

1 **Genetic ancestry inference from cancer-**
2 **derived molecular data across genomic**
3 **and transcriptomic platforms**

4
5 Pascal Belleau^{1,2}, Astrid Deschênes^{2,3}, Nyasha Chambwe⁴, David A. Tuveson^{2,3}, and
6 Alexander Krasnitz^{1,2}

7 ¹Simons Center for Quantitative Biology, Cold Spring Harbor Laboratory, Cold Spring Harbor, New
8 York, USA; ²Cancer Center, Cold Spring Harbor Laboratory, Cold Spring Harbor, New York, USA;
9 ³Lustgarten Foundation Pancreatic Cancer Research Laboratory, Cold Spring Harbor, New York,
10 USA; ⁴Institute of Molecular Medicine, Feinstein Institutes for Medical Research, Northwell Health,
11 Manhasset, New York, USA

12
13 Corresponding Author: Alexander Krasnitz, Simons Center for Quantitative Biology, Cold Spring Harbor Laboratory,
14 Cold Spring Harbor, New York, USA, 11724; E-mail: krasnitz@cshl.edu; Telephone: 516-3676863

15
16 **Running Title**

17 Ancestry inference from cancer-derived molecular data

18
19 **Keywords**

20 GENETICS OF CANCER RISK & OUTCOME/GENETICS OF CANCER RISK & OUTCOME, COMPUTATIONAL
21 METHODS/Algorithms, COMPUTATIONAL METHODS/Sequence analysis, COMPUTATIONAL METHODS/Software,
22 genetic ancestry

23
24 **Conflicts of interest**

25 The authors declare that they have no competing interests.

26
27 **Word count**

28 5042

29
30 **Total number of figures**

31 5

32
33 **Total number of tables**

34 2

35 **Significance statement**

36

37 The development of a computational approach that enables accurate and robust ancestry inference from cancer-
38 derived molecular profiles without matching cancer-free data provides a valuable methodology for genetic ancestry-
40 oriented cancer research.

41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61

Abstract Genetic ancestry-oriented cancer research requires the ability to perform accurate and robust genetic ancestry inference from existing cancer-derived data, including whole exome sequencing, transcriptome sequencing, and targeted gene panels, very often in the absence of matching cancer-free genomic data. Here we examined the feasibility and accuracy of computational inference of genetic ancestry relying exclusively on cancer-derived data. A data synthesis framework was developed to optimize and assess the performance of the ancestry inference for any given input cancer-derived molecular profile. In its core procedure, the ancestral background of the profiled patient is replaced with one of any number of individuals with known ancestry. The data synthesis framework is applicable to multiple profiling platforms, making it possible to assess the performance of inference specifically for a given molecular profile and separately for each continental-level ancestry; this ability extends to all ancestries, including those without statistically sufficient representation in the existing cancer data. The inference procedure was demonstrated to be accurate and robust in a wide range of sequencing depths. Testing of the approach in four representative cancer types and across three molecular profiling modalities showed that continental-level ancestry of patients can be inferred with high accuracy, as quantified by its agreement with the gold standard of deriving ancestry from matching cancer-free molecular data. This study demonstrates that vast amounts of existing cancer-derived molecular data are potentially amenable to ancestry-oriented studies of the disease without requiring matching cancer-free genomes or patient self-reported ancestry.

62 Introduction

63 There is ample epidemiological evidence that race and/or ethnicity are important determinants of
64 incidence, clinical course and outcome in multiple types of cancer (1-5). As such, these categories
65 must be taken into account in the analysis of molecular data derived from cancer. A number of
66 recently published large-scale genomic studies of cancer point to differences in the molecular make-
67 up of the disease among groups of different ancestral background and to the need for more
68 molecular data to power discovery of such differences (6-11).

69 Ancestry annotation of cancer-derived data largely draws on two sources: patient's self-identified
70 race and/or ethnicity (SIRE) and patient's cancer-free genotype. SIRE is often missing, sometimes
71 inaccurate and usually incomplete. As a recent analysis (12) of PubMed database entries since 2010
72 reveals, patients' SIRE is massively under-reported in genome and exome sequencing studies of
73 cancer, with only 37% of these reporting race, and 17% reporting ethnicity. Furthermore, SIRE is
74 not always consistent with genetic ancestry. Finally, a self-declaring patient is often given a choice
75 from a small number of broad racial or ethnic categories, which fail to capture complete ancestral
76 information, especially in cases of mixed ancestry (13).

77 A far more accurate and detailed ancestral characterization may be obtained by genotyping a
78 patient's DNA from a cancer-free tissue. Powerful methods exist for ancestry inference from
79 germline DNA sequence (14-17). These methods were recently used to determine ancestry of
80 approximately 10,000 patients profiled by The Cancer Genome Atlas (TCGA) (7,11). However,
81 genotyping of DNA from patient-matched cancer-free specimens is not part of standard clinical
82 practice, where the purpose of DNA profiling is often identification of mutations with known
83 oncogenic effects, such as those in the Catalog Of Somatic Mutations In Cancer (COSMIC)
84 database (18). As a result, it is not performed routinely outside academic clinical centers or major
85 research projects. There also are studies yielding sequence data from tumors, whose purpose does
86 not require germline profiling. RNA sequencing (RNA-seq) for expression quantification is in this
87 category. Finally, peripheral blood is most often the source of germline DNA in the clinic, but this is
88 not always the case for diseases of the hematopoietic system, such as leukemia, wherein cancer
89 cells are massively present in circulation. In summary, matched germline DNA sequence is not
90 universally available for cancer-derived molecular data. In such cases, it is necessary to infer ancestry
91 from the nucleic acid sequence of the tumor itself.

92 Standard methods of ancestry inference commonly rely on population specificity of germline
93 single-nucleotide variants (SNV). Whole-genome (WGS) or whole-exome sequences (WES), at depths
94 sufficient for reliably calling single-nucleotide variants, and readouts from genotyping microarrays,
95 are therefore data types most suitable for this purpose. However, such detailed DNA profiling is
96 often not performed in molecular studies of cancer. In such cases, it is necessary to infer ancestry
97 from other types of tumor-derived data, including RNA sequence and DNA sequence for a small
98 panel of genes, e.g., FoundationOne® CDx (19).

99 For all types of tumor-derived sequence, accurate inference of ancestry is a potential challenge.
100 Tumor genome is often replete with somatic alterations, including loss of heterozygosity (LOH),
101 copy number variants (CNV), translocations, microsatellite instabilities and SNV. These alterations
102 interfere with germline genotyping of the patient that is used as input for inference of genetic
103 ancestry. Structural variants, especially LOH and CNV, are the most likely to affect the germline
104 genotyping, and thereby the genetic ancestry calls. This effect is especially clearly seen in the case of
105 LOH, as a result of which heterozygous genotypes are transformed into homozygous, but other
106 types of alterations also are, to various degrees, potential obstacles to accurate ancestry inference.
107 Tumor RNA-seq presents additional challenges, namely, extremely uneven coverage of the
108 transcript due to a broad range of RNA expression levels and distortions due to allele-specific
109 expression. Gene panels represent a very small fraction of the genome, whose sufficiency for
110 ancestry inference is not clear and may vary from panel to panel. In addition, cancer gene panels
111 are enriched in cancer driver genes, which tend to undergo somatic alteration more frequently than
112 other parts of the genome.

113 Important recent publications on ancestral effects in cancer reported patient ancestry inferred

114 from matching cancer-free DNA (7,8,11). At the same time, there has been much less work on
115 ancestry inference from tumor-derived nucleic acids (7,11,20-23). Collectively, this work
116 demonstrates the feasibility of accurate genetic ancestry inference from cancer-derived DNA
117 profiled by SNP arrays or by high-coverage gene panels, such as the FoundationOne® CDx gene
118 panel (19). However, to our knowledge, no systematic computational framework for ancestry
119 inference from cancer-derived molecular data, across assay and cancer types, has been developed
120 to date. There is presently no ability to assess the inference accuracy specifically for a given input
121 tumor-derived molecular profile with all its attendant properties, including the data quality and the
122 depth of coverage. Reliable and accurate ancestry inference from tumor-derived nucleic acids thus
123 represents an unmet need, which the present work aims to address.

124 For this purpose, we designed an inference procedure having in mind a scenario, likely to occur in
125 studies of existing data or of archived tissue specimens, with an input molecular profile of a tumor
126 from a single patient, and no matching cancer-free sequence available. The profile in question may
127 have its unique set of sequence properties. These include the target sequence and uniformity of its
128 coverage depth, read length and sequencing quality. These profile-specific properties may be vastly
129 dissimilar from those in the available public data sets with reliably known genetic ancestry of the
130 patients. Furthermore, not all ancestries are equally easy to infer: for example, an American
131 ancestral category is sometimes difficult to distinguish either from African or from European
132 ancestry. This profile specificity would make it impossible to confidently assess the accuracy of the
133 inference procedure for the input profile from its performance with the public cancer-derived data in
134 aggregate. In order to overcome this difficulty, we developed a computational technique, which is
135 described schematically in *Figure 1*, wherein the ancestral background of the patient is supplanted in
136 the input profile by one of an unrelated individual with known ancestry. A similar data synthesis
137 procedure was employed in our prior work in a different genomic context (24). We next apply
138 established methods of ancestry inference to this synthetic profile and compare the result to that
139 known ancestry. Generating multiple such synthetic profiles allows us to assess how accurate the
140 ancestry inference is for the patient, both overall and as a function of the profile's continental-level
141 ancestry. Furthermore, using synthetic data, we are able to optimize the inference procedure with
142 respect to parameters on which it depends. Importantly, this assessment and optimization
143 procedure does not require the profile in question to be part of a larger data set from a cohort of
144 patients with a similar diagnosis. Very often in existing cancer-derived data, such cohorts do not
145 provide statistically meaningful representation of non-European ancestries. This insufficiency is not
146 an impediment to the application of our methodology.

147 In the following, we assess the accuracy of global ancestry calls from tumor exomes, narrowly
148 targeted gene panels and RNA sequences, in comparison to such calls from matching germline
149 genotypes, as profiled by exome sequencing or genotyping microarrays. We do so for four cancer
150 types, namely, pancreatic adenocarcinoma (PDAC), ovarian cystadenocarcinoma (OV) and breast
151 carcinoma (BRCA) as representative types of epithelial tumors, and acute myeloid leukemia (AML), as
152 an example of hematopoietic malignancy. Each of these data sets was chosen because it presents a
153 challenge for patients' ancestry inference and/or an opportunity to test our approach. Specifically, OV
154 is characterized by massive copy number alterations, often spanning much of the genome. Our
155 PDAC data originate from patient-derived organoid (PDO) models of the disease (25). In PDO, near-
156 100% tumor purity is achieved, exacerbating effects of copy number loss and loss of heterozygosity
157 on the sequence. In BRCA, a large patient cohort size makes it possible for us to choose an
158 ancestrally diverse subset of the data for testing our methods. In AML the peripheral blood, the
159 usual source of cancer-free DNA, may be severely contaminated by the cancer.

160 **Methods and Materials**

161 **Data sets and pre-processing**

162 The data sets used in this work originate from four sources: TCGA collection for ovarian
163 cystadenocarcinoma (26) (TCGA-OV), an ancestrally diverse subset of TCGA collection for breast
164 carcinoma (27) (TCGA-BRCA), Beat AML clinical trial (28) (Beat AML), and a study of pancreatic ductal
165 adenocarcinoma using patient-derived organoids (25) (PDAC). For all four, the data used are

166 summarized in the form of Venn diagrams in **Figure 2A-D** and tabulated in Supplementary Table S1.
167 These data include cancer DNA (whole-exome or whole-genome) sequence, cancer RNA sequence and
168 matching normal DNA (whole-exome or whole-genome) sequence. As explained in the following,
169 genetic ancestry inferred from the latter was used as the ground truth in assessing the performance
170 of ancestry inference from the cancer-derived data cohort-wide for each of the four cohorts. Also
171 available for comparison was the donor SIRE, as depicted in **Figure 2E**. In addition, published
172 genetic ancestry calls from matching cancer-free genotypes, representing a consensus of five
173 inference pipelines (C5), were available for comparison with our findings for the TCGA-OV and
174 TCGA-BRCA cohorts (7).

175 Throughout the study, we used the 1000 Genomes (1KG) data set, with no relatives for the
176 individuals included (29-31), as reference, against which patient molecular data were compared to
177 infer continental-level global ancestry. The latter is defined as a categorical variable taking five values:
178 African (AFR), East Asian (EAS), European (EUR), American (AMR) and South Asian (SAS). These are
179 called super-populations in the 1KG terminology. Each super-population comprises a number of
180 subcontinental-level populations, as explained in the 1000 Genomes consortium publications (31).
181 The composition of the 1KG data, as used in this study, is summarized in Supplementary Table S2.

182 In all cases, read data mapped to the hg38 version of the human genome were used. In order to
183 study ancestry inference from targeted panels, the cancer-derived whole-exome data were reduced
184 to reads mapping to the FoundationOne® CDx cancer-related gene panel (19). The pre-processing
185 is illustrated in the first part of the **Figure 3**. Reads in the cancer patient-derived data were filtered
186 for quality using a cutoff phred score of 20. Following this filter, single-nucleotide substitutions were
187 called at all positions with read coverage of at least 10, using snp-pileup in FACETS (32) and Varscan
188 version 2.4.4 (33). This set of positions is called the high-confidence substitution (HCS) set in the
189 following. From the 1000 Genomes (1KG) variant call data in the Variant Call Format (VCF) (34),
190 genomic positions where substitution variants occur at a frequency of at least 0.01 in at least one of
191 the super-populations comprising 1KG were selected as a basis for the ancestry inference. This set
192 is referred to as the high-frequency substitution (HFS) set in the following. The genotype was called
193 at the HFS positions in the cancer-derived profile with the coverage above 10. This subset of the HFS
194 positions is referred to as high-confidence genotype (HCG) set in the following. In the HCG set, the
195 total read count and the read counts for the reference and the alternative (according to HFS) alleles
196 were determined. A genotype at an HCG position was considered undetermined if the excess of the
197 total read count over the sum of the reference and alternative counts was inconsistent with the error
198 of 0.001 at the $p = 0.001$ level of significance. The same rule was used to call a heterozygous
199 genotype. The HCG genomic positions were pruned to reduce correlation between neighboring
200 genotypes using Bioconductor SNPRelate package version 1.22.0 (35), resulting in the pruned high-
201 confidence genotype (PHCG) set of positions.

202 **Ancestry inference**

203 **Figure 3** lays out the workflow for ancestry inference. For a given cancer-derived profile, principal
204 component analysis of the 1KG genotypes reduced to the PHCG was performed, and D top principal
205 components retained. The patient genotype reduced to PHCG was projected onto the subspace
206 spanned by these D components. Within this subspace, the patient's ancestry was called as that of
207 the 1KG super-population with the highest number of 1KG individuals among K nearest neighbors of
208 the patient's genotype, using Euclidean distance in the D -dimensional subspace. If two or more
209 super-populations were found tied in the nearest-neighbor count, no ancestry call was made for the
210 patient. Only two such ties were observed in this work.

212 **Measures of performance**

213 We evaluate the performance of the ancestry inference by comparison to the ancestry inferred from
214 the matching cancer-free data, wherever the latter are available. This is the case for the entirety of
215 Beat AML, TCGA-OV and TCGA-BRCA data. For all three, we infer the ancestry from the matching
216 cancer-free exome profiles. In the case of TCGA-OV and TCGA-BRCA data, we also compare the
217 results to the consensus ancestry calls (7).

218 In the case of PDAC matching cancer-free WGS data are available for 22 patient cases (**Figure 2**),
219 and our assessment of accuracy is based on this subset of the data. We compute, for each dataset,
220 the 5×5 confusion matrix (CM) for the 1KG superpopulation calls from the cancer-derived and
221 cancer-free data sources. From the CM, the call accuracy is computed as the sum of the diagonal
222 terms divided by that of the whole CM. Since the ancestral composition of all data sets considered
223 here except TCGA-BRCA is heavily skewed towards the European super-population, we also
224 compute the multi-class version of the area under the receiver operating characteristic curve
225 (AUROC) (36). AUROC is a measure of the call quality which compensates for the asymmetry in
226 the class sizes. We use an R package pROC (CRAN version 1.16.2) (37) for this purpose, and
227 compute both the class-specific AUROC for each super-population and the 5-class overall AUROC.
228 In the class-specific case, we use a version DeLong's algorithm (38,39) as implemented in the pROC
229 package to compute the AUROC confidence intervals. In the overall 5-class case the confidence
230 intervals are computed using bootstrap with 100-fold sampling.

231 **Data synthesis**

232 Data synthesis is defined here as replacement of PHCG genotypes in a cancer-derived profile P by
233 those found in the genome of an unrelated individual U . Ingredients required for this procedure are:
234 (a) allele fraction (AF) estimates in P , as explained in detail in the Supplementary Methods and
235 illustrated in Figure S1; and (b) the haplotype of U in the portion of the genome covered by P . With this
236 knowledge, the procedure, depicted in **Figure 4**, consists of the following steps. First, sequence reads
237 comprising P are distributed at random among the alleles with probabilities equal to the observed
238 allele fractions. Second, in each haplotype block in the genome of U that is covered by P , allele
239 assignment is made at random, yielding variant and reference read counts for each PHCG
240 substitution in the genome of U within the scope of P .

242 **Inference parameter optimization using synthetic data**

243 In order to optimize ancestry inference parameters D and K for a given cancer-derived molecular
244 profile, we generate a synthetic data set by repeatedly pairing the profile with 1KG genomes. A
245 subset of 780 1KG genomes is set aside for this purpose by drawing at random 30 genomes from
246 each of the 26 ancestral populations represented in 1KG. Genetic ancestry is then inferred for each of
247 the 780 synthetic profiles following the procedure described in the Ancestry Inference subsection,
248 each time with the 1KG genome used for synthesis removed from the reference data set. The
249 inference performance is then assessed as the 5-class AUROC, as explained in the Measures of
250 Performance subsection. AUROC is computed for the D, K pairs in a range of values of these
251 parameters, and the optimal D, K pairs yielding the highest accuracy are identified. Throughout this
252 work, AUROC was computed for all D and K in the rectangle $3 \leq D \leq 11$; $3 \leq K \leq 15$. For all
253 combinations of data sources and profiling modalities considered, a set of D, K pairs was found
254 where the performance was optimal or differed from the optimum by no more than 3% (**Figure 5**).

255 **Down-sampling of sequence data**

256 In order to down-sample the sequence data to a desired fraction f of the original coverage, we
257 sampled reads from the original patient profile P with the Bernoulli probability f without
258 replacement. The ancestry inference procedure was then performed with the resulting sample of
259 reads.

261 **Software used in making figures**

262 All diagrams were made using draw.io version 15.7.3 (<http://www.diagrams.net>). The Venn diagrams
263 in **Figure 2** were produced with CRAN packages VennDiagram version 1.7.3 (40) and multipanelfigure
264 version 2.1.2 (41). The bar plot in **Figure 2** and the plots in **Figure 5** were made using packages ggplot2
265 (version 3.3.6, RRID: SCR_014601) and cowplot (version 1.1.1, RRID: SCR_018081).

267 **Software and data availability**

268 Ancestry inference methods introduced in this work are implemented in an R language package
269 RAIDS (Robust Ancestry Inference using Data Synthesis) is publicly available, under the Apache-2.0

270 license, at <https://github.com/KrasnitzLab/RAIDS>. Documentation for this software is available at
271 <https://krasnitzlab.github.io/RAIDS/>. The data analyzed in this study were obtained from the
272 National Center for Biotechnology (NCBI) database of Genotypes and Phenotypes (dbGaP)
273 archive under accession numbers [phs001611.v1.p1](#), [phs001657.v1.p1](#) and [phs000178.v11.p8](#).
274

275 Results

276 We assessed the performance of genetic ancestry inference from three genomic data types: whole
277 exomes, gene panels targeting exomes of several hundred cancer-related genes each and RNA
278 sequences. Our assessment relied on molecular data collected from four patient cohorts, each
279 representing a cancer type, namely, tissue donors to the Cold Spring Harbor Laboratory (CSHL)
280 pancreatic ductal adenocarcinoma (PDAC) library of patient-derived organoids; acute myeloid
281 leukemia (AML) patients enrolled in Beat AML clinical trial; patients comprising TCGA ovarian cancer
282 cohort (TCGA-OV) (26) and a subset of TCGA breast cancer cohort (TCGA-BRCA). Throughout the
283 study we used the 1000 Genomes (1KG) genotype collection as our population reference.

284 As explained in detail in the Methods and Materials section, for inference of genetic ancestry we
285 employed principal-component analysis (PCA) in combination with K -nearest-neighbor
286 classification. For a subset of patients in each cohort we individually assessed the performance of
287 the ancestry inference, as a function of the parameters K and D , the number of principal dimensions
288 retained. We relied on data synthesis for this assessment. Both super-population-specific and overall
289 AUROC values were computed in a range of D , K pairs, as illustrated in **Figure 5** for 10 PDAC
290 patients and AMR-specific AUROC and in **Figure S2** for all other cohorts and super-populations.
291 Optimal D , K pairs maximizing the overall AUROC were chosen. From this subset of patients we
292 observed, for each cancer type considered and for each of the three molecular profiling modalities,
293 an optimal range of D and K parameters where the performance of inference was consistently high in
294 the subset and only weakly dependent on these parameters (**Figure S2**). For all four tumor types, our
295 overall performance findings using data synthesis are summarized in **Tables S3-S6**. We then
296 selected and used, for the remainder of the patients with this cancer type and for this profiling
297 modality, a pair D and K values from within the optimal range. As an additional validation of our
298 parameter optimization procedure, we applied it to cancer-free WES profiles of TCGA-OV and
299 TCGA-BRCA patients included in this study. Comparing the resulting ancestry calls to the consensus
300 calls (C5) by TCGA (7), we find the two to be in good agreement (**Tables S7-S10**).

301 We also assessed the cohort-wide performance of our ancestry calls from the original cancer-
302 derived molecular data, by comparison to the gold standard of ancestry as determined from the
303 matching cancer-free genotypes. For Beat AML, TCGA-OV and TCGA-BRCA patients, we performed
304 ancestry inference from cancer-free patient exomes, using the same methodology as we did for the
305 cancer-derived sequences of these patients. In the case of PDAC, cancer-free whole-genome
306 sequencing data were available, and used for the same purpose for a portion of the patient cohort.
307 For all four cohorts, we summarize our cohort-wide findings in **Table 1**. We also used the C5
308 ancestry calls (7) in our performance assessment for TCGA-OV and TCGA-BRCA and found close
309 agreement for both these cohorts (**Tables S7-S10**).

310 We note that in all patient cohorts we analyze here except TCGA-BRCA (**Table 2** and **Table**
311 **S11**) the sampling of patients with non-European ancestries is statistically insufficient for a purely
312 cohort-based assessment of performance (**Table S12-S14**). We therefore report cohort-wide overall
313 but not super-population specific AUROC values for Beat AML, TCGA-OV and TCGA-BRCA. Using data
314 synthesis, we are able to compensate for this data shortfall in non-European ancestries and
315 estimate super-population specific AUROC, as explained above (**Tables S15-S18** and **Figure S2**). We do
316 report super-population-specific AUROC for TCGA-BRCA and for the aggregate of all four cohorts.

317 The results of our analysis as presented in **Tables S15-S18**, lead to the following key
318 observations. First, we demonstrate a consistently high performance of our inference procedure
319 across all cohorts and profiling modalities. Second, the super-population specific performance was
320 the highest for the European and both Asian super populations. The slightly lower accuracy as
321 observed for the African and American super-populations is likely due to a greater genetic variability
322 within the African super-population and to a higher degree of (the predominantly European)

323 admixture in both super-populations. Third, the optimal choice of the D, K inference parameters, in
324 general, depends on an individual cancer-derived molecular profile, even within the same cancer
325 type and profiling modality (Figure S2 B,G,L). Full results of our inferential analysis for the patients
326 in all four cohorts are compiled in Table S19.

327 In order to examine whether our inference procedure is robust against variation in the sequence
328 target coverage, we re-computed the ancestry calls for a subset of ten TCGA-OV patients, with the
329 cancer-derived whole-exome and RNA sequences of these patients down-sampled to between 75%
330 and 10% of the original coverage. The results, presented in (Figure S3) exhibit no substantial
331 sensitivity of the inference accuracy to the depth of coverage in this range.

332 Discussion

333 With this work, we introduce a systematic approach to ancestry inference from cancer-derived
334 molecular data. The approach is rooted in a combination of an established, extensively used PCA-
335 based technique of ancestry inference with a central idea of inference parameter optimization using
336 data synthesized *in silico*. Crucially, this combination permits a statistically rigorous assessment of
337 inference accuracy for an individual cancer-derived molecular profile, with its unique biological (e.g.
338 cancer type) and technical (e.g., sequencing depth and quality) properties. Synthetic data here
339 are used as a substitute for a real-world set of molecular profiles sharing these properties and
340 with known ground-truth genetic ancestry. It is unrealistic to expect such a real-world set to be
341 available in all cases. Our tests of the resulting computational methodology on a representative
342 subset of cancer-derived data demonstrate its accurate and robust performance. As we describe in
343 detail in the Methods section, our data synthesis method relies on heuristic components for an
344 estimate of the allele fractions throughout the cancer-derived profile. This estimate can be made
345 more rigorous by using haplotypes in future implementations of the method, but the present version
346 produces allele fractions in good agreement with published allele fractions (ASCAT2 results in
347 (42,43)).

348 A line of research and development initiated with this work must be extended in several
349 directions. First, the performance of the methods presented must be examined more
350 comprehensively across cancer types, and sequence properties, such as quality and depth. This
351 task is computing-intensive but feasible given extensive, well annotated repositories of cancer-
352 derived data, such as those resulting from TCGA Research Network (44) and International Cancer
353 Genome Consortium (ICGC) (45) projects. For these, the genetic ancestry of the patients either is
354 known or can be readily established using matching cancer-free molecular data. Second, an
355 extension of our approach to additional profiling modalities should be examined. Chief among these
356 are low-coverage whole-genome sequences commonly used for copy-number analysis, single-
357 molecule, long-read sequences, chromatin-accessibility profiles (ATAC-seq) and cytosine-converted
358 sequences used for methylation profiling. Each of these presents unique challenges and
359 opportunities for the ancestry inference. For example, in the low-coverage whole-genome profiles
360 the sparsity of coverage is compensated by its whole-genome breadth, whereas in the long-read
361 sequences the trade-off is between the high sequence error rate and the long-distance phasing
362 afforded by the read length. Third, while the present work relied on PCA followed by nearest-
363 neighbor classification for ancestry assessment, alternatives including UMAP for the former and
364 Random Forest or Support Vector Machine for the latter exist and should be evaluated. Third,
365 future method development should be extended beyond inference of global ancestry to that of local
366 ancestry and ancestral admixture. Such an extension is particularly important in the study of cancer
367 in strongly admixed super-populations, such as AFR and AMR, and may require more extensive
368 reference data, in addition to the 1KG reference used here. Finally, beyond cancer, our
369 methodology can be applied to any molecular data from which ancestry inference is challenging.
370 Examples include RNA-seq of non-cancer origin and sequences originating in any kind of
371 fragmentary or damaged nucleic-acid specimens, such as those encountered in forensic,
372 archaeological or paleontological contexts.

373 We anticipate the computational approach described here to have a major, two-fold, impact on
374 investigation of links between ancestry and cancer. First, it will become possible to massively boost
375 the statistical power of such studies by leveraging existing tumor-derived molecular data sets without

376 matching germline sequences or ancestry annotation. Our search of the Gene Expression Omnibus
377 (GEO) database alone has identified over 1,250 such data sets, containing RNA expression data for
378 nearly 48,000 cancer tissue specimens. Such resources dwarf those of fully annotated repositories,
379 such as TCGA (44) and ICGC (45). Other molecular data repositories are likely to contain resources
380 of this category on a similar order of magnitude. Second, hundreds of thousands of tumor tissue
381 specimens stored at multiple clinical centers constitute another major resource for ancestry-aware
382 molecular studies of cancer. Here again, matching normal tissue specimens are often absent, and
383 so is ethnic or racial annotation for the patients. According to a recent estimate (46), such annotation
384 is missing in electronic health records (EHR) of over 50% of patients. Where the donor SIRE is
385 provided by the EHR, it can be used to guide the initial specimen collection for a study of ancestral
386 effects in cancer, with a subsequent genetic ancestry validation using methods developed in this
387 work. In summary, inferential tools presented here will make massive resources of archival tissues
388 available for ancestry-oriented cancer research.

389 Multiple directions of exploratory and correlative analysis are open to pursuit with the accurate
390 ancestry annotation made possible by the methods described here, even in the absence of matching
391 cancer-free molecular data. Single-nucleotide and other small-scale somatic alterations may be
392 identified in cancer-only exomes, both whole and restricted to specialized gene panels, using
393 methods developed for this purpose (47) alongside databases of frequent somatic variants in cancer
394 (18) and of frequent germline variants like gnomAD (48) and 1KG (31). Copy number variants and
395 losses of heterozygosity in cancer exomes are overwhelmingly somatic and may be determined
396 computationally (49,50). Cancer RNA expression quantification is feasible in the absence of the
397 germline genotype of the patient, including allele- and isoform-specific analysis. These and similar
398 genomic and transcriptional properties may be explored for associations with ancestral background
399 of the patients.

400

401 **Acknowledgments**

402 DAT is a distinguished scholar of the Lustgarten Foundation and Director of the Lustgarten
403 Foundation-designated Laboratory of Pancreatic Cancer Research. DAT is also supported by the
404 Cold Spring Harbor Laboratory Association, the New York Genome Center Polyethnic 1000 Project, the
405 Simons Foundation (552716), and the NIH (P30CA45508, P20CA192996, U01CA224013,
406 U01CA210240, R01CA188134, R01CA249002, and R01CA229699). DAT also acknowledges support
407 from The Pershing Square Foundation, William Ackman, and Neri Oxman. AK's work is supported by
408 the New York Genome Center Polyethnic-1000 Project, Simons Foundation award # 519054 and by
409 the Simons Center for Quantitative Biology at Cold Spring Harbor Laboratory and by the Lustgarten
410 Foundation. The results published here are in part based upon data generated by TCGA Research
411 Network: <https://www.cancer.gov/tcga>. We thank Adam Siepel, Lloyd Trotman, Jeffrey Boyd, W.
412 Richard McCombie, Thomas Gingeras, Justin Kinney, Camila dos Santos, Michael Schatz, Louis Staudt,
413 Michael Berger, David Solit and Samuel Aparicio for illuminating discussions.

Tables

Study	D	K	Accuracy	95% CI	AUROC	95% CI
TCGA-OV WES	5	13	0.998	0.994-1	0.993	0.992-0.994
TCGA-OV Panel	4	12	0.984	0.972-0.996	0.966	0.965-0.967
TCGA-OV RNA-seq	7	12	0.993	0.983-1	0.977	0.975-0.979
BeatAML WES	5	13	0.989	0.978-1	0.978	0.976-0.980
BeatAML Panel	4	13	0.991	0.981-1	0.999	0.999-0.999
BeatAML RNA-seq	4	13	0.992	0.981-1	0.999	0.999-0.999
PDAC WES	8	13	1	NA	NA	NA
PDAC Panel	6	5	0.952	0.861-1	0.958	NA
PDAC RNA-seq	4	13	1	NA	NA	NA
TCGA-BRCA WES	4	9	1	NA	NA	NA
TCGA-BRCA Panel	4	9	0.995	0.984-1	0.995	0.994-0.996
TCGA-BRCA RNA-seq	4	9	0.995	0.984-1	0.995	0.994-0.996
Aggregate WES	-	-	0.993	0.981-1	0.997	0.997-0.998
Aggregate Panel	-	-	0.988	0.972-1	0.987	0.986-0.988
Aggregate RNA-seq	-	-	0.993	0.981-1	0.993	0.993-0.994

Table 1. Overall cohort-wide performance measures for super-population calls from cancer-derived molecular data, as compared to the matching cancer-free WES or (in the case of PDAC) WGS. A reliable estimate of the confidence intervals (CI) was not possible in the case of PDAC, due to the small number of cases with matching cancer-free genotypes. The *D* and *K* values shown provide consistently high performance in each respective data set.

(a) TCGA-BRCA WES						(b) Aggregate WES							
		Inferred							Inferred				
pop		EAS	EUR	AFR	AMR	SAS	pop		EAS	EUR	AFR	AMR	SAS
Cancer-free WES	EAS	47	0	0	0	0	EAS	69	0	0	0	0	
	EUR	0	56	0	0	0	EUR	0	732	0	6	0	
	AFR	0	0	51	0	0	AFR	0	0	96	0	0	
	AMR	0	0	0	25	0	AMR	0	1	0	70	0	
	SAS	0	0	0	0	4	SAS	0	0	0	0	14	

(c) TCGA-BRCA Panel						(b) Aggregate Panel							
		Inferred							Inferred				
pop		EAS	EUR	AFR	AMR	SAS	pop		EAS	EUR	AFR	AMR	SAS
Cancer-free WES	EAS	47	0	0	0	0	EAS	69	0	0	0	0	
	EUR	0	56	0	0	0	EUR	0	733	0	5	0	
	AFR	0	0	51	0	0	AFR	0	0	95	1	0	
	AMR	0	0	0	24	1	AMR	0	5	0	65	1	
	SAS	0	0	0	0	4	SAS	0	0	0	0	14	

(a) TCGA-BRCA RNA						(b) Aggregate RNA							
		Inferred							Inferred				
pop		EAS	EUR	AFR	AMR	SAS	pop		EAS	EUR	AFR	AMR	SAS
Cancer-free WES	EAS	47	0	0	0	0	EAS	62	0	0	0	0	
	EUR	0	56	0	0	0	EUR	0	521	0	2	0	
	AFR	0	0	51	0	0	AFR	0	0	83	0	0	
	AMR	0	0	0	24	1	AMR	1	1	0	59	1	
	SAS	0	0	0	0	4	SAS	0	0	0	0	10	

Table 2. Confusion matrices comparing TCGA-BRCA or aggregate of all patients' super-population calls from the cancer-derived molecular profiles for the three profiling modalities (rows) to those from the matching cancer-free WES.

References

1. Ashktorab H, Kupfer SS, Brim H, Carethers JM. Racial Disparity in Gastrointestinal Cancer Risk. *Gastroenterology* **2017**;153:910-23
2. Cronin KA, Lake AJ, Scott S, Sherman RL, Noone AM, Howlader N, *et al.* Annual Report to the Nation on the Status of Cancer, part I: National cancer statistics. *Cancer* **2018**;124:2785-800
3. Siegel RL, Miller KD, Jemal A. Cancer statistics, 2020. *CA Cancer J Clin* **2020**;70:7-30
4. Tan DS, Mok TS, Rebbeck TR. Cancer Genomics: Diversity and Disparity Across Ethnicity and Geography. *J Clin Oncol* **2016**;34:91-101
5. Huang BZ, Stram DO, Le Marchand L, Haiman CA, Wilkens LR, Pandol SJ, *et al.* Interethnic differences in pancreatic cancer incidence and risk factors: The Multiethnic Cohort. *Cancer Med* **2019**;8:3592-603
6. Bhatnagar B, Kohlschmidt J, Mrozek K, Zhao Q, Fisher JL, Nicolet D, *et al.* Poor Survival and Differential Impact of Genetic Features of Black Patients with Acute Myeloid Leukemia. *Cancer Discov* **2021**;11:626-37
7. Carrot-Zhang J, Chambwe N, Damrauer JS, Knijnenburg TA, Robertson AG, Yau C, *et al.* Comprehensive Analysis of Genetic Ancestry and Its Molecular Correlates in Cancer. *Cancer Cell* **2020**;37:639-54 e6
8. Carrot-Zhang J, Soca-Chafre G, Patterson N, Thorner AR, Nag A, Watson J, *et al.* Genetic Ancestry Contributes to Somatic Mutations in Lung Cancers from Admixed Latin American Populations. *Cancer Discov* **2021**;11:591-8
9. Mahal BA, Alshalalfa M, Kensler KH, Chowdhury-Paulino I, Kantoff P, Mucci LA, *et al.* Racial Differences in Genomic Profiling of Prostate Cancer. *N Engl J Med* **2020**;383:1083-5
10. Sinha S, Mitchell KA, Zingone A, Bowman E, Sinha N, Schäffer AA, *et al.* Higher prevalence of homologous recombination deficiency in tumors from African Americans versus European Americans. *Nature Cancer* **2020**;1:112-21
11. Yuan J, Hu Z, Mahal BA, Zhao SD, Kensler KH, Pi J, *et al.* Integrated Analysis of Genetic Ancestry and Genomic Alterations across Cancers. *Cancer Cell* **2018**;34:549-60 e9
12. Nugent A, Conatser KR, Turner LL, Nugent JT, Sarino EMB, Ricks-Santi LJ. Reporting of race in genome and exome sequencing studies of cancer: a scoping review of the literature. *Genet Med* **2019**;21:2676-80
13. Mersha TB, Abebe T. Self-reported race/ethnicity in the age of genomic research: its potential impact on understanding health disparities. *Hum Genomics* **2015**;9:1
14. Alexander DH, Novembre J, Lange K. Fast model-based estimation of ancestry in unrelated individuals. *Genome Res* **2009**;19:1655-64
15. Diaz-Papkovich A, Anderson-Trocmé L, Ben-Eghan C, Gravel S. UMAP reveals cryptic population structure and phenotype heterogeneity in large genomic cohorts. *PLOS Genetics* **2019**;15:e1008432
16. Price AL, Patterson NJ, Plenge RM, Weinblatt ME, Shadick NA, Reich D. Principal components analysis corrects for stratification in genome-wide association studies. *Nat Genet* **2006**;38:904-9
17. Pritchard JK, Stephens M, Donnelly P. Inference of population structure using multilocus genotype data. *Genetics* **2000**;155:945-59
18. Tate JG, Bamford S, Jubb HC, Sondka Z, Beare DM, Bindal N, *et al.* COSMIC: the Catalogue Of Somatic Mutations In Cancer. *Nucleic Acids Res* **2019**;47:D941-d7
19. Frampton GM, Fichtenholtz A, Otto GA, Wang K, Downing SR, He J, *et al.* Development and validation of a clinical cancer genomic profiling test based on massively parallel DNA sequencing. *Nat Biotechnol* **2013**;31:1023-31
20. Dutil J, Chen Z, Monteiro AN, Teer JK, Eschrich SA. An Interactive Resource to Probe Genetic Diversity and Estimated Ancestry in Cancer Cell Lines. *Cancer Res* **2019**;79:1263-73
21. Huang Q, Baudis M. Enabling population assignment from cancer genomes with SNP2pop. *Sci Rep* **2020**;10:4846
22. Kessler MD, Bateman NW, Conrads TP, Maxwell GL, Dunning Hotopp JC, O'Connor TD. Ancestral characterization of 1018 cancer cell lines highlights disparities and reveals gene expression and mutational differences. *Cancer* **2019**;125:2076-88
23. Arora K, Tran TN, Kemel Y, Mehine M, Liu YL, Nandakumar S, *et al.* Genetic Ancestry Correlates with Somatic Differences in a Real-World Clinical Cancer Sequencing Cohort. *Cancer Discov* **2022**
24. Krasnitz A, Kendall J, Alexander J, Levy D, Wigler M. Early Detection of Cancer in Blood Using Single-Cell Analysis: A Proposal. *Trends Mol Med* **2017**;23:594-603
25. Tiriác H, Belleau P, Engle DD, Plenker D, Deschênes A, Somerville TDD, *et al.* Organoid Profiling Identifies Common Responders to Chemotherapy in Pancreatic Cancer. *Cancer Discov* **2018**;8:1112-29
26. Integrated genomic analyses of ovarian carcinoma. *Nature* **2011**;474:609-15
27. Koboldt DC, Fulton RS, McLellan MD, Schmidt H, Kalicki-Veizer J, McMichael JF, *et al.* Comprehensive molecular portraits of human breast tumours. *Nature* **2012**;490:61-70
28. Tyner JW, Tognon CE, Bottomly D, Wilmot B, Kurtz SE, Savage SL, *et al.* Functional genomic landscape of acute myeloid leukaemia. *Nature* **2018**;562:526-31
29. Abecasis GR, Altshuler D, Auton A, Brooks LD, Durbin RM, Gibbs RA, *et al.* A map of human genome variation from population-scale sequencing. *Nature* **2010**;467:1061-73

30. Byrska-Bishop M, Evani US, Zhao X, Basile AO, Abel HJ, Regier AA, *et al.* High coverage whole genome sequencing of the expanded 1000 Genomes Project cohort including 602 trios. *bioRxiv* **2021**:2021.02.06.430068
31. Fairley S, Lowy-Gallego E, Perry E, Flicek P. The International Genome Sample Resource (IGSR) collection of open human genomic variation resources. *Nucleic acids research* **2020**;48:D941-D7
32. Shen R, Seshan VE. FACETS: allele-specific copy number and clonal heterogeneity analysis tool for high-throughput DNA sequencing. *Nucleic Acids Res* **2016**;44:e131
33. Koboldt DC, Chen K, Wylie T, Larson DE, McLellan MD, Mardis ER, *et al.* VarScan: variant detection in massively parallel sequencing of individual and pooled samples. *Bioinformatics* **2009**;25:2283-5
34. Lowy-Gallego E, Fairley S, Zheng-Bradley X, Ruffier M, Clarke L, Flicek P. Variant calling on the GRCh38 assembly with the data from phase three of the 1000 Genomes Project. *Wellcome Open Res* **2019**;4:50
35. Zheng X, Levine D, Shen J, Gogarten SM, Laurie C, Weir BS. A high-performance computing toolset for relatedness and principal component analysis of SNP data. *Bioinformatics* **2012**;28:3326-8
36. Hand DJ, Till RJ. A simple generalisation of the area under the ROC curve for multiple class classification problems. *Mach Learn* **2001**;45:171-86
37. Robin X, Turck N, Hainard A, Tiberti N, Lisacek F, Sanchez JC, *et al.* pROC: an open-source package for R and S+ to analyze and compare ROC curves. *BMC Bioinformatics* **2011**;12:77
38. DeLong ER, DeLong DM, Clarkepearson DI. Comparing the Areas under 2 or More Correlated Receiver Operating Characteristic Curves - a Nonparametric Approach. *Biometrics* **1988**;44:837-45
39. Sun X, Xu WC. Fast Implementation of DeLong's Algorithm for Comparing the Areas Under Correlated Receiver Operating Characteristic Curves. *Ieee Signal Proc Let* **2014**;21:1389-93
40. Chen H, Boutros PC. VennDiagram: a package for the generation of highly-customizable Venn and Euler diagrams in R. *Bmc Bioinformatics* **2011**;12
41. Graumann J, Cotton R. multipanelfigure: Simple Assembly of Multiple Plots and Images into a Compound Figure. *J Stat Softw* **2018**;84:1-10
42. Grossman RL, Heath AP, Ferretti V, Varmus HE, Lowy DR, Kibbe WA, *et al.* Toward a Shared Vision for Cancer Genomic Data. *N Engl J Med* **2016**;375:1109-12
43. Heath AP, Ferretti V, Agrawal S, An M, Angelakos JC, Arya R, *et al.* The NCI Genomic Data Commons. *Nat Genet* **2021**;53:257-62
44. Gao GF, Parker JS, Reynolds SM, Silva TC, Wang LB, Zhou W, *et al.* Before and After: Comparison of Legacy and Harmonized TCGA Genomic Data Commons' Data. *Cell Syst* **2019**;9:24-34.e10
45. Zhang J, Bajari R, Andric D, Gerthoffert F, Lepsa A, Nahal-Bose H, *et al.* The International Cancer Genome Consortium Data Portal. *Nat Biotechnol* **2019**;37:367-9
46. Polubriaginof FCG, Ryan P, Salmasian H, Shapiro AW, Perotte A, Safford MM, *et al.* Challenges with quality of race and ethnicity data in observational databases. *J Am Med Inform Assoc* **2019**;26:730-6
47. Sun JX, He Y, Sanford E, Montesion M, Frampton GM, Vignot S, *et al.* A computational approach to distinguish somatic vs. germline origin of genomic alterations from deep sequencing of cancer specimens without a matched normal. *PLoS Comput Biol* **2018**;14:e1005965
48. Karczewski KJ, Francioli LC, Tiao G, Cummings BB, Alfoldi J, Wang Q, *et al.* The mutational constraint spectrum quantified from variation in 141,456 humans. *Nature* **2020**;581:434-43
49. Oh S, Geistlinger L, Ramos M, Morgan M, Waldron L, Riester M. Reliable Analysis of Clinical Tumor-Only Whole-Exome Sequencing Data. *JCO Clin Cancer Inform* **2020**;4:321-35
50. Riester M, Singh AP, Brannon AR, Yu K, Campbell CD, Chiang DY, *et al.* PureCN: copy number calling and SNV classification using targeted short read sequencing. *Source Code Biol Med* **2016**;11:13

Figure Legends

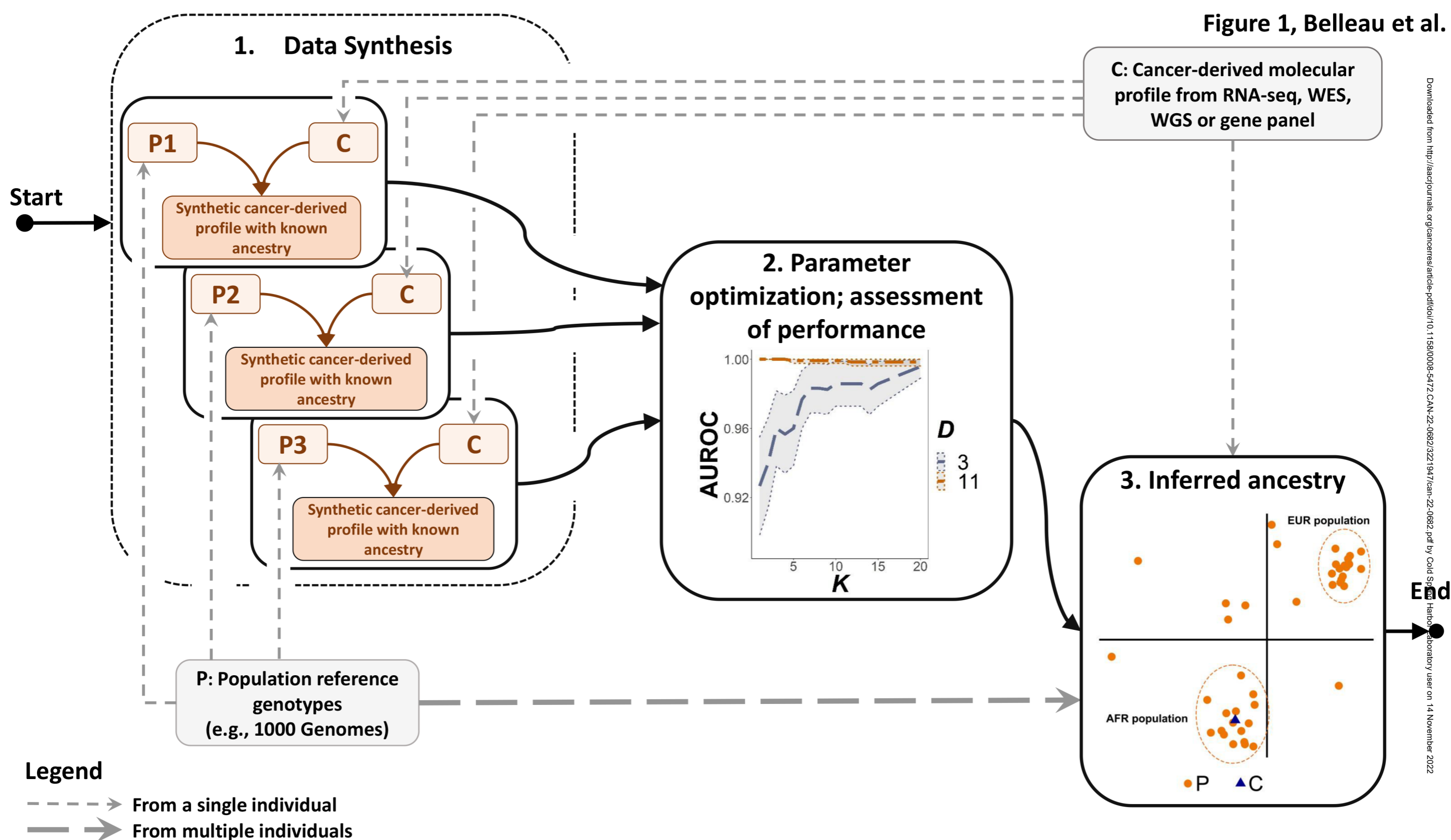
Figure 1. An overview of genetic ancestry inference from cancer-derived molecular data using data synthesis.

Figure 2. Summary of the molecular data used in this study. These originate from four patient cohorts: **A)** donors to TCGA ovarian cancer collection **B)** Beat AML clinical trial **C)** pancreatic ductal adenocarcinoma patients donating to CSHL patient-derived organoid collection **D)** a subset of donors to TCGA breast cancer collection. **E)** SIRE composition for the TCGA-OV, Beat AML, PDAC and TCGA-BRCA cohorts and in aggregate over all four cohorts. UNK means not reported or unknown.

Figure 3. A flowchart of the inference of genetic ancestry.

Figure 4. An overview of the data synthesis.

Figure 5. Dependence of AMR-specific AUROC on the inference parameters D and K , computed using data synthesis for 10 PDAC patients and the three profiling modalities: WES, RNA-seq and FoundationOne[®] CDx panels. The central AUROC values are shown in solid, and the 95% CI in dashed, lines.



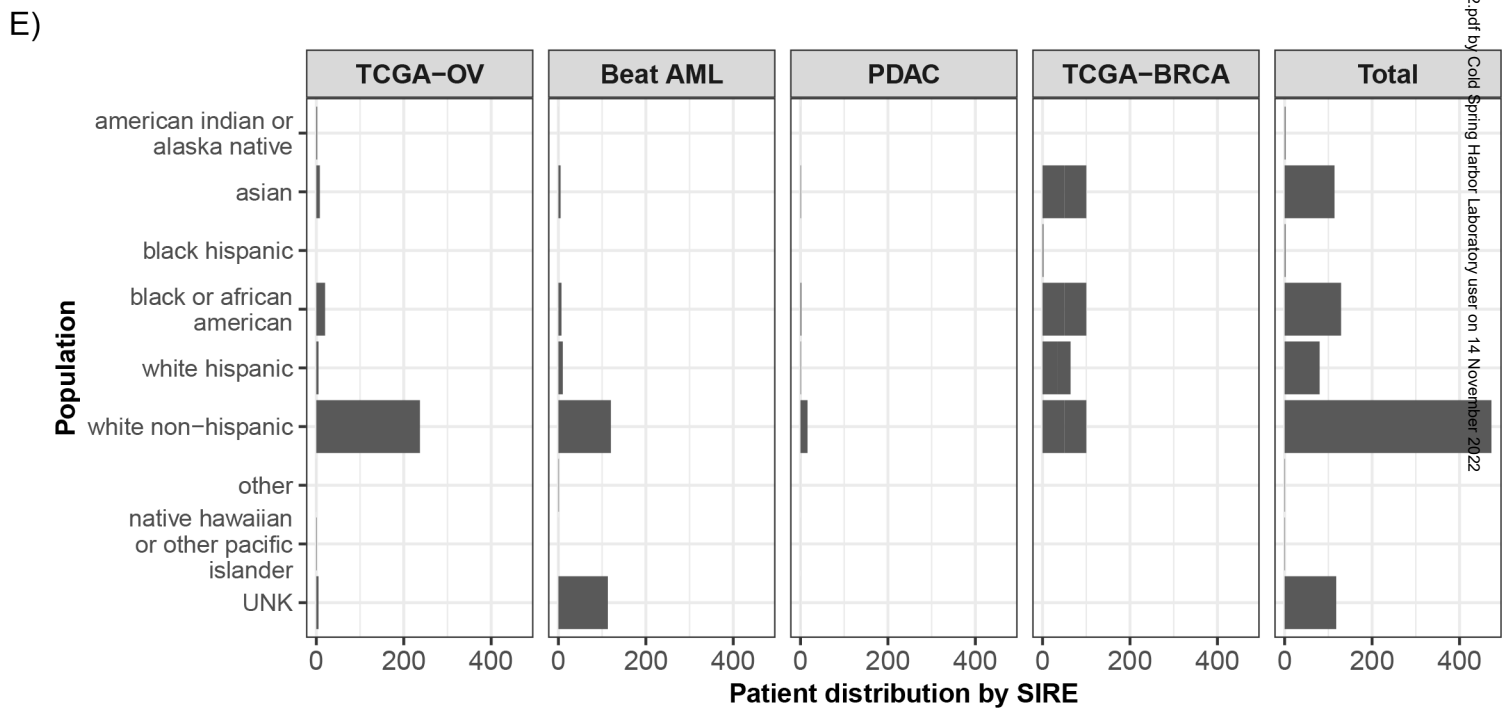
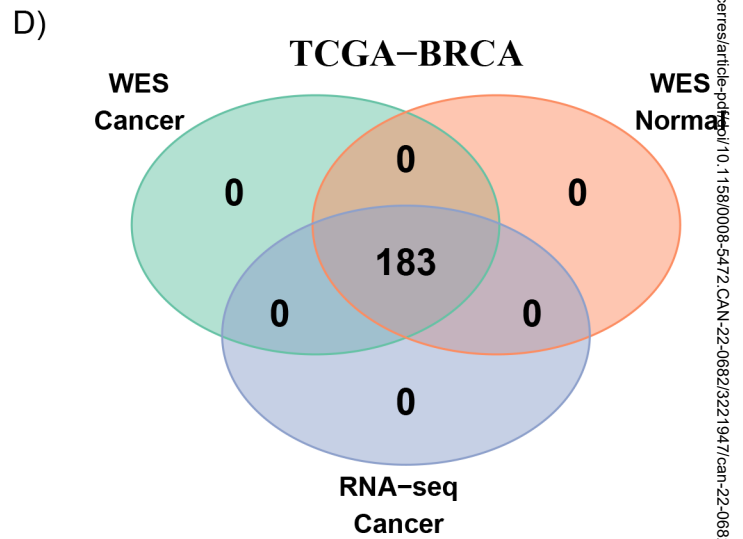
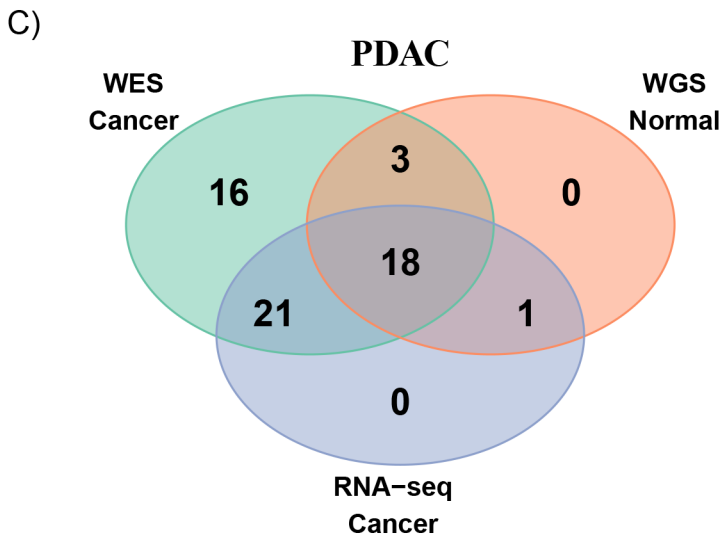
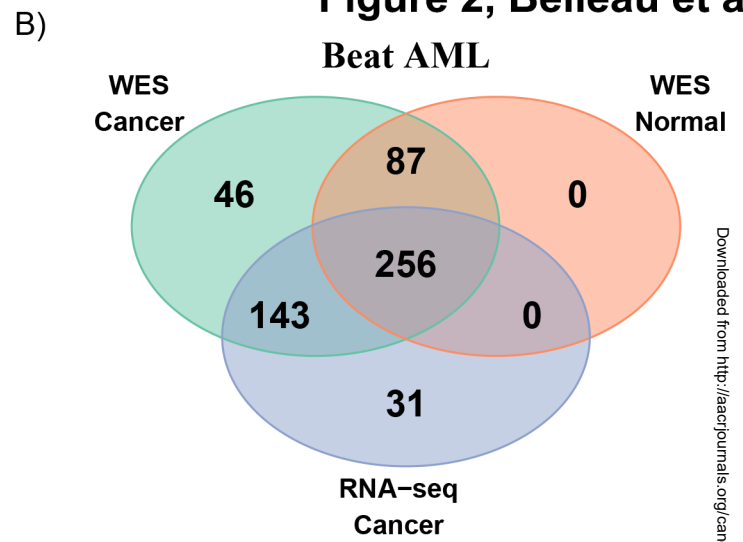
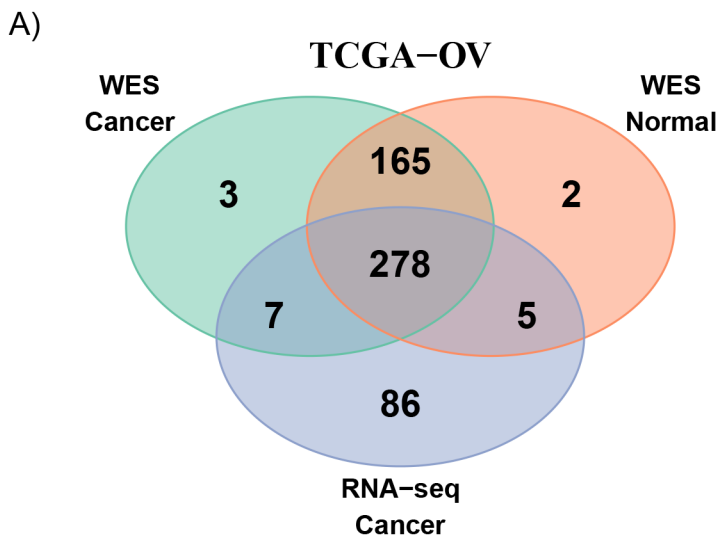


Figure 3, Belleau et al.

Legend:

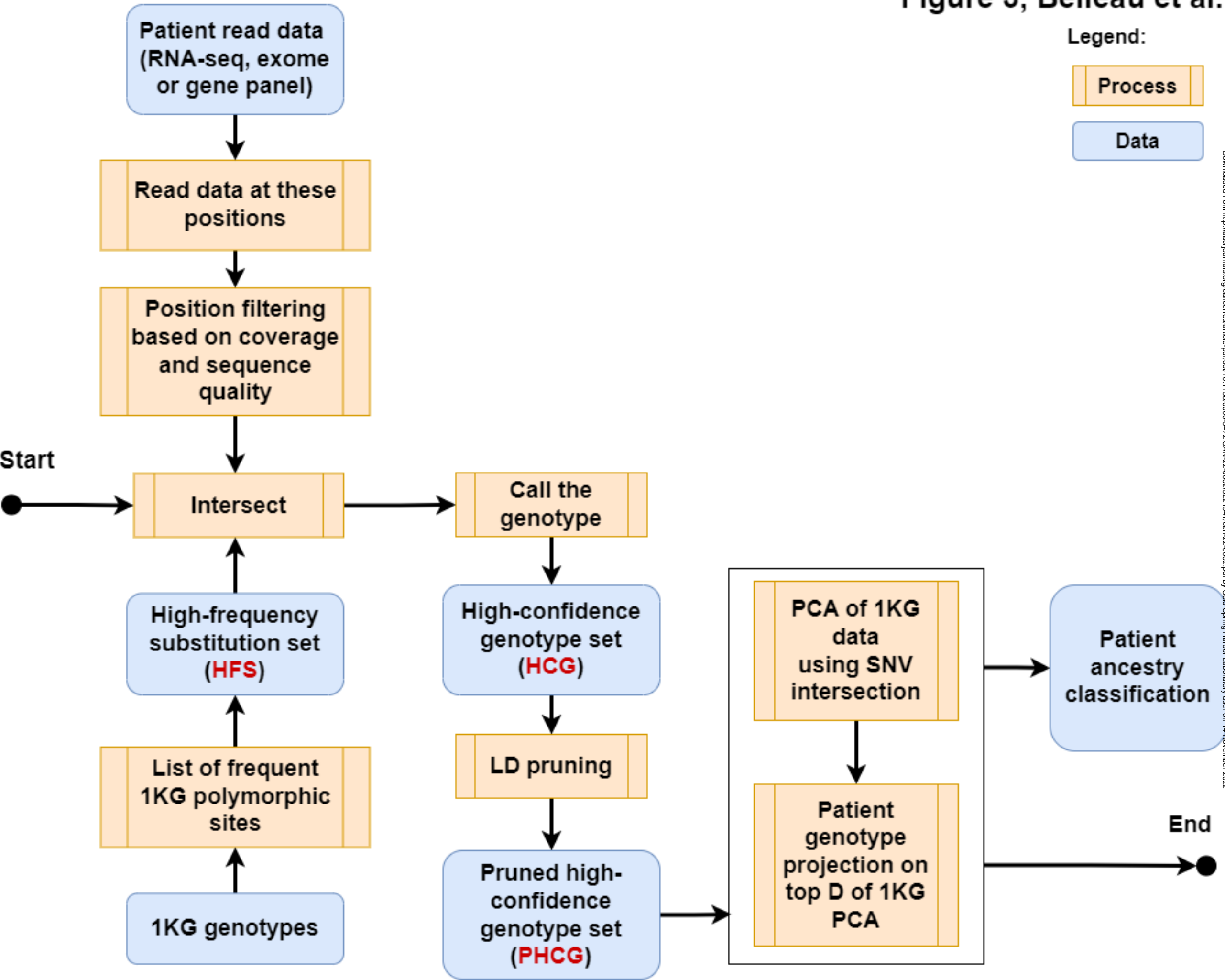
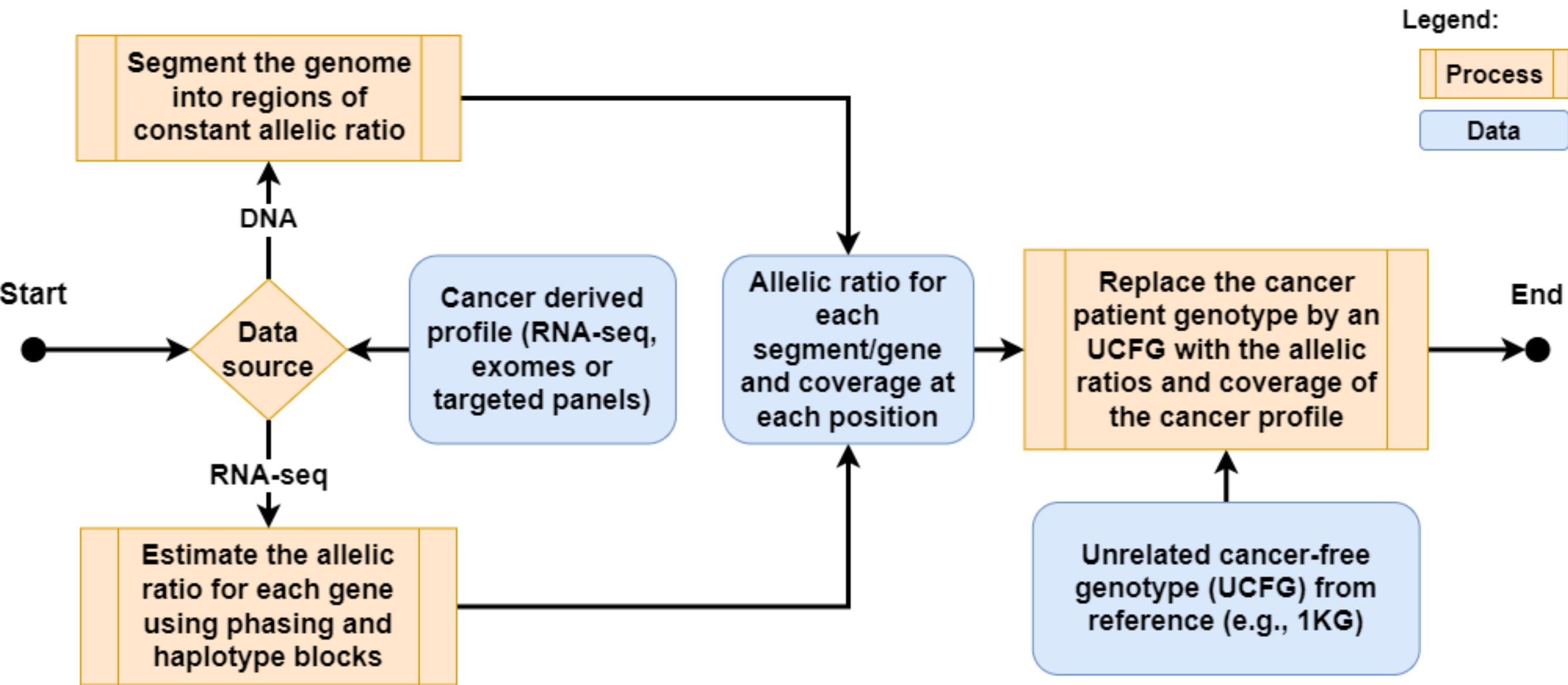
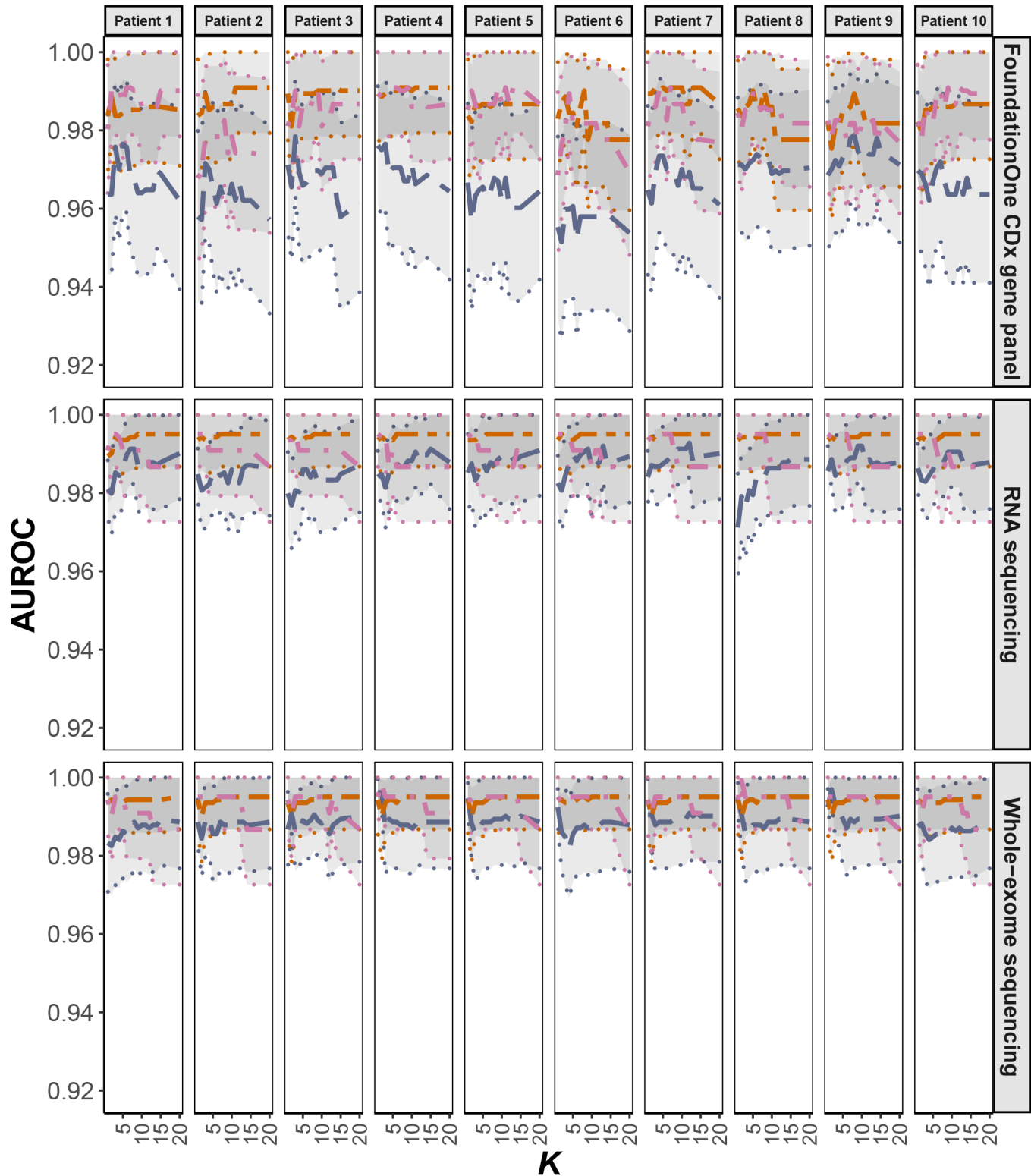


Figure 4, Belleau et al.





D