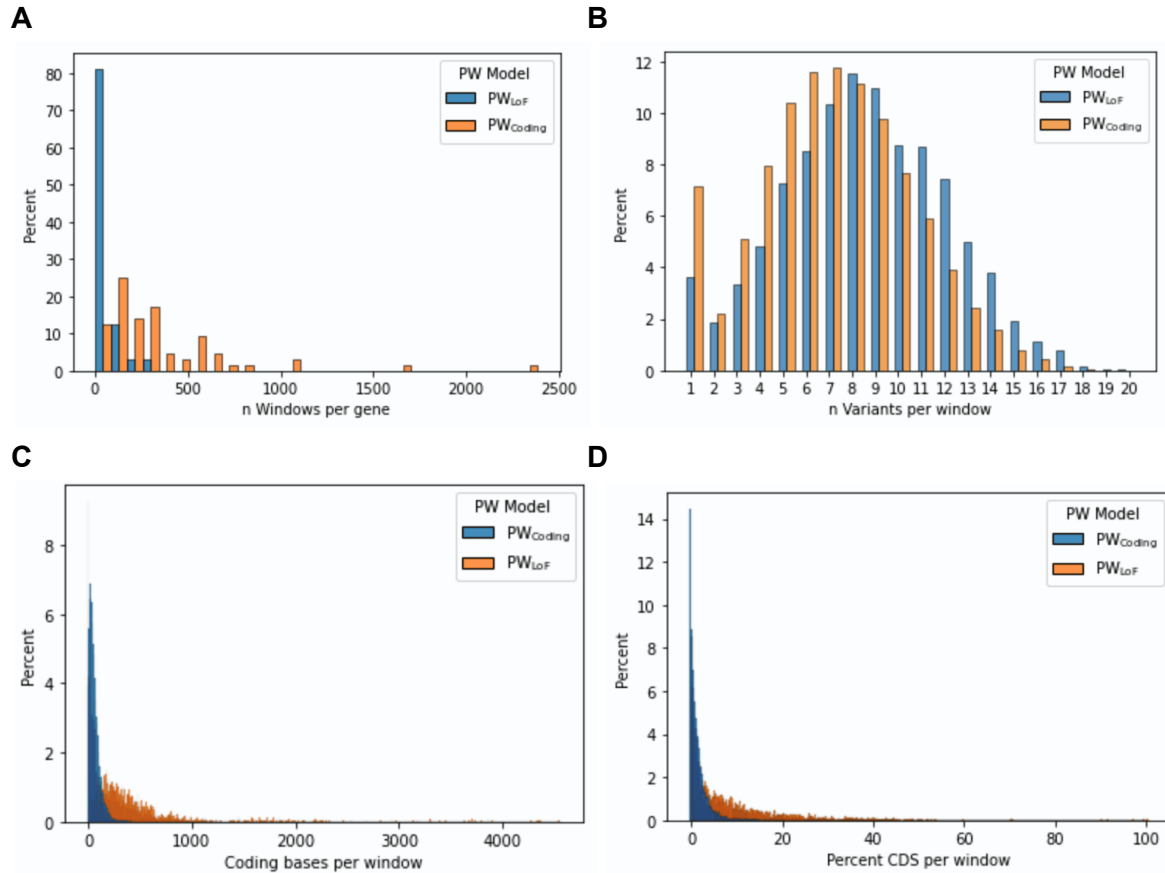


**HGGA, Volume 5**

**Supplemental information**

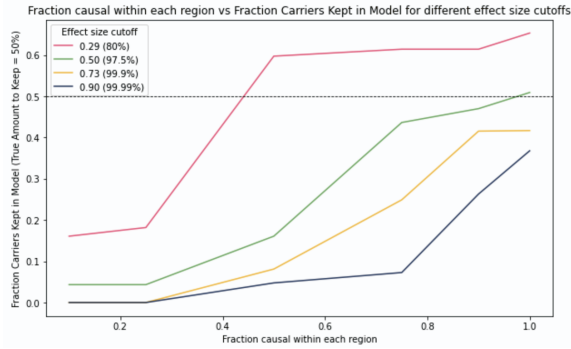
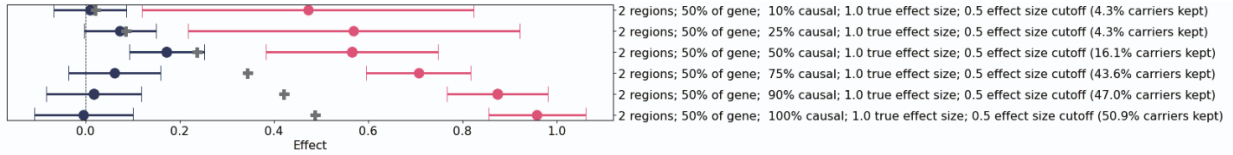
**A power-based sliding window approach to evaluate  
the clinical impact of rare genetic variants  
in the nucleotide sequence or the spatial position of the folded protein**

**Elizabeth T. Cirulli, Kelly M. Schiabor Barrett, Alexandre Bolze, Daniel P. Judge, Pamala A. Pawloski, Joseph J. Grzymalski, William Lee, and Nicole L. Washington**

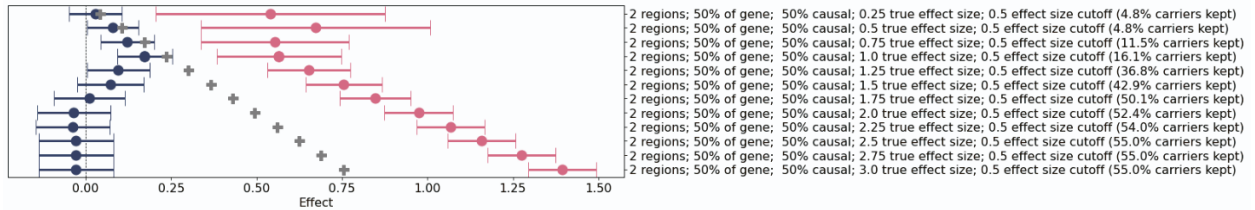


**Figure S1. Stats for power windows.** Because of the size of *TTN*, the stats for this gene are not displayed and instead are written separately in the legend. A) Histogram of number of windows analyzed per gene. The mean number of windows analyzed per gene was 357 / 67 for  $PW_{Coding}$  /  $PW_{LoF}$  (range 30-2385 and 2-302). *TTN* = 10,034 / 1362 for  $PW_{Coding}$  /  $PW_{LoF}$ . B) Histogram of number of variants included in each window. The mean number of variants analyzed per window was 7 / 8 for  $PW_{Coding}$  /  $PW_{LoF}$  (range 1-20 and 1-20). *TTN* = 7 (1-18) / 11 (1-18) for  $PW_{Coding}$  /  $PW_{LoF}$ . C) Histogram of number of coding bases included in each window; when a window only included one variant or only included variants at the same site, then the length is 1. The mean number of coding bases included per window was 61 / 404 for  $PW_{Coding}$  /  $PW_{LoF}$  (range 1-1396 and 1-4549). *TTN* = 74 (1-2847) / 943 (1-2632) for  $PW_{Coding}$  /  $PW_{LoF}$ . D) Histogram of the percent of each CDS for the gene that was included in each window. The mean percent of the CDS of the gene that was analyzed in each window was 2% / 10% for  $PW_{Coding}$  /  $PW_{LoF}$  (range <1-36% and <1-100%). *TTN* = 0.07% (0.001%-2.6%) / 0.87% (0.001%-2.4%) for  $PW_{Coding}$  /  $PW_{LoF}$ .

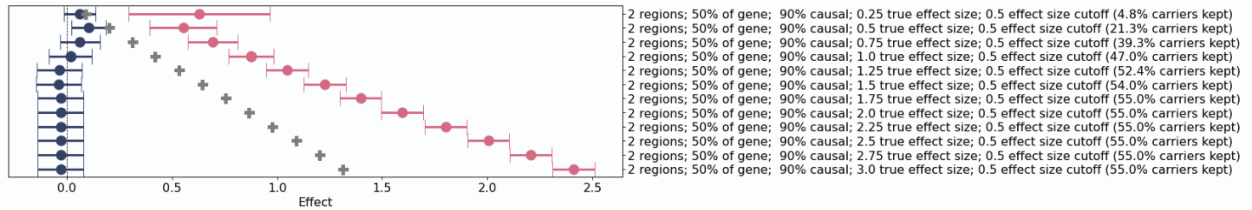
### A Different % causal variants within 2 regions



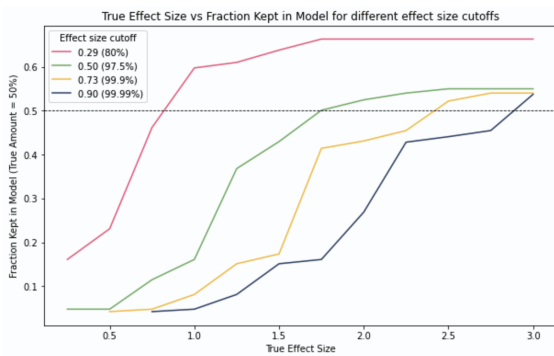
### B Different effect sizes: 50% causal variants within 2 regions



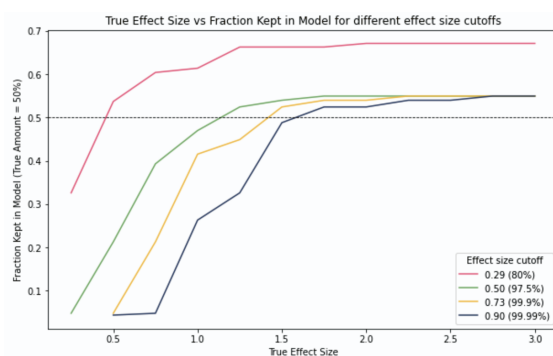
### Different effect sizes: 90% causal variants within 2 regions

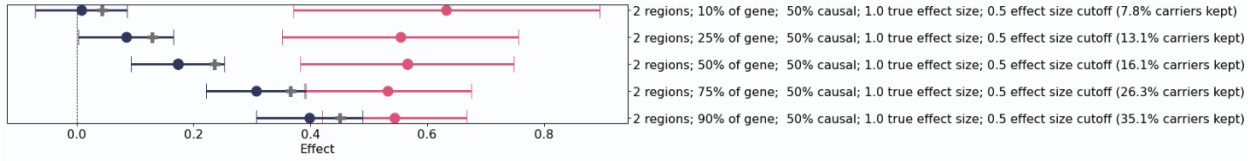
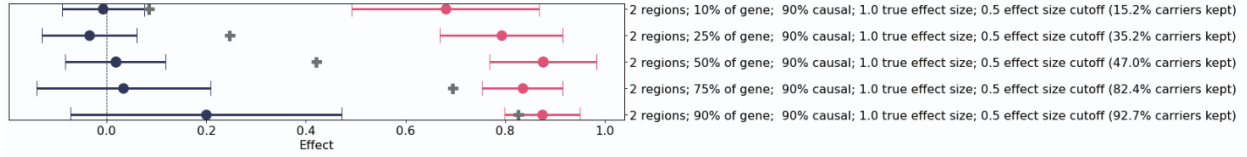
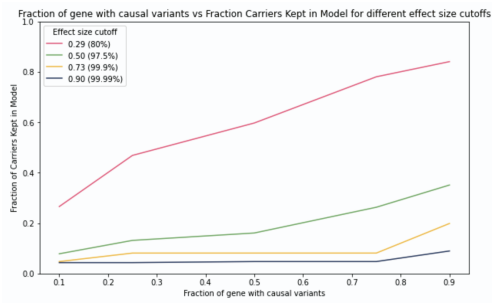
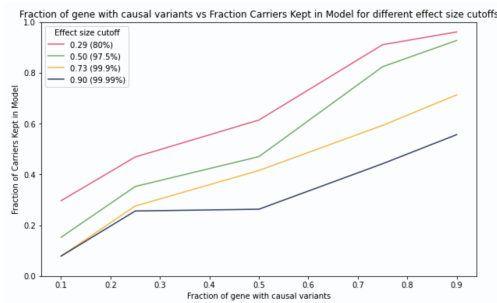
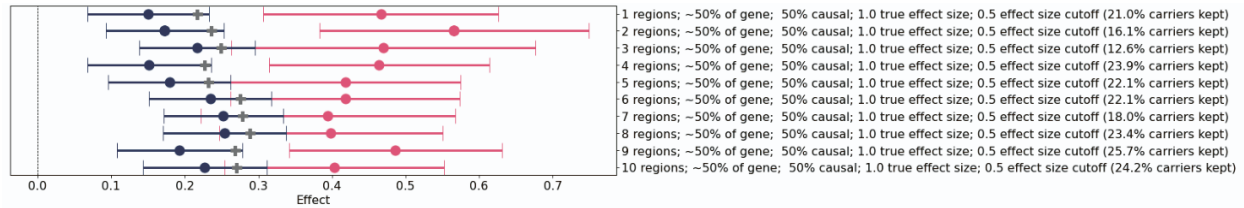
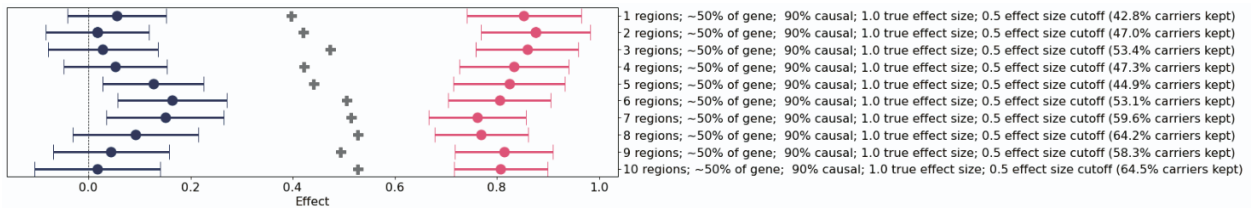
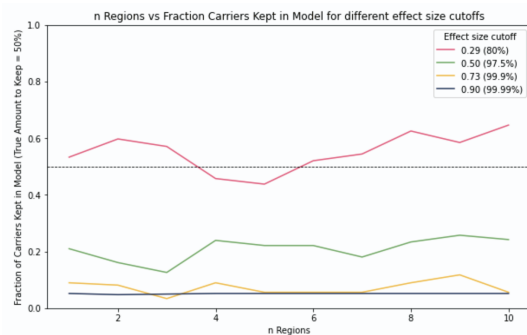
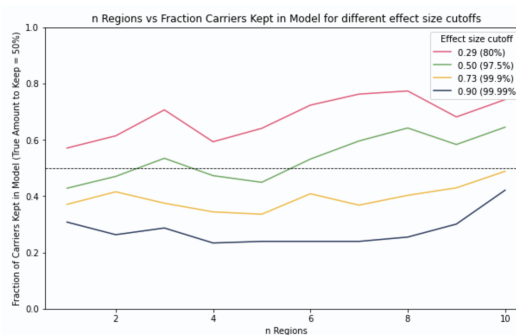


### 50% causal variants within 2 regions

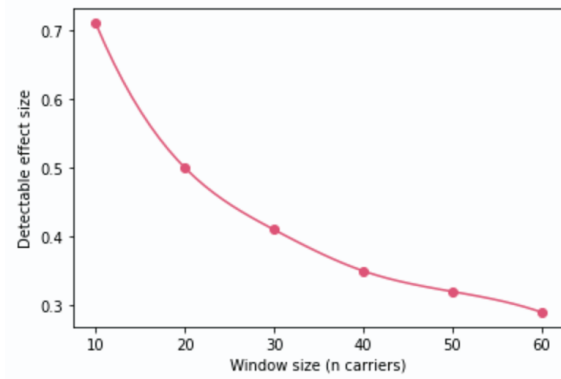
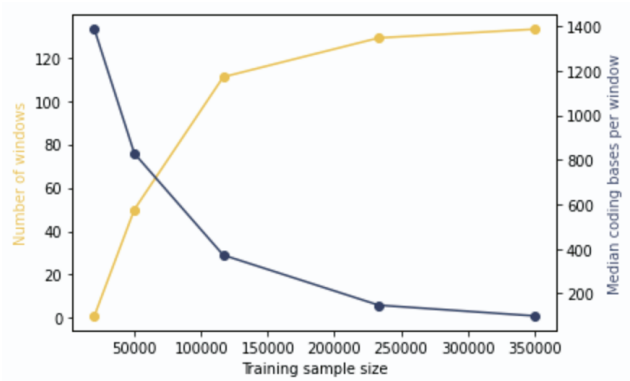


### 90% causal variants within 2 regions



**C****Different percents of gene with causal variants: 2 regions in which 50% are causal variant****Different percents of gene with causal variants: 2 regions in which 90% are causal variants****50% causal variants within 2 regions****90% causal variants within 2 regions****D****Different number of causal regions per gene: 50% causal variants in each region (effect size 1)****Different number of causal regions per gene: 90% causal variants in each region (effect size 1)****50% causal variants within each region****90% causal variants within each region**

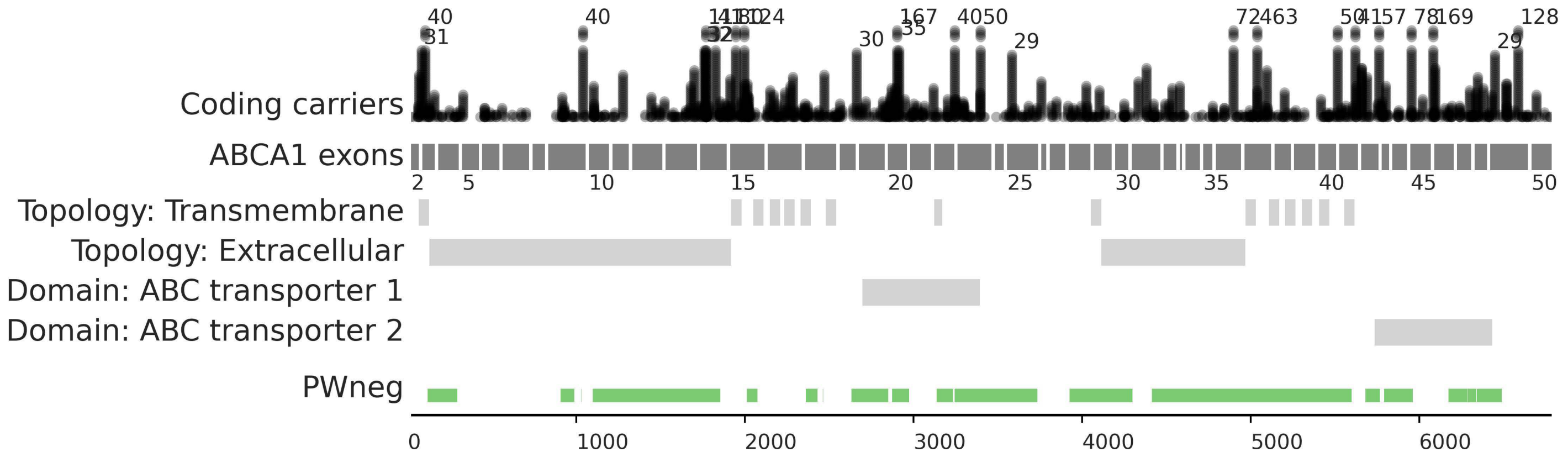
**Figure S2. Simulations showing the effect of tuning different parameters for Power Window in the example gene GCK and phenotype glucose with a coding model.** For the simulations, causal regions were defined within GCK and glucose levels for each carrier were reassigned according to the parameters specified. The figures show the model performance in the independent replication set of 117k individuals, as in **Figure 3**. Effect size values for gene regions are defined by inclusion in a PW model (PW, pink dot) or exclusion from a PW model (non-PW, blue dot) are plotted against the score for the whole gene (grey cross), together with 95% confidence intervals. Each row has different parameters, as labeled, with the effect size cutoff for inclusion into the model set at 0.5 just like in the main manuscript. All models shown are significant at  $p < 0.0005$  in the test set unless otherwise specified. Each panel also includes a plot showing the relationship between the specified parameter and the percent of carriers kept in the model, with different effect size cutoffs for inclusion into the model (rejecting an effect size of 0 with 80% confidence, 97.5% (as in the main manuscript), 99.9%, and 99.99%). For all panels, “percent causal variants” is equivalent to “of the individuals with variants in the region, the percent whose variant is causal”, i.e., putting aside that one variant may be seen in 1 individual and another in 5 individuals. **A)** Parameter changed: percent causal variants within associated regions of the gene. The simulation was set with 2 implicated regions taking up 50% of the gene in which causal variants had an effect size of 1, and the percent causal variants within those regions, ranged from 0% to 100%. To capture all of the associated regions with a cutoff of 0.5 when the effect size is 1, at least 75% of the variants in implicated regions must be causal. The only non-significant model was the one with 10% of the variants in the region being causal ( $p = 0.01$ ). **B)** Parameter changed: effect size of causal variants. The simulation was set with 2 implicated regions taking up 50% of the gene, in which either 50% or 90% of the variants were causal, and the effect size of the causal variants ranged from 0.25 to 3. To capture all of the associated regions with a cutoff of 0.5 when only 50% of the variants in implicated regions are truly associated, the true effect size of causal variants must be at least  $\sim 1.5$ ; if 90% of the variants are truly associated, then the effect size of the causal variants must be at least  $\sim 1$ . The only non-significant models were the ones with an effect size of 0.25. **C)** Parameter changed: percent of gene implicated. The simulation was set with 2 implicated regions in which either 50% or 90% of the variants were causal with effect sizes of 1, with the percent of the gene implicated ranging from 10% to 90%. When 50% of the variants in the regions were causal, PW generally missed approximately half of the regions; when 90% of the variants in the regions were causal, PW identified the correct regions but also sometimes included some additional regions. **D)** Parameter changed: number of implicated regions. The simulation was set with the implicated regions taking up  $\sim 50\%$  of the gene (note the actual percentages varied from 46-54% due to variations in numbers of variants carriers in different regions), in which either 50% or 90% of the variants were causal with an effect size of 1, and the number of regions ranged from 1 to 10. When 50% of the variants in the regions were causal, PW performed similarly regardless of the number of regions but missed more carriers when there were fewer regions; when 90% of the variants in the regions were causal, PW correctly captured them when they were split into 1 or 2 regions but allowed in more carriers than should have been included as the number of implicated regions in the gene grew.

**A****B**

**Figure S3. Properties of Power Window size. A)** For a quantitative trait in a training set of 350k samples, the window size (number of carriers in the window) is compared to the effect size for which there is power to identify (with 97.5% confidence) that the beta is not 0. **B)** For the example gene *GCK* with a coding model (which in our main model had 133 windows, with a mean window length of 99 coding bases), the training sample size is compared to the physical length of each window in terms of coding position, for a set window size of 20 carriers per window. A larger sample size results in being able to home in on more specific regions for analysis.



**Figure S4. Percent CDS kept in Power Window models.** For each gene ( $n=65$ ), the percent of the CDS for the gene that was retained by the PW model was evaluated for (A)  $PW_{\text{Coding}}$  and (B)  $PW_{\text{LoF}}$ . Within each model, the phenotypes were grouped into quantitative (pink) or binary (yellow) traits, and the genes were grouped based on the original genome-wide association as follows: **LoF**: original whole-gene associations had an absolute beta at least 3x as high in the LoF model as the Coding model; **Coding and LoF**: original whole-gene associations had a beta  $<3x$  as high for LoF as for Coding (no gene-phenotype combinations had a whole-gene Coding model absolute beta that was at least 3x higher than the whole-gene LoF beta). Two genes had 0 windows included in the final model for  $PW_{\text{Coding}}$  (*ASXL1* and *NF1*), and 2 for  $PW_{\text{LoF}}$  (*GFI1B* and *JAK2*). Dotted horizontal line indicates the median percent carriers kept.

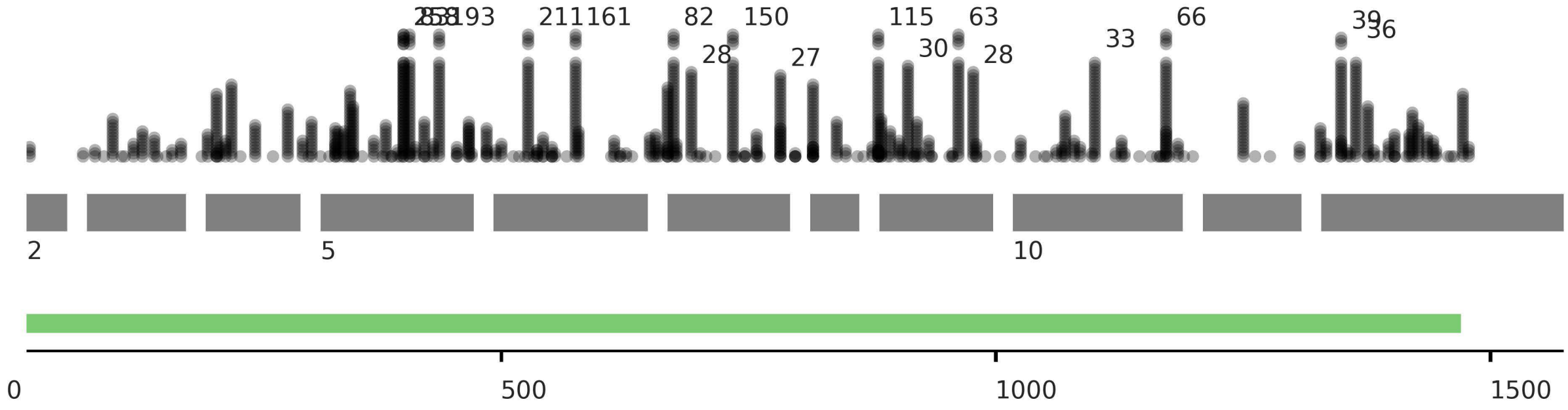


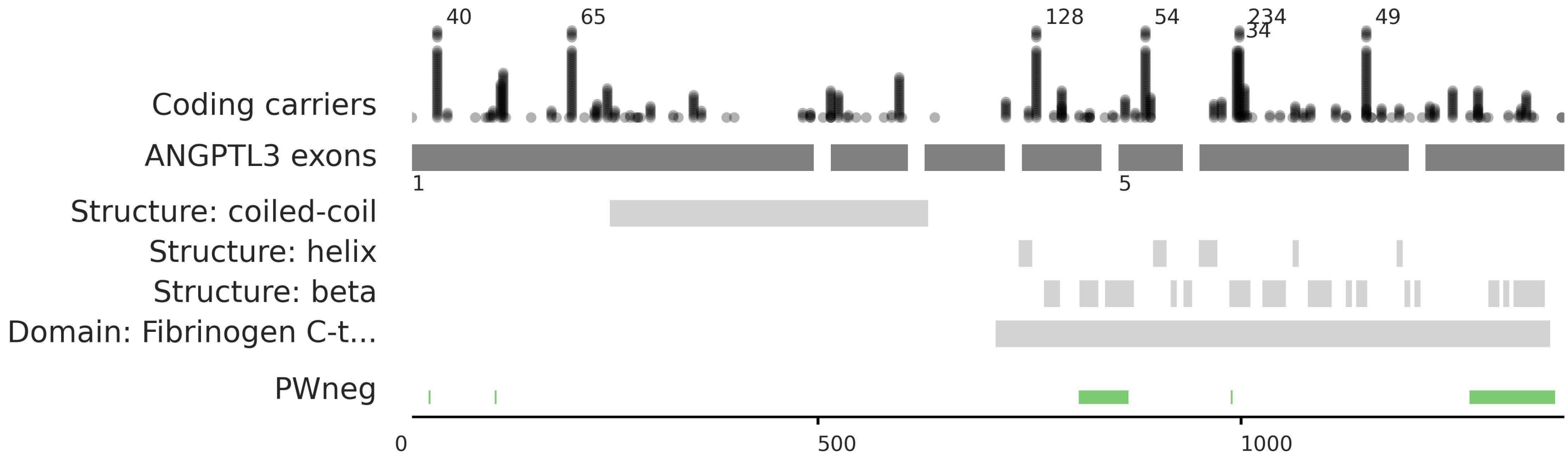


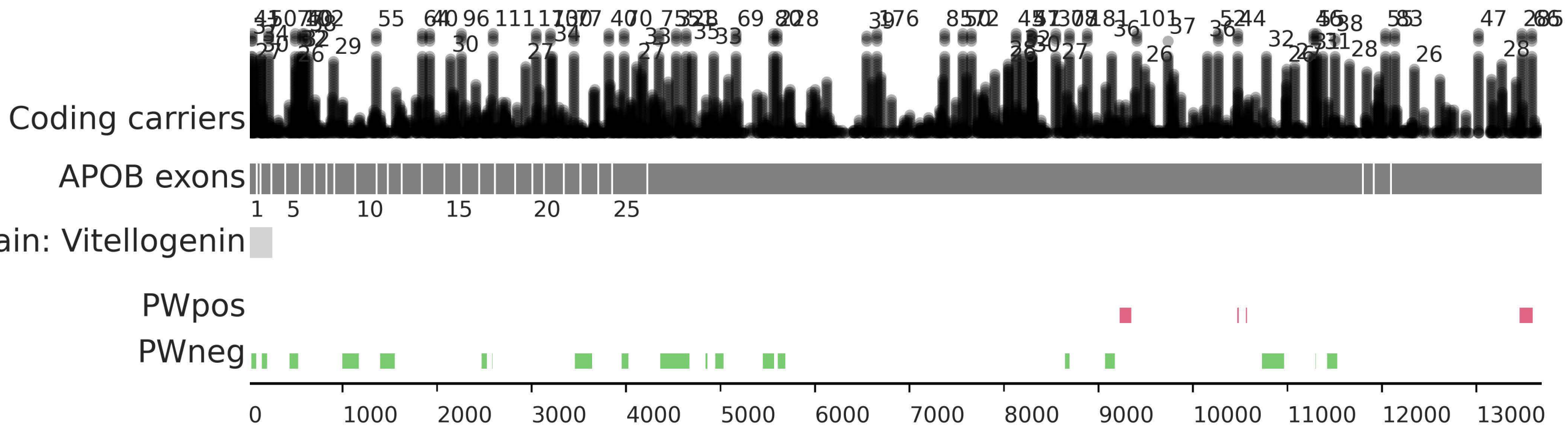
Coding carriers

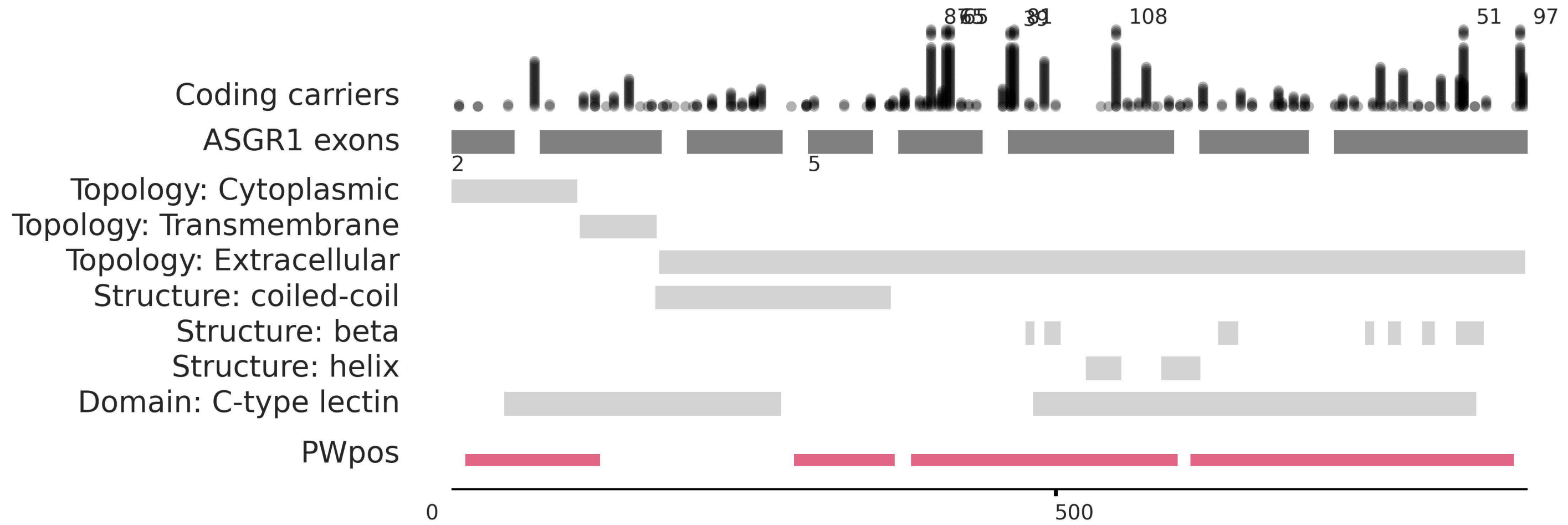
ALPL exons

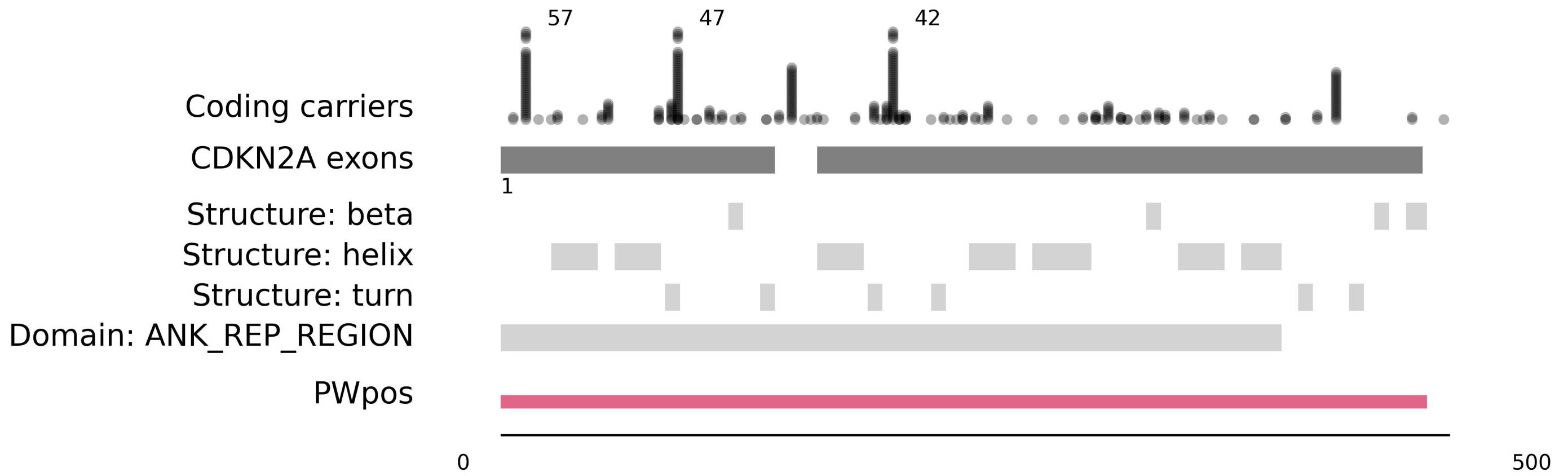
PWneg



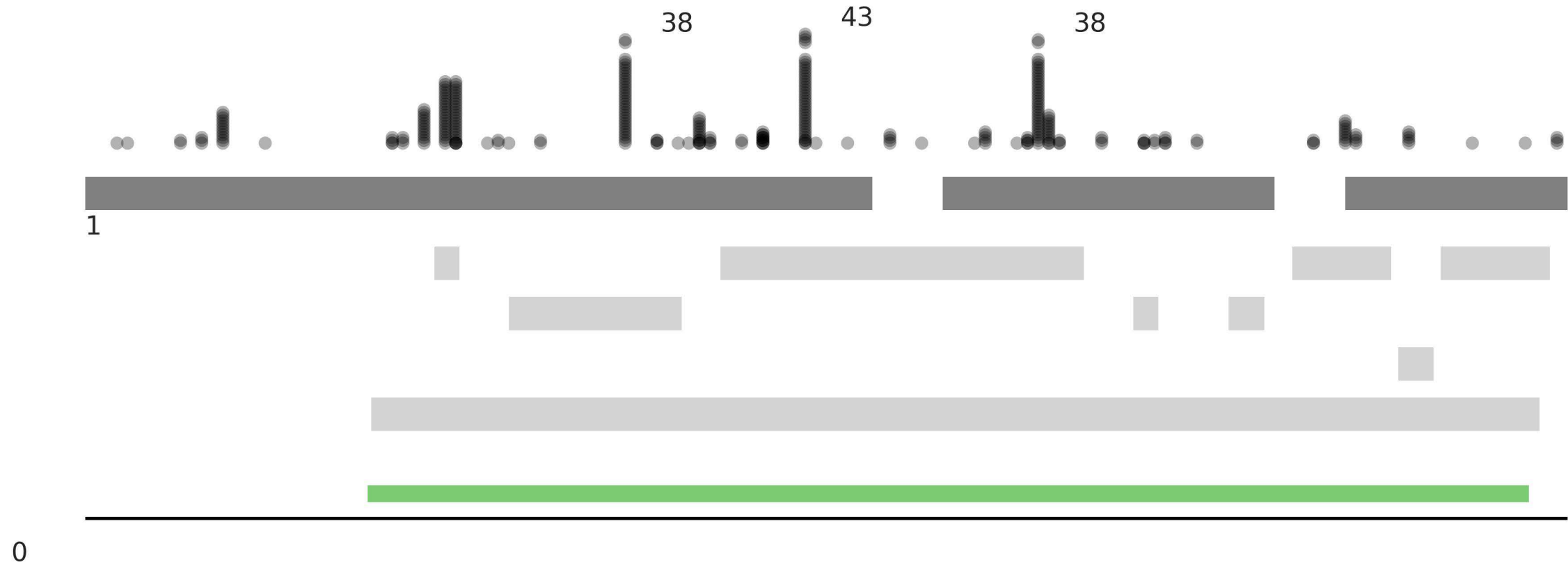


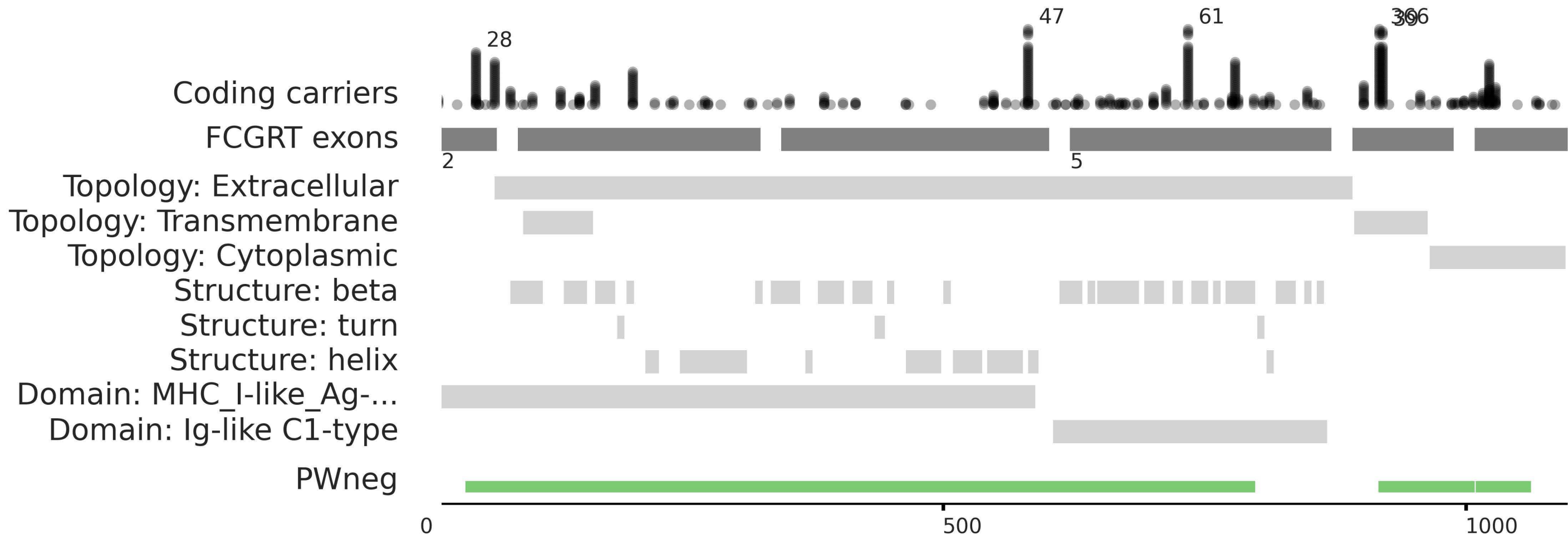


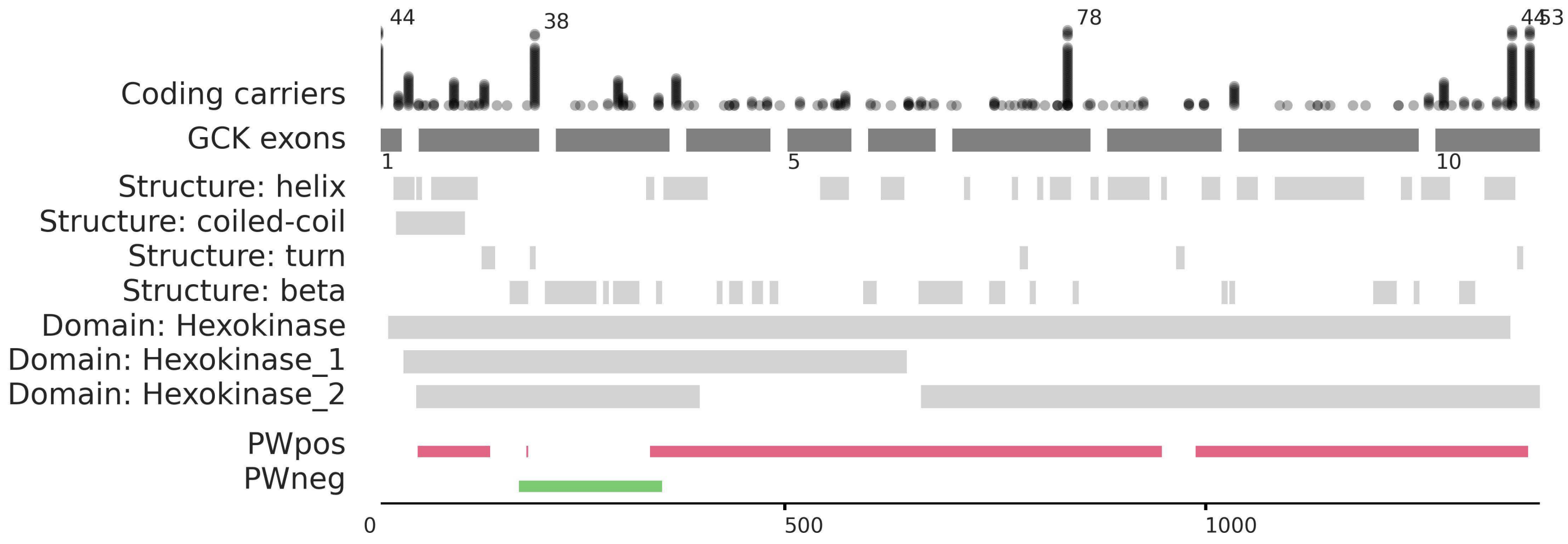




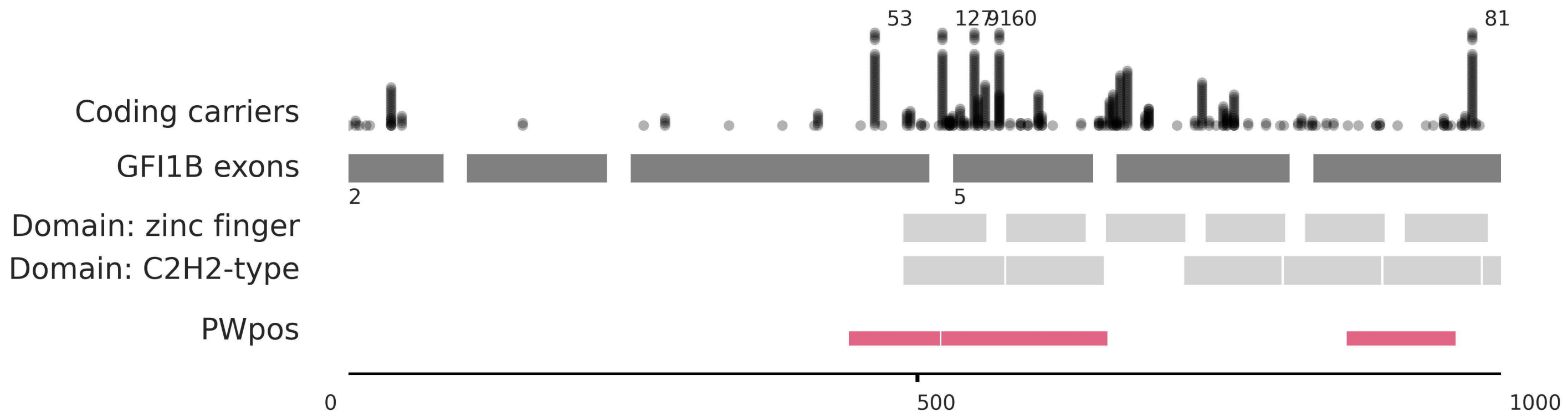
Coding carriers  
CST3 exons  
Structure: beta  
Structure: helix  
Structure: turn  
Domain: Cystatin  
PWneg











Coding carriers

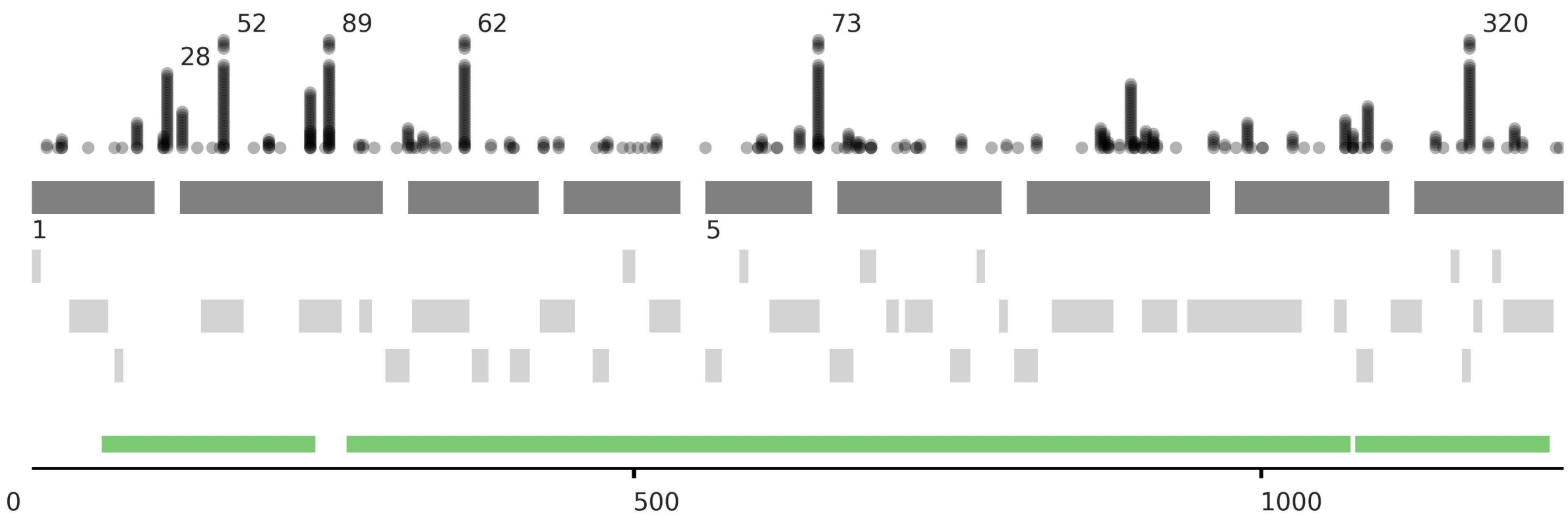
GOT1 exons

Structure: turn

Structure: helix

Structure: beta

PWneg



Coding carriers

GP1BB exons

Topology: Extracellular

Topology: Transmembrane

Topology: Cytoplasmic

Structure: beta

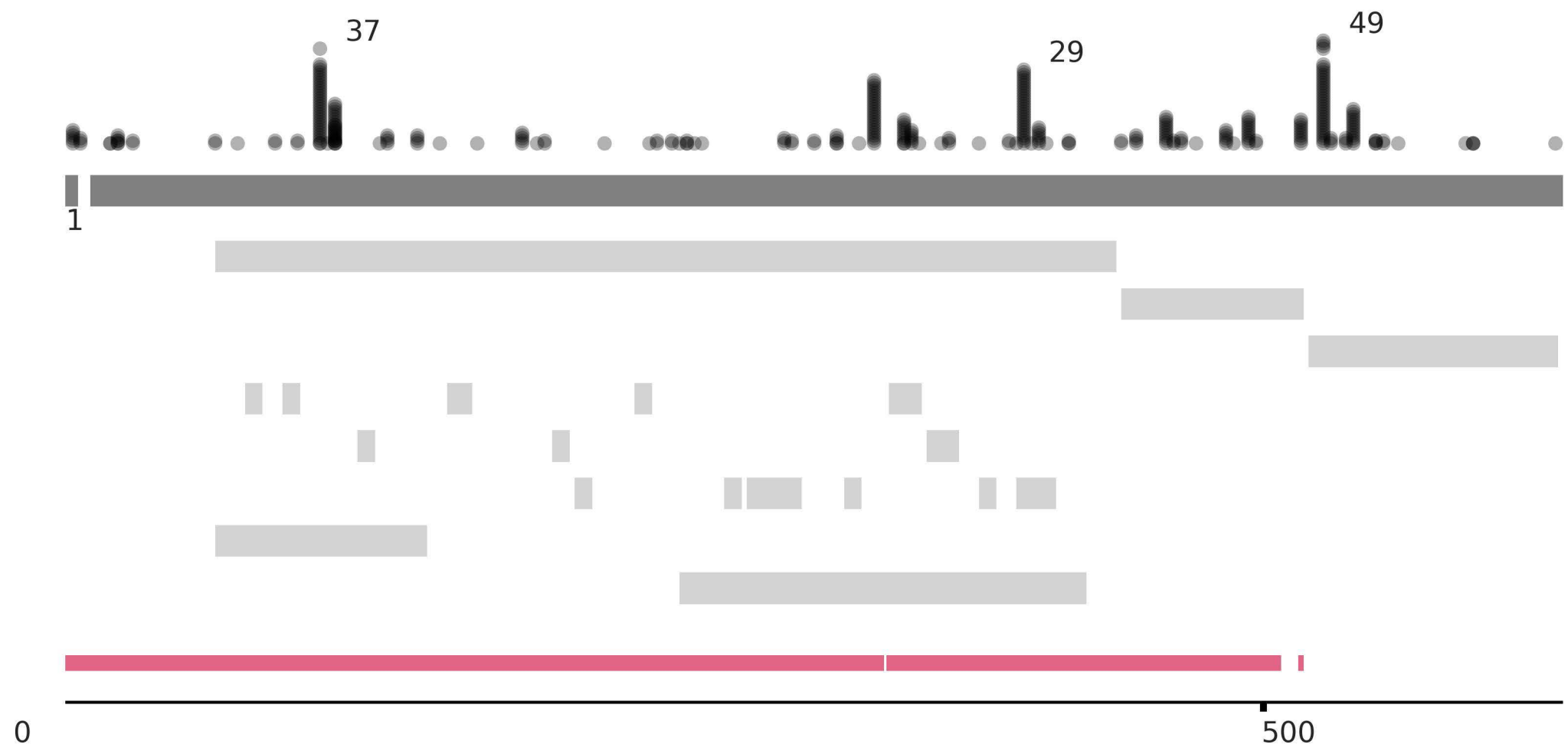
Structure: turn

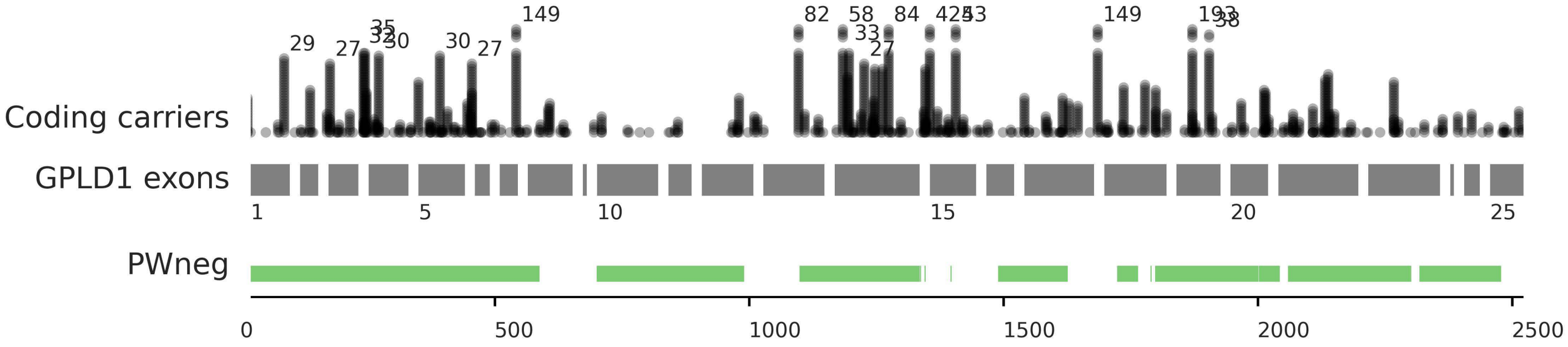
Structure: helix

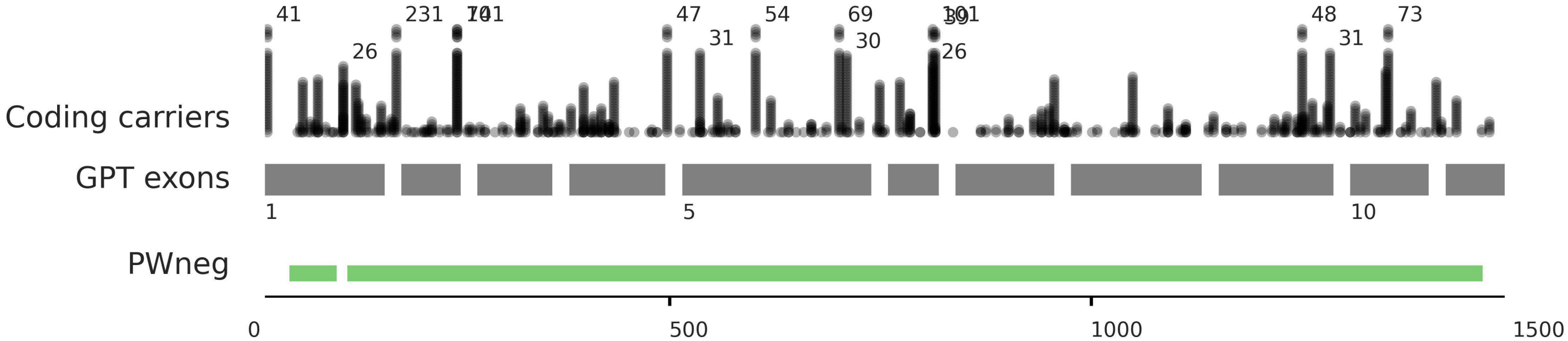
Domain: LRRNT

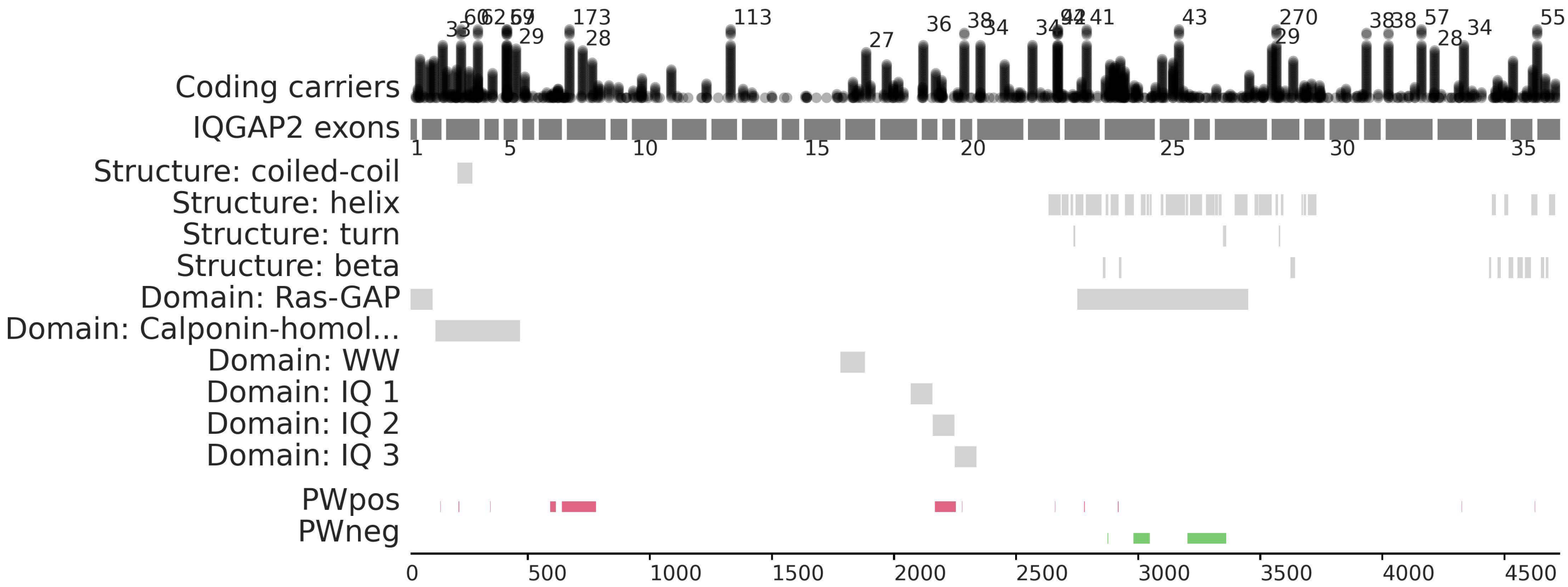
Domain: LRRCT

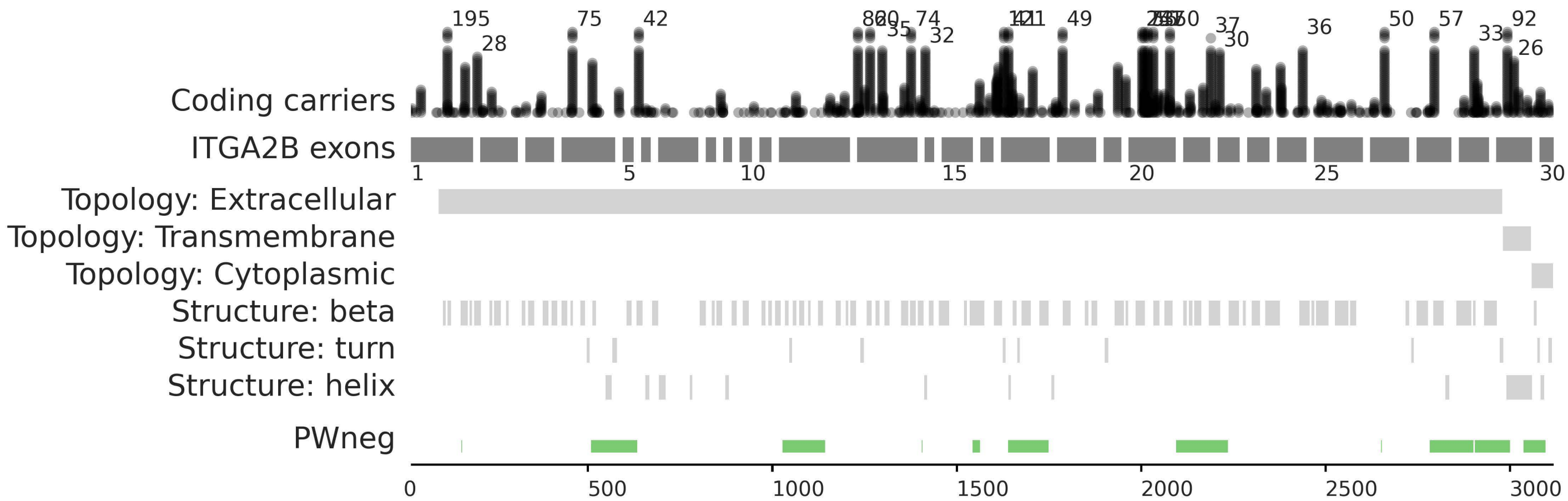
PWpos

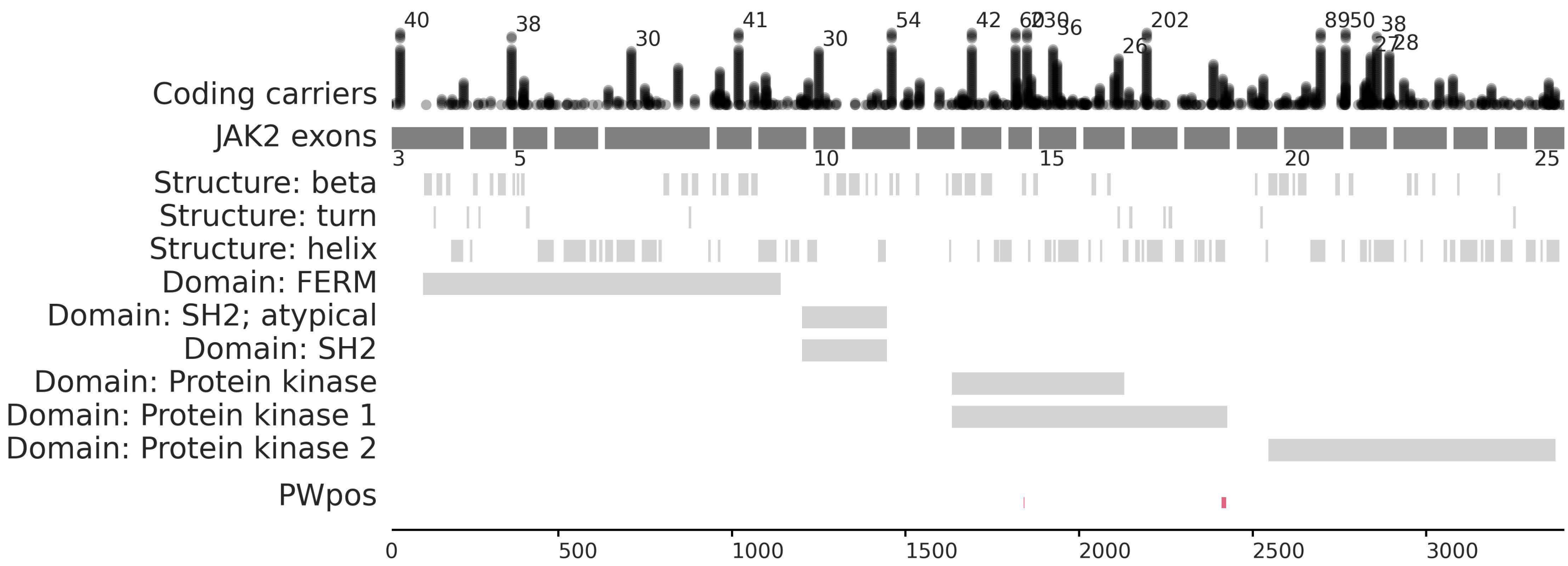




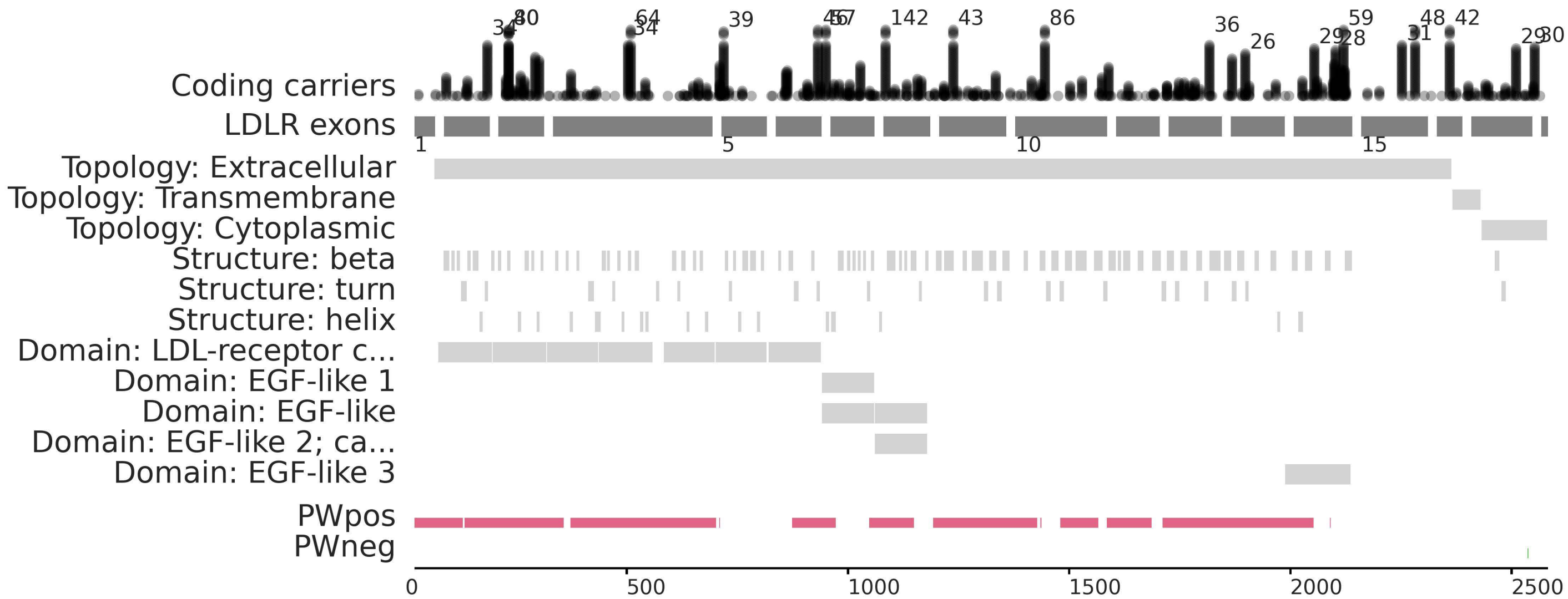












Coding carriers

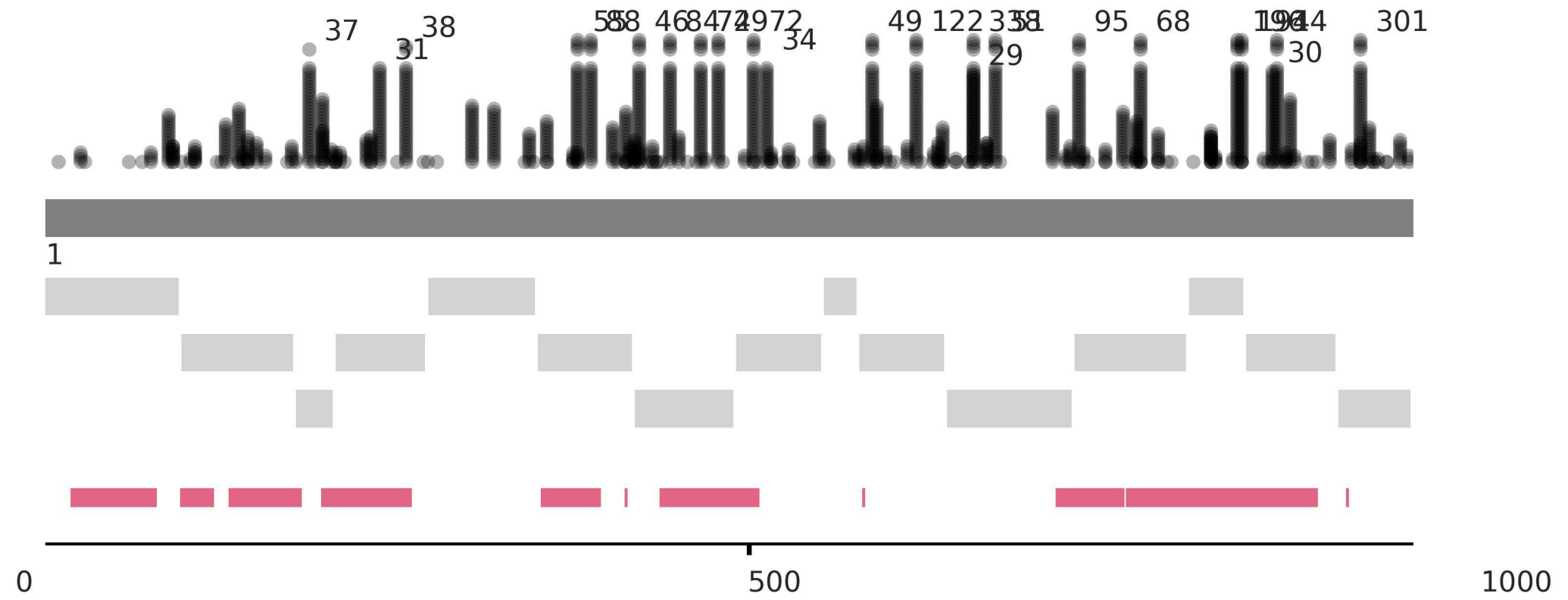
MC1R exons

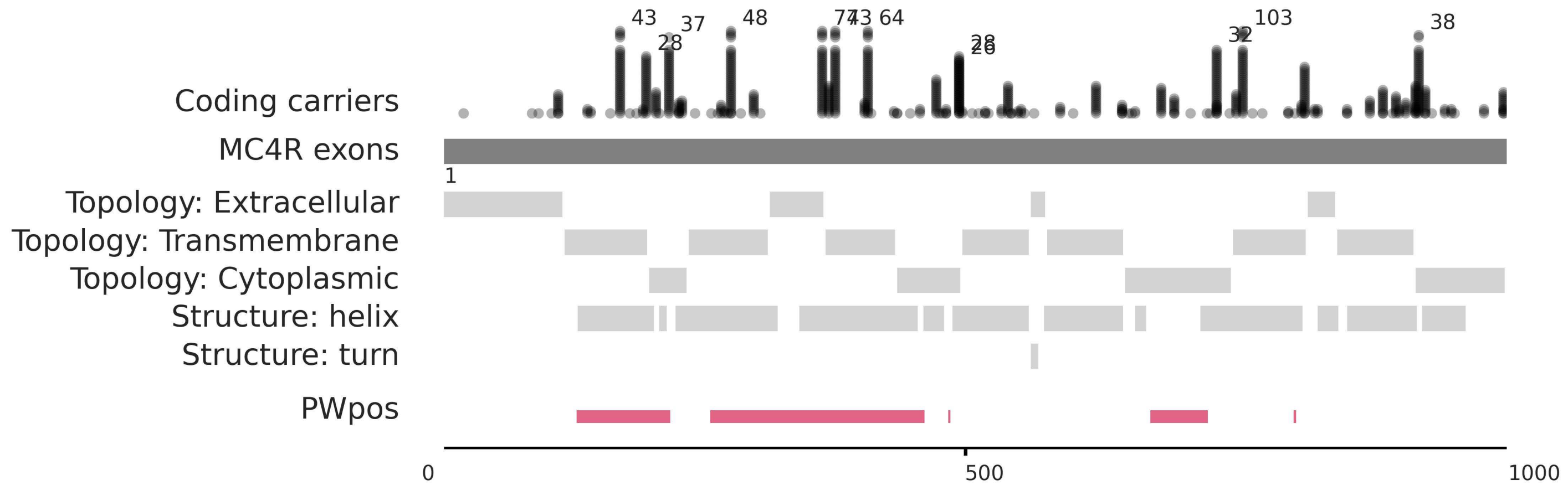
Topology: Extracellular

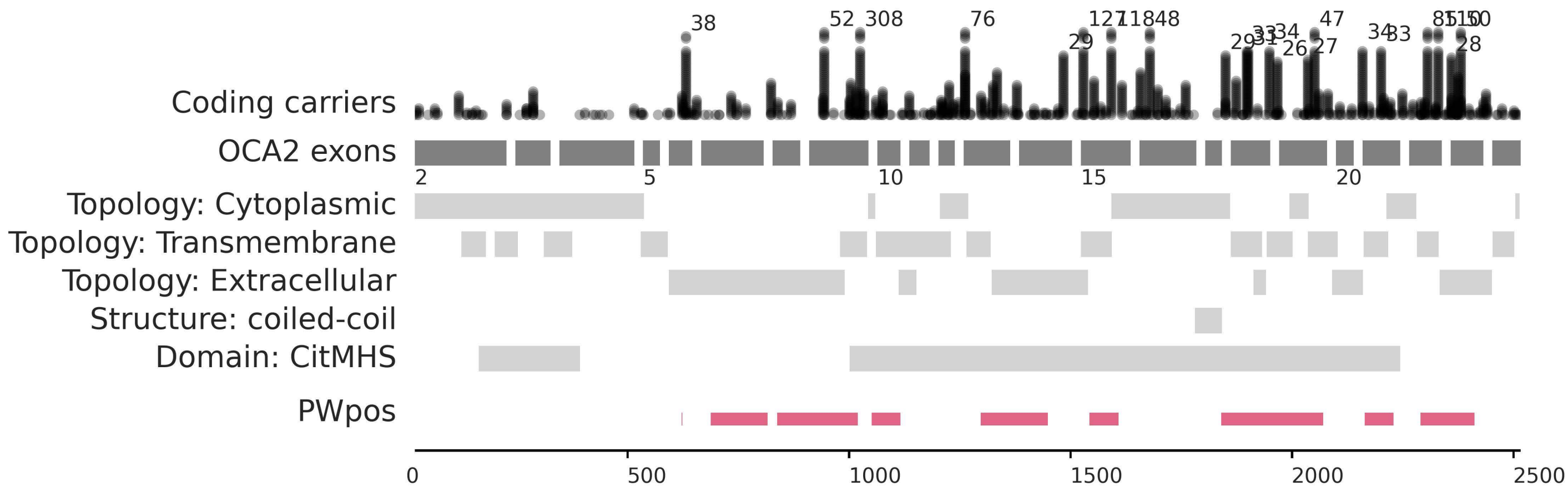
Topology: Transmembrane

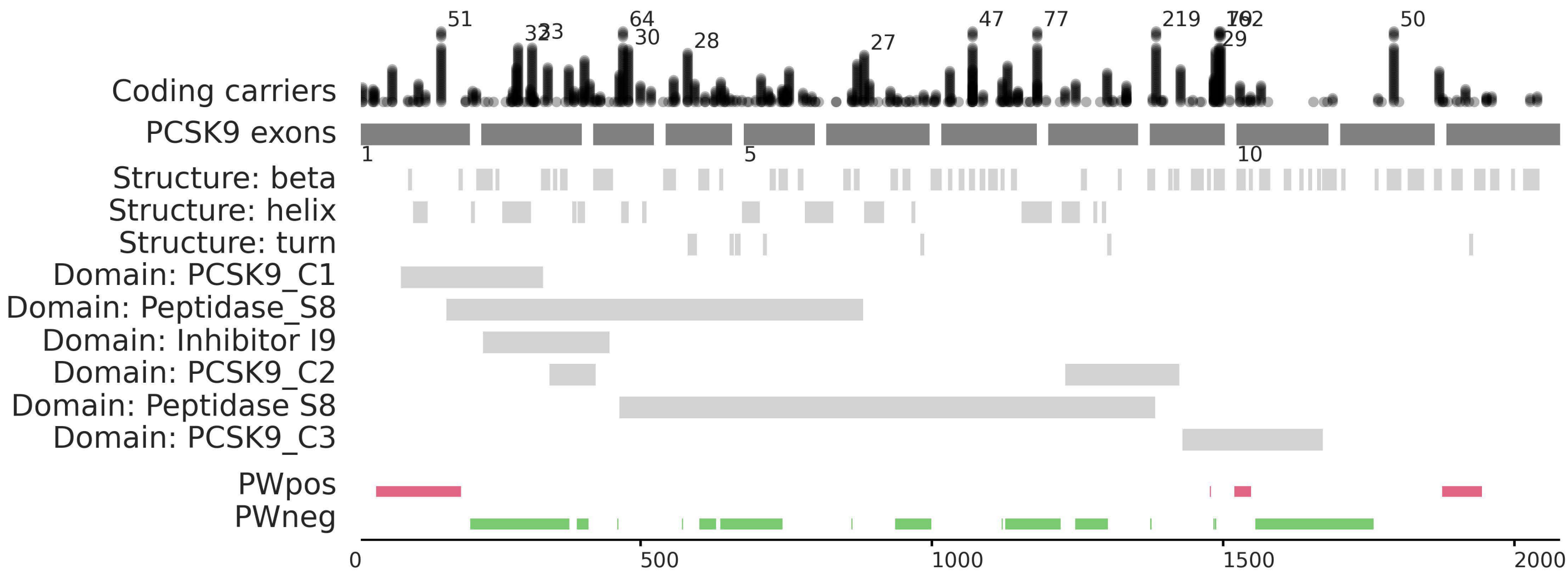
Topology: Cytoplasmic

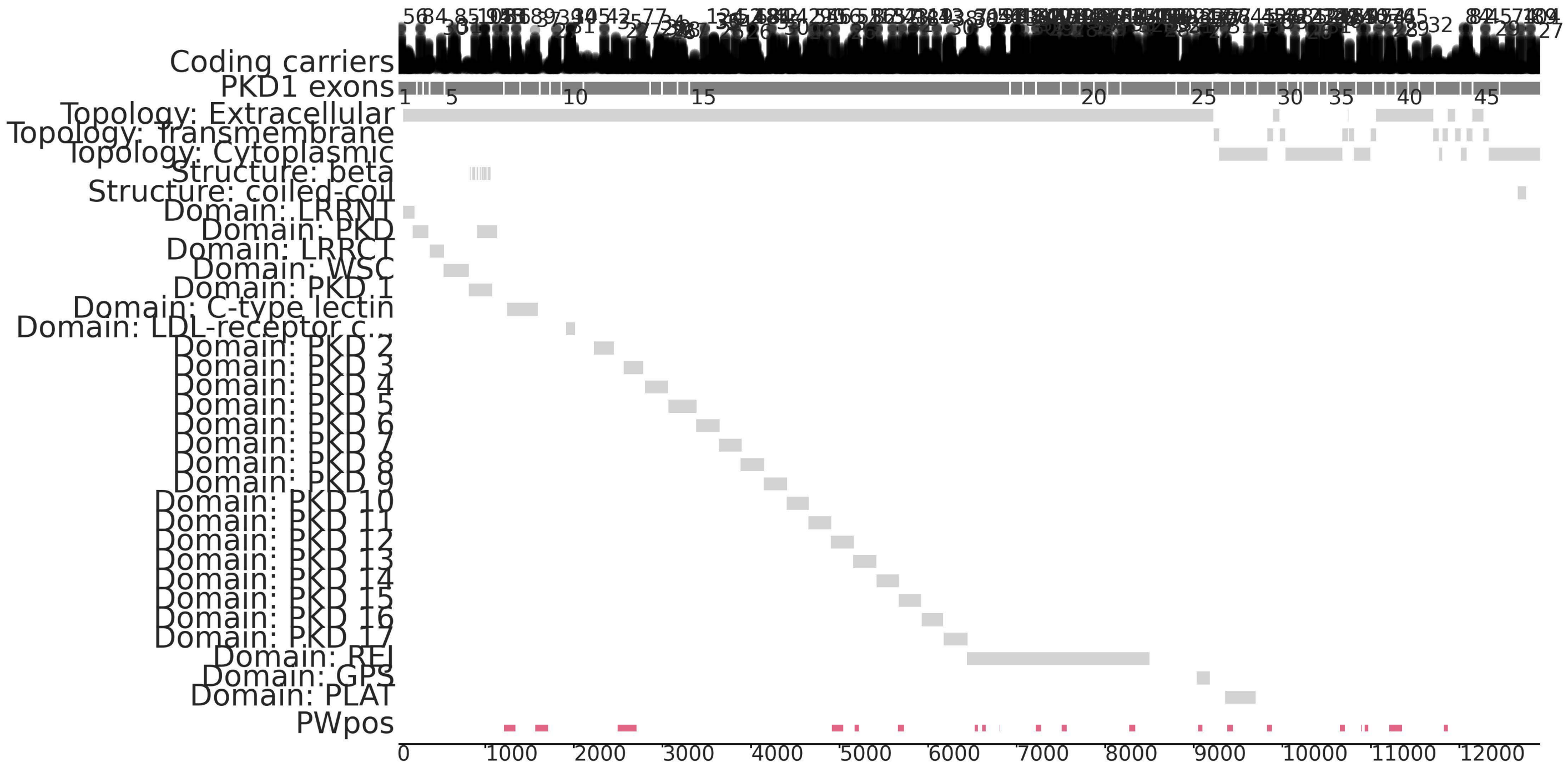
PWpos

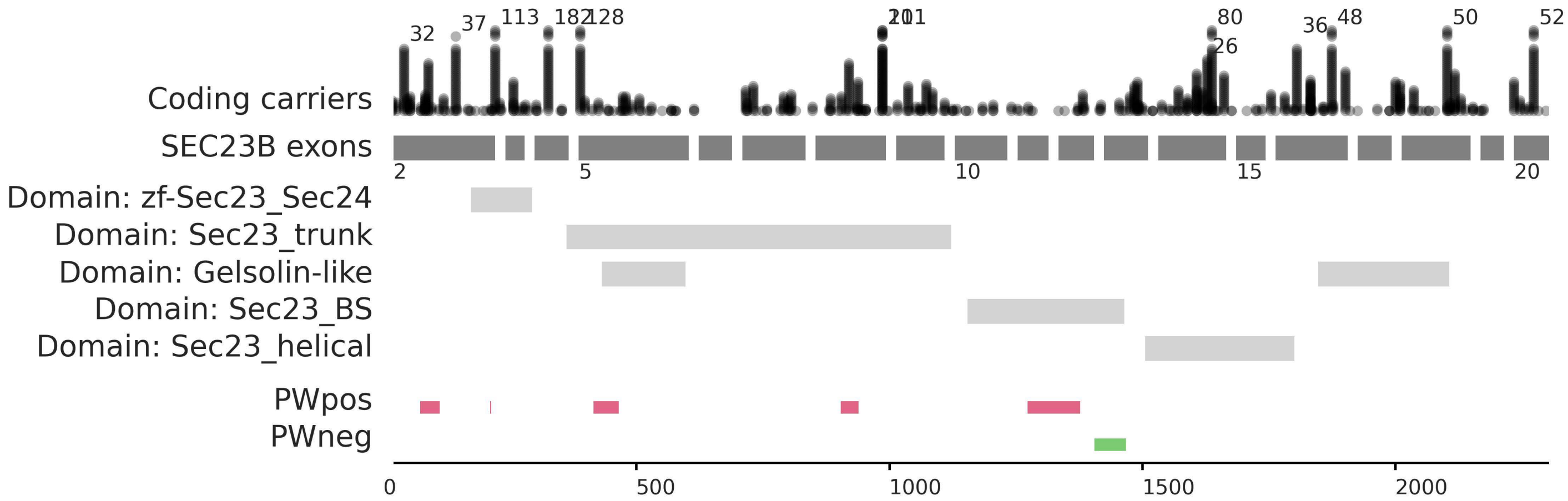


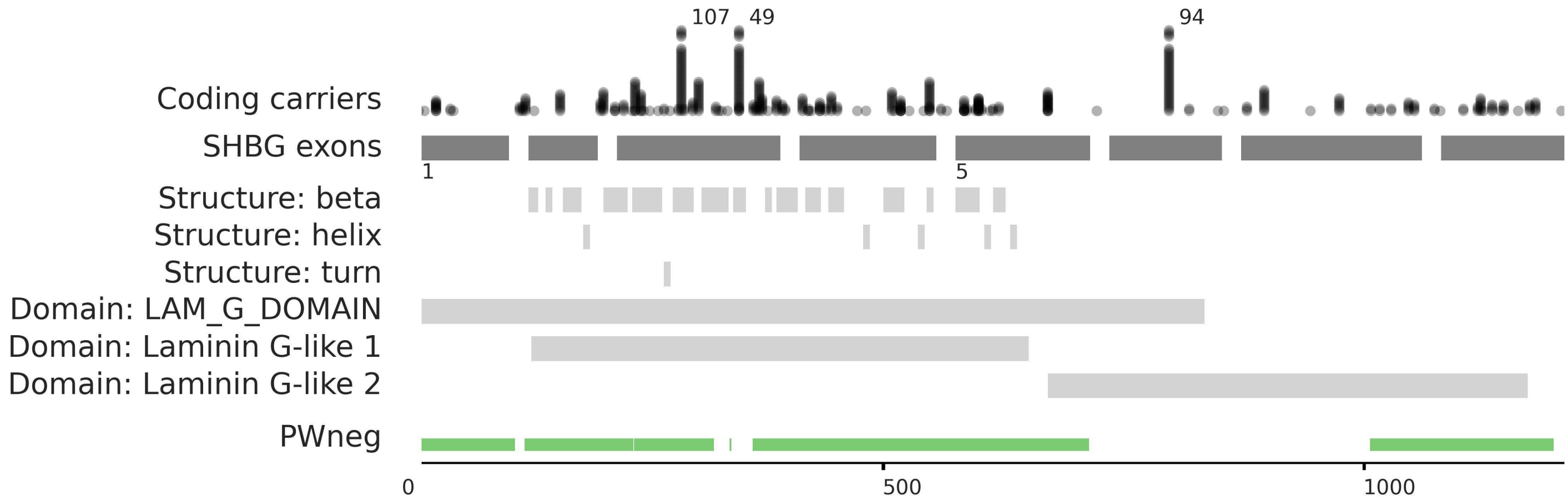




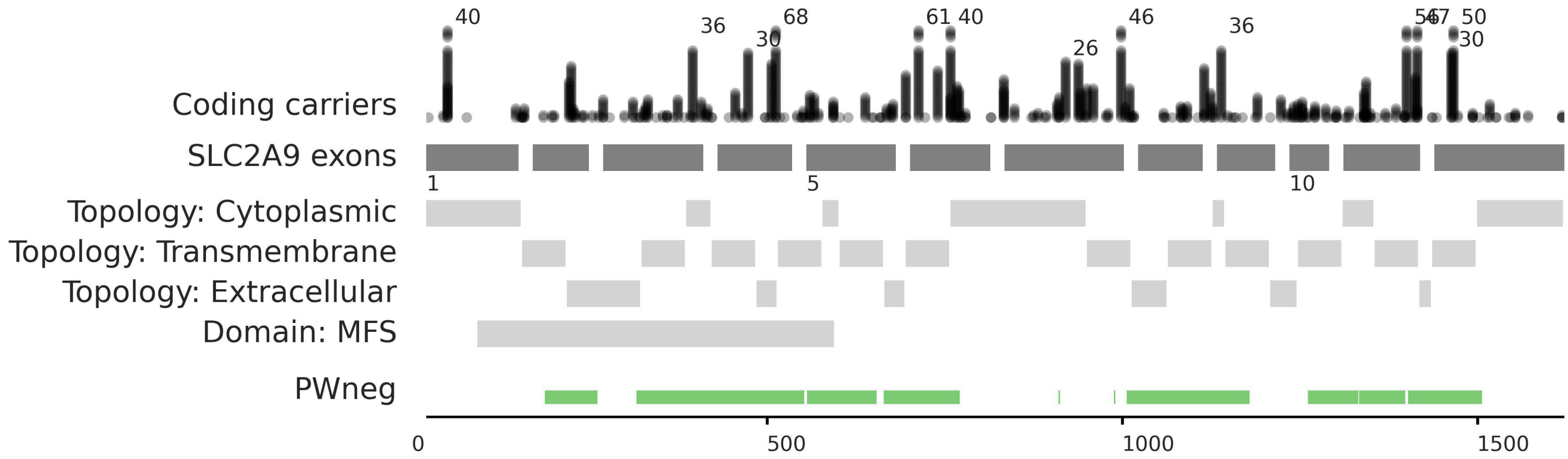


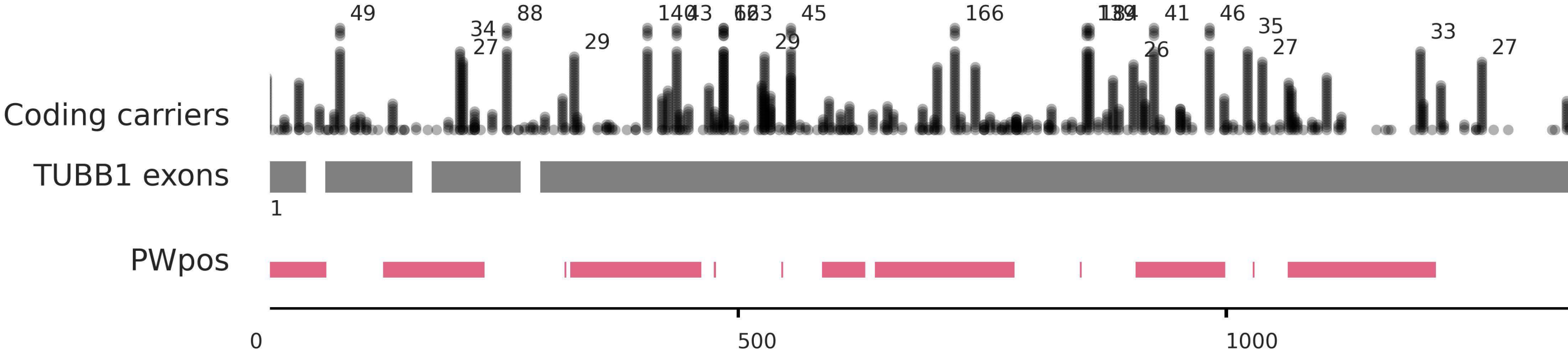


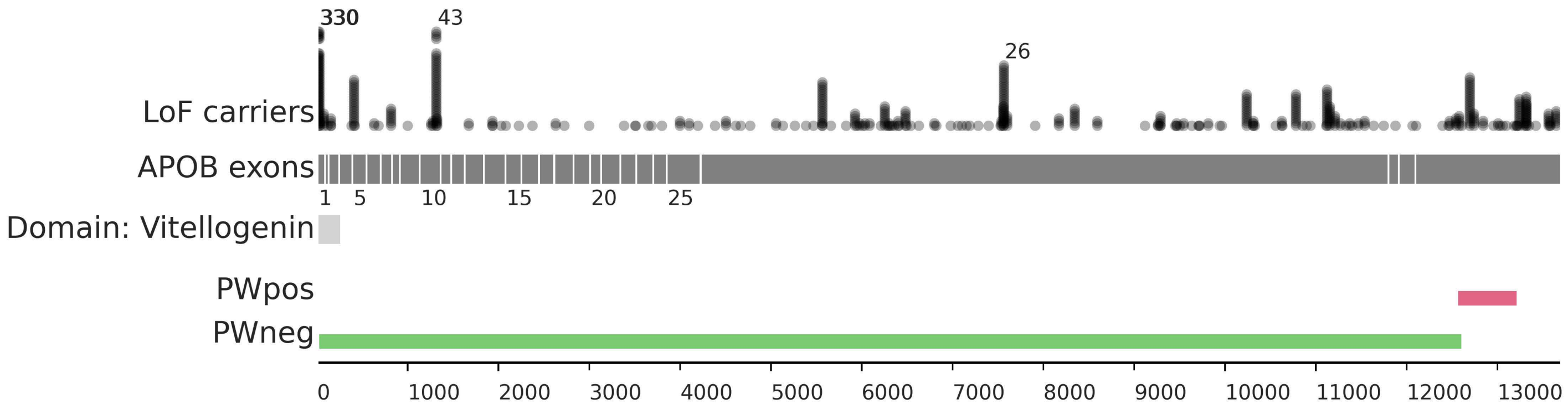


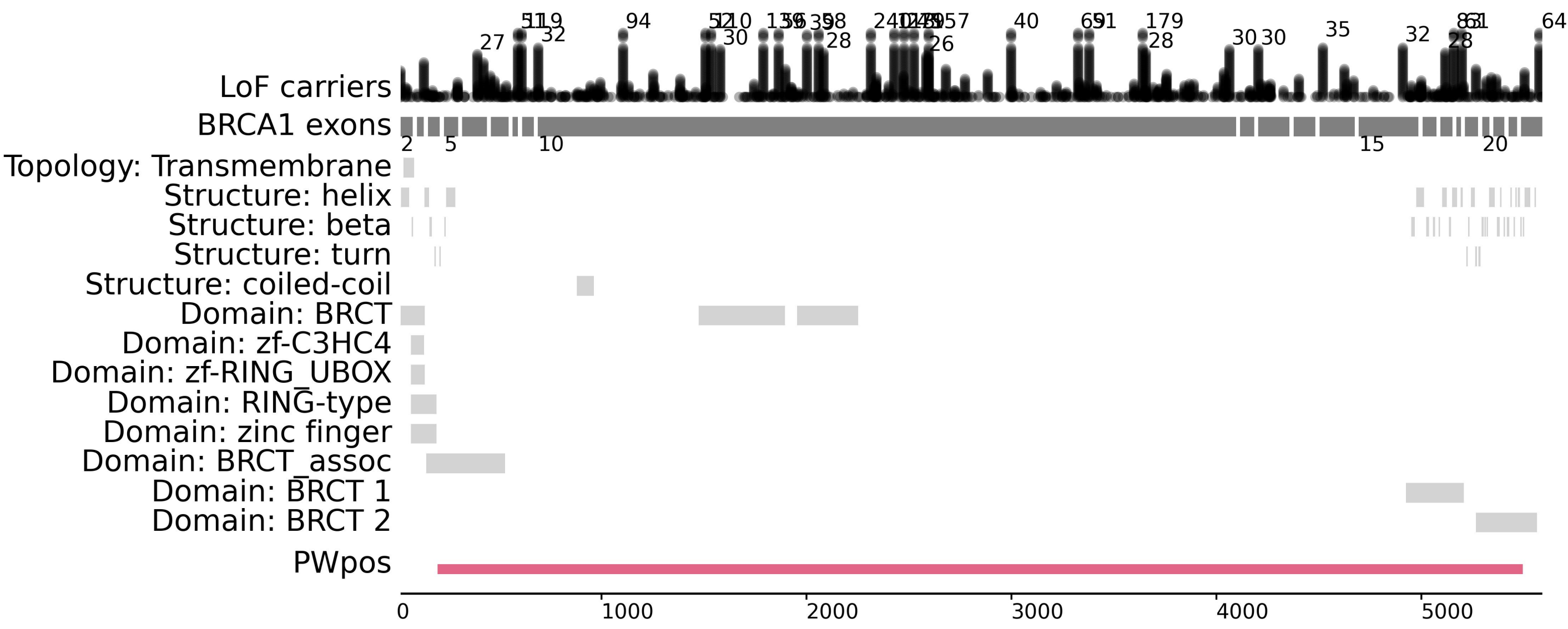


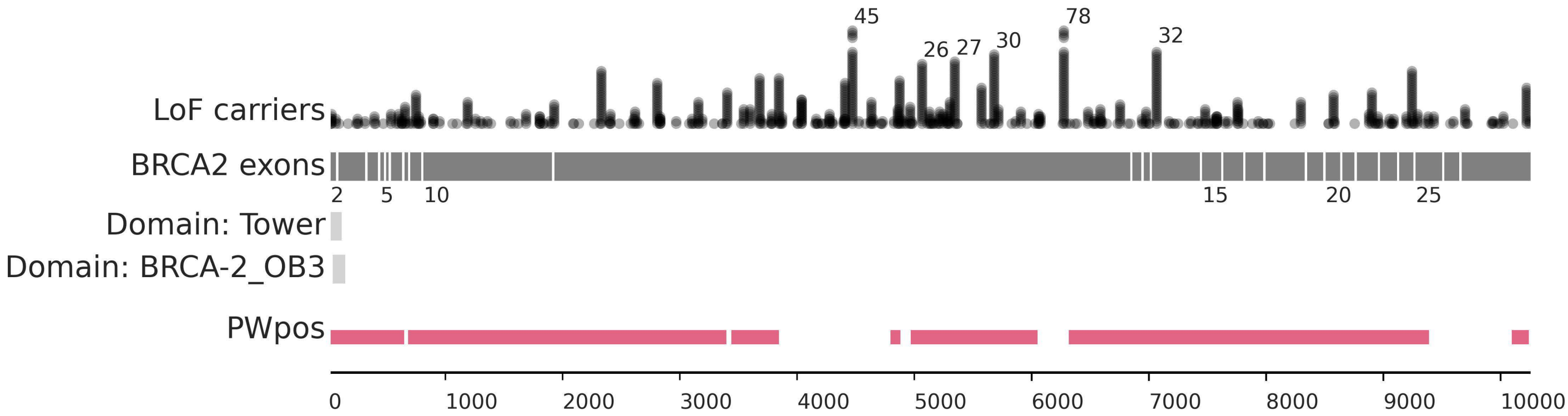


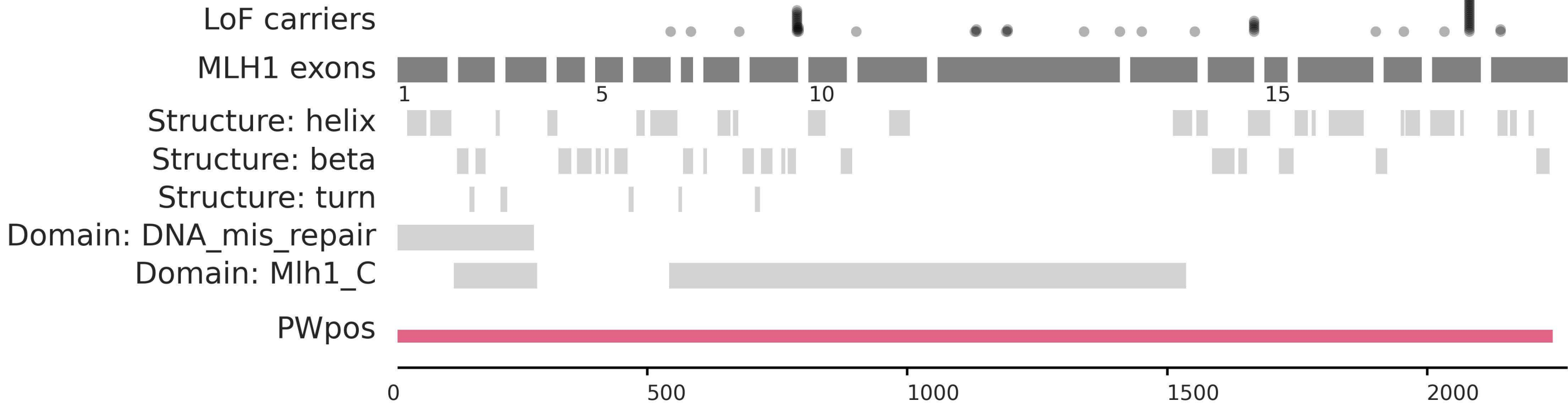












78

LoF carriers

MSH2 exons

Structure: helix

Structure: beta

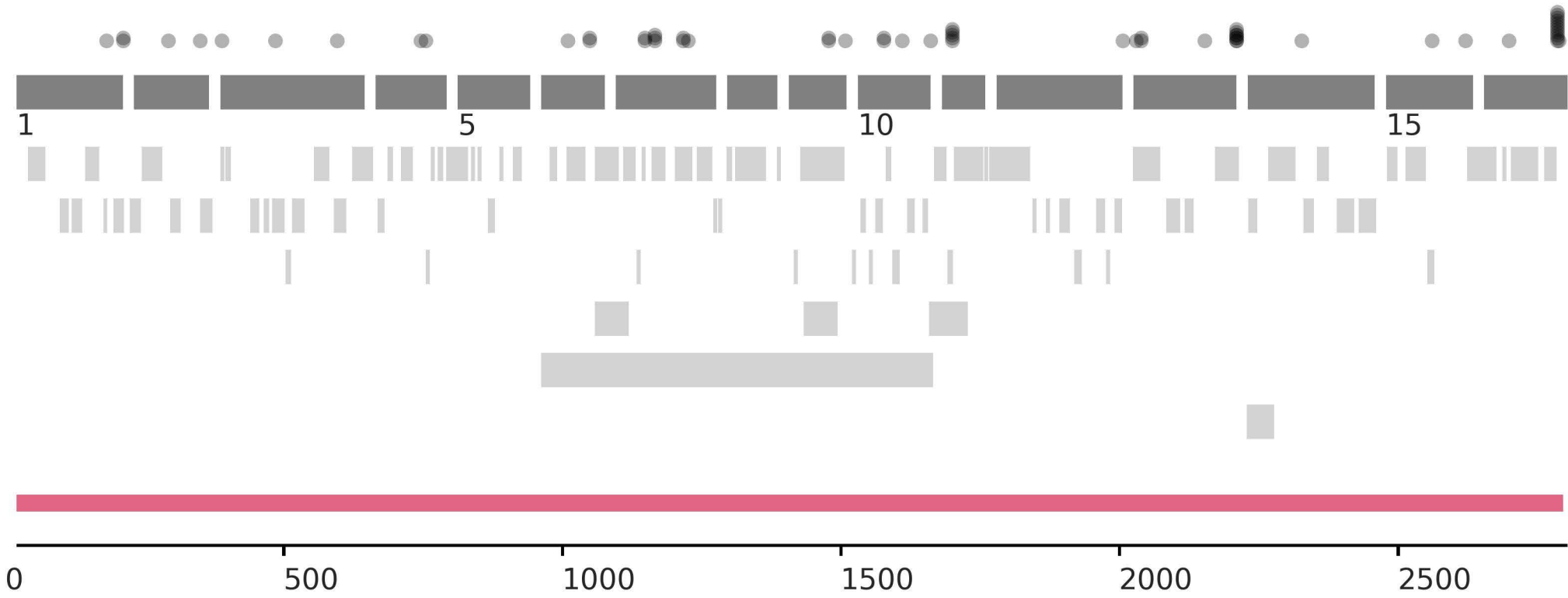
Structure: turn

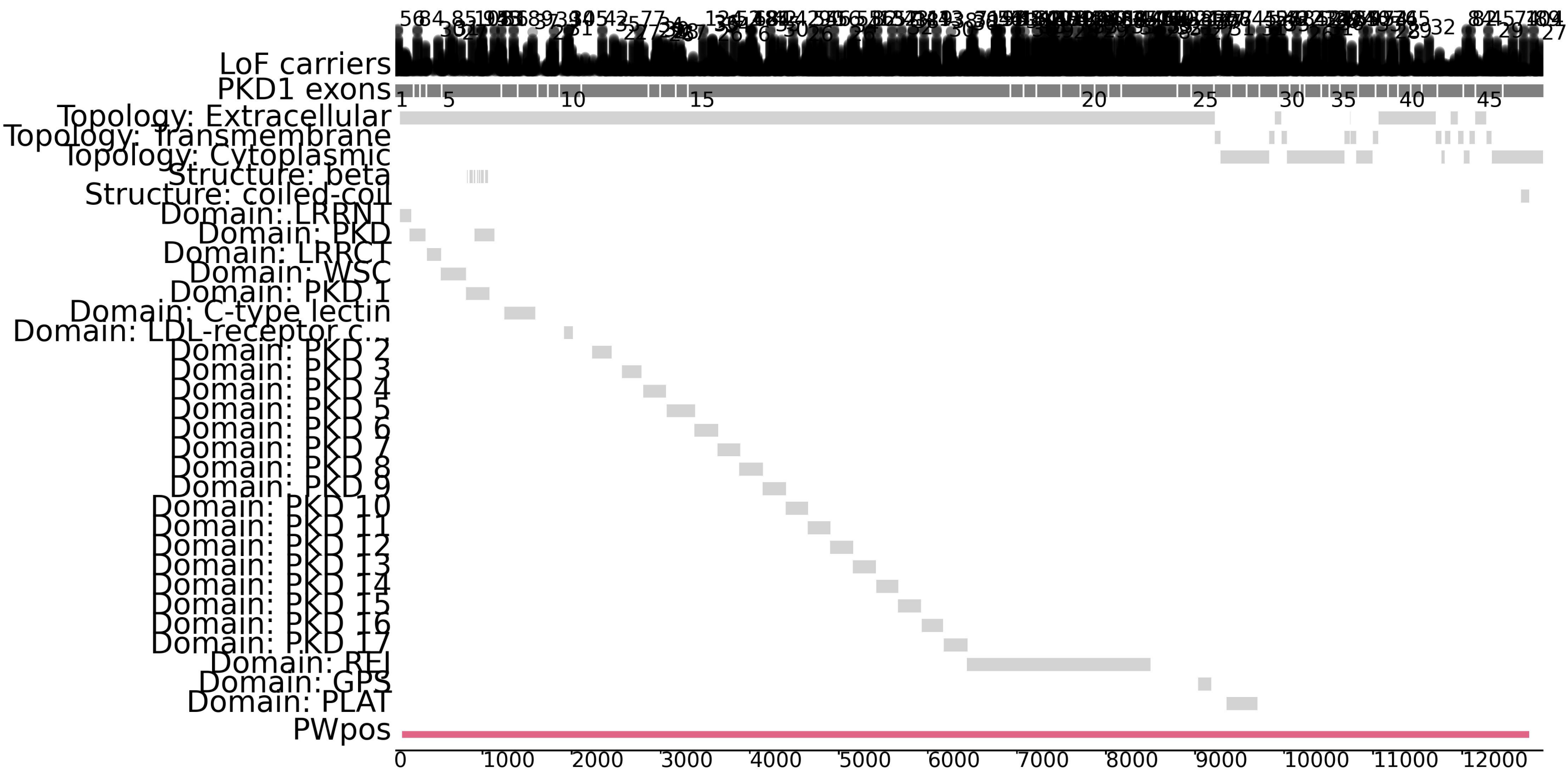
Structure: coiled-coil

Domain: MUTSd

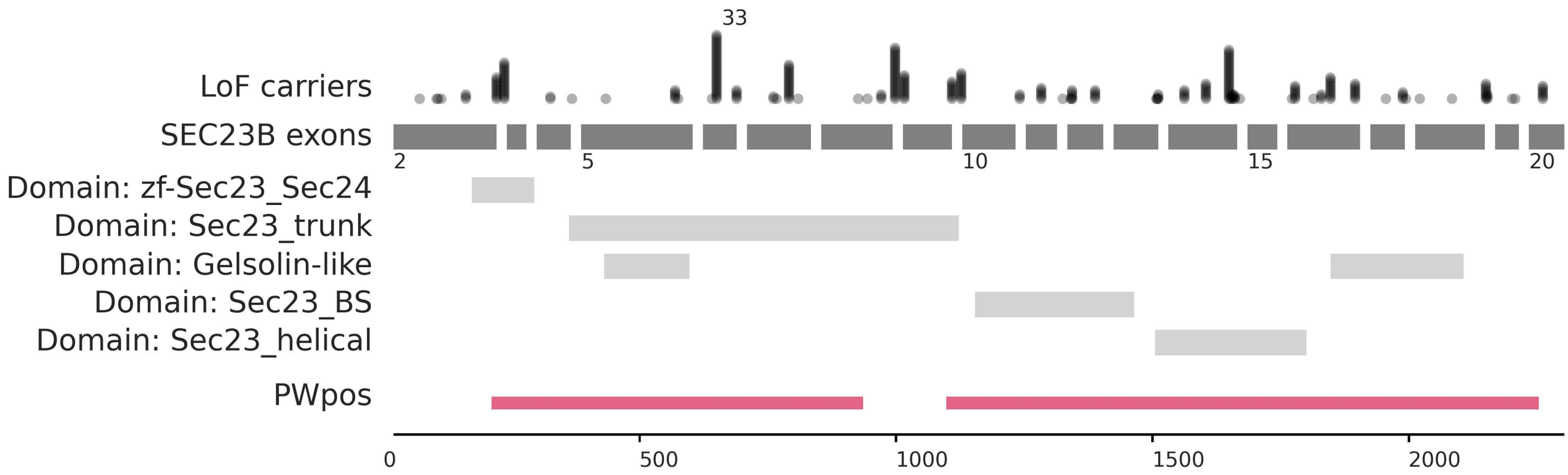
Domain: DNA\_MISMATCH\_R...

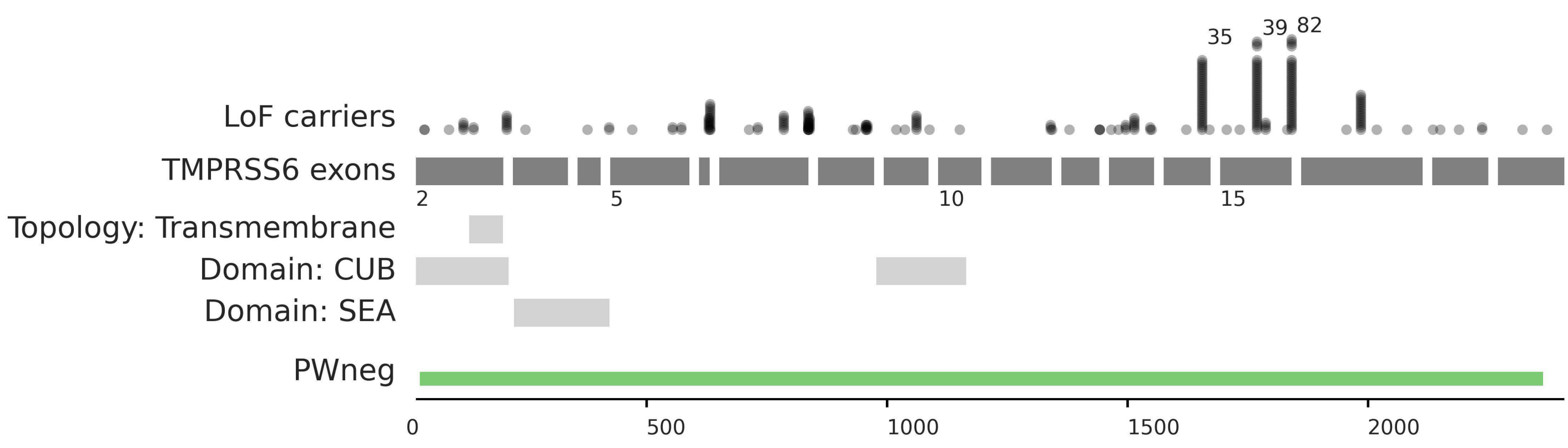
PWpos

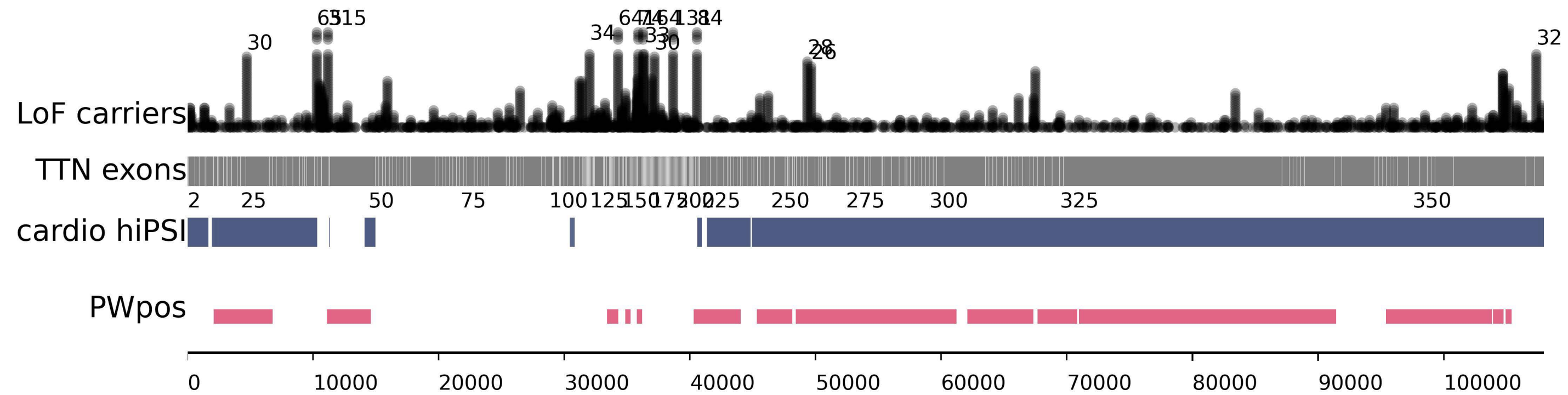






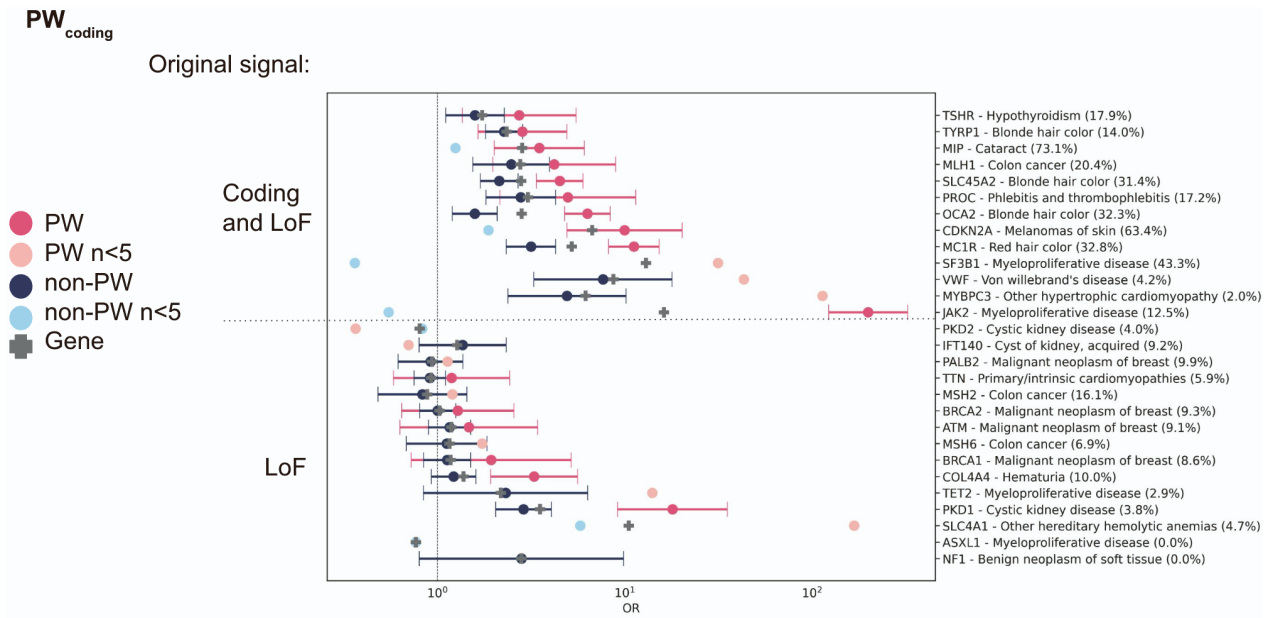




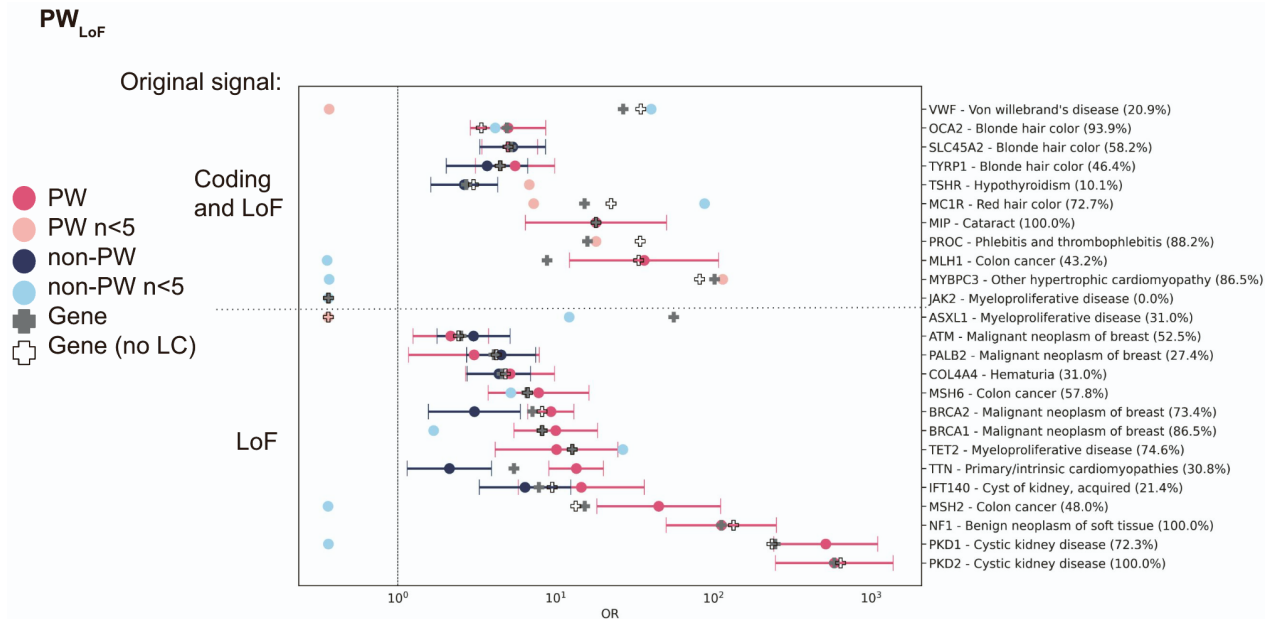


**Figure S5. Power Window results for the 37 significant models.** Tracks are drawn against each gene's canonical coding transcript. Coding position is shown at the bottom scale. Carriers: each dot represents an individual carrier and are stacked for each carrier at a given position. For brevity, carriers are trimmed to 35 and total number of carriers is indicated when total carriers>40. Exons: exons (dark grey to scale; introns not to scale). Exon number indicated below exon track. Topology, secondary structures and major domains are annotated according to UniProt. PWpos: merger of all significant windows with a positive direction of effect ( $\beta > 0.5$  or  $OR > \text{cutoff}$  for that gene; pink). PWneg: merger of all significant windows with a negative direction of effect ( $\beta < -0.5$ ; green). Significant associations are determined as indicated in the main text.

**A**

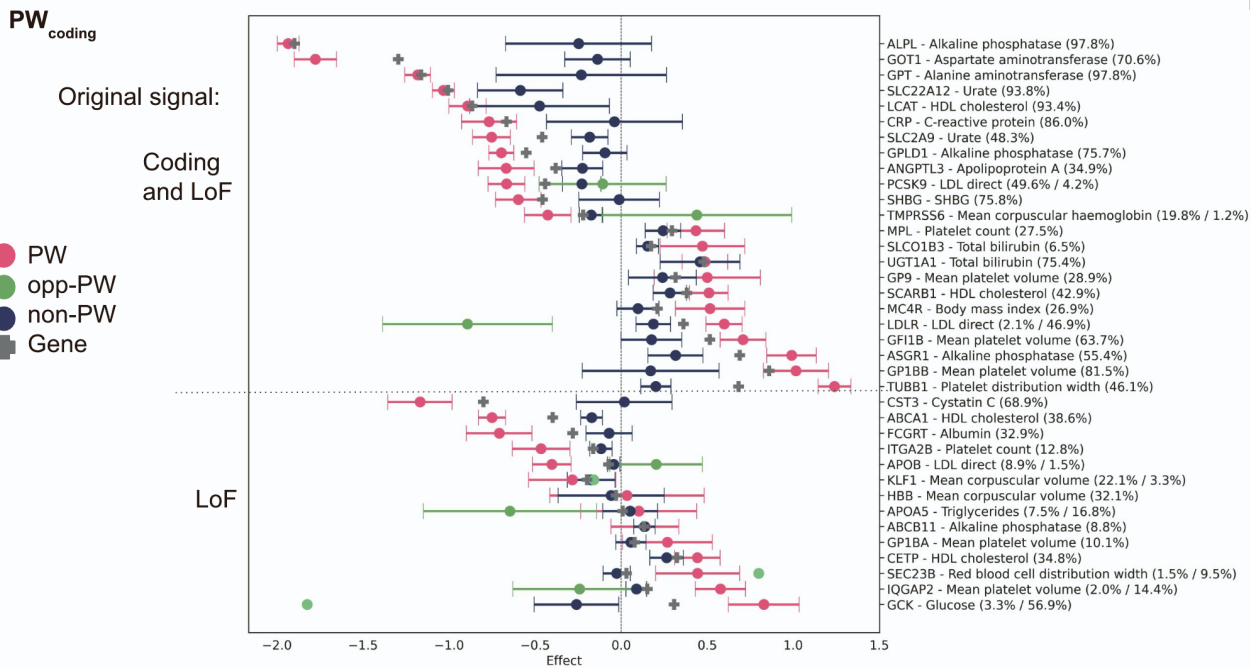


**B**

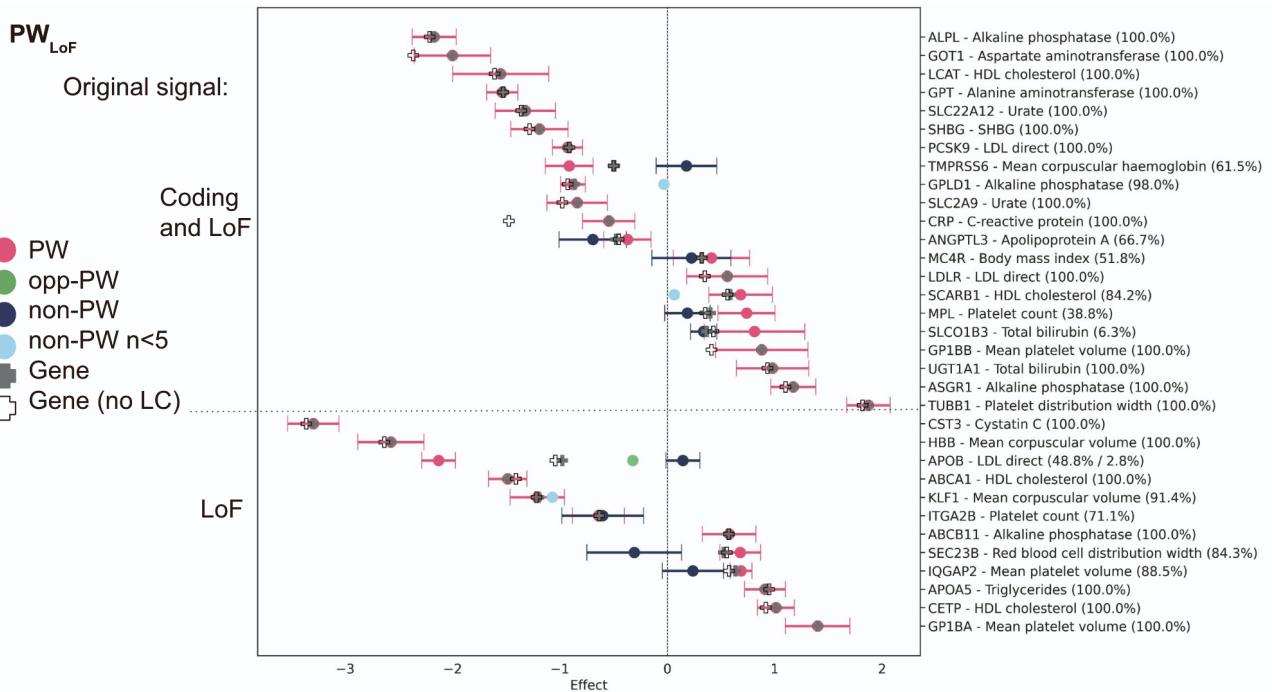


**Figure S6. Performance of Power Window for all binary traits in the 117k UKB test set.** A) Coding model B) LoF model. The odds ratios for PW models are shown in pink, and the excluded regions are shown in blue. The gene-phenotype label includes in parentheses an indication of what percent of rare variant carriers in the gene were included in the PW model. For brevity, 95% confidence intervals are only shown when there are at least 5 case carriers: datapoints from <5 carriers are not well supported and require additional data to confirm. The genes were grouped based on the original genome-wide association as follows: **LoF**: original whole-gene associations had an absolute beta at least 3x as high in the LoF model as the Coding model; **Coding and LoF**: original whole-gene associations had a beta <3x as high for LoF as for Coding (no gene-phenotype combinations had a whole-gene Coding model absolute beta that was at least 3x higher than the whole-gene LoF beta). Two whole-gene models are shown: gray filled cross, which includes all qualifying variants in the gene, and black empty cross, in which LOFTEE LC variants are excluded.

**A**



**B**

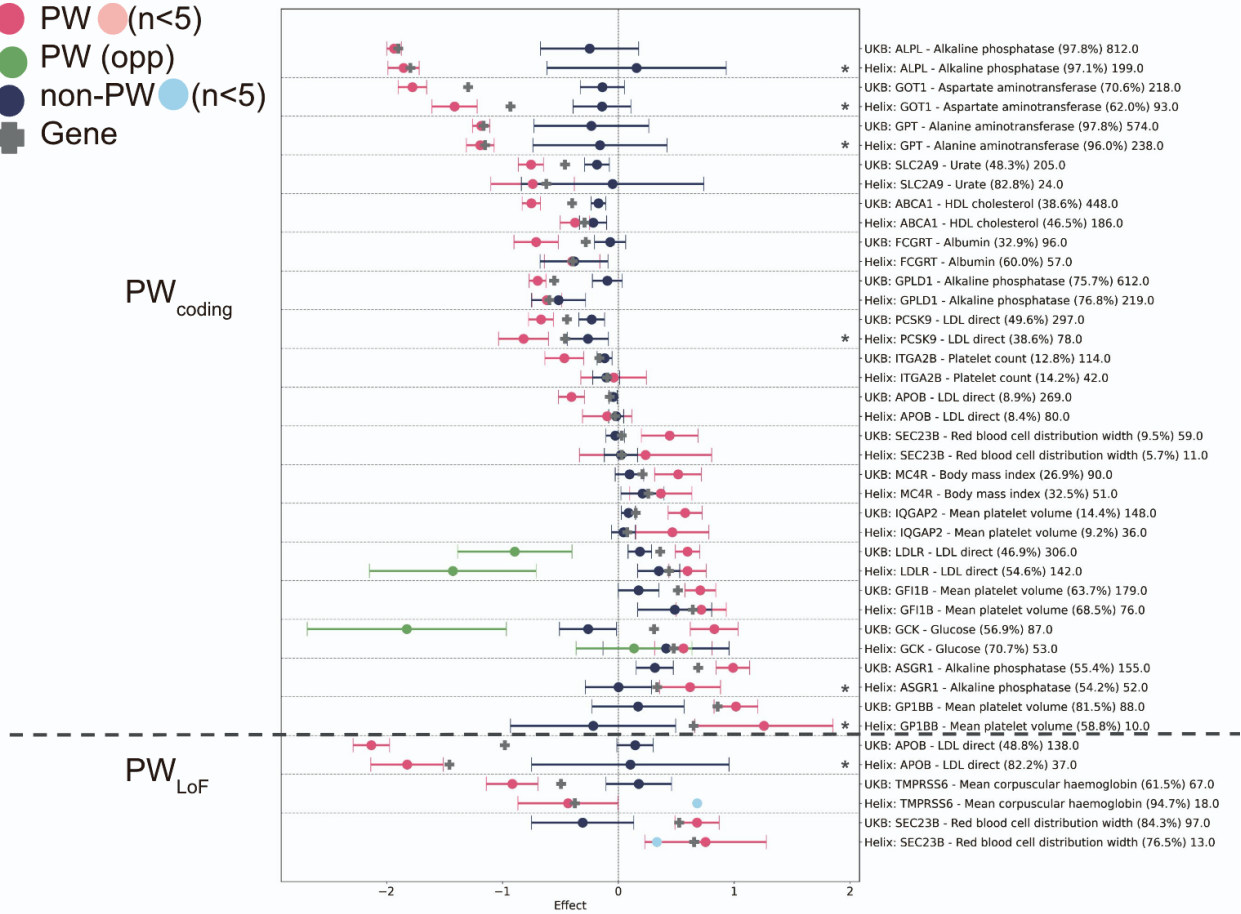


**Figure S7. Performance of Power Window for all quantitative traits in the 117k UKB test set.** A) Coding model B) LoF model. The effect sizes (normalized phenotypes) for PW models are shown in pink, and the excluded regions are shown in blue. The gene-phenotype label includes in parentheses an indication of what percent of rare variant carriers in the gene were included in the PW model (when there are 2 models for a gene, the percent shown is for negative model / positive model). When a model for the opposite direction of the main effect was built (opp-PW), it is shown in green. For brevity, 95% confidence intervals are only shown when there are at least 5 carriers: datapoints from <5 carriers are not well supported and require additional data to confirm. The genes were grouped based on the original genome-wide association as follows: **LoF**: original whole-gene associations had an absolute beta at least 3x

as high in the LoF model as the Coding model; **Coding and LoF**: original whole-gene associations had a beta  $< 3x$  as high for LoF as for Coding (no gene-phenotype combinations had a whole-gene Coding model absolute beta that was at least  $3x$  higher than the whole-gene LoF beta). Two whole-gene models are shown: gray filled cross, which includes all qualifying variants in the gene, and black empty cross, in which LOFTEE LC variants are excluded.

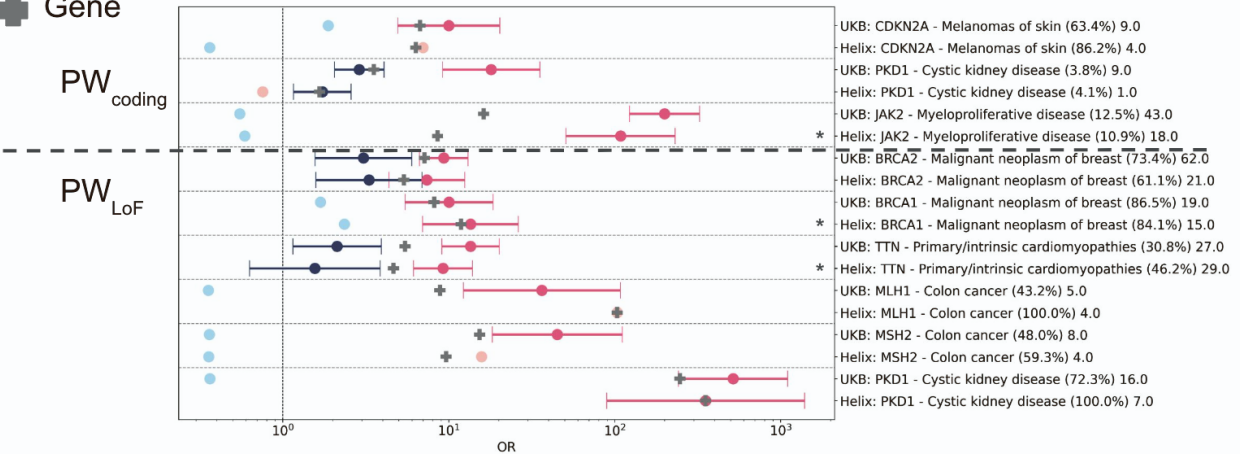
**A**

- PW (n<5)
- PW (opp)
- non-PW (n<5)
- Gene



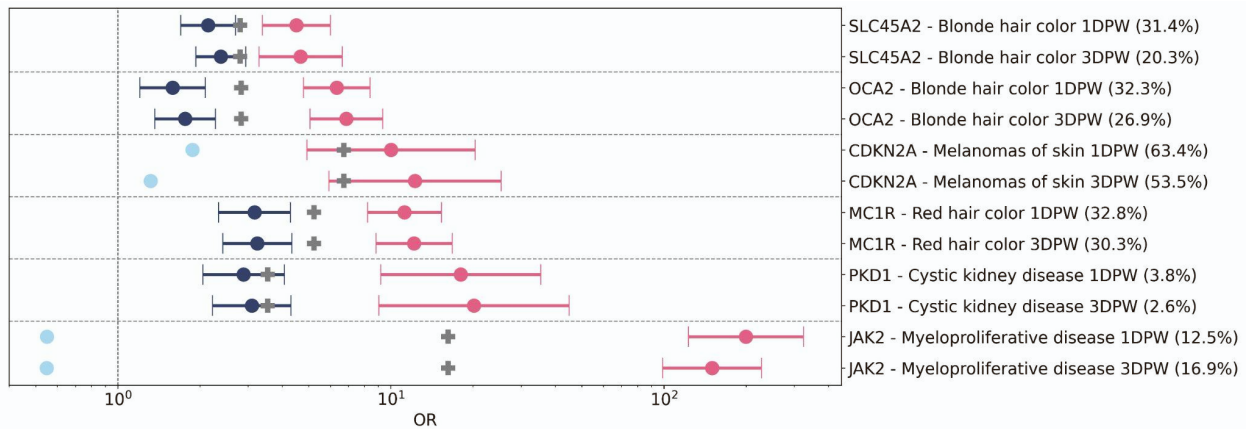
**B**

- PW (n<5)
- PW (opp)
- non-PW (n<5)
- Gene



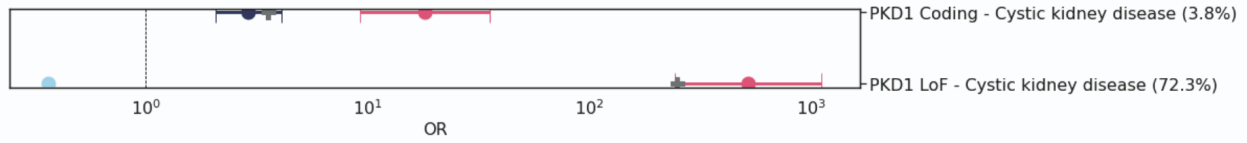


**Figure S8. Replication in the Helix cohorts.** A) Quantitative traits. The effect sizes (normalized phenotypes) for PW models are shown in pink, and the excluded regions (non-PW) are shown in blue. When a model for the opposite direction of the main effect was built (PW-opp), it is shown in green. Significant models ( $p < 0.002$  with Bonferroni correction for multiple tests) are marked with an asterisk. B) Binary traits. The odds ratios for PW models are shown in pink, and the excluded regions (non-PW) are shown in blue. The gene-phenotype label includes in parentheses an indication of what percent of rare variant carriers in the gene were included in the PW model (when there are 2 models for a gene, the percent shown is for negative model / positive model). For brevity, 95% confidence intervals are only shown when there are at least 5 case carriers. Only models that were significant in the UKB117k test cohort are tested. The UKB data shown are from the 117k test cohort. The total number of PW variant carriers (or case carriers for binary traits) is shown at the end of each row. The results for *PKD1* coding and LoF as well as *MLH1* LoF were not included in the final counts as they did not have case carriers in both PW and non-PW regions for comparison.

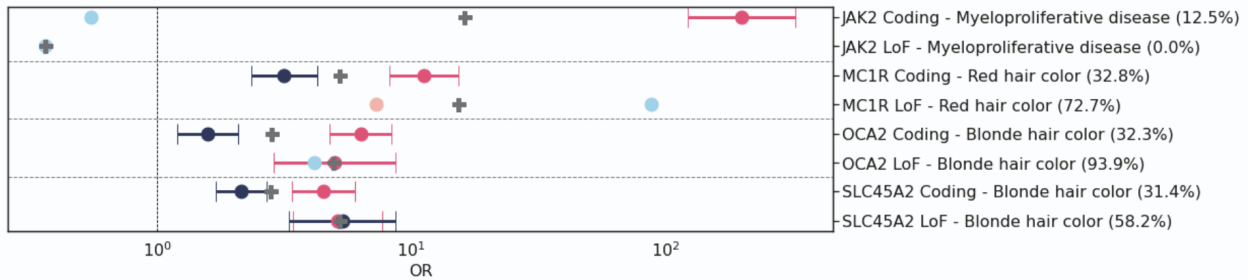


**Figure S9. 3D Power Window for binary traits.** As in Figure 5, PW models are shown in the test set of 117k individuals for significant  $PW_{\text{coding}}$  models. Each gene-phenotype pair is shown twice, for the 1D and 3DPW models, with 95% confidence intervals. The model results are very similar in terms of percent of carriers kept in the model and final effect sizes. The percentage to the right of the phenotype is the percent carriers retained in the model.

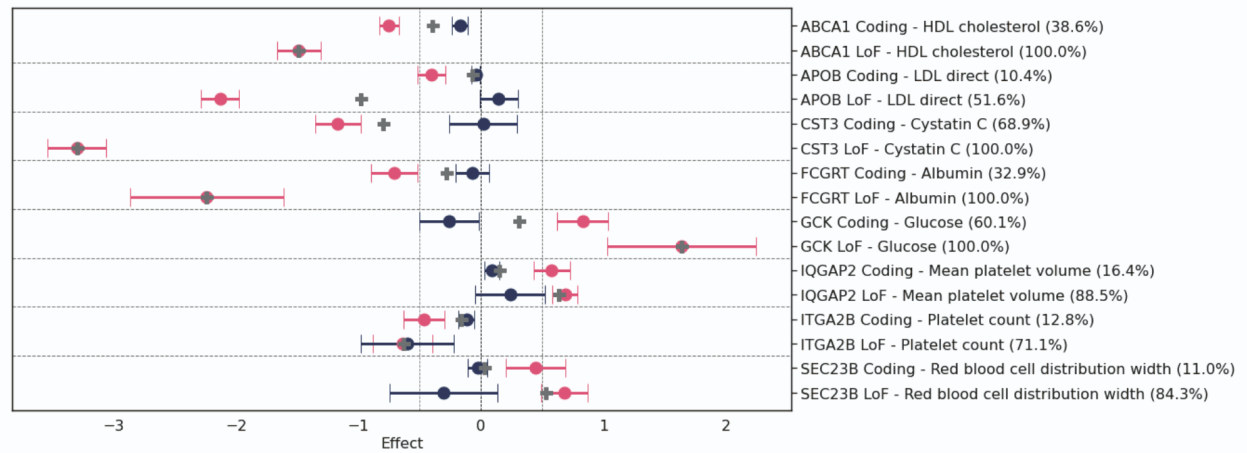
A - Binary traits- original signal primarily for LoF model



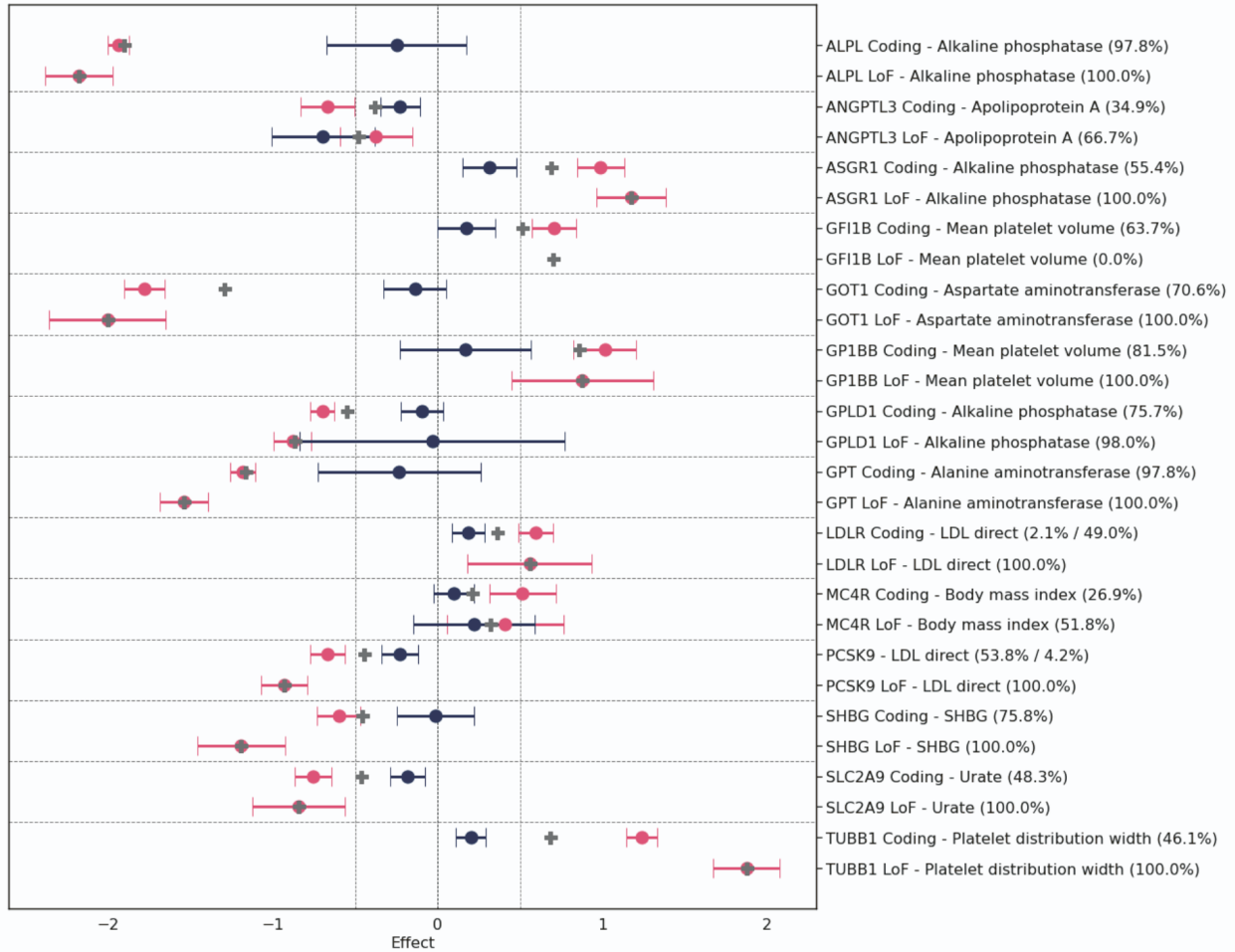
B - Binary traits- original signal for both models



C - Quantitative traits- original signal primarily for LoF model

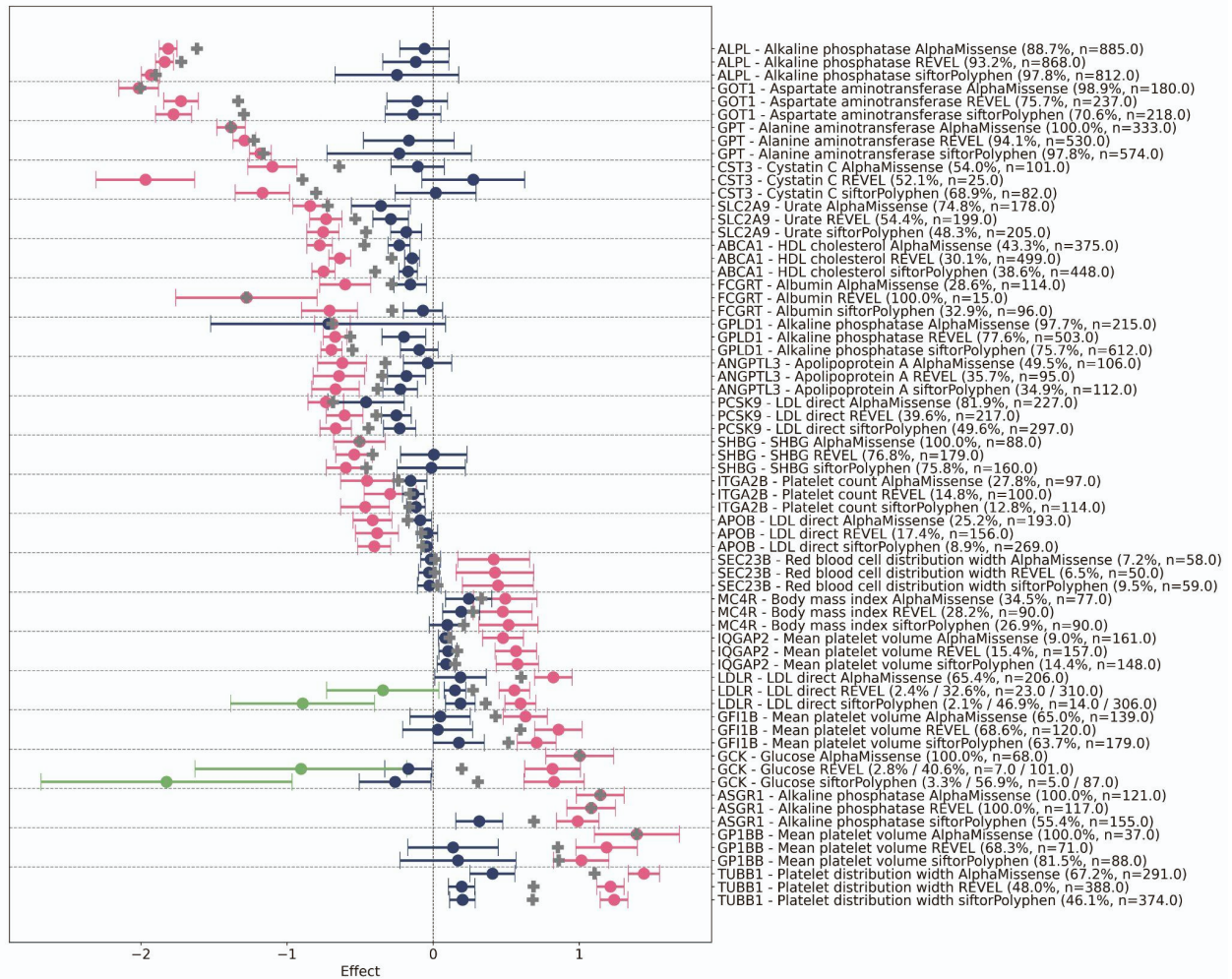


D - Quantitative traits- original signal for both models

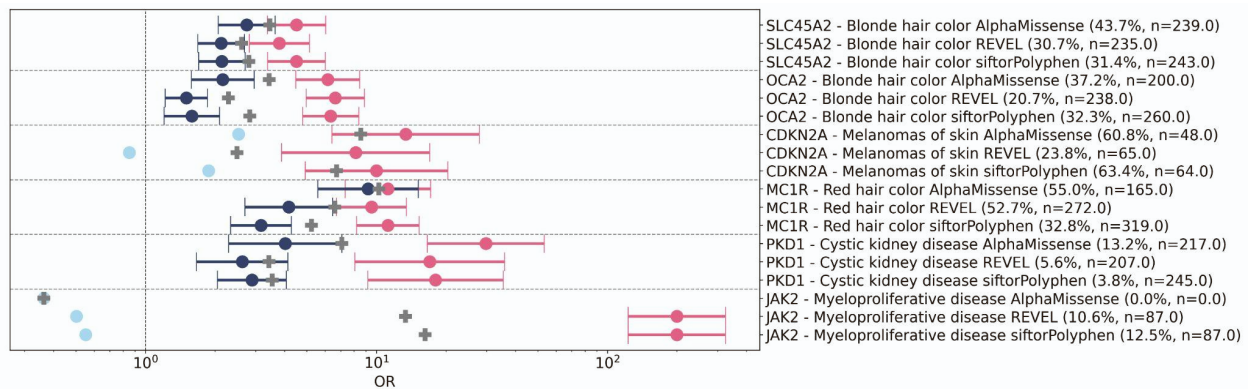


**Figure S10. Performance of significant  $PW_{\text{coding}}$  models compared to  $PW_{\text{LoF}}$  models for the same phenotype.** A) Binary traits where the original whole-gene associations were primarily for the LoF model (LoF beta >3x Coding beta); B) binary traits where the original whole-gene associations were for both Coding and LoF models (beta <3x as high in LoF as in Coding) C) Quantitative traits where the original whole-gene associations were primarily for the LoF model (LoF beta >3x Coding beta); D) quantitative traits where the original whole-gene associations were for both Coding and LoF models (beta <3x as high in LoF as in Coding). Data are shown for the 117k UKB test set, with 95% confidence intervals. These  $PW_{\text{coding}}$  models fit the criteria for statistical significance and are also displayed in Figure 3. For the  $PW_{\text{coding}}$  models, the effect size of the  $PW_{\text{coding}}$  model is similar to that of the  $PW_{\text{LoF}}$  model for 1 of the 3 binary models (A) and 3 of 8 quantitative models (C) where the whole-gene association was primarily LoF; this was also true for all 4 binary models (B, excluding *JAK2* where there was no LoF signal) and 14 of 19 quantitative models (D) where the whole-gene association was for both coding and LoF. In these cases, the model was able to identify coding variants with what appear to be complete LoF effects. For the remaining models, the  $PW_{\text{coding}}$  effect size is a substantial improvement over the whole-gene coding model, but is still not as extreme as those with LoF variants. These may reflect variants that reduce function without completely losing it.

A



B



**Figure S11. Comparison of significant PWcoding models built with REVEL, AlphaMissense, or sift/Polyphen.**

As in Figure 3, PW models are shown in the test set of 117k individuals with 95% confidence intervals for significant PW<sub>coding</sub> models for A) quantitative and B) binary traits. Each gene-phenotype pair is shown three times, for when REVEL was used as the bioinformatic predictor of what is damaging (cutoff 0.25), when AlphaMissense was, and when sift/Polyphen were (default model used in the paper, remove missense variants that were benign by both of these). The model results are very similar in terms of percent of carriers kept in the model and final effect sizes. The

percentage to the right of the phenotype is the percent carriers retained in the model; when there is both a negative and positive model, the percent carriers retained in the negative model is shown first.