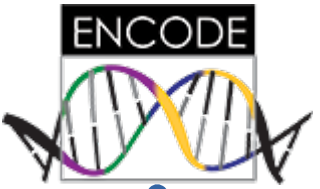
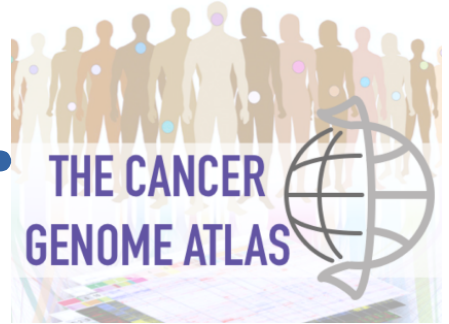
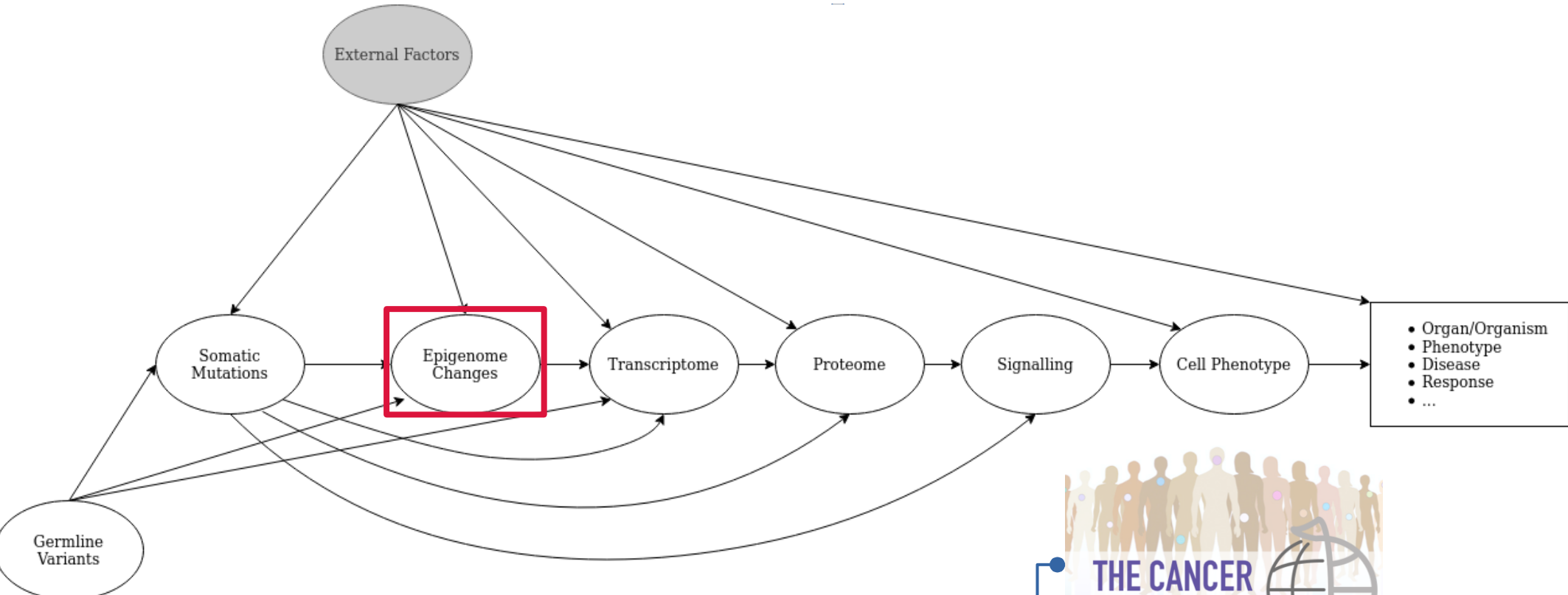


Getting personal with epigenetics: Towards machine-learning-assisted precision epigenomics

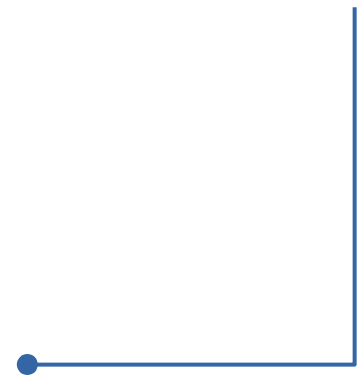
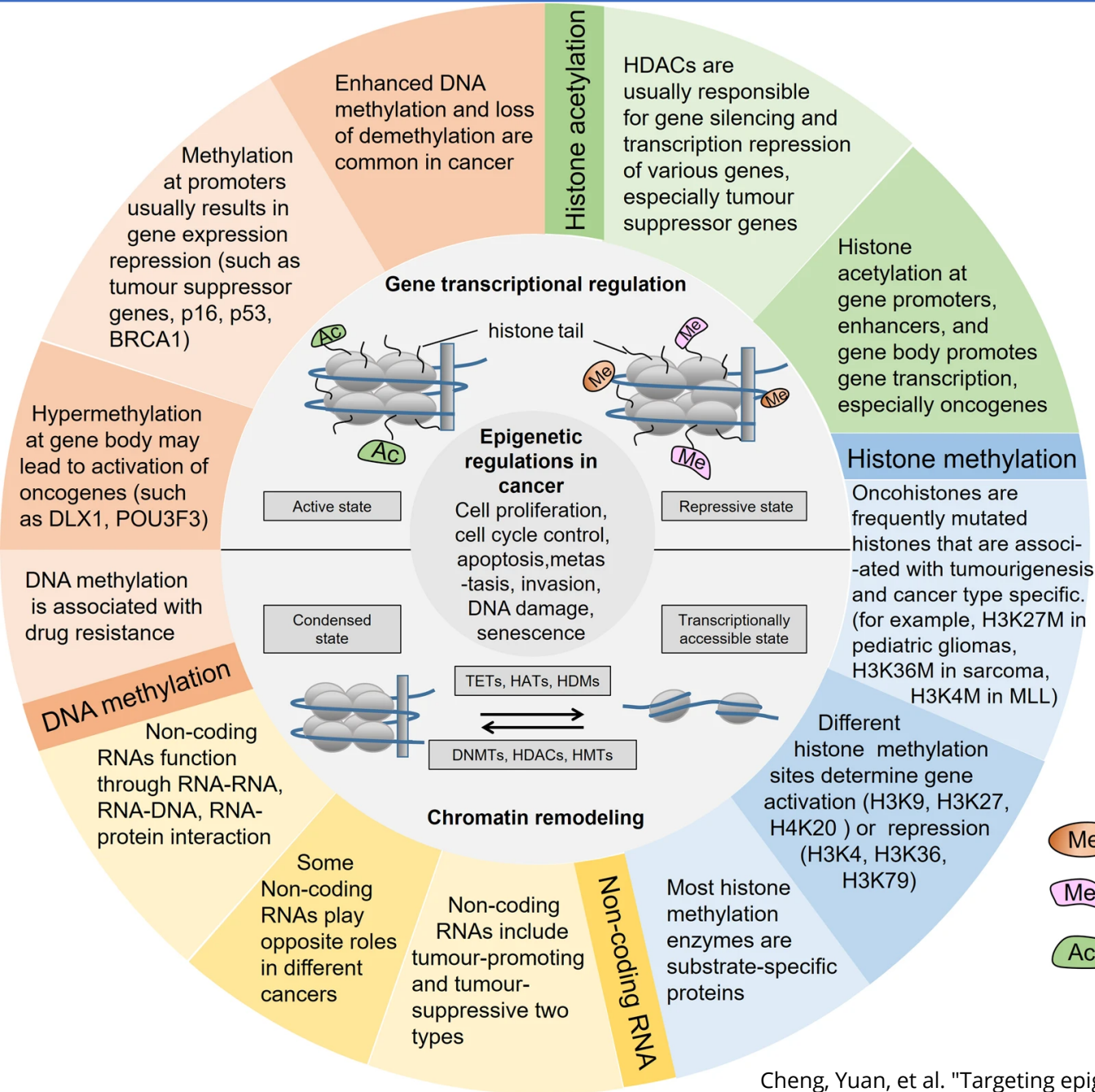
Giovanni Visonà


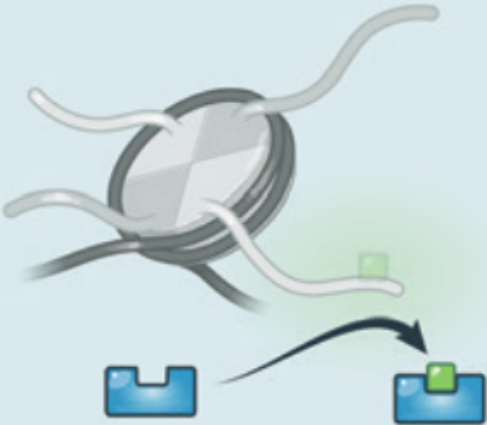


Large scale efforts conducted by several consortia have been key to the study of cancer and cell biology in general

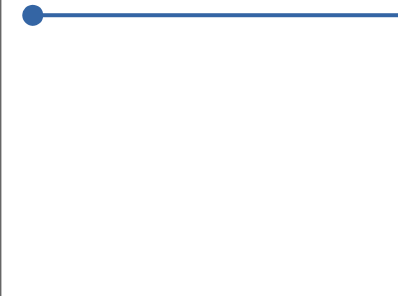


Epigenetic modifications are **involved in many aspects of carcinogenesis and cancer progression**, which offer **compelling therapeutic targets**



Category	Epigenetic Regulators	Function	FDA-Approved Drug
Writers 	DNMT1, 3A, and 3B	Methylates cytosines on DNA, and mutation can lead to aberrant methylation	Azacitidine, decitabine
	EZH2	Methylates histone H3K27	Tazemetostat
	DOT1L	Methylates histone H3K79	
	KMT2A–D, SETD2, NSD1	Methylate histone lysines	
	EP300, CREBBP	Acetylate histone lysines	
Erasers 	TET2	Is the first step in cytosine demethylation; is inhibited by 2-hydroxyglutarate (2-HG)	Azacitidine, decitabine
	IDH1, IDH2	Mutated protein produces 2-HG from isocitrate that inhibits TET2 and lysine demethylases	Ivosidenib, enasidenib
	HDAC1–3, 8 HDAC6	Deacetylase removes acetyl groups from histone lysines	Vorinostat, belinostat, panobinostat, romidepsin
	KDM1A, KDM6A (UTX)	Demethylates histone lysines	

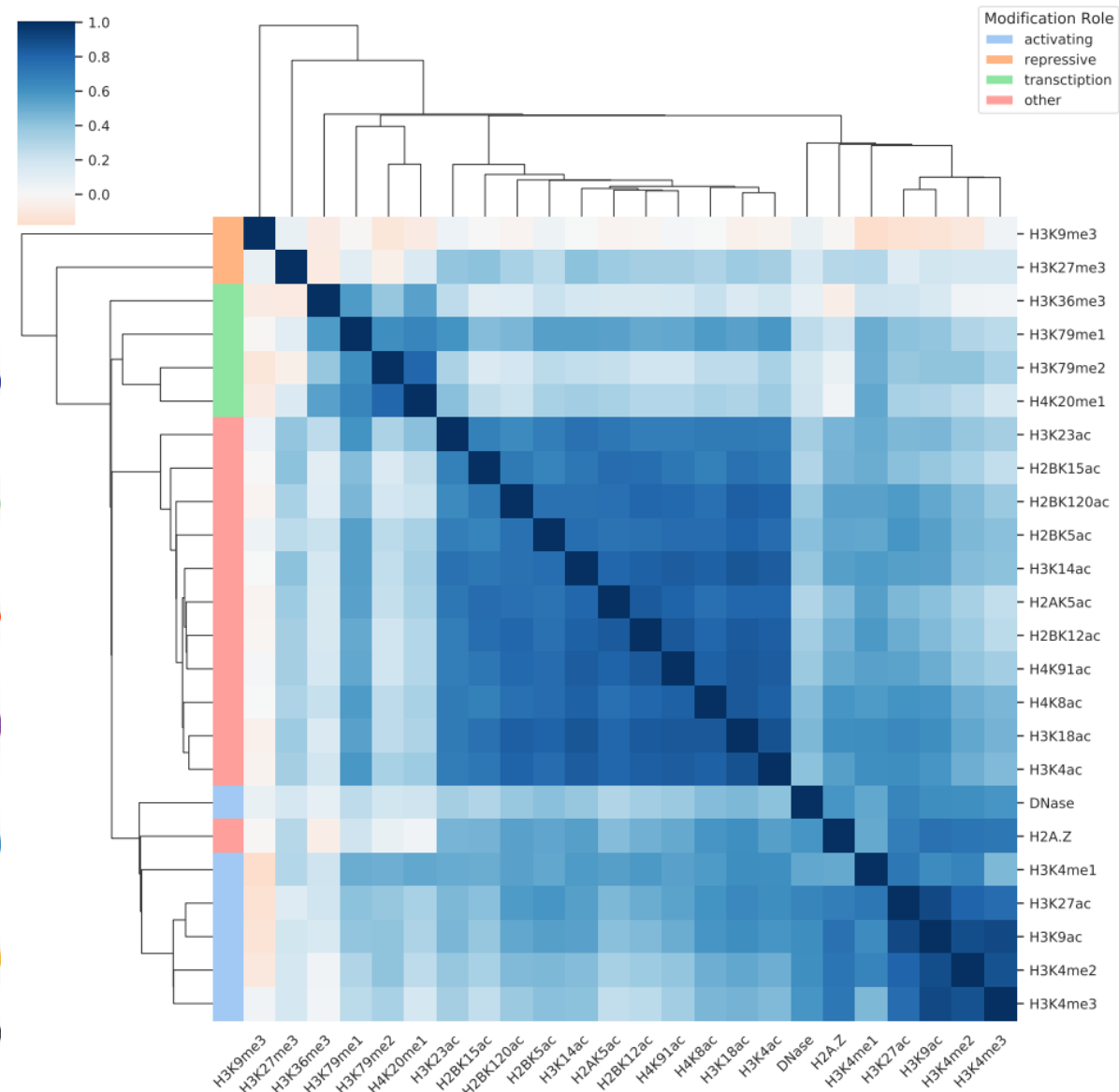
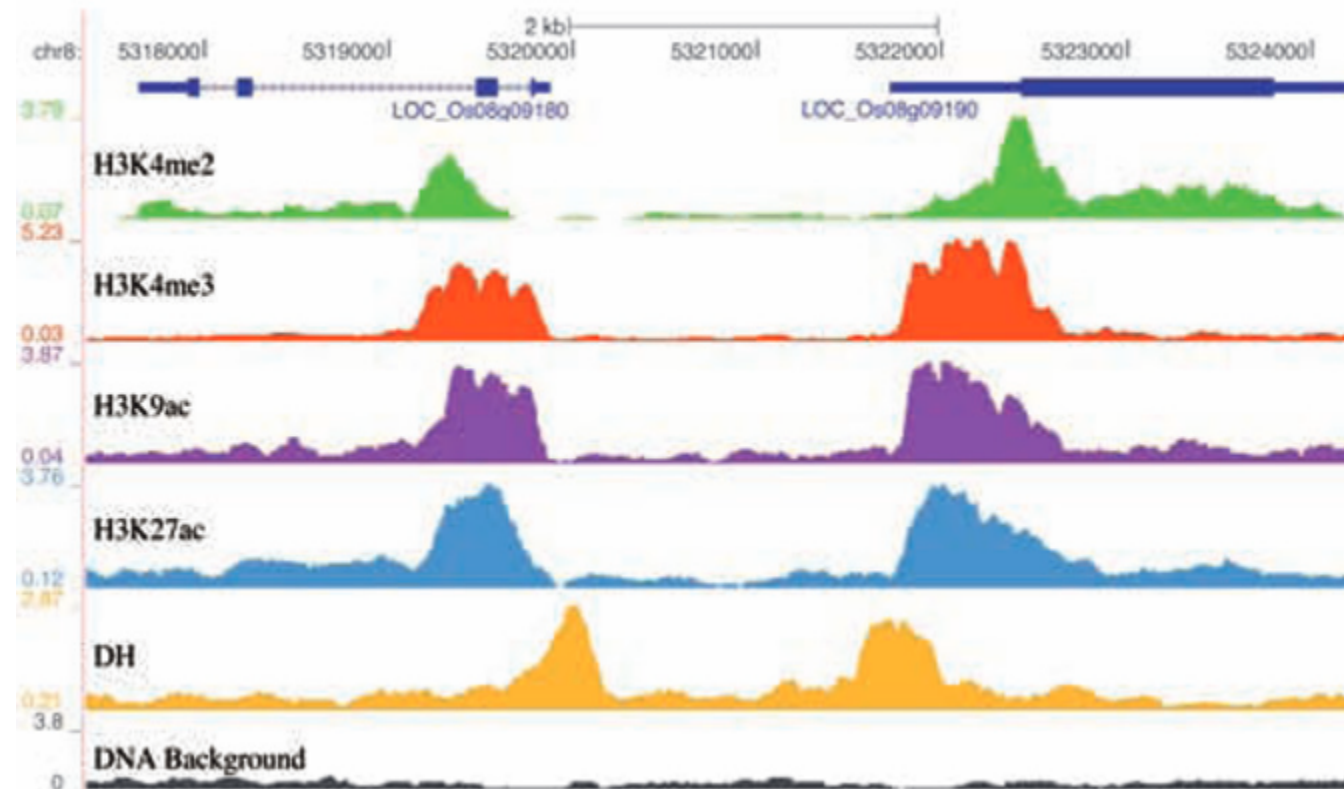
	Bisulfite-Seq	MeDIP-Seq	MRE-Seq	RRBS	DNaseI	DGF	mRNA-Seq	smRNA-Seq	ChIP-input	H3K27me3	H3K36me3	H3K4me1	H3K4me3	H3K9ac	H3K9me3	H3K27ac	H2AK5ac	H2AK9ac	H2AZ	H2BK120ac	H2BK12ac	H2BK15ac	H2BK20ac	H3K14ac	H3K18ac	H3K23ac	H3K4ac	H3K4me2	H3K56ac	H3K79me1	H3K79me2	H4K20me1	H4K5ac	H4K8ac	H4K12ac	H4K91ac	H3K23me2	H2BK5ac	H3K9me1	H3T11ph								
ADRENAL-Fetal																																																
BRAIN																																																
BRAIN-Fetal																																																
BREAST																																																
ES CELLS																																																
ES-derived cells																																																
Exocrine-Endocrine																																																
FAT-Adult																																																
GI-Adult																																																
GI-Fetal																																																
GU-Adult																																																
HEART-Adult																																																
HEART-Fetal																																																
Hematopoietic Stem																																																
CD34, Primary Cells																																																
CD34, Mobilized Primary Cells																																																
CD34, Cultured Cells																																																
iPS CELLS																																																
KIDNEY-Fetal																																																
LUNG-Adult																																																
LUNG-Fetal																																																
MUSCLE-Adult																																																
MUSCLE-Fetal																																																
PLACENTA-Fetal																																																
REPRODUCTIVE-Adult																																																
REPRODUCTIVE-Fetal																																																
SKIN-Fetal																																																
SPLEEN-Fetal																																																
STROMAL-CONNECTIVE																																																
THYMUS-Fetal																																																
White Blood																																																



Can we leverage machine learning to **fill the gaps?**

Epigenetic modifications are **cell-type specific**

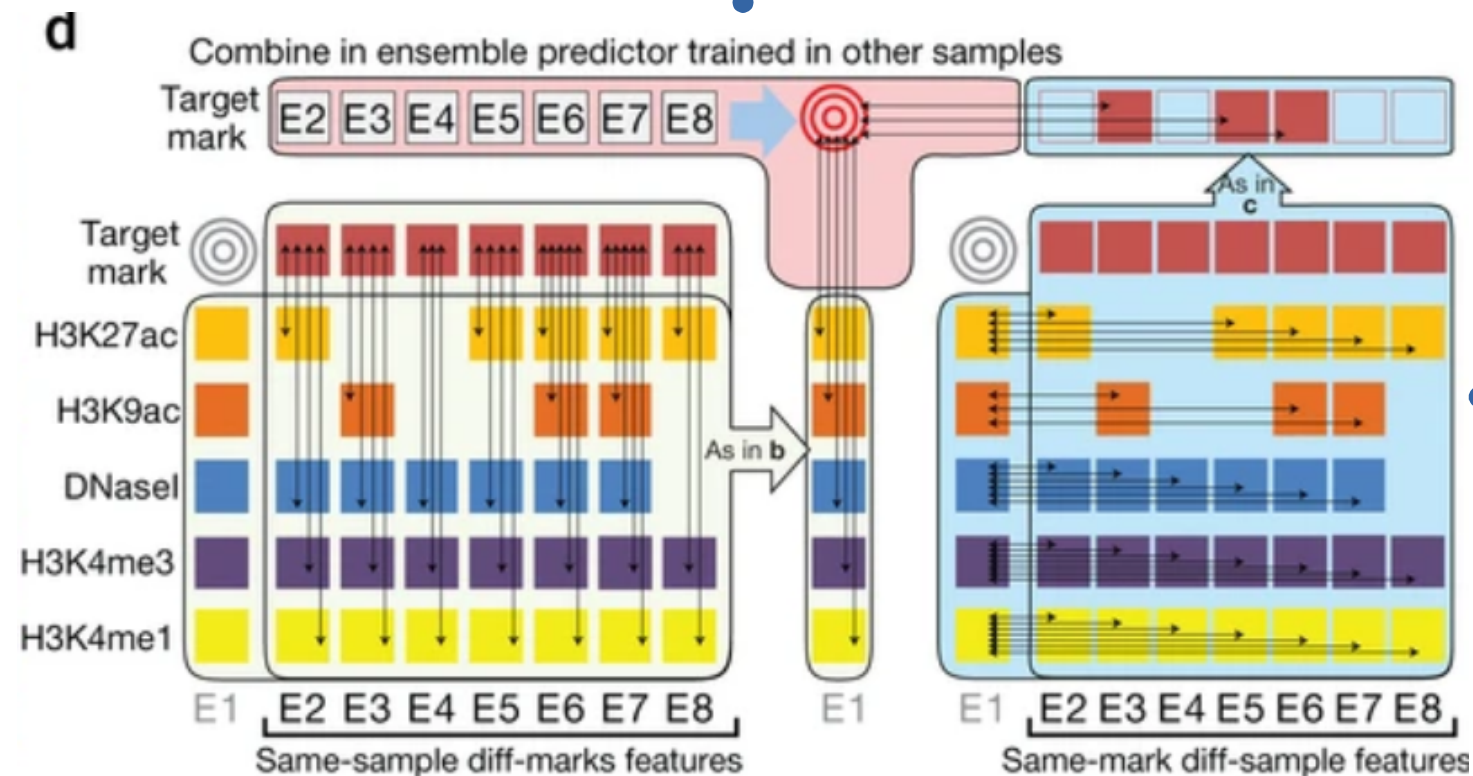
Epigenetic marks are highly **correlated**



Previous work: ChromImpute

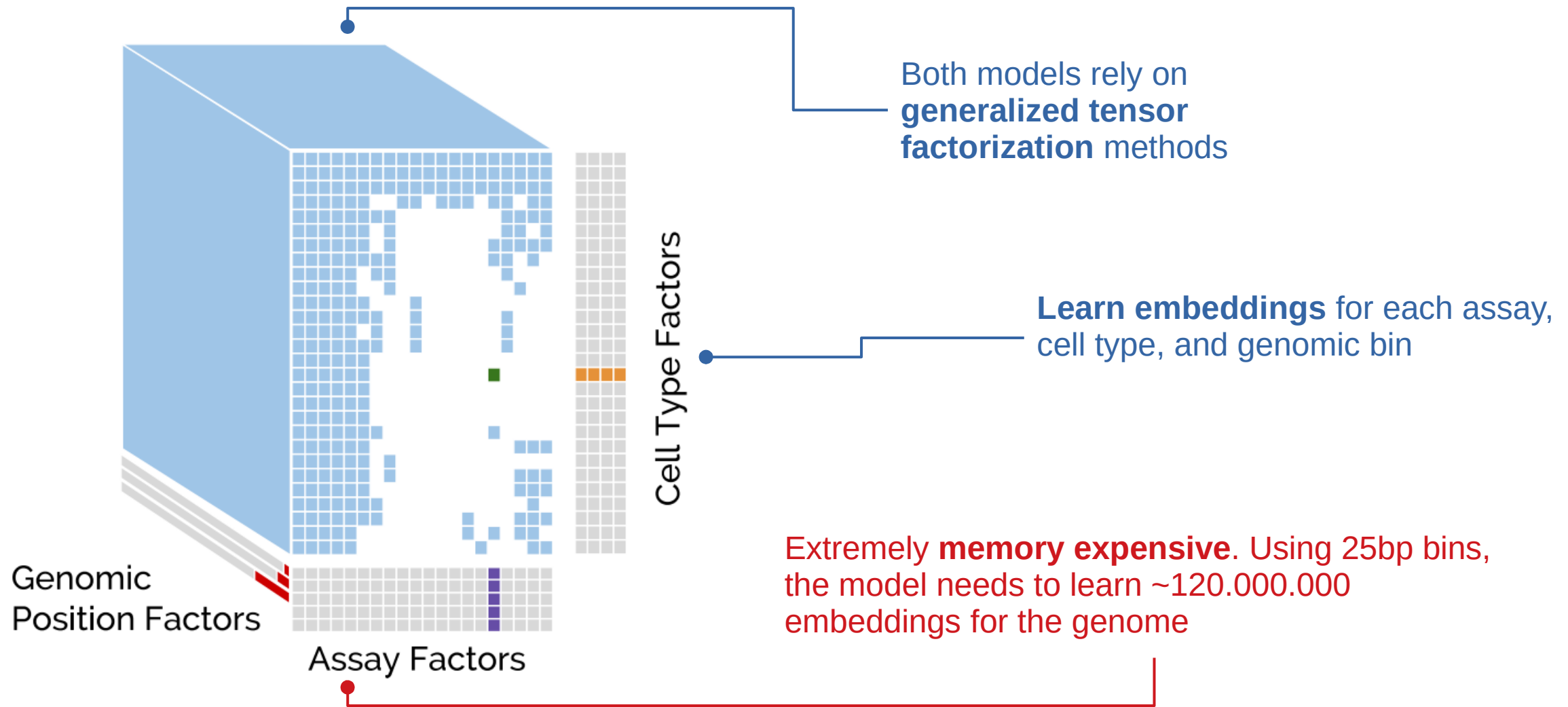
Trains ensembles of regression trees

Significant performance improvement compared to all previous methods



For each cell type-assay query CI requires the training of a **new ensemble** of regression trees

Previous work: PREDICTD and Avocado



Durham, Timothy J., et al. "PREDICTD parallel epigenomics data imputation with cloud-based tensor decomposition." Nature communications 9.1 (2018)

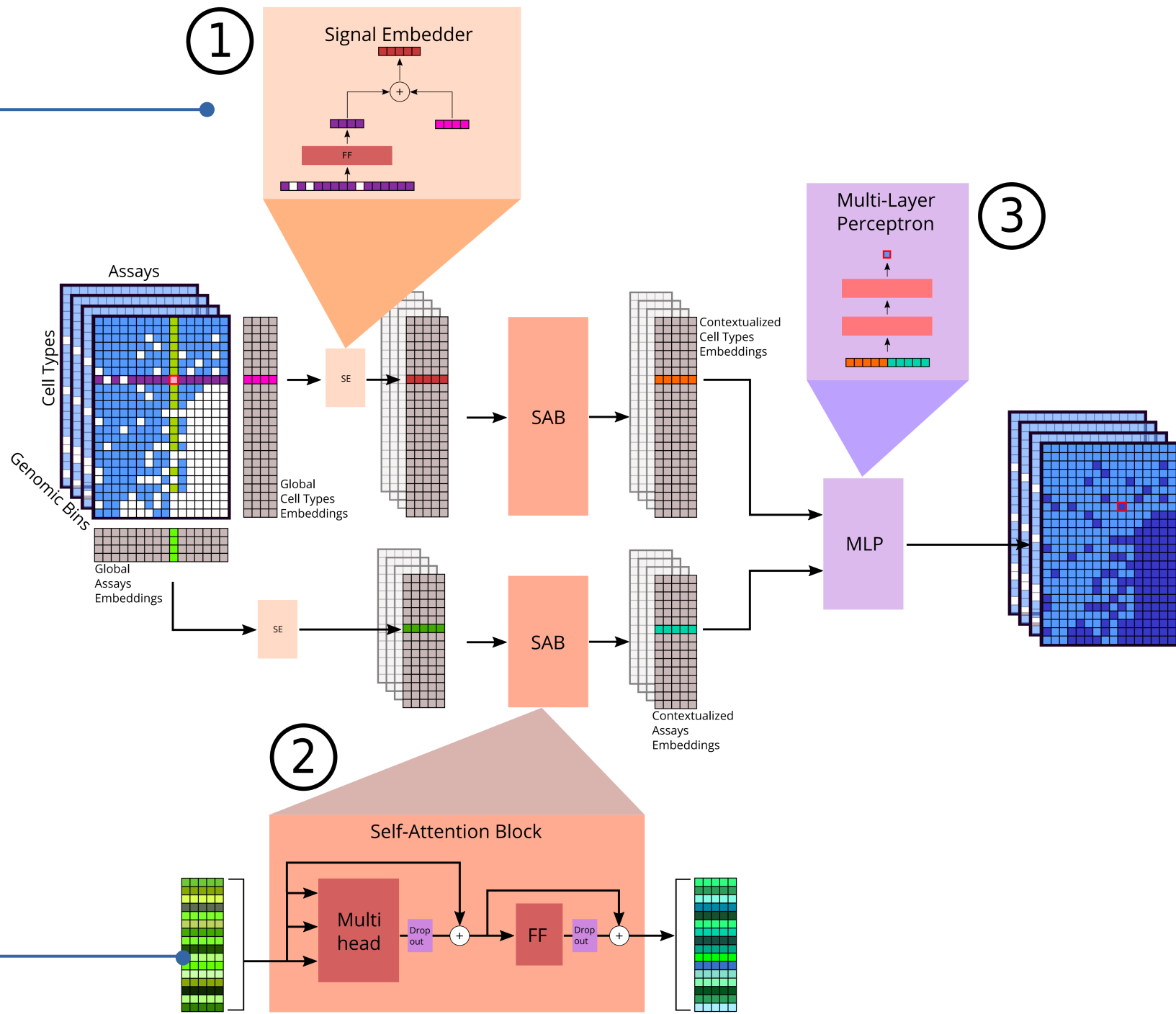
Schreiber, Jacob, et al. "Avocado: a multi-scale deep tensor factorization method learns a latent representation of the human epigenome." Genome biology 21 (2020)

eDICE

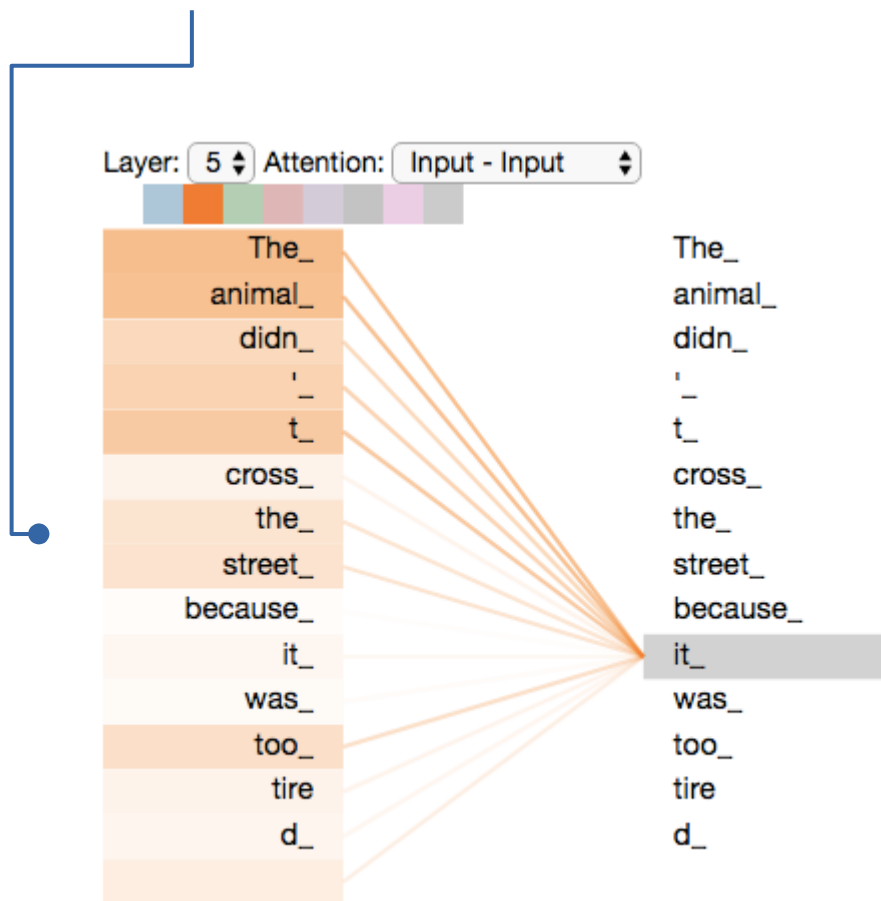
epigenomic Data Imputation through Contextualized Embeddings

Local signal captures the information of each genomic bin

Self-Attention includes the interaction between assays and between cell types in the model

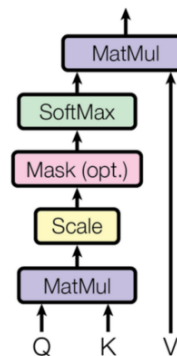


What is attention?



<https://jalammar.github.io/illustrated-transformer/>

Scaled Dot-Product Attention



Multi-Head Attention

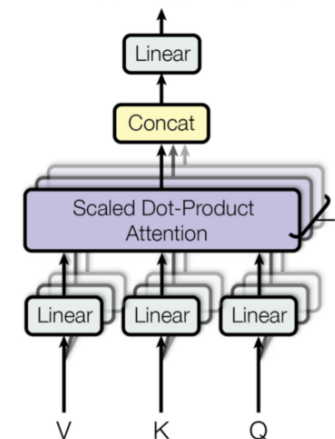
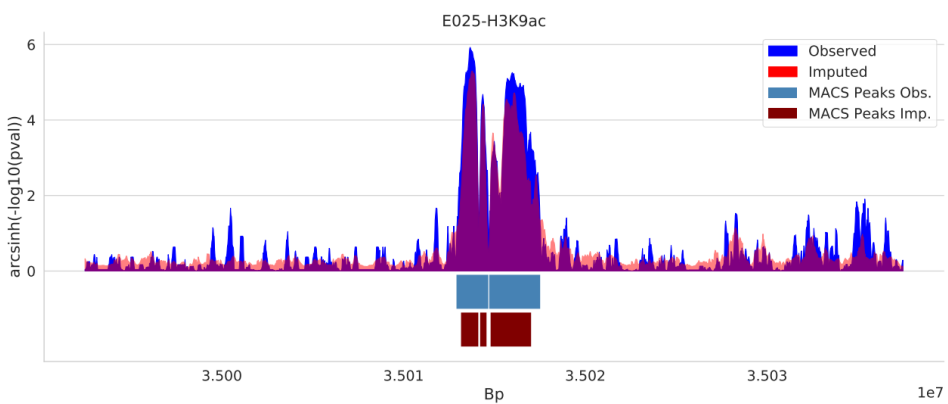
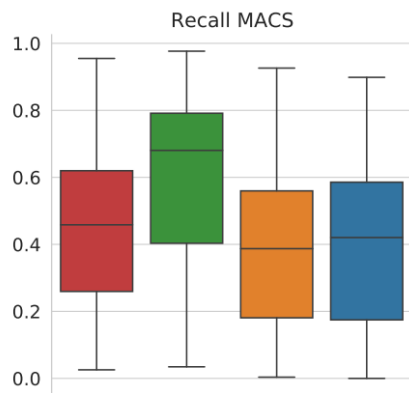
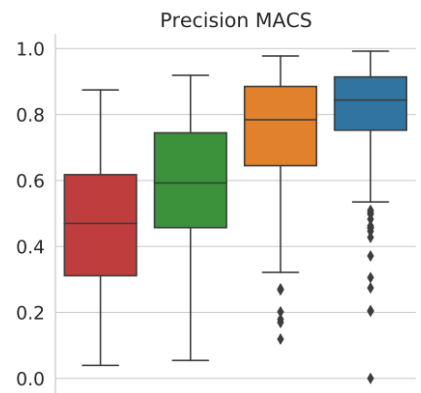
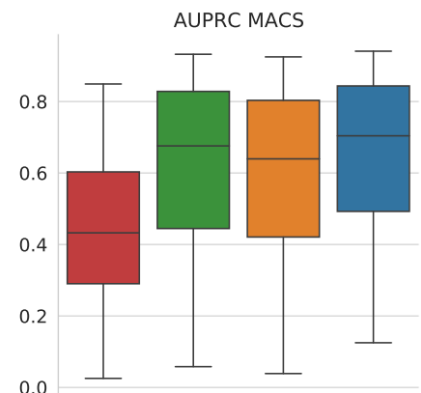
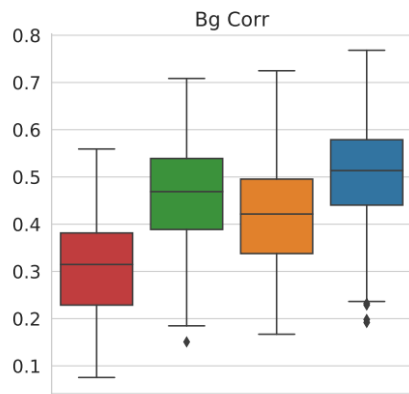
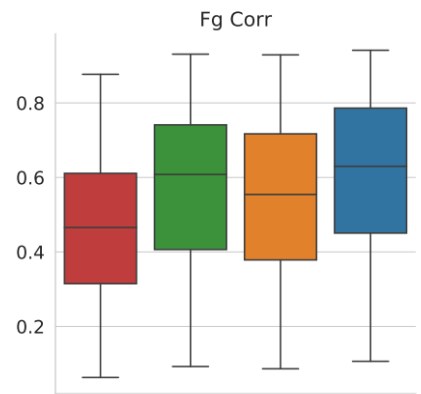
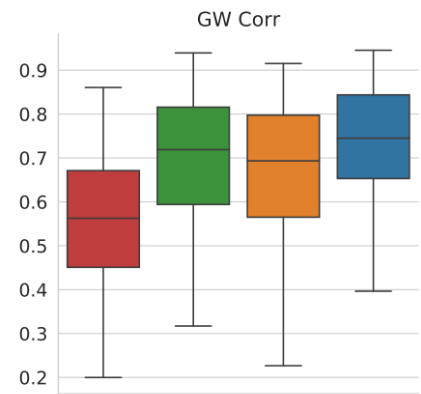
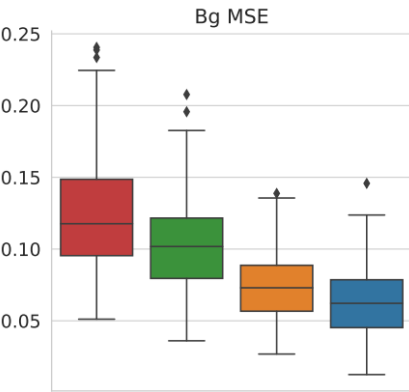
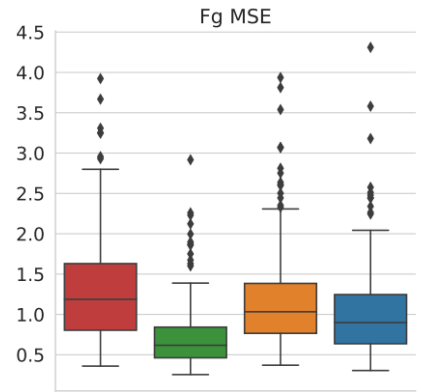
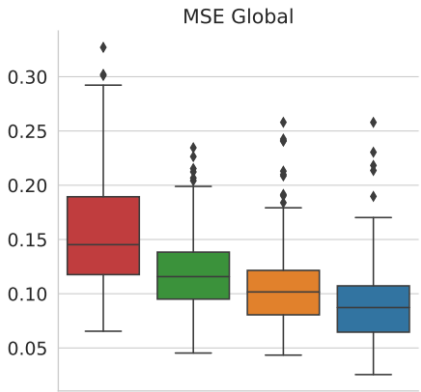


Figure 2: (left) Scaled Dot-Product Attention. (right) Multi-Head Attention consists of several attention layers running in parallel.

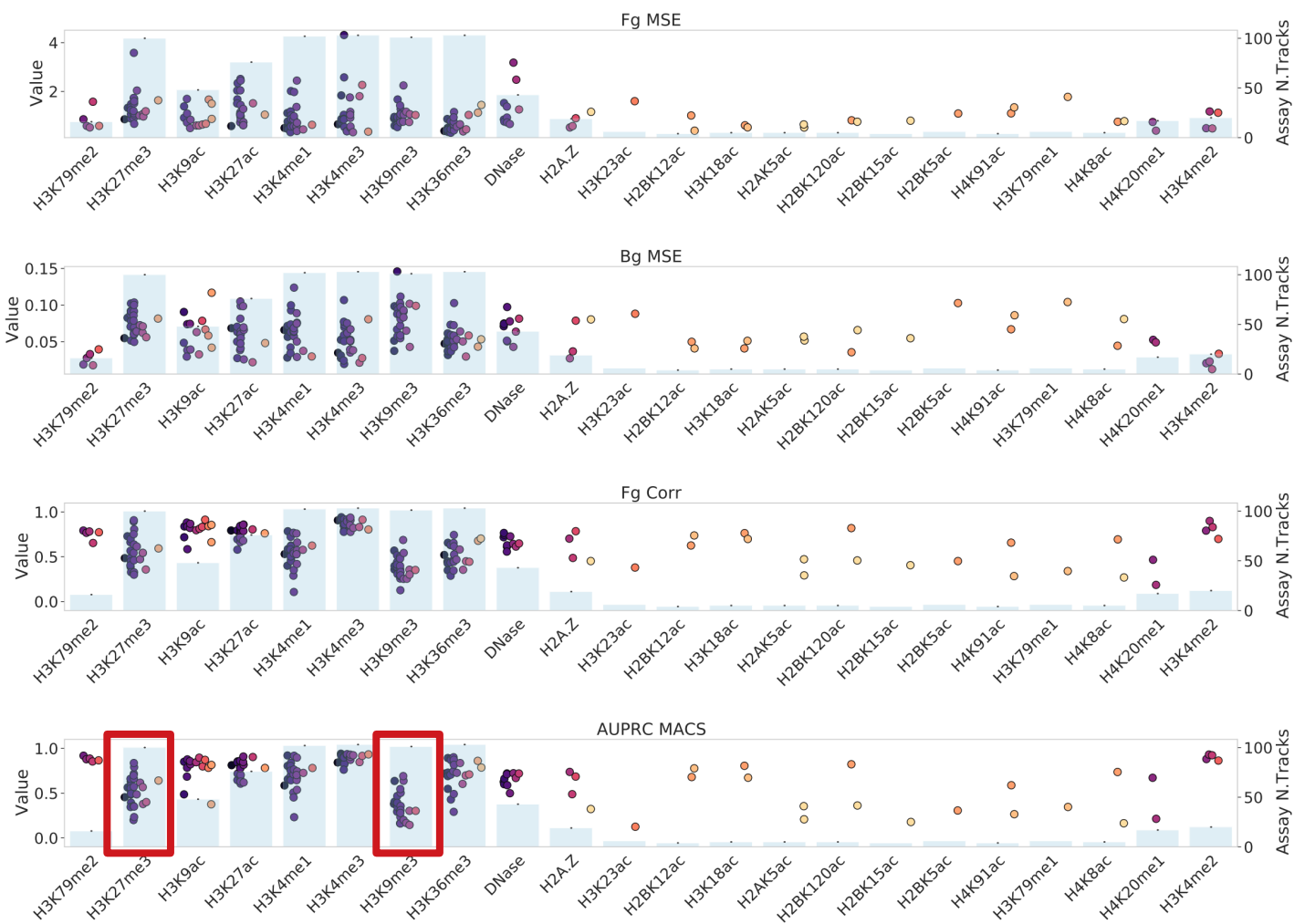
The **Transformer** is one of the most successful attention models

Tested on the Roadmap dataset
(chromosome 21),
**eDICE outperforms ChromImpute
and PREDICTD on almost all metrics**

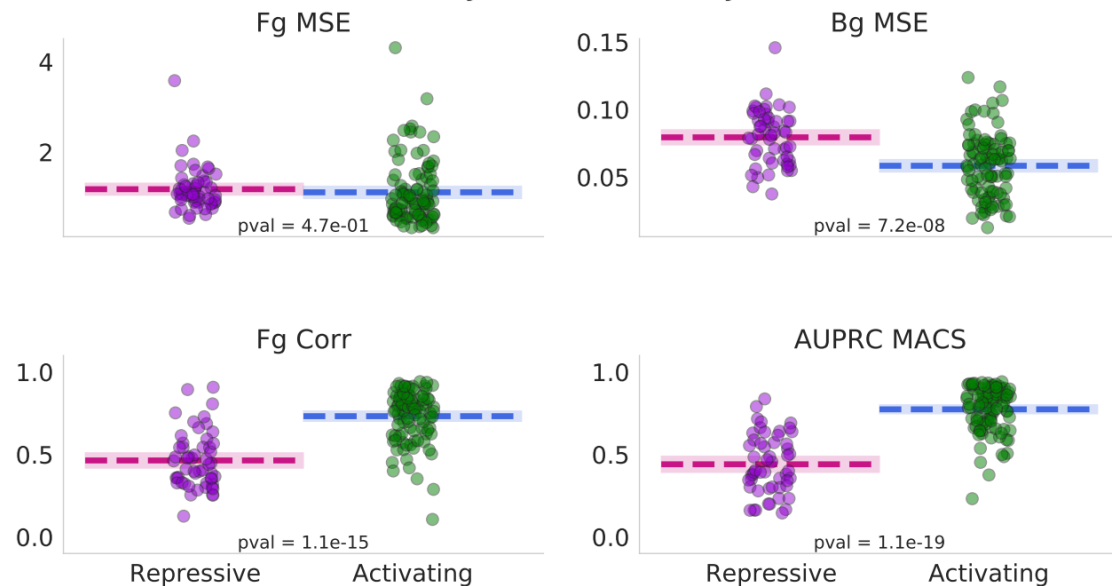


Imputations work better on some assays.
This problem is not unique to eDICE.

eDICE - Assay-level Performance

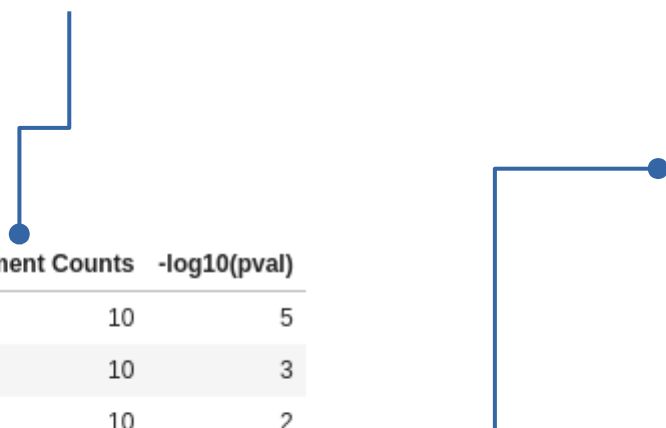


eDICE - Assay Performance by Function

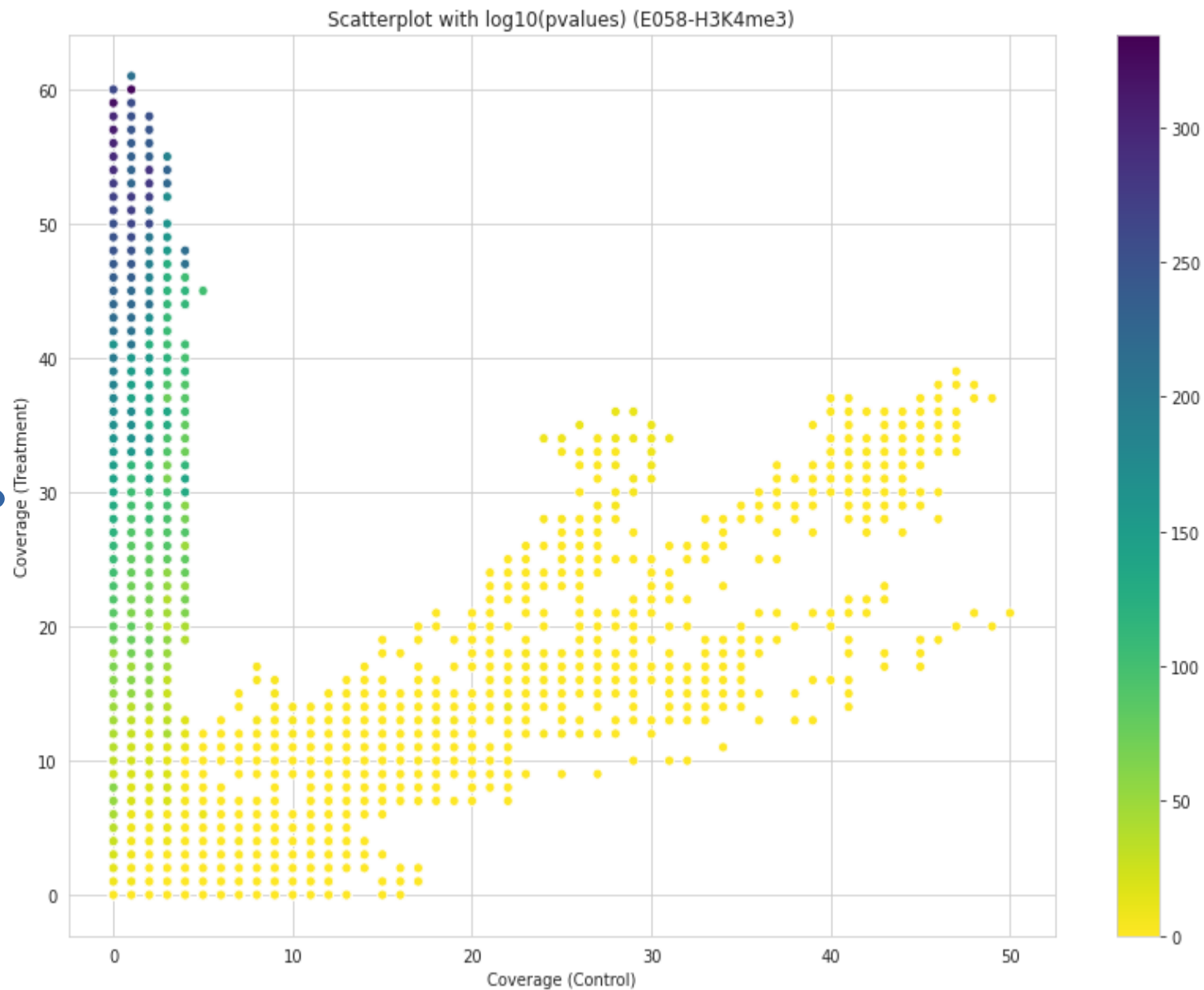


The discrepancies are possibly due to biases in the sequencing of heterochromatin-associated marks.

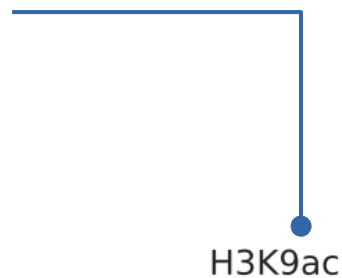
The processing pipeline is affected by **multiple sources of bias**: small number of replicates, low quality control samples, different sequencing platforms.



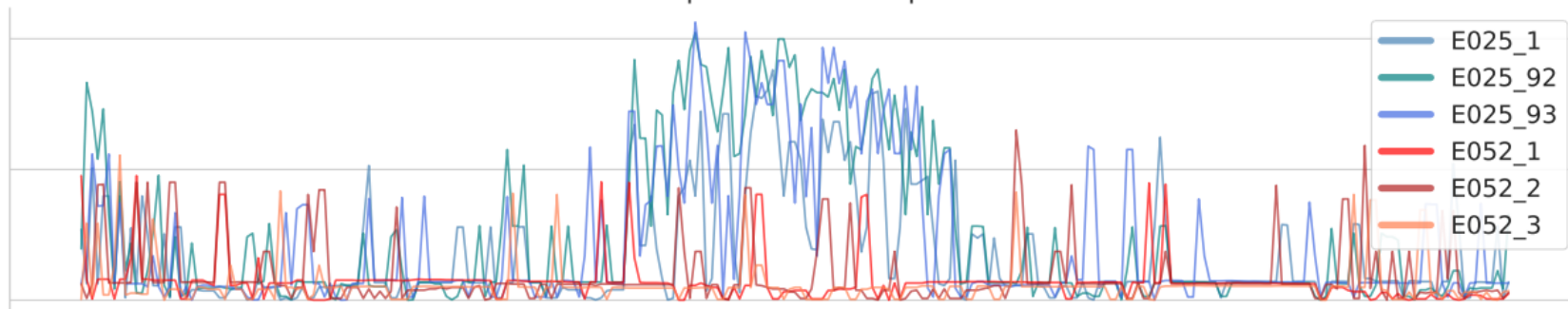
Control Counts	Treatment Counts	$-\log_{10}(\text{pval})$
2	10	5
3	10	3
4	10	2
8	100	72
9	100	67
10	100	63



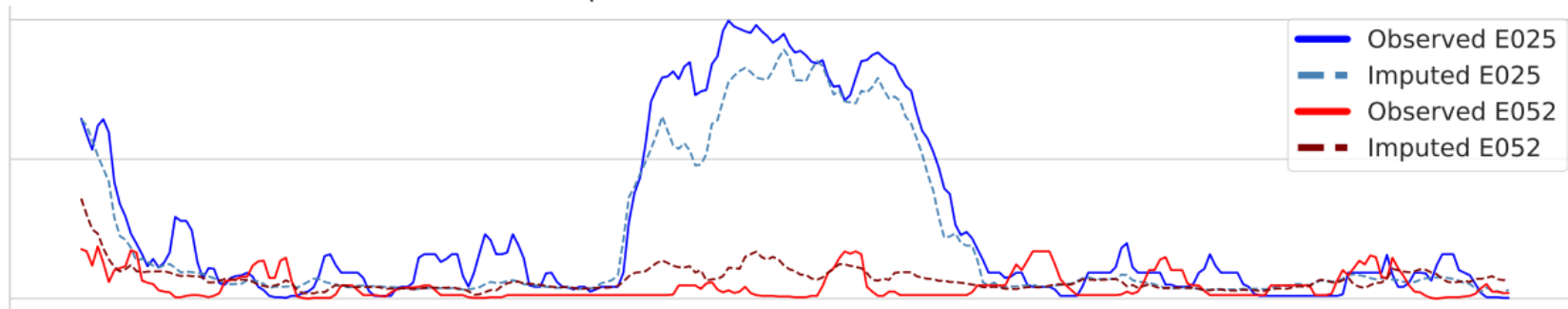
Identifying meaningful differences between biological samples is crucial to progress our understanding of the regulatory mechanisms of the genome



Tissue-specific Peak - Replicates

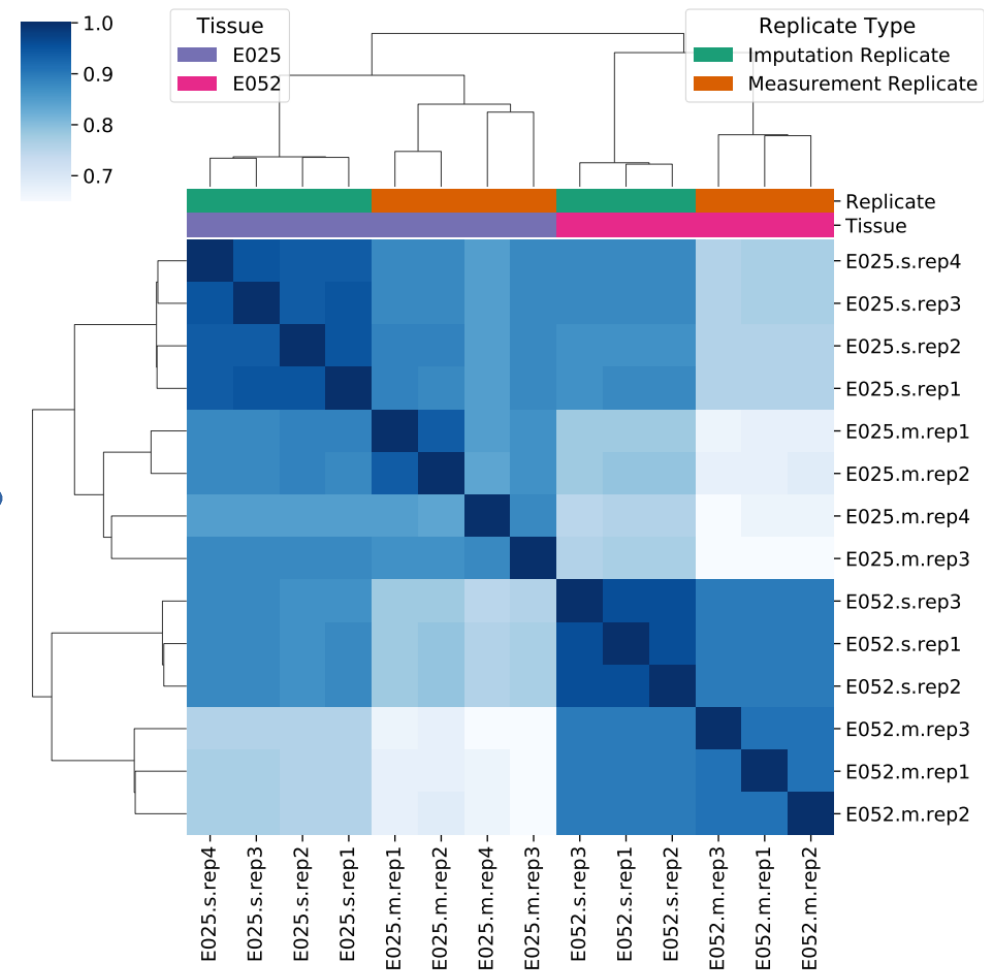
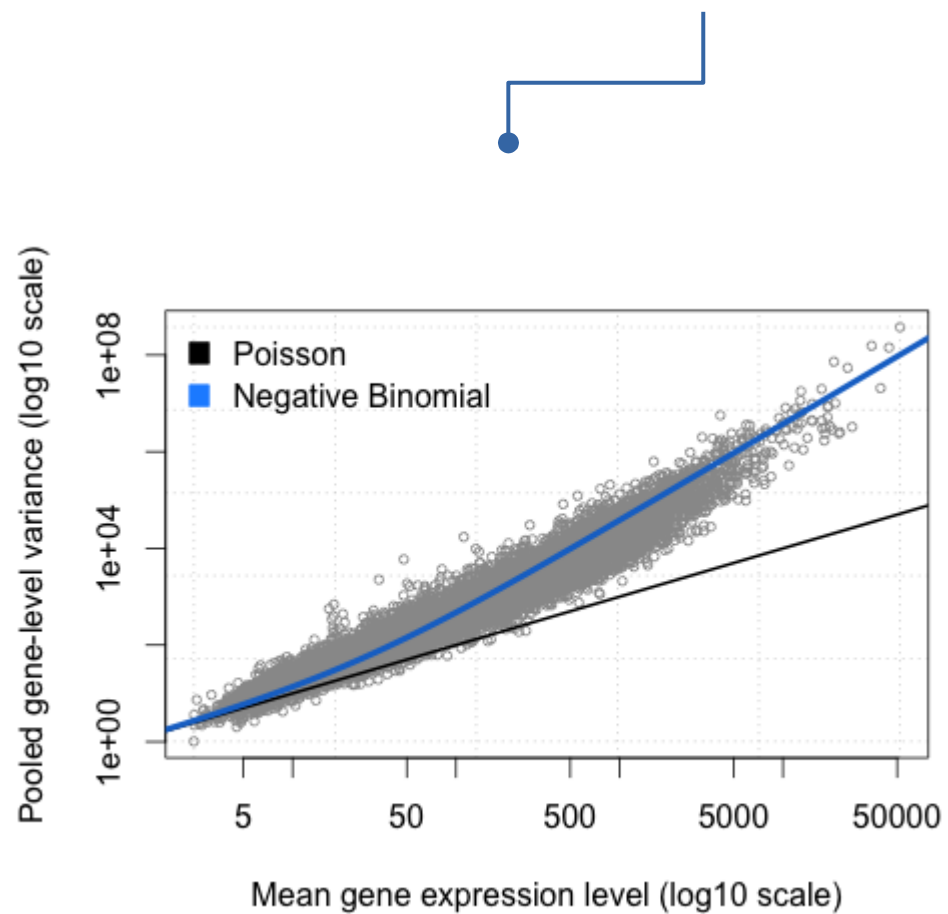


Tissue-specific Peak - Reference P-value tracks

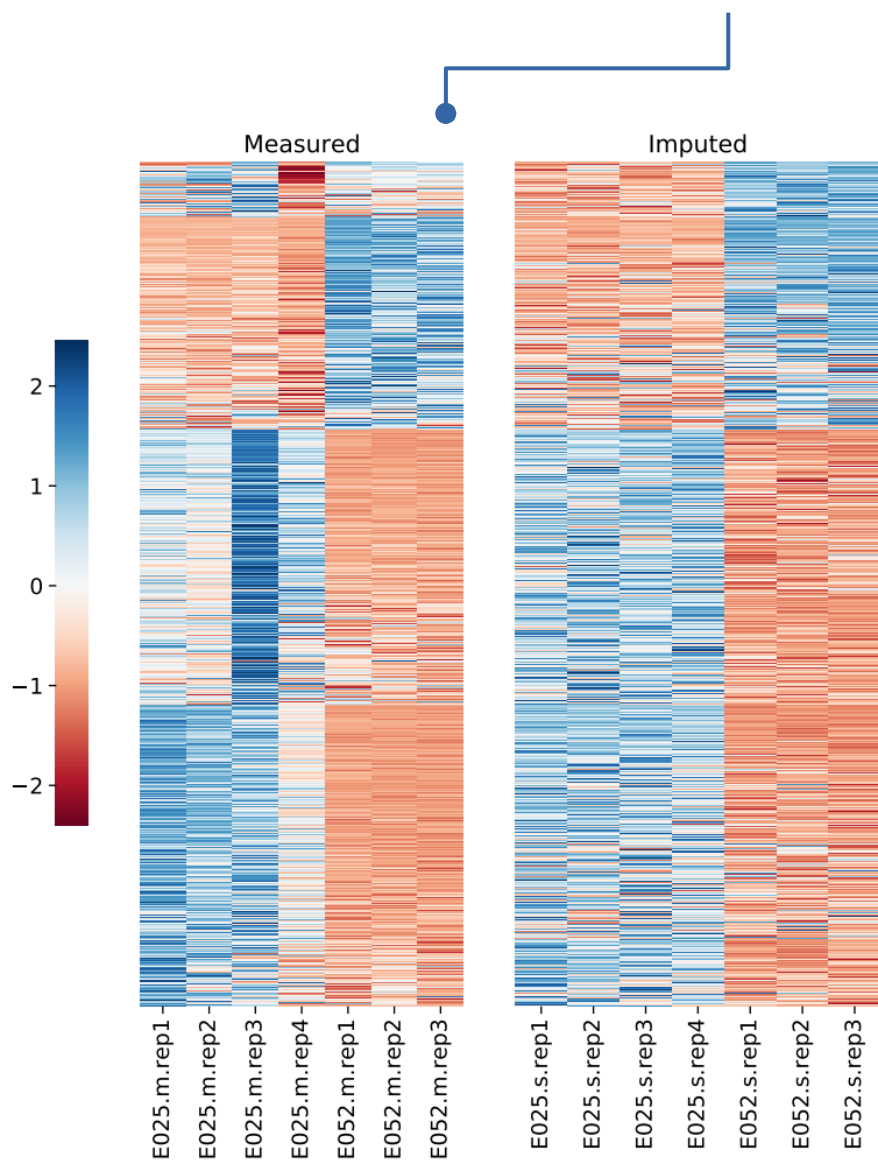


The use of multiple replicates is fundamental for robust analysis

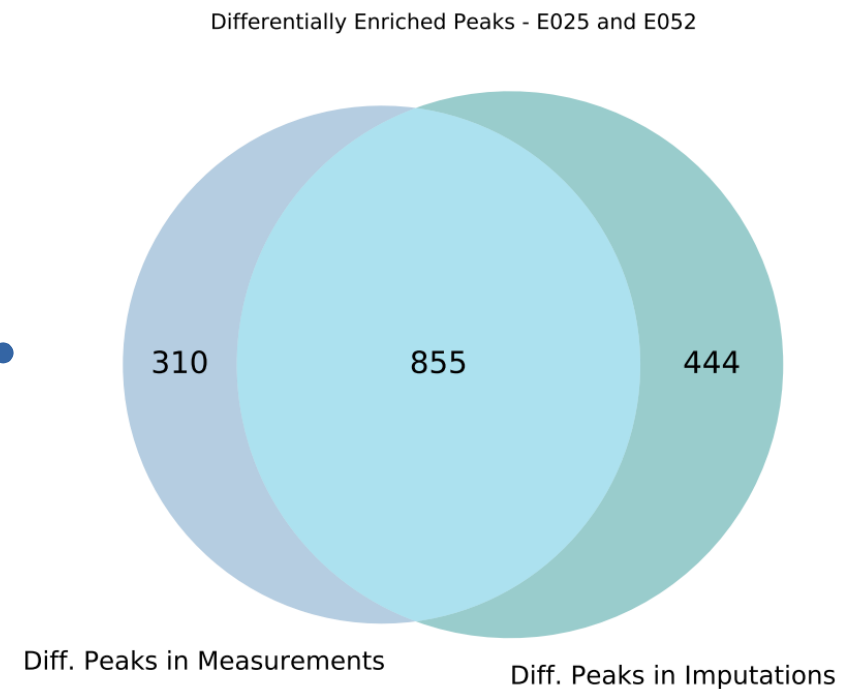
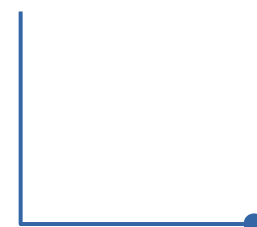
We simulate pseudo-replicates for two tissues using parameters estimated from the imputations. (Next iteration will explicitly predict the variance of the signal)



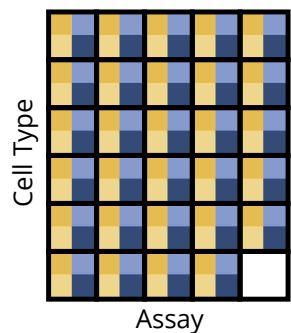
Binding affinity scores for the imputed replicates reflect the pattern found in the measurements.



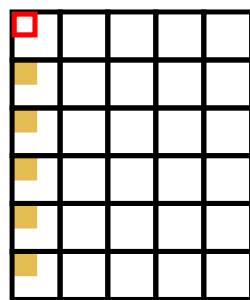
This differential analysis retrieves most of the meaningful differences between sets of replicates. (PPV ~70%)



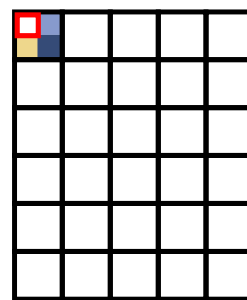
We compare to two baselines: averaging over tissues (AVG), and averaging over individuals (TrackAVG)



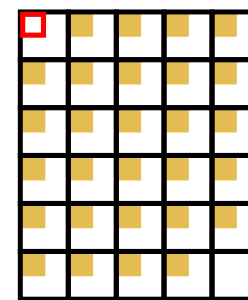
AVG



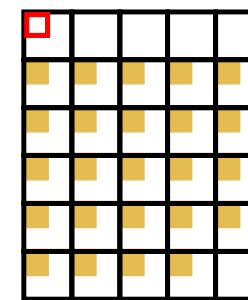
TrackAVG

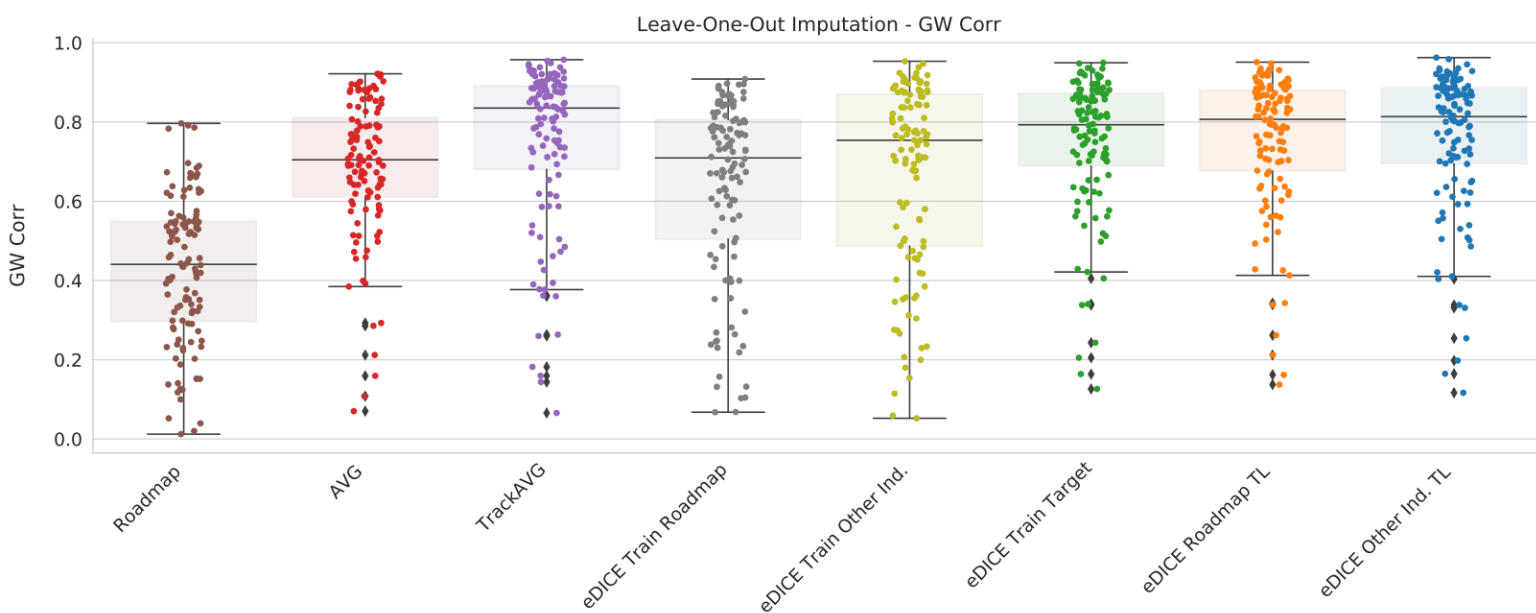
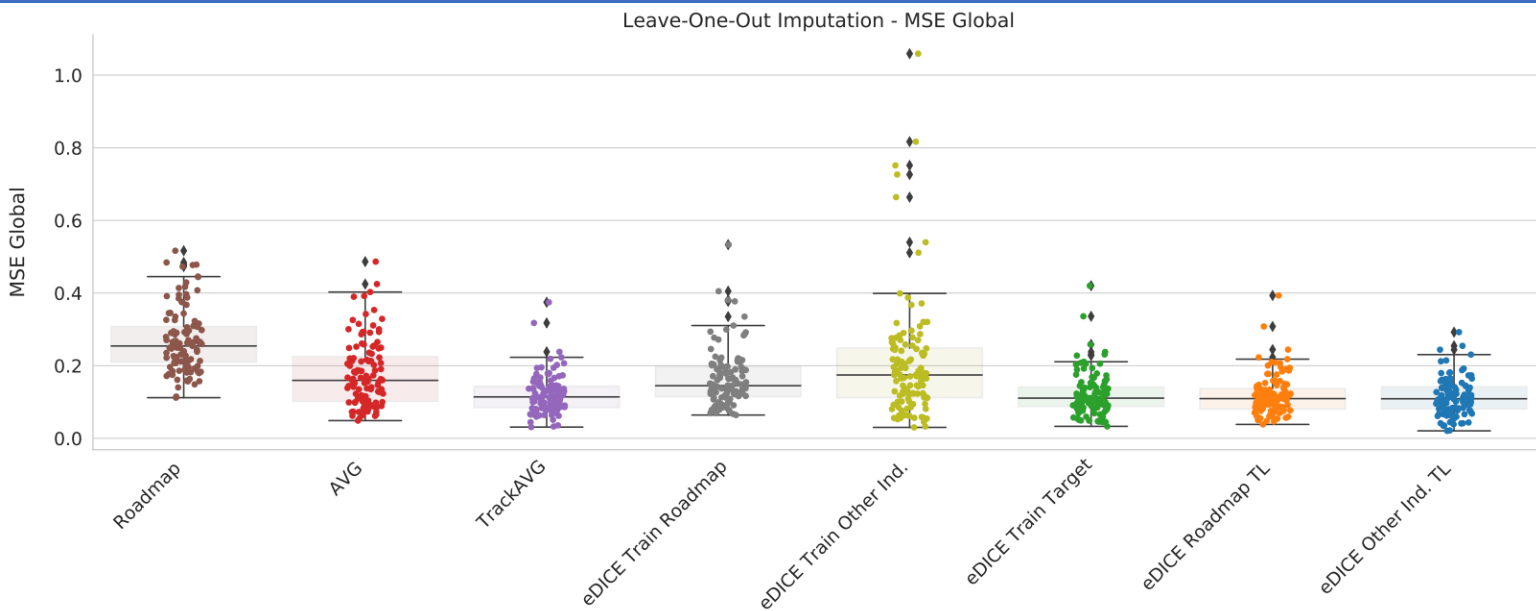


eDICE LOO support



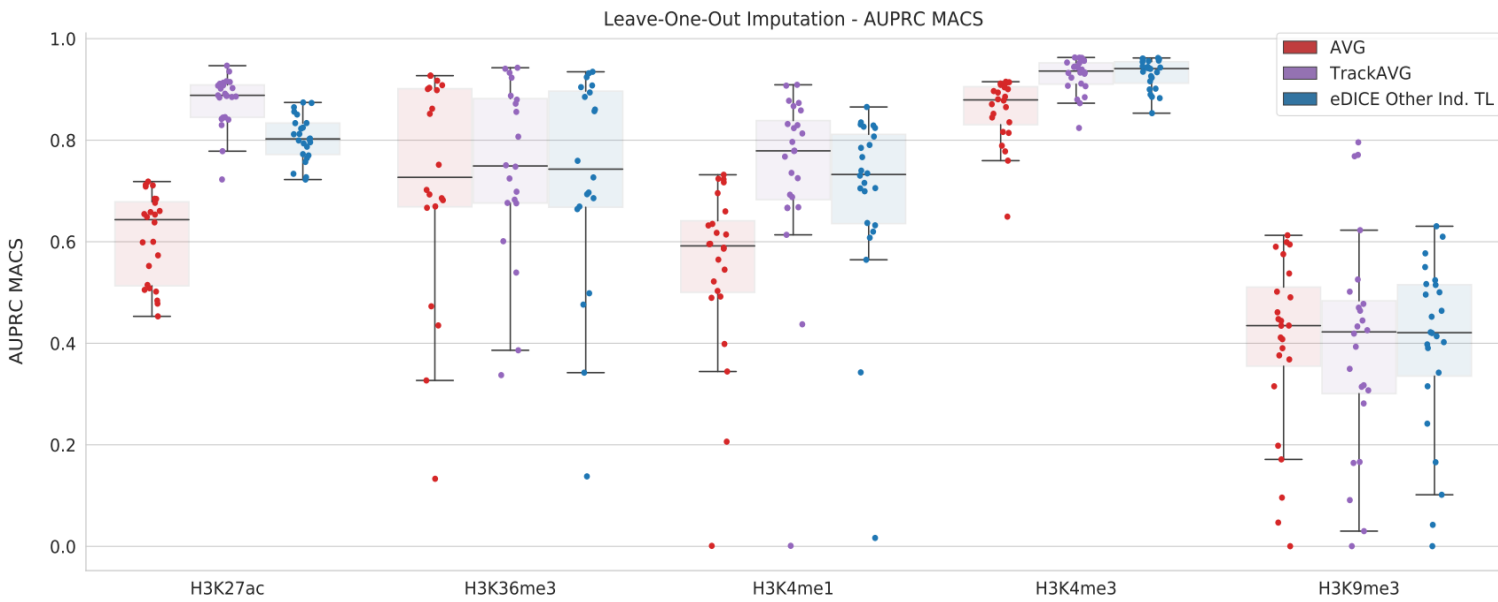
eDICE generalisation to unseen cell type



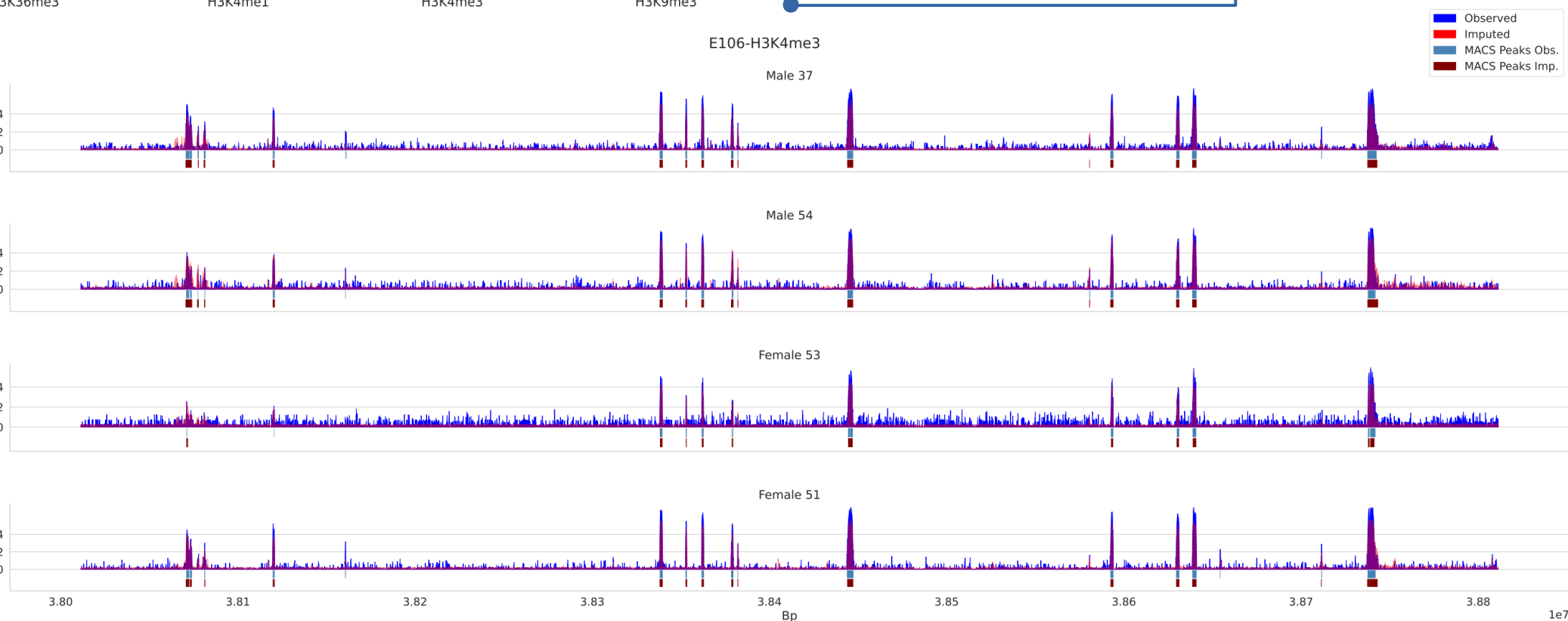


Among eDICE models, transfer learning from other individuals works best

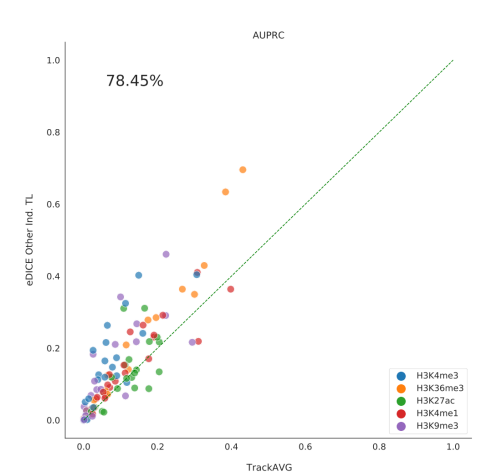
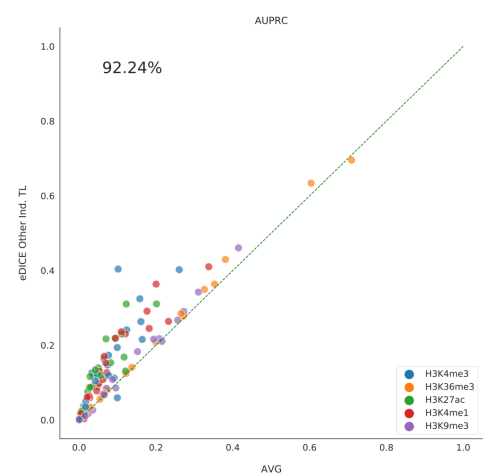
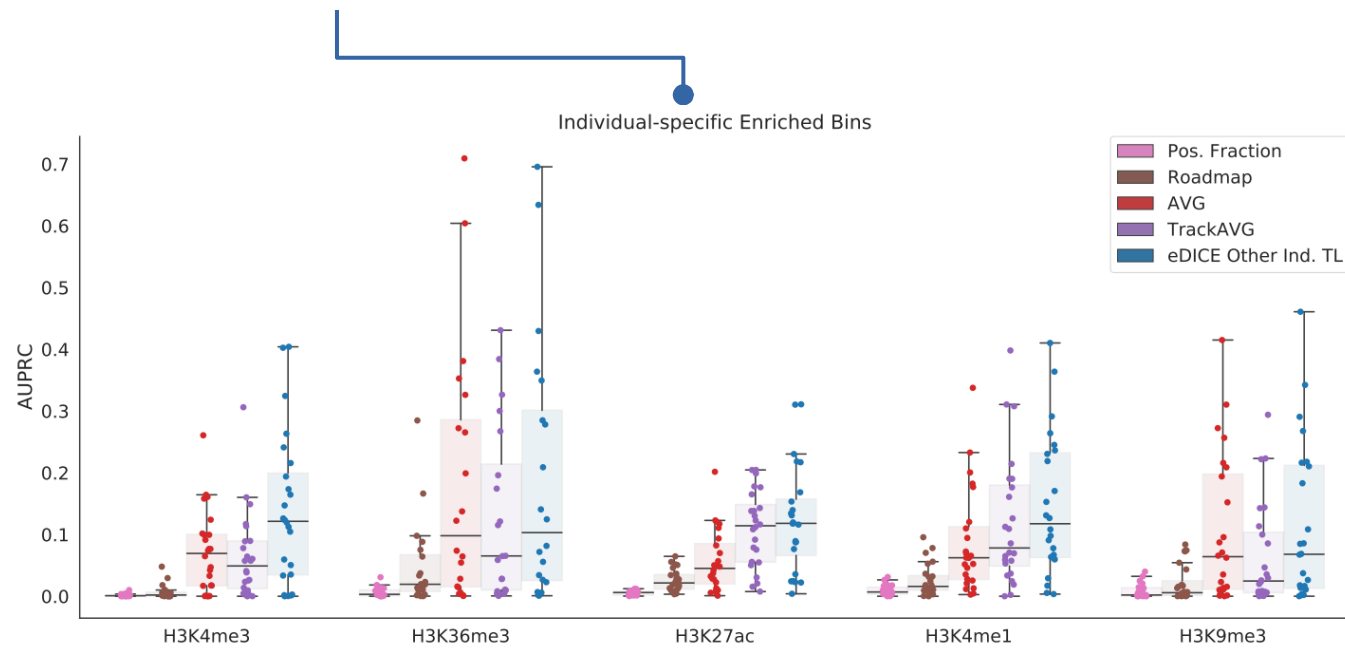
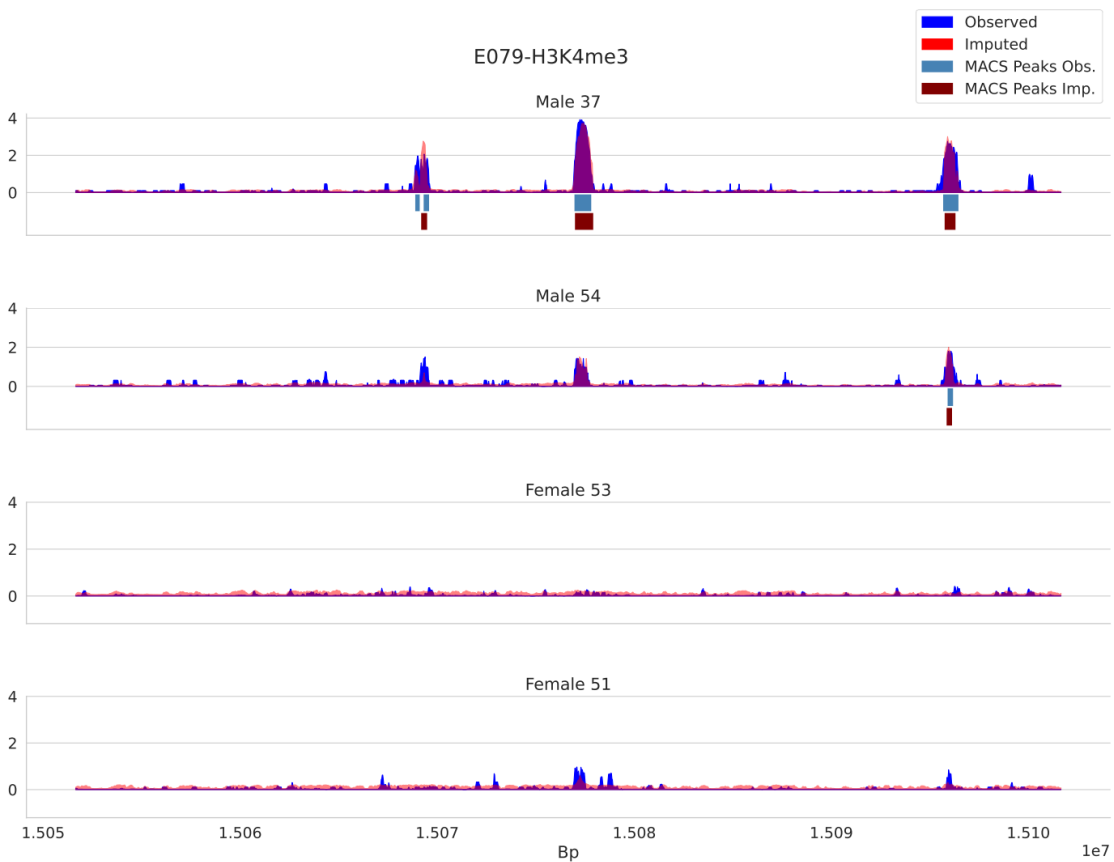
It also offers the lowest computational cost



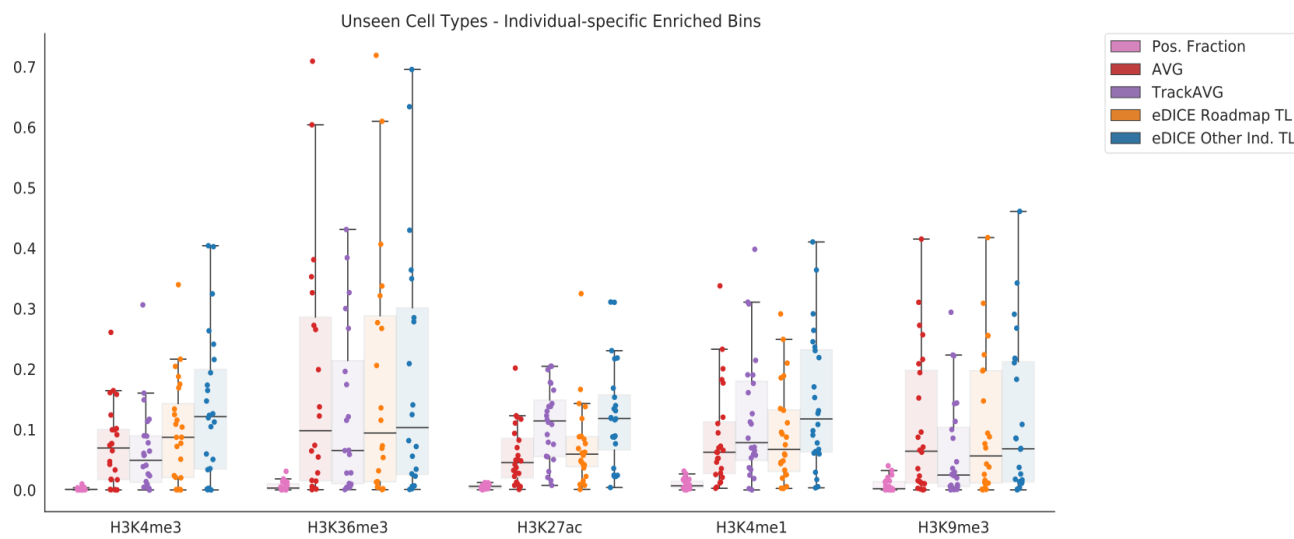
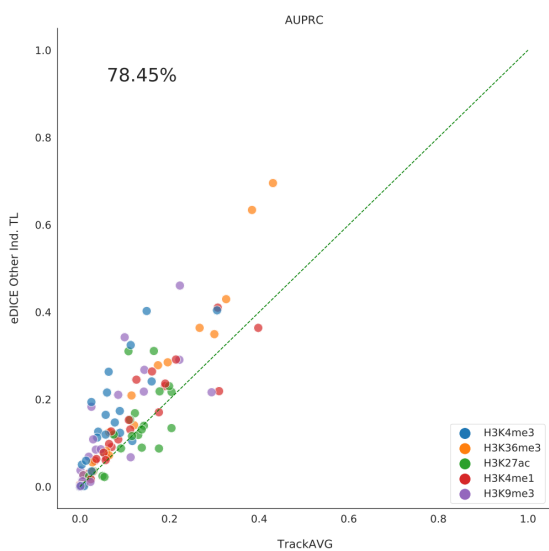
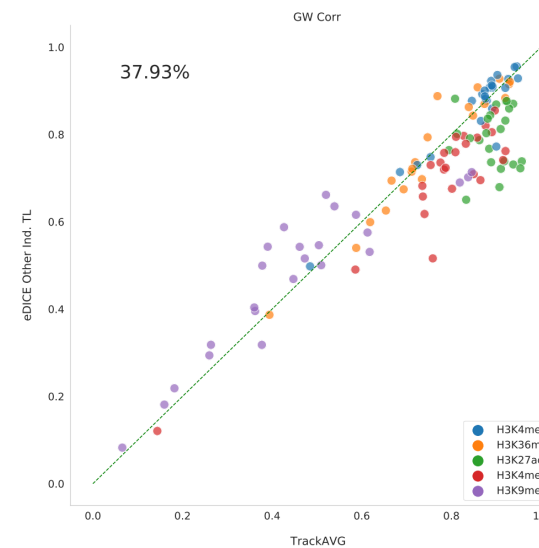
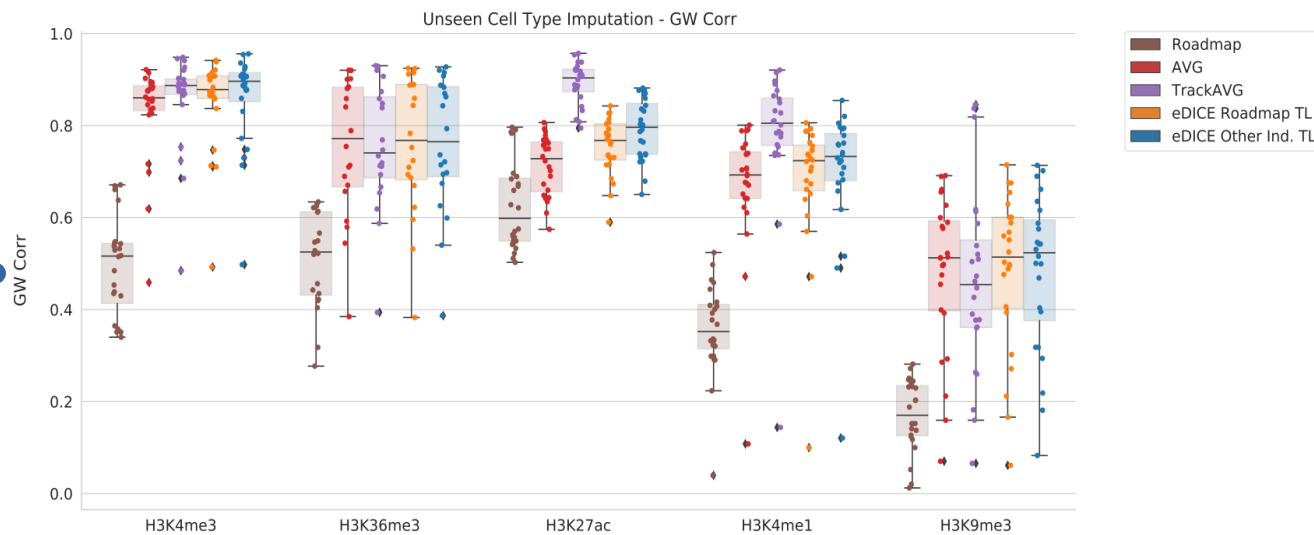
Averaging over individuals finds most of the enriched regions, which are conserved between patients



By construction, TrackAVG does not contain personalised information. We examined the performance of the imputations after masking out the shared regions.

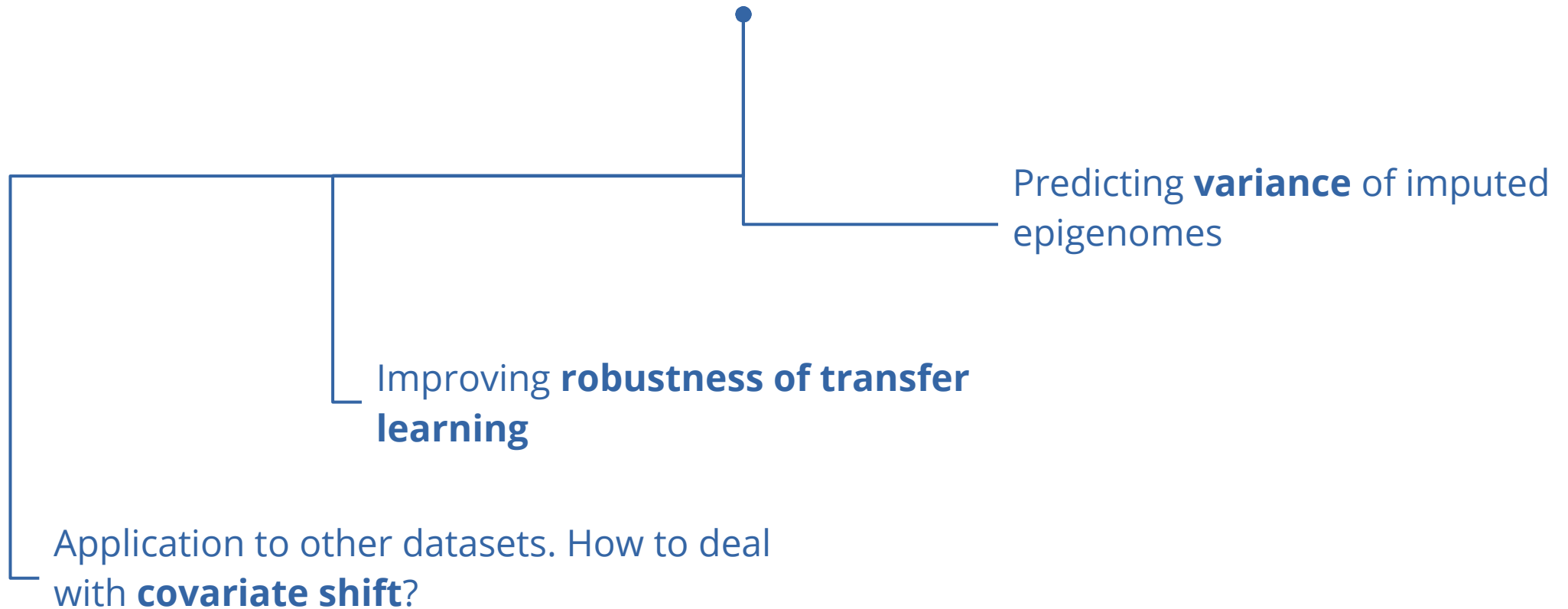


The transfer learning process allows generalisation to unseen tissues



The OOD imputations capture individual-specific enrichment better than the baseline

Open Questions and Future Work

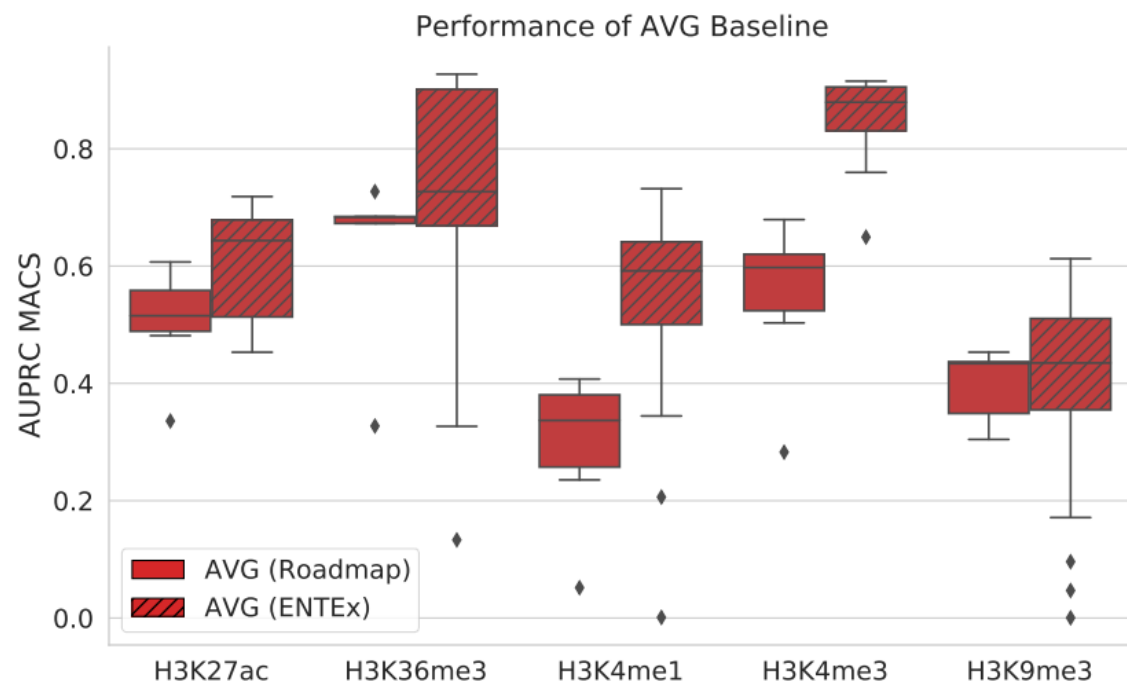


Acknowledgements

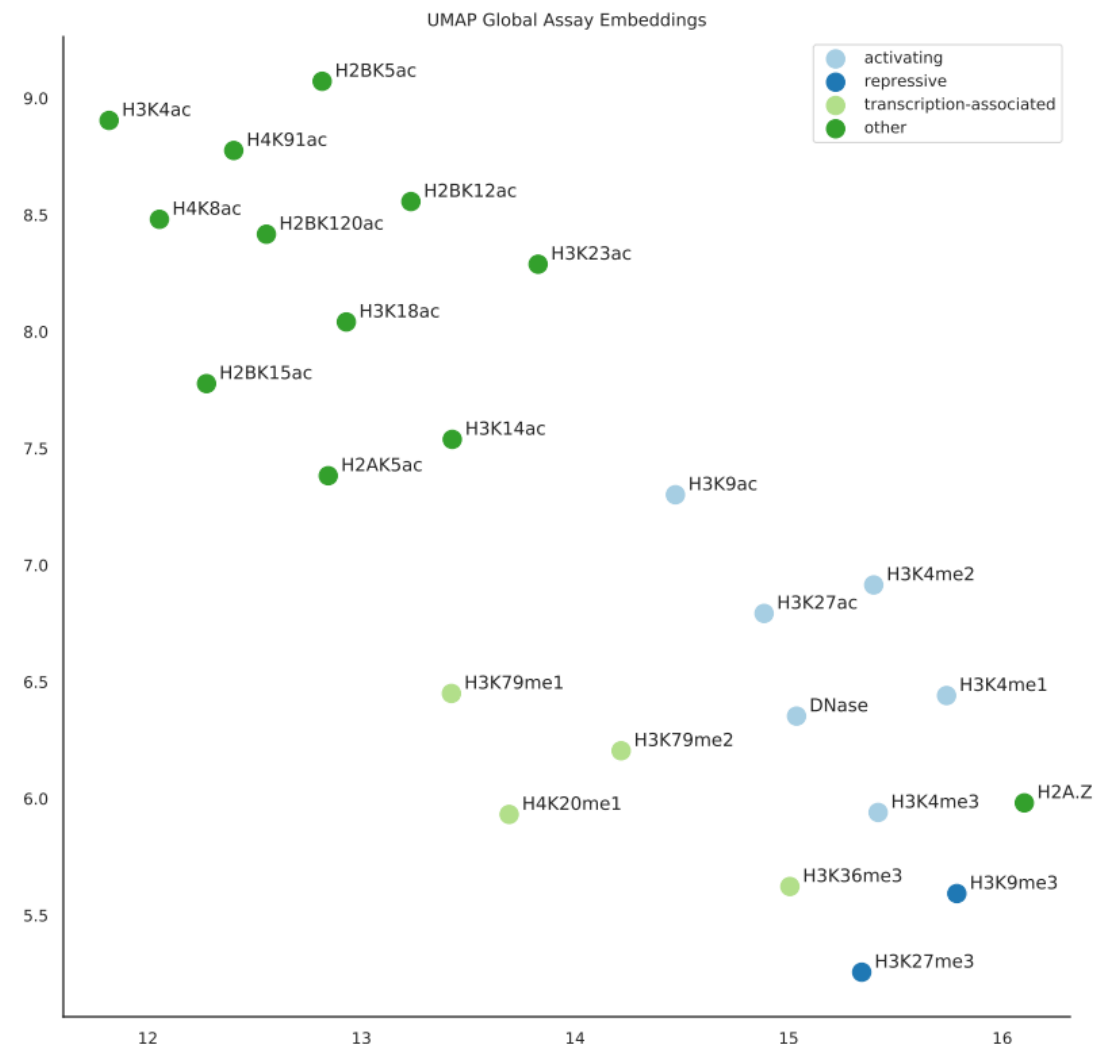
- Alex Hawkins-Hooker, UCL
- Tanmayee Narendra, University of Tübingen
- Mateo Rojas-Carulla, Laker AI
- Bernhard Schölkopf, Max-Planck Institute for Intelligent Systems
- Gabriele Schweikert, University of Dundee

Thank you for your attention

Systematic dataset shifts, differences between assays, and hidden confounders are difficult challenges for personalised imputation



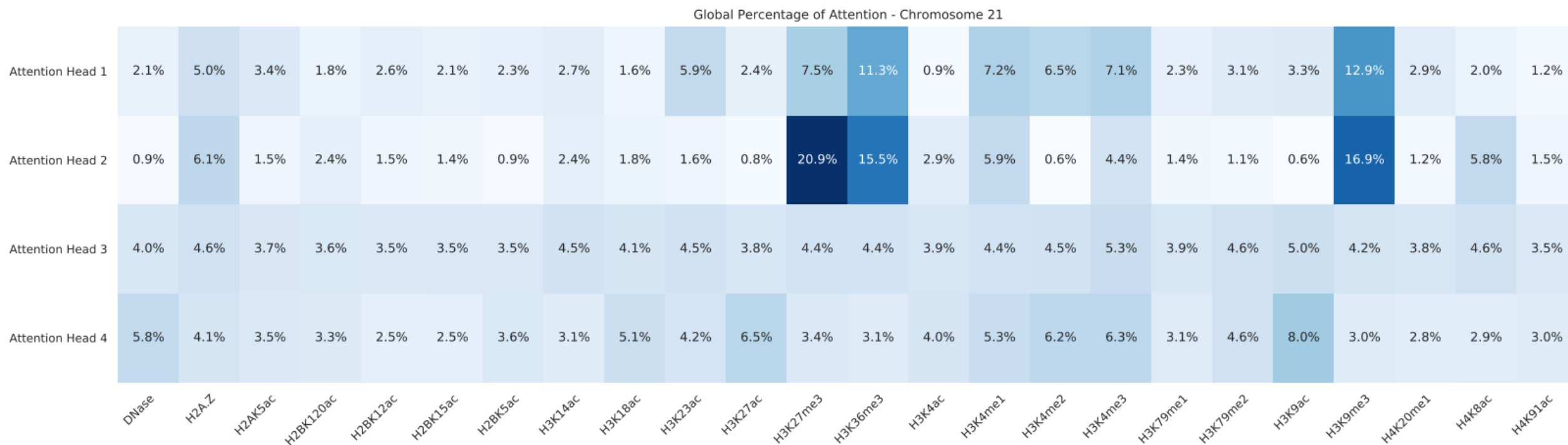
Global embeddings capture tissue similarity and general epigenetic mark function



Self-Attention opens up possibilities for **interpretation** of the model

Percentage of attention

$$p_h(\alpha) = \frac{\sum_{g \in G} \sum_{a \in A} w_{a\alpha}^{(g,h)}}{\sum_{g \in G} \sum_{a \in A} \sum_{a' \in A} w_{aa'}^{(g,h)}}$$



How does this portion of attention **shift** within functional regions of the genome?



$$d_h^{(R)}(\alpha) = \frac{p_h(\alpha)|_{G \equiv R} - p_h(\alpha)}{p_h(\alpha)}$$

