



Most protein domains exist as variants with distinct functions across cells, tissues and diseases

Vitting-Seerup, Kristoffer

Published in:
NAR Genomics and Bioinformatics

Link to article, DOI:
[10.1093/nargab/lqad084](https://doi.org/10.1093/nargab/lqad084)

Publication date:
2023

Document Version
Publisher's PDF, also known as Version of record

[Link back to DTU Orbit](#)

Citation (APA):
Vitting-Seerup, K. (2023). Most protein domains exist as variants with distinct functions across cells, tissues and diseases. *NAR Genomics and Bioinformatics*, 5(3), Article lqad084. <https://doi.org/10.1093/nargab/lqad084>

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

Most protein domains exist as variants with distinct functions across cells, tissues and diseases

Kristoffer Vitting-Seerup *

The Bioinformatics Section, Department of Health Technology, The Technical University of Denmark (DTU), Denmark

*To whom correspondence should be addressed. Email: k.vitting.seerup@gmail.com

ABSTRACT

Protein domains are the active subunits that provide proteins with specific functions through precise three-dimensional structures. Such domains facilitate most protein functions, including molecular interactions and signal transduction. Currently, these protein domains are described and analyzed as invariable molecular building blocks with fixed functions. Here, I show that most human protein domains exist as multiple distinct variants termed ‘domain isotypes’. Domain isotypes are used in a cell, tissue and disease-specific manner and have surprisingly different 3D structures. Accordingly, domain isotypes, compared to each other, modulate or abolish the functionality of protein domains. These results challenge the current view of protein domains as invariable building blocks and have significant implications for both wet- and dry-lab workflows. The extensive use of protein domain isotypes within protein isoforms adds to the literature indicating we need to transition to an isoform-centric research paradigm.

INTRODUCTION

Most proteins have, throughout evolution, been created using protein domains as molecular building blocks (1). The importance of these protein domains is hard to overstate as they encode the core functionality needed for most if not all, cellular functions. Noteworthy examples include signal transduction (2) and the ability of proteins to bind to DNA, RNA and other proteins (3). Current state-of-the-art annotation of protein domains (e.g. Pfam (4)) relies on encapsulating known domain sequences into a reference model and matching those models to a sequence of interest. Domain annotation tools are used in many types of analysis and range from mechanistic studies of proteins to interpreting the impact of mutations and are cited by thousands of papers each year.

While protein domains are defined, described and analyzed as protein subunits with fixed functionality (4,5) that might not be the case. Numerous mutational and cancer studies find that removing just a tiny part of a protein domain can either eliminate or modify domain function (6–8). Such findings lead to an intriguing hypothesis: just like most genes produce protein isoforms with distinct functions (9,10), naturally occurring protein domain variants could exist. Such domain variants would originate from alternative splicing or evolution, resulting in different isoforms or genes (paralogs) with variants of the same protein domain.

Here, I term this phenomenon ‘domain isotypes’ and show that in humans, they are ubiquitous, used in a cell, tissue and disease-specific manner, have distinct 3D structures and can modify the biological function of a protein domain.

MATERIALS AND METHODS

Data and statistical analysis

All analysis was done in R 4.0.4 unless explicitly stated. Significance is in all analyses defined as Benjamini-Hochberg (BH) (11) (also known as False Discovery Rate, FDR) corrected *P*-values smaller than 0.05.

Domain analysis

All domains analyzed in this article were identified by running Pfam’s (4) `pfam_scan.pl` tool on the corresponding amino acid sequences only using the manually curated pfamA database (4) and default parameters. All results in this article are only based on protein domains defined by the ‘type’ column of the `pfam_scan` output (meaning no families, repeats, etc., were considered).

Calculating domain variation sizes

To analyze domain sizes, I use the information available in the output of `pfam_scan.pl` that summarizes the alignment between a reference model (hmm) and a query sequence (seq). Specifically the hmm alignment start and end coordinates (hmm_start and hmm_end), the corresponding sequencing alignment start and end coordinates (seq_start and seq_end) as well as the reference domain length (hmm_length). Truncation sizes for each end of the domain are calculated as respectively `hmm_start - 1` and `hmm_end - hmm_length`, and the largest value is used as the truncation size. The truncation fraction is calculated by dividing the truncation size by the `hmm_length`.

Indel sizes are extracted by comparing the length of the aligned sequence with the reference model. Specifically, the length of the aligned reference model is calculated as `hmm_end - hmm_start + 1`, and the length of the aligned se-

Received: June 14, 2023. Revised: August 9, 2023. Editorial Decision: September 2, 2023. Accepted: September 5, 2023

© The Author(s) 2023. Published by Oxford University Press on behalf of NAR Genomics and Bioinformatics.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

quence is calculated as $\text{seq_end} - \text{seq_start} + 1$. The indel size is calculated by subtracting the aligned sequence length from the aligned model length, and the indel fraction is calculated by dividing by the model length (hmm_length).

The overall domain variation is defined as the largest absolute fraction when comparing the truncation fraction and indel fraction. For the cluster analysis, I for each domain analyzed the lengths using the `Modes()` function from the `LaplacesDemon` R package modifying the ‘`min.size`’ argument as indicated in the figure (0.1–0.3). The `Modes` function finds the number of local maxima in the density distribution of the variables while ensuring a mode corresponds to a least the fraction of observations indicated by the ‘`min.size`’ argument (meaning at least 10–30% of observation in each cluster).

pfamAnalyzeR

I here implemented and used `pfamAnalyzeR` v. 0.1.0 (available at Bioconductor (12) <https://bioconductor.org/packages/pfamAnalyzeR/>), which currently has three functionalities: (i) it reads the fixed-with-file (fwf) format produced by `pfam_scan.pl` or the Pfam webserver into R taking known format problems and data versions into account; (ii) domain data is augmented with the domain truncation/indel calculations described above; (iii) each domain is assigned one of five mutually exclusive isotype categories. First, a domain is defined as having a truncation if the truncation fraction ≥ 0.1 (Supplementary Figure S1). Secondly, a domain is defined as having an indel if the absolute indel fraction is ≥ 0.1 (Supplementary Figure S1). Then the five mutually exclusive categories are defined as:

- 1) ‘Reference’: No truncation and no indel
- 2) ‘Complex’: Both truncation and indel
- 3) ‘Truncation’: Truncation but no indel
- 4) ‘Insertion’: No truncation but indel and indel fraction > 0
- 5) ‘Deletion’: No truncation but indel and indel fraction < 0

These classes are further simplified to ‘reference’ (#1 from above) and ‘non-reference’ (#2–#5 from above).

TCGA analysis

The TCGA isoform switches were directly obtained from the supplementary material of Vitting-Seerup *et al.* (13). I extended `IsoformSwitchAnalyzeR` (14) to incorporate results from `pfamAnalyzeR` and to compare the simplified isotype classification (reference versus non-reference) as a new functional consequence (available in `IsoformSwitchAnalyzeR` v1.17.06). The number of genes with isoform switches was extracted with the `extractConsequenceSummary()` function. All isoform visualizations were created with the `switchPlot()` and `switchPlotTranscript()` functions. Data is available in Supplementary Table S1.

Species analysis

The amino acid fasta file for the manual curated SwissProt subset of UniProt (15) was downloaded from UniProt’s website, downloading only ‘canonical sequences’ (meaning one ‘representative’ sequence per gene). The human file was downloaded on 11 October 2021. The *Saccharomyces cerevisiae*, *Mus musculus* and *Drosophila melanogaster* files were downloaded on 12 May 2022. All files were analyzed with

`pfam_scan` and `pfamAnalyzeR` as described above. Domains identified, along with their isotype, are available in Supplementary Table S2.

Tissue expression analysis

The GTEx transcript TPM expression for all 55 human tissues was downloaded from www.gtexportal.org using v8 data (16). For each tissue, the median transcript expression was calculated. Since the GTEx v8 data is quantified via GENCODE (17) v23 annotation, I obtained the corresponding amino acid from genencodegenes.org and analyzed it with `pfam_scan` and `pfamAnalyzeR` as described above (Supplementary Table S8). The major transcript is defined as the most expressed protein-coding transcript from a given gene. For Figure 1G, I analyzed genes with two different transcripts defined as major in at least one tissue. Domain switches were then identified as those with differences in which domain isotypes were encoded in the transcripts.

Cell type expression analysis

The differential transcript usage analysis across 9 human immune cell types identified from full-length single-cell data was downloaded from the supplementary data Hagemann-Jensen *et al.* (18) (the ‘noHEK’ analysis). Transcripts with usage changes were extracted by filtering for adjusted *P*-values < 0.05 . For each significant transcript, I identified the cell type where it was most expressed. For Figure 1G, I analyzed genes where two different transcripts are defined as major in two different cell types. Domain switches were then identified by analyzing the isotypes of the domains encoded in the corresponding GENCODE v23 transcripts.

Structural analysis

I used PDB’s (19) `batch_download.sh`, script to download all human pdb files 13 October 2021. From these pdb files, I extracted the amino acid sequence of the first chain (not considering X) using the `read.pdb()` function from the `bio3d` R package (20). The resulting fasta file was analyzed with `pfam_scan` and `pfamAnalyzeR` as described above. Domains identified, along with their isotype, are available in Supplementary Table S3.

To analyze the similarity of the 3D structures of protein domains I used pairwise comparisons. It was not computationally feasible to analyze all possible within-domain combinations as some had thousands of occurrences. Instead, for each domain (`hmm_name`), I made all pairwise comparisons of max 9 randomly selected reference isotype domains and max 21 randomly selected non-reference isotype domains. The selected reference isotype domains were also compared to each other to estimate the reference variability.

All structural comparisons were made using the `bio3d` R package (20) v 2.4.2. The pdb files were loaded into R using the `read.pdb()`. The correct chain and the subset of the structure containing the domain were extracted using the `atom.select()` function with the Pfam domain alignment coordinates. The structures were compared using the `struct.aln()` function without trimming optimization (`max.cycles = 0`). The `struct.aln()` function works by first making a sequence alignment and afterward using the aligned sub-sequence to rotate and transpose the 3D structures ensuring they are comparable. From this structural alignment, the `struct.aln()` function also calculates the backbone RMSD, and I manually ex-

tended the function to calculate the backbone TM-scores as defined in equations 1 and 5 in Zhang *et al.* (21). The sequence alignment was run with the following parameters: gapOpening = 40, gapExtension = 4, substitutionMatrix = BLOSUM50 (from the Biostrings package) to encourage an alignment better reflecting the expected input where large parts were missing from ends (truncation) or in the middle of the sequence (insertion/deletion).

To estimate the similarity expected from un-related domains, I randomly selected one 3D structure for 150 different domains, all classified as the reference isotype, and compared the structures in a pairwise manner considering all possible comparisons.

I only analyzed comparisons where both of the aligned sequences were > 20 amino acids and where the length of the non-reference aligned sub-sequence corresponded to >80% of the original sequence length ($n = 48\ 369$) (Supplementary Tables S4 and S5). To obtain the average structural differences I for each domain (hmm_name) and isotype extracted the mean RMSD and TM-score for the cases where at least 3 structural comparisons passed the filtering described above. Average TM-scores and RMSD values are highly negatively correlated (Spearman's rho: 0.891, P -value < 2.22e-16). The interval of unrelated RMSD/TM values was defined as the 5–95 percentile of the 11 026 unrelated comparisons described above.

For the His_Phos_1 example 3 reference isotype structures (PDB ID: 1K6M (22), 5HTK (23), 6HVH (24)) and 2 truncation isotype structures (PDB ID: 1YFK (25), 2A9J (26)) were manually selected. Note that these structures are from 5 different genes, and 'only' 5 structures were chosen because this was the maximal number, I could get the multiple alignment described below to work with. The pdb files and domain subsets were handled as described above. All 5 domain structures were aligned with the pdbaln() function using the abovementioned parameters. This function works the same as struct.aln(), except it uses multiple alignment to align 2+ sequences. Open-Source PyMOL (<https://github.com/schrodinger/pymol-open-source>) was used to visualize the aligned structures by saving them with pymol.pdbs() and afterward opening it in PyMOL. The catalytic residues were manually annotated by comparison to figure 3 in Rigden *et al.* (27) using the 1K6M structure as reference since it is included in both this and their analyses.

Domain interaction analysis

Protein domains interacting in solved 3D structures were downloaded from the newest version (2020_01) of the 3DID database (3) (<https://3did.irbbarcelona.org>). DNA and RNA binding domains are defined as domains returned by searching the Pfam website (<https://pfam.xfam.org/>) for, respectively, 'DNA binding' and 'RNA binding' keywords. See Supplementary Table S6.

For the interaction analysis, I only analyzed domains found in the human SwissProt database. For each interaction type, I used Fisher's exact test to test the overlap between each domain isotype and the domains associated with the interaction type.

Binding motif analysis

For the domain-motif interaction (DMI), I downloaded the 3DID DMI (3) from <https://3did.irbbarcelona.org>. The DMI data was reduced to the interaction containing the human pdb files and chains I used for the structural analysis described above. The data was further reduced to domains uniquely overlapping (in protein coordinates) of the same type (hmm_name), and the domain isotype was transferred from the pdb data to the DMI data.

For each domain (hmm_name), I clustered motifs using the distance matrix produced by stringdistmatrix() from the stringDist R package (28) as input to standard hierarchical clustering with hclust(). Clusters were defined using the cutree() function specifying $b = 4$. The interpretation is that clusters are sets of motifs where > 4 modifications (deletions, insertions, substitutions and transpositions) must be made to make the motifs similar to any motif outside the cluster. Considering that the median motif consists of six amino acids, this results in clusters with very different motifs.

To identify domains where isotypes were associated with distinct motifs clusters, I only considered motif clusters where for each domain (hmm_name), at least one domain isotype was identified at least three times. A motif cluster was defined as isotype-specific if 80% of motifs within the cluster interacted with either reference or non-reference domain isotypes. Of the 98 domains (hmm_name) with at least two motif clusters, I found 10 where reference and non-reference isotypes were specific to distinct motif clusters. A Fisher test was used to test the statistical significance of the Trypsin example using the 4×2 contingency table of observed interaction counts where '4' indicate the 4 motif clusters and '2' is reference vs. non-reference groups.

Differences in presence, absence, or isotype of interaction domains

To determine changes in the presence, absence, or isotype of interaction domains, I used Pfam to identify protein domains (as described above) in the GENCODE v38 (17) annotation. I then extracted genes containing the interaction protein domains (as defined above). I then counted how many genes had isoforms both with and without a protein domain ($n = 6638$) and how many genes had protein domains (same hmm_name) that overlapped in genomic coordinates but were classified as different domain isotypes by pfamAnalyzeR (as described above).

Gene-set association

I used enrichment analysis to associate specific domains (given by their hmm_name) with biological functions. First, I extracted gene sets from MSigDB (29,30) explicitly using the following collections: (i) hallmarks, C6 oncogenic, GO:BP Biological Processes part of gene ontology, KEGG pathways, Wikipathways and CPG collection containing gene sets describing chemical and genetic perturbations. From these, only gene sets with at least 10 and, at most, 1000 gene symbols were used. The domains identified in SwissProt were reduced to those found in at least 10 genes. Both gene sets and domains were reduced to genes found in both datasets and gene sets with <10 genes were again removed. For each

domain (hmm_name), I use the Fisher.test() function with alternative='greater' to evaluate the enrichment of the genes in which the domain was found among the genes in a gene set. For each domain, *P*-values were corrected for multiple FDR testing, and significant associations were defined as FDR <0.05.

For each gene set significantly associated with a domain, I tested if the same gene set was also significantly associated with either the reference or the non-reference isotype of that domain as follows. For each gene and domain (hmm_name), I calculated the average overall domain variation (see Domain Analysis section for that definition), where larger values mean larger deviation from the reference isotype. For each domain, I then extracted the subset of genes with that domain and ranked them according to the average overall domain variation in decreasing order. I then used fgsea (31) to test if the significant gene set was overrepresented at the top or bottom of the list (skew). The decreasing ranking means that positive NES scores indicate associations with non-reference isotypes, and negative NES scores indicate associations with reference isotypes. Gene sets were only tested if there was at least one gene in the ranked list that was not in the gene set. For each domain, *P*-values were corrected for multiple testing using FDR, and significant associations were defined as FDR <0.05. The gene sets associated with each domain (and isotype skew) are available in Supplementary Table S7.

Disease association

I obtained the Pfam domain associated with diseases from Supplementary Table 4 of Savojardo *et al.* (32) by filtering for associations > 0. Protein domains associated with cancers (oncodomains) were obtained from the supplementary data of Peterson *et al.* (33) and added to the list of disease-associated domains. The disease-associated domains were reduced to those found in the human SwissProt data. For each disease group, I then used a Fisher's exact test to test the overlap between each isotype and the domain associated with the disease group. Similarly, another test was performed where all non-reference isotypes were jointly analyzed. *P*-values were corrected for multiple testing using FDR, and significant associations were defined as FDR < 0.05.

RESULTS

Domain isotypes exist in humans

To investigate if alternative splicing creates domain isotypes in humans, I first analyzed dysregulated splicing in human cancers. Specifically, I focused on cases where tumors, compared to adjacent healthy tissue, had upregulated one isoform while downregulated another (often referred to as isoform switches) (13). I re-analyzed these isoform switches for cases where the difference between the isoforms was a partial domain loss/gain. Across the 12 cancer types analyzed, I identified 492 genes where an isoform switch resulted in substantial changes to, but not complete gain/removal of, the protein domain (Supplementary Figure S2, Supplementary Table S1). Notably, 204 of these genes were found in more than one cancer type (Supplementary Figure S3), highlighting their functional importance.

A prominent example is the isoform switch in the *STK4* gene (frequently called *MST1*) found in three cancer types:

thyroid, kidney and colorectal (Supplementary Figure S4–S6). *STK4* is a key component in the Hippo pathway (34) which functions as a tumor suppressor by phosphorylating and inactivating the YAP1 oncogene (35). Inactivation of *STK4* has been linked to poor patient prognosis in many cancer types, including colon and kidney cancer (36–38). The isoform switch I identified in *STK4* results in the central part of the Pkinase domain, including the active sites, being removed from the cancer isoform (Figure 1A). Since the Pkinase domain is responsible for the phosphorylation of YAP1, the cancer isoform most likely cannot inhibit YAP1, potentially leading to cancer progression, and suggests domain isotypes, created via alternative splicing, are important in human cancers.

Encouraged by these examples, I next asked if domain isotypes are also created through the evolution of genes. I, therefore, analyzed the length of all protein domains identified in the SwissProt canonical database (11076 domain-containing proteins, one isoform per gene) (15). Compared to the reference domain length annotated in the Pfam database (4), I found substantial length differences in the domains identified in human proteins (Figure 1B). Intriguingly many protein domains seem to follow individual patterns in how they deviate from the reference domain length (Supplementary Figure S7). To investigate this, I used the fact that most protein domains are found many times in the human proteome and extracted the length of the domain for each of these occurrences. For each protein domain, I then used cluster analysis to look for patterns in the length distribution of individual domains. For most protein domains (58.3%), the length of the domains separates into two or more clusters (Figures 1C, S8), suggesting domain isotypes are frequently created through the evolution of genes.

In summary, I find that alternative splicing and evolution frequently creates domain isotypes.

Enabling systematic analysis of domain isotypes

To systematically analyze domain isotypes, I developed pfamAnalyzeR, which enables detection and in-depth analysis of domain isotypes (available via Bioconductor (12) at <https://bioconductor.org/packages/pfamAnalyzeR/>). pfamAnalyzeR uses a highly stringent approach to detect domain changes (Supplementary Figure S1) and can distinguish between 5 categories of domain isotypes (Figure 1D). These isotypes are the reference isotype and four isotypes that, compared to the reference isotype, are best described as a truncation, an insertion, a deletion, or combinations thereof ('complex') (See schematic illustration in Figure 1D). Importantly the stringent cutoff used here also means that my results cannot be explained by erroneous Pfam domain identification such as those identified by Triant *et al.* (39) as those, if truly erroneous, represent <0.1% of the data I here analyze.

Domain isotypes are ubiquitous

I used pfamAnalyzeR to analyze all protein domains in the manually curated SwissProt database of multiple species (15). Intriguingly I find that most protein domains, as a consequence of evolution, exist as both reference and non-reference iso-types (Figure 1E, Supplementary Figures S9 and S10, Supplementary Table S2).

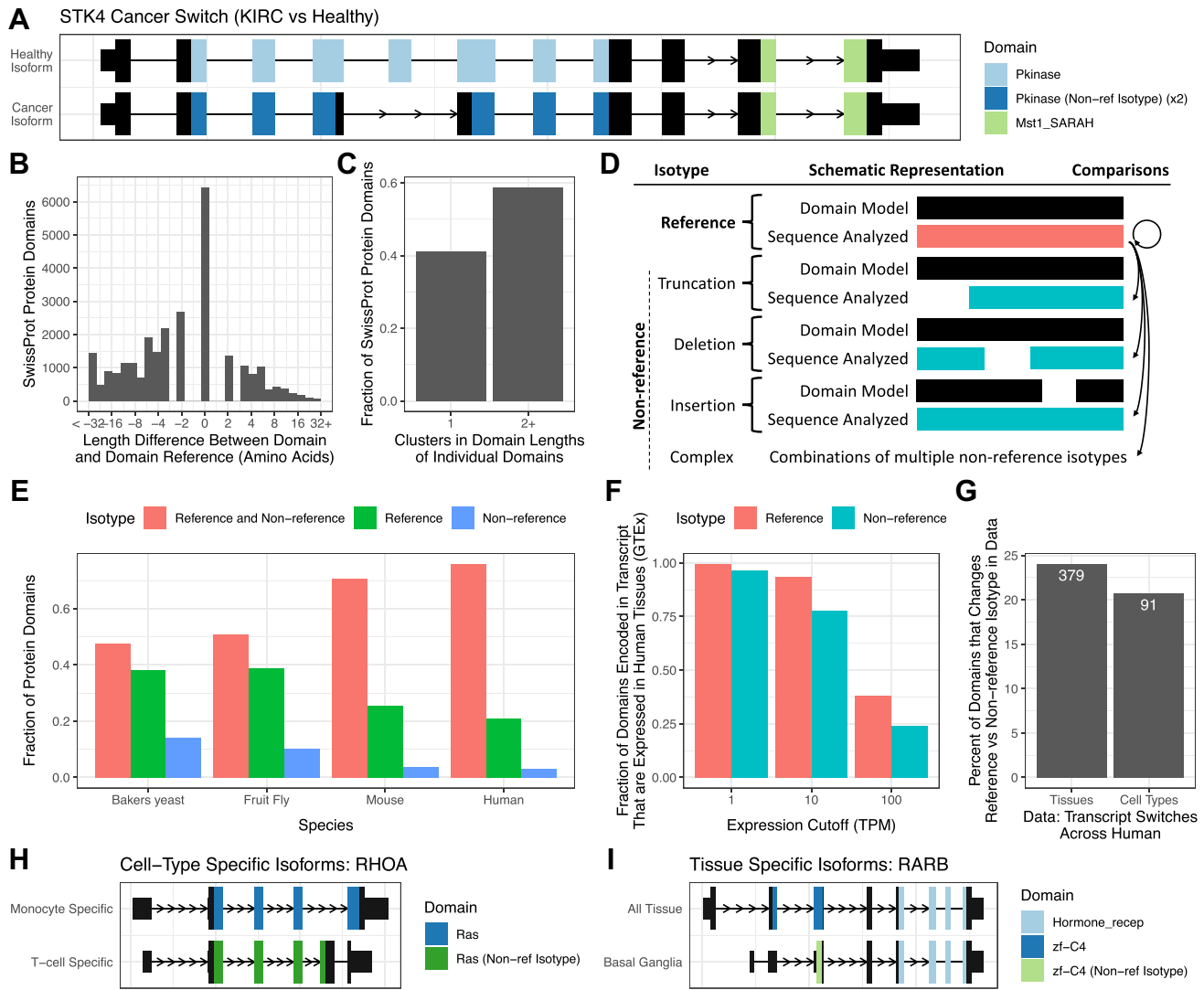


Figure 1. Domain isotypes are ubiquitous. **(A)** The isoform switch in STK4 identified by comparing Kidney Renal Cell Carcinoma (KIRC) to adjacent healthy tissue. In the Switch Plot, black boxes represent exons where the wider sections annotate the coding regions. Protein domains are indicated by color. The isoform usage is indicated on the left. Note the differences in the Pkinase domain. **(B)** The difference between the actual length of protein domains and their corresponding reference length across all human domains identified in SwissProt. **(C)** The number of clusters found in the length distribution of each human protein domain. **(D)** Schematic representation of the five different protein domain isotypes identified. Black boxes represent the reference model, and red/turquoise boxes indicate the sequence analyzed (red being the reference isotype and turquoise being non-reference isotypes). Arrows indicate comparisons analyzed in this paper; the non-reference isotypes are highlighted to the left. **(E)** The fraction of protein domains that exist as reference and non-reference isotypes in SwissProt of various species. **(F)** The fraction (y-axis) of protein domain isotypes (color) contained in transcripts expressed above a cutoff (x-axis) in the average expression profile of at least one human tissue. **(G)** The most expressed transcript was identified from each gene in each human tissue, and genes where the major transcript differed across tissues or cell types (x-axis) were extracted. From these transcripts, I calculated the percent of encoded protein domains (Y-axis) where there was a domain isotype switch due to the transcript switch. Numbers on the bars indicate the actual number of domains identified. **(H, I)** Switch plot as described in A, showing tissue- **(I)** and cell-type specific **(H)** domain isotype usage respectively found in the RARB **(J)** and RHOA **(I)**.

Domain isotypes are used in a cell and tissue-specific manner

To assess the biological relevance of domain isotypes, I examined the expression of mRNA transcripts containing either reference or non-reference domain isotypes across the 55 human tissues in the GTEx data (40 849 domain-contain isoforms from 11141 genes) (16). I find that reference and non-reference domain isotypes are found within expressed transcripts at comparable rates across expression thresholds (Figure 1F), suggesting non-reference domain isotypes are frequently used in human tissues. Accordingly, many non-

reference isotype domains were also found when only considering the most expressed (major) transcripts from each gene (Supplementary Figure S11). Interestingly within the major transcripts, where there are differences in which transcript is the most expressed across tissues, I find 379 (24%) protein domains where the most used domain isotype also changes across tissues (Figure 1G). Similarly, I found 91 (20.7%) domains with isotype switches across a small selection of 9 human immune cell types profiled at single-cell resolution (18) (Figure 1G). Intriguing examples are the cell-type specific use of the Ras domain isotypes (Figure 1H) and tissue specific use

of zink-finger domain isotypes (Figure 1I). RHOA is a GTPase important for a wide range of functions in the immune system (40). Essential for this function is the Ras domain used in most immune cell types. Interestingly we find that T-cells utilize a truncated domain isotype which lacks the last 30 amino acids including the 'G5' and the 'hypervariable region' responsible for many interactions (41). It is therefore an intriguing hypothesis that the Ras isotype differences are responsible for some of the already known cell-type specific functions of RHOA (40). RARB is a hormone receptor important for cell development and differentiation (42). RARB works by using its zinc-finger ('zf-C4') domain to bind DNA allowing the 'Hormone_recep' domain to work as a (co)factor to facilitate transcriptional changes. It is therefore peculiar that RARB seem to utilize a truncated zinc-finger domain isotype in basal ganglia which presumable would lead to a severely reduced ability to bind DNA. One fascinating possibility could be that this domain isotype enhance or enable DNA independent functions of RARB (43).

Thus, domain isotypes, created by alternative splicing, are not an exception but the norm and seem to be used in a cell and tissue-specific manner thereby contributing to cell-type and tissue specialization.

Domain isotypes have distinct 3D structures

Having shown that domain isotypes are prevalent, I next asked if different domain isotypes have different biological functions. Following the 'structure is function' axiom (44), I started with the three-dimensional (3D) structure of the protein domain. I analyzed all human protein domains within the experimentally solved 3D structures in the Protein Data Bank (PDB) (19) (Supplementary Table S3). For each protein domain, I compared the 3D structure of a reference domain isotype to the structure of other reference and non-reference isotypes (as illustrated by arrows in Figure 1D). In other words, I investigated if the 3D structure of protein domain isotypes had evolved through evolution. An interesting example is found in the structure of the histidine phosphatase domain 'His_Phos_1'. This domain removes phosphorylations from other proteins via its essential histidine residue ('His₈') aided by a few other residues, including 'His₁₀₈' (27). While the 'His_Phos_1' protein domains were virtually identical to each other within isotype category (mean RMSD_{aligned-backbone} < 1), the reference and truncated isotypes were strikingly different (mean RMSD_{aligned-backbone} > 11) (Figures 2A, S12). This difference is exacerbated by the absence of the conserved catalytical residues 'His₁₀₈' (27) from the truncated isotype (shown in blue in the reference isotype domains of Figures 2A and S12). The absence of a conserved catalytical residue naturally implies that the truncated isotype modifies the functionality of the domain(45). Importantly this is not just an artifact of the domain detection, as the residue is missing even when considering an additional 50 amino acids downstream of the domain boundary (Supplementary Figure S13). Encouraged by this example, I summarized the structural differences of domain isotypes across the 952 protein domains analyzed (8741 PDB files, Supplementary Table S4). I find the 3D structures of different domain isotypes are excessively different (Figures 2B, S14) (median RMSD_{aligned-backbone} > 3× higher than within reference isotype comparison, $P < 2.14 \times 10^{-2}$, Wilcoxon tests). Surprisingly specific categories of domain isotypes are associated with larger structural differences

that approach what is expected when comparing random domains (shaded area in Figure 2B, below dashed line Supplementary Figure S14, Supplementary Table S5). In summary, the 3D structure of domain isotypes, created through evolution, is surprisingly different, thus challenging the current view of domains as invariable evolutionary building blocks.

Domain isotypes affect biological function

Since interaction between proteins and other macromolecules (e.g. protein, DNA, RNA) requires a specific 3D structure, the structural differences of domain isoforms could modify such interactions (both positively and negatively). In support of this, I find that in SwissProt non-reference protein domain isotypes are highly enriched for protein domains that facilitate interactions with DNA, RNA and proteins (3) (overall odds ratio per type of interaction >1.93, P -value < 3.23×10^{-10}) (Figure 2C)(Supplementary Table S6). Indeed, by analyzing the specific amino acid motif that interacts with protein domains in solved 3D structures from DPB (3), I found 10 cases (10.2% of analyzed) where domain isotypes were significantly associated with different clusters of binding motifs (see Materials and Methods for details). One example is the trypsin binding domain, where the reference and truncation isotype bind two distinct motif clusters with high specificity (FDR = 8.68×10^{-18} , fisher test) (Figure 2D). Thus, protein domain isotypes created through evolution are important for determining the interaction between proteins and other macromolecules, making them particularly interesting from a drug-developmental point of view (46).

This is further emphasized since domain isotypes are also created through splicing. Indeed, the GENCODE v38 (17) annotation suggests the human genome contains at least 7735 protein-coding genes (38.71%) that encode protein isoforms, which differ in the presence, absence, or isotype of interaction domains (Supplementary Figures S15 and S16). Thus, protein domain isotypes likely also modulate the function of individual genes through alternative splicing.

The structure is function axiom also suggests that protein domain isotypes could have distinct biological functions. We, therefore, analyzed if individual SwissProt protein domains were overrepresented in the genes annotated to specific pathways or go-terms. I found a significant association between 552 protein domains and 6018 gene sets (FDR < 0.05, Supplementary Table S7), suggesting the protein domains are critical for the associated pathways. Next, I tested whether the significant gene-sets were preferentially enriched for reference or non-reference domain isotypes (see methods). I found that 77 protein domains had at least one gene set significantly skewed towards an isotype (Figure 2E, Supplementary Table S7), suggesting that a particular domain isotype was preferred in that pathway. One of these is the Ras protein domain, which is significantly associated with 198 gene-sets, of which 20 were significantly skewed towards a domain isotype (Figure 2F). Interestingly the genes with the reference isotype of the Ras domain seem to be preferentially involved in vesicle formation and release, whereas the genes containing the non-reference isotypes were associated with cell movement (Figure 2F). Thus, domain isotypes might facilitate different biological functions.

Next, I utilized that some protein domains have already been associated with human diseases through overrepresentation of either mutations or naturally oc-

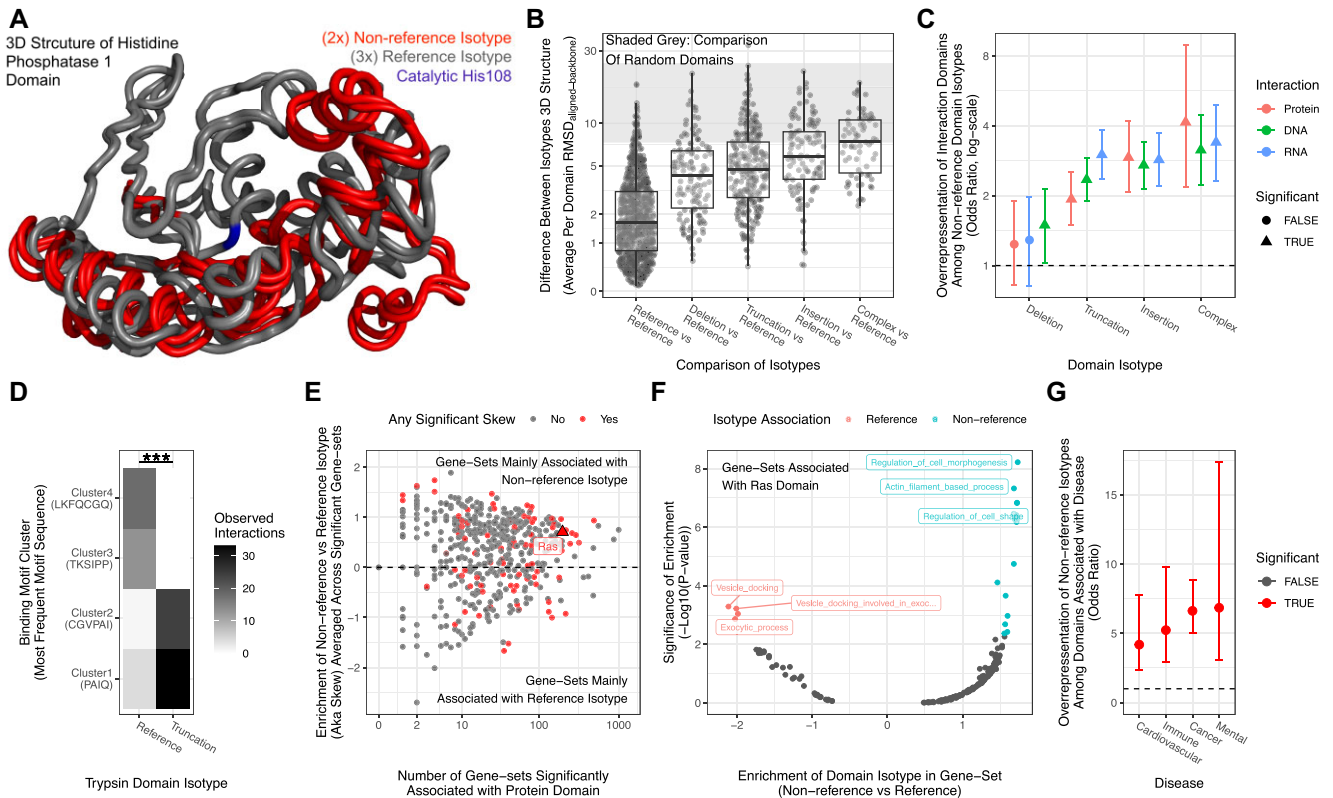


Figure 2. Domain isotypes are functionally important. **(A)** Structural alignment of three references (grey) and two truncated (red) isotypes of the His_Phos_1 protein domain. The catalytic residue His108 is indicated in blue. **(B)** For each protein domain, the 3D structures of reference isotypes were compared to the structure of all identified isotypes (x-axis). The structural difference was quantified RMSD and averaged per domain (Y-axis) (higher means more different). The grey area indicates the RMSD expected for random domains by showing where the middle 90% of random comparisons fall. **(C)** The enrichment (y-axis) of domains involved in interactions (color) for each domain isotype (x-axis). Triangles indicate FDR < 0.05. **(D)** The number of interactions (color) found between trypsin domain isotype (x-axis) and clusters of amino acid binding motifs (y-axis) in experimentally solved 3D structures. For each cluster, the most frequent motif is indicated in brackets. **(E)** For each protein domain, the number of significantly associated gene sets (x-axis)(FDR < 0.05) as a function of the median skew towards reference or non-reference domain isotype (y-axis). Color indicates if at least one gene set was significantly skewed towards an isotype (FDR < 0.05). **(F)** For each-gene set significantly associated with the RAS domain, the skew towards reference or non-reference domain isotype (x-axis) and the associated certainty of the shift (P-value, y-axis). Color indicates significant skew (FDR < 0.05), and top 3 skewed gene-sets in each direction are highlighted. **(G)** The enrichment of non-reference domain isotypes (y-axis) among domains associated with various disease classes (x-axis). Color denotes FDR < 0.05.

curing disease-related genetic variations (SNPs) (32,33). I find that non-reference protein domain isotypes are highly enriched for protein domains associated with all 17 disease groups tested (Supplementary Figure S17), with prominent examples being cardiovascular diseases, mental disorders, immune pathologies and cancer (all having odds ratio >4, P-value < 2.33e-07) (Figure 2G). Importantly this suggests that many non-reference domain isotypes are required for homeostasis and, when disrupted, can lead to a wide range of human diseases.

DISCUSSION

I show that protein domain isotypes are created through both alternative splicing and gene evolution (via paralogs). These domain isotypes are ubiquitous and modulate or disrupt the functionality of protein domains in a cell, tissue and disease-specific manner. Despite the natural dichotomy between alternative splicing and evolution, I here present a joint analysis of both mechanisms to comprehensively examine domain isotypes. Surprisingly these results highlight the synergy of con-

sidering both mechanisms since both show domain isotypes are prevalent and functionally important.

However, long-read technology consistently demonstrates that the current databases underrepresent protein diversity (47,48) meaning I probably underestimate the prevalence of domain isotypes. Additionally, naturally occurring genetic variation (SNPs) affect splicing in most human protein-coding genes (16), likely leading to additional domain isotype usage. Jointly this indicates that my results underestimate the true diversity and importance of domain isotypes. Therefore, considerable additional computational and experimental work is needed to characterize the functional role of domain isotypes, both as groups and for individual protein domains.

This article challenges the current perception of protein domains as building blocks with a fixed function. Instead, my findings suggest that protein domain isotypes are via alternative splicing or through evolution used to modulate domain function, thereby increasing the flexibility of protein adaptation. These findings also suggest caution when transferring conclusions about a protein domain from one setting to another (both other proteins, cells/tissues and isoforms) as it is likely that different domain isotypes are used, whereby func-

tional differences are expected. This is especially important when using protein domains for enzyme discovery (49) or drug design (46,50). Likewise extra caution is needed when using domain data to transfer information between species. This is especially relevant for resources such as Gene Ontologies which relies heavily on such information transfer (51).

The widespread use of protein domain isoforms also highlights, along with many other considerations (13,16,52,53), the need to change how we think about molecular biology. Although it is a considerable effort, we must progress from the current 'gene-centric' research paradigm to a more nuanced 'isoform-centric' paradigm.

DATA AVAILABILITY

All data analyzed and produced here are freely available. Data analyzed here are described and linked to in the appropriate methods sections. Data produced here are available in the supplementary tables.

SUPPLEMENTARY DATA

Supplementary Data are available at NARGAB Online.

ACKNOWLEDGEMENTS

Albin Sandlin for the encouraging discussion of the initial idea. For insightful feedback on manuscript drafts, Elena Papaleo, Lars Rønn Olsen, reviewers and Bo Porse.

FUNDING

Funded by a generous grant from the Lundbeck Foundation [R413-2022-878].

Conflict of interest statement

None declared.

REFERENCES

- Aziz, M.F. and Caetano-Anollés, G. (2021) Evolution of networks of protein domain organization. *Sci. Rep.*, **11**, 12075.
- Roskoski, R. (2004) Src protein-tyrosine kinase structure and regulation. *Biochem. Biophys. Res. Commun.*, **324**, 1155–1164.
- Mosca, R., Céol, A., Stein, A., Olivella, R. and Aloy, P. (2014) 3did: a catalog of domain-based interactions of known three-dimensional structure. *Nucleic Acids Res.*, **42**, D374–D379.
- Mistry, J., Chuguransky, S., Williams, L., Qureshi, M., Salazar, G.A., Sonnhammer, E.L.L., Tosatto, S.C.E., Paladin, L., Raj, S., Richardson, L.J., et al. (2020) Pfam: the protein families database in 2021. *Nucleic Acids Res.*, **49**, gkaa913.
- Jones, P., Binns, D., Chang, H.-Y., Fraser, M., Li, W., McAnulla, C., McWilliam, H., Maslen, J., Mitchell, A., Nuka, G., et al. (2014) InterProScan 5: genome-scale protein function classification. *Bioinformatics*, **30**, 1236–1240.
- Klimovich, B., Merle, N., Neumann, M., Elmshäuser, S., Nist, A., Mernberger, M., Kazdal, D., Stenzinger, A., Timofeev, O. and Stiewe, T. (2022) p53 partial loss-of-function mutations sensitize to chemotherapy. *Oncogene*, **41**, 1011–1023.
- Oren, M. and Rotter, V. (2010) Mutant p53 gain-of-function in cancer. *Csh. Perspect. Biol.*, **2**, a001107.
- Kato, S., Han, S.-Y., Liu, W., Otsuka, K., Shibata, H., Kanamaru, R. and Ishioka, C. (2003) Understanding the function–structure and function–mutation relationships of p53 tumor suppressor protein by high-resolution missense mutation analysis. *Proc. Natl. Acad. Sci. U.S.A.*, **100**, 8424–8429.
- Marasco, L.E. and Kornblihtt, A.R. (2023) The physiology of alternative splicing. *Nat. Rev. Mol. Cell Biol.*, **24**, 242–254.
- Wright, C.J., Smith, C.W.J. and Jiggins, C.D. (2022) Alternative splicing as a source of phenotypic diversity. *Nat. Rev. Genet.*, **23**, 697–710.
- Benjamini, Y. and Hochberg, Y. (1995) Controlling the false discovery rate: a practical and powerful approach to multiple testing. *J. Roy. Stat. Soc. Ser. B Methodol.*, **57**, 289–300.
- Huber, W., Carey, V.J., Gentleman, R., Anders, S., Carlson, M., Carvalho, B.S., Bravo, H.C., Davis, S., Gatto, L., Girke, T., et al. (2015) Orchestrating high-throughput genomic analysis with Bioconductor. *Nat. Methods*, **12**, 115–121.
- Vitting-Seerup, K. and Sandelin, A. (2017) The landscape of isoform switches in human cancers. *Mol. Cancer Res.*, **15**, 1206–1220.
- Vitting-Seerup, K. and Sandelin, A. (2019) IsoformSwitchAnalyzeR: analysis of changes in genome-wide patterns of alternative splicing and its functional consequences. *Bioinformatics*, **35**, 4469–4471.
- UniProt-Consortium, T., Bateman, A., Martin, M.-J., Orchard, S., Magrane, M., Agivetova, R., Ahmad, S., Alpi, E., Bowler-Barnett, E.H., Britto, R., et al. (2020) UniProt: the universal protein knowledgebase in 2021. *Nucleic Acids Res.*, **49**, D480–D489.
- GTEX-Consortium, T. (2020) The GTEx Consortium atlas of genetic regulatory effects across human tissues. *Science*, **369**, 1318–1330.
- Frankish, A., Diekhans, M., Jungreis, J., Lagarde, J., Loveland, J.E., Mudge, J.M., Sisu, C., Wright, J.C., Armstrong, J., Barnes, J., et al. (2020) GENCODE 2021. *Nucleic Acids Res.*, **49**, D916–D923.
- Hagemann-Jensen, M., Ziegenhain, C., Chen, P., Ramsköld, D., Hendriks, G.-J., Larsson, A.J.M., Faridani, O.R. and Sandberg, R. (2020) Single-cell RNA counting at allele and isoform resolution using Smart-seq3. *Nat. Biotechnol.*, **38**, 708–714.
- Berman, H.M., Westbrook, J., Feng, Z., Gilliland, G., Bhat, T.N., Weissig, H., Shindyalov, I.N. and Bourne, P.E. (2000) The Protein Data Bank. *Nucleic Acids Res.*, **28**, 235–242.
- Grant, B.J., Skjærven, L. and Yao, X. (2020) The Bio3D packages for structural bioinformatics. *Protein Sci.*, **30**, 20–30.
- Zhang, Y. and Skolnick, J. (2004) Scoring function for automated assessment of protein structure template quality. *Proteins Struct. Funct. Bioinform.*, **57**, 702–710.
- Lee, Y.-H., Li, Y., Uyeda, K. and Hasemann, C.A. (2003) Tissue-specific structure/function differentiation of the liver isoform of 6-phosphofructo-2-kinase/fructose-2,6-bisphosphatase*. *J. Biol. Chem.*, **278**, 523–530.
- Crochet, R.B., Kim, J., Lee, H., Yim, Y., Kim, S., Neau, D. and Lee, Y. (2017) Crystal structure of heart 6-phosphofructo-2-kinase/fructose-2,6-bisphosphatase (PFKFB2) and the inhibitory influence of citrate on substrate binding. *Proteins Struct. Funct. Bioinform.*, **85**, 117–124.
- Boutard, N., Bialas, A., Sabiniarz, A., Guzik, P., Banaszak, K., Biela, A., Bielnik, M., Buda, A., Bugaj, B., Cieluch, E., et al. (2019) Discovery and structure–activity relationships of N-aryl 6-aminoquinoxalines as potent PFKFB3 kinase inhibitors. *ChemMedChem*, **14**, 169–181.
- Wang, Y., Wei, Z., Liu, L., Cheng, Z., Lin, Y., Ji, F. and Gong, W. (2005) Crystal structure of human B-type phosphoglycerate mutase bound with citrate. *Biochem. Biophys. Res. Commun.*, **331**, 1207–1215.
- Wang, Y., Liu, L., Wei, Z., Cheng, Z., Lin, Y. and Gong, W. (2006) Seeing the process of histidine phosphorylation in human bisphosphoglycerate mutase. *J. Biol. Chem.*, **281**, 39642–39648.
- Rigden, D.J. (2007) The histidine phosphatase superfamily: structure and function. *Biochem. J.*, **409**, 333–348.
- Loo, M. and van der, P.J. (2014) The stringdist package for approximate string matching. *R J.*, **6**, 111.

29. Liberzon,A., Subramanian,A., Pinchback,R., Thorvaldsdóttir,H., Tamayo,P. and Mesirov,J.P. (2011) Molecular signatures database (MSigDB) 3.0. *Bioinformatics*, **27**, 1739–1740.
30. Liberzon,A., Birger,C., Thorvaldsdóttir,H., Ghandi,M., Mesirov,J.P. and Tamayo,P. (2015) The molecular signatures database hallmark gene set collection. *Cell Syst.*, **1**, 417–425.
31. Korotkevich,G., Sukhov,V., Budin,N., Shpak,B., Artyomov,M.N. and Sergushichev,A. (2021) Fast gene set enrichment analysis. bioRxiv doi: <https://doi.org/10.1101/060012>, 01 February 2021, preprint: not peer reviewed .
32. Savojardo,C., Babbi,G., Martelli,P.L. and Casadio,R. (2021) Mapping OMIM disease-related variations on protein domains reveals an association among variation type, Pfam models, and disease classes. *Front. Mol. Biosci.*, **8**, 617016.
33. Peterson,T.A., Gauran,I.I.M., Park,J., Park,D. and Kann,M.G. (2017) Oncodomains: a protein domain-centric framework for analyzing rare variants in tumor samples. *PLoS Comput. Biol.*, **13**, e1005428.
34. Thompson,B.J. and Sahai,E. (2015) MST kinases in development and disease. *J. Cell Biol.*, **210**, 871–882.
35. Patel,S.H., Camargo,F.D. and Yimlamai,D. (2017) Hippo signaling in the liver regulates organ size, cell fate, and Carcinogenesis. *Gastroenterology*, **152**, 533–545.
36. Cinar,B., Alp,E., Al-Mathkour,M., Boston,A., Dwead,A., Khazaw,K. and Gregory,A. (2021) The Hippo pathway: an emerging role in urologic cancers. *Am. J. Clin. Exp. Urol.*, **9**, 301–317.
37. Han,Y. (2019) Analysis of the role of the Hippo pathway in cancer. *J. Transl. Med.*, **17**, 116.
38. Rybarczyk,A., Klacz,J., Wronska,A., Matuszewski,M., Kmiec,Z. and Wierzbicki,P.M. (2017) Overexpression of the YAP1 oncogene in clear cell renal cell carcinoma is associated with poor outcome. *Oncol. Rep.*, **38**, 427–439.
39. Triant,D.A. and Pearson,W.R. (2015) Most partial domains in proteins are alignment and annotation artifacts. *Genome Biol.*, **16**, 99.
40. Bros,M., Haas,K., Moll,L. and Grabbe,S. (2019) RhoA as a key regulator of innate and adaptive immunity. *Cells*, **8**, 733.
41. Schaefer,A., Reinhard,N.R. and Hordijk,P.L. (2014) Toward understanding RhoGTPase specificity: structure, function and local activation. *Small GTPases*, **5**, e968004.
42. Hauksdottir,H., Farhoud,B. and Privalsky,M.L. (2003) Retinoic acid receptors β and γ do not repress, but instead activate target gene transcription in both the absence and presence of hormone ligand. *Mol. Endocrinol.*, **17**, 373–385.
43. Aranda,A. and Pascual,A. (2001) Nuclear hormone receptors and gene expression. *Physiol. Rev.*, **81**, 1269–1304.
44. Abbot,E.S. (1916) The causal relations between structure and function in biology. *Am J Psychology*, **27**, 245.
45. Ma,R., Kandera,E., Sundh,U.B., Geng,M., Ek,P., Zetterqvist,Ö. and Li,J.-P. (2005) Mutational study of human phosphohistidine phosphatase: effect on enzymatic activity. *Biochem. Biophys. Res. Commun.*, **337**, 887–891.
46. Scott,D.E., Bayly,A.R., Abell,C. and Skidmore,J. (2016) Small molecules, big targets: drug discovery faces the protein–protein interaction challenge. *Nat. Rev. Drug Discov.*, **15**, 533–550.
47. Glinos,D.A., Garborcauskas,G., Hoffman,P., Ehsan,N., Jiang,L., Gokden,A., Dai,X., Aguet,F., Brown,K.L., Garimella,K., et al. (2022) Transcriptome variation in human tissues revealed by long-read sequencing. *Nature*, **608**, 353–359.
48. Gupta,I., Collier,P.G., Haase,B., Mahfouz,A., Joglekar,A., Floyd,T., Koopmans,F., Barres,B., Smit,A.B., Sloan,S.A., et al. (2018) Single-cell isoform RNA sequencing characterizes isoforms in thousands of cerebellar cells. *Nat. Biotechnol.*, **36**, 1197–1202.
49. Robinson,S.L., Piel,J. and Sunagawa,S. (2021) A roadmap for metagenomic enzyme discovery. *Nat. Prod. Rep.*, **38**, 1994–2023.
50. Doğan,T., Güzelcan,E.A., Baumann,M., Koyas,A., Atas,H., Baxendale,I.R., Martin,M. and Cetin-Atalay,R. (2021) Protein domain-based prediction of drug/compound–target interactions and experimental validation on LIM kinases. *PLoS Comput. Biol.*, **17**, e1009171.
51. Gaudet,P., Livstone,M.S., Lewis,S.E. and Thomas,P.D. (2011) Phylogenetic-based propagation of functional annotations within the Gene Ontology consortium. *Brief Bioinform.*, **12**, 449–462.
52. Ji,Y., Mishra,R.K. and Davuluri,R.V. (2020) In silico analysis of alternative splicing on drug-target gene interactions. *Sci. Rep.*, **10**, 134.
53. Barnkob,M.B., Vitting-Seerup,K. and Olsen,L.R. (2022) Target isoforms are an overlooked challenge and opportunity in chimeric antigen receptor cell therapy. *Immunother. Adv.*, **2**, ltac009.