avocado

# Deep tensor factorization characterizes the human epigenome through imputation of thousands of epigenomic and transcriptomic experiments

Jacob Schreiber
Paul G. Allen School of Computer Science and Engineering
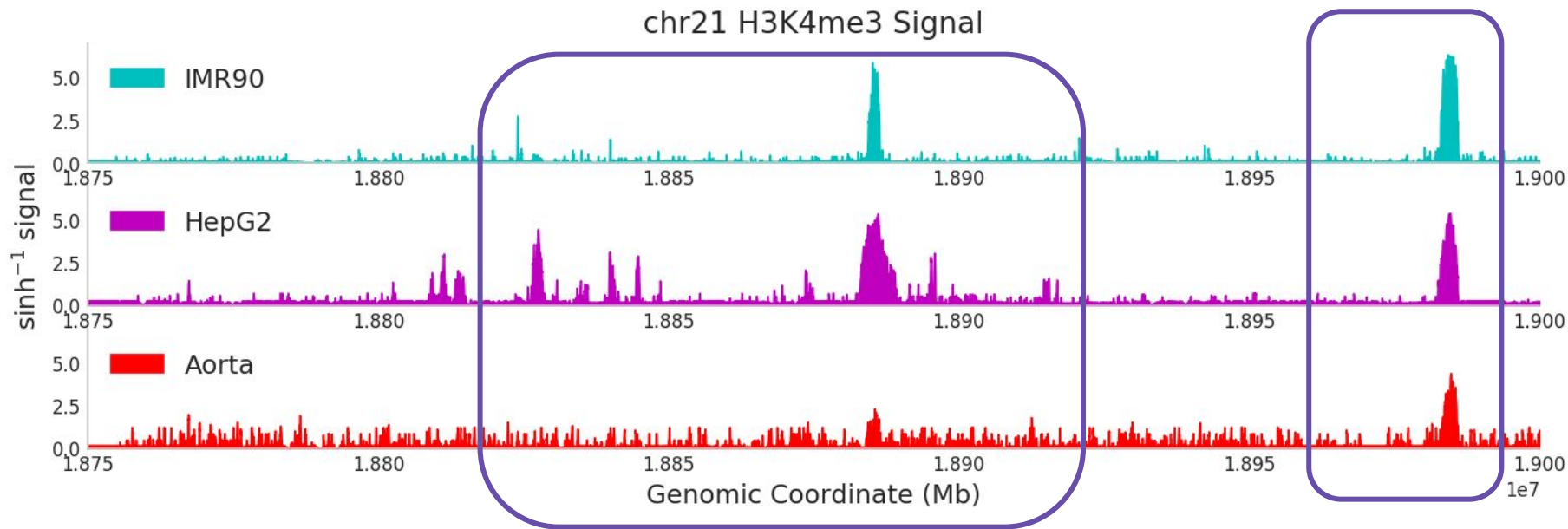University of Washington
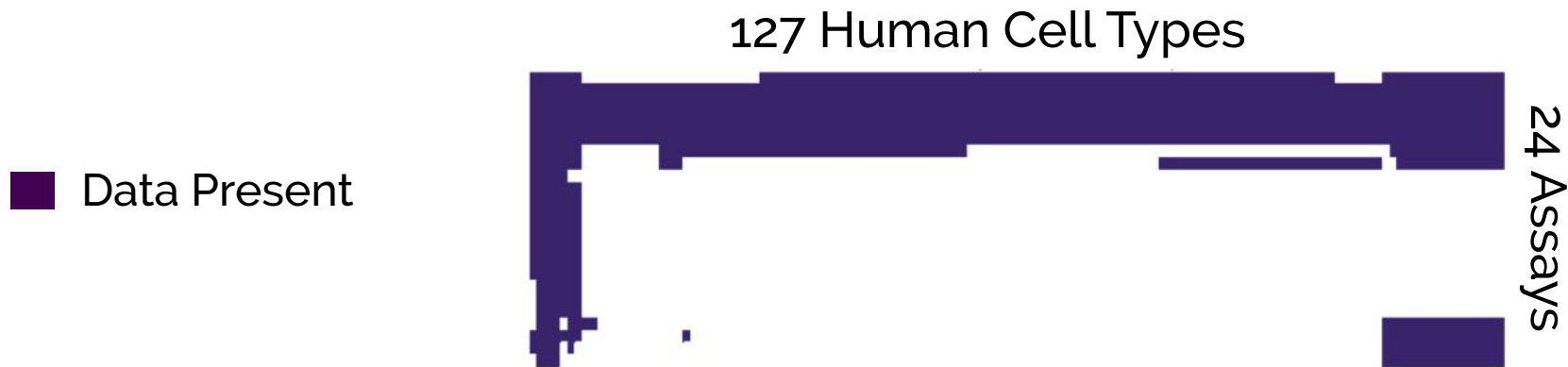
jmschreiber91

@jmschrei

@jmschreiber91

1

# The signal of epigenomic assays vary across cell types



chr21 H3K4me3 Signal

127 Human Cell Types
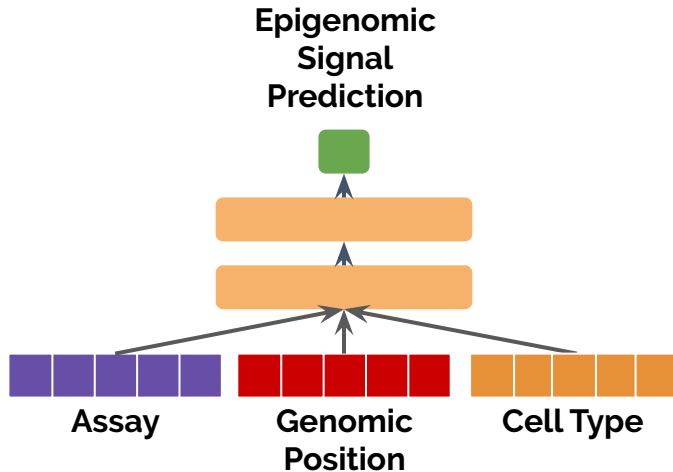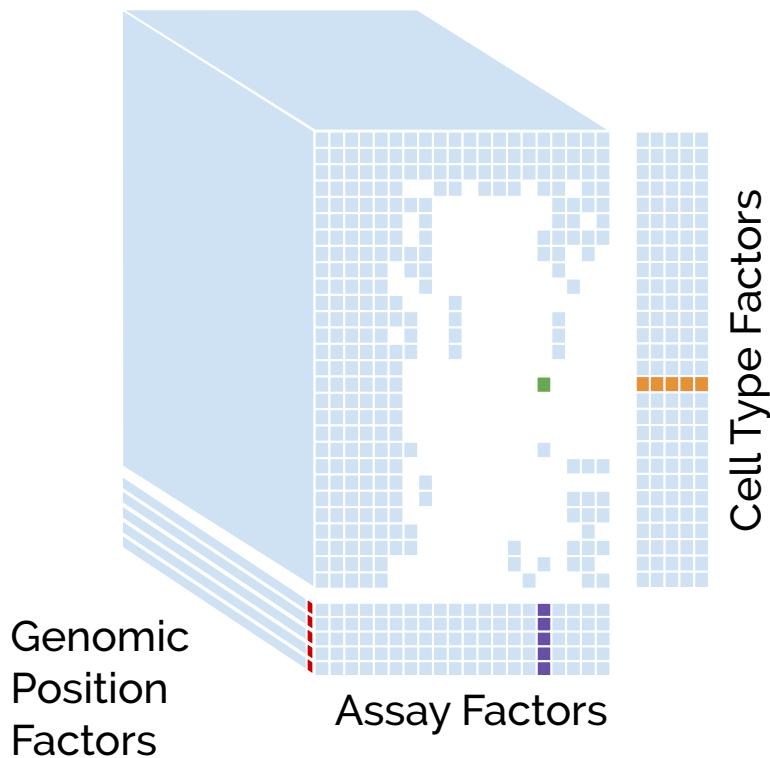
24 Assays



Data Present

Unfortunately the Roadmap compendium is incomplete. Previous work sought to fill in the matrix through imputing all potential experiments (ChromImpute[1], PREDICTD[2])

1. Ernst, et al. *Nature Methods, 2015*
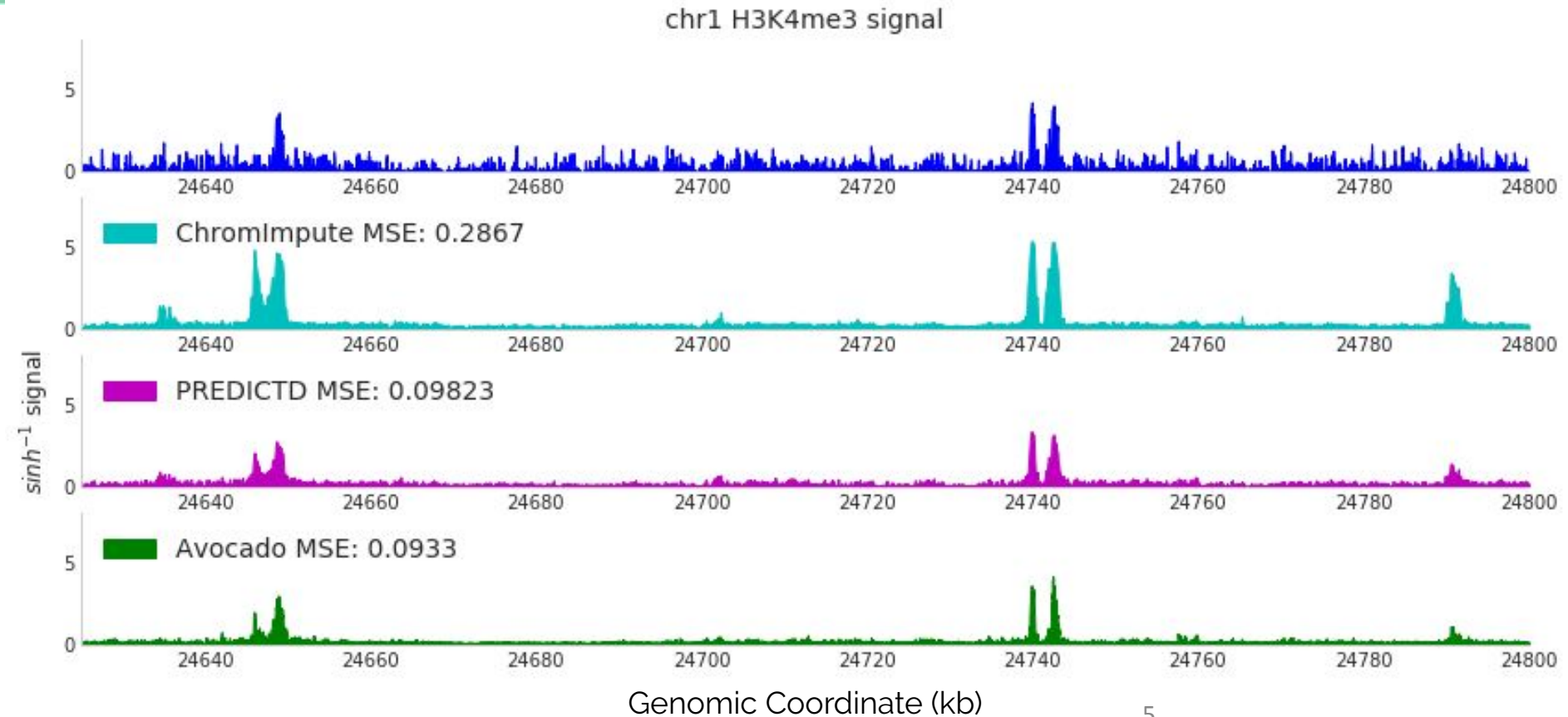2. Durham, et al. *Nature Communications, 2018*

3

# Avocado is a deep tensor factorization approach

Genomic Coordinate (kb)

5

# Avocado performs well genome-wide

| MSE- | global | 1obs | 1imp | Prom | Gene | Enh |
|---|---|---|---|---|---|---|
| ChromImpute | 0.113 | **0.941** | 1.09 | 0.3246 | 0.1494 | 0.3164 |
| PREDICTD | **0.1** | 1.76 | 0.897 | 0.2576 | **0.1295** | 0.267 |
| Avocado | **0.1** | 1.66 | **0.845** | **0.249** | **0.1295** | **0.26** |

**MSE-global:** Mean squared error (MSE) across the full length of the genome

**MSE-1obs:** MSE at the top 1% of genomic positions ranked by experimental signal

**MSE-1imp:** MSE at the top 1% of genomic positions ranked by imputed signal

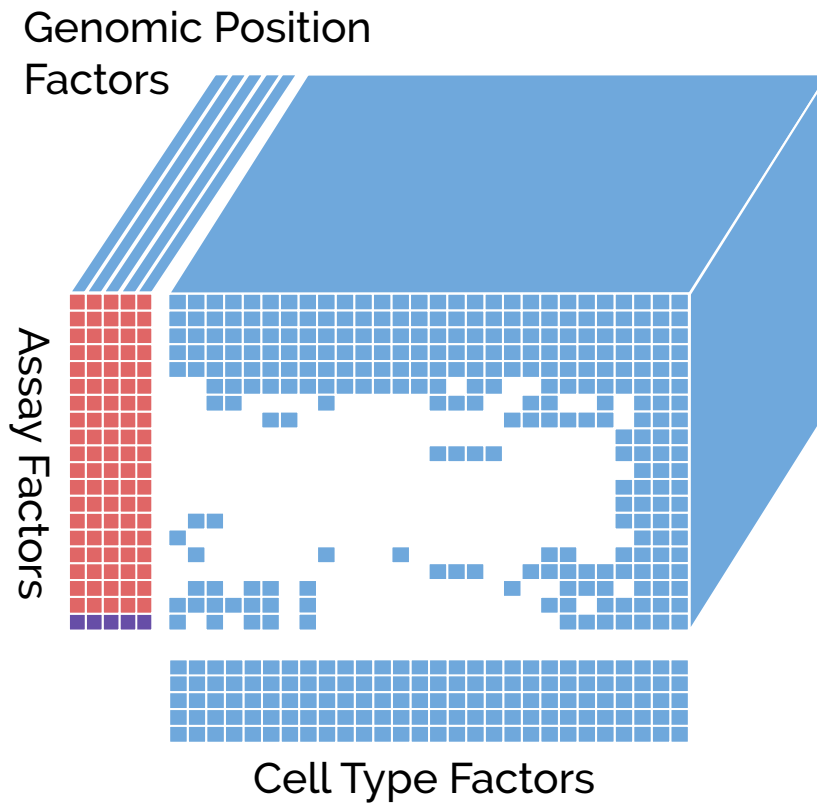**MSE-Prom:** MSE at promoter regions defined by GENCODE

**MSE-Gene:** MSE at gene bodies defined by GENCODE

**MSE-Enh:** MSE at enhancer regions defined by FANTOM5

# Okay, so have we characterized human epigenomics now?



Genomic Position Factors

**Histone Modification ChIP-seq**
**Chromatin Accessibility**

Assay Factors

Cell Type Factors

**# Cell Types: from 127 to 400**

**Histone Modification ChIP-seq**
**Chromatin Accessibility**



Genomic Position Factors

Assay Factors

Cell Type Factors

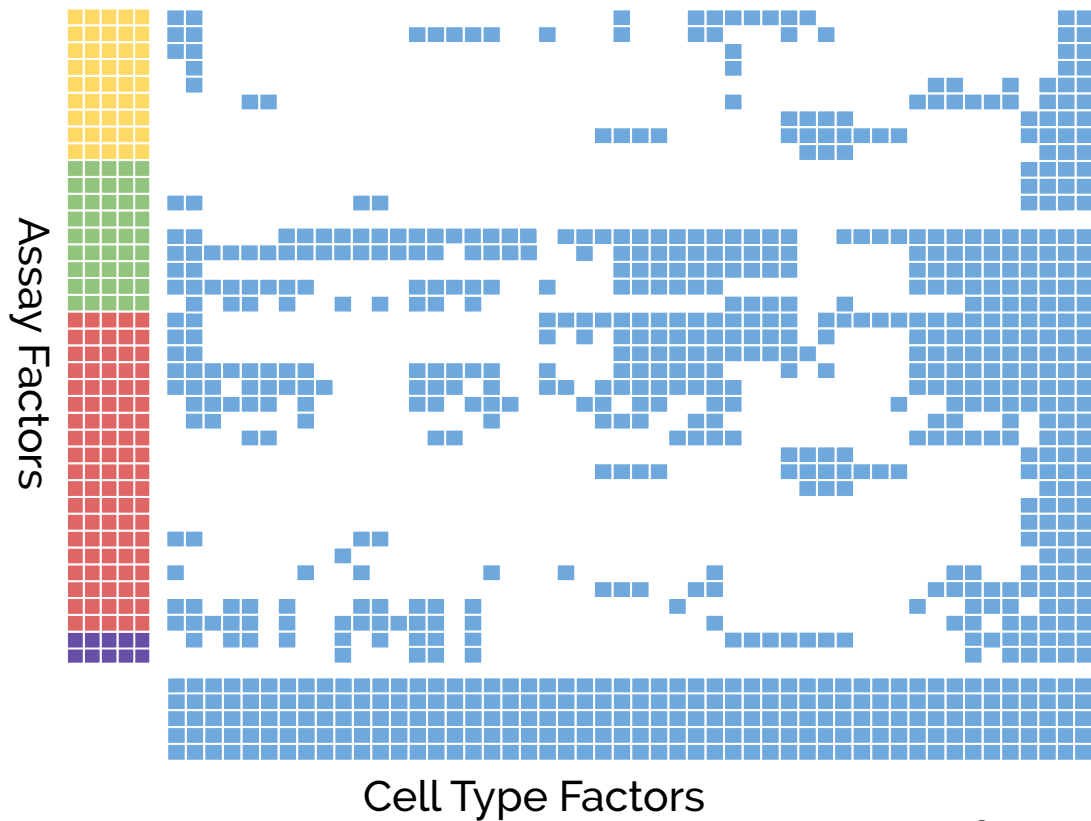# … and more assays

**# Cell Types: from 127 to 400**
**# Assays: from 24 to 76**

**Histone Modification ChIP-seq**
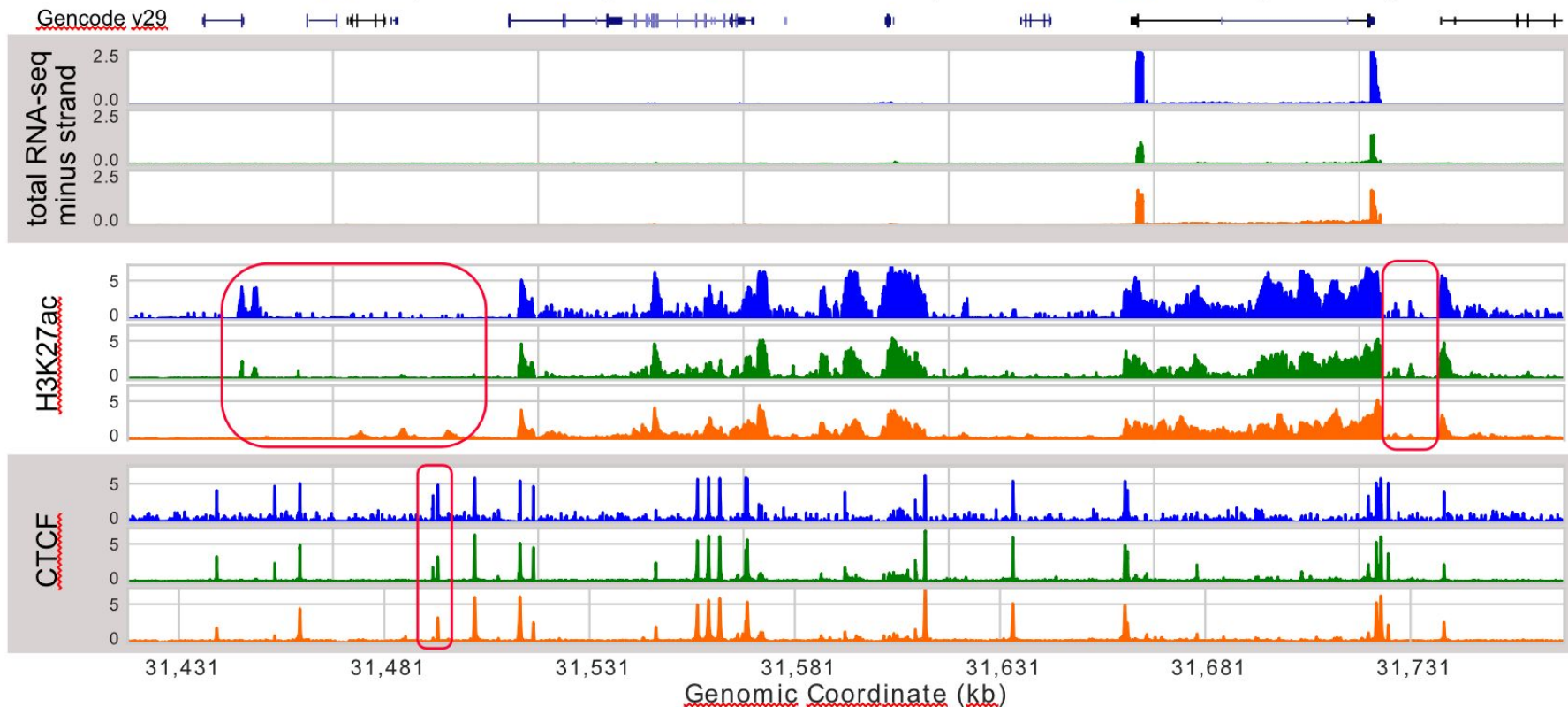**Chromatin Accessibility**
**Gene Transcription**
**Transcription Factor ChIP-seq**

Assay Factors

Cell Type Factors

# Avocado can jointly model many forms of activity

# Avocado imputes TF binding better than the participants in the ENCODE-DREAM challenge*

| Biosample | iPSC | PC-3 | liver | liver | liver | liver | liver | liver | liver |
|---|---|---|---|---|---|---|---|---|---|
| Assay | CTCF | CTCF | EGR1 | FOXA1 | GABPA | JUND | MAX | REST | TAF1 |
| Method | | | | | | | | | |
| Yuanfang Guan | 0.729 | 0.600 | 0.397 | 0.282 | 0.353 | 0.533 | 0.441 | 0.319 | 0.281 |
| dxquang | 0.866 | 0.783 | 0.274 | 0.400 | 0.347 | 0.260 | 0.330 | 0.312 | 0.264 |
| autosome.ru | 0.778 | 0.486 | 0.331 | 0.243 | 0.342 | 0.416 | 0.384 | 0.264 | 0.221 |
| J-TEAM | **0.812** | 0.747 | 0.363 | **0.462** | 0.344 | 0.415 | 0.377 | 0.196 | 0.272 |
| Avocado | 0.723 | **0.791** | **0.530** | 0.354 | **0.396** | **0.660** | **0.574** | **0.477** | **0.384** |
| Similar Biosample | — | — | 0.363 | 0.389 | 0.226 | 0.568 | 0.446 | 0.408 | — |
| Same Biosample | 0.741 | 0.878 | 0.648 | 0.716 | 0.573 | 0.731 | 0.622 | 0.622 | 0.556 |
| Average Activity | 0.574 | 0.735 | 0.240 | 0.299 | 0.253 | 0.223 | 0.349 | 0.124 | 0.140 |

*Performance metric is auPR (average precision)*

\* read about the caveats in our preprint

# Okay so now have we fully characterized human epigenomics?

No; the ENCODE compendium does not include hundreds of protein binding assays or a number of cell states, diseases, and mutations.
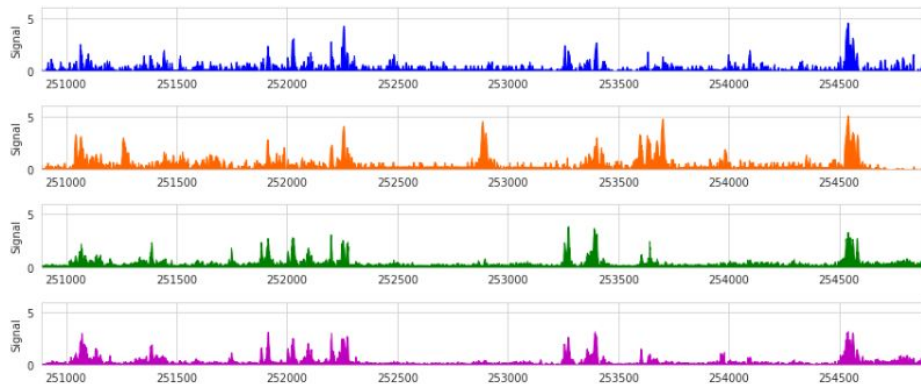
However:
- New biosamples and assays can be added to a pre-trained model with as little as a single experiment

- We are exploring zero-shot imputation approaches that precalculate assay embeddings using protein similarity and interaction networks

# Leveraging the large amount of human data enables zero-shot imputation of TF binding across species

```
Average Activity:                  0.09677
Mouse + 3,814 Human Experiments:   0.09252
Mouse + 6,870 Human Experiments:   0.08570
```

ELF1                                      MAX



Genomic Position                          Genomic Position

# GitHub repo, pretrained models, and preprints online!



**https://github.com/jmschrei/avocado**

# Acknowledgements

Timothy Durham   Deepthi Hedge        Jeffrey Bilmes   William Noble

eScience Institute
ADVANCING DATA-INTENSIVE DISCOVERY IN ALL FIELDS

# The learned latent representations capture known associations