

Predicting gene expression from plasma cell-free DNA using both the fragment length and fragment position

John St John¹, Erik Gafni¹, Brandon White¹, Ajay Kannan¹, Loren Hansen¹, Artur Jaroszewicz¹, Anshul Kundaje², Nathan Boley¹

¹Freenome; ²Stanford University

BACKGROUND

- The ability to use a blood sample to determine the transcriptional state of cells that are releasing DNA into the bloodstream of a patient may be helpful in a variety of clinical applications, including the early detection of cancer
- Cell-free DNA (cfDNA) contains epigenetic signatures of the cells from which it was produced. As a result, cfDNA can be used to predict the gene expression state of cfDNA-producing cells. To date, the two published approaches used to predict gene expression largely ignore cfDNA fragment size^{1,2}
- V-plots are a powerful visualization technique originally applied to MNase-seq data. These plots show the density of fragments of each length at a particular genomic location, and can provide single base pair resolution of nucleosomes as well as other proteins that protect DNA from digestion³
- cfDNA closely approximates MNase-seq data^{4,5}; we therefore used V-plots as an information-rich input to our gene expression prediction model

OBJECTIVES

- To develop a gene expression prediction model that uses cfDNA fragment coverage and length to predict which genes are highly or lowly expressed in cfDNA-producing cells
- To apply this model of gene expression prediction to a set of colon-specific genes in order to detect colon cancer and adjacent colon-derived cfDNA, which is expected to be present in patients with advanced colorectal cancer (CRC)

METHODS

Sample collection

- De-identified plasma samples from patients with CRC (n=532) and non-cancer controls (n=234) were obtained from academic medical centers and commercial biobanks. CRC stage information was as follows: stage I (n=169), stage II (n=256), stage III (n=97), stage IV (n=6) and unknown stage information (n=8)

Figure 1. Model architecture from gene expression to disease prediction

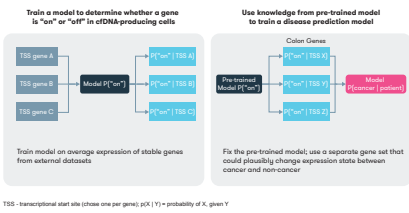


Figure 2. V-plots made from cfDNA capture DNA-protein associations and reflect transcriptional state

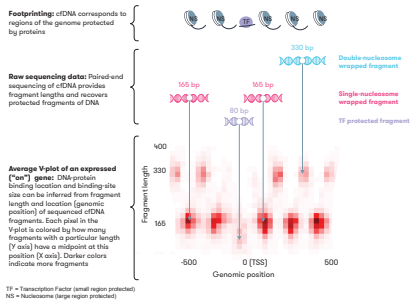
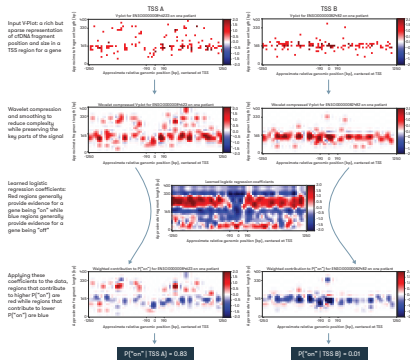
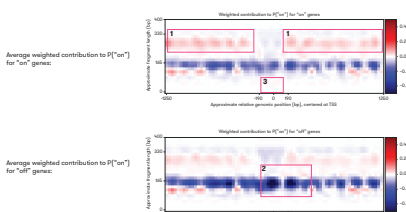


Figure 3. Predicting gene expression from cfDNA-derived V-plots around TSS regions



RESULTS

Figure 4. Interpretation of what the gene expression model learned



- Strongest drivers of predicting "on" are dinucleosome peaks flanking the TSS both up and downstream (Figure 4 Region 1) and a relatively weak mononucleosome band
- Strongest drivers predicting "off" are mononucleosome positions (especially Figure 4 Region 2) and a relatively weak dinucleosome band
- Although not always present, short sub-mononucleosome fragments at the TSS support an "on" prediction (Figure 4 Region 3) [see the logistic regression coefficients in Figure 3 for more evidence of this]

Figure 5. Classifiers using representations of fragment length and position accurately categorize "on" and "off" genes

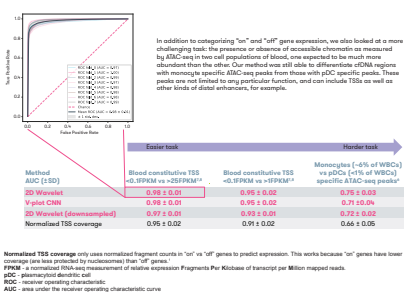
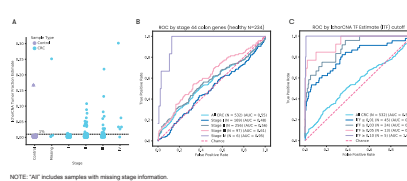
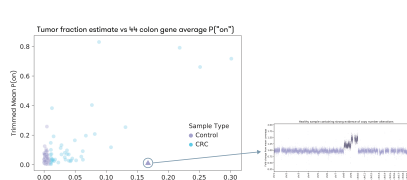


Figure 6. Tumor-targeted gene set enables classification of cancer samples as a function of cfDNA tumor fraction rather than stage



- For this approach we used 44 genes expressed in colon and not in blood cells as measured in the Roadmap Epigenomics Project⁶. We expect that colon genes will be expressed in colon cancer, as well as adjacent healthy colon tissue, which does not contribute substantial quantities of material to cfDNA in healthy individuals⁶
- IchatRNA-based TF estimates (ITF) increase with stage but most stage I-III CRC have low estimated ITF (~1%) (Figure 6A)
- Classification performance increases more strongly by tumor fraction than stage (Figure 6B & 6C)

Figure 7. Average gene expression prediction can augment CNV-based tumor fraction estimation



- A high ITF non-cancer control displayed a low average probability of expression P("on") of the 44 colon genes, differentiating it from high ITF CRC samples (Figure 7). These copy number changes may either be germline or somatic but do not originate from colon/CRC DNA. Plausible sources include DNA from blood-cells or from a non-CRC tumor

CONCLUSIONS

- 2D representations of fragment length and location can be used to accurately predict extremes in gene expression (Figure 5)
- The method presented here can accurately predict whether a patient has cancer with high fractions of tumor-derived cfDNA, which are typically observed in later stages but can be observed at any stage of disease (Figure 6)
- Despite limited sensitivity in patients with low tumor fractions, one practical use for this method is in identifying cases where observed CNVs in cfDNA do not originate from the cancer of interest (Figure 7)

NEXT STEPS

- This approach could be used with different cell-type-specific gene sets to predict the tissue of origin of a cancer
- We are in the process of evaluating and verifying this approach on immune-derived signals, as well as combining with other analytes, for early cancer detection

ACKNOWLEDGEMENT

The authors gratefully acknowledge Dr. Andrew Godwin and the University of Kansas Cancer Center's Bioprescreening Repository Core Facility staff funded in part by the National Cancer Institute Cancer Center Support Grant (P30 CA68524), National Health Services Research Scotland, Tayside Biorepository, Genetisat Inc., iSpecimen Inc., and Individuum for support of this research by providing the identified plasma samples. We also thank Signe Fransen, Girish Pulato and David Weinberg for their extensive suggestions, feedback, and editorial support.

REFERENCES

- Ulz et al. Nature Genetics, 2016.
- Snyder et al. Cell, 2016.
- Hancock et al. PNAS, 2011.
- Imvor et al. BMC Genomics, 2015.
- Zhang et al. Clinical Chemistry, 2017.
- Calderon et al. bioRxiv, 2018.
- Roadmap Epigenomics Consortium. Nature, 2015.
- ENCODE Project Consortium. Nature, 2012.
- Adalsteinsson et al. Nat. Commun, 2017.
- Moss et al. Nat. Commun, 2018.