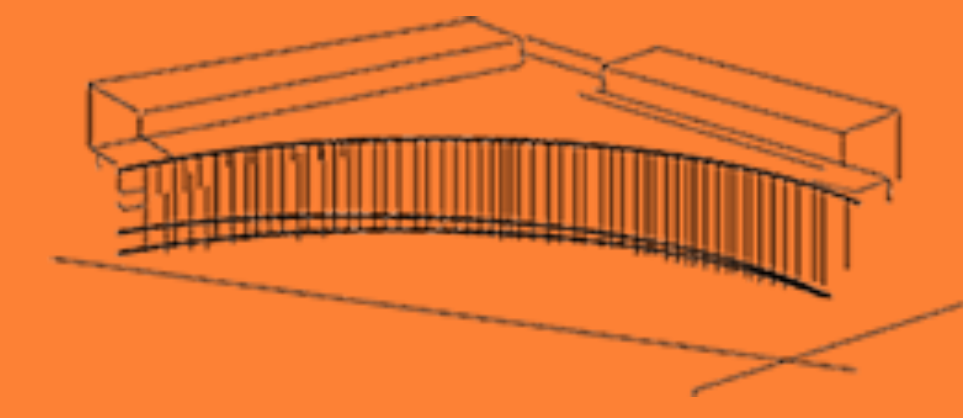


# Identifying Functional Protein Domains Positions Using Population Variation



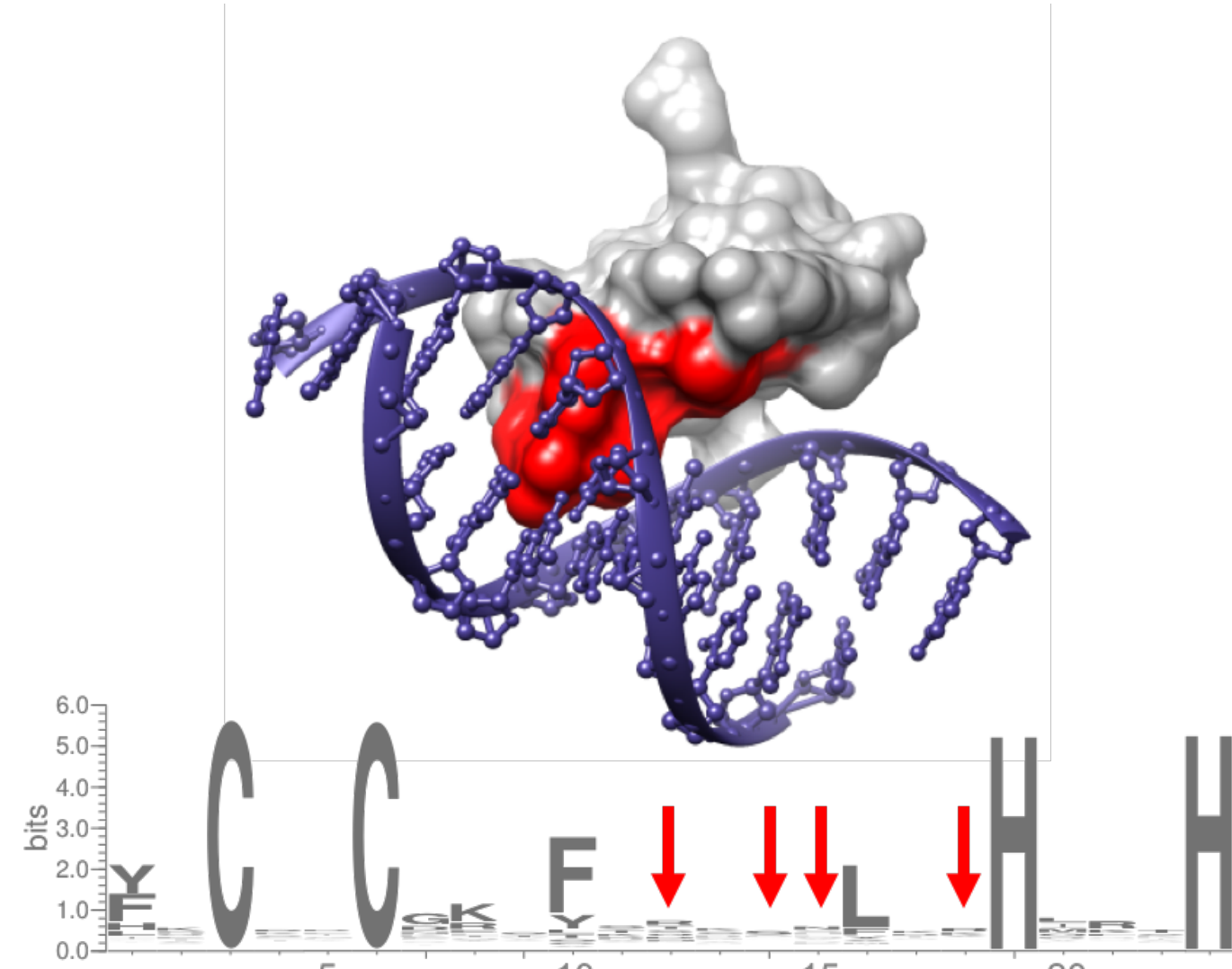
## Abstract

Identifying essential protein residues is a major challenge in computational biology. We developed a computational framework that incorporates variation from large-scale population sequencing data, to highlight functional positions of protein domains.

## Background

### Challenges in identifying proteins "key" residue positions

- Structure-based methods - limited to ~30% of human proteins with a solved structure.
- Sequence-based methods - often miss non-conserved specificity-determining positions (SDPs):

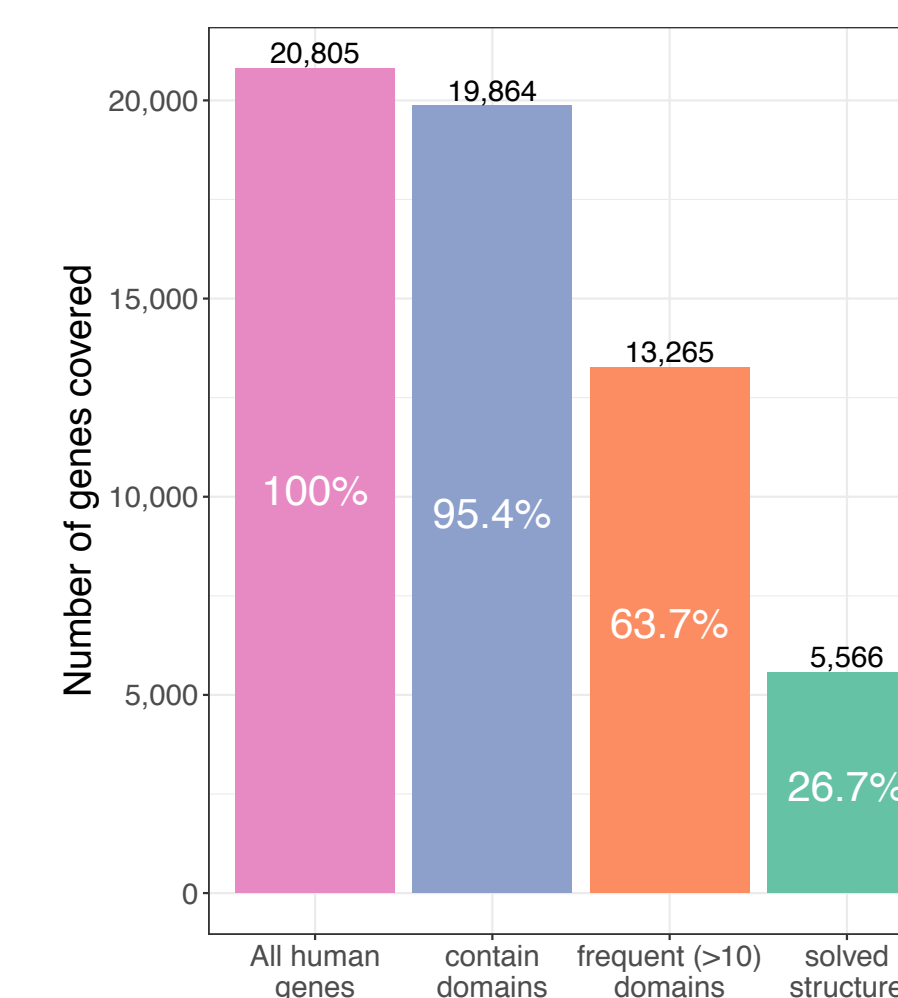


**Fig 1: Non-conserved SDPs exemplified in the Cys<sub>2</sub>-His<sub>2</sub> Zinc finger domain.** Known<sup>1</sup> Functionally-important positions are colored in red on the domain surface<sup>2</sup>, and pointed by red arrows on the domain sequence logo<sup>3</sup>. They are critical for the DNA-binding specificity, and are not conserved.

## Results

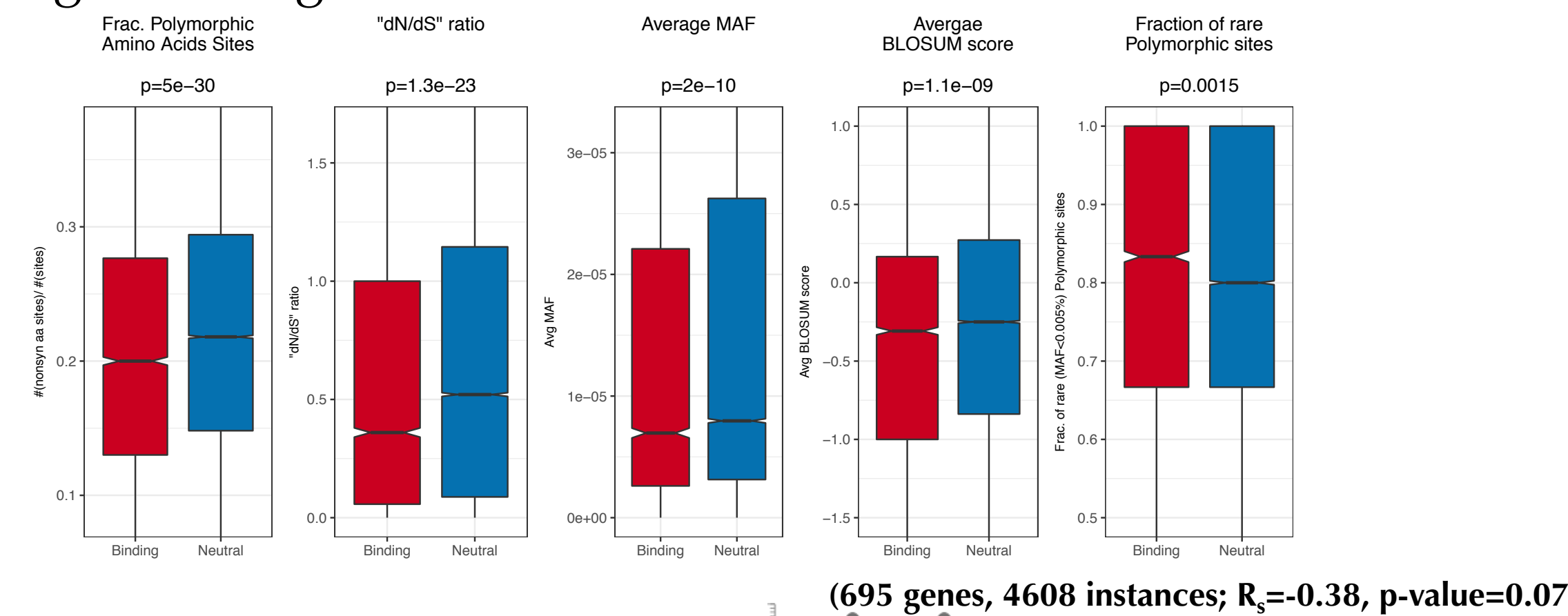
### Fig 3: Human genome coverage.

Most of the human genes have at least one protein domain (> 95%)<sup>5</sup>. Large fraction of the genes are covered by frequent domains (~ $\frac{2}{3}$ ). Methods that rely on structure have significantly lower coverage (<  $\frac{1}{3}$ )<sup>5</sup>.



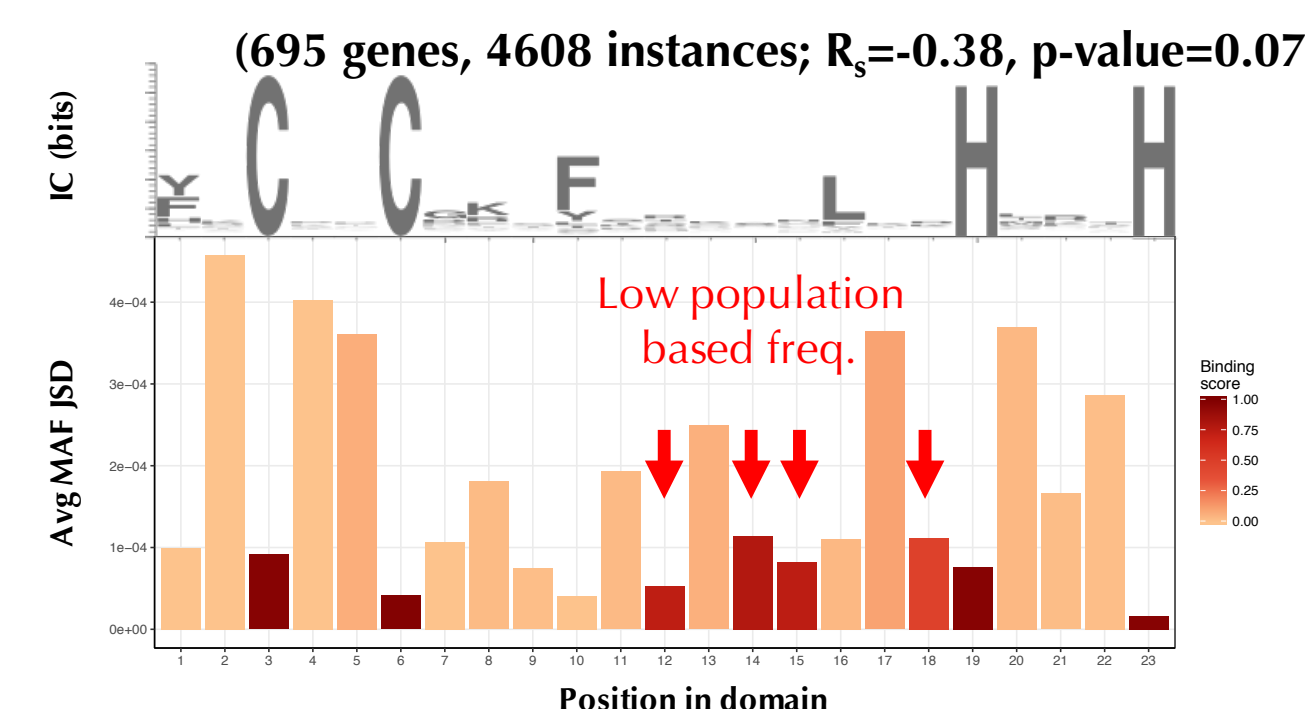
### Fig 4: Population-based positional features.

Features calculated for protein domains (≥ 10 instances) positions: 409 domains, 5,842 binding and 44,692 neutral. Positional labels according to Kobren & Singh binding scores<sup>6</sup>. P-values are Wilcoxon one-tailed test.

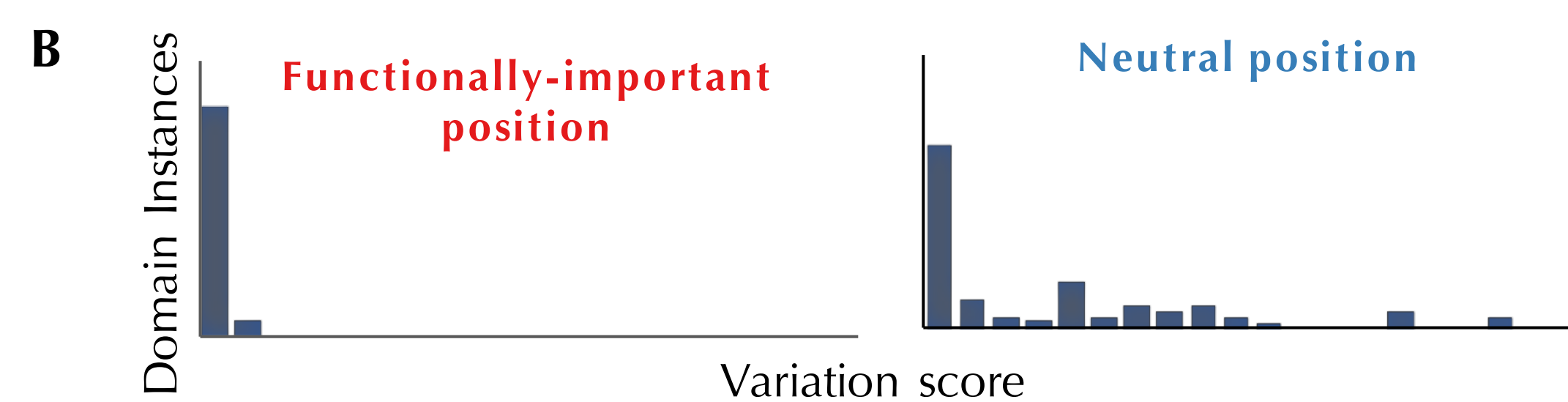
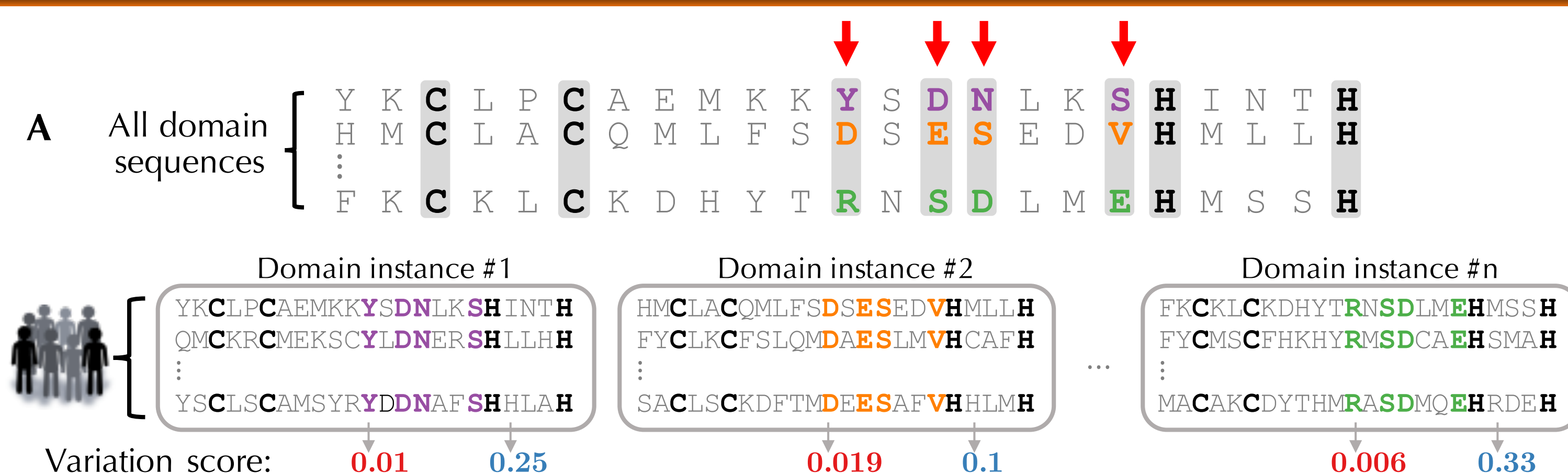


### Fig 5: Predicting binding potential.

Jensen-Shannon Divergence (JSD) per domain position, based on population allele frequency. The known SDPs in the C2H2-zinc finger domain have low JSD, as expected.



## Methods



**Fig 2: Outline of approach using the C2H2-zinc finger domain as an example.** (A) Aggregation of SNPs population variation score across domain instances per domain position. (B) The distributions for two such positions are represented here as histograms.

### Population-based domain-centered method

Predict functional importance from natural variations within sequencing data of 60,000 healthy individuals<sup>4</sup>:

- Single nucleotide polymorphisms (SNPs) in each proteins domain, from all protein sequences across all individuals are aggregated by domain position (Fig2A).
- Domain position distributions are analyzed to predict functional importance.

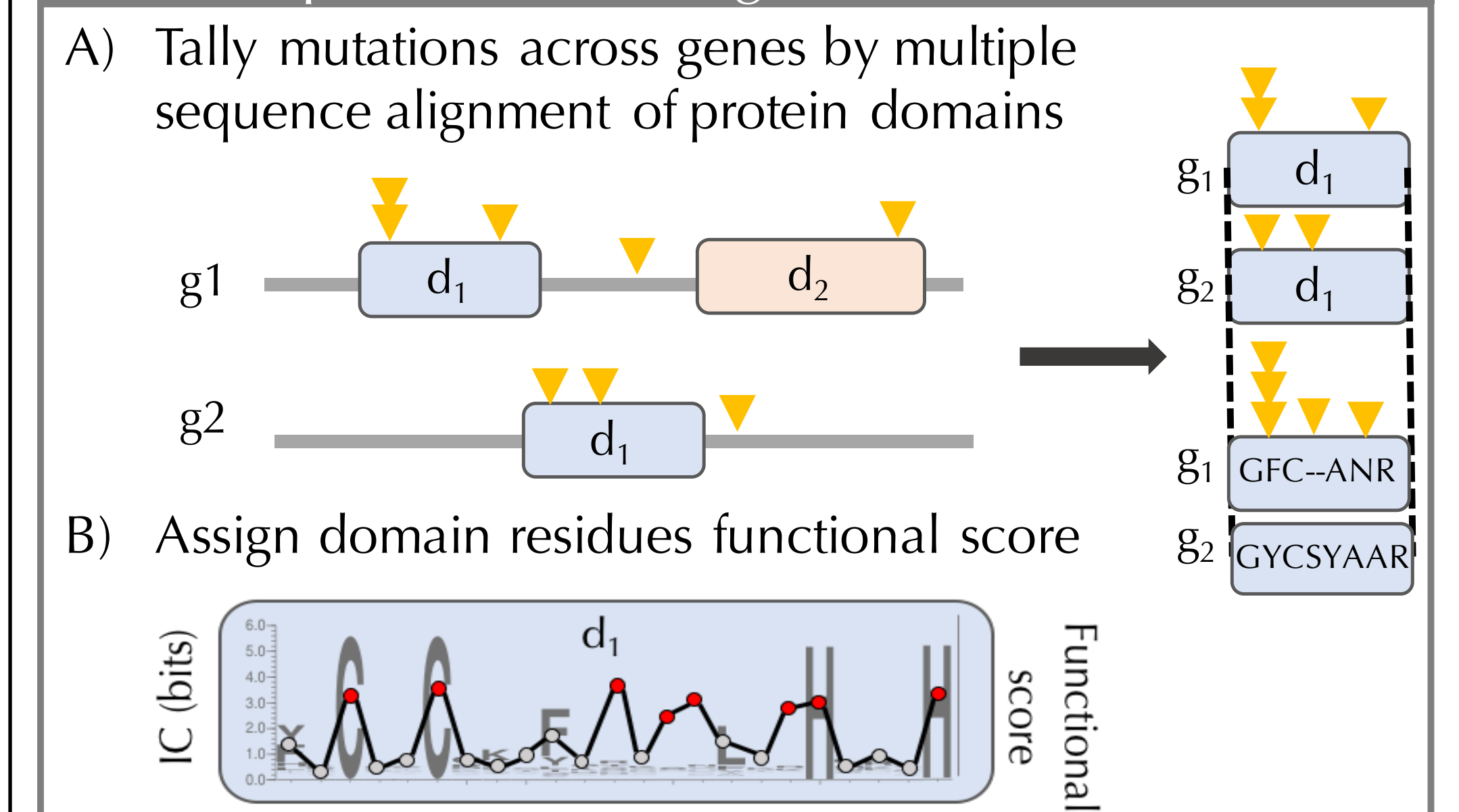
\*crucial protein domain positions are more conserved across the healthy population, even if they vary across proteins (Fig2B).

## Significance

Systematic identification of "key" residue positions is critical for understanding protein **functions**, predicting **interactions**, and characterizing the **effect of mutations** in the context of **disease**.

### Future plans

#### Map mutations, assign functional scores



### Domains enrichment

- C) Determine domains mutation burden: mutations count weighted by functional score
- D) Identify domains enriched for mutations in functional positions: Permutation test

### Functional hotspots

- E) Identify functionally-important hotspots of mutations in domains: Binomial test

**Fig 6: Proposed method for using domains positional functional scores in cancer mutations analysis.**

## References

- [1] Wolfe et al. Ann. Rev. Biochem. (2001) [2] Pettersen et al. J. Comput. Chem. (2004) [3] Crooks et al. Genome Res. (2004) [4] Lek et al. Nature (2016) [5] Yates et al. Nucleic Acids Res. (2016) [6] Kobren and Singh (Unpublished).

## Acknowledgements

We would like to thank Shilpa N. Kobren for helping with the computational framework, and the rest of the Singh lab for their helpful insights. Funded by T32 HG003284, NIH CA208148, and the Forese Family Fund for Innovation.