

Nucleosome Patterns in Circulating Tumor DNA Reveal Transcriptional Regulation of Advanced Prostate Cancer Phenotypes



Navonil De Sarkar^{1,2,3,4}, Robert D. Patton^{1,2}, Anna-Lisa Doebley^{1,2,5,6}, Brian Hanratty², Mohamed Adil¹, Adam J. Kreitzman^{1,2}, Jay F. Sarthy⁷, Minjeong Ko^{1,2}, Sandipan Brahma⁷, Michael P. Meers⁷, Derek H. Janssens⁷, Lisa S. Ang², Ilsa M. Coleman², Arnab Bose², Ruth F. Dumpit², Jared M. Lucas², Talina A. Nunez², Holly M. Nguyen⁸, Heather M. McClure⁹, Colin C. Pritchard^{10,11}, Michael T. Schweizer^{3,12}, Colm Morrissey⁸, Atish D. Choudhury^{9,13}, Sylvan C. Baca^{9,13}, Jacob E. Berchuck⁹, Matthew L. Freedman^{9,13}, Kami Ahmad⁷, Michael C. Haffner^{2,3,10}, R. Bruce Montgomery¹², Eva Corey⁸, Steven Henikoff^{7,14}, Peter S. Nelson^{2,3,8,11,12}, and Gavin Ha^{1,2,11,15}

ABSTRACT

Advanced prostate cancers comprise distinct phenotypes, but tumor classification remains clinically challenging. Here, we harnessed circulating tumor DNA (ctDNA) to study tumor phenotypes by ascertaining nucleosome positioning patterns associated with transcription regulation. We sequenced plasma ctDNA whole genomes from patient-derived xenografts representing a spectrum of androgen receptor active (ARPC) and neuroendocrine (NEPC) prostate cancers. Nucleosome patterns associated with transcriptional activity were reflected in ctDNA at regions of genes, promoters, histone modifications, transcription factor binding, and accessible chromatin. We identified the activity of key phenotype-defining transcriptional regulators from ctDNA, including AR, ASCL1, HOXB13, HNF4G, and GATA2. To distinguish NEPC and ARPC in patient plasma samples, we developed prediction models that achieved accuracies of 97% for dominant phenotypes and 87% for mixed clinical phenotypes. Although phenotype classification is typically assessed by IHC or transcriptome profiling from tumor biopsies, we demonstrate that ctDNA provides comparable results with diagnostic advantages for precision oncology.

SIGNIFICANCE: This study provides insights into the dynamics of nucleosome positioning and gene regulation associated with cancer phenotypes that can be ascertained from ctDNA. New methods for classification in phenotype mixtures extend the utility of ctDNA beyond assessments of somatic DNA alterations with important implications for molecular classification and precision oncology.

¹Division of Public Health Sciences, Fred Hutchinson Cancer Center, Seattle, Washington. ²Division of Human Biology, Fred Hutchinson Cancer Center, Seattle, Washington. ³Division of Clinical Research, Fred Hutchinson Cancer Center, Seattle, Washington. ⁴Department of Pathology and Prostate Cancer Center of Excellence, Medical College of Wisconsin, Milwaukee, Wisconsin. ⁵Molecular and Cellular Biology Graduate Program, University of Washington, Seattle, Washington. ⁶Medical Scientist Training Program, University of Washington, Seattle, Washington. ⁷Division of Basic Sciences, Fred Hutchinson Cancer Center, Seattle, Washington. ⁸Department of Urology, University of Washington, Seattle, Washington. ⁹Dana-Farber Cancer Institute, Boston, Massachusetts. ¹⁰Department of Laboratory Medicine and Pathology, University of Washington, Seattle, Washington. ¹¹Brotman Baty Institute for Precision Medicine, Seattle, Washington. ¹²Division of Oncology, Department of Medicine, University of Washington, Seattle, Washington. ¹³Broad Institute of MIT and Harvard, Cambridge, Massachusetts. ¹⁴Howard Hughes Medical Institute, Chevy

Chase, Maryland. ¹⁵Department of Genome Sciences, University of Washington, Seattle, Washington.

Note: N. De Sarkar and R.D. Patton contributed equally to this work. P.S. Nelson and G. Ha share senior authorship of this work.

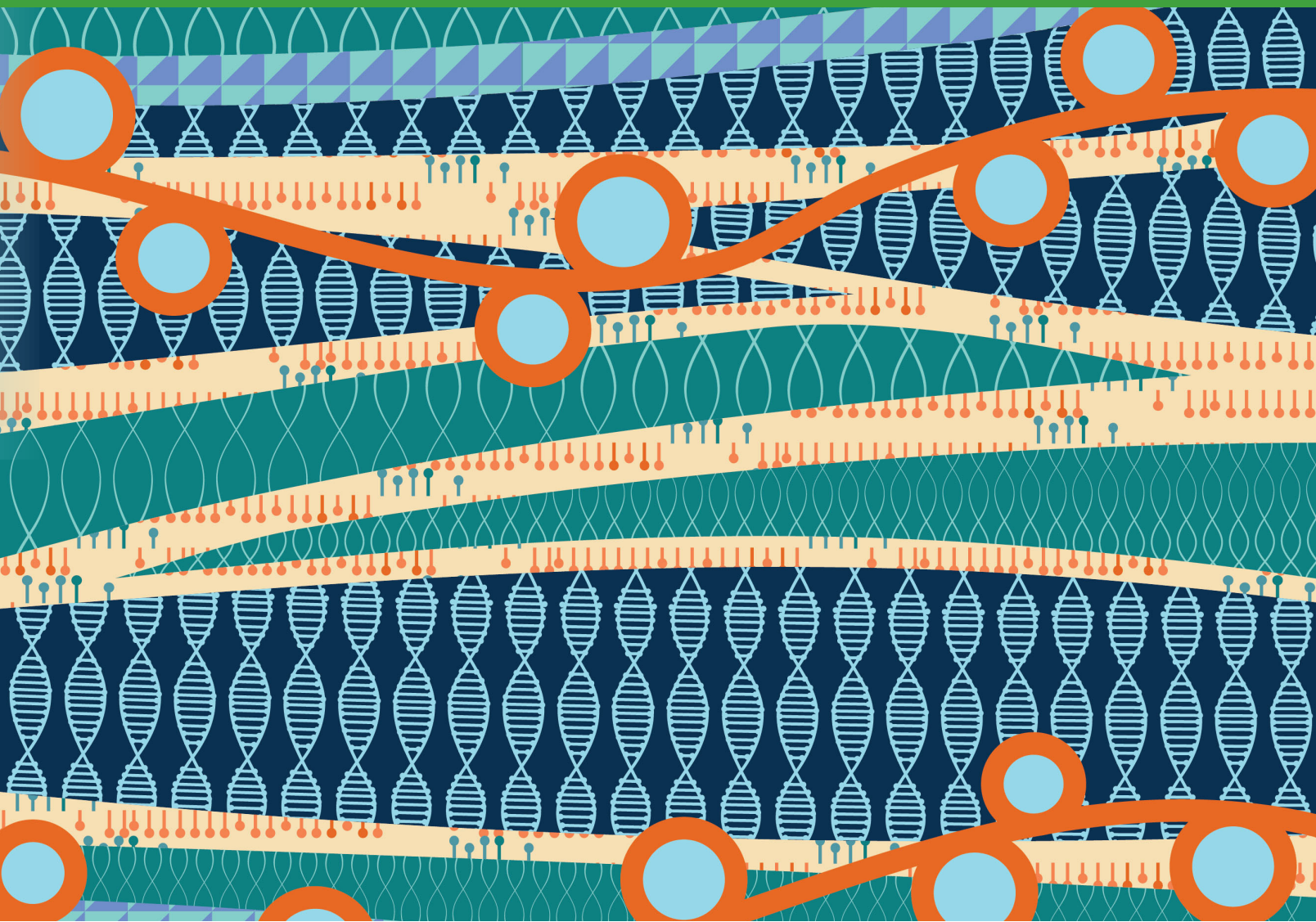
Corresponding Authors: Gavin Ha, Fred Hutchinson Cancer Center, 1100 Fairview Avenue North, Seattle, WA 98109. Phone: 206-667-2802; E-mail: gha@fredhutch.org; and Peter S. Nelson, Fred Hutchinson Cancer Center, 1100 Fairview Avenue North, Seattle, WA 98109. Phone: 206-667-3377; E-mail: pnelson@fredhutch.org

Cancer Discov 2023;13:632–53

doi: 10.1158/2159-8290.CD-22-0692

This open access article is distributed under the Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International (CC BY-NC-ND 4.0) license.

©2022 The Authors; Published by the American Association for Cancer Research



INTRODUCTION

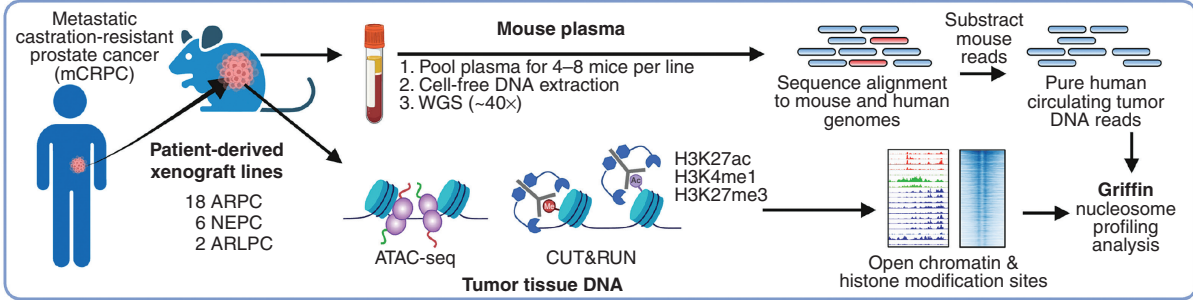
Metastatic castration-resistant prostate cancer (mCRPC) describes the stage in which the disease has developed resistance to androgen ablation therapies and is lethal (1). Androgen receptor signaling inhibitors (ARSI), designed for the treatment of CRPC, repress androgen receptor (AR) activity and improve survival, but these therapies eventually fail (2, 3). Since the adoption of ARSI as standard of care for mCRPC, there has been a prominent increase in the frequency of treatment-resistant tumors with neuroendocrine (NE) differentiation and features of small-cell carcinomas (4–7). These aggressive tumors may develop through a resistance mechanism of transdifferentiation from AR-positive adenocarcinoma (ARPC) to NE prostate cancer (NEPC) that lacks AR activity (4, 7–10). Additional phenotypes can also arise based on the expression of AR activity and NE genes, including AR-low prostate cancer (ARLPC) and double-negative prostate cancer (DNPC; AR-null/NE-null; refs. 5, 11–13). Distinguishing prostate cancer subtypes has clinical relevance in view of differential responses to therapeutics, but the need for a biopsy to diagnose tumor

histology can be challenging: invasive procedures are expensive and accompanied by morbidity, a subset of tumors are not accessible to biopsy, and bone sites pose particular challenges with respect to sample quality (7, 14).

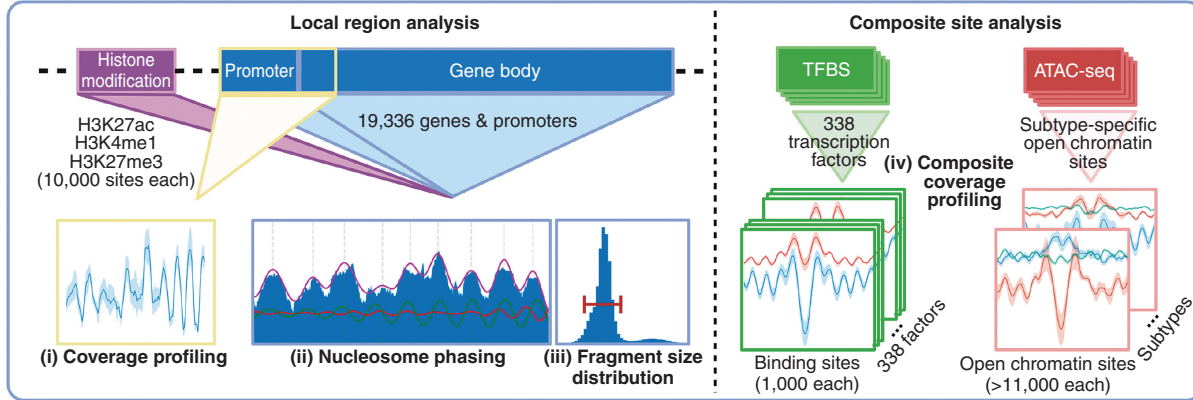
Circulating tumor DNA (ctDNA) released from tumor cells into the blood as cell-free DNA (cfDNA) is a noninvasive “liquid biopsy” solution for accessing tumor molecular information. The analysis of ctDNA to detect mutation and copy-number alterations has served to classify genomic subtypes of CRPC tumors (4, 15–21). However, the defining losses of *TP53* and *RBI* in NEPC do not always lead to NE transdifferentiation (7, 22). Rather, ARPC and NEPC tumors are associated with distinct reprogramming of transcriptional regulation (8, 9, 23). Methylation analysis of cfDNA in mCRPC to profile the epigenome shows promise for distinguishing phenotypes, but requires specialized assays such as bisulfite conversion, enzymatic treatment, or immunoprecipitation (24–27).

The majority of cfDNA represents DNA protected by nucleosomes when released from dying cells into circulation, leading to DNA fragmentation that is reflective of the nonrandom

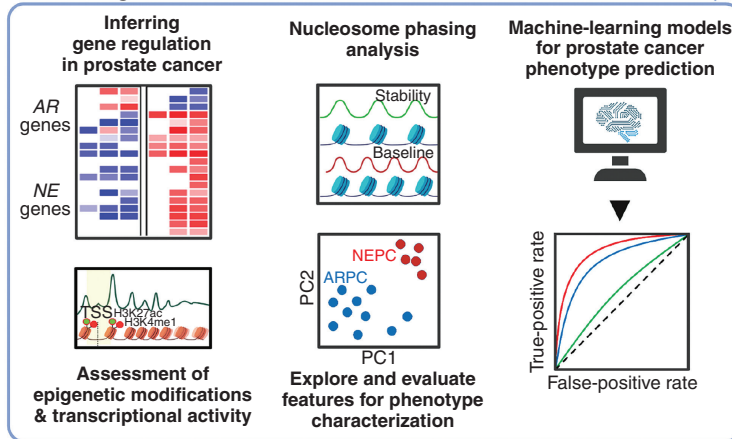
A PDX Plasma and tumor sequencing



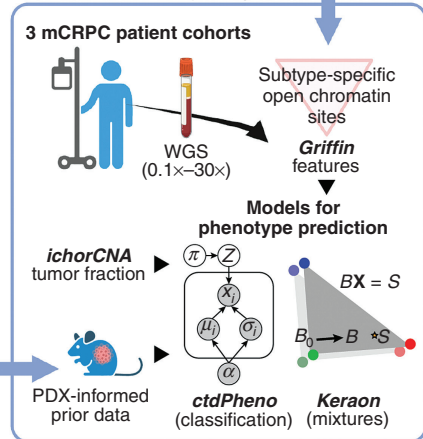
Circulating tumor DNA feature discovery



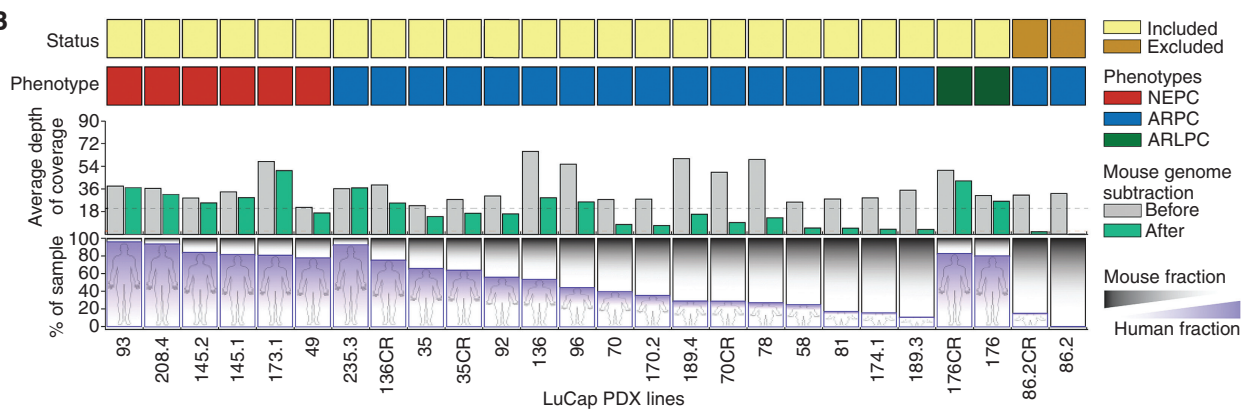
Circulating tumor DNA characterization



Patient phenotyping



B



enzymatic cleavage by nucleases (28, 29). Emerging approaches to analyze cfDNA fragmentation patterns from plasma for studying cancer can be performed directly from the standard whole-genome sequencing (WGS; refs. 30–35). cfDNA fragments primarily have a characteristic size of ~167 bp, consistent with protection by a single core nucleosome octamer and histone linkers, but the size distribution may vary between healthy individuals and patients with cancer (36–39). Recent studies have demonstrated that the nucleosome occupancy in cfDNA at the transcription start site (TSS) and transcription factor binding site (TFBS) can be used to infer gene expression and transcription factor (TF) activity from cfDNA (40–43). However, nucleosome positioning and spacing are dynamic in active and repressed gene regulation (44–46). A detailed understanding of the nucleosome patterns and accessible chromatin associated with transcriptional regulation in tumor phenotypes has not been fully explored in cfDNA.

The objective of this study is to determine if ctDNA could be used to accurately classify tumor phenotypes in men with mCRPC. A major challenge for ctDNA analysis is the low tumor content (tumor fraction) in patient plasma samples. By contrast, plasma from patient-derived xenograft (PDX) models may contain nearly pure human ctDNA after bioinformatic exclusion of mouse DNA reads (37, 39, 47). This provides a resource that is ideal for studying the properties of ctDNA, developing new analytic tools, and validating both genetic and phenotypic features by comparison with matching tumors. In this study, we performed WGS of ctDNA from mouse plasma across 24 CRPC PDX lines with diverse phenotypes. Applying newly developed computational methods, we comprehensively interrogated the nucleosome patterns in ctDNA across genes, regulatory loci, TFBSs, TSSs, and open chromatin sites to reveal transcriptional regulation associated with mCRPC phenotypes. Finally, we designed two probabilistic models to accurately classify treatment-resistant tumors into divergent phenotypes and to estimate the phenotype heterogeneity within a ctDNA sample. We then validated the performance of these models in 159 plasma samples from three mCRPC patient cohorts. Overall, these results highlight that transcriptional regulation of tumor phenotypes can be ascertained from ctDNA and has potential utility for diagnostic applications in cancer precision medicine.

RESULTS

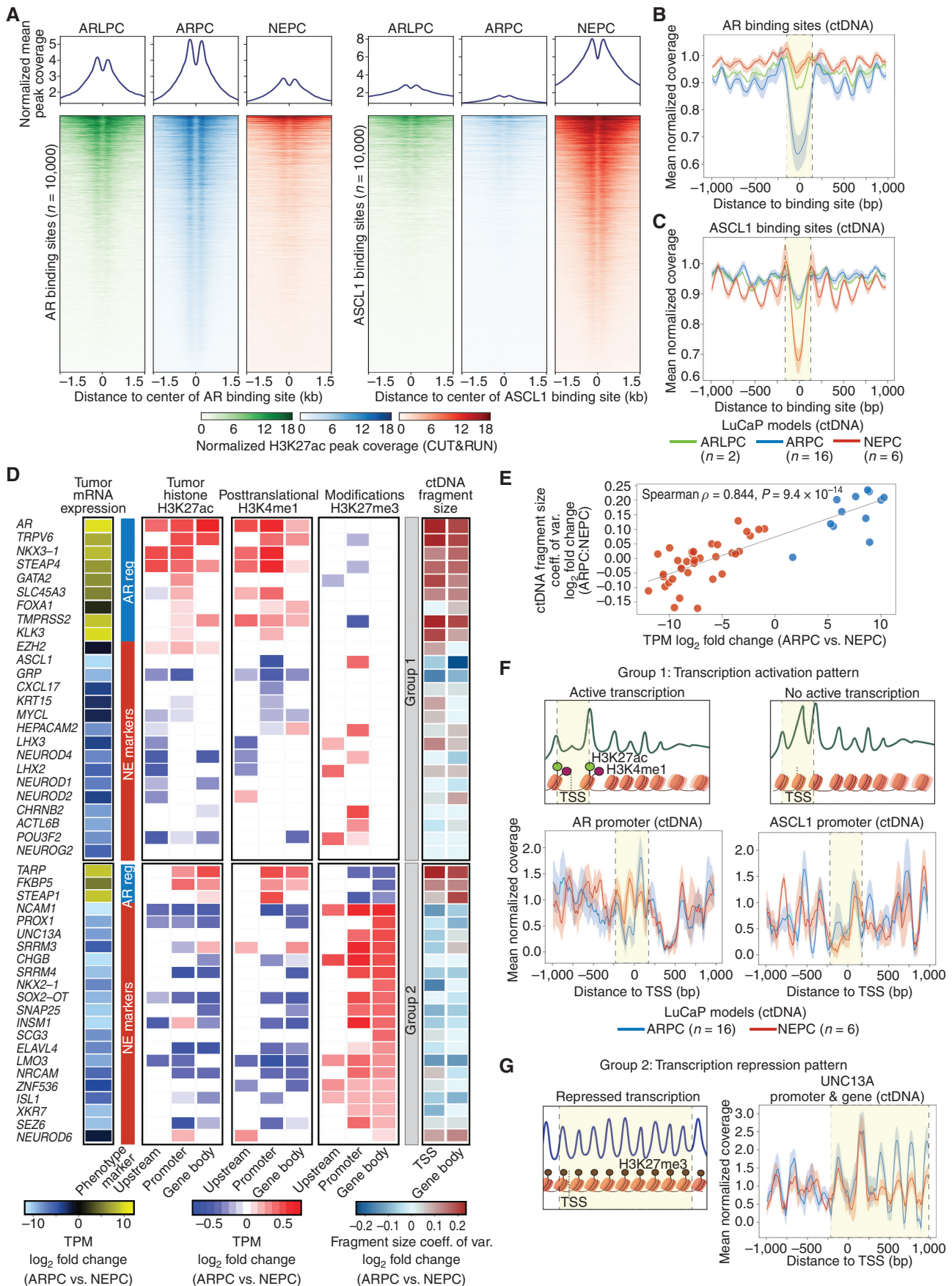
Comprehensive Resource of Matched Tumor and Liquid Biopsies from PDX Models of Advanced Prostate Cancer

To develop approaches for the accurate classification of mCRPC using ctDNA, we evaluated 26 models from the

LuCaP PDX series of advanced prostate cancer with well-defined phenotypes determined by whole-transcriptome RNA-sequencing (RNA-seq) and IHC assays for protein expression (48). There were 18 models classified as ARPC, two classified as AR-low and NE-negative prostate cancer (ARLPC), and six classified as NEPC (Fig. 1A). For each PDX line, we pooled mouse plasma (1.9–3.0 mL) from four to eight mice (mean tumor volume range 393–1,239 mm³), extracted cfDNA, and performed deep WGS (mean 38.4× coverage, range 21–85×; Methods; Supplementary Table S1). We used bioinformatic subtraction of mouse-sequenced reads to obtain nearly pure human ctDNA data (Methods). We observed that 25 lines had human ctDNA comprising more than 10% of the sample (mean 52.9%; range, 10.6%–96%) with NEPC samples having significantly higher human fractions (mean 85.1%; range, 77.1%–96%, two-tailed Mann-Whitney *U* test $P = 9.6 \times 10^{-4}$; Fig. 1B; Supplementary Table S1). After subsequent filtering by human ctDNA sequencing coverage, 24 PDX lines remained for further analysis (16 ARPC, 6 NEPC, 2 ARLPC; mean 20.5×; range, 3.8–50.6×; Supplementary Table S1). In the matching tumors, we performed Cleavage Under Targets and Release using Nuclease (CUT&RUN) to profile H3K27ac, H3K4me1, and H3K27me3 histone post-translational modifications (PTM; refs. 49, 50; Supplementary Fig. S1). We hypothesized that nucleosome organization inferred from ctDNA reflects the transcriptional activity state regulated by histone PTMs (51).

To study transcriptional regulation in mCRPC phenotypes from ctDNA, we interrogated four different features: (i) local promoter coverage, (ii) nucleosome positioning, (iii) fragment size analysis, and (iv) composite TFBSs plus open chromatin sites analysis using the Griffin framework (ref. 52; Fig. 1A; Methods). First, we analyzed three different local regions within ctDNA: all gene promoters and gene bodies and sites of histone PTMs guided by CUT&RUN analysis. For each of the three local regions, we extracted features of nucleosome periodicity using a nucleosome phasing approach and computed the fragment size variability. For promoter regions, we also computed the coverage at the TSS. Next, we analyzed ctDNA at TFBSs and open chromatin regions. For each TF, ctDNA coverage at TFBSs was aggregated into composite profiles representing the inferred activity (42, 52). Similarly, features in the composite profiles of phenotype-specific open chromatin regions were extracted for analyzing the signatures of chromatin accessibility in ctDNA. Altogether, we assembled a multiomic sequencing dataset from matching tumor and plasma for a total of 24 PDX lines, making this a unique molecular resource and platform for developing transcriptional regulation signatures of tumor phenotype prediction from ctDNA.

Figure 1. Characterizing advanced prostate cancer through matched tumor and liquid biopsies from PDX models. **A**, Top, both blood and tissue samples were taken from 26 PDX mouse models with tumors originating from mCRPC with AR-positive adenocarcinoma (ARPC), neuroendocrine (NEPC), AR-low neuroendocrine-negative (ARLPC) phenotypes. cfDNA was extracted from pooled plasma collected from 4 to 8 mice and WGS was performed. Following bioinformatic mouse read subtraction, pure human ctDNA reads remained. From PDX tissue, ATAC-seq and CUT&RUN (targeting H3K27ac, H3K4me1, and H3K27me3) data were generated. Middle, four distinct ctDNA features were analyzed at five genomic region types using Griffin (52) or nucleosome phasing methods developed in this study (Methods). Bottom (left), overview of PDX ctDNA features profiled to characterize the mCRPC pathways, transcriptional regulation, and nucleosome positioning. ctDNA features were evaluated for phenotype classification. Bottom (right), phenotype classification using probabilistic and analytic models that accounted for ctDNA tumor content and were informed by PDX features were applied to 159 samples in three patient cohorts. **B**, PDX phenotypes and mouse plasma sequencing. Inclusion status based on final mean depth after mouse read subtraction (< 3× coverage was excluded; red dotted line). Phenotype status, including 6 NEPC, 18 ARPC (2 excluded), and 2 ARLPC. Average depth of coverage before and after mouse subtraction (mean coverage 20.5×; dotted line). Percentage of the cfDNA sample that contains human ctDNA after mouse read subtraction.



Characterizing Transcriptional Activity of AR and ASCL1 in PDX Phenotypes through Analysis of Tumor Histone Modifications and ctDNA

We sought to further characterize the transcriptional activity in different tumor phenotypes by studying epigenetic regulation via histone PTMs. We identified broad peak regions for H3K4me1 (median of 17,643 regions; range, 1,894–64,934), H3K27ac (median 7,093; range, 1,610–34,047), and H3K27me3 (median 8,737; range, 2,024–42,495) in the tumors of the 24 PDX lines and an additional nine LuCaP PDX lines where only tumors were available (total of 25 ARPC, 2 ARLPC, and 6 NEPC; Methods; Supplementary Fig. S1; Supplementary Table S2). Using unsupervised clustering and principal components analysis (PCA), we identified putative active regulatory regions of enhancers and promoters (H3K27ac and H3K4me1) and gene repressive heterochromatic marks (H3K27me3) that were specific to ARPC, ARLPC, and NEPC phenotypes (ref. 53; Supplementary Fig. S2A).

AR and ASCL1 are two key differentially expressed TFs with known regulatory roles in ARPC and NEPC phenotypes, respectively (9, 54–56). When inspecting AR binding sites in ARPC tumors, we observed increased signals from flanking nucleosomes with H3K27ac PTMs compared with the other phenotypes (area under mean peak profile of 18.46 vs. 15.08 in ARLPC and 10.63 in NEPC, Fig. 2A; Supplementary Fig. S2B; Methods). We also observed the strongest signals at the nucleosome depleted region (NDR) in ARPC for H3K27ac (1.54 coverage decrease vs. 0.78 for ARLPC and 0.41 for NEPC). Conversely, in NEPC tumors, we observed stronger signals at nucleosomes with H3K27ac PTMs flanking ASCL1 binding sites (area under mean peak profile 62.65 vs. 29.18 for ARLPC and 10.83 for ARPC), and stronger NDR signals (2.26 coverage decrease vs. 0.19 for ARPC and 0.37 for ARLPC). We observed similar trends for H3K4me1 PTMs in the LuCaP PDX lines (Supplementary Fig. S2C).

We analyzed the ctDNA composite coverage profiles at 1,000 consensus TFBSs to evaluate nucleosome accessibility, where lower normalized central (± 30 bp window) mean coverage across these sites suggests more nucleosome depletion (Methods). For AR TFBSs, we observed the strongest signal for nucleosome depletion in ARPC, as indicated by the lowest mean central coverage (average 0.64, $n = 16$), compared with moderate signals for ARLPC (average 0.88, $n = 2$), and weakest signals for NEPC (average 0.95, $n = 6$; Fig. 2B). Conversely, the composite coverage profile at ASCL1 TFBSs showed the

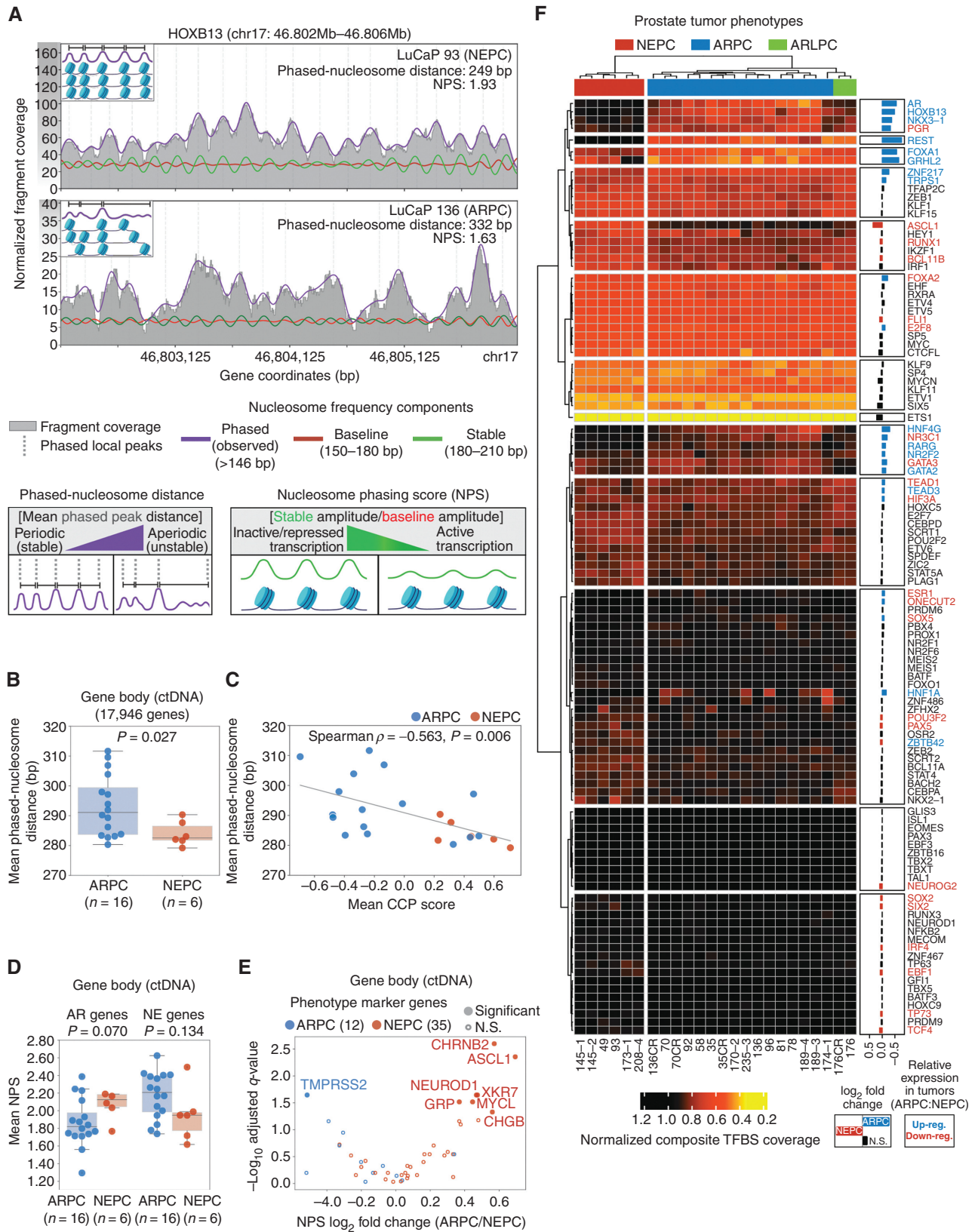
strongest nucleosome depletion for NEPC samples (mean central coverage 0.69) compared with ARLPC (0.86) and ARPC (0.88; Fig. 2C). These observations were consistent with the differential binding activity by AR and ASCL1 in their respective phenotypes from tumor tissue (Fig. 2A). We confirmed the same differential binding activity trends when analyzing TFBSs identified from other primary tissue sources (refs. 9, 57, 58; Supplementary Fig. S3A and S3B). We also noted that the composite TFBS coverage patterns in ctDNA resembled the NDR flanked by nucleosomes with H3K27ac and H3K4me1 modifications inferred by CUT&RUN (Fig. 2A; Supplementary Fig. S2B and S2C). Together, these results suggest that the nucleosome depletion in ctDNA at AR and ASCL1 binding sites represents active TF binding and regulatory activity in specific prostate PDX tumor phenotypes.

Nucleosome Patterns at Gene Promoters Inferred from ctDNA Are Consistent with Transcriptional Activity for Phenotype-Specific Genes

We selected 47 genes comprising 12 ARPC and 35 NEPC lineage markers established previously (4, 5, 59, 60) and confirmed their phenotype associations by RNA-seq from the PDX tumors (Fig. 2D; Supplementary Table S3; Methods). To assess the activity of these genes from ctDNA, we analyzed the ctDNA fragment size in TSSs (± 1 kb window) and gene bodies and found that the differential size variability between phenotypes was positively correlated with relative expression (Spearman $r = 0.844$, $P = 9.4 \times 10^{-14}$; Fig. 2E; Supplementary Fig. S4; Supplementary Table S2; Methods). However, closer inspection of ctDNA coverage patterns at promoters revealed consistent nucleosome organization for transcription activity and repression (refs. 40, 61–63; Fig. 2D). Therefore, we grouped the genes based on differential signals in H3K27me3 histone PTMs, which are linked with polycomb repressive complex mediated regulation and chromatin compaction (64).

For 25 genes without differential H3K27me3 peaks (group 1), including *AR*, *KLK3*, and *ASCL1*, we observed nucleosome depletion at the TSS consistent with the presence of active PTMs, such as for *AR* (mean coverage 0.47, $n = 16$) in ARPC and *ASCL1* (0.30, $n = 6$) in NEPC samples (Fig. 2F; Supplementary Fig. S5). By contrast, we observed increased coverage at the TSS of *AR* (1.08) in NEPC and *ASCL1* (0.42) in ARPC, which supports nucleosome depletion in the absence of PTMs and inactive transcription. For 22 genes with differential H3K27me3 peaks (group 2),

Figure 2. Analysis of tumor histone modifications and ctDNA reveals nucleosome patterns consistent with transcriptional regulation in CRPC phenotype-specific genes. **A**, H3K27ac peak signals between ARLPC, ARPC, and NEPC PDX tumor phenotypes at 10,000 AR binding sites (left) and at ASCL1 binding sites (right). Binding sites were selected from the GTRD (ref. 71; Methods). **B** and **C**, Composite coverage profiles at 1,000 AR (**B**) and ASCL1 (**C**) binding sites in ctDNA analyzed using Griffin for 140–250 bp fragments (Methods). Coverage profile means (lines) and 95% confidence interval computed using 1,000 bootstraps for a subset of sites (shading) are shown. The region ± 150 bp is indicated with a vertical dotted line and yellow shading. **D**, Heat map of \log_2 fold change in 47 key genes upregulated and downregulated between ARPC and NEPC established through RNA-seq (left) grouped by the type of histone modification which dictates translation levels: Group 1 shows gene activity attributable to H3K27ac or H3K4me1 PTM marks in the gene promoters or putative distal enhancers and lacking H3K27me3 heterochromatic mark in the gene body; group 2 features gene body spanning H3K27me3 repression marks. Central columns show differential peak intensity for each of the assayed histone modifications, separated by whether they appear upstream or in the promoter or the body of each gene. On the right, the \log_2 fold change between ARPC and NEPC lines' ctDNA fragment size CV is shown for TSS ± 1 KB windows and respective gene bodies. **E**, Comparison of the \log_2 fold change (ARPC/NEPC) of mean mRNA expression vs. mean CV in the 47 phenotypic lineage marker genes' promoter regions. **F**, Top, illustrations of expected ctDNA coverage profiles for group 1 genes with and without H3K27ac or H3K4me1 modification leading to active and inactive transcription, respectively. Bottom, $\pm 1,000$ bp surrounding the promoter region for AR and ASCL1 in ARPC and NEPC. Shown are coverage profile means (lines) and 95% confidence interval computed using 1,000 bootstraps for a subset of sites (shading). Decreased coverage is reflective of increased nucleosome accessibility and thus increased transcription. Dotted line and yellow shading highlight the TSS at -230 bp and $+170$ bp. **G**, Illustration of expected ctDNA coverage profiles for group 2 genes with repressed transcription caused by H3K27me3 modifications in the gene body. Neuronal gene UNCL13A has increased nucleosome phasing in the ctDNA of ARPC samples compared with NEPC.



including *INSM1*, *CHGB*, and *SRRM4*, we observed a relatively consistent increase in nucleosome occupancy and phasing in the TSS as well as in the gene body for ~50% of the genes (Fig. 2G; Supplementary Fig. S6). The neuronal signaling genes in this group, such as *UNC13A* and *INSM1*, had reduced signals for the stable nucleosome dyad position, consistent with the heterogeneous (“fuzzy”) nucleosome patterns described for actively transcribed genes (45, 65). Interestingly, although *UNC13A* was active in NEPC tumors, we did not detect H3K27ac nor H3K4me1 PTM marks in the regulatory loci of this gene (Supplementary Fig. S7A and S7B). These results illustrate that ctDNA analysis can reveal patterns that are consistent with different modalities of transcriptional regulation by histone modifications for key genes that define prostate cancer phenotypes.

Phasing Analysis in ctDNA Reveals Nucleosome Periodicity Associated with Transcriptional Activity between CRPC Phenotypes

Regions of inactive or repressed transcription are expected to have stably bound nucleosomes, resulting in more periodic phasing in the gene body (62, 66, 67). Conversely, actively transcribed regions may exhibit overall disordered phasing in the gene body due to fast nucleosome turnover, resulting in relatively aperiodic patterns with highly varied protection from nucleases along the gene (68). To systematically quantify internucleosomal spacing and predict nucleosome phasing, we developed TritonNP, a tool utilizing Fourier transforms and band-pass filters on guanine-cytosine-corrected ctDNA coverage to isolate frequency components corresponding to phased nucleosomes (Fig. 3A; Supplementary Fig. S8A and S8B; Methods). This approach allows for calling phased nucleosome dyad positions to generate an average internucleosome distance from the originating cells, encapsulating potential heterogeneity in nucleosome occupancy and stability. In PDX ctDNA, we observed a larger mean phased-nucleosome distance across 17,946 genes in the ARPC lines compared with the NEPC lines (median 291.1 bp vs. 282.6 bp, $P = 0.027$; two-tailed Mann-Whitney U test; Fig. 3B). The phased nucleosome distance was also negatively correlated with the mean cell-cycle progression (CCP) score (Spearman $\rho = -0.563$, $P = 0.006$; Fig. 3C; Methods). These results suggest that increased nucleosome periodicity in NEPC ctDNA may reflect the condensed chromatin in hyperchromatic nuclei of NE cells (14), and the phasing

analysis may have potential utility for assessing tumor proliferation and aggressiveness (69).

To model the relationship between nucleosome phasing and transcriptional activity more directly, we further extracted the frequency components corresponding to the interdyad distances of “stable” nucleosomes (180–210 bp) and a “baseline” component (150–180 bp) for normalization between samples of differing depths (70). We then computed the ratio of the mean frequency amplitudes between these components, which we designated the nucleosome phasing score (NPS), where a higher score corresponded to more ordered nucleosome phasing and repressed transcription. As an example, *HOXB13*, which is transcriptionally inactive in NEPC, had higher overall GC-corrected coverage (mean 56.85 depth) and a phased nucleosome distance of 249 bp with a 1.93 NPS in the LuCaP 93 NEPC PDX (Fig. 3A). By contrast, *HOXB13* is actively transcribed in ARPC and had lower coverage (mean 13.54 depth) and a more disordered phased-nucleosome distance of 332 bp with a 1.63 NPS in the LuCaP 136 ARPC PDX. When assessing the 47-prostate cancer phenotype marker genes, we observed that the mean NPS for the 35 NE genes was lower in NEPC lines compared with ARPC (median NPS 1.95 vs. 2.21, $P = 0.134$; two-tailed Mann-Whitney U test; Fig. 3D); although this was not statistically significant, it was consistent with their active transcription. Conversely, the mean NPS for the 12 AR-regulated genes was lower in ARPC lines compared with NEPC (median NPS 1.82 vs. 2.13, $P = 0.070$; two-tailed Mann-Whitney U test). In particular, 26 (74%) of the NE genes had lower NPS in NEPC compared with ARPC [\log_2 fold change (ARPC:NEPC) > 0], including seven genes (*ASCL1*, *CHGB*, *CHRN2*, *GRP*, *MYCL*, *XKR7*, and *NEUROD1*) that were statistically significant ($P < 0.05$); 10 (83%) of the AR-regulated genes had lower NPS in ARPC (\log_2 fold change < 0), with *TMPRSS2* being statistically significant (Fig. 3E; Supplementary Table S3). These results illustrate that the NPS captured signals distinguishing key lineage-specific gene markers.

Inferred TF Activity from Analysis of Nucleosome Accessibility at TFBSs in ctDNA Confirms Key Regulators of Tumor Phenotypes

To characterize the regulation of prostate tumor phenotype lineages, we considered nucleosome accessibility at TFBSs in PDX ctDNA for 338 TFs from the Gene Transcription Regulation Database (GTRD; ref. 71; Methods). First,

Figure 3. Phasing analysis in ctDNA recapitulates nucleosome stability and trends in transcriptional activity between CRPC phenotypes. **A**, Illustration of nucleosome phasing analysis using TritonNP for *HOXB13*, which is expressed in ARPC but not NEPC. Fourier transform and a band-pass filter-based smoothing method was used to identify phased peaks (gray dotted lines). Frequency components corresponding to nucleosome dyads (wavelength > 146 bp) are shown in purple. The mean internucleosome distance was computed from all peaks in the gene body: lower values represent more periodic and stable nucleosomes. NPS is defined as the ratio of the mean amplitudes between frequency components 180–210 bp (“stable,” green curve) and 150–180 bp (“baseline,” red curve). **B**, Boxplot of mean phased-nucleosome distance in 17,946 gene bodies per ctDNA sample for ARPC and NEPC PDX lines. Two-tailed Mann-Whitney U test P value is shown. **C**, Comparison of the mean phased-nucleosome distance and the mean CCP score (estimated from RNA-seq) for 16 ARPC and 6 NEPC PDX lines. **D**, Boxplot of NPS in gene bodies of 47 phenotype-defining genes (35 NE-regulated and 12 AR-regulated) between ARPC and NEPC lines. Two-tailed Mann-Whitney U test P values are shown. **E**, Volcano plot of NPS \log_2 fold change (ARPC/NEPC) in the 47 phenotype-defining genes. Genes with significantly higher NPS scores (solid-colored dots; two-tailed Mann-Whitney U test, Benjamini-Hochberg adjusted FDR at $P < 0.05$) and nonsignificant genes (open circle) are shown. **F**, Hierarchical clustering of the normalized composite central mean coverage at TFBSs from the Griffin analysis of ctDNA for 108 TFs in LuCaP PDX lines of ARPC ($n = 16$), NEPC ($n = 6$), and ARLPC ($n = 2$) phenotypes. This list of TFs was initially selected as having differential expression between ARPC and NEPC from LuCaP PDX RNA-seq analysis. Heat map colors indicate increased accessibility (low values; yellow, orange, red) and decreased accessibility (higher values; black) in ctDNA. TFs with increased accessibility in NEPC samples (\log_2 fold change > 0.05 , Mann-Whitney U test $P < 0.05$) are indicated with red bars; increased accessibility in ARPC (\log_2 fold change < -0.05 , $P < 0.05$) are indicated with blue bars. Text color indicates relative expression between ARPC and NEPC PDX tumors by RNA-seq shown for TFs with significant differential accessibility.

we identified 108 TFs out of the 338 that were differentially expressed between ARPC and NEPC PDX tumors by RNA-seq (Supplementary Fig. S9; Supplementary Table S3; Methods). Through unsupervised hierarchical clustering of composite TFBS central coverage values for these 108 TFs, we observed distinct groups of TFs in PDX ctDNA (Fig. 3F). Of these 108 TFs, 38 had significantly different accessibility in ctDNA between ARPC and NEPC phenotypes (two-tailed Mann-Whitney U test, Benjamini-Hochberg adjusted $P < 0.05$; Supplementary Table S3). Most of these TFs [27/38 (71%)] had differential inferred accessibility in ctDNA that was consistent with their upregulation in the same phenotype by tumor mRNA expression, although some TFs [11/38 (29%)] did not show this trend (Fig. 3F; Supplementary Fig. S10). A comparison of TFBS between paralogous TFs revealed that the binding sites used in the analysis had limited overlap (median 18.3%; range, 0–81.2%), suggesting that many of the TFs may have some independent inferred accessibility (Supplementary Fig. S11; Supplementary Table S3). For paralogs with high TFBS overlap ($\geq 19\%$), such as AR, NR3C1, and PGR, we noted only a subset of TFs were expressed in one phenotype.

REST had the largest difference in accessibility as supported by a decrease in coverage within ARPC models compared with NEPC (\log_2 fold change = -0.77 , adjusted $P = 5.7 \times 10^{-4}$; Supplementary Fig. S12A; Supplementary Table S3). FOXA1 and GRHL2 binding sites were significantly more accessible in ARPC (and ARLPC) samples compared with NEPC (\log_2 fold change < -0.57 ; adjusted $P < 1.3 \times 10^{-3}$). AR, HOXB13, and NKX3-1 had higher accessibility in ARPC compared with NEPC (\log_2 fold change < -0.37 , adjusted $P < 1.3 \times 10^{-3}$), but with only moderate accessibility in ARLPC, as expected. We also observed a group of TFs that followed a similar trend, including nuclear hormone receptors (NR2F2 and RARG), pioneer factor GATA2, and nuclear factors HNF4G and HNF1A (\log_2 fold change < -0.10 , adjusted $P < 0.027$; Supplementary Fig. S12A).

For factors that had higher accessibility in NEPC models compared with ARPC and ARLPC, ASCL1 had the largest TFBS coverage difference (\log_2 fold change 0.36; adjusted $P = 5.7 \times 10^{-4}$; Figs. 2C and 3F). Other TFs, including RUNX1, BCL11B, POU3F2, NEUROG2, and SOX2, also had sites with higher accessibility in NEPC (\log_2 fold change > 0.06 ; adjusted $P < 0.048$; Supplementary Fig. S12B), although the difference was modest. Other notable factors such as MYC and ETS transcription family genes (ETV4, ETV5, ETS1, and ETV1) had high accessibility across all phenotypes, whereas NEUROD1, RUNX3, and TP63 sites were inaccessible in nearly all samples. Furthermore, we considered restricting the analysis to 20 TFs with TFBSs that were observed in prostatic tissue and cell lines and were also differentially expressed in the PDX tumors by RNA-seq (Methods). However, although hierarchical clustering distinguished PDX tumor phenotypes, key NEPC-defining markers, such as ASCL1, were omitted from this analysis as ChIP-seq for many NEPC-defining markers had not been performed on prostate lineage samples in GTRD (Supplementary Fig. S13). Overall, we identified the accessibility of known prostate cancer regulators, including ASCL1, HNF4G, HNF1A, GATA2, and SOX2 (72–74), that have not been shown before from ctDNA analysis in these tumor phenotypes.

Phenotype-Specific Open Chromatin Regions (ATAC-Seq) in PDX Tumor Tissue Are Reflected in ctDNA Profiles of Nucleosome Accessibility

Nucleosome profiling from cfDNA sequencing analysis has shown agreement with overall chromatin accessibility in tumor tissue (38, 42, 75); however, its application for distinguishing tumor phenotypes has been limited. We hypothesized that due to lack of protection from nucleases, regions of open chromatin would be underrepresented in ctDNA assays. We investigated the use of ATAC-seq data from tumor tissue for 10 LuCaP PDX lines (5 ARPC and 5 NEPC) to inform phenotype-related differences in chromatin accessibility (9). We defined an initial set of 28,765 ARPC and 21,963 NEPC differential consensus open chromatin regions that we further restricted to those that overlapped TFBSs for 338 TFs, resulting in 15,879 ARPC and 11,692 NEPC sites (Methods; Fig. 4A). For ARPC-specific open chromatin sites, we observed decreased overall composite site coverage (± 1 kb window) and central coverage (± 30 bp) in the ctDNA for ARPC PDX lines (mean central coverage 0.75, $n = 16$) compared with NEPC lines (mean 0.96, $n = 6$) and cfDNA from healthy human donors (mean 0.97, $n = 14$; Fig. 4B; Supplementary Table S3, Methods). Conversely, for NEPC-specific open chromatin sites, coverage was decreased in ctDNA for NEPC lines (mean 0.89) compared with ARPC lines (mean 1.01) and healthy donor cfDNA (mean 1.00; Fig. 4C; Supplementary Table S3). Coverage patterns were discernible between phenotypes for as few as 100 sites, suggesting that even a smaller subset of open chromatin regions may still be informative (Supplementary Fig. S14A and S14B). These results confirmed that tumor tissue chromatin accessibility can be corroborated in ctDNA and that ARPC and NEPC phenotypes have distinct ctDNA coverage profiles at these sites.

Comprehensive Evaluation of ctDNA Features across Genomic Contexts for CRPC Phenotype Classification

To assess the utility of ctDNA nucleosome profiling for informing prostate cancer phenotype classification, we systematically evaluated four groups of global genome-wide ctDNA features: phasing, fragment sizes, local coverage profiling, and composite site coverage profiling (Fig. 1A). From PCA, we observed distinct feature signals between ARPC and NEPC phenotypes for composite TFBS coverage of TFs, NPS of the 47 phenotype marker genes, and fragment size variability at global sites of PTMs (Fig. 4D; Supplementary Fig. S15A; Supplementary Table S4; Methods). In addition to these features, we also included previously reported features, including short-long fragment ratio and local coverage patterns at the TSS (max wave height between -120 bp and 195 bp; refs. 30, 41; Methods).

We then quantitatively evaluated all combinations of coverage, phasing, and fragment size features for different genomic contexts to investigate their potential to classify ARPC and NEPC phenotypes. For each feature set, we conducted 100 iterations of stratified cross-validation using a supervised machine-learning classifier (XGBoost) on ctDNA samples from the ARPC and NEPC models and computed the area under the receiver operating characteristic (AUC) curve (Methods). First, we evaluated an established set of 10 genes associated with AR activity (5, 12). We observed that the

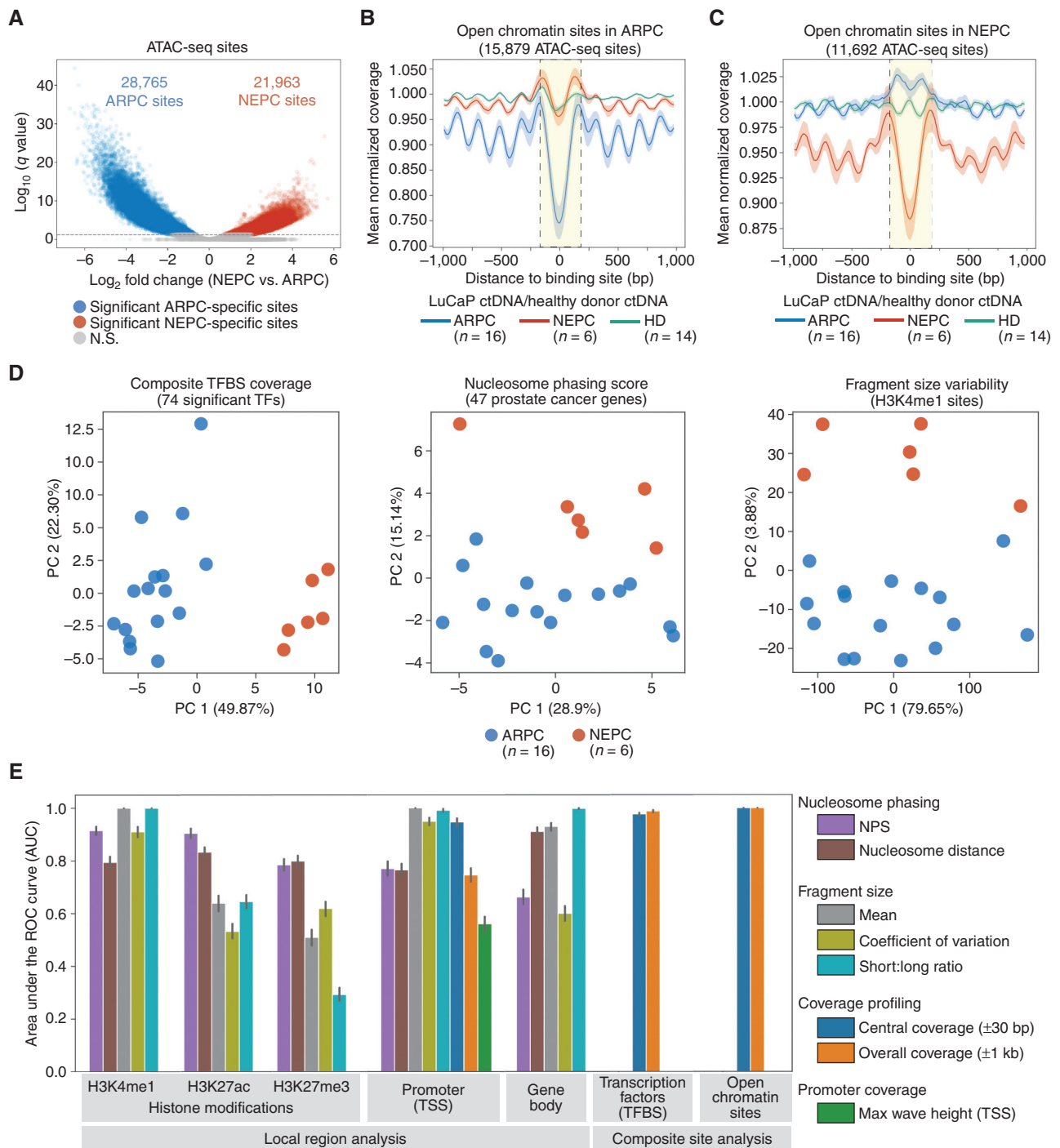
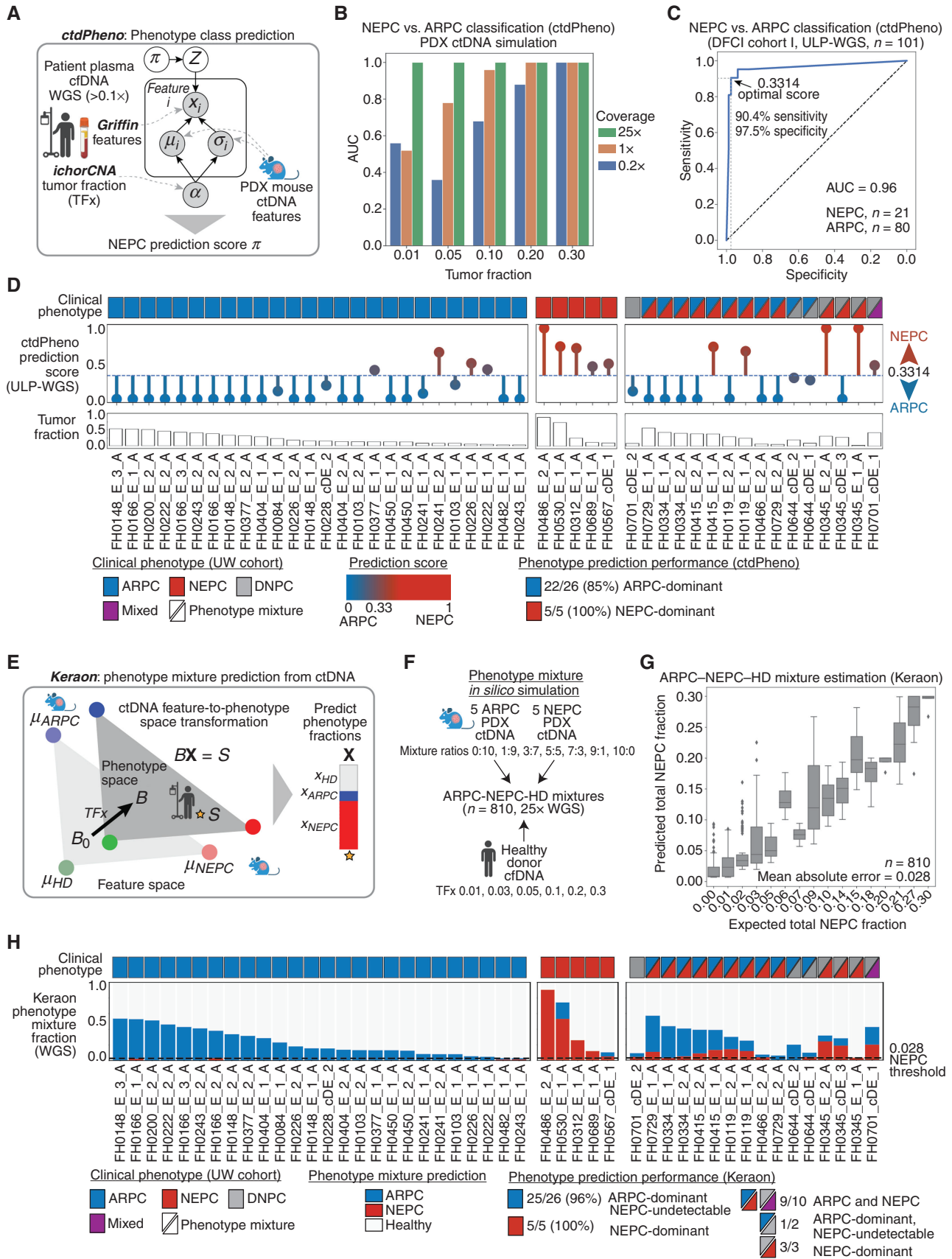


Figure 4. Comprehensive evaluation of ctDNA features throughout the genome for CRPC phenotype classification in PDX models. **A**, Volcano plot of log_2 fold change of ATAC-seq peak intensity between 5 ARPC and 5 NEPC lines; the dotted line demarcates sites by q -value < 0.05 . **B** and **C**, Composite coverage profiles at open chromatin sites specific to ARPC (**B**) and NEPC (**C**) PDX tumors analyzed by Griffin. Sites from **A** were filtered for overlap with known TFBSs in 338 factors from GTRD (71). Coverage profile means (lines) and 95% confidence interval with 1,000 bootstraps (shading) are shown. The region ± 150 bp is indicated with a vertical dotted line and yellow shading. **D**, PCAs of ctDNA features demonstrates grouping between ARPC and NEPC phenotypes. Left, composite central coverage of TFBSs significant for 74 TFs with differential accessibility out of 338 factors between ARPC and NEPC (Supplementary Table S4). Center, NPS in the gene bodies of the 47 phenotype-defining genes. Right, fragment size variability (CV) at H3K4me1 histone modification sites ($n = 9,750$). **E**, Performance of classifying ARPC vs. NEPC PDX from ctDNA using supervised machine learning (XGBoost) in various region types (all genes, TFBSs, and open regions; Methods). Area under the ROC curve (AUC) with 95% confidence interval (100 repeats of stratified cross-validation) is shown for the performance of all feature types.



phased nucleosome distance at H3K27ac sites and the central coverage at TSSs had moderate predictive performance (AUC 0.88; Supplementary Fig. S15B; Supplementary Table S4). For the set of 47 phenotype markers, the NPS of gene bodies was most predictive (AUC 0.98; Supplementary Fig. S15C; Supplementary Table S4). When considering all PTM sites, promoters, genes, TFs, and open chromatin regions, the best-performing features included mean fragment size at H3K4me1 sites ($n = 9,750$, AUC 1.0) and promoter TSSs ($n = 17,946$, AUC 1.0), and both open chromatin composite site features (AUC 1.0; Fig. 4E; Supplementary Table S4).

Accurate Classification of ARPC and NEPC Phenotypes from Patient Plasma Using a Probabilistic Model Informed by PDX ctDNA Analysis

An important consideration and challenge in analyzing plasma from patients is the presence of cfDNA released by hematopoietic cells, which leads to a lower ctDNA fraction (i.e., tumor fraction). Furthermore, the small patient cohorts with available tumor phenotype information make supervised machine-learning approaches suboptimal. Therefore, we developed ctdPheno, a probabilistic model to classify ARPC and NEPC from an individual plasma sample, accounting for the tumor fraction (Fig. 5A; Methods). We focused on the phenotype-specific open chromatin composite site features and used the PDX plasma ctDNA signals (27,571 total sites, Fig. 4B and C; Supplementary Table S3) to inform the model. The model produces a normalized prediction score that represents the estimated signature of ARPC (lower values) and NEPC (higher values). We applied this method to benchmarking datasets generated by simulating varying tumor fractions and sequencing coverages using five ARPC and NEPC PDX ctDNA samples each, and healthy donor plasma cfDNA (Supplementary Fig. S15D; Methods). We achieved a 1.0 AUC at 25 \times coverage down to 0.01 tumor fraction, 1.0 AUC at 1 \times down to 0.2 tumor fraction, and 1.0 AUC at 0.2 \times coverage at 0.3 tumor fraction, suggesting a possible upper-bound performance for classifying samples with lower tumor fraction in plasma (Fig. 5B; Supplementary Table S4).

To test the performance of ctdPheno on patient samples, we analyzed a published dataset of ultra-low-pass WGS (ULP-WGS) of plasma cfDNA (mean coverage 0.52 \times , range, 0.28–0.92 \times) from 101 patients with mCRPC, comprising 80 adenocarcinoma (ARPC) and 21 NEPC samples [Dana-Farber Cancer Institute (DFCI) cohort I; ref. 25]. Using ctdPheno,

which was unsupervised and used parameters informed only by the PDX analysis, we achieved an overall AUC of 0.96 (Fig. 5C; Supplementary Table S5). The performance was 0.97 AUC and 0.76 AUC when considering samples with high (≥ 0.1) and low (< 0.1) tumor fractions, respectively, and 0.83 AUC when using only 2,000 sites for analysis (Supplementary Fig. S16A and S16B). We identified an optimal overall performance at 97.5% specificity (ARPC) and 90.4% sensitivity (NEPC), which corresponded to the prediction score of 0.3314 (Fig. 5C). These results were concordant (92.1%) with phenotype classification by cfDNA methylation on the same plasma samples (Supplementary Fig. S16C; Supplementary Table S5). In another published dataset of 11 mCRPC samples from 6 patients who had high prostate-specific antigen (PSA), treatment with ARSI, or both (DFCI cohort II; refs. 76, 77), the model correctly classified patients as ARPC in 8 (73%) ULP-WGS ($\sim 0.1\times$) samples when using the optimal score cutoff (Supplementary Fig. S16D; Supplementary Table S5).

Next, we analyzed 61 clinical plasma samples from 31 patients with CRPC with ARPC, NEPC, and mixed phenotypes that are representative of typical clinical histories [University of Washington (UW) Cohort; Supplementary Table S5]. We performed ULP-WGS of cfDNA and selected 47 samples (26 ARPC, 5 NEPC, and 16 mixed phenotype) from 27 patients based on having greater than 3% estimated tumor fraction (Supplementary Table S5; Methods). For the 26 samples with ARPC clinical phenotype, ctdPheno correctly classified 22 (85%) samples with ARPC-dominant clinical phenotype and all five (100%) samples with NEPC-dominant clinical phenotype using the score cutoff of 0.3314 (Fig. 5D). For the remaining 16 samples with clinical histories or tumor histologies that reflected mixed phenotypes such as a tumor with AR-positive adenocarcinoma intermixed with NEPC, the classification results were variable (Fig. 5D; Supplementary Table S5; Supplementary Fig. S17). Overall, we achieved an accuracy of 87% for ULP-WGS data of ctDNA samples with dominant clinical phenotypes, but the variable predictions for mixed-phenotype samples underscore the complexities associated with tumor heterogeneity in the setting of metastatic disease.

Quantifying ARPC and NEPC Phenotype Heterogeneity within Individual Patient Plasma ctDNA

Phenotype heterogeneity may arise in the clinical setting, particularly when transdifferentiation can lead to a mixture of ARPC and NEPC cells or lesions. To account for and

Figure 5. Accurate classification and estimation of prostate cancer in patient plasma samples. **A**, Schematic illustration of the ctdPheno classification method. Griffin-derived features and ichorCNA tumor fraction estimates from patient plasma samples are combined in a probabilistic framework informed by PDX models to predict the presence of NEPC. **B**, Performance for classification on admixtures samples using ctdPheno. Five ctDNA admixtures were generated for each phenotype from PDX lines, each at various sequencing coverages and tumor fractions. In total, 125 admixtures were evaluated. The mean AUC across the 5 admixtures is shown for each configuration. **C**, ROC curve for 101 patients with mCRPC (DFCI cohort I) with ULP-WGS data. The optimal performance of 90.4% sensitivity (for predicting NEPC) and 97.5% specificity (for predicting ARPC) corresponding to a prediction score cutoff of 0.3314 is indicated with horizontal and vertical dotted lines, respectively. **D**, Prediction scores from ctdPheno for 47 ULP-WGS plasma samples with clinical phenotypes comprising 26 ARPC (blue), 5 NEPC (red), and 16 mixed or ambiguous phenotypes (purple, triangles), including DNPC (gray). The 0.3314 score cutoff threshold (dotted line) was used for classifying NEPC and ARPC. Tumor fractions were estimated by ichorCNA from WGS data. **E**, Schematic illustration of the Keraon mixture estimation method. Griffin-derived features from PDX lines and healthy donors define a known feature space, which is transformed based on Griffin features and ichorCNA tumor fraction estimates for each patient plasma sample. Based on the patient's location in the transformed phenotype space, fractions of each phenotype are inferred directly. **F**, Illustration of mixture simulations. Five ARPC and five NEPC PDX samples were combined in the ratios shown with a single healthy donor at the tumor fractions shown, for a total of 810 mixed-phenotype samples at 25 \times for evaluating mixture proportions with Keraon. **G**, Boxplot of predicted total NEPC fraction in 810 simulated mixed-phenotype samples using Keraon, Pearson $r = 0.884$. MAE was computed as the median absolute difference between estimated and expected NEPC fraction across all samples. **H**, Fractional phenotype estimates for 47 WGS plasma samples with clinical phenotypes comprising 26 ARPC (blue), 5 NEPC (red), and 16 mixed or ambiguous phenotypes (purple, triangles), including DNPC; gray). The 2.8% NEPC fraction threshold indicates the predicted presence of NEPC (dotted line).

predict phenotype mixtures within a patient ctDNA sample, we developed Keraon, an analytic model that estimates the proportions of phenotypes from WGS using the same ctDNA features as ctdPheno (Fig. 5E; Methods). First, we evaluated Keraon using a benchmark dataset generated for simulating varying tumor fractions and proportions of ARPC–NEPC mixtures at 25× coverage using PDX ctDNA and healthy donor cfDNA data (Fig. 5F; Methods). In 810 simulated phenotype mixtures, we observed the estimated total NEPC fraction was consistent with expected proportions (Pearson $r = 0.884$) with a mean absolute error (MAE) of 0.028, highlighting the method's potential for accurate estimation of emergent phenotypes in mixed histology samples (Fig. 5G; Supplementary Table S5). Next, we evaluated Keraon for classifying NEPC in DFCI cohort I and observed the highest performance (0.96 AUC) using all 27,571 open chromatin sites, with decreased performance (0.84 AUC) when using only 2,000 sites (Supplementary Fig. S16D). Applying Keraon to analyze DFCI cohort II, we correctly estimated dominant ARPC with undetectable NEPC phenotype in 10 (91%) samples with WGS (mean coverage, 27×; Supplementary Fig. S18; Supplementary Table S5).

We performed deeper WGS (22.13× mean coverage; range, 15.15×–31.79×) for the UW cohort ctDNA samples and applied Keraon to classify the presence of NEPC and to estimate the proportions of ARPC and NEPC phenotypes (Fig. 5H). Keraon correctly estimated the dominant phenotype (≥ 0.5 relative phenotype fraction) in 25 of 26 (96%) samples with ARPC clinical phenotype and in 5 of 5 (100%) NEPC samples. For 10 samples with the presence of ARPC and NEPC phenotypes reported in the clinical histories, Keraon correctly detected both phenotypes in nine samples (NEPC fraction ≥ 0.028 , ARPC fraction ≥ 0.06). In two samples with ARPC–DNPC phenotypes, one was estimated to be ARPC-dominant (0.20 fraction), and in three samples with NEPC–DNPC phenotypes, all three were estimated as being NEPC-dominant (≥ 0.028 fraction). In 14 (82%) out of 17 patients with multiple plasma collected, the predicted phenotypes were consistent across all ctDNA samples. Overall, we observed an accuracy of 97% for correctly classifying ARPC- and NEPC-dominant phenotypes and 87% for estimating NEPC fractions in samples with admixed clinical phenotypes from ctDNA.

DISCUSSION

The development of minimally invasive blood-based assays of ctDNA to define tumor subtypes has dramatically changed the landscape of clinical oncology. To date, the majority of these assays characterize genomic alterations in oncogenes such as *EGFR* or tumor suppressors, such as *BRCA2*, that inform outlier responses to specific therapeutics. However, tumor classification determined by gene-expression analyses, such as the PAM50 subtyping of breast carcinoma and the transcript-based classification of urothelial cancers, is also informative of clinical trajectories. Consequently, the ability to characterize tumor phenotype using blood-based assays has the potential to add relevant information for guiding treatment allocation.

In the present study, we analyzed multiple features of DNA to infer the activity of gene-expression programs corresponding to distinct prostate cancer phenotypes. A key component of the work that allowed for the development of optimized methods and the identification of the most informative ctDNA

features was the use of PDX models. The sequencing of mouse plasma provided a unique opportunity to comprehensively interrogate the epigenetic nucleosome patterns in ctDNA from well-characterized tumor models. We developed and applied computational methodologies to evaluate a multitude of ctDNA features, each of which was associated with transcriptional regulation across CRPC tumor phenotypes. The use of PDX mouse plasma overcomes the challenge of low ctDNA content or incomplete knowledge of the tumor when studying patient samples. Using features learned from the PDX ctDNA, we developed models to accurately classify ARPC and NEPC and to estimate their proportions in phenotypically heterogeneous samples from patient plasma in three clinical cohorts. Although these data were focused on ARPC and NEPC phenotypes, the approaches may serve as a framework for the use of ctDNA to subtype malignancies arising in other organ sites based on distinctive gene-expression programs.

The analysis of the LuCaP PDX ctDNA sequencing data confirmed the activity of key regulators between ARPC and NEPC phenotypes, including a set of 47 established differentially expressed genes that associate with cell lineage. Although gene-expression inference from ctDNA has been shown in proof-of-concept studies (34, 41), the PDX ctDNA allowed for a detailed dissection of nucleosome organization associated with the transcriptional activity of individual genes that define the tumor phenotypes. Previous analytic approaches have profiled nucleosome occupancy from cfDNA (38, 75). However, our assessment of nucleosome stability by means of the NPS is the first to capture the highly variable spacing, positioning, and turnover of the nucleosome arrays associated with transcription and tumor aggressiveness (44, 68, 69, 78).

In addition to the existing molecular profiling available for these models, we now provide a characterization of histone PTMs in LuCaP PDX tumors using CUT&RUN. At regions with these PTMs on histone tails, we observed expected nucleosome patterns inferred in ctDNA that were consistent with active or repressed gene transcription. To our knowledge, this is the first time that ctDNA analysis has been performed in the context of histone PTMs and will provide a blueprint to develop new approaches for studying additional epigenetic alterations using PDX plasma.

Although the regulation of key factors such as AR, HOXB13, NKX-3.1, FOXA1, and REST has been shown from ctDNA in CRPC (35, 42), we report the differential activity of other key factors in CRPC from ctDNA analysis. This included nuclear factors HNF4G and HNF1A and pioneering factor GATA2, which are associated with prostate adenocarcinoma (ARPC; refs. 72, 74, 79). ASCL1 is a pioneer TF with roles in neuronal differentiation and was recently described to be active during NE transdifferentiation and in NEPC (9, 56). To our knowledge, this study is the first to demonstrate ASCL1 binding site accessibility and provide a detailed characterization of its transcriptional activity in NEPC from plasma ctDNA.

We show an expansive analysis of TFBSs for 338 factors in each plasma sample without the need for chromatin immunoprecipitation or other epigenetic assays. However, we did not find a significant difference in accessibility for 70 out of the 108 TFs in ctDNA, which may be consistent with TF activity not necessarily being correlated with its own expression level (80). On the other hand, the accessibility of TFBSs may not necessarily indicate true

TF activity, as other cobound TFs or coactivators/corepressors influence gene regulation. Moreover, our analyses were based on TFBSs obtained from public databases, including for a limited number of prostate-specific TFs; however, expanded phenotype-specific TF cistrome data may improve this approach.

We applied state-of-the-art computational approaches built on existing and new concepts of ctDNA data analysis to extract tumor-specific features, including the representation of nucleosome phasing, periodicity, and spacing associated with transcriptional activity. Other approaches have also considered regions, such as TSSs, TFBSs, and DNase hypersensitivity sites (33, 38, 41, 42); however, after a systematic evaluation, we found that ctDNA features in open chromatin sites derived from ATAC-seq of PDX tissue (9) provided the highest performance for distinguishing CRPC phenotypes. We presented ctdPheno, which is a probabilistic model that classifies ARPC and NEPC from ULP-WGS data, and Keraon, an analytic model that estimates the proportion of ARPC and NEPC from WGS of patient plasma. Both models are unsupervised and utilize a statistical framework informed directly by parameters from the LuCaP PDX ctDNA analysis. These models do not require training on patient samples but do require tumor fraction estimates (ichorCNA; ref. 81) and in the case of ctdPheno a prediction score cutoff determined from DFCI cohort I. Both frameworks can also be extended to model additional phenotypes. Insights from additional datasets such as single-cell nucleosome and accessibility profiling (82, 83) of PDX tumors and clinical samples may improve the resolution for ctDNA analysis. Although we observed optimal performance analyzing all open chromatin sites, a smaller subset was still informative, which may be useful when considering targeted assays for clinical applications.

Applying the prediction models to patient datasets with definitive clinical phenotypes yielded high performance even when using low depth of coverage sequencing. In particular, our performance for the DFCI cohort I was also consistent with the reported phenotype classification results using ctDNA methylation in the same patients (25). Similarly, in the UW cohort, samples with well-defined clinical phenotypes had near-perfect concordance from WGS data. We established the lower limits of phenotype classification performance to be at 8% tumor fraction for ctdPheno (ULP-WGS) and 3% for Keraon (WGS). These results support a strategy whereby ULP-WGS is performed for screening using ctdPheno, along with clinical assessments, and followed up with standard WGS for more accurate and comprehensive phenotype characterization using Keraon. Although this framework may have limited performance for low (<3%) ctDNA levels, it may be optimal at the initial assessment of metastatic disease and at tumor progression on therapy, which is when the clinical decision points are most critical.

Tumor heterogeneity and coexistence of different molecular phenotypes are common in mCRPC where treatment-induced phenotypic plasticity may vary within and between tumors in an individual patient. In real data simulations and patients with cases of mixed clinical phenotypes, Keraon accurately detected the contributions of mixed phenotypes with a detection limit of 2.8% NEPC, providing the first approach to directly quantify phenotype proportions and heterogeneity from ctDNA. In this study, estimation of phenotype heterogeneity using Keraon required standard depths

of WGS. Larger studies with a comprehensive assessment of the tumor histologies will be needed for evaluating these models as potential biomarkers of treatment response.

In summary, this study illustrates that analysis of ctDNA from PDX mouse plasma at scale can facilitate a detailed investigation of tumor regulation. These results, together with the suite of computational methods presented here, highlight the utility of ctDNA for surveying transcriptional regulation of tumor phenotypes and its potential diagnostic applications in cancer precision medicine.

METHODS

PDX Mouse Models

The establishment and characterization of the LuCaP PDX models were described previously (84). PDXs were propagated *in vivo* in male NOD/SCID IL2R-gamma-null (NSG) mice (cat. #005557). The collection of tumors for the establishment of PDX lines was approved by the UW Human Subjects Division Institutional Review Board (IRB #2341). PDX lines were evaluated using histopathology by at least two expert pathologists, and histologic phenotypic subtype annotations were orthogonally validated based on transcriptome-derived signature marker expression scores to define phenotypes (4, 5, 22): ARPC, NEPC, and ARLPC. Resected PDX tumors (300–800 mm³) were divided into ~50 mg to ~100 mg pieces and stored at –80°C. Animal studies were approved by the Fred Hutchinson Cancer Center (FHCC) Institutional Animal Care and Use Committee (protocol 1618) and performed in accordance with the NIH guidelines. For the current study, blood was collected by cardiac puncture from animals bearing PDX tumors (measurable size 300–1,400 mm³).

Human Subjects

UW Cohort. Blood samples were collected from men with mCRPC at the UW (collected under UW Human Subjects Division IRB protocol number CC6932 between years 2014 and 2021). Patients in this study have provided written informed consent for research participation. In this study, 61 plasma samples from 31 patients were analyzed. After initial ULP-WGS analysis, 47 plasma samples from 27 patients with sufficient tumor fraction (>3%, based on initial ichorCNA analysis using GRCh37 genome build) and three additional samples not meeting the threshold but with clear AR amplification seen in manual curation (FH0243_E_1_A, FH0345_E_1_A, and FH0482_E_1_A) were retained for further high depth of coverage WGS analysis. All samples were deidentified prior to ctDNA analysis and we used a double-blinded approach for evaluating clinical phenotype predictions.

DFCI Cohort I. Plasma was collected from men diagnosed with mCRPC and treated at the DFCI, Brigham and Women's Hospital, or Weill Cornell Medicine (WCM) between April 2003 and August 2021. All patients provided written informed consent for research participation and genomic analysis of their biospecimen and blood. The use of samples was approved by the DFCI IRB (#01-045 and 09-171) and WCM (1305013903) IRBs. The ULP-WGS data at mean coverage 0.5× (range 0.3×–0.9×) for 101 patients were published previously (25).

DFCI Cohort II. Plasma samples in this cohort were collected from men diagnosed with mCRPC and treated at the DFCI. All patients provided written informed consent for blood collection and the analysis of their clinical and genetic data for research purposes (DFCI protocol #01-045 and 11-104). WGS data at mean coverage 27× (range, 11×–44×; ref. 76), and ULP-WGS data at mean coverage 0.13× (range, 0.07×–0.18×; refs. 77, 81) were downloaded from dbGAP accession phs001417. Eleven samples from six patients had matching WGS and ULP-WGS with paired-end reads, necessary for

analysis by Griffin. PSA (ng/mL) values and treatment at the time of the blood draw were previously published (77). The six patients were treated for adenocarcinoma using abiraterone, enzalutamide, or bicalutamide, or the patients had detectable levels of PSA.

Healthy donor plasma cfDNA WGS data used in this study were obtained from previously published studies. Two samples (HD45 and HD46, both male) with coverage of 13× and 15×, respectively, were accessed from dbGAP under accession phs001417 (76, 81). These donors were consented under DFCI protocol IRB (# 03-022). Plasma cfDNA WGS data from thirteen healthy donors [12 male: NPH002, 03, 06, 07, 12, 18, 23, 26, 33, 34, 35, 36; 1 female (used in admixtures): NPH004] with coverages between 13.5× and 27.6× were obtained from the European Phenome Archive (EGA) under accession EGAD00001005343 (42).

PDX Plasma Processing

Blood samples were collected from NSG mice bearing subcutaneous PDX tumors at the time of sacrifice. The PDX lines were maintained at vivaria in the UW and FHCC. The blood was processed following methods described for human plasma DNA processing for subsequent DNA isolation. Blood was collected in Sarstedt Micro sample tube K3 EDTA tubes and processed within 4 hours. All blood samples were sequentially double spun, first at 2,500 × g for 10 minutes followed by a 16,000 × g centrifugation of the plasma fraction for 10 minutes at room temperature. For each PDX line, 4 to 8 mouse plasma samples were pooled. Processed plasma samples were preserved in clean, screw-capped cryo-microfuge tubes and stored at −80°C prior to cfDNA isolation.

cfDNA Isolation

The QIAamp Circulating Nucleic Acid Kit was used to isolate cfDNA from PDX mouse-derived plasma using the recommended protocol. The pooled plasma samples from 4 to 8 mice for each PDX line contained 1.9 to 3 mL total plasma volume for each line. The filter retention-based cfDNA kit method does not implement any fragment size class enrichment. Isolated cfDNA was quantified using the Qubit dsDNA HS assay (Invitrogen) and the cfDNA fragment size profiles were analyzed using TapeStation HS D5000 and HS D1000 assays (Agilent).

cfDNA Library Preparation and Sequencing

For LuCaP PDX mouse plasma samples, NGS libraries were prepared with 50 ng input cfDNA. Illumina NGS sequencing libraries were prepared with the KAPA hyperprep kit, adopting nine cycles of amplification, and purified using lab-standardized SPRI beads. We used KAPA UDI dual-indexed library adapters. Library concentrations were balanced and pooled for multiplexing and sequenced using the Illumina HiSeq 2500 at the Fred Hutchinson Genomics Shared Resources (200 cycles) and Illumina NovaSeq platform at the Broad Institute Genomics Platform Walkup-Seq Services using S4 flow cells (300 cycles). To match with Illumina HiSeq 2500 data, truncated 200 cycles FASTQ files were generated (100 bp paired-end reads).

Clinical patient plasma samples collected at the University of Washington (UW cohort) were submitted to the Broad Institute Blood Biopsy Services. Briefly, cfDNA was extracted from 2 mL double-spun plasma and ULP-WGS to approximately 0.2× coverage was performed. The ichorCNA pipeline was used to estimate tumor DNA content (i.e., tumor fraction; see below). Forty-seven samples (from 31 patients) had either ≥5% tumor fraction or ≥2% tumor fraction with AR amplification observed in ichorCNA and were subsequently sequenced to deeper WGS coverage (~20×).

cfDNA Sequencing Analysis and Mouse Subtraction

All cfDNA sequencing data used in this study were realigned to the hg38 (GRCh38) human reference genome (<http://hgdownload.soe.ucsc.edu/goldenPath/hg38/bigZips/hg38.fa.gz>). FASTQ files

were realigned using BWA (v0.7.17) mem (85). The complete alignment pipeline including configuration settings may be accessed at https://github.com/GavinHaLab/fastq_to_bam_paired_snakemake.

For PDX ctDNA WGS data, we performed mouse genome subtraction following the protocol described previously (86), wherein reads were aligned using BWA mem to a concatenated reference consisting of both human (hg38) and mouse (mm10, GRCm38.p6, http://igenomes.illumina.com.s3-website-us-east-1.amazonaws.com/Mus_musculus/NCBI/GRCm38/Mus_musculus_NCBI_GRCm38.tar.gz) reference genomes. Read pairs where both reads aligned to the human reference genome were retained and all other read pairs were removed. Then, the remaining reads were realigned to the human-only reference. Finally, the GATK best practices workflow was applied to each sample (87); the complete mouse subtraction pipeline used in this study, including tool versions and parameters, can be accessed at https://github.com/GavinHaLab/PDX_mouseSubtraction. Following mouse subtraction, samples with <3× depth were removed for downstream analysis.

CCP Score Calculation

The 31-gene CCP signature (69) was computed from RNA-seq data using GSVA (88). The single-sample enrichment scores were calculated with default parameters using genome-wide log₂ FPKM values as input for the 31 genes.

Differential mRNA Expression Analysis

RNA isolation of 102 tumors from 46 LuCaP PDX samples was performed as described previously (11). RNA concentration, purity, and integrity were assessed by NanoDrop (Thermo Fisher Scientific Inc) and Agilent TapeStation, and RNA integrity number (RIN) ≥ 8 was retained for library preparation. RNA-seq libraries were constructed from 1 µg of total RNA using the Illumina TruSeq Stranded mRNA LT Sample Prep Kit according to the manufacturer's protocol. Barcoded libraries were pooled and sequenced by Illumina NovaSeq 6000 or Illumina HiSeq 2500, generating 50 bp paired-end reads. Sequencing reads were mapped to the hg38 human reference genome and mm10 mouse reference genomes using STAR.v2.7.3a (89). All subsequent analyses were performed in R-4.1.0. Sequences aligning to the mouse genome and therefore derived from potential contamination with mouse tissue were removed from the analysis using XenofilteR (v1.6; ref. 90). Gene-level abundance was quantitated using the R package GenomicAlignments v1.32.0 summarizeOverlaps function using mode = IntersectionStrict, restricted to primary aligned reads. We used refSeq gene annotations for transcriptome analysis. Transcript abundances (FPKM) were input to edgeR v3.38.1 (91), filtered for a minimum expression level using the filterByExpr function with default parameters, and then limma v3.52.1 voom was used for differential expression analysis of NEPC versus ARPC and ARLPC versus ARPC. We then filtered the results using a list of 1,635 human TFs published previously (92), which resulted in 514 genes with FDR < 0.05 and log₂ fold change > 1.58. Out of these 514, deregulation of gene expression for 404 TF genes delineated ARPC from NEPC.

CUT&RUN

CUT&RUN is an antibody-targeted enzyme tethering chromatin profiling assay in which controlled cleavage by micrococcal nuclease releases specific protein-DNA complexes into the supernatant for paired-end DNA sequencing analysis. We performed CUT&RUN assays for three histone modifications, H3K27ac, H3K4me1, and H3K27me3, according to published protocols (49). We performed CUT&RUN on LuCaP PDX tumors using ~75 mg flash-frozen tissue pieces.

Paired-end (50 bp) sequencing was performed and reads were aligned using bowtie2 v2.4.2 (93) to the hg38 human reference assembly. Aligned reads were processed as described in the SEACR protocol (<https://github.com/FredHutch/SEACR#preparing-input-bedgraph-files>).

Peaks were called using SEACR version 1.3 (50) using “stringent” settings and with reference to paired IgG controls. BigWig files were prepared using bamCoverage in deepTools 3.5.0 (94). Genome-wide peak heat map, targeted heat map, and respective profiles were plotted using deepTools v3.5.0. bigwig formatted files for each phenotype were obtained using the mean function in wiggletools 1.2.8. and deepTools computeMatrix. Phenotype-specific informative region coordinates were obtained from diffBind v3.5.0, and the top 10,000 most significant regions (all with FDR < 0.05) differentially open between ARPC and NEPC lines were used for downstream feature analyses (see Gene body and promoter region selection for additional subsetting criteria applied on a feature-by-feature basis). For heat maps and profiles, the plotHeatmap function was used. We utilized the “Peak Center” option to derive desired heat maps. These steps were all performed for H3K27ac, H3K4me1, and H3K27me3 antibodies. Scaled heat map profiles’ AUC (± 1.5 kbp) and peak height at the profile center were estimated using deepStats v0.4 (<https://zenodo.org/record/3668336>; comparable profiles are scaled to 10 units).

Differential Histone PTM Analysis

Differential PTM analysis was performed with the DiffBind version 2.16.0 package (95) in R-4.0.1 using standard parameters (<https://bioconductor.org/packages/3.0/bioc/html/DiffBind.html>). ARPC, NEPC, and ARLPC samples were grouped by histopathologic and transcriptome signature-defined phenotypes described in the “PDX mouse models” section (Supplementary Table S2). Samples were loaded with the dba function, reads counted with the dba.count function, and contrast specified as phenotype with dba.contrast and a minimum members of 2. Differential peak sites were computed with the dba.analyze function with default settings. Differential peak binding of NEPC and ARLPC was computed against ARPC samples. Unique binding sites in NEPC and ARLPC were cataloged using bedtools v2.29.2 (96). Intergroup differentially bound peaks were annotated using ChIPseeker 1.28.3 (97) and TxDb.Hsapiens.UCSC.hg38.knownGene 3.2.2 in R 4.1.0.

ATAC-seq Analysis

ATAC-seq sequence data for 15 tumor samples from 10 PDX lines were published previously (9). These lines included LuCaP PDX lines with ARPC (23.1, 77, 78, 81, and 96) and NEPC (three replicates of 173.1, two replicates each of 49, 93, 145.1, and one replicate of 145.2) phenotypes. Paired-end reads were aligned using bowtie2 v2.4.2 (93) to the UCSC hg38 human reference assembly with the “very-sensitive” “-k 10” settings. Peaks were called using Genrich version 0.6.1 (<https://github.com/jsh58/Genrich>). Differential binding analysis was performed using DiffBind version 3.5.0 package in R version 4.1.0. ENCODE blacklisted regions were excluded using hg38-blacklist.v2 (ref. 98; <https://github.com/Boyle-Lab/Blacklist>). Phenotype-specific regions were isolated by first selecting for positive fold change open chromatin enrichment and then using Intervene 0.6.5 (99), where regions were considered overlapping if they shared at least 1 bp with another phenotype. Regions with FDR adjusted $P < 0.05$ were then subset to those overlapping the 3,380,000 established TFBSs (338 TFs \times 10,000 binding sites, see Griffin analysis for site selection) by at least 1 bp using BedTools v2.30.0 Intersect. Only regions that overlapped an established TFBS from those lists were retained. For analyses restricted to 10,000, 1,000, or 100 sites, sites were ranked and chosen by adjusted P value.

Nucleosome Profiling of ctDNA

Griffin is a method for profiling nucleosome protection and accessibility on predefined genomic loci (52). For this study, Griffin (v0.1.0) was used and can be found on GitHub (<https://github.com/adoebley/Griffin/releases/tag/v0.1.0>). The analysis was performed as follows: First, GC bias was quantified for each sample using an approach described previously (100). Briefly, for each possible fragment length and GC content,

the number of reads in a bam file and the number of genomic positions with that specific length and GC content were counted. The GC bias for each fragment length and GC content was calculated by dividing the number of observed reads by the number of observed genomic positions for that fragment length and GC content. The GC bias for all possible GC contents at a given fragment length was then normalized to a mean bias of 1. GC biases were then smoothed by taking the median of values for fragments with similar lengths and GC contents (k nearest neighbors smoothing) to generate smoothed GC bias values.

After GC correction, nucleosome profiling was performed in each sample. For each mappable site of interest, fragments aligning to the region $\pm 5,000$ bp from the site were fetched from the bam file. Fragments were filtered to remove duplicates and low-quality alignments (<20 mapping quality) and by fragment length. Nucleosome size fragments (140–250 bp) were retained and used in all downstream Griffin analyses. Fragments were then GC corrected by assigning each fragment a weight of $1/\text{GC_bias}$ for that given fragment length and GC content. The fragment midpoint was identified, and the number of weighted fragment midpoints in 15 bp bins across the site was counted. For composite sites, all sites of a given type (such as all sites for a given transcription factor) were summed together to generate a single coverage profile. Individual or composite coverage profiles were normalized to a mean coverage of 1 in the $\pm 5,000$ bp region surrounding the site. Finally, sites were smoothed using a Savitzky-Golay filter with a window length of 165 bp and a polynomial order of 3. The window $\pm 1,000$ bp around the site was retained for plotting and feature extraction when plotting sites; shading illustrates the 95% confidence interval within sample groups. Features extracted from individual or composite sites included:

- “mean central coverage,” the mean coverage between -30 and 30 bp relative to the site center,
- “mean window coverage,” the mean coverage between -990 and 990 bp relative to the site center, and
- “max wave height,” the absolute difference between the minimum coverage within the window from -120 to 30 bp and maximum coverage in the window from 31 to 195 bp relative to the TSS.

TFBS Selection from GTRD

TFBS identified using ChIP-seq were downloaded from the GTRD database version 19.10 (https://grtd.biouml.org/downloads/19.10/chip-seq/Homo%20sapiens_meta_clusters.interval.gz). This database contains binding sites (meta-clusters) that were observed in one or more ChIP-seq experiments. Low mappability sites were excluded by examining the mean mappability score in a window around each site ($\pm 5,000$ bp). Mappability information (hg38 Umap multiread mappability for 50 bp reads) was obtained from the UCSC genome browser (ref. 101; <https://hgdownload.soe.ucsc.edu/gbdb/hg38/hoffmanMappability/k50.Umap.MultiTrackMappability.bw>). Highly mappable sites (>0.95 mean mappability) were retained for further analysis. 338 TFs were selected for analysis using three criteria: (i) TF was contained in GTRD, (ii) had at least 10,000 highly mappable binding sites on autosomes (chr1–22) in GTRD, and (iii) TF was present in the CIS-BP database (ref. 102; CIS-BP v2.00 downloaded from <http://cisbp.cabr.utoronto.ca/bulk.php>) and had a known binding motif (“TF_status” is not N). Unless otherwise noted, analyses utilized the top 1,000 TFBSs ranked by the highest “peak.count” across all experiments as computed by GTRD (71). In addition, in the case of AR and ASCL1, we also compared the top 1,000 with the top 10,000 sites chosen with the same “peak.count” criterion.

After intersecting these 338 TFs with the 404 differentially expressed TFs identified through RNA-seq, 108 remained. On both the 108 and prostate-specific 41 TFs (described below) we performed unsupervised hierarchical clustering of central window mean values (see Griffin analysis). Hierarchical clustering was performed using the Ward.D2 method

with Euclidean distance and complete linkage settings; the groupings were determined using `cutree_cols = 2` for columns (LuCaP CRPC phenotypes) and `cutree_rows = 13` for rows (TFs) on the dendrograms.

To generate a prostate lineage-specific TF set, we first merged GTRD metadata (file; <http://gtrd.biouml.org:8888/downloads/current/metadata/ChIP-seq.metadata.txt> and http://gtrd.biouml.org:8888/downloads/current/metadata/cell_types_and_tissues.metadata.txt). We identified human prostate lineage-specific experiments by restricting the “species” field to “Homo sapiens” and the “title” (tissue or cell type) field by performing a string match of the following {“Prostate,” “prostate,” “LNCaP,” “DU145,” “PrEC”}. This resulted in a list of 1,086 prostate lineage ChIP-seq experiments. Then, we selected metapeaks from the “Homo_sapiens_meta_clusters.interval” file that had been observed in at least one of the prostate lineage experiments using the “exp.set” field. This resulted in a set of 82 TFs. We then filtered the peaks by mappability and kept only highly mappable peaks (as described above). We excluded any TF that was not in the initial set of 338 TFs (this removed ChIP targets that were not true TFs, lacked a known binding site, or did not have 10,000 total autosomal peaks in GTRD). Of the remaining TFs, we analyzed those with 1,000 highly mappable peaks on autosomes in prostate lineage experiments, resulting in 41 TFs. Twenty out of these 41 TFs overlapped the list of 108 differentially expressed TFs by RNA-seq of the PDX tumors. Note that the top 1,000 sites for each of the 41 TFs were different than in the same TFs of the 338 set because sites must meet the criteria of being derived from at least one experiment involving prostate tissue or cell lines.

TFBS Selection from Other Sources

For AR we further considered 17,619 sites identified through ChIP-seq by Pomerantz and colleagues (ref. 57; which overlapped 10.9% of the GTRD top 1,000 using bedtools), and 41,633 sites identified by Severson and colleagues (58) across four metastatic tumors (which overlapped 99.4% of the GTRD top 1,000). For ASCL1, we obtained 11,124 ChIP-seq sites from Cejas and colleagues (ref. 9; which overlapped 60.9% of the GTRD top 1,000). All of these site lists were lifted over from genome build GRCh37 to GRCh38. No mappability filtering was applied so that all possible sites from these prostate experiments and studies were considered.

Phenotype Lineage-Specific Gene Marker Selection

We selected 47 genes comprising 12 ARPC and 35 NEPC lineage markers established previously (4, 5, 59, 60) and confirmed by differential expression analysis from PDX tumor RNA-seq data (Supplementary Table S3). In tissues, AR and NE activities were measured on lineage-determinant signature gene’s mRNA expression (GSVA score; ref. 88). The 47 selected gene list comprises the majority of these signature sets of genes defining mPC characteristic phenotypes or phenotypic activities.

Gene Body and Promoter Region Selection

For individual gene body and promoter analyses, Ensembl BioMart v104 (hg38; ref. 103) was used to directly retrieve protein-coding transcript start (TSS) and end (TES) coordinates. For promoter region analysis, the window $\pm 1,000$ bp relative to the TSS was considered. For gene body analysis, the region between the TSS and TES was considered. In the case of genes with multiple transcripts, analyses were limited to the longest transcript, resulting in 19,336 regions. In a downstream analysis of LuCaP PDX cfDNA, if any lines did not meet specific criteria in a region (including differentially open histone modification regions), that feature/region combination was excluded from analysis, leading to a variable lower number of regions considered based on the feature. These criteria included requiring at least 10 total fragments in a region for all fragment size analysis (see below) and a nonzero number of “short” and “long” fragments for the short-long

ratio; short-long ratios less than 0.01 or greater than 10.0 were also excluded as outliers. For phasing analysis (see below), we also excluded amplitude components and thus NPS where individual components were 0, or where the ratio was less than 0.01 or greater than 10.0, indicative of insufficient coverage. In the case of mean phased nucleosome distance, if no peaks were identified or the value in a region exceeded 500 (indicative of highly irregular/sparse pileups also from low coverage), those regions were also excluded. Any region with no coverage in a line was excluded from all analyses. This resulted in gene lists that differed in numbers between genomic contexts and feature types.

cfDNA Fragment Size Analysis

Fragments were first filtered to remove duplicates and low-quality alignments (<20 mapping quality) and by fragment length (15–500 bp). In individual genomic loci/windows, we computed the fragment short-long ratio (FSLR) as the ratio of short (15–120 bp) to long (140–250 bp) fragments. We also calculated the mean, median absolute deviation (MAD; $\text{median}(|X_i - \text{median}(X)|)$), and coefficient of variation (CV; σ/μ , where σ = standard deviation, μ = mean) of the fragment length distribution for each selected window. The fragment size analysis code and implementation used in this study can be accessed at <https://github.com/GavinHaLab/CRPCSubtypingPaper/tree/main/FragmentAnalysis>.

Nucleosome Phasing Analysis (TritonNP)

Fragments were first filtered to remove duplicates and low-quality alignments (<20 mapping quality) and by fragment length (nucleosome-sized: 140–250 bp). Next, we performed fragment-level GC bias correction utilizing the same preprocessing method defined in Griffin. A band-pass filter was then applied to the corrected coverage in each region of interest by taking the fast Fourier transform (FFT; `scipy.fft v1.8.0`; ref. 104) and removing high-frequency components corresponding to frequency components <146 bp before reconstructing the signal. This cutoff was chosen to ensure that periodic fit signal for downstream evaluation must come from the minimum possible internucleosome distance, thus excluding peak pileups that would not indicate an overall trend in nucleosome phasing. Local peak calling was then done on the smoothed signal to infer the average internucleosome distance or “phased nucleosome distance” by finding maxima directly. To quantify the clarity of overall phasing, we took the average frequency amplitude in two bands corresponding to stably bound, well-phased nucleosomes (180–210 bp) and a baseline (150–180 bp), with the former measuring the strength of typically aligned nucleosomes and the latter giving a measure of the underlying signal strength not coming from either high-frequency noise or low-frequency shifts in total coverage. The ratio of these two amplitude averages forms the NPS. Because peak locations are assumed to be independent of copy-number alterations (CNA) or depth, and the NPS by virtue of being a ratio divides out any confounding DNA/depth variation between sites, both features are taken as agnostic of CNAs or variable depth. Code and implementation of the method can be found at <https://github.com/denniepatton/TritonNP>.

ctDNA Tumor–Normal Admixtures and Benchmarking

Admixtures for evaluating benchmarking performance were constructed using 5 ARPC (LuCaP 35, 35CR, 58, 92, and 136CR) and 5 NEPC (LuCaP 49, 93, 145.2, 173.1, and 208.4) lines mixed to 1%, 5%, 10%, 20%, and 30% tumor fraction with a single healthy donor plasma line (NPH023, EGAD00001005343) for use in binary classification (50 admixes), and in mixtures of 1%, 3%, 5%, 10%, 20%, and 30% tumor fraction at ARPC:NEPC ratios, of 0.0, 0.1, 0.3, 0.5, 0.7, 0.9, and 1.0 in all possible combinations (810 admixes) for mixture model evaluation. All admixes were mixed at $\sim 25\times$ mean coverage, assuming 100% tumor fraction in post-mouse subtracted

PDX sequencing data. After extracting chromosomal DNA (chr1–22, X, Y) with SAMtools v1.14 (105) and removing duplicates with Picard (<https://broadinstitute.github.io/picard/>), SAMtools was used to merge BAM files. To evaluate the ULP-WGS performance, admixtures were then downsampled using SAMtools to the number of reads corresponding to 1× or 0.2×. During unsupervised benchmarking of each admixture, the healthy donor and the LuCaP line used in the admixture were excluded from the generation of feature distributions to ensure the model would not learn from the lines being interrogated. The admixture pipeline used in this study can be accessed at https://github.com/GavinHaLab/Admixtures_snakemake.

Supervised Binary Classification of ARPC and NEPC

Binary classification of ARPC and NEPC subtypes using individual region and feature combinations was conducted using XGBoost v1.4.2 “XGBClassifier” implemented in Python with default parameters. Features included NPS and mean phased nucleosome distance (see Phasing analysis) in histone modification regions, promoters, and gene bodies; fragment size mean, short-long ratio, and CV (see Fragment size analysis) in histone modification regions, promoters, and gene bodies; central and window coverage (see Griffin analysis) in promoters, composite TFBSs, and composite differentially open chromatin regions identified through ATAC-seq; and max wave height (see Griffin analysis) in promoters. We applied stratified 6-fold cross-validation where two ARPC samples and one NEPC sample were held out in each fold. This was repeated 100 times and performance was computed using AUC and 95% confidence intervals for each individual feature and region combination. Code and implementation of the method can be found at <https://github.com/GavinHaLab/CRPCSubtypingPaper/tree/main/SupervisedLearning>.

Tumor Fraction Estimation

Tumor fractions from patient plasma samples were assessed using ichorCNA (81) with binSize 1,000,000 bp and both GRCh37 and GRCh38 reference genomes. Default tumor fraction estimates reported by ichorCNA were used. See https://github.com/GavinHaLab/CRPCSubtypingPaper/tree/main/ichorCNA_configuration for complete configuration settings.

Phenotype Class Prediction Model (ctdPheno)

We developed a probabilistic model to classify the mCRPC phenotype (ARPC or NEPC) in an individual patient plasma ctDNA sample. This is a generative mixture model that is unsupervised—it does not train on the patient cohort of interest. However, the model accepts the preestimated tumor fraction from ichorCNA for the given patient ctDNA sample, as well as the precomputed ctDNA features values from the LuCaP PDX ctDNA and healthy donor ctDNA as prior information. For each patient ctDNA sample, it fits specific feature values against the pure PDX LuCaP models, shifted toward healthy based on the estimated tumor fraction. The expected feature values (mean μ and standard deviation σ) from each phenotype k for feature i were taken from the mean of LuCaP PDX samples ($\mu_{i,k}$) or taken from the mean of a panel of normals H ($\mu_{i,H}$, male only, $n = 14$; see Human subjects: Healthy donor samples). Assuming a Gaussian distribution, feature values were shifted such that the shifted $\mu'_{i,k}$, $\sigma'_{i,k}$ took the form:

$$\begin{aligned}\mu'_{i,k} &= \alpha\mu_{i,k} + (1-\alpha)\mu_{i,H} \\ \sigma'_{i,k} &= \sqrt{\alpha\sigma_{i,k}^2 + (1-\alpha)\sigma_{i,H}^2}\end{aligned}$$

where α is the tumor fraction estimate for each test sample. In the final model, four features were used: composite open chromatin regions (central and window mean coverage) for specific phenotypes (ARPC and NEPC) identified from the LuCaP PDX ATAC-seq analysis using Griffin (see Griffin analysis). For each feature i , we then

found the probability that the observed sample came from a mixture of the tumor fraction-corrected Gaussian distributions, where θ is the NEPC mixture weight:

$$p_i(x|\theta) = \theta p(x|k = \text{NEPC}) + (1-\theta)p(x|k = \text{ARPC})$$

The θ parameter is estimated by maximizing the joint log-likelihood L for a given patient sample:

$$\theta' = \arg \max_{\theta} [L(x|\theta)]$$

where $L(x|\theta) = \sum_i \ln[p_i(x|\theta)]$

θ has range $[0,1]$, where higher values indicate an increased probability of the sample having an NEPC phenotype and was used as the NEPC prediction score metric. Code and implementation of the method can be found at <https://github.com/GavinHaLab/CRPCSubtypingPaper/tree/main/ctdPheno>.

Phenotype Heterogeneity Prediction and Quantification (Keraon)

We developed an analytic model to directly estimate the contributing fractions of ctDNA from different mCRPC phenotypes (ARPC and NEPC) in individual patient plasma ctDNA samples. Like ctdPheno, this model is unsupervised and does not require training on the patient cohort of interest. However, the model accepts the preestimated tumor fraction from ichorCNA for the given patient ctDNA sample, as well as the precomputed ctDNA features values from the LuCaP PDX ctDNA and healthy donor ctDNA as prior information (see Class phenotype prediction model).

As a preprocessing step, the model first computes the mean vector μ_i and covariance matrix Σ_i for each anchor class i in K , under the assumption that each subtype (including healthy) fits a multivariate Gaussian distribution. Based on model constraints, $K - 1$ noncorrelated features fully specify the system, and so for ARPC:NEPC:healthy ($K = 3$) fraction estimation, we limited analyses to sets of two features of interest ($F = 2$).

Next, for each sample defined by some location in feature space v and estimated tumor fraction t , we first performed a change of basis to translate the sample's location from feature space to class space, where each (not necessarily orthogonal) axis defined a single phenotype, and the origin represented pure healthy. If $F = K - 1$, this was accomplished by solving the determined, linear matrix equation for the shifted basis components \mathbf{X} :

$$\mathbf{B}\mathbf{X} = \mathbf{S}$$

where $\mathbf{B} = [\mu_{i \neq HD} - \mu_{HD}]$ is the matrix defining all basis vectors from the healthy mean anchor to each phenotype mean anchor, and \mathbf{S} is the vector from the healthy mean anchor to the sample of interest, $\mathbf{S} = v - \mu_{HD}$. If the system is overdetermined ($F > K - 1$), least squares was used to estimate the approximate solution. This step allows us to learn where in the class space the sample lies, which determined how estimates were evaluated:

1. Anchor space: If all basis components are positive, then the sample lies within the volume of order $K - 1$, which has vertices defined by the class means. The relative ratio of basis component magnitudes in the direction of each class is corrected by estimated tumor fraction directly: $BC_{i \neq HD} = \frac{X_i}{\sum X} t$.
2. Contra space: If all basis components are negative, then the sample lies within the volume of order $K - 1$, which forms a reflection of that formed by the class vertices about healthy. Component fractions for each basis are computed to capture the inverse distance from the healthy anchor, such that $BC_{i \neq HD} = \frac{X_i + 1}{\sum (X + 1)} t$.

3. Extra space: If some basis components are positive but others are negative, the sample lies in some space outside of the anchor or contra space. In this case, only positive contributions are considered, such that $BC_{i+HD} = \frac{X_i}{\sum X} t$ for all i such that $X_i > 0$.

The tumor fraction-normalized basis component estimates BC have a range [0,1], where values directly correspond to the total fraction of each class in the sample.

The code and implementation of the method can be found at <https://github.com/denniepatton/Keraon>.

Analysis and Classification of Clinical Patient Samples

After establishing feature distributions using the LuCaP PDX lines and normal panel as described above, both models were applied to three clinical patient cohorts (see Human subjects for cohort information).

Binary Class Prediction. Initial scoring using ctdPheno was run on DFCI cohort I, consisting of 101 ULP-WGS samples with paired-end reads. Tumor fraction estimates predicted by ichorCNA and tumor phenotype classifications were obtained from the original study (25). A prediction score threshold of 0.3314 for calling NEPC was chosen because it offered an optimal performance for sensitivity (90%) and specificity (97.5%), where sensitivity is the true positive rate for identifying NEPC samples $\left(\frac{TP}{TP + FN}\right)$ and specificity is the true negative rate for identifying ARPC samples $\left(\frac{TN}{TN + FP}\right)$. Alternative thresholds maximizing sensitivity and specificity were 0.1077, at which 95% sensitivity was achieved with a lower specificity of 93.8%, and 0.3769 with a lower sensitivity of 81.0% but higher specificity of 98.8%. To compare these predictions with cfDNA methylation (cfMeDIP-seq) classification on the same plasma samples in DFCI cohort I, the concordance was computed between the ctdPheno NEPC prediction score and the cfMeDIP NEPC score obtained from the original study using a 0.15 threshold (25).

We then validated the model on two cohorts, beginning with the already published DFCI cohort II (76, 77, 81). We restricted our analysis to 11 samples from 6 patients with matched ULP-WGS and WGS data with paired-end reads. Tumor fraction estimates from ichorCNA were obtained from the original study (81). All samples were considered adenocarcinoma (ARPC) based on clinical histories (see Human subjects). The scoring threshold of 0.3314, determined from DFCI cohort I, was used for phenotype classification.

For the UW cohort, consisting of 47 samples from 27 patients (average 22.13× depth of coverage sequencing), ichorCNA was used to estimate sample tumor fractions as described above (GRCh38), whereas clinical phenotype was determined from clinical histories and expert chart review. We evaluated model performance on matched ULP-WGS and WGS data for unambiguous clinical phenotypes of ARPC and NEPC. The chosen scoring threshold of 0.3314 was used, and the fraction of correctly predicted ARPC ($n = 26$) and NEPC ($n = 5$) was computed. The remaining 16 samples with mixed histologies were not evaluated for performance in ctdPheno.

Phenotype Prediction and Proportion Estimation. Keraon does not require *de novo* threshold selection, so all clinical cohorts were treated as validation sets. Based on the MAE of 2.8% for estimating NEPC fraction garnered in the heterogeneous mixture benchmarking, this value was chosen as the minimum NEPC fraction threshold for calling the presence of NEPC in WGS cohorts. The same tumor fraction estimates used by ctdPheno in ULP were utilized by Keraon, with standard classification conducted on pure clinical phenotypes. The 16 samples with mixed phenotypes in the UW cohort were evaluated both qualitatively and based on the 2.8% threshold in the absence of quantifiable burden estimates from histories.

Statistical Analysis

Quantification of and statistical approaches for high-throughput sequencing data analysis are described in the Methods. When nonparametric distributions (not normally distributed) of numerical values of a particular parameter in a population were compared (using boxplots or in tables), the two-tailed Mann-Whitney U test (also known as the Wilcoxon Rank Sum test; `scipy.stats.mannwhitneyu`; ref. 104) was used to test if any two distributions being compared were significantly different, with Benjamini-Hochberg (`statsmodels.stats.multitest.fdr_correction`; <https://www.statsmodels.org>) correction applied in multiple testing scenarios. All boxplots represent the median with a centerline, interquartile range (IQR) with a box, and first quartile - 1.5 IQR and third quartile + 1.5 IQR with whiskers. PCA was conducted in Python (`sklearn.decomposition.PCA`; <https://scikit-learn.org>).

Data Availability

The LuCaP PDX plasma ctDNA-seq data generated in this study can be accessed under NCBI BioProject accession PRJN900550 (<https://www.ncbi.nlm.nih.gov/bioproject/PRJN900550>). The processed patient plasma data can be accessed at <https://github.com/GavinHaLab/CRPCSubtypingPaper/tree/main/Data>. The raw sequencing data generated for the UW cohort are not publicly available because patients did not consent to genomic data sharing but are available upon reasonable request from the corresponding authors. This paper also analyzes existing, publicly available data, including LuCaP PDX RNA-seq (GSE199596) and ATAC-seq data (GSE156292). The CUT&RUN processed data can be accessed at https://github.com/nielOnav/LuCaP_nucleosome_profile. Published data for DFCI cohort I was obtained from the authors (25) after establishing a data use agreement with the Dana-Farber Cancer Institute.

Any additional information required to reanalyze the data reported in this paper is available from the corresponding authors upon request.

Authors' Disclosures

N. De Sarkar reports a patent for a cell-free DNA sequence data analysis method to examine nucleosome protection and chromatin accessibility (PCT/US2022/024082) pending and provisional patent 22-172-US-PSP pending to FHCC. R.D. Patton reports a patent for 63/353,331 pending to FHCC and a patent for PCT/US2022/024082 pending to FHCC. A. Doebly reports a patent for PCT/US2022/024082 issued. M.T. Schweizer reports personal fees from Sanofi, AstraZeneca, PharmaIn, and Resverlogix, and grants from Zenith Epigenetics, Bristol Myers Squibb, Merck, Immunomedics, Janssen, AstraZeneca, Pfizer, Hoffmann-La Roche, Tmunity, SignalOne Bio, and Ambrx outside the submitted work. A.D. Choudhury reports grants and personal fees from Bayer, and personal fees from Clovis, Dendreon, AstraZeneca, Pfizer, Astellas, Blue Earth, Janssen, Eli Lilly, and Tolmar outside the submitted work. S.C. Baca reports personal fees and other support from Precede Biosciences outside the submitted work. J.E. Berchuck reports grants from the Department of Defense during the conduct of the study, as well as a patent for methods for identifying neuroendocrine prostate cancer with tissue-informed cell-free DNA methylation analysis pending. M.L. Freedman reports personal fees from Nuscan outside the submitted work. E. Corey reports grants from the NIH during the conduct of the study, as well as grants from Bayer, Forma Therapeutics, GSK, AbbVie, Kronos, Foghorn, AstraZeneca, and MacroGenics outside the submitted work. S. Henikoff reports grants from the Howard Hughes Medical Institute during the conduct of the study. P.S. Nelson reports personal fees from Bristol Myers Squibb and Merck, and grants from Janssen outside the submitted work, as well as a patent for a cell-free DNA sequence data analysis method to examine nucleosome protection and chromatin accessibility pending. G. Ha reports a provisional patent (63/353,331) pending to FHCC and a patent for a cell-free DNA sequence data analysis method to examine nucleosome protection and chromatin accessibility

(PCT/US2022/024082) pending to FHCC. No disclosures were reported by the other authors.

Authors' Contributions

N. De Sarkar: Conceptualization, resources, data curation, software, formal analysis, validation, investigation, visualization, methodology, writing—original draft, project administration, writing—review and editing. **R.D. Patton:** Conceptualization, data curation, software, formal analysis, validation, investigation, visualization, methodology, writing—original draft, project administration, writing—review and editing. **A.-L. Doebley:** Data curation, software, formal analysis, investigation, methodology, writing—review and editing. **B. Hanratty:** Software, formal analysis, investigation, visualization. **M. Adil:** Software, formal analysis. **A.J. Kreitzman:** Software, formal analysis, investigation. **J.F. Sarthy:** Formal analysis, investigation, visualization. **M. Ko:** Formal analysis, visualization. **S. Brahma:** Resources, data curation. **M.P. Meers:** Resources, data curation. **D.H. Janssens:** Resources, data curation. **L.S. Ang:** Resources. **I.M. Coleman:** Formal analysis. **A. Bose:** Investigation. **R.F. Dumpit:** Resources. **J.M. Lucas:** Investigation. **T.A. Nunez:** Resources. **H.M. Nguyen:** Resources. **H.M. McClure:** Resources. **C.C. Pritchard:** Writing—review and editing. **M.T. Schweizer:** Writing—review and editing. **C. Morrissey:** Resources, data curation, writing—review and editing. **A.D. Choudhury:** Data curation, writing—review and editing. **S.C. Baca:** Resources. **J.E. Berchuck:** Resources. **M.L. Freedman:** Resources. **K. Ahmad:** Writing—review and editing. **M.C. Haffner:** Data curation, visualization, writing—review and editing. **R.B. Montgomery:** Writing—review and editing. **E. Corey:** Resources, data curation, writing—review and editing. **S. Henikoff:** Resources, writing—review and editing. **P.S. Nelson:** Conceptualization, resources, data curation, supervision, funding acquisition, visualization, methodology, writing—original draft, project administration, writing—review and editing. **G. Ha:** Conceptualization, resources, data curation, software, formal analysis, supervision, funding acquisition, validation, investigation, visualization, methodology, writing—original draft, project administration, writing—review and editing.

Acknowledgments

This research was supported by the Pacific Northwest Prostate Cancer SPORE grant (P50 CA097186) and the Department of Defense Idea Development Award (W81XWH-21-1-0513). This work was also supported by the Brotman Baty Institute for Precision Medicine grants (to G. Ha, R.D. Patton, and N. De Sarkar), Prostate Cancer Foundation Young Investigator Awards (G. Ha and N. De Sarkar), National Institutes of Health (K22 CA237746, R21 CA264383, and DP2 CA280624 to G. Ha; R01 CA234715 and R01 CA266452 to P.S. Nelson; P01 CA163227 to P.S. Nelson, E. Corey, and C. Morrissey; R01 CA251555 to M.L. Freedman; K99 GM138920 to S.C. Baca; K99 GM140251 to M.P. Meers), Department of Defense (W81XWH-18-1-0406 to P.S. Nelson; W81XWH-17-1-0380 to N. De Sarkar; W81XWH-18-0756, W81XWH-18-1-0356, PC170510, PC170503P2, and PC200262P to C.C. Pritchard; W81XWH-20-1-0118 to J.E. Berchuck; W81XWH-20-1-0084 to A. Bose). Support was also provided by the Prostate Cancer Foundation, the Institute for Prostate Cancer Research, V Foundation Scholar Grants (to G. Ha and M.C. Haffner), Fund for Innovation in Cancer Informatics Major Grant (to G. Ha); Doris Duke Charitable Foundation and Safeway Foundation (to M.C. Haffner); Wong Family Award in Translational Oncology and DFCI Medical Oncology grant (to A.D. Choudhury); H.L. Snyder Medical Research Foundation, the Cutler Family Fund for Prevention and Early Detection, and Claudia Adams Barr Program for Innovative Cancer Research (to M.L. Freedman); ASCO Young Investigator Award, Kure It Cancer Research Foundation, and PhRMA Foundation (to S.C. Baca). This research was also supported in part by the NIH/NCI Cancer Center Support Grant (P30 CA015704) and Scientific Computing Infrastructure (ORIP Grant

S10OD028685). We thank the many patients and their families for their altruistic contributions to this study. We thank the Fred Hutchinson Cancer Center Genomics Shared Resources Core members, the Institute for Prostate Cancer Research clinicians and staff that support the University of Washington rapid autopsy program and the PDX program. We thank Patricia Galipeau and members of the Ha and Nelson Laboratories for critically reading this manuscript.

The publication costs of this article were defrayed in part by the payment of publication fees. Therefore, and solely to indicate this fact, this article is hereby marked “advertisement” in accordance with 18 USC section 1734.

Note

Supplementary data for this article are available at Cancer Discovery Online (<http://cancerdiscovery.aacrjournals.org/>).

Received June 16, 2022; revised October 1, 2022; accepted November 16, 2022; published first November 18, 2022.

REFERENCES

- Karantanos T, Corn PG, Thompson TC. Prostate cancer progression after androgen deprivation therapy: mechanisms of castrate resistance and novel therapeutic approaches. *Oncogene* 2013;32:5501–11.
- Ryan CJ, Smith MR, de Bono JS, Molina A, Logothetis CJ, de Souza P, et al. Abiraterone in metastatic prostate cancer without previous chemotherapy. *N Engl J Med* 2013;368:138–48.
- Scher HI, Fizazi K, Saad F, Taplin M-E, Sternberg CN, Miller K, et al. Increased survival with enzalutamide in prostate cancer after chemotherapy. Cabot RC, Harris NL, Rosenberg ES, Shepard J-AO, Cort AM, Ebeling SH, et al., editors. *N Engl J Med* 2012;367:1187–97.
- Beltran H, Prandi D, Mosquera JM, Benelli M, Puca L, Cyrta J, et al. Divergent clonal evolution of castration-resistant neuroendocrine prostate cancer. *Nat Med* 2016;22:298–305.
- Bluemn EG, Coleman IM, Lucas JM, Coleman RT, Hernandez-Lopez S, Tharakan R, et al. Androgen receptor pathway-independent prostate cancer is sustained through FGF signaling. *Cancer Cell* 2017;32:474–89.
- Conteduca V, Oromendia C, Eng KW, Bareja R, Sigouros M, Molina A, et al. Clinical features of neuroendocrine prostate cancer. *Eur J Cancer* 2019;121:7–18.
- Aggarwal R, Huang J, Alumkal JJ, Zhang L, Feng FY, Thomas GV, et al. Clinical and genomic characterization of treatment-emergent small-cell neuroendocrine prostate cancer: a multi-institutional prospective study. *JCO* 2018;36:2492–503.
- Baca SC, Takeda DY, Seo J-H, Hwang J, Ku SY, Arafeh R, et al. Reprogramming of the FOXA1 cisrome in treatment-emergent neuroendocrine prostate cancer. *Nat Commun* 2021;12:1979.
- Cejas P, Xie Y, Font-Tello A, Lim K, Syamala S, Qiu X, et al. Subtype heterogeneity and epigenetic convergence in neuroendocrine prostate cancer. *Nat Commun* 2021;12:3775.
- Spetsieris N, Boukvala M, Patsakis G, Alafis I, Efstathiou E. Neuroendocrine and aggressive-variant prostate cancer. *Cancers* 2020;12:3792.
- Labrecque MP, Coleman IM, Brown LG, True LD, Kollath L, Lakely B, et al. Molecular profiling stratifies diverse phenotypes of treatment-refractory metastatic castration-resistant prostate cancer. *J Clin Invest* 2019;129:4492–505.
- Labrecque MP, Alumkal JJ, Coleman IM, Nelson PS, Morrissey C. The heterogeneity of prostate cancers lacking AR activity will require diverse treatment approaches. *Endocrine-related cancer*. Bioscientifica Ltd; 2021;28:T51–66.
- Liu Y, Horn JL, Banda K, Goodman AZ, Lim Y, Jana S, et al. The androgen receptor regulates a druggable translational regulon in advanced prostate cancer. *Sci Transl Med* 2019;11:eaaw4993.
- Epstein JI, Amin MB, Beltran H, Lotan TL, Mosquera J-M, Reuter VE, et al. Proposed morphologic classification of prostate cancer with neuroendocrine differentiation. *Am J Surg Pathol* 2014;38:756–67.

15. Annala M, Taavitsainen S, Khalaf DJ, Vandekerkhove G, Beja K, Sipola J, et al. Evolution of castration-resistant prostate cancer in ctDNA during sequential androgen receptor pathway inhibition. *Clin Cancer Res* 2021;27:4610–23.
16. Aparicio AM, Shen L, Tapia ELN, Lu J-F, Chen H-C, Zhang J, et al. Combined tumor suppressor defects characterize clinically defined aggressive variant prostate cancers. *Clin Cancer Res* 2016;22:1520–30.
17. Carreira S, Romanel A, Goodall J, Grist E, Ferraldeschi R, Miranda S, et al. Tumor clone dynamics in lethal prostate cancer. *Sci Transl Med* 2014;6:254ra125.
18. Du M, Tian Y, Tan W, Wang L, Wang L, Kilari D, et al. Plasma cell-free DNA-based predictors of response to abiraterone acetate/prednisone and prognostic factors in metastatic castration-resistant prostate cancer. *Prostate Cancer Prostatic Dis* 2020;23:705–13.
19. Sumanasuriya S, Seed G, Parr H, Christova R, Pope L, Bertan C, et al. Elucidating prostate cancer behaviour during treatment via low-pass whole-genome sequencing of circulating tumour DNA. *Eur Urol* 2021;80:243–53.
20. Ulz P, Belic J, Graf R, Auer M, Lafer I, Fischereeder K, et al. Whole-genome plasma sequencing reveals focal amplifications as a driving force in metastatic prostate cancer. *Nat Commun* 2016;7:12008.
21. Wyatt AW, Annala M, Aggarwal R, Beja K, Feng F, Youngren J, et al. Concordance of circulating tumor DNA and matched metastatic tissue biopsy in prostate cancer. *J Natl Cancer Inst* 2018;110:78–86.
22. Nyquist MD, Corella A, Coleman I, De Sarkar N, Kaipainen A, Ha G, et al. Combined TP53 and RB1 loss promotes prostate cancer resistance to a spectrum of therapeutics and confers vulnerability to replication stress. *Cell Rep* 2020;31:107669.
23. Berger A, Brady NJ, Bareja R, Robinson B, Conteduca V, Augello MA, et al. N-Myc-mediated epigenetic reprogramming drives lineage plasticity in advanced prostate cancer. *J Clin Invest* 2019;129:3924–40.
24. Beltran H, Romanel A, Conteduca V, Casiraghi N, Sigouros M, Franceschini GM, et al. Circulating tumor DNA profile recognizes transformation to castration-resistant neuroendocrine prostate cancer. *J Clin Invest* 2020;130:1653–68.
25. Berchuck JE, Baca SC, McClure HM, Korthauer K, Tsai HK, Nuzzo PV, et al. Detecting neuroendocrine prostate cancer through tissue-informed cell-free DNA methylation analysis. *Clin Cancer Res* 2022;28:928–38.
26. Shen SY, Singhanian R, Fehringer G, Chakravarthy A, Roehrl MHA, Chadwick D, et al. Sensitive tumour detection and classification using plasma cell-free DNA methylomes. *Nature* 2018;563:579–83.
27. Wu A, Cremaschi P, Wettterskog D, Conteduca V, Franceschini GM, Klefogiannis D, et al. Genome-wide plasma DNA methylation features of metastatic prostate cancer. *J Clin Invest* 2020;130:1991–2000.
28. Heitzer E, Auinger L, Speicher MR. Cell-free DNA and apoptosis: how dead cells inform about the living. *Trends Mol Med* 2020;26:519–28.
29. Lo YMD, Han DSC, Jiang P, Chiu RWK. Epigenetics, fragmentomics, and topology of cell-free DNA in liquid biopsies. *Science* 2021;372:eaaw3616.
30. Cristiano S, Leal A, Phallen J, Fiksel J, Adleff V, Bruhm DC, et al. Genome-wide cell-free DNA fragmentation in patients with cancer. *Nature* 2019;570:385–9.
31. Jiang P, Sun K, Peng W, Cheng SH, Ni M, Yeung PC, et al. Plasma DNA end-motif profiling as a fragmentomic marker in cancer, pregnancy, and transplantation. *Cancer Discov* 2020;10:664–73.
32. Mathios D, Johansen JS, Cristiano S, Medina JE, Phallen J, Larsen KR, et al. Detection and characterization of lung cancer using cell-free DNA fragmentomes. *Nat Commun* 2021;12:5060.
33. Peneder P, Stütz AM, Surdez D, Krumbholz M, Semper S, Chicard M, et al. Multimodal analysis of cell-free DNA whole-genome sequencing for pediatric cancers with low mutational burden. *Nat Commun* 2021;12:3230.
34. Zhu G, Guo YA, Ho D, Poon P, Poh ZW, Wong PM, et al. Tissue-specific cell-free DNA degradation quantifies circulating tumor DNA burden. *Nat Commun* 2021;12:2229.
35. Herberts C, Annala M, Sipola J, Ng SWS, Chen XE, Nurminen A, et al. Deep whole-genome ctDNA chronology of treatment-resistant prostate cancer. *Nature* 2022;608:199–208.
36. Jiang P, Chan CWM, Chan KCA, Cheng SH, Wong J, Wong VW-S, et al. Lengthening and shortening of plasma DNA in hepatocellular carcinoma patients. *Proc Natl Acad Sci U S A* 2015;112:E1317–25.
37. Moulriere F, Chandrananda D, Piskorz AM, Moore EK, Morris J, Ahlborn LB, et al. Enhanced detection of circulating tumor DNA by fragment size analysis. *Sci Transl Med* 2018;10:eaat4921.
38. Snyder MW, Kircher M, Hill AJ, Daza RM, Shendure J. Cell-free DNA comprises an in vivo nucleosome footprint that informs its tissues-of-origin. *Cell* 2016;164:57–68.
39. Underhill HR, Kitzman JO, Hellwig S, Welker NC, Daza R, Baker DN, et al. Fragment length of circulating tumor DNA. *PLOS Genet* 2016;12:426–37.
40. Ramachandran S, Ahmad K, Henikoff S. Transcription and remodeling produce asymmetrically unwrapped nucleosomal intermediates. *Molecular Cell*. Cell Press; 2017;68:1038–53.
41. Ulz P, Thallinger GG, Auer M, Graf R, Kashofer K, Jahn SW, et al. Inferring expressed genes by whole-genome sequencing of plasma DNA. *Nat Genet* 2016;48:1273–8.
42. Ulz P, Perakis S, Zhou Q, Moser T, Belic J, Lazzeri I, et al. Inference of transcription factor binding from cell-free DNA enables tumor subtype prediction and early detection. *Nat Commun* 2019;10:4666.
43. Esfahani MS, Hamilton EG, Mehrmohamadi M, Nabet BY, Alig SK, King DA, et al. Inferring gene expression from cell-free DNA fragmentation profiles. *Nat Biotechnol* 2022;40:585–97.
44. Brahma S, Henikoff S. Epigenome regulation by dynamic nucleosome unwrapping. *Trends Biochem Sci* 2020;45:13–26.
45. Lai WKM, Pugh BF. Understanding nucleosome dynamics and their links to gene expression and DNA replication. *Nat Rev Mol Cell Biol* 2017;18:548–62.
46. Yen K, Vinayachandran V, Batta K, Koerber RT, Pugh BF. Genome-wide nucleosome specificity and directionality of chromatin remodelers. *Cell* 2012;149:1461–73.
47. Rao S, Han AL, Zukowski A, Kopin E, Sartorius CA, Kabos P, et al. Transcription factor–nucleosome dynamics from plasma cfDNA identifies ER-driven states in breast cancer. *Sci Adv* 2022;8:eabm4358.
48. Nguyen HM, Vessella RL, Morrissey C, Brown LG, Coleman JM, Higano CS, et al. LuCaP prostate cancer patient-derived xenografts reflect the molecular heterogeneity of advanced disease and serve as models for evaluating cancer therapeutics. *Prostate* 2017;77:654–71.
49. Skene PJ, Henikoff S. An efficient targeted nuclease strategy for high-resolution mapping of DNA binding sites. *Reinberg D, editor. eLife* 2017;6:e21856.
50. Meers MP, Tenenbaum D, Henikoff S. Peak calling by sparse enrichment analysis for CUT&RUN chromatin profiling. *Epigenetics Chromatin* 2019;12:42.
51. Zhou VW, Goren A, Bernstein BE. Charting histone modifications and the functional organization of mammalian genomes. *Nat Rev Genet* 2011;12:7–18.
52. Doebley AL, Ko M, Liao H, Cruikshank AE, Santos K, Kikawa C, et al. A framework for clinical cancer subtyping from nucleosome profiling of cell-free DNA. *Nat Commun* 2022 Dec;13:7475.
53. Soares LM, He PC, Chun Y, Suh H, Kim T, Buratowski S. Determinants of histone H3K4 methylation patterns. *Mol Cell* 2017;68:773–85.
54. Brady NJ, Bagadion AM, Singh R, Conteduca V, Van Emmenis L, Arceci E, et al. Temporal evolution of cellular heterogeneity during the progression to advanced AR-negative prostate cancer. *Nat Commun* 2021;12:3372.
55. Wang YA, Sfakianos J, Tewari AK, Cordon-cardo C, Kyprianou N. Molecular tracing of prostate cancer lethality. *Oncogene* 2020;39:7225–38.
56. Rapa I, Ceppi P, Bollito E, Rosas R, Cappia S, Bacillo E, et al. Human ASH1 expression in prostate cancer with neuroendocrine differentiation. *Mod Pathol* 2008;21:700–7.
57. Pomerantz MM, Qiu X, Zhu Y, Takeda DY, Pan W, Baca SC, et al. Prostate cancer reactivates developmental epigenomic programs during metastatic progression. *Nat Genet* 2020;52:790–9.
58. Severson TM, Zhu Y, De Marzo AM, Jones T, Simons JW, Nelson WG, et al. Epigenetic and transcriptional analysis reveals a core transcriptional program conserved in clonal prostate cancer metastases. *Mol Oncol* 2021;15:1942–55.

59. Labrecque MP, Brown LG, Coleman IM, Lakely B, Brady NJ, Lee JK, et al. RNA splicing factors SRRM3 and SRRM4 distinguish molecular phenotypes of castration-resistant neuroendocrine prostate cancer. *Cancer Res* 2021;81:4736–50.
60. Tsai HK, Lehrer J, Alshalhafa M, Erho N, Davicioni E, Lotan TL. Gene expression signatures of neuroendocrine prostate cancer and primary small cell prostatic carcinoma. *BMC Cancer* 2017;17:759.
61. Jiang Z, Zhang B. On the role of transcription in positioning nucleosomes. *PLoS Comput Biol* 2021;17:e1008556.
62. Klemm SL, Shipony Z, Greenleaf WJ. Chromatin accessibility and the regulatory epigenome. *Nat Rev Genet* 2019;20:207–20.
63. Oruba A, Saccani S, van Essen D. Role of cell-type specific nucleosome positioning in inducible activation of mammalian promoters. *Nat Commun* 2020;11:1075.
64. Guo Y, Zhao S, Wang GG. Polycomb gene silencing mechanisms: PRC2 chromatin targeting, H3K27me3 “Readout”, and phase separation-based compaction. *Trends Genet* 2021;37:547–65.
65. Jiang C, Pugh BF. Nucleosome positioning and gene regulation: advances through genomics. *Nat Rev Genet* 2009;10:161–72.
66. Saxton DS, Rine J. Nucleosome positioning regulates the establishment, stability, and inheritance of heterochromatin in *Saccharomyces cerevisiae*. *Proc Natl Acad Sci U S A* 2020;117:27493–501.
67. Valouev A, Johnson SM, Boyd SD, Smith CL, Fire AZ, Sidow A. Determinants of nucleosome organization in primary human cells. *Nature* 2011;474:516–20.
68. Deal RB, Henikoff JG, Henikoff S. Genome-wide kinetics of nucleosome turnover determined by metabolic labeling of histones. *Science* 2010;328:1161–4.
69. Cuzick J, Swanson GP, Fisher G, Brothman AR, Berney DM, Reid JE, et al. Prognostic value of an RNA expression signature derived from cell cycle proliferation genes in patients with prostate cancer: a retrospective study. *Lancet Oncol* 2011;12:245–55.
70. Chereji RV, Bryson TD, Henikoff S. Quantitative MNase-seq accurately maps nucleosome occupancy levels. *Genome Biol* 2019;20:198.
71. Yevshin I, Sharipov R, Kolmykov S, Kondrakhin Y, Kolpakov F. GTRD: a database on gene transcription regulation—2019 update. *Nucleic Acids Res* 2019;47:D100–5.
72. AroraVK, SchenkeinE, MuraliR, SubudhiSK, WongvipatJ, BalbasMD, et al. Glucocorticoid receptor confers resistance to antiandrogens by bypassing androgen receptor blockade. *Cell* 2013;155:1309–22.
73. Mu P, Zhang Z, Benelli M, Karthaus WR, Hoover E, Chen C-C, et al. SOX2 promotes lineage plasticity and antiandrogen resistance in TP53- and RB1-deficient prostate cancer. *Science* 2017;355:84–8.
74. Shukla S, Cyrta J, Murphy DA, Walczak EG, Ran L, Agrawal P, et al. Aberrant activation of a gastrointestinal transcriptional circuit in prostate cancer mediates castration resistance. *Cancer Cell* 2017;32:792–806.
75. Sun K, Jiang P, Cheng SH, Cheng THT, Wong J, Wong VWS, et al. Orientation-aware plasma cell-free DNA fragmentation analysis in open chromatin regions informs tissue of origin. *Genome Res* 2019;29:418–27.
76. Viswanathan SR, Ha G, Hoff AM, Wala JA, Carrot-Zhang J, Whelan CW, et al. Structural alterations driving castration-resistant prostate cancer revealed by linked-read genome sequencing. *Cell* 2018;174:433–47.
77. Choudhury AD, Werner L, Francini E, Wei XX, Ha G, Freeman SS, et al. Tumor fraction in cell-free DNA as a biomarker in prostate cancer. *JCI Insight* 2018;3:e122109.
78. Klein DC, Hainer SJ. Genomic methods in profiling DNA accessibility and factor localization. *Chromosome Res* 2020;28:69–85.
79. Chaytor L, Simcock M, Nakjang S, Heath R, Walker L, Robson C, et al. The pioneering role of GATA2 in androgen receptor variant regulation is controlled by bromodomain and extraterminal proteins in castrate-resistant prostate cancer. *Mol Cancer Res* 2019;17:1264–78.
80. Corces MR, Granja JM, Shams S, Louie BH, Seoane JA, Zhou W, et al. The chromatin accessibility landscape of primary human cancers. *Science* 2018;362.
81. Adalsteinsson VA, Ha G, Freeman SS, Choudhury AD, Stover DG, Parsons HA, et al. Scalable whole-exome sequencing of cell-free DNA reveals high concordance with metastatic tumors. *Nat Commun* 2017;8:1324.
82. Fang R, Preissl S, Li Y, Hou X, Lucero J, Wang X, et al. Comprehensive analysis of single-cell ATAC-seq data with SnapATAC. *Nat Commun* 2021;12:1337.
83. Wu SJ, Furlan SN, Mihalas AB, Kaya-Okur HS, Feroze AH, Emerson SN, et al. Single-cell CUT&Tag analysis of chromatin modifications in differentiation and tumor progression. *Nat Biotechnol* 2021;39:819–24.
84. Lam H-M, Nguyen HM, Corey E. Generation of prostate cancer patient-derived xenografts to investigate mechanisms of novel treatments and treatment resistance. In: Culig Z, editor. *Prostate cancer: methods and protocols*. New York, NY: Springer; 2018, pp. 1–27.
85. Li H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. arXiv: 2013, 13033997 [q-bio] [cited 2022 Mar 22]. Available from: <http://arxiv.org/abs/1303.3997>.
86. Jo S-Y, Kim E, Kim S. Impact of mouse contamination in genomic profiling of patient-derived models and best practice for robust analysis. *Genome Biol* 2019;20:231.
87. DePristo MA, Banks E, Poplin R, Garimella KV, Maguire JR, Hartl C, et al. A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nat Genet* 2011;43:491–8.
88. Hänzelmann S, Castelo R, Guinney J. GSEA: gene set variation analysis for microarray and RNA-Seq data. *BMC Bioinf* 2013;14:7.
89. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* 2013;29:15–21.
90. Kluijn RJC, Kemper K, Kuilman T, de Ruiter JR, Iyer V, Forment JV, et al. XenofilteR: computational deconvolution of mouse and human reads in tumor xenograft sequence data. *BMC Bioinf* 2018;19:366.
91. Robinson MD, McCarthy DJ, Smyth GK. edgeR: a Bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics* 2010;26:139–40.
92. Lambert SA, Jolma A, Campitelli LF, Das PK, Yin Y, Albu M, et al. The human transcription factors. *Cell* 2018;172:650–65.
93. Langmead B, Wilks C, Antonescu V, Charles R. Scaling read aligners to hundreds of threads on general-purpose processors. *Bioinformatics* 2019;35:421–32.
94. Ramírez F, Ryan DP, Grüning B, Bhardwaj V, Kilpert F, Richter AS, et al. deepTools2: a next-generation web server for deep-sequencing data analysis. *Nucleic Acids Res* 2016;44:W160–5.
95. Ross-Innes CS, Stark R, Teschendorff AE, Holmes KA, Ali HR, Dunning MJ, et al. Differential oestrogen receptor binding is associated with clinical outcome in breast cancer. *Nature* 2012;481:389–93.
96. Quinlan AR, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics* 2010;26:841–2.
97. Yu G, Wang L-G, He Q-Y. ChIPseeker: an R/Bioconductor package for ChIP peak annotation, comparison and visualization. *Bioinformatics* 2015;31:2382–3.
98. Amemiya HM, Kundaje A, Boyle AP. The ENCODE blacklist: identification of problematic regions of the genome. *Sci Rep* 2019;9:9354.
99. Khan A, Mathelier A. Intervene: a tool for intersection and visualization of multiple gene or genomic region sets. *BMC Bioinf* 2017;18:287.
100. Benjamini Y, Speed TP. Summarizing and correcting the GC content bias in high-throughput sequencing. *Nucleic Acids Res* 2012;40:e72–.
101. Karimzadeh M, Ernst C, Kundaje A, Hoffman MM. Umap and Bismap: quantifying genome and methylome mappability. *Nucleic Acids Res* 2018;46:e120.
102. Weirauch MT, Yang A, Albu M, Cote AG, Montenegro-Montero A, Drewe P, et al. Determination and inference of eukaryotic transcription factor sequence specificity. *Cell* 2014;158:1431–43.
103. Howe KL, Achuthan P, Allen J, Allen J, Alvarez-Jarreta J, Amode MR, et al. Ensembl 2021. *Nucleic Acids Res* 2021;49:D884–91.
104. Virtanen P, Gommers R, Oliphant TE, Haberland M, Reddy T, Cournapeau D, et al. SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nat Methods* 2020;17:261–72.
105. Danecek P, Bonfield JK, Liddle J, Marshall J, Ohan V, Pollard MO, et al. Twelve years of SAMtools and BCFtools. *GigaScience* 2021;10:giab008.