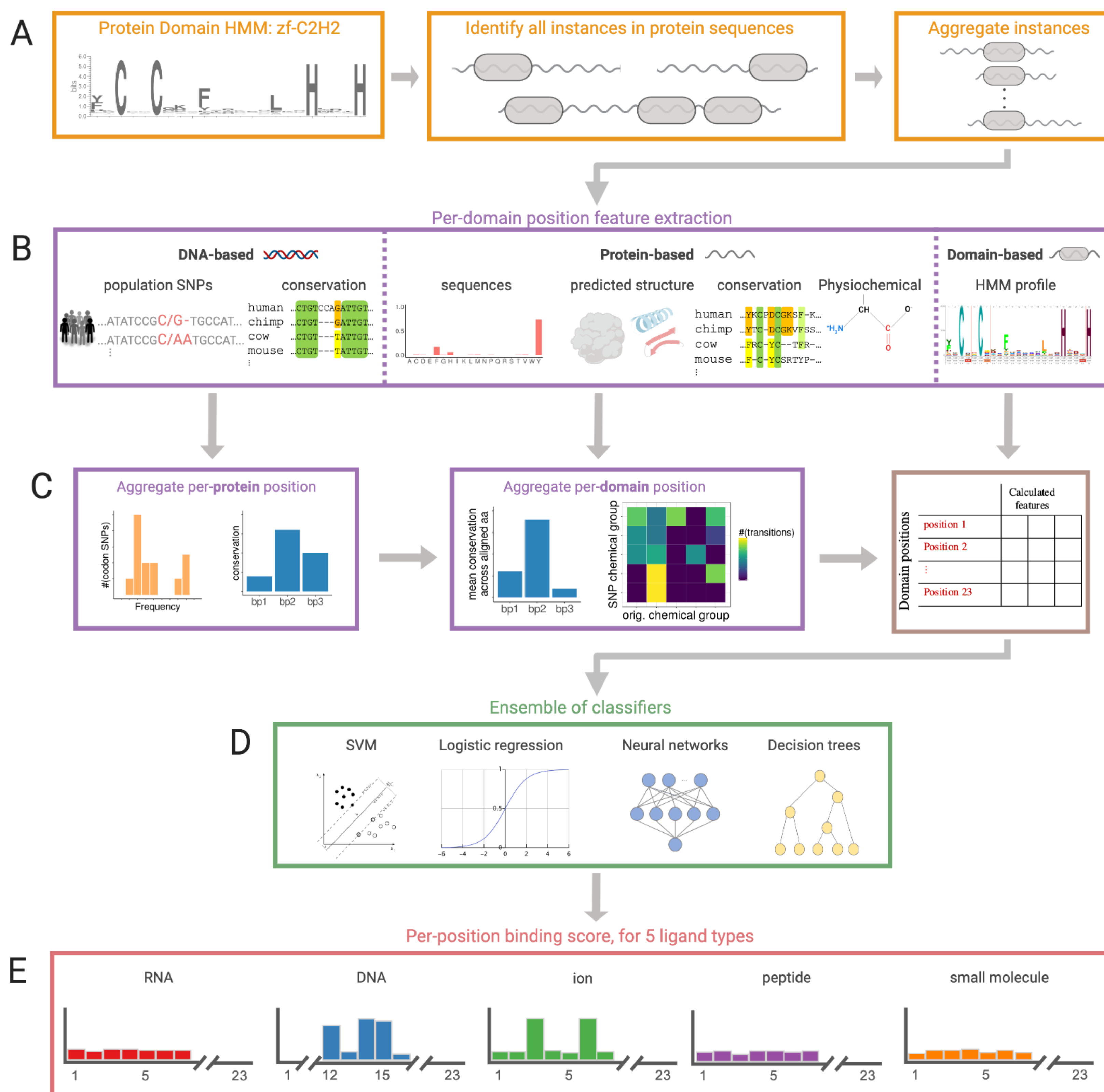


## 1 Abstract

We present **dSPRINT**: domain Sequence-based **P**rediction of **I**nteraction sites, an ensemble of machine learning classifiers using a novel stacking architecture, that predict binding positions within protein domains.

## 3 Methods: pipeline overview

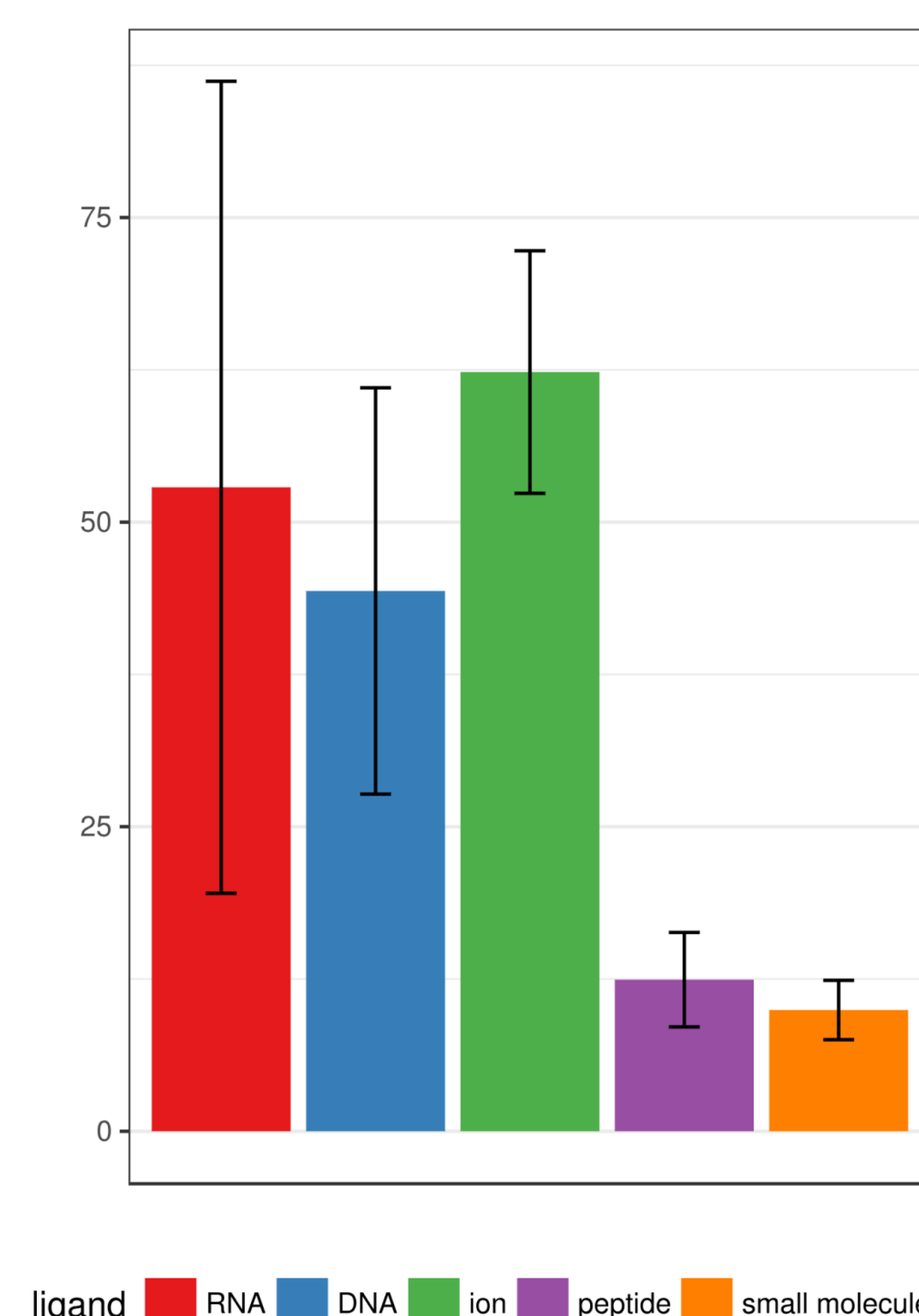


**Fig 2: dSPRINT workflow for domain-centered per-position prediction.**

## 6 Results: global evaluation

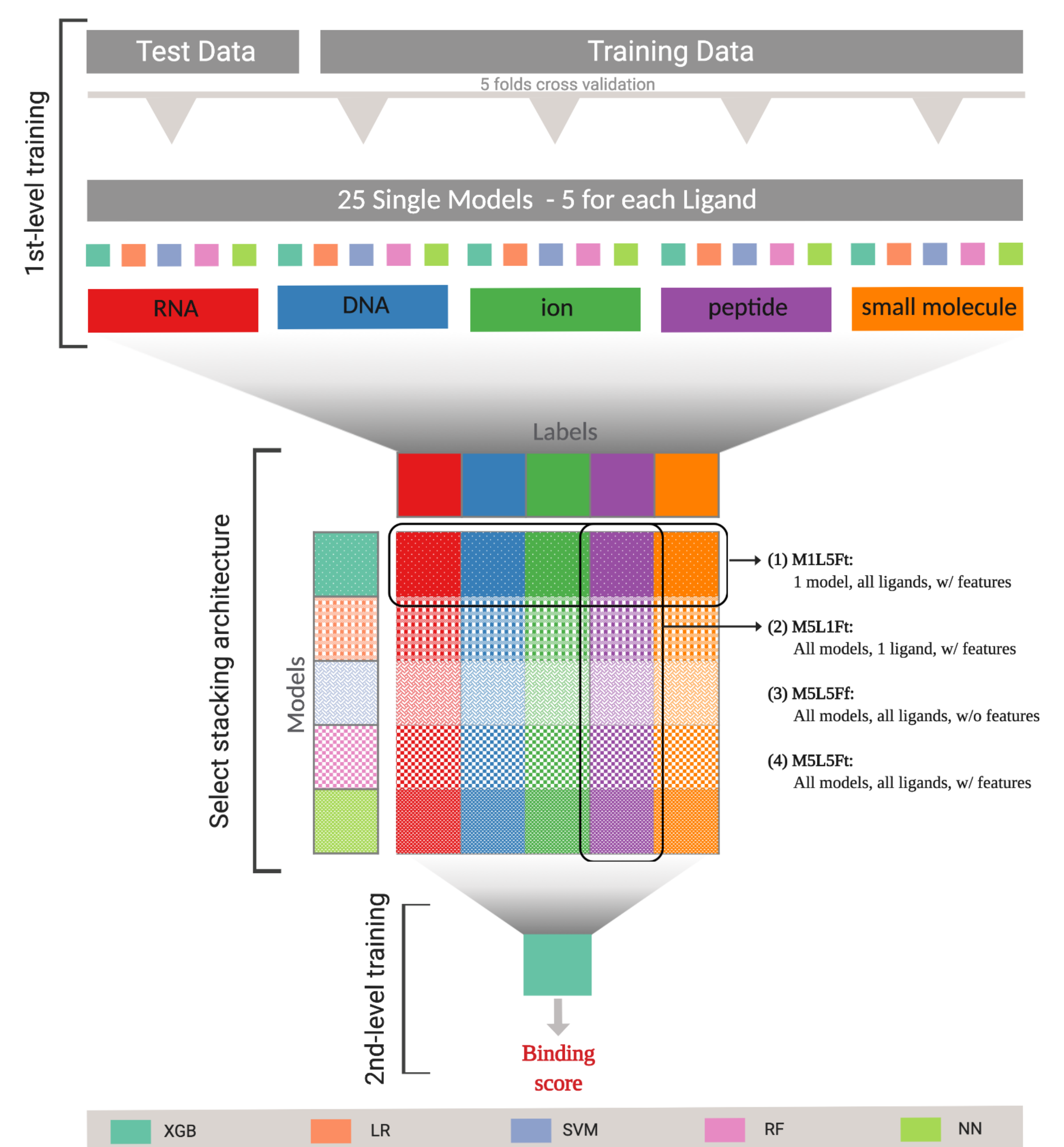
	Number of most reliable predictions chosen per domain	Number of domains with at least one Correct prediction	Fraction of domains With correct predictions
RNA	1	13	62%
	3	16	76%
	5	18	86%
DNA	1	22	67%
	3	27	82%
	5	29	88%
ion	1	47	52%
	3	60	66%
	5	64	70%
peptide	1	29	40%
	3	44	61%
	5	48	67%
small molecule	1	72	55%
	3	96	73%
	5	104	79%

**Table 1: The most reliable prediction(s) in each binding domain.**



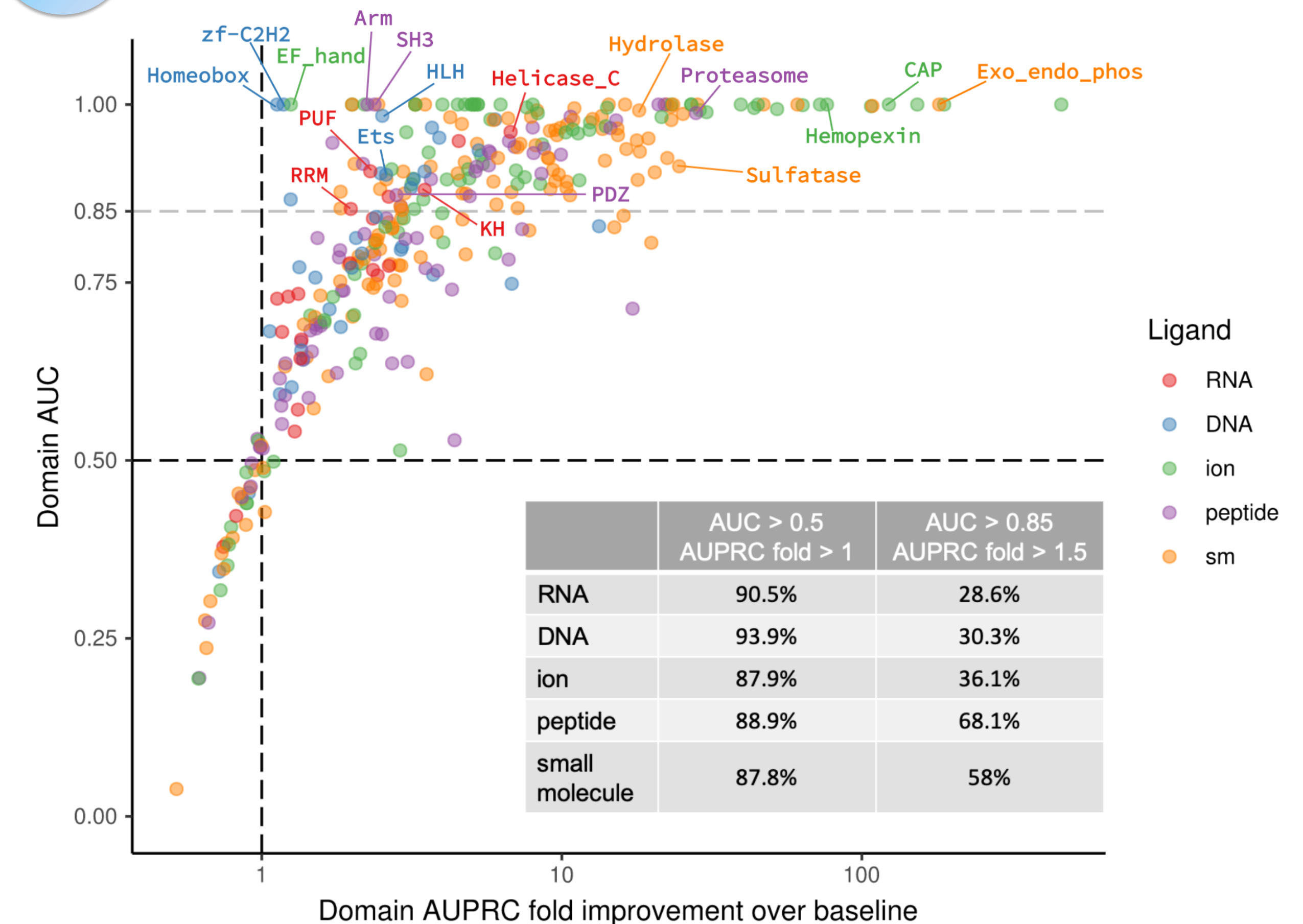
**Fig 6: Ligand AUPRC fold improvement.** The ratio of the AUPRC to a baseline corresponding to the fraction of binding positions of that ligand at that CV fold.

## 4 Methods: ML stacked architecture



**Fig 3: Ligand-combined classifier stacking architecture.** Five base-models are trained for each ligand in 5-fold cross validation. This results in 25 base-models that are used in combinations as illustrated in the colorful grid. The chosen stacking architecture is used as input to a meta-classifier in the 2<sup>nd</sup> stacking level.

## 5 Results: per-domain evaluation



**Fig 5: Performance evaluation on ligand-binding domains.** The table represents domain-ligand pairs with performance exceeding that of the random baselines (the dashed lines) of AUC=0.5 and AUPRC fold ratio=1.

## 7 Significance

**Systematic identification of ligand-binding residues** would have a farther-reaching applications:

- Identify the functional impact of coding variants
- Explore the variation in protein interaction network
- Characterize mutations' effect in the context of disease
- Suggest molecular targets for therapeutic intervention

### 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] Kobren and Singh, NAR (2018)

### Acknowledgements

We would like to thank Shilpa N. Kobren, and the Singh lab for their helpful insights. Funded by NIH GM076275 (MS), T32 HG003284 (AE), and CA208148 (MS).