every branch of the *PADI2* tree, we detected 20 individual relevant substitutions in different parts of the tree. We removed a few candidates for positively selected amino acids detected by CodeML in *Loxodonta africana* due to the alignment of a non-homologous region (238–278), possibly due to an assembly error in that species. The residues identified in the different species or branches are shown in [Table 1](#) and [Figure 2](#).

### Positively selected residues in peptidyl arginine deiminase 2 may mediate protein-protein interactions

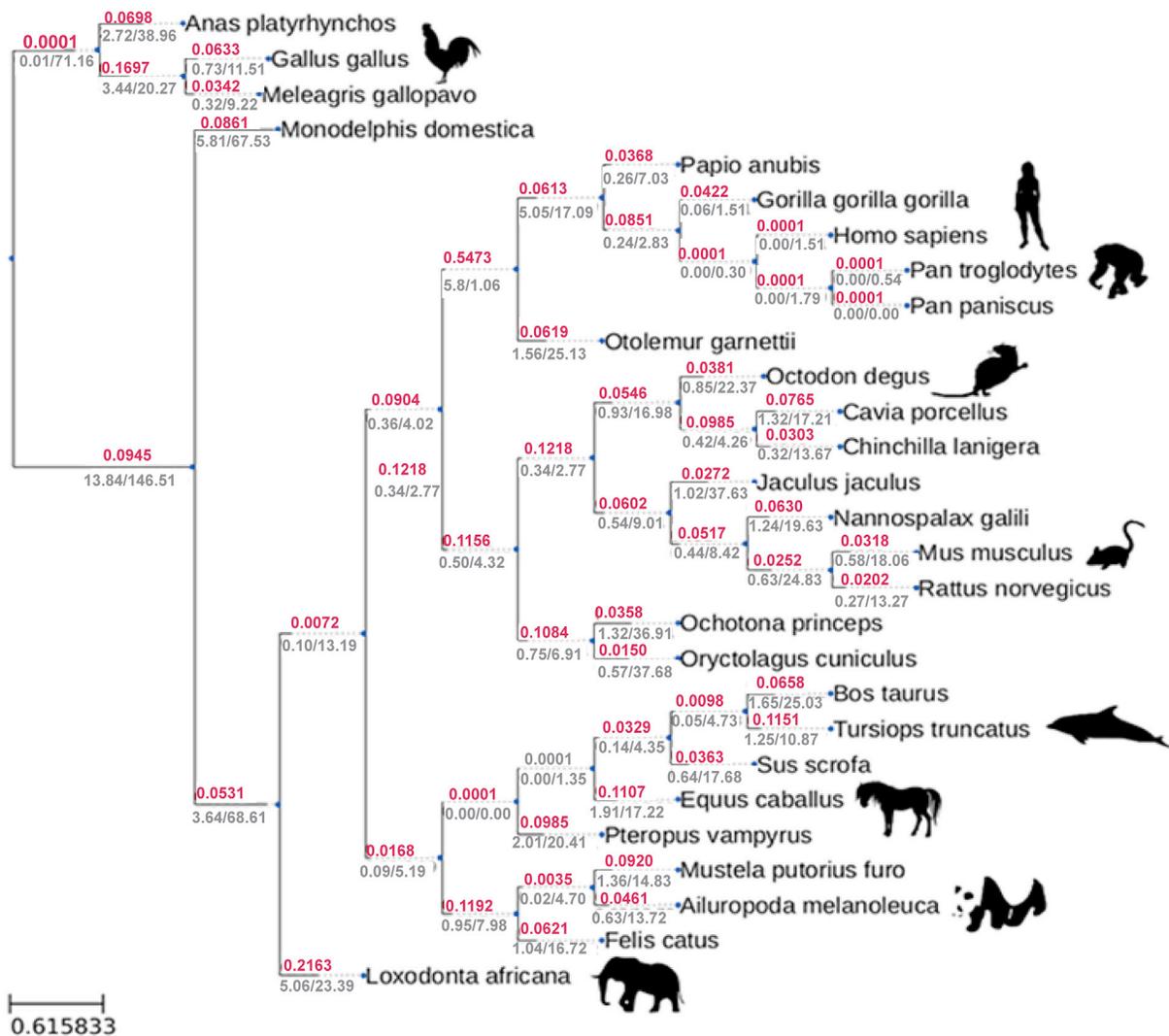
We next integrated our evolutionary insights into *PADI2* with known structural information to study the properties of positively selected residues. To pursue this analysis, we mapped all the positively selected amino acids onto the structure of *PADI2*. 12 of the 20 positively selected amino acids mapped onto the *PADI2*\_N and *PADI2*\_M domains, while eight mapped to the *PADI2*\_C domain. The latter contains the catalytic pocket and is responsible for the citrullination of arginine residues. ([Table 1](#); [Figures 3A](#) and [3B](#)). Of note, only one of the positively selected amino acids (L25) is fully buried (amino acids are numbered as shown in [Table 1](#), [File S2](#), [Video S1](#)); most of the other amino acids are fully exposed except for L342, V538, and M663, which are partially buried. In the secondary structure, approximately 50% of the amino acids are located in loops, while the others are either in beta strands or alpha helices (including helix-capping regions) ([Table 1](#)).

After evaluating the features of *PADI2*\_N and *PADI2*\_M domains and analyzing each particular amino acid and its potential role based on structural data, we concluded that the likely roles of E16, V53, E60, K136, L162, V201, S245, S249, S267, and T289 are to mediate interactions with putative binding partners. These amino acids are fully exposed, and their chemical properties correspond with classical protein interfaces.<sup>56,57</sup> We thus characterized these positively selected amino acids using Multi-VORFFIP (MVORFFIP), a tool that predicts protein-, peptide-, DNA-, and RNA-binding sites.<sup>56</sup> Notably, the MVORFFIP predictions for all except V201 and L25 yielded very high scores (>0.7, on a

A



B



**Figure 1. PADI2 conservation across species**

(A) Diagram showing the main steps involved in phylogenetic reconstruction and evolutionary analysis. Each box contains a list of the tools used. (B) The PADI2 gene tree showing the dN/dS ratios (red) for every branch as well as dN and dS (gray) for every branch. The common branch to all primates is highlighted in blue. Branch lengths are proportional to nucleotide substitutions, as calculated by ete3 (related to Figures S1–S4 and File S1).

**Table 1. Positively selected sites across species under the branch site (CodeML) and mapped onto the PADI2 human structure**

Domain	Alignment position	Sequence position (human)	Secondary structure	Flanking residues	MVORFFIP scores (0–1)	Accessibility
PADI_N	70	E16	Beta strand	RVEAV	0.8	Exposed
PADI_N	79	L25	Beta strand	TYLWT	0.0	Buried
PADI_N	107	V53	Beta strand	EVRD	0.5	Exposed
PADI_N	114	E60	Loop	AEEVA	0.7	Exposed
PADI_M	190	K136	Loop	NPKKA	0.6	Exposed
PADI_M	216	L162	Loop	PWLPK	0.7	Exposed
PADI_M	260	V201	Beta strand	EIVLY	0.2	Exposed
PADI_M	304	S245	Loop	GGSAE	0.6	Exposed
PADI_M	308	L249	Beta strand	ELLFF	0.8	Exposed
PADI_M	322	S263	Loop	GFSGL	0.7	Exposed
PADI_M	326	S267	Beta strand	LVSIH	0.7	Exposed
PADI_M	350	T289	Beta strand	TDTVI	0.6	Exposed
PADI_C	401	L342	Loop	QYLNR	0.7	Partially buried
PADI_C	439	F380	Loop	KDFPV	0.8	Exposed
PADI_C	460	S401	Loop	FESVT	0.8	Exposed
PADI_C	511	K452	Helix	FLKAQ	0.7	Exposed
PADI_C	567	D507	Helix capping	QKDGH	0.5	Exposed
PADI_C	596	S536	Helix	NESLV	0.7	Exposed
PADI_C	598	V538	Helix	SLVQE	0.8	Partially buried
PADI_C	723	M663	Helix capping	WHMVP	0.7	Partially buried

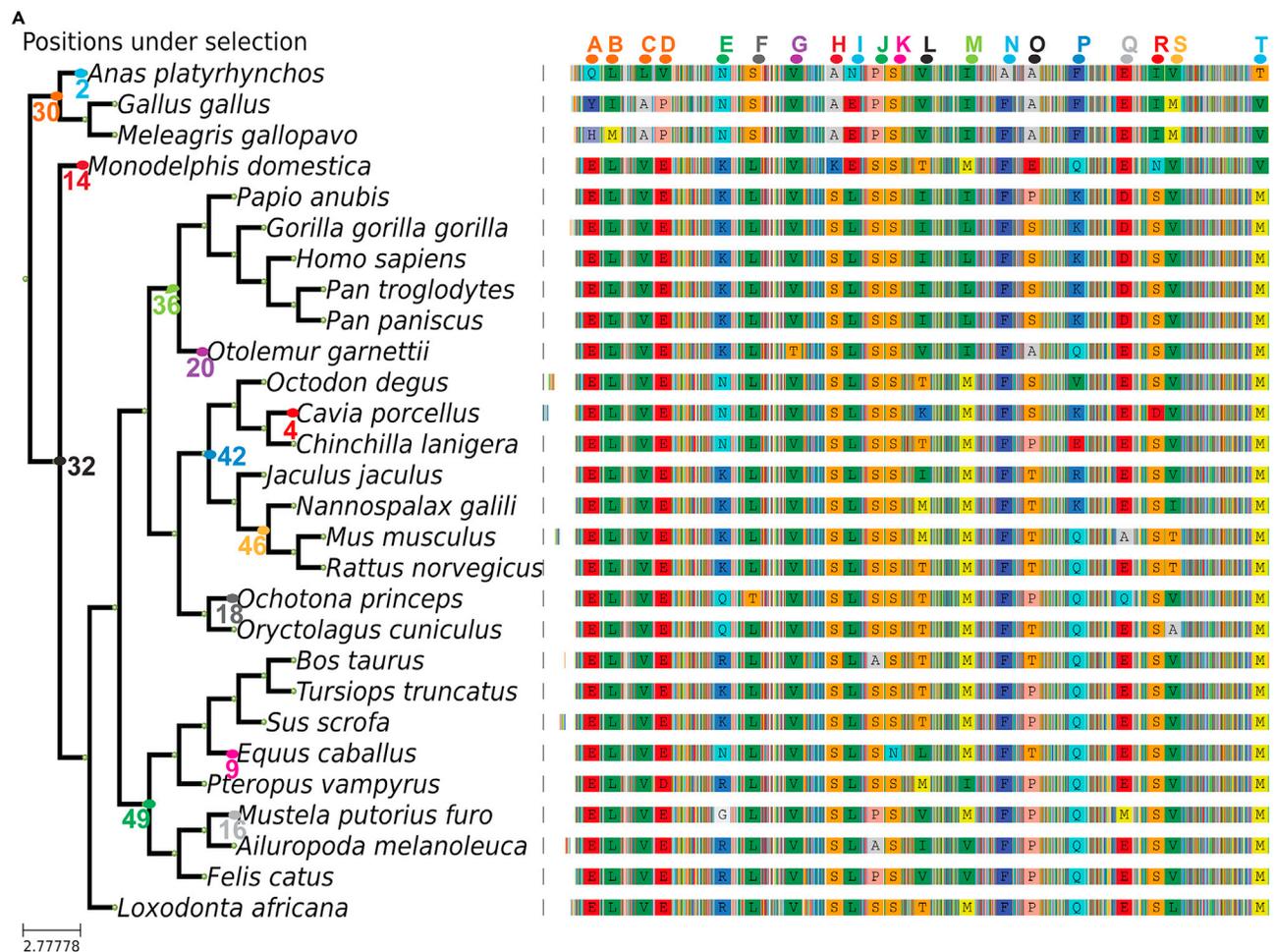
The position of the amino acid in the domain, MSA, and PDB file (Supplementary material [File S2](#)) are shown in the domain, sequence position, and PDB number, respectively. The flanking residues are shown with the given positively selected amino acid depicted in bold. The type of secondary structure and solvent accessibility (i.e., exposure) is shown in the secondary structure and accessibility column. Finally, the MVORFFIP scores (see [STAR Methods](#)) are also indicated.

scale of 0–1), supporting the hypothesis that they are interface residues ([Figures 3C and 3D](#); [Table 1](#)). S263 is located in a long loop close to the interface between domains, suggesting that it plays a role in hinge motions and geometrical orientation between domains. Given its position in the structure, the likely role of L25 (the only positively selected amino acid that is completely buried) may be to prevent the distortion of a nearby beta-sandwich. Indeed, a larger hydrophobic amino acid at this position would change the packing of the beta-strands.

We found that three positively selected residues in the PADI2\_C domain, L342, F380, and S401, are close to the catalytic pocket. These residues are likely to play a role in substrate specificity (L342) and in the dynamics of catalysis (F380 and S401). K452, S536, and V538, however, are located in a helix that is distant from the active site pocket. Judging from the structural microenvironment, these residues could also play a role in protein-protein interactions, in line with the high scores assigned by MVORFFIP ([Figure 3D](#); [Table 1](#)).<sup>56</sup> Finally, D507 and M663 are both located in helix-capping positions. Conservation in helix capping is important for the stability and integrity of the helix. Residue M663 may also play a role in the packing of the C-terminal tail. A CAPS Database search<sup>58</sup> showed that other helices with the same capping structure present a conserved small hydrophobic patch in the same position by including M, V, or L residues. It is noteworthy that M663 is indeed strictly conserved across all human PADIs. These observations suggest that most of these positively selected amino acids may function in facilitating and stabilizing protein-protein interactions.

### The middle domain of peptidyl arginine deiminase 2 maintains interactions with the positive transcription elongation factor b complex

We next tested whether residues that have been positively selected in evolution but are not linked to the catalytic domain of PADI2, might interact with other proteins. For this analysis, we examined the P-TEFb kinase complex, which we previously found to interact with PADI2.<sup>30</sup> The structures of PADI2<sup>25</sup> and CDK9-CCNT1<sup>59</sup> are known individually, and therefore it was possible to derive a structural model of the PADI2-CDK9/CCNT1 complex by protein docking. After deriving the docking ensemble of PADI2-CDK9/CCNT1, we ranked the models. Importantly, we used a swarm-based approach (among other metrics) to select the putative model of interaction, reasoning that determining the region of PADI2 would be relevant to its interaction with CDK9/CCNT1 ([Figure S4D](#), [File S3](#), [Video S2](#)). Moreover, we derived the structural model using AlphaFold-Multimer.<sup>60,61</sup> As shown in [Figure S4E](#) ([File S4](#), [Video S3](#)), interacting regions in the structural model from AlphaFold-Multimer overlap substantially with those identified by docking ([Figure S4D](#)). More specifically, the predicted interface between PADI2 and CDK9-CCNT1 contained PADI2\_M, and the hinge region was overrepresented by the number of positively selected residues. Importantly, L162 is present on the highly exposed loop ([Figure 4A](#)).



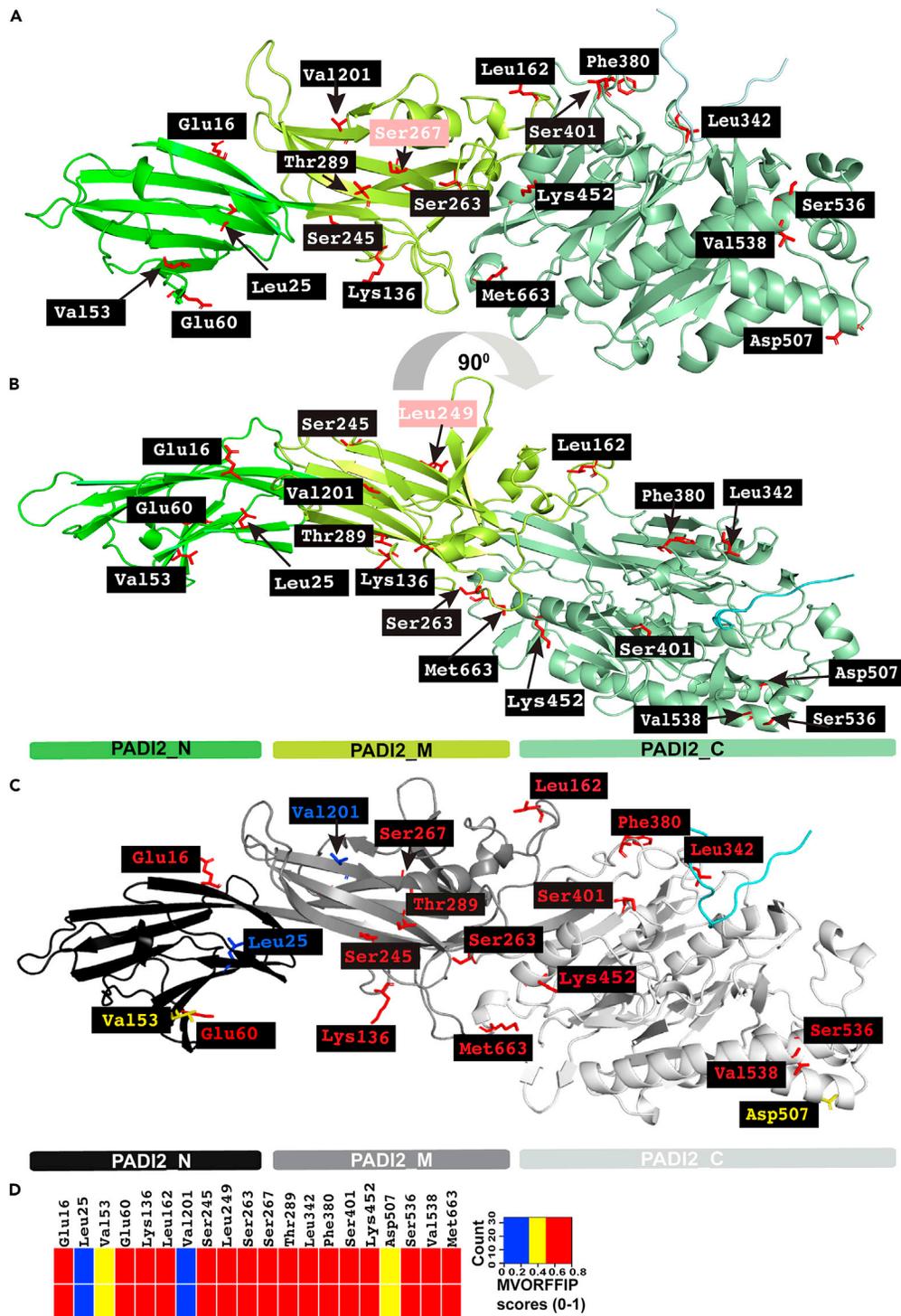
**Figure 2. Summary of the branch site test results for internal and terminal branches**

(A) Summary of the multiple sequence alignment, showing only columns with positions that are detected as being under positive selection with the branch site test. The colored-blurred positions in between represent the rest of the alignment (illustrated in [File S1](#)). The node number correspondence can be found in [Figure S1](#) (related to [Table 1](#)).

While AlphaFold-Multimer generates a single model, docking is represented by a range of docking poses. This range of models allowed us to identify regions of the predicted interface, i.e., PADI2\_M and PADI2\_C ([Figure 4A](#)), that are over-represented in the docking space. The distribution of the top 200 docking poses (represented using a unique point depicting the center of mass of CDK9/CCNT1) revealed that the region around the hinge between the PADI2\_M and PADI2\_C was over-represented ([Figure 4B](#), [File S5](#), [Video S4](#)). This particular region includes L162, a residue overrepresented in the docking conformers, along with T159 and W161. These three residues are located in a highly exposed region. We therefore selected these residues for experimental validation. Note that L162 in PADI2 was present among the top-ranking interface residues likely to mediate the interaction with CDK9/CCNT1 and is also a positively selected residue, suggesting that it is important for maintaining this interaction.

### The loop encompassing positively selected Leu162 is important for cell proliferation and maintaining peptidyl arginine deiminase 2 interactions with the positive transcription elongation factor b complex

Since the P-TEFb kinase complex is involved in transcription elongation and cell proliferation,<sup>62–64</sup> we experimentally tested whether the L162-encompassing loop in the middle domain of PADI2 affects cell proliferation. In addition to L162, we tested its neighboring amino acids, W161 and T159, which also ranked highly in our evolutionary and interaction prediction analysis and are located in a highly exposed loop region structurally close to the boundary between the PADI2\_M and the catalytic domain. We selectively expressed the green fluorescence protein (GFP)-tagged wild-type (WT) PADI2 as well as the single (T159A or W161A or L162A), double (L162A/W161A), or triple (L162A/W161A/T159A) mutant of PADI2. The GFP-positive cells were sorted by FACS to ensure positive cell selection ([Figures S5A](#) and [S5B](#)). We monitored cell proliferation in HeLa cells expressing the GFP-tagged PADI2 WT and mutants (single, T159A or W161A or L162A, double, L162A/W161A, and triple, L162A/W161A/T159A). Strikingly, we observed that among the three single PADI2 mutants, L162A significantly decreased cell

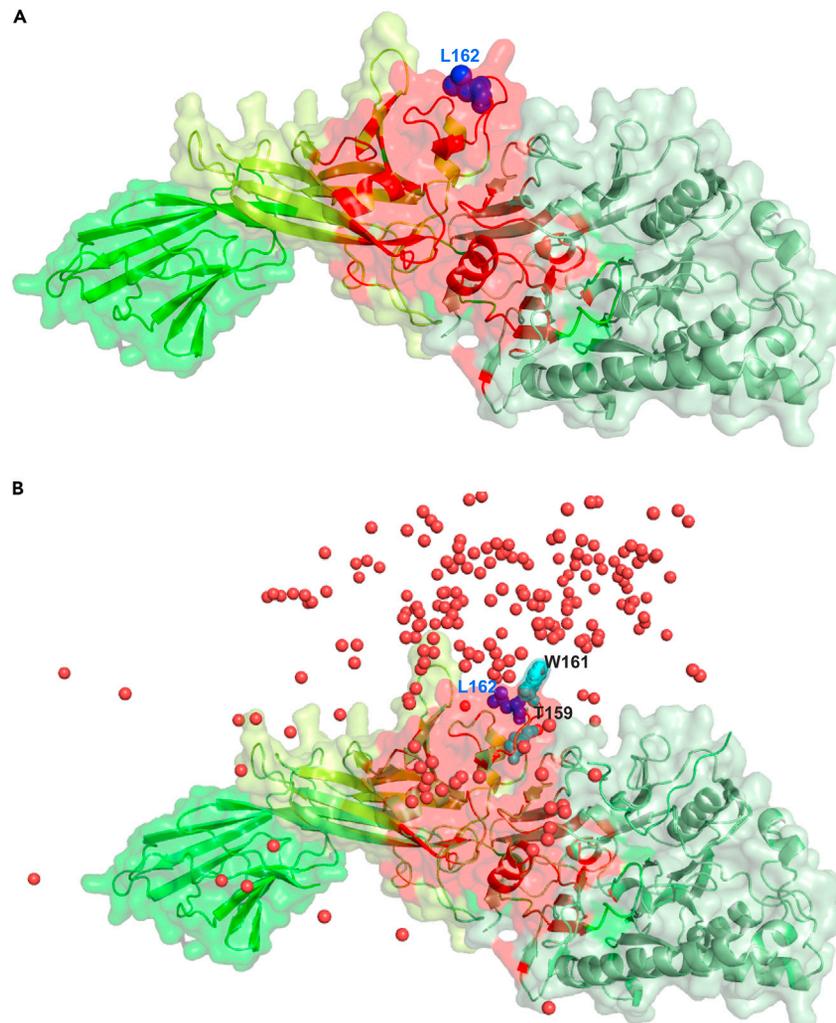


**Figure 3. Mapping of positively selected amino acids onto the PADI2 structure**

(A and B) Ribbon representation of PADI2 showing the N-terminal domain (PADI2\_N), the middle domain (PADI2\_M), and the C-terminal domain (PADI2\_C) in different shades of green (in green, bright green, and pale green, respectively). (B) The same as (A) but applying a 90-degree rotation over the longitudinal axis. Pink highlighted residues are specific to the respective orientation. PyMol session of PADI2 and the highlighted amino acids are available in Supplementary Material [File S2](#).

(C) Cartoon representation of PADI2 with positively selected amino acids colored according to MVORFFIP scores as shown in the heatmap in panel D; PADI2\_N, PADI2\_M and PADI2\_C domains colored in black, gray, and white, respectively.

(D) The heatmap is based on the MVORFFIP scores as mentioned in [Table 1](#) (related to [Table 1](#), [File S2](#), and [Video S1](#)).



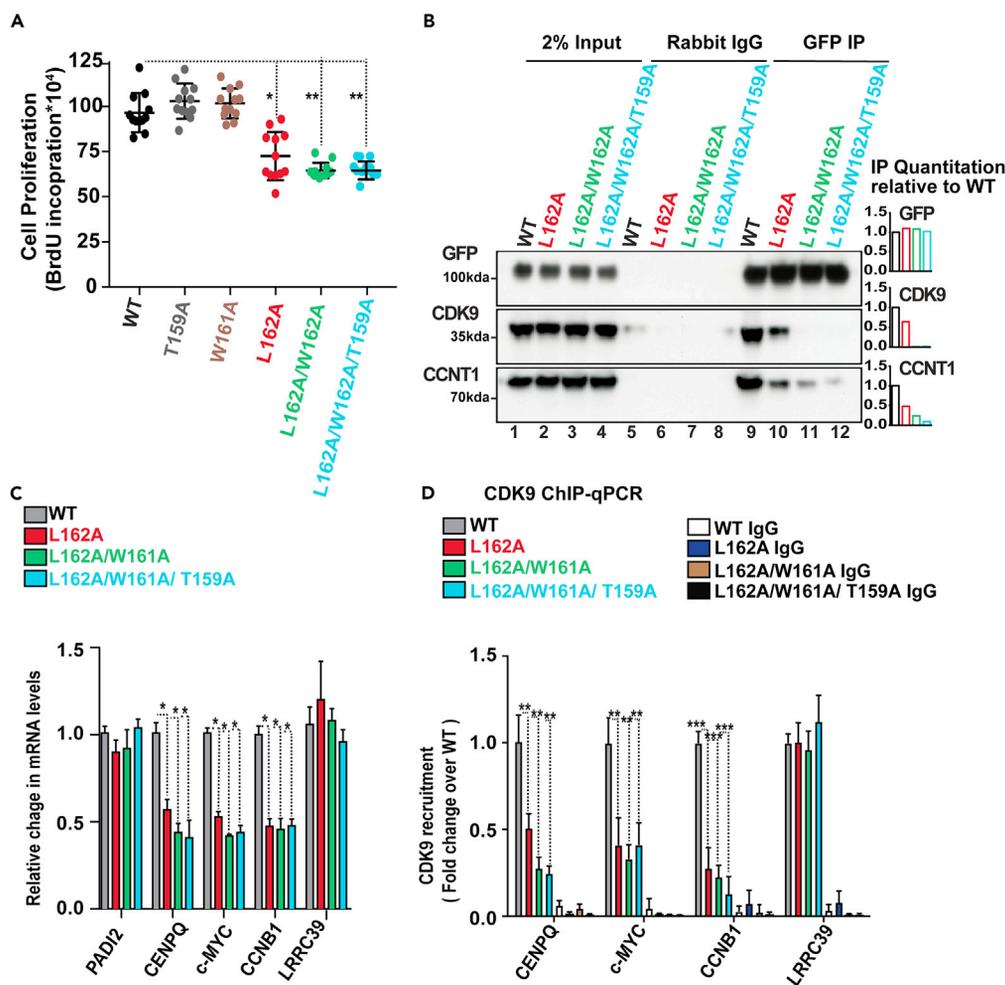
**Figure 4. The predicted interface between PADI2 and the P-TEFb complex**

(A) Surface representation of PADI2 with the same color scheme as in Figure 3 with the consensus interface with P-TEFb complex derived from docking and AlphaFold-Multimer structural models highlighted in red and L162 shown as a blue sphere.

(B) The same representation as in (A) including the top 200 docking poses represented as red spheres depicting the center of mass of CDK9/CCNT1; T159 and W161 are also shown in cyan as sphere representation. The PyMol session of PADI2 with the P-TEFb complex with docking and AlphaFold-Multimer shared region is available in Supplementary Material File S5 (related to File S5, Video S4, Figure S4).

proliferation compared to WT, highlighting the functional relevance of L162 in the middle domain of PADI2 (Figures 5A and S5C). In the same manner, the PADI2 double mutant L162A/W161A, as well as the triple mutant L162A/W161A/T159A, also showed a significant reduction in cell proliferation. These observations highlight the functional role of L162A in this remarkably exposed loop in the middle domain of PADI2.

Therefore, we focused on the L162A, L162A/W161A, and L162A/W161A/T159A mutants along with WT PADI2 for further analysis. In the immunoprecipitation assay, using a GFP-tagged antibody, we observed that the P-TEFb complex (containing CCNT1 and CDK9) from HeLa cells was efficiently immunoprecipitated with GFP-WT PADI2 but not with GFP-L162A-PADI2 (Figure 5B). These results confirmed the functional role of the positively selected L162 in maintaining PADI2 interactions with the P-TEFb complex. Likewise, the P-TEFb complex was only weakly immunoprecipitated with either the double or triple PADI2 mutant, highlighting the function of the highly exposed region in the middle domain of PADI2 (PADI2\_M) in this protein-protein interaction. Considering that the proper recruitment of the P-TEFb complex is required for efficient transcription of highly expressed genes relevant to cell proliferation, we validated the reduced expression levels of the *c-MYC*, *CENPQ*, and *CCNB1* genes in cells expressing the single, double, or triple PADI2 mutant (Figure 5C). Of note, the levels of PADI2 expression did not differ significantly. We observed that the L162A single mutant (as well as double and triple mutants) significantly and specifically reduced the expression of highly expressed genes, including *c-MYC*, *CCNB1*, and *CENPQ* without affecting the levels of a control gene (the low-expressed *LRR39* gene). Next, we examined if the L162A mutant in PADI2 can mediate CDK9 recruitment to the promoter region of the target genes. We performed a chromatin immunoprecipitation (ChIP) assay using a CDK9-specific antibody in HeLa cells specifically



**Figure 5. Positively selected L162 encompassing loop contributes to its interaction with the P-TEFb complex**

(A) Cell proliferation of HeLa cells specifically expressing the WT or mutant PADI2 (single, T159A or W161A or L162A; double, L162A/W161A or triple, L162A/W161A/T159A). Data represent mean  $\pm$  SEM of at least eight biological experiments. \*p value <0.05; \*\*p value <0.01.

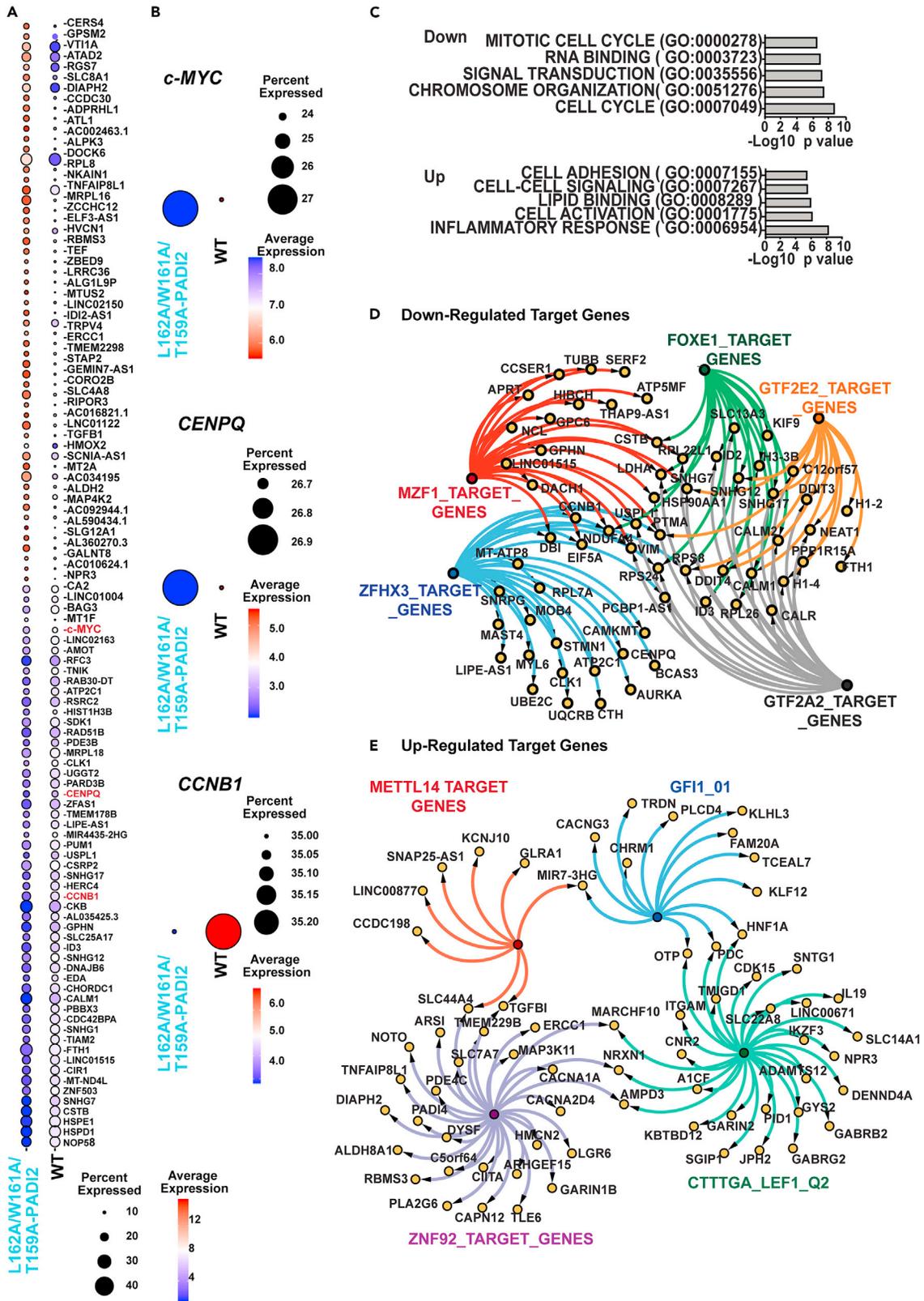
(B) Representative image of immunoprecipitation with GFP-specific antibody or non-immune rabbit IgG of GFP-positive HeLa cells nuclear extracts expressing wild-type (WT) or mutant (single, L162A; double, L162A/W161A or triple, L162A/W161A/T159A) PADI2 followed by western blot with the indicated antibodies. The relative quantification is shown as a bar plot of the two biological replicates.

(C) Quantitative RT-qPCR validation in HeLa cells selectively expressing the WT or mutant PADI2 (as in B). Changes in mRNA levels were normalized to GAPDH mRNA. Data represent the mean  $\pm$  SEM of  $n \geq 3$  biological experiments for all plots in the figure. Two-tailed unpaired Student's t-test was used to determine the statistical significance between the groups.

(D) ChIP-qPCR assay performed in HeLa cells selectively expressing the WT or mutant PADI2 (as in B) with CDK9 antibody. Non-immune IgG was used as a negative control. Y axis: fold change over the input samples. Data represent mean  $\pm$  SEM of three biological experiments, \*p value <0.05; \*\*p value <0.01 (related to Figure S5).

expressing either WT PADI2, the L162A mutant, or the double and triple mutants. We found that CDK9 occupancy decreases in the presence of mutants compared to the WT PADI2 control (Figure 5D), suggesting that the positively selected L162 residue, along with W161A/T159A, is important to maintaining PADI2 interaction with the P-TEFb complex.

Given the fact that the L162-encompassing loop is important to maintaining PADI2 interaction with the P-TEFb complex, we next explored the effects of the PADI2 triple-mutant (L162A/W161A/T159A) on the transcriptome using single-cell RNA sequencing experiments. We compared the single-cell transcriptomics atlas of PADI2 triple-mutant to that of WT PADI2-expressing HeLa cells (Figures S6A and S6B). After quality control, 49,503 cells were subjected to downstream analysis, of which, 17,139 and 32,364 cells were derived from the WT and triple mutant PADI2 samples, respectively. We did not observe any significant differences in cell populations through unsupervised clustering and dimensionality reduction. TSNE plots organized the data into three clusters (Figure S6A). However, these clusters are not significantly different in triple mutants in comparison to WT (Figure S6B). Among the differentially expressed genes, we found *c-MYC* (Figures 6A and 6B). Within the top genes that are differentially expressed, we also found *CCNB1* and *CENPQ*, which are related to cell cycle regulation



**Figure 6. The L162 enclosed loop regulates the transcription regulation**

(A). Dotplot depicting the top 55 up- and downregulated genes in L162A/W161A/T159A-PADI2 mutants vs. WT. Selected were the  $p_{\text{adj}} < 0.05$  and highest FC (fold change) according to the Find Markers function in Seurat (4.3.0,  $\log_2\text{FC.threshold} = 0.25$ ,  $\text{min. pct} = 0.1$ ).

(B). Dot plots depicting the expression of *c-MYC* and *CCNB1* in L162A/W161A/T159A-PADI2 mutants vs. WT. The circle size depicts the percentage of cells expressing the respective gene, while the color encodes the average expression.

(C). Gene set enrichment analysis (GSEA) for molecular functions and biological processes. Representative processes are presented (Upper- Downregulated genes, Lower- Up-regulated genes). The X axis shows the  $-\log_{10}$  transformed p values. GO, Gene Ontology.

(D and E). The connection map of the target genes from the representative regulatory elements in WT/(L162A/W161A/T159A) PADI2 (D) upregulated genes (E) downregulated genes (related to Figure S6).

(Figures 6A and 6B). Remarkably, in global differential expression (DEseq) analysis, the PADI2 triple mutant (L162A/W161A/T159A) affected the expression of over 469 genes ( $\log_2\text{FC} > 0.58$  or  $\log_2\text{FC} < -0.58$ ,  $p$  value  $< 0.05$ ). Of these, 162 genes were downregulated and 307 up-regulated (Table S2). Gene ontology analysis of the top down-regulated genes ( $\log_2\text{FC} < -1.0$ ,  $p$  value  $< 0.01$ ) revealed enrichment in the mitochondrial function, cell cycle, RNA binding, and chromosome organization. The up-regulated genes were instead enriched for genes linked to the inflammatory response, cell-cell signaling, and cell adhesion processes ( $\log_2\text{FC} > 1.0$ ,  $p$  value  $< 0.01$ ) (Figure 6C; Table S3). We also investigated the regulatory gene sets using the gene ontology tool. The promoter regions of downregulated genes were significantly enriched for GTF2A2 (UniProt: P52657), GTF2E2 (UniProt: P29084), FOXE1 (UniProt: O00358), MZF1 (UniProt: P28698), ZFH3\_3 (UniProt: Q15911) binding motifs (Figure 6D; Table S4). Meanwhile, the upregulated genes showed significant enrichment in the promoter regions for motifs that can bind ZNF92 (UniProt: Q03936), METTL14 (UniProt: Q9HCE5), GFI1 and CTTTGA motif (Figure 6E; Table S4). This analysis highlights a connection between regulatory elements and differentially expressed (WT vs. PADI2 triple mutant) genes.

Importantly, *c-MYC*, *CCNB1*, and *CENPQ* overexpression have been linked to oncogenesis.<sup>65–70</sup> These results support the functional relevance of the positively selected L162-encompassing loop in the middle domain of PADI2 in the expression of important oncogenic genes. Modulation of *c-MYC* expression by epigenetic regulation in cancer cells has been suggested to increase heterogeneity in the transcriptional landscape that promotes tumorigenesis.<sup>71</sup> These observations suggest a potential function for the L162-encompassing loop of PADI2 in the regulation of the *c-MYC/CCNB1/CENPQ* axis which has been linked to tumorigenesis.

**DISCUSSION**

We characterized the evolution of the PADI gene family using the most complete dataset of mammalian orthologous sequences to date. The three major extant groups in the class Mammalia are organized into two subclasses: Prototheria which includes monotremes (platypus and echidnas), and Theria, which includes the infraclasses Methatheria (marsupials), and Eutheria (placental mammals). According to Jones et al. 2009,<sup>72</sup> there are 5416 mammalian species divided into 29 orders. Most known species belong to the 7 biggest orders (4865 species), all of which have at least one representative in our phylogenetic analysis.

The functional relevance of PADI2 is highlighted by its high sequence conservation across species and very low dN/dS values. We took advantage of the natural experiments performed by evolution and applied sensitive methodologies such as the branch-site test of positive selection. Comparing the amino acid sequences of 25 mammalian species, we identified 20 positively selected residues predominantly located in the non-catalytic domain of the PADI2 in different species and phylogenetic branches (Figures 1, 2 and S1–S4). Since these selected residues could have a strong functional impact on PADI2, we studied them further using the human sequence.

By examining the location of the positively selected residues in the structure of human PADI2, modeled in a complex with R1810 in the CTD of RNAPII, we determined that the majority of positively selected PADI2 residues were structurally highly exposed in the non-catalytic domains of PADI2 (Figure 3; Table 1). Notably, the highly exposed position and chemical properties of these positively selected residues suggested they play important functions in cellular processes. In addition, our MVORFFIP analysis revealed that these interface residues may play a significant role in protein-protein interactions (Figure 3D). Therefore, it could be postulated that the appearance of these positively selected residues contributed to the evolution of PADI2 function by modulating its interaction with an essential cognate set of proteins. Further work will be required to investigate these possibilities.

Specifically, the P-TEFb complex, comprising CCNT1 and CDK9, interacts with PADI2 and facilitates the effects of the citrullination of RNAPII-R1810 on transcription.<sup>30</sup> Hence, we used structure modeling approaches to derive the tertiary structure of the PADI2–CCNT1/CDK9 complex. To identify the PADI2 residues that contribute to forming the PADI2–CCNT1/CDK9 complex, we analyzed structural models and identified a highly exposed loop in the PADI2\_M domain as a potential interacting region. Strikingly, this loop includes the positively selected residue L162, along with T159 and W161. This observation supports the notion that positively evolved residues tend to coincide with the important residues on the protein surface, and that positively selected structural clusters are important for cellular function (Figure S6C).

Our data in HeLa cells overexpressing the GFP-tagged WT and a mutant PADI2 also unveiled an essential function for the L162-containing loop in interaction with the P-TEFb kinase complex and cell proliferation. Notably, the mutation of L162, alone or in combination with mutated W161 and/or T159, reduced (i) cell proliferation, (ii) PADI2 interaction with P-TEFb kinase complex, and (iii) reduced the PADI2-dependent expression of highly expressed genes that are relevant for cell growth.

We focused on the interactions between PADI2 and the P-TEFb complex, specifically with the positively selected L162-encompassing loop. However, in the future, it would be interesting to investigate the functional role of other positively selected PADI2 amino acids in cellular processes. Nonetheless, our analysis shows the importance of using a multi-disciplinary approach. Here, we used comparative genomics, evolutionary analysis, structural modeling, and genetic perturbations to analyze the functional relevance of a recently evolved non-catalytic

domain in a highly conserved enzyme family. Overall, this evolutionary approach led to the identification of a PADI2 structural loop that is under positive selection and specifically supports an important role for PADI2 in the modulation of transcription elongation.

### Limitations of the study

Note that one limitation of our study was that we could not mutate the endogenous PADI2 gene due to the hyperpolyploid nature of the HeLa genome. Therefore, we based our analysis on the overexpression of GFP-tagged wild-type and mutant versions of PADI2 in HeLa cells.

### STAR★METHODS

Detailed methods are provided in the online version of this paper and include the following:

- **KEY RESOURCES TABLE**
- **RESOURCE AVAILABILITY**
  - Lead contact
  - Materials availability
  - Data and code availability
- **EXPERIMENTAL MODEL AND SUBJECT DETAILS**
  - Cell lines
- **METHOD DETAILS**
  - PADI2-GFP plasmids and fluorescence-activated cell sorting (FACS) sorting
  - Cell proliferation assay
  - Incucyte® proliferation assays for live-cell analysis
  - RNA extraction and RT-qPCRs
  - Single-cell samples, library preparation, and sequencing
  - Single-cell RNA sequencing data processing and analysis
  - Gene ontology (GO) and regulatory gene set analysis
  - GFP-tagged PADI2 immunoprecipitation, western blot
  - Multiple sequence alignments
  - Phylogenetic reconstruction
  - Estimation of dN/dS and positive selection
  - Characterization of positively selected residues using MVORFFIP
  - Structural modeling of PADI2-CDK9/CCNT1 and selection of putative interface residues
  - AlphaFold-multimer modeling of PADI2-CDK9/CCNT1 complex
- **QUANTIFICATION AND STATISTICAL ANALYSIS**

### SUPPLEMENTAL INFORMATION

Supplemental information can be found online at <https://doi.org/10.1016/j.isci.2024.109584>.

### ACKNOWLEDGMENTS

We thank Fátima Gebauer and François Le Dily from CRG for their constructive criticism and advice on this article. We acknowledge Veronica A. Raker for article editing. We thank the CNAG (National Center for Genomic Analysis, Barcelona) Genomic facility and CRG flow Cytometry unit for all technical support. This work was supported by grants to P.S. from the French National Research Agency (ANR) Young Investigator grant (ANR-21-CE12-0010), French Cancer Research Foundation (ARCPJA22021050003683), La Ligue Foundation for Cancer, France (Haute-Garonne, 285458), National Institute of Health and Medical Research (INSERM) Young Recruitment Support (U1194SHA), Cancéropôle Grand Sud-Ouest collaboration grant (R21031FF), Centre National de la Recherche Scientifique (CNRS), Université Toulouse III - Paul Sabatier, Toulouse, France. This work was also supported by a grant to M.B. from the European Research Council Synergy Grant “4DGenome” (609989).

### AUTHOR CONTRIBUTIONS

Conceptualization, P.S.; methodology, J.L.V., N.F.F., and P.S.; investigation, J.L.V., N.F.F., B.O., M.B., and P.S.; formal analysis, J.L.V., N.F.F., and P.S.; data curation, J.L.V., N.F.F., M.E.R., C.T., F.P.P., B.O., C.N., and P.S.; project administration, P.S., writing-original draft, P.S. writing-review and editing, J.L.V., N.F.F., B.O., M.E.R., C.T., F.P.P., M.B., C.N., and P.S.; funding acquisition, P.S.; supervision, P.S.

### DECLARATION OF INTERESTS

The authors declare no competing interests.

Received: December 21, 2023

Revised: February 13, 2024

Accepted: March 25, 2024

Published: March 27, 2024

## REFERENCES

- van Venrooij, W.J., and Puij, G.J.M. (2000). Citrullination: A small change for a protein with great consequences for rheumatoid arthritis. *Arthritis Res.* 2, 249–251. <https://doi.org/10.1186/ar95>.
- Wang, S., and Wang, Y. (2013). Peptidylarginine deiminases in citrullination, gene regulation, health and pathogenesis. *Biochim Biophys Acta - Gene Regul Mech* 1829, 1126–1135. <https://doi.org/10.1016/j.bbagr.2013.07.003>.
- Fuhrmann, J., Clancy, K.W., and Thompson, P.R. (2015). Chemical Biology of Protein Arginine Modifications in Epigenetic Regulation. *Chem. Rev.* 115, 5413–5461. <https://doi.org/10.1021/acs.chemrev.5b00003>.
- Christophorou, M.A., Sharma, P., Zhang, X., and Christophorou, M.A. (2023). Citrullination : New Tricks for an Old Mod, pp. 4–7.
- Christophorou, M.A. (2022). The virtues and vices of protein citrullination. *R. Soc. Open Sci.* 9, 220125–220195. <https://doi.org/10.1098/rsos.220125>.
- Vossenaar, E.R., Zendman, A.J.W., Van Venrooij, W.J., and Puij, G.J.M. (2003). PAD, a growing family of citrullinating enzymes: Genes, features and involvement in disease. *Bioessays* 25, 1106–1118. <https://doi.org/10.1002/bies.10357>.
- Tanikawa, C., Ueda, K., Suzuki, A., Iida, A., Nakamura, R., Atsuta, N., Tohnai, G., Sobue, G., Saichi, N., Momozawa, Y., et al. (2018). Citrullination of RGG Motifs in FET Proteins by PAD4 Regulates Protein Aggregation and ALS Susceptibility. *Cell Rep.* 22, 1473–1483. <https://doi.org/10.1016/j.celrep.2018.01.031>.
- Christophorou, M.A., Castelo-Branco, G., Halley-Stott, R.P., Oliveira, C.S., Loos, R., Radzishewska, A., Mowen, K.A., Bertone, P., Silva, J.C.R., Zermicka-Goetz, M., et al. (2014). Citrullination regulates pluripotency and histone H1 binding to chromatin. *Nature* 507, 104–108. <https://doi.org/10.1038/nature12942>.
- Sharma, P., Azebi, S., England, P., Christensen, T., Møller-Larsen, A., Petersen, T., Batsché, E., and Muchardt, C. (2012). Citrullination of Histone H3 Interferes with HP1-Mediated Transcriptional Repression. *PLoS Genet.* 8, 10029344–e1003015. <https://doi.org/10.1371/journal.pgen.1002934>.
- Chang, X., Xia, Y., Pan, J., Meng, Q., Zhao, Y., and Yan, X. (2013). PAD2 is significantly associated with rheumatoid arthritis. *PLoS One* 8, e81259. <https://doi.org/10.1371/journal.pone.0081259>.
- Ishigami, A., Ohsawa, T., Hiratsuka, M., Taguchi, H., Kobayashi, S., Saito, Y., Murayama, S., Asaga, H., Toda, T., Kimura, N., and Maruyama, N. (2005). Abnormal accumulation of citrullinated proteins catalyzed by peptidylarginine deiminase in hippocampal extracts from patients with Alzheimer's disease. *J. Neurosci. Res.* 80, 120–128. <https://doi.org/10.1002/jnr.20431>.
- Arif, M., and Kato, T. (2009). Increased expression of PAD2 after repeated intracerebroventricular infusions of soluble Aβ25–35 in the Alzheimer's disease model rat brain: Effect of memantine. *Cell. Mol. Biol. Lett.* 14, 703–714. <https://doi.org/10.2478/s11658-009-0029-x>.
- Cherrington, B.D., Zhang, X., Mcelwee, J.L., Morency, E., Anguish, L.J., and Coonrod, S.A. (2012). Potential Role for PAD2 in Gene Regulation in Breast Cancer Cells. *PLoS One* 7, e41242. <https://doi.org/10.1371/journal.pone.0041242>.
- Mohanani, S., Cherrington, B.D., Horibata, S., Mcelwee, J.L., Thompson, P.R., and Coonrod, S.A. (2012). Potential Role of Peptidylarginine Deiminase Enzymes and Protein Citrullination in Cancer Pathogenesis. *Biochem. Res. Int.* 2012, 895343. <https://doi.org/10.1155/2012/895343>.
- Guo, W., Zheng, Y., Xu, B., Ma, F., Li, C., Zhang, X., Wang, Y., and Chang, X. (2017). Investigating the expression , effect and tumorigenic pathway of PAD2 in tumors. *OncoTargets Ther.* 10, 1475–1485.
- Wang, L., Song, G., Zhang, X., Feng, T., Pan, J., Chen, W., Yang, M., Bai, X., Pang, Y., Yu, J., et al. (2017). PAD2-mediated citrullination promotes prostate cancer progression. *Cancer Res.* 77, 5755–5768. <https://doi.org/10.1158/0008-5472.CAN-17-0150>.
- Horibata, S., Rogers, K.E., Sadegh, D., Anguish, L.J., Mcelwee, J.L., Shah, P., Thompson, P.R., and Coonrod, S.A. (2017). Role of peptidylarginine deiminase 2 (PAD2) in mammary carcinoma cell migration. *BMC Cancer* 17, 378. <https://doi.org/10.1186/s12885-017-3354-x>.
- Song, S., and Yu, Y. (2019). Progression on citrullination of proteins in gastrointestinal cancers. *Front. Oncol.* 9, 15–16. <https://doi.org/10.3389/fonc.2019.00015>.
- Yuzhalin, A.E. (2019). Citrullination in cancer. *Cancer Res.* 79, 1274–1284. <https://doi.org/10.1158/0008-5472.CAN-18-2797>.
- Beato, M., and Sharma, P. (2020). Peptidyl arginine deiminase 2 (PAD2)-mediated arginine citrullination modulates transcription in cancer. *Int. J. Mol. Sci.* 21, 1351–1416. <https://doi.org/10.3390/ijms21041351>.
- Gao, B.S., Rong, C.S., Xu, H.M., Sun, T., Hou, J., and Xu, Y. (2020). Peptidyl Arginine Deiminase, Type II (PAD2) Is Involved in Urothelial Bladder Cancer. *Pathol. Oncol. Res.* 26, 1279–1285. <https://doi.org/10.1007/s12253-019-00687-0>.
- Xue, T., Liu, X., Zhang, M., Qiukai, E., Liu, S., Zou, M., Li, Y., Ma, Z., Han, Y., Thompson, P., and Zhang, X. (2021). PAD2-Catalyzed MEK1 Citrullination Activates ERK1/2 and Promotes IGF2BP1-Mediated SOX2 mRNA Stability in Endometrial Cancer. *Adv. Sci.* 8, 2002831–2002917. <https://doi.org/10.1002/adv.202002831>.
- Darrah, E. (2012). Peptidylarginine deiminase 2, 3 and 4 have distinct specificities against cellular substrates: Novel insights into autoantigen selection in rheumatoid arthritis. *Erika. Ann. Rheum. Dis.* 71, 92–98. <https://doi.org/10.1136/ard.2011.151712>.
- Dreyton, C.J., Knuckley, B., Jones, J.E., Lewallen, D.M., and Thompson, P.R. (2014). Mechanistic studies of protein arginine deiminase 2: Evidence for a substrate-assisted mechanism. *Biochemistry* 53, 4426–4433. <https://doi.org/10.1021/bi500554b>.
- Slade, D.J., Fang, P., Dreyton, C.J., Zhang, Y., Fuhrmann, J., Rempel, D., Bax, B.D., Coonrod, S.A., Lewis, H.D., Guo, M., et al. (2015). Protein arginine deiminase 2 binds calcium in an ordered fashion: Implications for inhibitor design. *ACS Chem. Biol.* 10, 1043–1053. <https://doi.org/10.1021/cb500933j>.
- Zhang, X., Bolt, M., Guertin, M.J., Chen, W., Zhang, S., Cherrington, B.D., Slade, D.J., Dreyton, C.J., Subramanian, V., Bicker, K.L., et al. (2012). Peptidylarginine deiminase 2-catalyzed histone H3 arginine 26 citrullination facilitates estrogen receptor  $\alpha$  target gene activation. *Proc. Natl. Acad. Sci. USA* 109, 13331–13336. <https://doi.org/10.1073/pnas.1203280109>.
- Guertin, M.J., Zhang, X., Anguish, L., Kim, S., Varticovski, L., Lis, J.T., Hager, G.L., and Coonrod, S.A. (2014). Targeted H3R26 Deimination Specifically Facilitates Estrogen Receptor Binding by Modifying Nucleosome Structure. *PLoS Genet.* 10, 10046133–e1004712. <https://doi.org/10.1371/journal.pgen.1004613>.
- Xiao, S., Lu, J., Sidrah, B., Cao, X., Yu, P., Zhao, T., Chen, C.C., McDermott, D., Sloopman, L., Wang, Y., et al. (2017). SMARCD1 Contributes to the Regulation of Naive Pluripotency by Interacting with Histone Citrullination. *Cell Rep.* 18, 3117–3128. <https://doi.org/10.1016/j.celrep.2017.02.070>.
- Falcão, A.M., Meijer, M., Scaglione, A., Rinwa, P., Agirre, E., Liang, J., Larsen, S.C., Heskol, A., Frawley, R., Klingener, M., et al. (2019). PAD2-Mediated Citrullination Contributes to Efficient Oligodendrocyte Differentiation and Myelination. *Cell Rep.* 27, 1090–1102.e10. <https://doi.org/10.1016/j.celrep.2019.03.108>.
- Sharma, P., Lioutas, A., Fernandez-Fuentes, N., Quilez, J., Carbonell-Caballero, J., Wright, R.H.G., Di Vona, C., Le Dily, F., Schüller, R., Eick, D., et al. (2019). Arginine Citrullination at the C-Terminal Domain Controls RNA Polymerase II Transcription. *Mol. Cell.* 73, 84–96.e7. <https://doi.org/10.1016/j.molcel.2018.10.016>.
- Cummings, T.F.M., Gori, K., Sanchez-Pulido, L., Gavriilidis, G., Moi, D., Wilson, A.R., Murchison, E., Dessimoz, C., Ponting, C.P., and Christophorou, M.A. (2022). Citrullination Was Introduced into Animals by Horizontal Gene Transfer from Cyanobacteria. *Mol. Biol. Evol.* 39, msab317. <https://doi.org/10.1093/molbev/msab317>.
- Crisp, A., Boschetti, C., Perry, M., Tunnacliffe, A., and Mickleth, G. (2015). Expression of multiple horizontally acquired genes is a hallmark of both vertebrate and invertebrate genomes. *Genome Biol.* 16, 50. <https://doi.org/10.1186/s13059-015-0607-3>.
- Emamalipour, M., Seidi, K., Zununi Vahed, S., Jahanban-Esfahlan, A., Jaymand, M., Majidi,

- H., Amoozgar, Z., Chitkushev, L.T., Javaheri, T., Jahanban-Esfahlan, R., and Zare, P. (2020). Horizontal Gene Transfer: From Evolutionary Flexibility to Disease Progression. *Front. Cell Dev. Biol.* 8, 229. <https://doi.org/10.3389/fcell.2020.00229>.
34. Soucy, S.M., Huang, J., and Gogarten, J.P. (2015). Horizontal gene transfer: Building the web of life. *Nat. Rev. Genet.* 16, 472–482. <https://doi.org/10.1038/nrg3962>.
35. Husnik, F., and McCutcheon, J.P. (2018). Functional horizontal gene transfer from bacteria to eukaryotes. *Nat. Rev. Microbiol.* 16, 67–79. <https://doi.org/10.1038/nrmicro.2017.137>.
36. Goulas, T., Mizgalska, D., Garcia-Ferrer, I., Kantyka, T., Guevara, T., Szmigielski, B., Sroka, A., Millán, C., Usón, I., Veillard, F., et al. (2015). Structure and mechanism of a bacterial host-protein citrullinating virulence factor, *Porphyromonas gingivalis* peptidylarginine deiminase. *Sci. Rep.* 5, 11969–12017. <https://doi.org/10.1038/srep11969>.
37. Touz, M.C., Rópolo, A.S., Rivero, M.R., Vranych, C.V., Conrad, J.T., Svard, S.G., and Nash, T.E. (2008). Arginine deiminase has multiple regulatory roles in the biology of *Giardia lamblia*. *J. Cell Sci.* 121, 2930–2938. <https://doi.org/10.1242/jcs.026963>.
38. Mistry, J., Chuguransky, S., Williams, L., Qureshi, M., Salazar, G.A., Sonnhammer, E.L.L., Tosatto, S.C.E., Paladini, L., Raj, S., Richardson, L.J., et al. (2021). Pfam: The protein families database in 2021. *Nucleic Acids Res.* 49, D412–D419. <https://doi.org/10.1093/nar/gkaa913>.
39. Asaga, H., and Ishigami, A. (2001). Protein deimination in the rat brain after kainate administration: Citrulline-containing proteins as a novel marker of neurodegeneration. *Neurosci. Lett.* 299, 5–8. [https://doi.org/10.1016/S0304-3940\(00\)01735-3](https://doi.org/10.1016/S0304-3940(00)01735-3).
40. Rogers, G., Winter, B., McLaughlan, C., Powell, B., and Nesci, T. (1997). Peptidylarginine deiminase of the hair follicle: Characterization, localization, and function in keratinizing tissues. *J. Invest. Dermatol.* 108, 700–707. <https://doi.org/10.1111/1523-1747.ep12292083>.
41. Rus'd, A.A., Ikejiri, Y., Ono, H., Yonekawa, T., Shiraiwa, M., Kawada, A., and Takahara, H. (1999). Molecular cloning of cDNAs of mouse peptidylarginine deiminase type I, type III and type IV, and the expression pattern of type I in mouse. *Eur. J. Biochem.* 259, 660–669. <https://doi.org/10.1046/j.1432-1327.1999.00083.x>.
42. Beltrao, P., Bork, P., Krogan, N.J., and Van Noort, V. (2013). Evolution and functional cross-talk of protein post-translational modifications. *Mol. Syst. Biol.* 9, 714–813. <https://doi.org/10.1002/msb.201304521>.
43. Kumar, S., Stecher, G., Suleski, M., and Hedges, S.B. (2017). TimeTree: A Resource for Timelines, Timetrees, and Divergence Times. *Mol. Biol. Evol.* 34, 1812–1819. <https://doi.org/10.1093/molbev/msx116>.
44. Corden, J.L. (2019). An Arginine Nexus in the RNA Polymerase II CTD. *Mol. Cell.* 73, 3–4. <https://doi.org/10.1016/j.molcel.2018.12.013>.
45. Marshall, N.F., and Price, D.H. (1995). Purification of P-TEFb, a transcription factor required for the transition into productive elongation. *J. Biol. Chem.* 270, 12335–12338. <https://doi.org/10.1074/jbc.270.21.12335>.
46. Adelman, K., and Lis, J.T. (2012). Promoter-proximal pausing of RNA polymerase II: emerging roles in metazoans. *Nat. Rev. Genet.* 13, 720–731. <https://doi.org/10.1038/nrg3293>.
47. Jonkers, I., and Lis, J.T. (2015). Getting up to speed with transcription elongation by RNA polymerase II. *Nat. Rev. Mol. Cell Biol.* 16, 167–177. <https://doi.org/10.1038/nrm3953>.
48. Core, L., and Adelman, K. (2019). Promoter-proximal pausing of RNA polymerase II: A nexus of gene regulation. *Genes Dev.* 33, 960–982. <https://doi.org/10.1101/gad.325142.119>.
49. Chen, F.X., Smith, E.R., and Shilatifard, A. (2018). Born to run: Control of transcription elongation by RNA polymerase II. *Nat. Rev. Mol. Cell Biol.* 19, 464–478. <https://doi.org/10.1038/s41580-018-0010-5>.
50. Martin, R.D., Hébert, T.E., and Tanny, J.C. (2020). Therapeutic targeting of the general RNA polymerase II transcription machinery. *Int. J. Mol. Sci.* 21, 3354. <https://doi.org/10.3390/ijms21093354>.
51. Muniz, L., Nicolas, E., and Trouche, D. (2021). RNA polymerase II speed: a key player in controlling and adapting transcriptome composition. *EMBO J.* 40, 1057400–e105821. <https://doi.org/10.15252/emboj.2020105740>.
52. Altenhoff, A.M., Glover, N.M., Train, C.M., Kaleb, K., Warwick-Vesztry, A., Dylus, D., de Farias, T.M., Zile, K., Stevenson, C., Long, J., et al. (2018). The OMA orthology database in 2018: Retrieving evolutionary relationships among all domains of life through richer web and programmatic interfaces. *Nucleic Acids Res.* 46, D477–D485. <https://doi.org/10.1093/nar/gkx1019>.
53. Yang, Z., and Bielawski, J. (2000). Statistical methods for detecting molecular adaptation. *Trends Ecol. Evol.* 15, 496–503. [https://doi.org/10.1016/S0169-5347\(00\)01994-7](https://doi.org/10.1016/S0169-5347(00)01994-7).
54. Zhang, J., Nielsen, R., and Yang, Z. (2005). Evaluation of an improved branch-site likelihood method for detecting positive selection at the molecular level. *Mol. Biol. Evol.* 22, 2472–2479. <https://doi.org/10.1093/molbev/msi237>.
55. Kosiol, C., Vinař, T., Da Fonseca, R.R., Hubisz, M.J., Bustamante, C.D., Nielsen, R., and Siepel, A. (2008). Patterns of positive selection in six mammalian genomes. *PLoS Genet.* 4, e1000144. <https://doi.org/10.1371/journal.pgen.1000144>.
56. Segura, J., Jones, P.F., and Fernandez-Fuentes, N. (2012). A holistic in silico approach to predict functional sites in protein structures. *Bioinformatics* 28, 1845–1850. <https://doi.org/10.1093/bioinformatics/bts269>.
57. Jones, S., Heyningen, P.V., Berman, H.M., and Thornton, J.M. (1999). *Arcturus\_Summer\_Hba*.
58. Segura, J., Oliva, B., and Fernandez-Fuentes, N. (2012). CAPS-DB: A structural classification of helix-capping motifs. *Nucleic Acids Res.* 40, 479–485. <https://doi.org/10.1093/nar/gkr879>.
59. Baumli, S., Lolli, G., Lowe, E.D., Troiani, S., Rusconi, L., Bullock, A.N., Debreczeni, J.E., Knapp, S., and Johnson, L.N. (2008). The structure of P-TEFb (CDK9/cyclin T1), its complex with flavopiridol and regulation by phosphorylation. *EMBO J.* 27, 1907–1918. <https://doi.org/10.1038/emboj.2008.121>.
60. Evans, R., O'Neill, M., Pritzel, A., Antropova, N., Senior, A.G.T., Židek, A., Bates, R., Blackwell, S., Yim, J., Ronneberger, O.B.S., et al. (2022). Protein complex prediction with AlphaFold-Multimer. Preprint at bioRxiv. <https://doi.org/10.1101/2021.10.04.463034>.
61. Jumper, J., Evans, R., Pritzel, A., Green, T., Figurnov, M., Ronneberger, O., Tunyasuvunakool, K., Bates, R., Židek, A., Potapenko, A., et al. (2021). Highly accurate protein structure prediction with AlphaFold. *Nature* 596, 583–589. <https://doi.org/10.1038/s41586-021-03819-2>.
62. Lin, F.R., Niparko, J.K., and Ferrucci, L. (2014). 基因的改变 NIH Public Access. *Bone* 23, 1–7. <https://doi.org/10.1146/annurev-biochem-052610-095910.RNA>.
63. Kohoutek, J. (2009). P-TEFb- The final frontier. *Cell Div.* 4, 19. <https://doi.org/10.1186/1747-1028-4-19>.
64. Manuscript, A. (2009). multi-tasking P-TEFb 20, 334–340. <https://doi.org/10.1016/j.ccb.2008.04.008>.
65. Gabay, M., Li, Y., and Felsher, D.W. (2014). MYC activation is a hallmark of cancer initiation and maintenance. *Cold Spring Harb. Perspect. Med.* 4, a014241. <https://doi.org/10.1101/cshperspect.a014241>.
66. Mollaoglu, G., Guthrie, M.R., Böhm, S., Brägelmann, J., Can, I., Ballieu, P.M., Marx, A., George, J., Heinen, C., Chalishazar, M.D., et al. (2017). MYC Drives Progression of Small Cell Lung Cancer to a Variant Neuroendocrine Subtype with Vulnerability to Aurora Kinase Inhibition. *Cancer Cell* 31, 270–285. <https://doi.org/10.1016/j.ccell.2016.12.005>.
67. Qiu, X., Boufaied, N., Hallal, T., Feit, A., de Polo, A., Luoma, A.M., Alahmadi, W., Larocque, J., Zadra, G., Xie, Y., et al. (2022). MYC drives aggressive prostate cancer by disrupting transcriptional pause release at androgen receptor targets. *Nat. Commun.* 13, 2559.
68. Cui, Z., Xiao, L., Chen, F., Wang, J., Lin, H., Li, D., and Wu, Z. (2021). High mRNA Expression of CENPL and Its Significance in Prognosis of Hepatocellular Carcinoma Patients. *Dis. Markers* 2021, 9971799. <https://doi.org/10.1155/2021/9971799>.
69. Zhang, H., Zhang, X., Li, X., Meng, W.B., Bai, Z.T., Rui, S.Z., Wang, Z.F., Zhou, W.C., and Jin, X.D. (2018). Effect of CCNB1 silencing on cell cycle, senescence, and apoptosis through the p53 signaling pathway in pancreatic cancer. *J. Cell. Physiol.* 234, 619–631. <https://doi.org/10.1002/jcp.26816>.
70. Chen, E.B., Qin, X., Peng, K., Li, Q., Tang, C., Wei, Y.C., Yu, S., Gan, L., and Liu, T.S. (2019). HnRNPR-CCNB1/CENPF axis contributes to gastric cancer proliferation and metastasis. *Aging (Albany NY)* 11, 7473–7491. <https://doi.org/10.18632/aging.102254>.
71. Liu, C., Kudo, T., Ye, X., and Gascoigne, K. (2023). Cell-to-cell variability in Myc dynamics drives transcriptional heterogeneity in cancer cells. *Cell Rep.* 42, 112401. <https://doi.org/10.1016/j.celrep.2023.112401>.
72. Jones, K.E., Bielby, J., Cardillo, M., Fritz, S.A., O'Dell, J., Orme, C.D.L., Safi, K., Sechrest, W., Boakes, E.H., Carbone, C., et al. (2009). PanTHERIA: a species-level database of life history, ecology, and geography of extant and recently extinct mammals. *Ecology* 90, 2648. <https://doi.org/10.1890/08-1494.1>.
73. Subramanian, A., Tamayo, P., Mootha, V.K., Mukherjee, S., Ebert, B.L., Gillette, M.A., Paulovich, A., Pomeroy, S.L., Golub, T.R., Lander, E.S., and Mesirov, J.P. (2005). Gene set enrichment analysis: A knowledge-based approach for interpreting genome-wide expression profiles. *Proc. Natl. Acad. Sci. USA* 102, 15545–15550. <https://doi.org/10.1073/pnas.0506580102>.

74. Veidenberg, A., Medlar, A., and Löytynoja, A. (2016). Wasabi: An integrated platform for evolutionary sequence analysis and data visualization. *Mol. Biol. Evol.* 33, 1126–1130. <https://doi.org/10.1093/molbev/msv333>.
75. Katoh, K., Misawa, K., Kuma, K.I., and Miyata, T. (2002). MAFFT: A novel method for rapid multiple sequence alignment based on fast Fourier transform. *Nucleic Acids Res.* 30, 3059–3066. <https://doi.org/10.1093/nar/gkf436>.
76. Larsson, A. (2014). AliView: A fast and lightweight alignment viewer and editor for large datasets. *Bioinformatics* 30, 3276–3278. <https://doi.org/10.1093/bioinformatics/btu531>.
77. Capella-Gutiérrez, S., Silla-Martínez, J.M., and Gabaldón, T. (2009). trimAl: A tool for automated alignment trimming in large-scale phylogenetic analyses. *Bioinformatics* 25, 1972–1973. <https://doi.org/10.1093/bioinformatics/btp348>.
78. Suyama, M., Torrents, D., and Bork, P. (2006). PAL2NAL: Robust conversion of protein sequence alignments into the corresponding codon alignments. *Nucleic Acids Res.* 34, 609–612. <https://doi.org/10.1093/nar/gkl315>.
79. Stamatakis, A. (2014). RAxML version 8: A tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics* 30, 1312–1313. <https://doi.org/10.1093/bioinformatics/btu033>.
80. Huerta-Cepas, J., Serra, F., and Bork, P. (2016). ETE 3: Reconstruction, Analysis, and Visualization of Phylogenomic Data. *Mol. Biol. Evol.* 33, 1635–1638. <https://doi.org/10.1093/molbev/msw046>.
81. Pronk, S., Páll, S., Schulz, R., Larsson, P., Bjelkmar, P., Apostolov, R., Shirts, M.R., Smith, J.C., Kasson, P.M., van der Spoel, D., et al. (2013). GROMACS 4.5: A high-throughput and highly parallel open source molecular simulation toolkit. *Bioinformatics* 29, 845–854. <https://doi.org/10.1093/bioinformatics/btt055>.
82. Segura, J., Marín-López, M.A., Jones, P.F., Oliva, B., and Fernandez-Fuentes, N. (2015). VORFFIP-driven dock: V-D2OCK, a fast, accurate protein docking strategy. *PLoS One* 10, 1–12. <https://doi.org/10.1371/journal.pone.0118107>.
83. Gibson, D.G., Young, L., Chuang, R.Y., Venter, J.C., Hutchison, C.A., and Smith, H.O. (2009). Enzymatic assembly of DNA molecules up to several hundred kilobases. *Nat. Methods* 6, 343–345. <https://doi.org/10.1038/nmeth.1318>.
84. Vicent, G.P., Zaurin, R., Nacht, A.S., Li, A., Font-Mateu, J., Le Dily, F., Vermeulen, M., Mann, M., and Beato, M. (2009). Two chromatin remodeling activities cooperate during activation of hormone responsive promoters. *PLoS Genet.* 5, e1000567. <https://doi.org/10.1371/journal.pgen.1000567>.
85. Daly, M.J., Patterson, N., Mesirov, J.P., Golub, T.R., Tamayo, P., and Spiegelman, B. (2003). PGC-1 $\alpha$ -responsive genes involved in oxidative phosphorylation are coordinately downregulated in human diabetes. *Nat. Genet.* 34, 267–273.
86. Shannon, P., Markiel, A., Owen, O., Baliga, N.S., Jonathan, T., Wang, D.R., Amin, N., Schwikowski, B., and Cytoscape, T.I. (2003). A Software Environment for Integrated Models. *Genome Res.* 13, 426. <https://doi.org/10.1101/gr.1239303.metabolite>.
87. Villanueva-Cañas, J.L., Ruiz-Orera, J., Agea, M.I., Gallo, M., Andreu, D., and Albà, M.M. (2017). New genes and functional innovation in mammals. *Genome Biol Evol* 9, 1886–1900. <https://doi.org/10.1093/gbe/evx136>.
88. Villanueva-Cañas, J.L., Laurie, S., and Albà, M.M. (2013). Improving genome-wide scans of positive selection by using protein isoforms of similar length. *Genome Biol. Evol.* 5, 457–467. <https://doi.org/10.1093/gbe/evt017>.
89. Yang, Z. (2007). PAML 4: Phylogenetic analysis by maximum likelihood. *Mol. Biol. Evol.* 24, 1586–1591. <https://doi.org/10.1093/molbev/msm088>.
90. Feliu, E., Aloy, P., and Oliva, B. (2011). On the analysis of protein-protein interactions via knowledge-based potentials for the prediction of protein-protein docking. *Protein Sci.* 20, 529–541. <https://doi.org/10.1002/pro.585>.
91. Pei, J., and Grishin, N.V. (2001). AL2CO: Calculation of positional conservation in a protein sequence alignment. *Bioinformatics* 17, 700–712. <https://doi.org/10.1093/bioinformatics/17.8.700>.

## STAR★METHODS

### KEY RESOURCES TABLE

REAGENT or RESOURCE	SOURCE	IDENTIFIER
<b>Antibodies</b>		
Mouse monoclonal anti-GFP	Roche	11814460001
Rabbit polyclonal anti-CCNT1	Bethyl Labs	A303-499A
Mouse monoclonal anti-GFP	Merck	11814460001
Rabbit polyclonal anti-CDK9 (D-7)	Santa Cruz Biot.	SC-13130
Rabbit polyclonal anti-CCNT1	Bethyl Labs	A303-499A
IgG mouse	Merck	12-371
IgG Rabbit	Cell Signaling	2729S
<b>Chemicals, peptides, and recombinant proteins</b>		
Lipofectamine 3000	Invitrogen	L3000008
Proteinase K	ThermoFisher Scientific	AM2546
TURBO™ DNase	ThermoFisher Scientific	AM2239
DNase I (RNase-free)	ThermoFisher Scientific	AM2222
Cell proliferation ELISA BrdU Colorimetric assay	Roche	11647229001
Dynabeads™ M-280 Sheep Anti-Mouse IgG	ThermoFisher	11201D
Protein G Plus / Protein A Agarose	Millipore	IP05
<b>Deposited data</b>		
Single Cell RNA-seq	This study	GSE246420
<b>Experimental models: Cell lines</b>		
HeLa cells	ATCC	CCL-2
<b>Oligonucleotides</b>		
PADI2- T159A Forward 5'TGGTGAAGTGTGACCGAGAGGCACCCT GGTTGCCCAAGGAGGACTGCCGTGATG 3'	This study	
PADI2- T159A Reverse 5'CATCACGGCAGTCCTCCTTGGGCAACCA GGGTGCTCTCGGTACAGTTCACCA 3'.		
PADI2- W161A Forward 5'TGGTGAAGTGTGACCGAGAGACACCCGC GTTGCCCAAGGAGGACTGCCGTGATG 3'	This study	
PADI2- W161A Reverse 5'CATCACGGCAGTCCTCCTTGGGCAACGC GGGTGCTCTCGGTACAGTTCACCA 3'.		
PADI2- L162A Forward 5'TGGTGAAGTGTGACCGAGAGACACCCTG GGCACCAAGGAGGACTGCCGTGATG3'	This study	
PADI2- L162A Reverse 5'CATCACGGCAGTCCTCCTTGGGTGCCCA GGGTGCTCTCGGTACAGTTCACCA3'.		

(Continued on next page)

**Continued**

REAGENT or RESOURCE	SOURCE	IDENTIFIER
PADI2- L162A/ W161A Forward 5'TGGTGAAGTGTGACCGAGAGACACCCGCC GCACCCAAGGAGGACTGCCGTGATG 3'	This study	
PADI2- L162A/ W161A Reverse 5'CATCACGGCAGTCCTCCTTGGGTGCGGCGG GTGCTCTCGGTACAGTTCACCA 3'		
PADI2- L162A/W161A/T159A Forward 5'TGGTGAAGTGTGACCGAGAGGCA <del>C</del> CCC GCCGCACCCAAGGAGGACTGCCGTGATG -3'	This study	
PADI2- L162A/W161A/T159A Reverse 5'CATCACGGCAGTCCTCCTTGGGTGCGGCGG GTGCTCTCGGTACAGTTCACCA 3'		

**Recombinant DNA**

Human PADI2 wild type gene cloned in pCPR0032 vector	This study	WT-PADI2
L162A mutation in WT-PADI2	This study	L162A-PADI2
W161A mutation in L162A-PADI2	This study	L162A/W161A-PADI2
T159A mutation in L162A/W161A-PADI2	This study	L162A/W161A/T159A-PADI2

**Software and algorithms**

MVORFFIP	Segura et al. <sup>56,58</sup>	<a href="http://www.bioinsilico.org/MVORFFIP">http://www.bioinsilico.org/MVORFFIP</a>
OMA Database	Altenhoff et al. <sup>52</sup>	<a href="https://omabrowser.org/oma/home/">https://omabrowser.org/oma/home/</a>
GSEA	Subramanian et al. <sup>73</sup>	<a href="http://software.broadinstitute.org/gsea/index.jsp">http://software.broadinstitute.org/gsea/index.jsp</a>
Wasabi	Veidenberg et al. <sup>74</sup>	<a href="http://wasabiapp.org/">http://wasabiapp.org/</a>
MAFFT	Katoh et al. <sup>75</sup>	<a href="https://mafft.cbrc.jp/alignment/server/index.html">https://mafft.cbrc.jp/alignment/server/index.html</a>
AliView	Larsson et al. <sup>76</sup>	<a href="http://ormbunkar.se/aliview/">http://ormbunkar.se/aliview/</a>
trimAl	Capella-Gutiérrez et al. <sup>77</sup>	<a href="http://trimal.cgenomics.org/">http://trimal.cgenomics.org/</a>
PAL2NAL	Suyama et al. <sup>78</sup>	<a href="https://www.bork.embl.de/pal2nal/">https://www.bork.embl.de/pal2nal/</a>
RAXML	Stamatakis et al. <sup>79</sup>	<a href="https://cme.h-its.org/exelixis/web/software/raxml/">https://cme.h-its.org/exelixis/web/software/raxml/</a>
ETE 3	Huerta-Cepas et al. <sup>80</sup>	<a href="http://et toolkit.org/">http://et toolkit.org/</a>
GROMACS	Pronk et al. <sup>81</sup>	<a href="http://www.gromacs.org">http://www.gromacs.org</a>
VD2OCK	Segura et al. <sup>82</sup>	<a href="http://www.bioinsilico.org/cgi-bin/VD2OCK/staticHTML/home">http://www.bioinsilico.org/cgi-bin/VD2OCK/staticHTML/home</a>

**RESOURCE AVAILABILITY**

**Lead contact**

Further information and requests for resources and results should be directed to the lead contact Priyanka Sharma ([priyanka.sharma@ipbs.fr](mailto:priyanka.sharma@ipbs.fr)).

**Materials availability**

Further information and requests for plasmids generated in this study will be available upon reasonable request to the Lead Contact.

**Data and code availability**

- Single-cell RNA sequencing data performed and used in this study have been deposited at GEO under the accession code GSE246420 and publically available.
- No original code was generated in this paper.
- Any additional information required to reanalyze the data reported in this paper is available from the [lead contact](#) upon request.

## EXPERIMENTAL MODEL AND SUBJECT DETAILS

### Cell lines

HeLa cells (ATCC CCL-2) were grown in DMEM with 10% FBS and 100U/ml penicillin-streptomycin according to the ATCC's recommendations. Cells were transfected using Lipofectamine 3000 (Invitrogen) according to the manufacturer's instructions.

## METHOD DETAILS

### PADI2-GFP plasmids and fluorescence-activated cell sorting (FACS) sorting

PADI2 was cloned into the pCPR0032 GFP-tagged vector using the forward primer (F) 5'- AGAACCTGTA~~CTT~~CCAATCCATGCTGCGCG AGCGGAC-3' and the reverse primer (R) 5'- GATCCGTATCCACCTTTACTTTAGGGCACCATGTGCCACC-3' using Gibson assembly.<sup>83</sup> The selected plasmid sequence was verified by Sanger sequencing. Next, PADI2 wild-type (WT) was used to generate a single mutant (T159A or W161A or L162A), a double mutant (L162A/W161A), and a triple mutant (L162A/ W161A/ T159A), using the following primers:

PADI2- T159A Forward

5'- TGGTGAAGTGTGACCGAGAGGCACCCTGGTTGCCCAAGGAGGACTGCCGTGATG -3'

PADI2- T159A Reverse

5'- CATCACGGCAGTCCTCCTTGGGCAACCAGGGTGCCTCTCGGTCACAGTTCACCA -3'.

PADI2- W161A Forward

5'- TGGTGAAGTGTGACCGAGAGACACCCGCGTTGCCCAAGGAGGACTGCCGTGATG -3'

PADI2- W161A Reverse

5'- CATCACGGCAGTCCTCCTTGGGCAACGCGGGTGTCTCTCGGTCACAGTTCACCA -3'.

PADI2- L162A Forward

5'- TGGTGAAGTGTGACCGAGAGACACCCTGGGCACCCAAGGAGGACTGCCGTGATG-3'

PADI2- L162A Reverse

5'-CATCACGGCAGTCCTCCTTGGGTGCCCAGGGTGTCTCTCGGTCACAGTTCACCA-3'.

PADI2- L162A/ W161A Forward

5'- TGGTGAAGTGTGACCGAGAGACACCCGCCGCACCCAAGGAGGACTGCCGTGATG -3'

PADI2- L162A/ W161A Reverse

5'- CATCACGGCAGTCCTCCTTGGGTGCCGCGGGTGTCTCTCGGTCACAGTTCACCA -3'.

PADI2- L162A/W161A/T159A Forward

5'- TGGTGAAGTGTGACCGAGAGGCCACCCGCCGCACCCAAGGAGGACTGCCGTGATG -3'

PADI2- L162A/W161A/T159A Reverse

5'- CATCACGGCAGTCCTCCTTGGGTGCCGCGGGTGCCTCTCGGTCACAGTTCACCA-3'.

The corresponding fragments were generated by each pair of primers using oligonucleotide assembly and introduced by Gibson assembly.<sup>83</sup> All the generated mutants were confirmed by Sanger sequencing.

For transfection,  $2 \times 10^6$  HeLa cells were seeded in 10-cm plates, and 4  $\mu$ g of the plasmid being tested was transfected using Lipofectamine 3000 (Invitrogen) for 24 hours according to the manufacturer's instructions. Cells were trypsinized, and GFP-positive live cells were sorted using BD influx (Becton and Dickinson, San Jose, CA). Briefly, cells were stained with 1 $\mu$ g/mL concentration of DAPI (4', 6-diamidino-2-phenylindole) before FACS sorting. A SSC-H (side scatter height) versus FSC-H (forward scatter height), morphological-related parameters dot-plot was used to exclude debris by gating cells; doublets were then excluded using a FSC-H versus FSC-A (forward scatter area) by gating and dead cells were excluded using DAPI versus FSC-A dot-plot by gating living cells. GFP-positive cells were identified and isolated using a G4 gate in GFP versus autofluorescence (AF) dot-plot. Obtained data were analyzed using Flow Jo 10. 6. GFP-positive cells were used for experiments. Cells were centrifuged and stored as pellets at  $-80^\circ\text{C}$  prior to RNA extraction and immunoprecipitation experiments.

### Cell proliferation assay

#### *BrdU (5'-bromo-2'-deoxyuridine) cell proliferation assay*

HeLa cells ( $0.3 \times 10^3$ ) were transfected with either wild-type PADI2 (WT) or mutant PADI2 (with L642A or T159 or W161A or L642A/W161A, or L642A/W161A/T159A) using Lipofectamine 3000 (Invitrogen) in 96-well plates. The cell proliferation ELISA BrdU (5'-bromo-2'-deoxyuridine) colorimetric assay (Roche,11647229001) was performed as per the manufacturer's instructions. The experiments were performed on at least eight biological replicates.

### Incucyte<sup>®</sup> proliferation assays for live-cell analysis

HeLa cells seeded in 96-well plates with 300 cells per well, transfected with either wild-type PADI2 (WT) or mutant PADI2 (with L642A or T159 or W161A or L642A/W161A, or L642A/W161A/T159A). After 24hours of transfection, imaging was performed using the IncuCyte live cell imaging system (Essen BioScience). Scans at 4 $\times$  magnification were taken every 8 hours for 5 days. Cell confluence was calculated from microscopy images using the Incucyte software algorithm to generate a proliferation index corresponding to the change in confluence for each well. These measurements are the mean  $\pm$  SEM of at least six replicates.

### RNA extraction and RT-qPCRs

RNA from HeLa cells transfected with either wild-type PADI2 (WT) or mutant PADI2 (with L642A, or L642A/W161A, or L642A/W161A/T159A) was extracted using RNeasy (Qiagen) according to the manufacturer's instructions. Purified RNA (1 µg) was used for DNase treatment (Thermo Scientific) and quantified with a Qubit 3.0 Fluorometer (Life Technologies).

Reverse transcription of RNA was performed using a qScript™ cDNA Synthesis Kit (Quanta Bioscience 95047-100) according to the manufacturer's instructions. Complementary DNA was quantified by qPCR using Roche Lightcycler (Roche) as described.<sup>84</sup> For each gene product, relative RNA abundance was calculated using the standard curve method and expressed as relative RNA abundance after normalizing against the human *GAPDH* gene level. All gene expression data generated by RT-qPCR represented the average and ± SEM of at least three biological replicates. Primers used for RT-qPCR are listed in Table S1.

### Single-cell samples, library preparation, and sequencing

HeLa cells ( $0.6 \times 10^6$ ) were transfected with either wild-type PADI2 (WT) or the PADI2 triple mutant (with L642A/W161A/T159A) using Lipofectamine 3000 (Invitrogen) in 60mm plates. GFP-positive singlet cells and live cells were suspended in 1 ml of PBS+BSA 0.05%. Cells were centrifuged at 400 rcf for 5 min at 4°C in order to bring the cell concentration to 300-1000 cells/µl. Cell concentration and viability were determined by manual counting using a Neubauer chamber and staining the cells with Trypan blue. Cells were partitioned into Gel Bead-In-Emulsions (GEMs) using the Chromium Controller system (10X Genomics), with a target recovery of 5000 total cells per sample. cDNA sequencing libraries were prepared using the Next GEM Single Cell 3' Reagent Kits v3.1 (10X Genomics, PN-1000268), following the manufacturer's instructions. Briefly, after GEM-RT clean-up, cDNA was amplified for 11 cycles, and cDNA QC and quantification were performed on an Agilent Bioanalyzer High Sensitivity chip (Agilent Technologies). cDNA libraries were indexed by PCR using the PN-1000215 Dual Index Kit TT Set A Plate. Size distribution and concentration of 3' cDNA libraries were verified on an Agilent Bioanalyzer High Sensitivity chip (Agilent Technologies). Finally, sequencing of cDNA libraries was carried out using the Illumina NovaSeq 6000.

### Single-cell RNA sequencing data processing and analysis

Transcriptome-wide analysis of human L162A/W161A/T159A-PADI2 mutant and WT cells was performed at single-cell resolution using NovaSeq 6000 pipeline. Sequencing data were aligned to the human reference genome Grch38. Data with at least 500 unique molecular identifiers (UMIs), log<sub>10</sub> genes per UMI >0.8, >250 genes per cell and a mitochondrial ratio of less than 20% were extracted, normalized, and integrated using the Seurat package v4.0.3 in R4.0.2. After quality control and integration, we performed a modularity-optimized Louvain clustering. After quality control, 17,139 WT and 32,364 PADI2mut cells remained. Altogether, we used two biological replicates (PADI2\_mut1: 16,177 cells, PADI2mut2: 16,187 cells, WT1: 8,568 cells, WT2: 8,571 cells).

### Gene ontology (GO) and regulatory gene set analysis

GO Annotation and regulatory gene set analysis were performed using the online tool Gene Set Enrichment Analysis (GSEA, <http://software.broadinstitute.org/gsea/index.jsp>) collection database v5.<sup>73,85</sup> The significant cut-off p-value and FDR q-value < 0.05. Plots were done with the use of Prism (GraphPad Prism 10.0.3 for MacOS), Cytoscape<sup>86</sup> (version 3.9.1) and R4.0.2.

### ChIP-qPCRs

For ChIP assays,<sup>30</sup>  $4 \times 10^6$  of FACS sorted GFP-positive HeLa cells (WT-PADI2, L162A-PADI2, L162A/W161A-PADI2, L162A/W161A/T159A-PADI2) were cross-linked for 10 min with 1% formaldehyde at 37°C. The chromatin lysate was sonicated to a DNA fragment size range of 100-200bp using a Biorupter sonicator (Diagenode). CDK9 was immunoprecipitated with 15 µg of an anti-CDK9 antibody (D-7, sc-13130, lot no # B1422) or control mouse IgG (12-371, Merck) in IP Buffer with 2X SDS buffer (100mM NaCl, 50mM Tris-HCl, pH8, 5mM EDTA and 0.5% SDS) and 1X Triton buffer (100mM Tris-HCl, pH8.8, 100mM NaCl, 5mM EDTA and 5% Triton-X) with protease inhibitors (11836170001, Roche) for 16 hours at 4°C. This step was followed by incubation with 50 µl of Dynabeads® M-280 sheep anti-mouse IgG (11201D, Thermo Scientific) for 3 hours. Beads were washed 3 times with low salt buffer (140mM NaCl, 50mM HEPES, pH 7.4, 1% Triton-X 100), 2 times with high salt buffer (500 mM NaCl, 50mM HEPES, pH 7.4, 1% Triton-X 100) followed by single wash with LiCl Buffer (10mM Tris HCl pH 8.0, 250 mM LiCl, 1% NP-40, 1% sodium deoxycholic acid and 1mM EDTA) and 1 × TE buffer at 4°C. Subsequently, crosslinks were reversed at 65°C overnight, followed by RNase treatment for 1.5 hours, and bound DNA was purified by Phenol-Chloroform extraction. The resulting eluted DNA was quantified by Qubit 3.0 Fluorometer (Life Technologies) and followed by real-time qPCR analysis. Data are represented as fold-change over input fraction from at least 3 biological replicate experiments. Primers used for qPCR are listed in Table S1.

### GFP-tagged PADI2 immunoprecipitation, western blot

Briefly,  $3 \times 10^6$  FACS sorted GFP-positive (WT-PADI2, L162A-PADI2, L162A/W161A-PADI2, L162A/W161A/T159A-PADI2) cells were lysed on ice for 30 min in lysis buffer (1% Triton X-100 in 50mM Tris pH 7.4–7.6, 130 mM NaCl) containing protease inhibitors (11836170001, Roche) with rotation, followed by sonication for 7 min with every 30 sec on / 30 sec off. After centrifugation at 4°C and 13,000 rpm for 10 min, extracts were used for protein quantitation. For the immunoprecipitation (IP) assay, 2mg of extract was incubated for 12 hours with 100 µl Dynabeads Protein A (10002D, Thermo Scientific). The monoclonal antibody anti-GFP (polyclonal rabbit, A-11122, Invitrogen) or a control rabbit IgG (2729S, Cell Signaling) was coupled with Dynabeads before incubation with extract at 4°C. The samples were washed 10 times with lysis buffer and boiled

for 5 min in SDS gel sample buffer. Proteins were visualized by 4% to 12% SDS-PAGE gels and western blotting, using anti-GFP (11814460001, Roche), anti-CDK9 (sc-13130, Santa Cruz), or anti-CCNT1 (A303-499A, Bethyl Labs) were used for western blots.

### Multiple sequence alignments

All of the homolog PADI sequences available in the OMA database were downloaded.<sup>52</sup> The species for the study were selected using three criteria: (i) a good representation of the different mammalian lineages, (ii) good sequence quality, and (iii) the presence of a full PADI2 sequence. Three bird species (chicken, turkey, and duck) were used as outgroups. The protein sequences were aligned using the software PRANK<sup>74</sup> and a pruned mammalian guide tree with branch distances.<sup>87</sup> This program uses an evolutionary model to place insertions and deletions, minimizing the over-alignment of non-homologous regions, and has been shown to improve dN/dS estimates.<sup>88</sup> Local realignments of two regions were done using MAFFT<sup>75</sup> within AliView.<sup>76</sup> Another MSA using only one-to-one orthologous PADI2 sequences was also built. TrimAl<sup>77</sup> was applied to the PADI family MSA, enabling the -automated1 option, as recommended before a phylogenetic reconstruction. The PADI2 protein MSA was then converted into a nucleotide alignment using the script Pal2Nal. The corresponding cDNA sequences for each species were gathered from the OMA database and converted using the script Pal2Nal.<sup>78</sup> The MSAs generated and the trees used are available as Supplementary Material S1.

### Phylogenetic reconstruction

We reconstructed the PADI family tree using RaxML<sup>79</sup> with 200 bootstrap values (-N 200), with an optimal amino acid substitution model chosen automatically (-m PROTGAMMAAUTO), and a rapid bootstrap algorithm (-f a) along with reproducible seed values (-p 12345, -x 345). Two reconstructions were performed: one with the original PADI family MSA and another one with the trimmed version of the alignment.

### Estimation of dN/dS and positive selection

Estimates were made for both the number of non-synonymous substitutions per non-synonymous site (dN) and the number of synonymous substitutions per synonymous site (dS) using the free-ratio model in CodeML.<sup>89</sup> For each branch and leaf in the tree, we performed a branch-site test of positive selection,<sup>54</sup> implemented in the Phylogenetic Analysis by Maximum Likelihood (PAML) software package<sup>89</sup> and available in the environment for tree exploration (ETE3) framework as bsA/bsA1.<sup>80</sup> A likelihood ratio (LRT) was calculated as  $2 \times (L1 - L0)$ , where L1 and L0 are the maximum likelihood value for the alternative hypothesis and the null hypothesis respectively. A chi-squared distribution with 1 degree of freedom was used to calculate the p-values. Only the positions for which the p-value (< 0.05) is significant for positive selection are reported.

### Characterization of positively selected residues using MVORFFIP

The prediction of interface residues was done using MVORFFIP.<sup>56</sup> MVORFFIP is a structure-based prediction method that identifies protein-, peptide-, RNA- and DNA interfaces based on a range of structural, evolutionary, experimental, and energy-based information integrated by a Random Forest classifier. The structure of PADI2 was submitted to the MVORFFIP server (<http://www.bioinsilico.org/MVORFFIP>)<sup>56</sup> and protein interface scores were assigned to individual amino acids.

### Structural modeling of PADI2-CDK9/CCNT1 and selection of putative interface residues

The structure of the trimer complex, PADI2-CDK9, was derived using protein docking as follows. The crystal structure of PADI2 (PDB code: 4n2a)<sup>25</sup> and CDK9/CCNT1 (PDB code: 3blr)<sup>59</sup> were available. The protein docking was performed using VD<sup>2</sup>OOCK,<sup>82</sup> generating the ensemble of over 10000 docking poses. Docking conformers were clustered using the GROMACS<sup>81</sup> and sorted by the ES3DC potential.<sup>90</sup> The top 200 poses were selected to identify the top 10 putative interface residues: for this, a composite score was calculated using the docking information, the MVORFFIP score and sequence conservation. Among the top 200 docking poses, as discussed above, a normalized score (Z score) was calculated for each exposed residue in PADI2. The mean and standard deviations to calculate the Z score were computed over the entire docking space, i.e., considering all docking poses. The MVORFFIP score was also computed as shown above. Finally, the conservation of each residue was computed using al2co.<sup>91</sup>

### AlphaFold-multimer modeling of PADI2-CDK9/CCNT1 complex

Structural models of the trimer PADI2-CDK9/CCNT1 were generated using AlphaFold-Multimer<sup>60</sup> as follows. The sequences were retrieved as mentioned above, PADI2 (PDB code: 4n2a)<sup>25</sup> and CDK9/CCNT1 (PDB code: 3blr).<sup>59</sup> The structure was derived using model parameters corresponding to "model\_1\_multimer".

## QUANTIFICATION AND STATISTICAL ANALYSIS

For RT-qPCR, ChIP-qPCR, and cell proliferation experiments, a Two-tailed unpaired Student's t-test was used to determine statistical significance between the groups. For cell proliferation experiments, the significance between groups calculated by Wilcoxon-Mann-Whitney test. Plots and indicated statistical analysis were done with the use of Prism (GraphPad Prism 10.0.3 for MacOS) unless otherwise stated. If exact p-values are not shown or indicated in legend then p-values are represented in all figures as follows: \*, p-value  $\leq$  0.05; \*\*, p-value  $\leq$  0.01; \*\*\*, p-value  $\leq$  0.001; °, p-value > 0.05.