# 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 Hanratty[2], Mohamed Adil[1], Adam J. Kreitzman[1,2], Jay F. Sarthy[7], Minjeong Ko[1,2], Sandipan Brahma[7], Michael P. Meers[7], Derek H. Janssens[7], Lisa S. Ang[2], Ilsa M. Coleman[2], Arnab Bose[2], Ruth F. Dumpit[2], Jared M. Lucas[2], Talina A. Nunez[2], Holly M. Nguyen[8], Heather M. McClure[9], Colin C. Pritchard[10,11], Michael T. Schweizer[3,12], Colm Morrissey[8], Atish D. Choudhury[9,13], Sylvan C. Baca[9,13], Jacob E. Berchuck[9], Matthew L. Freedman[9,13], Kami Ahmad[6], Michael C. Haffner[2,3,10], R. Bruce Montgomery[12], Eva Corey[8], Steven Henikoff[7,14], Peter S. Nelson[2,3,8,11,12,†], Gavin Ha[1,2,11,15,†]

[1] Division of Public Health Sciences, Fred Hutchinson Cancer Center, 1100 Fairview Ave. N, Seattle, WA 98109
[2] Division of Human Biology, Fred Hutchinson Cancer Center, 1100 Fairview Ave. N, Seattle, WA 98109
[3] Division of Clinical Research, Fred Hutchinson Cancer Center, 1100 Fairview Ave. N, Seattle, WA 98109
[4] Department of Pathology and Prostate Cancer Center of Excellence, Medical College of Wisconsin, 8701 W Watertown Plank Road, Milwaukee, WI 53226
[5] Molecular and Cellular Biology Graduate Program, University of Washington, 1959 NE Pacific St, Seattle WA 98195
[6] Medical Scientist Training Program, University of Washington, 1959 NE Pacific St, Seattle WA 98195
[7] Division of Basic Sciences, Fred Hutchinson Cancer Center, 1100 Fairview Ave. N, Seattle, WA 98109
[8] Department of Urology, University of Washington, 1959 NE Pacific St, Seattle, WA, 98195
[9] Dana-Farber Cancer Institute, 450 Brookline Ave, Boston, MA 02215
[10] Department of Laboratory Medicine and Pathology, University of Washington, 1959 NE Pacific St, Seattle, WA 98195
[11] Brotman Baty Institute for Precision Medicine, 1959 NE Pacific St, Seattle, WA, 98195
[12] Division of Oncology, Department of Medicine, University of Washington
[13] Broad Institute of MIT and Harvard, 415 Main Street, Cambridge, MA 02142
[14] Howard Hughes Medical Institute, USA
[15] Department of Genome Sciences, University of Washington, 1959 Pacific St, Seattle, WA, 98195
* These authors contributed equally.
† Co-senior authors.

**RUNNING TITLE:** Cancer phenotype classification using ctDNA

**KEYWORDS:** Circulating tumor DNA, liquid biopsies, castration-resistant prostate cancer, patient-derived xenografts, fragmentomics

## ADDITIONAL INFORMATION

1

## Correspondence:

Gavin Ha, Ph.D.
Fred Hutchinson Cancer Center
1100 Fairview Ave. N,
Seattle, WA 98109
206-667-2802
gha@fredhutch.org

Peter S. Nelson, M.D.
Fred Hutchinson Cancer Center
1100 Fairview Ave. N,
Seattle, WA 98109
206-667-3377
pnelson@fredhutch.org

## Conflict of interest disclosure:

G.H., A-L.D., N.D.S., R.D.P., P.S.N.: These authors have filed pending patent applications on methodologies developed in this manuscript (PCT/US2002/024082 and provisional patent USSN: 63/353,331).

P.S.N.: Served as a paid consultant to Janssen, Astellas, Pfizer, and Bristol Myers Squibb in work unrelated to the present study.

B.M.: Has institutional funding from Clovis, Janssen, Astellas, BeiGene, and AstraZeneca.

E.C.: Received research funding under institutional SRA from Janssen Research and Development, Bayer Pharmaceuticals, KronosBio, Forma Pharmaceutics Foghorn, Gilead, Sanofi, AbbVie, and GSK for work unrelated to the present study.

M.L.F.: Serves as a consultant to and has equity in Nuscan Diagnostics. This activity is outside of the scope of this manuscript. M.L.F. has a pending patent for detecting NEPC using DNA methylation.

M.T.S.: Paid consultant and/or received Honoria from Sanofi, AstraZeneca, PharmaIn and Resverlogix. He has received research funding to his institution from Zenith Epigenetics, Bristol Myers Squibb, Merck,

Immunomedics, Janssen, AstraZeneca, Pfizer, Madison Vaccines, Hoffman-La Roche, Tmunity, SignalOne Bio and Ambrx, Inc.

All other authors declare no competing interests.

1 **ABSTRACT**

2 Advanced prostate cancers comprise distinct phenotypes, but tumor classification remains
3 clinically challenging. Here, we harnessed circulating tumor DNA (ctDNA) to study tumor
4 phenotypes by ascertaining nucleosome positioning patterns associated with transcription
5 regulation. We sequenced plasma ctDNA whole genomes from patient-derived xenografts
6 representing a spectrum of androgen receptor active (ARPC) and neuroendocrine (NEPC)
7 prostate cancers. Nucleosome patterns associated with transcriptional activity were reflected in
8 ctDNA at regions of genes, promoters, histone modifications, transcription factor binding, and
9 accessible chromatin. We identified the activity of key phenotype-defining transcriptional
10 regulators from ctDNA, including AR, ASCL1, HOXB13, HNF4G, and GATA2. To distinguish
11 NEPC and ARPC in patient plasma samples, we developed prediction models that achieved
12 accuracies of 97% for dominant phenotypes and 87% for mixed clinical phenotypes. While
13 phenotype classification is typically assessed by immunohistochemistry or transcriptome
14 profiling from tumor biopsies, we demonstrate that ctDNA provides comparable results with
15 diagnostic advantages for precision oncology.
16

17 **STATEMENT OF SIGNIFICANCE**

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

## 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 trans-differentiation from AR-positive adenocarcinoma (ARPC) to NE prostate cancer (NEPC) that lack AR activity (4,7–10). Additional phenotypes can also arise based on expression of AR activity and NE genes, including AR-low prostate cancer (ARLPC) and double-negative prostate cancer (DNPC; AR-null/NE-null) (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 non-invasive "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 *RB1* in NEPC do not always lead to NE trans-differentiation (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 non-random enzymatic cleavage by nucleases (28,29). Emerging approaches to analyze cfDNA fragmentation patterns from plasma for studying cancer can be performed directly from standard whole genome sequencing (WGS) (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 cancer patients (36–

5

54 39). Recent studies have demonstrated that the nucleosome occupancy in cfDNA at the
55 transcription start site (TSS) and transcription factor binding site (TFBS) can be used to infer
56 gene expression and transcription factor (TF) activity from cfDNA (40–42). However,
57 nucleosome positioning and spacing are dynamic in active and repressed gene regulation (43–
58 45). A detailed understanding of the nucleosome patterns and accessible chromatin associated
59 with transcriptional regulation in tumor phenotypes has not been fully explored in cfDNA.

60 The objective of this study is to determine if ctDNA could be used to accurately classify tumor
61 phenotypes in men with mCRPC. A major challenge for ctDNA analysis is the low tumor content
62 (tumor fraction) in patient plasma samples. By contrast, plasma from patient-derived xenograft
63 (PDX) models may contain nearly pure human ctDNA after bioinformatic exclusion of mouse
64 DNA reads (37,39,46). This provides a resource that is ideal for studying the properties of
65 ctDNA, developing new analytical tools, and validating both genetic and phenotypic features by
66 comparison to matching tumors. In this study, we performed WGS of ctDNA from mouse plasma
67 across 24 CRPC PDX lines with diverse phenotypes. Applying newly developed computational
68 methods, we comprehensively interrogated the nucleosome patterns in ctDNA across genes,
69 regulatory loci, TFBSs, TSSs, and open chromatin sites to reveal transcriptional regulation
70 associated with mCRPC phenotypes. Finally, we designed two probabilistic models to
71 accurately classify treatment-resistant tumors into divergent phenotypes and to estimate the
72 phenotype heterogeneity within a ctDNA sample. We then validated the performance of these
73 models in 159 plasma samples from three mCRPC patient cohorts. Overall, these results
74 highlight that transcriptional regulation of tumor phenotypes can be ascertained from ctDNA and
75 has potential utility for diagnostic applications in cancer precision medicine.

76 **RESULTS**

77 **Comprehensive resource of matched tumor and liquid biopsies from patient derived**
78 **xenograft (PDX) models of advanced prostate cancer**

79 To develop approaches for the accurate classification of mCRPC using ctDNA, we evaluated 26
80 models from the LuCaP PDX series of advanced prostate cancer with well-defined phenotypes
81 determined by whole transcriptome RNAseq and immunohistochemical assays for protein
82 expression (47). The models consisted of 18 classified as ARPC, two classified as AR-low and
83 NE-negative prostate cancer (ARLPC), and six classified as NEPC (**Figure 1A**). For each PDX
84 line, we pooled mouse plasma (1.9 – 3.0 mL) from four to eight mice (mean tumor volume range
85 393-1239 mm$^3$), extracted cfDNA, and performed deep whole genome sequencing (WGS; mean

6

86　38.4x coverage, range 21 – 85x) **(Methods, Supplementary Table S1)**. We used bioinformatic

87　subtraction of mouse sequenced reads to obtain nearly pure human ctDNA data **(Methods)**. We

88　observed that 25 lines had human ctDNA comprising more than 10% of the sample (mean

89　52.9%, range 10.6 – 96%) with NEPC samples having significantly higher human fractions

90　(mean 85.1%, range 77.1 – 96%, two-tailed Mann-Whitney U test p = 9.6 x $10^{-4}$) (**Figure 1B,**

91　**Supplementary Table S1**). After subsequent filtering by human ctDNA sequencing coverage,

92　24 PDX lines remained for further analysis (16 ARPC, 6 NEPC, 2 ARLPC; mean 20.5x, range

93　3.8 – 50.6x, **Supplementary Table S1**). In the matching tumors, we performed Cleavage Under

94　Targets and Release using Nuclease (CUT&RUN) to profile H3K27ac, H3K4me1, and

95　H3K27me3 histone post-translational modifications (PTMs) (48,49) (**Supplementary Fig. S1**).

96　We hypothesized that nucleosome organization inferred from ctDNA reflects the transcriptional

97　activity state regulated by histone PTMs (50).

98　To study transcriptional regulation in mCRPC phenotypes from ctDNA, we interrogated four

99　different features: (i) local promoter coverage, (ii) nucleosome positioning, (iii) fragment size

100　analysis, and (iv) composite TFBSs plus open chromatin sites analysis using the Griffin

101　framework (51) (**Figure 1A, Methods**). First, we analyzed three different local regions within

102　ctDNA: all gene promoters and gene bodies and sites of histone PTMs guided by CUT&RUN

103　analysis. For each of the three local regions, we extracted features of nucleosome periodicity

104　using a nucleosome phasing approach and computed the fragment size variability. For promoter

105　regions, we also computed the coverage at the transcription start site (TSS). Next, we analyzed

106　ctDNA at transcription factor binding sites (TFBSs) and open chromatin regions. For each

107　transcription factor (TF), ctDNA coverage at TFBSs were aggregated into composite profiles

108　representing the inferred activity (42,51). Similarly, features in the composite profiles of

109　phenotype-specific open chromatin regions were extracted for analyzing the signatures of

110　chromatin accessibility in ctDNA. Altogether, we assembled a multi-omic sequencing dataset

111　from matching tumor and plasma for a total of 24 PDX lines, making this a unique molecular

112　resource and platform for developing transcriptional regulation signatures of tumor phenotype

113　prediction from ctDNA.

**Characterizing transcriptional activity of AR and ASCL1 in PDX phenotypes through**
**analysis of tumor histone modifications and ctDNA**

116　We sought to further characterize the transcriptional activity in different tumor phenotypes by

117　studying epigenetic regulation via histone PTMs. We identified broad peak regions for H3K4me1

7

118   (median of 17,643 regions, range 1,894 – 64,934), H3K27ac (median 7,093, range 1610 -

119   34,047), and H3K27me3 (median 8,737, range 2,024 - 42,495) in the tumors of the 24 PDX

120   lines and an additional nine LuCaP PDX lines where only tumors were available (total of 25

121   ARPC, 2 ARLPC, and 6 NEPC) (**Methods**, **Supplementary Fig. S1**, **Supplementary Table S2**).

122   Using unsupervised clustering and principal components analysis (PCA), we identified putative

123   active regulatory regions of enhancers and promoters (H3K27ac, H3K4me1) and gene

124   repressive heterochromatic marks (H3K27me3) that were specific to ARPC, ARLPC, and NEPC

125   phenotypes (52) (**Supplementary Fig. S2A**).

126   AR and ASCL1 are two key differentially expressed TFs with known regulatory roles in ARPC

127   and NEPC phenotypes, respectively (9,53–55). When inspecting AR binding sites in ARPC

128   tumors, we observed increased signals from flanking nucleosomes with H3K27ac PTMs

129   compared to the other phenotypes (area under mean peak profile of 18.46 vs. 15.08 in ARLPC

130   and 10.63 in NEPC, **Figure 2A**, **Supplementary Fig. S2B**, **Methods**). We also observed the

131   strongest signals at the nucleosome depleted region (NDR) in ARPC for H3K27ac (1.54

132   coverage decrease vs. 0.78 for ARLPC and 0.41 for NEPC). Conversely, in NEPC tumors, we

133   observed stronger signals at nucleosomes with H3K27ac PTMs flanking ASCL1 binding sites

134   (area under mean peak profile 62.65 vs. 29.18 for ARLPC and 10.83 for ARPC), and stronger

135   NDR signals (2.26 coverage decrease vs. 0.19 for ARPC and 0.37 for ARLPC). We observed

136   similar trends for H3K4me1 PTMs in the LuCaP PDX lines (**Supplementary Fig. S2C**).

137   We analyzed the ctDNA composite coverage profiles at 1,000 consensus TFBSs to evaluate

138   nucleosome accessibility, where lower normalized central (±30 bp window) mean coverage

139   across these sites suggests more nucleosome depletion (**Methods**). For AR TFBSs, we

140   observed the strongest signal for nucleosome depletion in ARPC, as indicated by the lowest

141   mean central coverage (average 0.64, n=16), compared to moderate signals for ARLPC

142   (average 0.88, n=2), and weakest signals for NEPC (average 0.95, n=6) (**Figure 2B**).

143   Conversely, the composite coverage profile at ASCL1 TFBSs showed the strongest nucleosome

144   depletion for NEPC samples (mean central coverage 0.69) compared to ARLPC (0.86) and

145   ARPC (0.88) (**Figure 2C**). These observations were consistent with the differential binding

146   activity by AR and ASCL1 in their respective phenotypes from tumor tissue (**Figure 2A**). We

147   confirmed the same differential binding activity trends when analyzing TFBSs identified from

148   other primary tissue sources (9,56,57) (**Supplementary Fig. S3A-B**). We also noted that the

149   composite TFBS coverage patterns in ctDNA resembled the NDR flanked by nucleosomes with

150   H3K27ac and H3K4me1 modifications inferred by CUT&RUN (**Figure 2A**, **Supplementary Fig.**

8

151 **S2B-C**). Together, these results suggest that the nucleosome depletion in ctDNA at AR and
152 ASCL1 binding sites represents active TF binding and regulatory activity in specific prostate
153 PDX tumor phenotypes.

**Nucleosome patterns at gene promoters inferred from ctDNA are consistent with transcriptional activity for phenotype-specific genes**

156 We selected 47 genes comprising 12 ARPC and 35 NEPC lineage markers established
157 previously (4,5,58,59) and confirmed their phenotype associations by RNA-Seq from the PDX
158 tumors (**Figure 2D**, **Supplementary Table S3**, **Methods**). To assess the activity of these genes
159 from ctDNA, we analyzed the ctDNA fragment size in TSSs (± 1 kb window) and gene bodies
160 and found that the differential size variability between phenotypes was positively correlated with
161 relative expression (Spearman's r = 0.844, p = 9.4 x $10^{-14}$, **Figure 2E**, **Supplementary Fig. S4**,
162 **Supplementary Table S2**, **Methods**). However, closer inspection of ctDNA coverage patterns
163 at promoters revealed consistent nucleosome organization for transcription activity and
164 repression (40,60–62) (**Figure 2D**). Therefore, we grouped the genes based on differential
165 signals in H3K27me3 histone PTMs, which are linked with polycomb repressive complex
166 mediated regulation and chromatin compaction (63).

167 For 25 genes without differential H3K27me3 peaks (Group 1), including AR, KLK3 and ASCL1,
168 we observed nucleosome depletion at the TSS consistent with presence of active PTMs, such
169 as for AR (mean coverage 0.47, n=16) in ARPC and ASCL1 (0.30, n=6) in NEPC samples
170 (**Figure 2F, Supplementary Fig. S5**). By contrast, we observed increased coverage at the TSS
171 of AR (1.08) in NEPC and ASCL1 (0.42) in ARPC, which supports nucleosome depletion in the
172 absence of PTMs and inactive transcription. For 22 genes with differential H3K27me3 peaks
173 (Group 2), including INSM1, CHGB and SRRM4, we observed relatively consistent increase in
174 nucleosome occupancy and phasing in the TSS as well as in the gene body for ~50% of the
175 genes (**Figure 2G, Supplementary Fig. S6**). The neuronal signaling genes in this group, such
176 as UNC13A and INSM1, had reduced signals for the stable nucleosome dyad position,
177 consistent with the heterogeneous ('fuzzy') nucleosome patterns described for actively
178 transcribed genes (44,64). Interestingly, while UNC13A was active in NEPC tumors, we did not
179 detect H3K27ac nor H3K4me1 PTM marks in the regulatory loci of this gene (**Supplementary
180 Fig. S7A-B**). These results illustrate that ctDNA analysis can reveal patterns that are consistent
181 with different modalities of transcriptional regulation by histone modifications for key genes that
182 define prostate cancer phenotypes.

9

**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 (61,65,66). 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 (67). To systematically quantify inter-nucleosomal spacing and predict nucleosome phasing, we developed TritonNP, a tool utilizing Fourier transforms and band-pass filters on GC-corrected ctDNA coverage to isolate frequency components corresponding to phased nucleosomes (**Figure 3A**, **Supplementary Fig. S8A-B**, **Methods**). This approach allows for calling phased nucleosome dyad positions to generate an average inter-nucleosome 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 to the NEPC lines (median 291.1 bp vs. 282.6 bp, $p = 0.027$; two-tailed Mann-Whitney U test, **Figure 3B**). The phased nucleosome distance was also negatively correlated with the mean cell cycle progression (CCP) score (Spearman's rho = -0.563, $p = 0.006$, **Figure 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 (68).

To model the relationship between nucleosome phasing and transcriptional activity more directly, we further extracted the frequency components corresponding to the inter-dyad distances of "stable" nucleosomes (180 – 210 bp) and a "baseline" component (150 – 180 bp) for normalization between samples of differing depths (69). 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 (**Figure 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 to ARPC

10

216 (median NPS 1.95 vs. 2.21, p = 0.134; two-tailed Mann-Whitney U test, **Figure 3D**); although

217 this was not statistically significant, it was consistent with their active transcription. Conversely,

218 the mean NPS for the 12 AR-regulated genes was lower in ARPC lines compared to NEPC

219 (median NPS 1.82 vs 2.13, p = 0.070; two-tailed Mann-Whitney U test). In particular, 26 (74%)

220 of the NE genes had lower NPS in NEPC compared to ARPC ($\log_2$ fold-change [ARPC:NEPC] >

221 0), including seven genes (ASCL1, CHGB, CHRNB2, GRP, MYCL, XKR7, NEUROD1) that

222 were statistically significant (p < 0.05); ten (83%) of the AR-regulated genes had lower NPS in

223 ARPC ($\log_2$ fold-change < 0), with TMPRSS2 being statistically significant (**Figure 3E**,

224 **Supplementary Table S3**). These results illustrate that the NPS captured signals distinguishing

225 key lineage-specific gene markers.

**Inferred TF activity from analysis of nucleosome accessibility at TFBSs in ctDNA confirms key regulators of tumor phenotypes**

228 To characterize the regulation of prostate tumor phenotype lineages, we considered

229 nucleosome accessibility at TFBSs in PDX ctDNA for 338 TFs from the Gene Transcription

230 Regulation Database (GTRD)(70) (**Methods**). First, we identified 108 TFs out of the 338 that

231 were differentially expressed between ARPC and NEPC PDX tumors by RNAseq

232 (**Supplementary Fig. S9, Supplementary Table S3, Methods**). Through unsupervised

233 hierarchical clustering of composite TFBS central coverage values for these 108 TFs, we

234 observed distinct groups of TFs in PDX ctDNA (**Figure 3F**). Of these 108 TFs, 38 had

235 significantly different accessibility in ctDNA between ARPC and NEPC phenotypes (two tailed

236 Mann-Whitney U test, Benjamini-Hochberg adjusted p < 0.05, **Supplementary Table S3**). Most

237 of these TFs (27/38 [71%]) had differential inferred accessibility in ctDNA that was consistent

238 with their up-regulation in the same phenotype by tumor mRNA expression, although some TFs

239 (11/38, [29%]) did not show this trend (**Figure 3F**, **Supplementary Fig. S10**). A comparison of

240 TFBS between paralogous TFs revealed that the binding sites used in the analysis had limited

241 overlap (median 18.3%, range 0-81.2%), suggesting that many of the TFs may have some

242 independent inferred accessibility (**Supplementary Fig. S11, Supplementary Table S3**). For

243 paralogs with high TFBS overlap (≥ 19%), such as AR, NR3C1 and PGR, we noted only a

244 subset of TFs were expressed in one phenotype.

245 REST had the largest difference in accessibility as supported by a decrease in coverage within

246 ARPC models compared to NEPC ($\log_2$ fold-change -0.77, adjusted p = $5.7 \times 10^{-4}$,

247 **Supplementary Fig. S12A**, **Supplementary Table S3**). FOXA1, and GRHL2 binding sites were

11

248    significantly more accessible in ARPC (and ARLPC) samples compared to NEPC ($log_2$ fold-

249    change < -0.57, adjusted p < 1.3 x $10^{-3}$). AR, HOXB13, and NKX3-1 had higher accessibility in

250    ARPC compared to NEPC ($log_2$ fold-change < -0.37, adjusted p < 1.3 x $10^{-3}$), but with only

251    moderate accessibility in ARLPC, as expected.  We also observed a group of TFs that followed

252    a similar trend, including nuclear hormone receptors (NR2F2, RARG), pioneer factor GATA2,

253    and nuclear factors HNF4G and HNF1A ($log_2$ fold-change < -0.10, adjusted p < 0.027,

254    **Supplementary Fig. S12A**).

255    For factors that had higher accessibility in NEPC models compared to ARPC and ARLPC,

256    ASCL1 had the largest TFBS coverage difference ($log_2$ fold-change 0.36, adjusted p = 5.7 x $10^{-4}$,

257    **Figure 2C**, **Figure 3F)**. Other TFs, including RUNX1, BCL11B, POU3F2, NEUROG2, and

258    SOX2 also had sites with higher accessibility in NEPC ($log_2$ fold-change > 0.06, adjusted p <

259    0.048, **Supplementary Fig. S12B**), although the difference was modest. Other notable factors

260    such as MYC and ETS transcription family genes (ETV4, ETV5, ETS1, ETV1) had high

261    accessibility across all phenotypes, while NEUROD1, RUNX3, and TP63 sites were

262    inaccessible in nearly all samples. Furthermore, we considered restricting the analysis to 20 TFs

263    with TFBSs that were observed in prostatic tissue and cell lines and were also differentially

264    expressed in the PDX tumors by RNAseq (**Methods**). However, while hierarchical clustering

265    distinguished PDX tumor phenotypes, key NEPC-defining markers, such as ASCL1, were

266    omitted from this analysis as ChIP-seq for many NEPC-defining markers had not been

267    performed on prostate lineage samples in GTRD (**Supplementary Fig. S13**). Overall, we

268    identified the accessibility of known prostate cancer regulators, including ASCL1, HNF4G,

269    HNF1A, GATA2 and SOX2 (71–73), that have not been shown before from ctDNA analysis in

270    these tumor phenotypes.

271    **Phenotype-specific open chromatin regions (ATAC-Seq) in PDX tumor tissue are**

272    **reflected in ctDNA profiles of nucleosome accessibility**

273    Nucleosome profiling from cfDNA sequencing analysis has shown agreement with overall

274    chromatin accessibility in tumor tissue (38,42,74); however, its application for distinguishing

275    tumor phenotypes has been limited. We hypothesized that due to lack of protection from

276    nucleases, regions of open chromatin would be under-represented in ctDNA assays. We

277    investigated the use of ATAC-Seq data from tumor tissue for 10 LuCaP PDX lines (5 ARPC and

278    5 NEPC) to inform phenotype-related differences in chromatin accessibility (9). We defined an

279    initial set of 28,765 ARPC and 21,963 NEPC differential consensus open chromatin regions

12

280 which we further restricted to those that overlapped TFBSs for 338 TFs, resulting in 15,879

281 ARPC and 11,692 NEPC sites (**Methods**, **Figure 4A**). For ARPC-specific open chromatin sites,

282 we observed decreased overall composite site coverage (+/- 1 kb window) and central coverage

283 (+/- 30 bp) in the ctDNA for ARPC PDX lines (mean central coverage 0.75, n=16) compared to

284 NEPC lines (mean 0.96, n=6) and cfDNA from healthy human donors (mean 0.97, n=14)

285 (**Figure 4B**, **Supplementary Table S3, Methods**). Conversely, for NEPC-specific open

286 chromatin sites, coverage was decreased in ctDNA for NEPC lines (mean 0.89) compared to

287 ARPC lines (mean 1.01) and healthy donors cfDNA (mean 1.00) (**Figure 4C**, **Supplementary**

288 **Table S3**). Coverage patterns were discernable between phenotypes for as few as 100 sites,

289 suggesting that even a smaller subset of open chromatin regions may still be informative

290 (**Supplementary Fig. S14A-B**). These results confirmed that tumor tissue chromatin

291 accessibility can be corroborated in ctDNA and that ARPC and NEPC phenotypes have distinct

292 ctDNA coverage profiles at these sites.

293 **Comprehensive evaluation of ctDNA features across genomic contexts for CRPC**

294 **phenotype classification**

295 To assess the utility of ctDNA nucleosome profiling for informing prostate cancer phenotype

296 classification, we systematically evaluated four groups of global genome-wide ctDNA features:

297 phasing, fragment sizes, local coverage profiling, and composite site coverage profiling (**Figure**

298 **1A**). From principal components analysis (PCA), we observed distinct feature signals between

299 ARPC and NEPC phenotypes for composite TFBS coverage of TFs, NPS of the 47 phenotype

300 marker genes, and fragment size variability at global sites of PTMs (**Figure 4D**, **Supplementary**

301 **Fig. S15A, Supplementary Table S4**, **Methods**). In addition to these features, we also

302 included previously reported features, including short-long fragment ratio and local coverage

303 patterns at the TSS (max wave height between -120bp to 195bp) (30,41) (**Methods**).

304 We then quantitatively evaluated all combinations of coverage, phasing, and fragment size

305 features for different genomic contexts to investigate their potential to classify ARPC and NEPC

306 phenotypes. For each feature set, we conducted 100 iterations of stratified cross-validation

307 using a supervised machine learning classifier (XGBoost) on ctDNA samples from the ARPC

308 and NEPC models and computed the area under the receiver operating characteristic curve

309 (AUC) (**Methods**). First, we evaluated an established set of 10 genes associated with AR

310 activity (5,12). We observed that the phased nucleosome distance at H3K27ac sites and the

311 central coverage at TSSs had moderate predictive performance (AUC 0.88) (**Supplementary**

312   **Fig. S15B**, **Supplementary Table S4**). For the set of 47 phenotype markers, the NPS of gene

313   bodies was most predictive (AUC 0.98) (**Supplementary Fig. S15C, Supplementary Table S4**).

314   When considering all PTM sites, promoters, genes, TFs, and open chromatin regions, the best

315   performing features included mean fragment size at H3K4me1 sites (n=9,750, AUC 1.0) and

316   promoter TSSs (n=17,946, AUC 1.0), and both open chromatin composite site features (AUC

317   1.0) (**Figure 4E**, **Supplementary Table S4**).

318   **Accurate classification of ARPC and NEPC phenotypes from patient plasma using a**

319   **probabilistic model informed by PDX ctDNA analysis**

320   An important consideration and challenge in analyzing plasma from patients is the presence of

321   cfDNA released by hematopoietic cells, which leads to a lower ctDNA fraction (i.e., tumor

322   fraction). Furthermore, the small patient cohorts with available tumor phenotype information

323   make supervised machine learning approaches suboptimal. Therefore, we developed ctdPheno,

324   a probabilistic model to classify ARPC and NEPC from an individual plasma sample, accounting

325   for the tumor fraction (**Figure 5A, Methods**). We focused on the phenotype-specific open

326   chromatin composite site features and used the PDX plasma ctDNA signals (27,571 total sites,

327   **Figure 4B-C**, **Supplementary Table S3**) to inform the model. The model produces a

328   normalized prediction score that represents the estimated signature of ARPC (lower values) and

329   NEPC (higher values). We applied this method to benchmarking datasets generated by

330   simulating varying tumor fractions and sequencing coverages using five ARPC and NEPC PDX

331   ctDNA samples each, and healthy donor plasma cfDNA (**Supplementary Figure 15D**,

332   **Methods**). We achieved a 1.0 AUC at 25X coverage down to 0.01 tumor fraction, 1.0 AUC at

333   1X down to 0.2 tumor fraction, and 1.0 AUC at 0.2x coverage at 0.3 tumor fraction, suggesting a

334   possible upper-bound performance for classifying samples with lower tumor fraction in plasma

335   (**Figure 5B**, **Supplementary Table S4**).

336   To test the performance of ctdPheno on patient samples, we analyzed a published dataset of

337   ultra-low-pass whole genome sequencing (ULP-WGS) of plasma cfDNA (mean coverage 0.52X,

338   range 0.28-0.92X) from 101 mCRPC patients comprising 80 adenocarcinoma (ARPC) and 21

339   NEPC samples (DFCI cohort I) (25). Using ctdPheno, which was unsupervised and used

340   parameters informed only by the PDX analysis, we achieved an overall AUC of 0.96 (**Figure 5C**,

341   **Supplementary Table S5**). The performance was 0.97 AUC and 0.76 AUC when considering

342   samples with high (≥ 0.1) and low (< 0.1) tumor fraction, respectively, and 0.83 AUC when using

343   only 2000 sites for analysis (**Supplementary Fig. S16A-B**). We identified an optimal overall

14

performance at 97.5% specificity (ARPC) and 90.4% sensitivity (NEPC) which corresponded to the prediction score of 0.3314 (**Figure 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 PSA, treatment with ARSI, or both (DFCI cohort II) (75,76), the model correctly classified patients as ARPC in 8 (73%) ULP-WGS (~0.1x) samples when using the optimal score cutoff (**Supplementary Fig. S16D**, **Supplementary Table S5).**

Next, we analyzed 61 clinical plasma samples from 31 CRPC patients with ARPC, NEPC, and mixed phenotypes that are representative of typical clinical histories (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 (**Figure 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 (**Figure 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 trans-differentiation can lead to a mixture of ARPC and NEPC cells or lesions. To account for and predict phenotype mixtures within a patient ctDNA sample, we developed Keraon, an analytical model that estimates the proportions of phenotypes from WGS using the same ctDNA features as ctdPheno (**Figure 5E, Methods**). First, we evaluated Keraon using a benchmark dataset generated for simulating varying tumor fractions and proportions of ARPC-NEPC mixtures at 25x coverage using PDX ctDNA and healthy donor cfDNA data (**Figure 5F, Methods**). In 810 simulated phenotype mixtures, we observed the estimated total NEPC fraction was consistent with expected proportions (Pearson's r=0.884) with a mean absolute error (MAE) of 0.028,

15

376  highlighting the method's potential for accurate estimation of emergent phenotypes in mixed
377  histology samples (**Figure 5G, Supplementary Table S5**). Next, we evaluated Keraon for
378  classifying NEPC in DFCI Cohort I and observed the highest performance (0.96 AUC) using all
379  27,571 open chromatin sites, with decreased performance (0.84 AUC) when using only 2,000
380  sites (**Supplementary Fig. S16D**). Applying Keraon to analyze DFCI Cohort II, we correctly
381  estimated dominant ARPC with undetectable NEPC phenotype in 10 (91%) samples with WGS
382  (mean coverage, 27x) (**Supplementary Fig. 18**, **Supplementary Table S5**).

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

396  **DISCUSSION**

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

405  In the present study, we analyzed multiple features of DNA to infer the activity of gene
406  expression programs corresponding to distinct prostate cancer phenotypes. A key component of

16

407   the work that allowed for the development of optimized methods and the identification of the
408   most informative ctDNA features was the use of PDX models. The sequencing of mouse plasma
409   provided a unique opportunity to comprehensively interrogate the epigenetic nucleosome
410   patterns in ctDNA from well-characterized tumor models. We developed and applied
411   computational methodologies to evaluate a multitude of ctDNA features, each of which were
412   associated with transcriptional regulation across CRPC tumor phenotypes. The use of PDX
413   mouse plasma overcomes the challenge of low ctDNA content or incomplete knowledge of the
414   tumor when studying patient samples. Using features learned from the PDX ctDNA, we
415   developed models to accurately classify ARPC and NEPC and to estimate their proportions in
416   phenotypically heterogenous samples from patient plasma in three clinical cohorts. While these
417   data were focused on ARPC and NEPC phenotypes, the approaches may serve as a framework
418   for the use of ctDNA to subtype malignancies arising in other organ sites based on distinctive
419   gene expression programs.

420   The analysis of the LuCaP PDX ctDNA sequencing data confirmed the activity of key regulators
421   between ARPC and NEPC phenotypes, including a set of 47 established differentially
422   expressed genes that associate with cell lineage. While gene expression inference from ctDNA
423   has been shown in proof-of-concept studies (34,41), the PDX ctDNA allowed for a detailed
424   dissection of nucleosome organization associated with transcriptional activity of individual genes
425   that define the tumor phenotypes. Previous analytical approaches have profiled nucleosome
426   occupancy from cfDNA (38,74). However, our assessment of nucleosome stability by means of
427   the Nucleosome Phasing Score is the first to capture the highly variable spacing, positioning,
428   and turnover of the nucleosome arrays associated with transcription and tumor aggressiveness
429   (43,67,68,77).

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

437   While the regulation of key factors such as AR, HOXB13, NKX-3.1, FOXA1, and REST has
438   been shown from ctDNA in CRPC (35,42), we report the differential activity of other key factors

17

439    in CRPC from ctDNA analysis. This included nuclear factors HNF4G and HNF1A, and

440    pioneering factor GATA2, which are associated with prostate adenocarcinoma (ARPC)

441    (71,73,78). ASCL1 is a pioneer TF with roles in neuronal differentiation and was recently

442    described to be active during NE trans-differentiation and in NEPC (9,55). To our knowledge,

443    this study is the first to demonstrate ASCL1 binding site accessibility and provide a detailed

444    characterization of its transcriptional activity in NEPC from plasma ctDNA.

445    We show an expansive analysis of TFBSs for 338 factors in each plasma sample without the

446    need for chromatin immunoprecipitation or other epigenetic assays. However, we did not find a

447    significant difference in accessibility for 70 out of the 108 TFs in ctDNA, which may be

448    consistent with TF activity not necessarily being correlated with its own expression level (79).

449    On the other hand, the accessibility of TFBSs may not necessarily indicate true TF activity as

450    other co-bound TFs or co-activators/co-repressors influence gene regulation. Moreover, our

451    analyses were based on TFBSs obtained from public databases, including for a limited number

452    of prostate-specific TFs; however, expanded phenotype-specific TF cistrome data may improve

453    this approach.

454    We applied state-of-the-art computational approaches built on existing and new concepts of

455    ctDNA data analysis to extract tumor-specific features, including the representation of

456    nucleosome phasing, periodicity, and spacing associated with transcriptional activity. Other

457    approaches have also considered regions, such as TSSs, TFBSs, and DNase hypersensitivity

458    sites (33,38,41,42); however, after a systematic evaluation, we found that ctDNA features in

459    open chromatin sites derived from ATAC-Seq of PDX tissue (9) provided the highest

460    performance for distinguishing CRPC phenotypes. We presented ctdPheno, which is a

461    probabilistic model that classifies ARPC and NEPC from ULP-WGS data, and Keraon, an

462    analytical model which estimates the proportion of ARPC and NEPC from WGS of patient

463    plasma. Both models are unsupervised and utilize a statistical framework informed directly by

464    parameters from the LuCaP PDX ctDNA analysis. These models do not require training on

465    patient samples but do require tumor fraction estimates (ichorCNA (80)) and in the case of

466    ctdPheno a prediction score cutoff determined from DFCI cohort I. Both frameworks can also be

467    extended to model additional phenotypes. Insights from additional datasets such as single-cell

468    nucleosome and accessibility profiling (81,82) of PDX tumors and clinical samples may improve

469    the resolution for ctDNA analysis. While we observed optimal performance analyzing all open

470    chromatin sites, a smaller subset was still informative which may be useful when considering

471    targeted assays for clinical applications.

18

472    Applying the prediction models to patient datasets with definitive clinical phenotypes yielded

473    high performance even when using low depth of coverage sequencing. In particular, our

474    performance for the DFCI cohort I was also consistent with the reported phenotype classification

475    results using ctDNA methylation in the same patients (25). Similarly, in the UW cohort, samples

476    with well-defined clinical phenotypes had near-perfect concordance from WGS data. We

477    established the lower limits of phenotype classification performance to be at 8% tumor fraction

478    for ctdPheno (ULP-WGS) and 3% for Keraon (WGS). These results support a strategy whereby

479    ULP-WGS is performed for screening using ctdPheno, along with clinical assessments, and

480    followed-up with standard WGS for more accurate and comprehensive phenotype

481    characterization using Keraon. While this framework may have limited performance for low

482    (<3%) ctDNA levels, it may be optimal at initial assessment of metastatic disease and at tumor

483    progression on therapy, which is when the clinical decision points are most critical.

484    Tumor heterogeneity and co-existence of different molecular phenotypes are common in

485    mCRPC where treatment-induced phenotypic plasticity may vary within and between tumors in

486    an individual patient. In real data simulations and patients with cases of mixed clinical

487    phenotypes, Keraon accurately detected the contributions of mixed phenotypes with a detection

488    limit of 2.8% NEPC, providing the first approach to directly quantify phenotype proportions and

489    heterogeneity from ctDNA. In this study, estimation of phenotype heterogeneity using Keraon

490    required standard depths of WGS. Larger studies with comprehensive assessment of the tumor

491    histologies will be needed for evaluating these models as potential biomarkers of treatment

492    response.

493    In summary, this study illustrates that analysis of ctDNA from PDX mouse plasma at scale can

494    facilitate a detailed investigation of tumor regulation. These results, together with the suite of

495    computational methods presented here, highlight the utility of ctDNA for surveying

496    transcriptional regulation of tumor phenotypes and its potential diagnostic applications in cancer

497    precision medicine.

498    **ACKNOWLEDGEMENTS**

502 Washington rapid autopsy program and the PDX program. We thank Patricia Galipeau and
503 members of the Ha and Nelson Laboratories for critically reading this manuscript.

504 **AUTHOR CONTRIBUTIONS**

505 Conceptualization: N.D.S., R.D.P., G.H., P.S.N.

506 Methodology: R.D.P., N.D.S., A-L.D., P.S.N., G.H.

507 Software: R.D.P., A-L.D., N.D.S., B.H., A.J.K., G.H.

508 Formal Analysis: R.D.P., N.D.S., A-L.D., B.H., A.J.K., M.K., M.A., I.M.C., G.H.

509 Investigation: N.D.S., R.D.P., A-L.D., B.H., J.F.S., J.M.L., A.B., G.H.

510 Resources: N.D.S., J.S., S.B., M.P.M., D.H.J., L.S.A., R.F.D., T.A.N., H.M.M., S.C.B., J.E.B.,
511 M.L.F., C.M., H.M.N., E.C., S.H., P.S.N., G.H.

512 Data Curation: N.D.S., R.D.P., A-L.D., M.P.M., S.B., D.H.J., E.C., C.M., A.D.C., M.C.H., P.S.N.,
513 G.H.

514 Writing – Original Draft: R.D.P, N.D.S, P.S.N., G.H.

515 Writing – Review & Editing: N.D.S, R.D.P., A-L.D., C.C.P., C.M., A.D.C., M.T.S., R.B.M., M.C.H.,
516 E.C., K.A., S.H., P.S.N., G.H.

517 Visualization: R.D.P., N.D.S., B.H., M.K., M.C.H., P.S.N., G.H.

518 Supervision: P.S.N., G.H.

519 Funding Acquisition: G.H., P.S.N.

## MATERIALS AND METHODS

### *PDX mouse models*

The establishment and characterization of the LuCaP PDX models were described previously (83). 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 University of Washington Human Subjects Division IRB (IRB #2341). PDX lines were evaluated using histopathology by at least two expert pathologists, and histological phenotypic subtype annotations were orthogonally validated based on transcriptome-derived signature marker expression scores to define phenotypes (4,5,22): adenocarcinoma AR-positive (ARPC), neuroendocrine positive (NEPC), and AR-low, neuroendocrine negative (ARLPC). Resected PDX tumors (300-800 mm$^3$) were divided into ~50mg to ~100mg pieces and stored at -80°C. Animal studies were approved by the Fred Hutchinson Cancer Center (FHCC) IACUC (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 300-1400 mm$^3$).

### *Human subjects*

*UW cohort:* Blood samples were collected from men with metastatic castration resistant prostate cancer at the University of Washington (collected under University of Washington Human Subjects Division IRB protocol number CC6932 between years 2014-2021). 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 ultra-low pass whole genome sequencing (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, FH0482_E_1_A) were retained for further high depth of coverage whole genome sequencing (WGS) analysis. All samples were de-identified prior to ctDNA analysis and we employed a double blinded approach for evaluating clinical phenotype predictions.

*DFCI cohort I:* Plasma was collected from men diagnosed with mCRPC and treated at the Dana-Farber Cancer Institute (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

21

552 samples was approved by the DFCI IRB (#01-045 and 09-171) and WCM (1305013903) IRBs.
553 The ULP-WGS data at mean coverage 0.5x (range 0.3x – 0.9x) for 101 patients were published
554 previously (25).

555 DFCI cohort II: Plasma samples in this cohort were collected from men diagnosed with mCRPC
556 and treated at the Dana-Farber Cancer Institute (DFCI). All patients provided written informed
557 consent for blood collection and the analysis of their clinical and genetic data for research
558 purposes (DFCI Protocol # 01-045 and 11-104). WGS data at mean coverage 27x (range 11x –
559 44x) (75), and ULP-WGS data at mean coverage 0.13x (range 0.07x – 0.18x) (76,80) were
560 downloaded from dbGAP accession phs001417. Eleven samples from six patients had
561 matching WGS and ULP-WGS with paired-end reads, necessary for analysis by Griffin. Prostate
562 specific antigen (PSA, ng/mL) values and treatment at the time of the blood draw were
563 previously published (76). The six patients were treated for adenocarcinoma using Abiraterone,
564 Enzalutamide, or Bicalutamide, or the patients had detectable levels of PSA.

565 Healthy donor plasma cfDNA WGS data used in this study were obtained from previously
566 published studies. Two samples (HD45 and HD46, both male) with coverage of 13x and 15x,
567 respectively, were accessed from dbGAP under accession phs001417 (75,80). These donors
568 were consented under DFCI protocol IRB (# 03-022). Thirteen healthy donor plasma cfDNA
569 WGS data (12 male: NPH002, 03, 06, 07, 12, 18, 23, 26, 33, 34, 35, 36; 1 female (used in
570 admixtures): NPH004) with coverages between 13.5x – 27.6x were obtained from the European
571 Phenome Archive (EGA) under accession EGAD00001005343 (42).

572 ***PDX plasma processing***
573 Blood samples were collected from NSG mice bearing subcutaneous PDX tumors at the time of
574 sacrifice. The PDX lines were maintained at vivaria in the University of Washington and FHCC.
575 The blood was processed following methods described for human plasma DNA processing for
576 subsequent DNA isolation. Blood was collected in Sarstedt Micro sample tube K3 EDTA tubes
577 and processed within 4 hours. All blood samples were sequentially double spun, first at 2500g
578 for 10 minutes followed by a 16000g centrifugation of the plasma fraction for 10 minutes at room
579 temperature. For each PDX line, 4-8 mouse plasma samples were pooled. Processed plasma
580 samples were preserved in clean, screw-capped cryo-microfuge tubes and stored at -80°C prior
581 to cfDNA isolation.

582     *Cell-free DNA isolation*

583     The QIAamp Circulating Nucleic Acid Kit was used to isolate cfDNA from PDX mouse-derived

584     plasma using the recommended protocol. The pooled plasma samples from 4-8 mice for each

585     PDX line contained 1.9 to 3 mL total plasma volume for each line. The filter retention-based

586     cfDNA kit method does not implement any fragment size class enrichment. Isolated cfDNA was

587     quantified using the Qubit dsDNA HS assay (Invitrogen) and the cfDNA fragment size profiles

588     were analyzed using TapeStation HS D5000 and HS D1000 assays (Agilent).

589     *Cell-free DNA library preparation and sequencing*

590     For LuCaP PDX mouse plasma samples, NGS libraries were prepared with 50ng input cfDNA.

591     Illumina NGS sequencing libraries were prepared with the KAPA hyperprep kit, adopting nine

592     cycles of amplification, and purified using lab standardized SPRI beads. We used KAPA UDI

593     dual indexed library adapters. Library concentrations were balanced and pooled for multiplexing

594     and sequenced using the Illumina HiSeq 2500 at the Fred Hutch Genomics Shared Resources

595     (200 cycles) and Illumina NovaSeq platform at the Broad Institute Genomics Platform Walkup-

596     Seq Services using S4 flow cells (300 cycles). To match with Illumina HiSeq 2500 data,

597     truncated 200 cycles FASTQ files were generated (100 bp paired end reads).

598     Clinical patient plasma samples collected at University of Washington (UW cohort) were

599     submitted to the Broad Institute Blood Biopsy Services. Briefly, cfDNA was extracted from 2 mL

600     double-spun plasma and ultra-low-pass whole genome sequencing (ULP-WGS) to

601     approximately 0.2x coverage was performed. The ichorCNA pipeline was used to estimate

602     tumor DNA content (i.e., tumor fraction, see below). Forty-seven samples (from 31 patients) had

603     either ≥ 5% tumor fraction or ≥ 2% tumor fraction with AR amplification observed in ichorCNA

604     and were subsequently sequenced to deeper WGS coverage (~20x).

605     *Cell-free DNA sequencing analysis and mouse subtraction*

606     All cfDNA sequencing data used in this study were realigned to the hg38 (GRCh38) human

607     reference    genome    (http://hgdownload.soe.ucsc.edu/goldenPath/hg38/bigZips/hg38.fa.gz).

608     FASTQ files were realigned using BWA (v0.7.17) mem (84). The complete alignment pipeline

609     including       configuration       settings       may       be       access       at

610     https://github.com/GavinHaLab/fastq_to_bam_paired_snakemake.

611     For PDX ctDNA whole-genome sequence data, we performed mouse genome subtraction

612     following the protocol described previously (85), wherein reads were aligned using BWA mem to

613 a concatenated reference consisting of both human (hg38) and mouse (mm10, GRCm38.p6,

614 http://igenomes.illumina.com.s3-website-us-east-

615 1.amazonaws.com/Mus_musculus/NCBI/GRCm38/Mus_musculus_NCBI_GRCm38.tar.gz)

616 reference genomes. Read pairs where both reads aligned to the human reference genome were

617 retained and all other read pairs were removed. Then, remaining reads were re-aligned to the

618 human-only reference. Finally, the GATK best practices workflow was applied to each sample

619 (86); the complete mouse subtraction pipeline used in this study, including tool versions and

620 parameters, can be accessed at https://github.com/GavinHaLab/PDX_mouseSubtraction.

621 Following mouse subtraction samples with < 3X depth were removed for downstream analysis.

622 ***Cell cycle progression (CCP) score calculation***

623 The 31-gene cell cycle proliferation (CCP) signature (68) was computed from RNAseq data

624 using GSVA (87). The single-sample enrichment scores were calculated with default

625 parameters using genome-wide log2 FPKM values as input for the 31 genes.

626 ***Differential mRNA expression analysis***

627 RNA isolation of 102 tumors from 46 LuCaP PDX samples was performed as described

628 previously (11). RNA concentration, purity, and integrity was assessed by NanoDrop (Thermo

629 Fisher Scientific Inc) and Agilent TapeStation and RNA RIN >=8 was retained for library

630 preparation. RNA-Seq libraries were constructed from 1 ug of total RNA using the Illumina

631 TruSeq Stranded mRNA LT Sample Prep Kit according to the manufacturer's protocol.

632 Barcoded libraries were pooled and sequenced by Illumina NovaSeq 6000 or Illumina HiSeq

633 2500 generating 50 bp paired end reads. Sequencing reads were mapped to the hg38 human

634 reference genome and mm10 mouse reference genomes using STAR.v2.7.3a (88). All

635 subsequent analyses were performed in R-4.1.0. Sequences aligning to the mouse genome and

636 therefor derived from potential contamination with mouse tissue were removed from the analysis

637 using XenofilteR (v1.6) (89). Gene level abundance was quantitated using the R package

638 GenomicAlignments v1.32.0 summarizeOverlaps function using mode=IntersectionStrict,

639 restricted to primary aligned reads. We used refSeq gene annotations for transcriptome analysis.

640 Transcript abundances (FPKM) were input to edgeR v3.38.1 (90), filtered for a minimum

641 expression level using the filterByExpr function with default parameters, and then limma v3.52.1

642 voom was used for differential expression analysis of NEPC vs. ARPC and ARLPC vs. ARPC.

643 We then filtered the results using a list of 1,635 human transcription factors published previously

24

644 (91), which resulted in 514 genes with FDR<0.05 and $\log_2$ fold change > 1.58. Out of these 514,

645 deregulation of gene expression for 404 transcription factor genes delineated ARPC from NEPC.

### Cleavage Under Targets & Release Using Nuclease (CUT&RUN)

647 CUT&RUN is an antibody targeted enzyme tethering chromatin profiling assay in which

648 controlled cleavage by micrococcal nuclease releases specific protein-DNA complexes into the

649 supernatant for paired-end DNA sequencing analysis. We performed CUT&RUN assays for

650 three histone modifications, H3K27ac, H3K4me1, and H3K27me3, according to published

651 protocols (48). We performed CUT&RUN on LuCaP PDX tumors using ~75mg flash-frozen

652 tissue pieces.

653 Paired-end (50 bp) sequencing was performed and reads were aligned using bowtie2 v2.4.2

654 (92) to the hg38 human reference assembly. Aligned reads were processed as described in the

655 SEACR protocol (https://github.com/FredHutch/SEACR#preparing-input-bedgraph-files). Peaks

656 were called using SEACR version 1.3 (49) using "stringent" settings and with reference to paired

657 IgG controls. BigWig files were prepared using bamCoverage in deepTools 3.5.0 (93).

658 Genomewide peak heatmap, targeted heatmap, and respective profiles were plotted using

659 deepTools v3.5.0. bigwig formatted files for each phenotype were obtained using the mean

660 function in wiggletools 1.2.8. and deepTools computeMatrix. Phenotype-specific informative

661 region coordinates were obtained from diffBind v3.5.0, and the top 10,000 most significant

662 regions (all with FDR < 0.05) differentially open between ARPC and NEPC lines were used for

663 downstream feature analyses (see Gene body and promoter region selection for additional

664 subsetting criteria applied on a feature-by-feature basis). For heatmaps and profiles the

665 plotHeatmap function was used. We utilized the "Peak Center" option to derive desired

666 heatmaps. These steps were all performed for H3K27ac, H3K4me1 and H3K27me3 antibodies.

667 Scaled heatmap profiles' area under the curve (AUC; ±1.5kbp) and peak height at the profile

668 center were estimated using deepStats v0.4 (https://zenodo.org/record/3668336) (comparable

669 profiles are scaled to 10 units).

### Differential histone post-translational modification (PTM) analysis

671 Differential PTM analysis was performed with the Diffbind version 2.16.0 package (94) in R-

672 4.0.1 using standard parameters

673 (https://bioconductor.riken.jp/packages/3.0/bioc/html/DiffBind_.html). ARPC, NEPC and ARLPC

674 samples were grouped by histopathological and transcriptome signature defined phenotypes

675 described in the "PDX mouse models" section (**Supplementary Table S2**). Samples were

25

676    loaded with the dba function, reads counted with the dba.count function, and contrast specified

677    as phenotype with dba.contrast and a minimum members of 2. Differential peak sites were

678    computed with the dba.analyze function with default settings. Differential peak binding of NEPC

679    and ARLPC was computed against ARPC samples. Unique binding sites in NEPC and ARLPC

680    were catalogued using bedtools v2.29.2 (95). Intergroup differentially bound peaks were

681    annotated using ChIPseeker 1.28.3 (96) and TxDb.Hsapiens.UCSC.hg38.knownGene 3.2.2 in

682    R 4.1.0.

683    ***ATAC-Seq analysis***

684    ATAC-Seq sequence data for 15 tumor samples from 10 PDX lines were published previously

685    (9). These lines included LuCaP PDX lines with ARPC (23.1, 77, 78, 81, 96) and NEPC (three

686    replicates of 173.1, two replicates each of 49, 93, 145.1, and one replicate of 145.2) phenotypes.

687    Paired end reads were aligned using bowtie2 v2.4.2 (92) to the UCSC hg38 human reference

688    assembly with the "very-sensitive" "-k 10" settings. Peaks were called using Genrich version

689    0.6.1 (https://github.com/jsh58/Genrich). Differential binding analysis was performed using

690    Diffbind version 3.5.0 package in R version 4.1.0. ENCODE blacklisted regions were excluded

691    using hg38-blacklist.v2 (97) (https://github.com/Boyle-Lab/Blacklist). Phenotype specific regions

692    were isolated by first selecting for positive fold change open chromatin enrichment and then

693    using Intervene 0.6.5 (98) where regions were considered overlapping if they shared at least 1

694    bp with another phenotype. Regions with FDR adjusted p-values < 0.05 were then subset to

695    those overlapping the 3,380,000 established TFBSs (338 TFs x 10,000 binding sites, see Griffin

696    analysis for site selection) by at least 1 bp using BedTools v2.30.0 Intersect. Only regions that

697    overlapped an established TFBS from those lists were retained. For analyses restricted to

698    10,000, 1,000, or 100 sites, sites were ranked and chosen by adjusted p-value.

699    ***Nucleosome profiling of ctDNA***

700    Griffin is a method for profiling nucleosome protection and accessibility on predefined genomic

701    loci (51). For this study, Griffin (v0.1.0) was used and can be found on GitHub

702    (https://github.com/adoebley/Griffin/releases/tag/v0.1.0). The analysis was performed as

703    follows: First, GC bias was quantified for each sample using an approach described previously

704    (99). Briefly, for each possible fragment length and GC content, the number of reads in a bam

705    file and the number of genomic positions with that specific length and GC content were counted.

706    The GC bias for each fragment length and GC content was calculated by dividing the number of

707    observed reads by the number of observed genomic positions for that fragment length and GC

708    content. The GC bias for all possible GC contents at a given fragment length was then
709    normalized to a mean bias of 1. GC biases were then smoothed by taking the median of values
710    for fragments with similar lengths and GC contents (k nearest neighbors smoothing) to generate
711    smoothed GC bias values.

712    After GC correction, nucleosome profiling was performed in each sample. For each mappable
713    site of interest, fragments aligning to the region ± 5000 bp from the site were fetched from the
714    bam file. Fragments were filtered to remove duplicates and low-quality alignments (<20
715    mapping quality) and by fragment length. Nucleosome size fragments (140-250 bp) were
716    retained and used in all down-stream Griffin analyses. Fragments were then GC corrected by
717    assigning each fragment a weight of 1/GC_bias for that given fragment length and GC content.
718    The fragment midpoint was identified and the number of weighted fragment midpoints in 15bp
719    bins across the site were counted. For composite sites, all sites of a given type (such as all sites
720    for a given transcription factor) were summed together to generate a single coverage profile.
721    Individual or composite coverage profiles were normalized to a mean coverage of 1 in the ±
722    5000bp region surrounding the site.  Finally, sites were smoothed using a Savitsky-Golay filter
723    with a window length of 165bp and a polynomial order of 3. The window ± 1000 bp around the
724    site was retained for plotting and feature extraction when plotting sites, shading illustrates the
725    95% confidence interval within sample groups. Features extracted from individual or composite
726    sites included:

727    a) "mean central coverage," the mean coverage between -30 to 30 bp relative to the site
728       center,
729    b) "mean window coverage," the mean coverage between -990 to 990 bp relative to the site
730       center, and
731    c) "max wave height," the absolute difference between the minimum coverage within the
732       window from -120 to 30 bp and maximum coverage in the window from 31 to 195 bp
733       relative to the TSS.

734    ***Transcription factor binding site (TFBS) selection from GTRD***

735    TFBS identified using ChIP-seq were downloaded from the GTRD database version 19.10
736    (https://gtrd.biouml.org/downloads/19.10/chip-seq/Homo%20sapiens_meta_clusters.interval.gz).
737    This database contains binding sites (meta-clusters) that were observed in one or more ChIP
738    seq experiment. Low mappability sites were excluded by examining the mean mappability score
739    in a window around each site (+/- 5000 bp). Mappability information (hg38 Umap multi-read

740     mappability for 50bp reads) was obtained from UCSC genome browser (100)
741     (https://hgdownload.soe.ucsc.edu/gbdb/hg38/hoffmanMappability/k50.Umap.MultiTrackMappabi
742     lity.bw). Highly mappable sites (>0.95 mean mappability) were retained for further analysis. 338
743     TFs were selected for analysis using three criteria: (i) TF was contained in GTRD, (ii) had at
744     least 10,000 highly mappable binding sites on autosomes (chr1-22) in GTRD, (iii) TF was
745     present in the CIS-BP database (101) (CIS-BP v2.00 downloaded from
746     http://cisbp.ccbr.utoronto.ca/bulk.php) and had a known binding motif ('TF_status' is not N).
747     Unless otherwise noted, analyses utilized the top 1,000 TFBSs ranked by the highest
748     'peak.count' across all experiments as computed by GTRD (70). In addition, in the case of AR
749     and ASCL1 we also compared the top 1,000 vs the top 10,000 sites chosen with the same
750     'peak.count' criterion.

751     After intersecting these 338 TFs with the 404 differentially expressed TFs identified through
752     RNA-Seq 108 remained. On both the 108 and prostate-specific 41 TFs (described below) we
753     performed unsupervised hierarchical clustering of central window mean values (see Griffin
754     analysis). Hierarchical clustering was performed using the Ward.D2 method with Euclidean
755     distance and complete linkage settings; the groupings were determined using cutree_cols=2 for
756     columns (LuCaP CRPC phenotypes) and cutree_rows=13 for rows (TFs) on the dendrograms.

757     To generate a prostate lineage-specific TF set, we first merged GTRD metadata (file;
758     http://gtrd.biouml.org:8888/downloads/current/metadata/ChIP-seq.metadata.txt &
759     http://gtrd.biouml.org:8888/downloads/current/metadata/cell_types_and_tissues.metadata.txt).
760     We identified human prostate lineage-specific experiments by restricting the "species" field to
761     "Homo sapiens" and the "title" (tissue or cell type) field by performing a string match of the
762     following {"Prostate", "prostate", "LNCaP", "DU145", "PrEC"}. This resulted in a list of 1,086
763     prostate lineage ChIP seq experiments. Then, we selected metapeaks from the
764     "Homo_sapiens_meta_clusters.interval" file which had been observed in at least one of the
765     prostate lineage experiments using the "exp.set" field. This resulted in a set of 82 TFs. We then
766     filtered the peaks by mappability and kept only highly mappable peaks (as described above).
767     We excluded any TF that wasn't in the initial set of 338 TFs (this removed ChIP targets that
768     weren't true TFs, lacked a known binding site, or didn't have 10,000 total autosomal peaks in
769     GTRD). Of the remaining TFs, we analyzed those with 1,000 highly mappable peaks on
770     autosomes in prostate lineage experiments, resulting in 41 TFs. 20 out of 41 of these TFs
771     overlapped the list of 108 differentially expressed TFs by RNAseq of the PDX tumors. Note that
772     the top 1000 sites for each of the 41 TFs were different than in same TFs of the 338 set

28

773  because sites must meet the criteria of being derived from at least one experiment involving

774  prostate tissue or cell lines.

### Transcription factor binding site (TFBS) selection from other sources

776  For AR we further considered 17,619 sites identified through ChIP-seq by Pomerantz et al. (56)

777  (which overlapped 10.9% of the GTRD top 1,000 using bedtools), 41,633 sites identified by

778  Severson et al. (57) across four metastatic tumors (which overlapped 99.4% of the GTRD top

779  1,000). For ASCL1 we obtained 11,124 ChIP-seq sites from Cejas et al. (9) (which overlapped

780  60.9% of the GTRD 1,000). All of these site lists were lifted over from genome build GRCh37 to

781  GRCh38. No mappability filtering was applied so that all possible sites from these prostate

782  experiments and studies were considered.

### Phenotype-lineage specific gene marker selection

784  We selected 47 genes comprising 12 ARPC and 35 NEPC lineage markers established

785  previously (4,5,58,59) and confirmed by differential expression analysis from PDX tumor RNA-

786  Seq data (**Supplementary Table S3**). In tissues, AR and NE activities were measured on

787  lineage determinant signature gene's mRNA expression (GSVA score)(87). The 47 selected

788  gene list comprises the majority of these signature sets of genes defining mPC characteristic

789  phenotypes or phenotypic activities.

### Gene body and promoter region selection

791  For individual gene body and promoter analyses Ensembl BioMart v104 (hg38) (102) was used

792  to directly retrieve protein coding transcript start (TSS) and end (TES) coordinates. For promoter

793  region analysis the window ±1000 bp relative to the TSS was considered. For gene body

794  analysis, the region between the TSS and TES was considered. In the case of genes with

795  multiple transcripts, analyses were limited to the longest transcript resulting in 19,336 regions. In

796  downstream analysis of LuCaP PDX cfDNA, if any lines did not meet specific criteria in a region

797  (including differentially open histone modification regions) that feature/region combination was

798  excluded from analysis, leading to a variable lower number of regions considered based on the

799  feature. These criteria included requiring at least 10 total fragments in a region for all Fragment

800  size analysis (see below) and a non-zero number of "short" and "long" fragments for the short-

801  long ratio; short-long ratios less than 0.01 or greater than 10.0 were also excluded as outliers.

802  For Phasing analysis (see below) we also excluded amplitude components and thus NPS where

803  individual components were 0, or where the ratio was less than 0.01 or greater than 10.0,

804  indicative of insufficient coverage. In the case of mean phased nucleosome distance, if no

29

805     peaks were identified or the value in a region exceeded 500 (indicative of highly irregular/sparse

806     pileups also from low coverage) those regions were also excluded. Any region with no coverage

807     in a line was excluded from all analyses. This resulted in gene lists that differed in numbers

808     between genomic contexts and feature types.

809     ***Cell-free DNA fragment size analysis***

810     Fragments were first filtered to remove duplicates and low-quality alignments (<20 mapping

811     quality) and by fragment length (15-500 bp). In individual genomic loci/windows, we computed

812     the fragment short-long ratio (FSLR) as the ratio of short (15 - 120 bp) to long (140 - 250 bp)

813     fragments. We also calculated the mean, median absolute deviation (MAD: $median(|X_i -$

814     $median(X)|)$), and coefficient of variation (CV: $\frac{\sigma}{\mu}$ where σ = standard deviation, μ = mean) of the

815     fragment length distribution for each selected window. The fragment size analysis code and

816     implementation used in this study can be accessed at

817     https://github.com/GavinHaLab/CRPCSubtypingPaper/tree/main/FragmentAnalysis.

818     ***Nucleosome phasing analysis (TritonNP)***

819     Fragments were first filtered to remove duplicates and low-quality alignments (<20 mapping

820     quality) and by fragment length (nucleosome-sized: 140-250 bp). Next, we performed fragment-

821     level GC bias correction utilizing the same pre-processing method defined in Griffin. A band-

822     pass filter was then applied to the corrected coverage in each region of interest by taking the

823     Fast Fourier Transform (FFT) (scipy.fft v1.8.0) (103) and removing high-frequency components

824     corresponding to frequency components < 146 bp before reconstructing the signal. This cutoff

825     was chosen to ensure that periodic fit signal for downstream evaluation must come from the

826     minimum possible inter-nucleosome distance, thus excluding peak pileups that would not

827     indicate an overall trend in nucleosome phasing. Local peak calling was then done on the

828     smoothed signal to infer average inter nucleosome distance or "phased nucleosome distance"

829     by finding maxima directly. To quantify clarity of overall phasing we took the average frequency

830     amplitude in two bands corresponding to stably bound, well-phased nucleosomes (180-210 bp)

831     and a baseline (150-180 bp), with the former measuring the strength of typically aligned

832     nucleosomes and the latter giving a measure of the underlying signal strength not coming from

833     either high frequency noise or low frequency shifts in total coverage. The ratio of these two

834     amplitude averages forms the Nucleosome Phasing Score (NPS). Because peak locations are

835     assumed to be independent of copy number alterations or depth, and the NPS by virtue of being

836     a ratio divides out any confounding DNA/depth variation between sites, both features are taken

837    as agnostic of CNAs or variable depth. Code and implementation of the method can be found at

838    https://github.com/denniepatton/TritonNP.

839    ***ctDNA tumor-normal admixtures and benchmarking***

840    Admixtures for evaluating benchmarking performance were constructed using 5 ARPC (LuCaP

841    35, 35CR, 58, 92, 136CR) and 5 NEPC (LuCaP 49, 93, 145.2, 173.1, 208.4) lines mixed to 1%,

842    5%, 10%, 20%, and 30% tumor fraction with a single healthy donor plasma line (NPH023,

843    EGAD00001005343) for use in binary classification (50 admixes), and in mixtures of 1%, 3%,

844    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

845    1.0 in all possible combinations (810 admixes) for mixture model evaluation. All admixes were

846    mixed at ~25X mean coverage, assuming 100% tumor fraction in post-mouse subtracted PDX

847    sequencing data. After extracting chromosomal DNA (chr1-22, X, Y) with SAMtools v1.14 (104)

848    and removing duplicates with Picard (https://broadinstitute.github.io/picard/), SAMtools was

849    used to merge BAM files. To evaluate the ultra-low pass WGS performance, admixtures were

850    then down-sampled using SAMtools to the number of reads corresponding to 1X or 0.2X. During

851    unsupervised benchmarking of each admixture, the healthy donor and the LuCaP line used in

852    the admixture were excluded from the generation of feature distributions to ensure the model

853    would not learn from the lines being interrogated. The admixture pipeline used in this study can

854    be accessed at https://github.com/GavinHaLab/Admixtures_snakemake.

855    ***Supervised binary classification of ARPC and NEPC***

856    Binary classification of ARPC and NEPC subtypes using individual region and feature

857    combinations was conducted using XGBoost v1.4.2 'XGBClassifier' implemented in Python with

858    default parameters. Features included NPS and Mean Phased Nucleosome Distance (see

859    Phasing analysis) in histone modification regions, promoters, and gene bodies; fragment size

860    mean, short-long ratio, and coefficient of variation (see Fragment size analysis) in histone

861    modification regions, promoters, and gene bodies; central and window coverage (see Griffin

862    analysis) in promoters, composite TFBSs, and composite differentially open chromatin regions

863    identified through ATAC-Seq; and Max Wave Height (See Griffin analysis) in promoters. We

864    applied stratified 6-fold cross-validation where two ARPC samples and one NEPC sample were

865    held out in each fold. This was repeated 100 times and performance was computed using area

866    under the receiver operating characteristic (ROC) curve (AUC) and 95% confidence intervals for

867    each individual feature and region combination. Code and implementation of the method can be

868    found at https://github.com/GavinHaLab/CRPCSubtypingPaper/tree/main/SupervisedLearning.

31

869 ***Tumor fraction estimation***

870 Tumor fractions from patient plasma samples were assessed using ichorCNA (80) with binSize
871 1,000,000 bp and both GRCh37 and GRCh38 reference genomes. Default tumor fraction
872 estimates reported by ichorCNA were used. See
873 https://github.com/GavinHaLab/CRPCSubtypingPaper/tree/main/ichorCNA_configuration for
874 complete configuration settings.

875 ***Phenotype class prediction model (ctdPheno)***

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

$$\mu'_{i,k} = \alpha\mu_{i,k} + (1 - \alpha)\mu_{i,H}$$

$$\sigma'_{i,k} = \sqrt{\alpha\sigma^2_{i,k} + (1 - \alpha)\sigma^2_{i,H}}$$

887 where $\alpha$ is the tumor fraction estimate for each test sample. In the final model, four features
888 were used: composite open chromatin regions (central and window mean coverage) for specific
889 phenotypes (ARPC and NEPC) identified from the LuCaP PDX ATAC-Seq analysis using Griffin
890 (see Griffin analysis). For each feature $i$, we then found the probability that the observed sample
891 came from a mixture of the tumor-fraction-corrected Gaussian distributions, where $\theta$ is the
892 NEPC mixture weight:

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

893 The $\theta$ parameter is estimated by maximizing the joint log-likelihood $L$ for a given patient sample:

$$\theta' = \underset{\theta}{\mathrm{argmax}}[L(x|\theta)]$$

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

32

894 $\theta$ has range [0,1], where higher values indicate an increased probability of the sample having a

895 NEPC phenotype and was used as the NEPC prediction score metric. Code and implementation

896 of the method can be found at

897 https://github.com/GavinHaLab/CRPCSubtypingPaper/tree/main/ctdPheno.

898

899 ***Phenotype heterogeneity prediction and quantification (Keraon)***

900 We developed an analytical model to directly estimate the contributing fractions of ctDNA from

901 different mCRPC phenotypes (ARPC and NEPC) in individual patient plasma ctDNA samples.

902 Like ctdPheno this model is unsupervised and does not require training on the patient cohort of

903 interest. However, the model accepts the pre-estimated tumor fraction from ichorCNA for the

904 given patient ctDNA sample, as well as the pre-computed ctDNA features values from the

905 LuCaP PDX ctDNA and healthy donor ctDNA as prior information (see Class phenotype

906 prediction model).

907

908 As a pre-processing step, the model first computes the mean vector $\mu_i$ and covariance matrix $\Sigma_i$

909 for each anchor class *i* in K, under the assumption that each subtype (including healthy) fits a

910 multivariate Gaussian distribution. Based on model constraints, K - 1 non-correlated features

911 fully specify the system, and so for ARPC:NEPC:Healthy (K = 3) fraction estimation we limited

912 analyses to sets of two features of interest (F = 2).

913

914 Next, for each sample defined by some location in feature space ***v*** and estimated tumor fraction

915 *t* we first performed a change of basis to translate the sample's location from feature space to

916 class space, where each (not necessarily orthogonal) axis defined a single phenotype, and the

917 origin represented pure healthy. If F = K -1, this was accomplished by solving the determined,

918 linear matrix equation for the shifted basis components **X**:

$$BX = S$$

919 Where **B** = $[\mu_{i \neq HD} - \mu_{HD}]$ is the matrix defining all basis vectors from the healthy mean anchor to

920 each phenotype mean anchor, and **S** is the vector from the healthy mean anchor to the sample

921 of interest, **S** = ***v*** - $\mu_{HD}$. If the system is overdetermined (F > K – 1), least squares was used to

922 estimate the approximate solution. This step allows us to learn where in the class space the

923 sample lies, which determined how estimates were evaluated:

33

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 are 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 $BC_{i \neq HD} = \frac{X_i}{\sum X} t$ for all $i$ such that $X_i > 0$.

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

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,

34

954 the concordance was computed between the ctdPheno NEPC prediction score and the cfMeDIP

955 NEPC score obtained from the original study using a 0.15 threshold (25).

956 We then validated the model on two cohorts, beginning with the already published DFCI cohort

957 II (75,76,80). We restricted our analysis to eleven samples from six patients with matched ULP-

958 WGS and WGS data with paired-end reads. Tumor fraction estimates from ichorCNA were

959 obtained from the original study (80). All samples were considered adenocarcinoma (ARPC)

960 based on clinical histories (see Human subjects). The scoring threshold of 0.3314, determined

961 from DFCI cohort I was used for phenotype classification.

962 For the *UW cohort*, consisting of 47 samples from 27 patients (average 22.13X depth of

963 coverage sequencing), ichorCNA was used to estimate sample tumor fractions as described

964 above (GRCh38), while clinical phenotype was determined from clinical histories and expert

965 chart review. We evaluated model performance on matched ULP-WGS and WGS data for

966 unambiguous clinical phenotypes of ARPC and NEPC. The chosen scoring threshold of 0.3314

967 was used, and the fraction of correctly predicted ARPC (n=26) and NEPC (n=5) was computed.

968 The remaining 16 samples with mixed histologies were not evaluated for performance in

969 ctdPheno.

970 Phenotype prediction and proportion estimation: Keraon does not require de novo threshold

971 selection, and so all clinical cohorts were treated as validation sets. Based on the Mean

972 Absolute Error (MAE) of 2.8% for estimating NEPC fraction garnered in the heterogenous

973 mixture benchmarking, this value was chosen as the minimum NEPC fraction threshold for

974 calling presence of NEPC in WGS cohorts. The same tumor fraction estimates used by

975 ctdPheno in ULP were utilized by Keraon, with standard classification conducted on pure clinical

976 phenotypes. The 16 samples with mixed phenotypes in the UW cohort were evaluated both

977 qualitatively and based on the 2.8% threshold in the absence of quantifiable burden estimates

978 from histories.

979 **STATISTICAL ANALYSIS**

980 Quantification of and statistical approaches for high-throughput sequencing data analysis are

981 described in the methods above. When non-parametric distributions (not normally distributed) of

982 numerical values of a particular parameter in a population were compared (using boxplots or in

983 tables), the two-tailed Mann-Whitney U test (also known as the Wilcoxon Rank Sum test;

984 scipy.stats.mannwhitneyu, (103) was used to test if any two distributions being compared were

985 significantly different, with Benjamini-Hochberg (statsmodels.stats.multitest.fdrcorrection,

986 https://www.statsmodels.org) correction applied in multiple testing scenarios. All boxplots

987 represent the median with a centerline, interquartile range (IQR) with a box, and first quartile –

988 1.5 IQR and third quartile + 1.5 IQR with whiskers. PCA was conducted in Python

989 (sklearn.decomposition.PCA; https://scikit-learn.org)

990 **DATA AVAILABILITY**

991 The LuCaP patient derived xenograft (PDX) plasma ctDNA sequencing data generated in this

992 study can be accessed under NCBI BioProject accession PRJN900550

993 (https://www.ncbi.nlm.nih.gov/bioproject/PRJNA90055). The processed patient plasma data can

994 be accessed at https://github.com/GavinHaLab/CRPCSubtypingPaper/tree/main/Data. The raw

995 sequencing data generated for the UW cohort are not publicly available because patients did

996 not consent to genomic data sharing but are available upon reasonable request from the

997 corresponding author. This paper also analyzes existing, publicly available data, including

998 LuCaP PDX RNA-Seq (GSE199596) and ATAC-Seq data (GSE156292). The CUT&RUN

999 processed data can be accessed at https://github.com/nielOnav/LuCaP_nucleosome_profile.

1000 Published data for DFCI Cohort I was obtained from the authors (25) after establishing a data

1001 use agreement with the Dana-Farber Cancer Institute.

1002 Any additional information required to reanalyze the data reported in this paper is available from

1003 the lead contact upon request.

**REFERENCES**

1.  Karantanos T, Corn PG, Thompson TC. Prostate cancer progression after androgen deprivation therapy: mechanisms of castrate resistance and novel therapeutic approaches. Oncogene. Nature Publishing Group; 2013;32:5501–11.

2.  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. Massachusetts Medical Society; 2013;368:138–48.

3.  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. New England Journal of Medicine. 2012;367:1187–97.

4.  Beltran H, Prandi D, Mosquera JM, Benelli M, Puca L, Cyrta J, et al. Divergent clonal evolution of castration-resistant neuroendocrine prostate cancer. Nature Medicine. Nature Publishing Group; 2016;22:298–305.

5.  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. Elsevier; 2017;32:474-489.e6.

6.  Conteduca V, Oromendia C, Eng KW, Bareja R, Sigouros M, Molina A, et al. Clinical features of neuroendocrine prostate cancer. European Journal of Cancer. 2019;121:7–18.

7.  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. American Society of Clinical Oncology; 2018;36:2492–503.

8.  Baca SC, Takeda DY, Seo J-H, Hwang J, Ku SY, Arafeh R, et al. Reprogramming of the FOXA1 cistrome in treatment-emergent neuroendocrine prostate cancer. Nat Commun. Nature Publishing Group; 2021;12:1979.

9.  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;12:5775.

10. Spetsieris N, Boukovala M, Patsakis G, Alafis I, Efstathiou E. Neuroendocrine and Aggressive-Variant Prostate Cancer. Cancers. Multidisciplinary Digital Publishing Institute; 2020;12:3792.

11. 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. American Society for Clinical Investigation; 2019;129:4492–505.

12. 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.

37

1042   13.   Liu Y, Horn JL, Banda K, Goodman AZ, Lim Y, Jana S, et al. The androgen receptor
1043         regulates a druggable translational regulon in advanced prostate cancer. Science
1044         Translational Medicine. American Association for the Advancement of Science;
1045         2019;11:eaaw4993.

1046   14.   Epstein JI, Amin MB, Beltran H, Lotan TL, Mosquera J-M, Reuter VE, et al. Proposed
1047         Morphologic Classification of Prostate Cancer With Neuroendocrine Differentiation. The
1048         American Journal of Surgical Pathology. 2014;38:756–67.

1049   15.   Annala M, Taavitsainen S, Khalaf DJ, Vandekerkhove G, Beja K, Sipola J, et al. Evolution
1050         of Castration-Resistant Prostate Cancer in ctDNA during Sequential Androgen Receptor
1051         Pathway Inhibition. Clinical Cancer Research. 2021;27:4610–23.

1052   16.   Aparicio AM, Shen L, Tapia ELN, Lu J-F, Chen H-C, Zhang J, et al. Combined Tumor
1053         Suppressor Defects Characterize Clinically Defined Aggressive Variant Prostate Cancers.
1054         Clinical Cancer Research. 2016;22:1520–30.

1055   17.   Carreira S, Romanel A, Goodall J, Grist E, Ferraldeschi R, Miranda S, et al. Tumor clone
1056         dynamics in lethal prostate cancer. Science translational medicine. 2014;6:254ra125.

1057   18.   Du M, Tian Y, Tan W, Wang L, Wang L, Kilari D, et al. Plasma cell-free DNA-based
1058         predictors of response to abiraterone acetate/prednisone and prognostic factors in
1059         metastatic castration-resistant prostate cancer. Prostate Cancer Prostatic Dis. Nature
1060         Publishing Group; 2020;23:705–13.

1061   19.   Sumanasuriya S, Seed G, Parr H, Christova R, Pope L, Bertan C, et al. Elucidating
1062         Prostate Cancer Behaviour During Treatment via Low-pass Whole-genome Sequencing of
1063         Circulating Tumour DNA. European Urology. 2021;80:243–53.

1064   20.   Ulz P, Belic J, Graf R, Auer M, Lafer I, Fischereder K, et al. Whole-genome plasma
1065         sequencing reveals focal amplifications as a driving force in metastatic prostate cancer.
1066         Nat Commun. Institute of Human Genetics, Medical University of Graz, A-8010 Graz,
1067         Austria. Department of Urology, Medical University of Graz, A-8036 Graz, Austria.
1068         Department of Internal Medicine I, Hospital Barmherzige Schwestern Linz, A-4020 Linz,
1069         Austria. Departme; 2016;7:12008.

1070   21.   Wyatt AW, Annala M, Aggarwal R, Beja K, Feng F, Youngren J, et al. Concordance of
1071         Circulating Tumor DNA and Matched Metastatic Tissue Biopsy in Prostate Cancer. JNCI:
1072         Journal of the National Cancer Institute. Oxford University Press; 2018;110:78–86.

1073   22.   Nyquist MD, Corella A, Coleman I, De Sarkar N, Kaipainen A, Ha G, et al. Combined TP53
1074         and RB1 Loss Promotes Prostate Cancer Resistance to a Spectrum of Therapeutics and
1075         Confers Vulnerability to Replication Stress. Cell Reports. 2020;31:107669.

1076   23.   Berger A, Brady NJ, Bareja R, Robinson B, Conteduca V, Augello MA, et al. N-Myc–
1077         mediated epigenetic reprogramming drives lineage plasticity in advanced prostate cancer.
1078         J Clin Invest. American Society for Clinical Investigation; 2019;129:3924–40.

1079   24.   Beltran H, Romanel A, Conteduca V, Casiraghi N, Sigouros M, Franceschini GM, et al.
1080         Circulating tumor DNA profile recognizes transformation to castration-resistant

1081     neuroendocrine prostate cancer. J Clin Invest. American Society for Clinical Investigation;
1082     2020;130:1653–68.

1083  25.  Berchuck JE, Baca SC, McClure HM, Korthauer K, Tsai HK, Nuzzo PV, et al. Detecting
1084     Neuroendocrine Prostate Cancer Through Tissue-Informed Cell-Free DNA Methylation
1085     Analysis. Clinical Cancer Research. 2022;28:928–38.

1086  26.  Shen SY, Singhania R, Fehringer G, Chakravarthy A, Roehrl MHA, Chadwick D, et al.
1087     Sensitive tumour detection and classification using plasma cell-free DNA methylomes.
1088     Nature. Nature Publishing Group; 2018;563:579–83.

1089  27.  Wu A, Cremaschi P, Wetterskog D, Conteduca V, Franceschini GM, Kleftogiannis D, et al.
1090     Genome-wide plasma DNA methylation features of metastatic prostate cancer. J Clin
1091     Invest. American Society for Clinical Investigation; 2020;130:1991–2000.

1092  28.  Heitzer E, Auinger L, Speicher MR. Cell-Free DNA and Apoptosis: How Dead Cells Inform
1093     About the Living. Trends in Molecular Medicine. Elsevier Ltd; 2020;26:519–28.

1094  29.  Lo YMD, Han DSC, Jiang P, Chiu RWK. Epigenetics, fragmentomics, and topology of cell-
1095     free DNA in liquid biopsies. Science [Internet]. American Association for the Advancement
1096     of Science; 2021 [cited 2021 Apr 12];372. Available from:
1097     https://science.sciencemag.org/content/372/6538/eaaw3616

1098  30.  Cristiano S, Leal A, Phallen J, Fiksel J, Adleff V, Bruhm DC, et al. Genome-wide cell-free
1099     DNA fragmentation in patients with cancer. Nature. Nature Publishing Group;
1100     2019;570:385–9.

1101  31.  Jiang P, Sun K, Peng W, Cheng SH, Ni M, Yeung PC, et al. Plasma DNA End-Motif
1102     Profiling as a Fragmentomic Marker in Cancer, Pregnancy, and Transplantation. Cancer
1103     Discov. American Association for Cancer Research; 2020;10:664–73.

1104  32.  Mathios D, Johansen JS, Cristiano S, Medina JE, Phallen J, Larsen KR, et al. Detection
1105     and characterization of lung cancer using cell-free DNA fragmentomes. Nat Commun.
1106     2021;12:5060.

1107  33.  Peneder P, Stütz AM, Surdez D, Krumbholz M, Semper S, Chicard M, et al. Multimodal
1108     analysis of cell-free DNA whole-genome sequencing for pediatric cancers with low
1109     mutational burden. Nat Commun. Nature Publishing Group; 2021;12:3230.

1110  34.  Zhu G, Guo YA, Ho D, Poon P, Poh ZW, Wong PM, et al. Tissue-specific cell-free DNA
1111     degradation quantifies circulating tumor DNA burden. Nature Communications. Nature
1112     Publishing Group; 2021;12:2229.

1113  35.  Herberts C, Annala M, Sipola J, Ng SWS, Chen XE, Nurminen A, et al. Deep whole-
1114     genome ctDNA chronology of treatment-resistant prostate cancer. Nature. 2022;608:199–
1115     208.

1116  36.  Jiang P, Chan CWM, Chan KCA, Cheng SH, Wong J, Wong VW-S, et al. Lengthening and
1117     shortening of plasma DNA in hepatocellular carcinoma patients. Proceedings of the
1118     National Academy of Sciences of the United States of America. 2015;112:E1317-25.

1119 37. Mouliere F, Chandrananda D, Piskorz AM, Moore EK, Morris J, Ahlborn LB, et al.
1120      Enhanced detection of circulating tumor DNA by fragment size analysis. Science
1121      Translational Medicine. 2018;10:eaat4921.

1122 38. Snyder MW, Kircher M, Hill AJ, Daza RM, Shendure J. Cell-free DNA Comprises an in
1123      Vivo Nucleosome Footprint that Informs Its Tissues-Of-Origin. Cell. Elsevier Inc.;
1124      2016;164:57–68.

1125 39. Underhill HR, Kitzman JO, Hellwig S, Welker NC, Daza R, Baker DN, et al. Fragment
1126      Length of Circulating Tumor DNA. PLOS Genet. 2016;12:426–37.

1127 40. Ramachandran S, Ahmad K, Henikoff S. Transcription and Remodeling Produce
1128      Asymmetrically Unwrapped Nucleosomal Intermediates. Molecular Cell. Cell Press;
1129      2017;68:1038-1053.e4.

1130 41. Ulz P, Thallinger GG, Auer M, Graf R, Kashofer K, Jahn SW, et al. Inferring expressed
1131      genes by whole-genome sequencing of plasma DNA. Nature Genetics. Nature Publishing
1132      Group; 2016;48:1273–8.

1133 42. Ulz P, Perakis S, Zhou Q, Moser T, Belic J, Lazzeri I, et al. Inference of transcription factor
1134      binding from cell-free DNA enables tumor subtype prediction and early detection. Nature
1135      Communications. 2019;10:4666.

1136 43. Brahma S, Henikoff S. Epigenome Regulation by Dynamic Nucleosome Unwrapping.
1137      Trends in Biochemical Sciences. Elsevier; 2020;45:13–26.

1138 44. Lai WKM, Pugh BF. Understanding nucleosome dynamics and their links to gene
1139      expression and DNA replication. Nat Rev Mol Cell Biol. 2017;18:548–62.

1140 45. Yen K, Vinayachandran V, Batta K, Koerber RT, Pugh BF. Genome-wide Nucleosome
1141      Specificity and Directionality of Chromatin Remodelers. Cell. 2012;149:1461–73.

1142 46. Rao S, Han AL, Zukowski A, Kopin E, Sartorius CA, Kabos P, et al. Transcription factor–
1143      nucleosome dynamics from plasma cfDNA identifies ER-driven states in breast cancer.
1144      Science Advances. American Association for the Advancement of Science;
1145      2022;8:eabm4358.

1146 47. Nguyen HM, Vessella RL, Morrissey C, Brown LG, Coleman IM, Higano CS, et al. LuCaP
1147      Prostate Cancer Patient-Derived Xenografts Reflect the Molecular Heterogeneity of
1148      Advanced Disease an--d Serve as Models for Evaluating Cancer Therapeutics. The
1149      Prostate. John Wiley & Sons, Ltd; 2017;77:654–71.

1150 48. Skene PJ, Henikoff S. An efficient targeted nuclease strategy for high-resolution mapping
1151      of DNA binding sites. Reinberg D, editor. eLife. eLife Sciences Publications, Ltd;
1152      2017;6:e21856.

1153 49. Meers MP, Tenenbaum D, Henikoff S. Peak calling by Sparse Enrichment Analysis for
1154      CUT&RUN chromatin profiling. Epigenetics & Chromatin. 2019;12:42.

1155 50. Zhou VW, Goren A, Bernstein BE. Charting histone modifications and the functional
1156      organization of mammalian genomes. Nat Rev Genet. 2011;12:7–18.

40

1157    51.   Doebley A-L, Ko M, Liao H, Cruikshank AE, Kikawa C, Santos K, et al. Griffin: Framework
1158          for clinical cancer subtyping from nucleosome profiling of cell-free DNA. medRxiv.
1159          2021;2021.08.31.21262867.

1160    52.   Soares LM, He PC, Chun Y, Suh H, Kim T, Buratowski S. Determinants of Histone H3K4
1161          Methylation Patterns. Molecular Cell. 2017;68:773-785.e6.

1162    53.   Brady NJ, Bagadion AM, Singh R, Conteduca V, Van Emmenis L, Arceci E, et al. Temporal
1163          evolution of cellular heterogeneity during the progression to advanced AR-negative
1164          prostate cancer. Nat Commun. Nature Publishing Group; 2021;12:3372.

1165    54.   Wang YA, Sfakianos J, Tewari AK, Cordon-cardo C, Kyprianou N. Molecular tracing of
1166          prostate cancer lethality. Oncogene. Nature Publishing Group; 2020;39:7225–38.

1167    55.   Rapa I, Ceppi P, Bollito E, Rosas R, Cappia S, Bacillo E, et al. Human ASH1 expression in
1168          prostate cancer with neuroendocrine differentiation. Mod Pathol. Nature Publishing Group;
1169          2008;21:700–7.

1170    56.   Pomerantz MM, Qiu X, Zhu Y, Takeda DY, Pan W, Baca SC, et al. Prostate cancer
1171          reactivates developmental epigenomic programs during metastatic progression. Nat Genet.
1172          2020;52:790–9.

1173    57.   Severson TM, Zhu Y, De Marzo AM, Jones T, Simons JW, Nelson WG, et al. Epigenetic
1174          and transcriptional analysis reveals a core transcriptional program conserved in clonal
1175          prostate cancer metastases. Molecular Oncology. 2021;15:1942–55.

1176    58.   Labrecque MP, Brown LG, Coleman IM, Lakely B, Brady NJ, Lee JK, et al. RNA Splicing
1177          Factors SRRM3 and SRRM4 Distinguish Molecular Phenotypes of Castration-Resistant
1178          Neuroendocrine Prostate Cancer. Cancer Research. 2021;81:4736–50.

1179    59.   Tsai HK, Lehrer J, Alshalalfa M, Erho N, Davicioni E, Lotan TL. Gene expression
1180          signatures of neuroendocrine prostate cancer and primary small cell prostatic carcinoma.
1181          BMC Cancer. 2017;17:759.

1182    60.   Jiang Z, Zhang B. On the role of transcription in positioning nucleosomes. PLOS
1183          Computational Biology. Public Library of Science; 2021;17:e1008556.

1184    61.   Klemm SL, Shipony Z, Greenleaf WJ. Chromatin accessibility and the regulatory
1185          epigenome. Nature Reviews Genetics. Nature Publishing Group; 2019;20:207–20.

1186    62.   Oruba A, Saccani S, van Essen D. Role of cell-type specific nucleosome positioning in
1187          inducible activation of mammalian promoters. Nat Commun. 2020;11:1075.

1188    63.   Guo Y, Zhao S, Wang GG. Polycomb Gene Silencing Mechanisms: PRC2 Chromatin
1189          Targeting, H3K27me3 "Readout", and Phase Separation-Based Compaction. Trends in
1190          Genetics. Elsevier; 2021;37:547–65.

1191    64.   Jiang C, Pugh BF. Nucleosome positioning and gene regulation: advances through
1192          genomics. Nat Rev Genet. 2009;10:161–72.

1193    65.    Saxton DS, Rine J. Nucleosome Positioning Regulates the Establishment, Stability, and
1194           Inheritance of Heterochromatin in Saccharomyces cerevisiae. Proceedings of the National
1195           Academy of Sciences. Proceedings of the National Academy of Sciences;
1196           2020;117:27493–501.

1197    66.    Valouev A, Johnson SM, Boyd SD, Smith CL, Fire AZ, Sidow A. Determinants of
1198           nucleosome organization in primary human cells. Nature. Nature Publishing Group;
1199           2011;474:516–20.

1200    67.    Deal RB, Henikoff JG, Henikoff S. Genome-Wide Kinetics of Nucleosome Turnover
1201           Determined by Metabolic Labeling of Histones. Science. American Association for the
1202           Advancement of Science; 2010;328:1161–4.

1203    68.    Cuzick J, Swanson GP, Fisher G, Brothman AR, Berney DM, Reid JE, et al. Prognostic
1204           value of an RNA expression signature derived from cell cycle proliferation genes in
1205           patients with prostate cancer: a retrospective study. The Lancet Oncology. 2011;12:245–
1206           55.

1207    69.    Chereji RV, Bryson TD, Henikoff S. Quantitative MNase-seq accurately maps nucleosome
1208           occupancy levels. Genome Biology. 2019;20:198.

1209    70.    Yevshin I, Sharipov R, Kolmykov S, Kondrakhin Y, Kolpakov F. GTRD: a database on
1210           gene transcription regulation—2019 update. Nucleic Acids Res. Oxford Academic;
1211           2019;47:D100–5.

1212    71.    Arora VK, Schenkein E, Murali R, Subudhi SK, Wongvipat J, Balbas MD, et al.
1213           Glucocorticoid Receptor Confers Resistance to Antiandrogens by Bypassing Androgen
1214           Receptor Blockade. Cell. Elsevier; 2013;155:1309–22.

1215    72.    Mu P, Zhang Z, Benelli M, Karthaus WR, Hoover E, Chen C-C, et al. SOX2 promotes
1216           lineage plasticity and antiandrogen resistance in TP53- and RB1-deficient prostate cancer.
1217           Science. American Association for the Advancement of Science; 2017;355:84–8.

1218    73.    Shukla S, Cyrta J, Murphy DA, Walczak EG, Ran L, Agrawal P, et al. Aberrant Activation of
1219           a Gastrointestinal Transcriptional Circuit in Prostate Cancer Mediates Castration
1220           Resistance. Cancer Cell. Elsevier; 2017;32:792-806.e7.

1221    74.    Sun K, Jiang P, Cheng SH, Cheng THT, Wong J, Wong VWS, et al. Orientation-aware
1222           plasma cell-free DNA fragmentation analysis in open chromatin regions informs tissue of
1223           origin. Genome research. Cold Spring Harbor Laboratory Press; 2019;29:418–27.

1224    75.    Viswanathan SR, Ha G, Hoff AM, Wala JA, Carrot-Zhang J, Whelan CW, et al. Structural
1225           Alterations Driving Castration-Resistant Prostate Cancer Revealed by Linked-Read
1226           Genome Sequencing. Cell. Elsevier; 2018;174:433-447.e19.

1227    76.    Choudhury AD, Werner L, Francini E, Wei XX, Ha G, Freeman SS, et al. Tumor fraction in
1228           cell-free DNA as a biomarker in prostate cancer. JCI Insight [Internet]. American Society
1229           for Clinical Investigation; 2018 [cited 2019 Mar 1];3. Available from:
1230           https://insight.jci.org/articles/view/122109

42

1231    77.    Klein DC, Hainer SJ. Genomic methods in profiling DNA accessibility and factor
1232            localization. Chromosome Res. 2020;28:69–85.

1233    78.    Chaytor L, Simcock M, Nakjang S, Heath R, Walker L, Robson C, et al. The Pioneering
1234            Role of GATA2 in Androgen Receptor Variant Regulation Is Controlled by Bromodomain
1235            and Extraterminal Proteins in Castrate-Resistant Prostate Cancer. Mol Cancer Res.
1236            American Association for Cancer Research; 2019;17:1264–78.

1237    79.    Corces MR, Granja JM, Shams S, Louie BH, Seoane JA, Zhou W, et al. The chromatin
1238            accessibility landscape of primary human cancers. Science. 2018;362.

1239    80.    Adalsteinsson VA, Ha G, Freeman SS, Choudhury AD, Stover DG, Parsons HA, et al.
1240            Scalable whole-exome sequencing of cell-free DNA reveals high concordance with
1241            metastatic tumors. Nature Communications. 2017;8.

1242    81.    Fang R, Preissl S, Li Y, Hou X, Lucero J, Wang X, et al. Comprehensive analysis of single
1243            cell ATAC-seq data with SnapATAC. Nat Commun. Nature Publishing Group;
1244            2021;12:1337.

1245    82.    Wu SJ, Furlan SN, Mihalas AB, Kaya-Okur HS, Feroze AH, Emerson SN, et al. Single-cell
1246            CUT&Tag analysis of chromatin modifications in differentiation and tumor progression. Nat
1247            Biotechnol. 2021;39:819–24.

1248    83.    Lam H-M, Nguyen HM, Corey E. Generation of Prostate Cancer Patient-Derived
1249            Xenografts to Investigate Mechanisms of Novel Treatments and Treatment Resistance. In:
1250            Culig Z, editor. Prostate Cancer: Methods and Protocols [Internet]. New York, NY:
1251            Springer; 2018 [cited 2022 Mar 22]. page 1–27. Available from:
1252            https://doi.org/10.1007/978-1-4939-7845-8_1

1253    84.    Li H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
1254            arXiv:13033997 [q-bio] [Internet]. 2013 [cited 2022 Mar 22]; Available from:
1255            http://arxiv.org/abs/1303.3997

1256    85.    Jo S-Y, Kim E, Kim S. Impact of mouse contamination in genomic profiling of patient-
1257            derived models and best practice for robust analysis. Genome Biology. 2019;20:231.

1258    86.    DePristo MA, Banks E, Poplin R, Garimella KV, Maguire JR, Hartl C, et al. A framework for
1259            variation discovery and genotyping using next-generation DNA sequencing data. Nat
1260            Genet. Nature Publishing Group; 2011;43:491–8.

1261    87.    Hänzelmann S, Castelo R, Guinney J. GSVA: gene set variation analysis for microarray
1262            and RNA-Seq data. BMC Bioinformatics. 2013;14:7.

1263    88.    Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast
1264            universal RNA-seq aligner. Bioinformatics. 2013;29:15–21.

1265    89.    Kluin RJC, Kemper K, Kuilman T, de Ruiter JR, Iyer V, Forment JV, et al. XenofilteR:
1266            computational deconvolution of mouse and human reads in tumor xenograft sequence
1267            data. BMC Bioinformatics. 2018;19:366.

1268 90. Robinson MD, McCarthy DJ, Smyth GK. edgeR: a Bioconductor package for differential
1269      expression analysis of digital gene expression data. Bioinformatics. 2010;26:139–40.

1270 91. Lambert SA, Jolma A, Campitelli LF, Das PK, Yin Y, Albu M, et al. The Human
1271      Transcription Factors. Cell. 2018;172:650–65.

1272 92. Langmead B, Wilks C, Antonescu V, Charles R. Scaling read aligners to hundreds of
1273      threads on general-purpose processors. Bioinformatics. 2019;35:421–32.

1274 93. Ramírez F, Ryan DP, Grüning B, Bhardwaj V, Kilpert F, Richter AS, et al. deepTools2: a
1275      next generation web server for deep-sequencing data analysis. Nucleic Acids Research.
1276      2016;44:W160–5.

1277 94. Ross-Innes CS, Stark R, Teschendorff AE, Holmes KA, Ali HR, Dunning MJ, et al.
1278      Differential oestrogen receptor binding is associated with clinical outcome in breast cancer.
1279      Nature. 2012;481:389–93.

1280 95. Quinlan AR, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic features.
1281      Bioinformatics. 2010;26:841–2.

1282 96. Yu G, Wang L-G, He Q-Y. ChIPseeker: an R/Bioconductor package for ChIP peak
1283      annotation, comparison and visualization. Bioinformatics. 2015;31:2382–3.

1284 97. Amemiya HM, Kundaje A, Boyle AP. The ENCODE Blacklist: Identification of Problematic
1285      Regions of the Genome. Sci Rep. Nature Publishing Group; 2019;9:9354.

1286 98. Khan A, Mathelier A. Intervene: a tool for intersection and visualization of multiple gene or
1287      genomic region sets. BMC Bioinformatics. 2017;18:287.

1288 99. Benjamini Y, Speed TP. Summarizing and correcting the GC content bias in high-
1289      throughput sequencing. Nucleic Acids Research. 2012;40:e72–e72.

1290 100. Karimzadeh M, Ernst C, Kundaje A, Hoffman MM. Umap and Bismap: quantifying genome
1291       and methylome mappability. Nucleic Acids Research. 2018;46:e120.

1292 101. Weirauch MT, Yang A, Albu M, Cote AG, Montenegro-Montero A, Drewe P, et al.
1293       Determination and inference of eukaryotic transcription factor sequence specificity. Cell.
1294       2014;158:1431–43.

1295 102. Howe KL, Achuthan P, Allen J, Allen J, Alvarez-Jarreta J, Amode MR, et al. Ensembl 2021.
1296       Nucleic Acids Research. 2021;49:D884–91.

1297 103. Virtanen P, Gommers R, Oliphant TE, Haberland M, Reddy T, Cournapeau D, et al. SciPy
1298       1.0: fundamental algorithms for scientific computing in Python. Nat Methods. Nature
1299       Publishing Group; 2020;17:261–72.

1300 104. Danecek P, Bonfield JK, Liddle J, Marshall J, Ohan V, Pollard MO, et al. Twelve years of
1301       SAMtools and BCFtools. GigaScience. 2021;10:giab008.

**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 patient-derived xenograft (PDX) mouse models with tumors originating from metastatic castration-resistant prostate cancer (mCRPC) with AR-positive adenocarcinoma (ARPC), neuroendocrine (NEPC), AR-low neuroendocrine-negative (ARLPC) phenotypes. Cell-free DNA (cfDNA) was extracted from pooled plasma collected from 4-8 mice and whole genome sequencing (WGS) was performed. Following bioinformatic mouse read subtraction, pure human circulating tumor DNA (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 (51) 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 analytical 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 (< 3x coverage were 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 20.5x; dotted line). Percentage of the cfDNA sample that contains human ctDNA after mouse read subtraction.

**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 (70) (**Methods**).

(B-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 1000 bootstraps for a subset of sites (shading) are shown. The region ±150 bp is indicated with vertical dotted line and yellow shading.

(D) Heatmap of $\log_2$ fold change in 47 key genes up and down regulated between ARPC and NEPC established through RNA-Seq (**left**) grouped by the type of histone modification which dictates translation levels: Group 1 shows genes 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' cfDNA fragment size coefficient of variation (CV) is shown for TSS+/- 1KB windows and respective gene bodies.

(E) Comparison of the $\log_2$ fold change (ARPC/NEPC) of mean mRNA expression vs mean coefficient of variation (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**) ±1000 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 1000 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 transcription start site (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 UNC13A has increased nucleosome phasing in ctDNA of ARPC samples compared to NEPC.

46

**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 (grey dotted lines). Frequency components corresponding to nucleosome dyads (wavelength > 146 bp) are shown in purple. The mean inter-nucleosome distance was computed from all peaks in the gene body: lower values represent more periodic and stable nucleosomes. Nucleosome Phasing Score (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 shown.

**(C)** Comparison of the mean phased-nucleosome distance and the mean cell-cycle progression (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 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 non-significant 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. Heatmap 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 (log2-fold-change < -0.05, $p < 0.05$) are indicated with blue bars. Text color indicates relative expression between ARPC and NEPC PDX tumors by RNAseq shown for TFs with significant differential accessibility.

**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-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 (70). Coverage profile means (lines) and 95% confidence interval with 1000 bootstraps (shading) are shown. The region ±150 bp is indicated with 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 (coefficient of variation) 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 receiver operating characteristic curve (AUC) with 95% confidence interval (100 repeats of stratified cross validation) is shown for performance of all feature types.
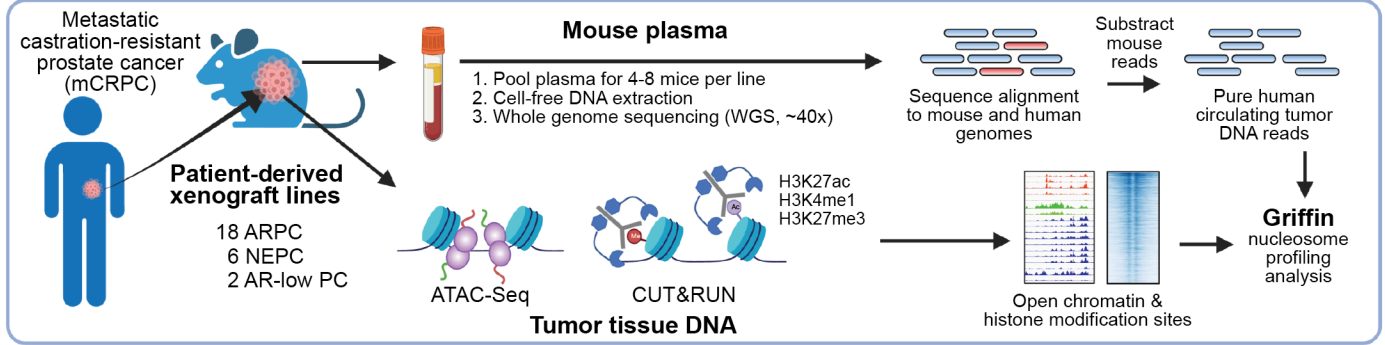
1415 **Figure 5. Accurate classification and estimation of prostate cancer in patient plasma**
1416 **samples**

1417 (A) Schematic illustration of the ctdPheno classification method. Griffin-derived features and
1418 ichorCNA tumor fraction estimates from patient plasma samples are combined in a
1419 probabilistic framework informed by PDX models to predict the presence of NEPC.

1420 (B) Performance for classification on admixtures samples using ctdPheno. Five ctDNA
1421 admixtures were generated for each phenotype from PDX lines, each at various
1422 sequencing coverages and tumor fractions. In total, 125 admixtures were evaluated. The
1423 mean AUC across the 5 admixtures is shown for each configuration.

1424 (C) Receiver operating characteristic (ROC) curve for 101 mCRPC patients (DFCI cohort I)
1425 with ultra-low-pass WGS (ULP-WGS) data. The optimal performance of 90.4%
1426 sensitivity (for predicting NEPC) and 97.5% specificity (for predicting ARPC)
1427 corresponding to a prediction score cutoff of 0.3314 is indicated with horizontal and
1428 vertical dotted lines, respectively.

1429 (D) Prediction scores from ctdPheno for 47 ULP-WGS plasma samples with clinical
1430 phenotypes comprising 26 ARPC (blue), 5 NEPC (red), and 16 mixed or ambiguous
1431 phenotypes (purple, triangles), including double-negative prostate cancer (DNPC; grey).
1432 The 0.3314 score cutoff threshold (dotted line) was used for classifying NEPC and
1433 ARPC. Tumor fractions were estimated by ichorCNA from WGS data.

1434 (E) Schematic illustration of the Keraon mixture estimation method. Griffin-derived features
1435 from PDX lines and healthy donors define a known feature space, which is transformed
1436 based on Griffin features and ichorCNA tumor fraction estimates for each patient plasma
1437 sample. Based on the patient's location in the transformed phenotype space, fractions of
1438 each phenotype are inferred directly.

1439 (F) Illustration of mixture simulations. 5 ARPC and 5 NEPC PDX samples were combined in
1440 the ratios shown with a single healthy donor at the tumor fractions shown, for a total of
1441 810 mixed phenotype samples at 25x for evaluating mixture proportions with Keraon.

1442 (G) Boxplot of predicted total NEPC fraction in 810 simulated mixed-phenotype samples
1443 using Keraon, Pearson's r = 0.884. Median absolute error (MAE) was computed as the
1444 median absolute difference between estimated and expected NEPC fraction across all
1445 samples.

1446 (H) Fractional phenotype estimates for 47 WGS plasma samples with clinical phenotypes
1447 comprising 26 ARPC (blue), 5 NEPC (red), and 16 mixed or ambiguous phenotypes
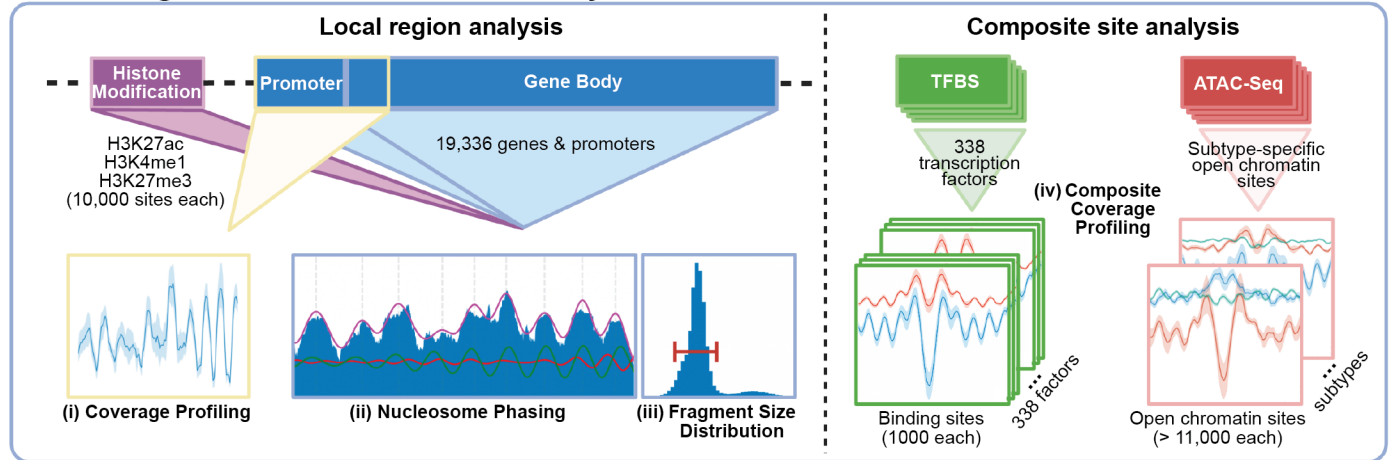
49

1448        (purple, triangles), including double-negative prostate cancer (DNPC; grey). The 2.8%

1449        NEPC fraction threshold indicates the predicted presence of NEPC (dotted line).

1450

**A**

## PDX Plasma and Tumor Sequencing

Metastatic castration-resistant prostate cancer (mCRPC)

Patient-derived xenograft lines

18 ARPC
6 NEPC
2 AR-low PC

**Mouse plasma**

1. Pool plasma for 4-8 mice per line
2. Cell-free DNA extraction
3. Whole genome sequencing (WGS, ~40x)

Sequence alignment to mouse and human genomes

Substract mouse reads

Pure human circulating tumor DNA reads

**Griffin** nucleosome profiling analysis

ATAC-Seq    CUT&RUN

H3K27ac
H3K4me1
H3K27me3

**Tumor tissue DNA**

Open chromatin & histone modification sites

## Circulating Tumor DNA Feature Discovery

### Local region analysis

Histone Modification    Promoter    Gene Body

H3K27ac
H3K4me1
H3K27me3
(10,000 sites each)

19,336 genes & promoters

**(i) Coverage Profiling**    **(ii) Nucleosome Phasing**    **(iii) Fragment Size Distribution**

### Composite site analysis

TFBS

338 transcription factors

**(iv) Composite Coverage Profiling**

ATAC-Seq

Subtype-specific open chromatin sites

Binding sites (1000 each)    338 factors

Open chromatin sites (> 11,000 each)    subtypes

## Circulating Tumor DNA Characterization

**Inferring gene regulation in prostate cancer**

AR genes
NE genes

**Assessment of epigenetic modifications & transcriptional activity**

**Nucleosome phasing analysis**

Stability
Baseline

NEPC
ARPC

PC2
PC1

**Explore and evaluate features for phenotype characterization**

**Machine-learning models for prostate cancer phenotype prediction**

True Positive Rate
False Positive Rate

## Patient Phenotyping

**3 mCRPC patient cohorts**

Whole genome sequencing (WGS, 0.1x - 30x)

Subtype-specific open chromatin sites

**Griffin** features

**Models for phenotype prediction**

***ichorCNA*** tumor fraction

$\pi \rightarrow Z$
$x_i$
$\mu_i \quad \sigma_i$
$\alpha$

PDX-informed prior data

***ctdPheno*** (classification)    ***Keraon*** (mixtures)

$BX = S$
$B_0 \rightarrow B \quad S$

**B**

Status: Included / Excluded

Phenotype: NEPC / ARPC / ARLPC

Average depth of coverage

Mouse genome subtraction: Before / After

% of sample

Mouse fraction / Human fraction

LuCaP PDX lines: 93, 208.4, 145.2, 145.1, 173.1, 49, 235.3, 136CR, 35, 35CR, 92, 136, 96, 70, 170.2, 189.4, 70CR, 78, 58, 81, 174.1, 189.3, 176CR, 176, 86.2CR, 86.2

**Figure 1**

**Figure 2**

**A**

HOXB13 (chr17: 46.802Mb - 46.806Mb)

LuCaP 93 (NEPC)
Phased-nucleosome distance: 249 bp
NPS: 1.93

LuCaP 136 (ARPC)
Phased-nucleosome distance: 332 bp
NPS: 1.63

Gene coordinates (bp)

**Fragment coverage** — Grey: Fragment coverage, Dotted: Phased local peaks

**Nucleosome frequency components**
- Phased (observed) (>146 bp)
- Baseline (150-180 bp)
- Stable (180-210 bp)

**Phased-Nucleosome Distance**
[Mean phased peak distance]
Periodic (stable) ← → Aperiodic (unstable)

**Nucleosome Phasing Score (NPS)**
[Stable amplitude / Baseline amplitude]
Inactive/repressed transcription ← → Active transcription

**B** Gene body (ctDNA) (17,946 genes)
p = 0.027
Mean phased-nucleosome distance (bp)
ARPC (n=16), NEPC (n=6)

**C**
Spearman's ρ = -0.563, p = 0.006
Mean phased-nucleosome distance (bp) vs Mean CCP score
ARPC, NEPC

**D** Gene body (ctDNA)
AR genes p=0.070, NE genes p=0.134
Mean NPS
ARPC (n=16), NEPC (n=6), ARPC (n=16), NEPC (n=6)

**E** Gene body (ctDNA)
Phenotype marker genes: ARPC (12), NEPC (35); Significant, N.S.
-log₁₀ adjusted q-value vs NPS log2-fold-change (ARPC / NEPC)
TMPRSS2, CHRNB2, ASCL1, NEUROD1, GRP, XKR7, MYCL, CHGB

**F**
Prostate tumor phenotypes: NEPC, ARPC, ARLPC
Normalized composite TFBS coverage
log2 fold change
Relative expression in tumors (ARPC:NEPC): Up-reg., Down-reg., N.S.

**Figure 3**

**A** ATAC-Seq sites

28,765 ARPC sites
21,963 NEPC sites

log₂ fold change (NEPC vs. ARPC)

- Significant ARPC-specific sites
- Significant NEPC-specific sites
- N.S.

**B** Open chromatin sites in ARPC
(15,879 ATAC-Seq sites)

Distance to binding site (bp)

LuCaP ctDNA / Healthy donor cfDNA
ARPC (n=16)  NEPC (n=6)  HD (n=14)

**C** Open chromatin sites in NEPC
(11,692 ATAC-Seq sites)

Distance to binding site (bp)

LuCaP ctDNA / Healthy donor cfDNA
ARPC (n=16)  NEPC (n=6)  HD (n=14)

**D**

Composite TFBS coverage (74 significant TFs)
PC 1 (49.87%) / PC 2 (22.30%)

Nucleosome Phasing Score (47 prostate genes)
PC 1 (28.9%) / PC 2 (15.14%)

Fragment size variability (H3K4me1 sites)
PC 1 (79.65%) / PC 2 (3.88%)

ARPC (n=16)  NEPC (n=6)

**E**

Area under the ROC curve (AUC)

Histone modifications: H3K4me1, H3K27ac, H3K27me3
Promoter (TSS), Gene body
Transcription factors (TFBS)
Open chromatin sites

Local region analysis    Composite site analysis

Nucleosome phasing
- NPS
- Nucleosome distance

Fragment size
- Mean
- Coefficient of variation
- Short:Long ratio

Coverage profiling
- Central coverage (+/- 30bp)
- Overall coverage (+/- 1kb)

Promoter coverage
- Max wave height (TSS)

**Figure 4**

**Figure 5**