

1 **Title: Pan-cancer detection and typing by mining patterns in large genome-wide cell-free**

2 **DNA sequencing datasets**

3

4 **Running head:** cfDNA data mining for cancer detection

5

6 **Authors:** Huiwen Che¹, Tatjana Jatsenko¹, Liesbeth Lenaerts², Luc Dehaspe³, Leen
7 Vancoillie³, Nathalie Brison³, Ilse Parijs³, Kris Van Den Bogaert³, Daniela Fischerova⁴, Ruben
8 Heremans⁵, Chiara Landolfo⁶, Antonia Carla Testa⁷, Adriaan Vanderstichele⁶, Lore Liekens⁸,
9 Valentina Pomella⁸, Agnieszka Wozniak⁹, Christophe Doods^{10,11}, Els Wauters^{10,11}, Sigrid
10 Hatse^{9,12}, Kevin Punie^{12,13}, Patrick Neven^{6,12}, Hans Wildiers^{12,13}, Sabine Tejpar⁸, Diether
11 Lambrechts¹⁴, An Coosemans¹⁵, Dirk Timmerman^{5,6}, Peter Vandenberghe^{16,17}, Frédéric
12 Amant^{2,6,18}, Joris Robert Vermeesch^{1,3*}

13 **Affiliations:**

14 ¹Department of Human Genetics, Laboratory for Cytogenetics and Genome Research, KU
15 Leuven, Leuven, Belgium.

16 ²Department of Oncology, Laboratory of Gynecological Oncology, KU Leuven, Leuven,
17 Belgium.

18 ³Centre for Human Genetics, University Hospitals Leuven, Leuven, Belgium.

19 ⁴Department of Obstetrics and Gynaecology, First Faculty of Medicine, Charles University and
20 General University Hospital in Prague, Prague, Czech Republic.

21 ⁵Department of Development and Regeneration, Woman and Child, KU Leuven, Leuven,
22 Belgium.

23 ⁶Department of Gynecology and Obstetrics, University Hospitals Leuven, Leuven, Belgium.

24 ⁷Department of Woman and Child Health, Fondazione Policlinico Universitario A. Gemelli,
25 IRCCS, Università Cattolica del Sacro Cuore Roma, Rome, Italy.

NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

26 ⁸Department of Oncology, Molecular Digestive Oncology, KU Leuven, Leuven, Belgium.

27 ⁹Department of Oncology, Laboratory of Experimental Oncology, KU Leuven, Leuven,
28 Belgium.

29 ¹⁰Department of Chronic Diseases and Metabolism, Laboratory of Respiratory Diseases and
30 Thoracic Surgery (BREATHE), KU Leuven, Leuven, Belgium.

31 ¹¹Department of Pneumology, University Hospitals Leuven, Leuven, Belgium.

32 ¹²Multidisciplinary Breast Centre, Leuven Cancer Institute, University Hospitals Leuven,
33 Leuven, Belgium.

34 ¹³Department of General Medical Oncology, Leuven Cancer Institute, University Hospitals
35 Leuven, Leuven, Belgium.

36 ¹⁴Department of Human Genetics, Laboratory of Translational Genetics, VIB-KU Leuven,
37 Leuven, Belgium.

38 ¹⁵Department of Oncology, Laboratory of Tumor Immunology and Immunotherapy, Leuven
39 Cancer Institute, KU Leuven, Leuven, Belgium.

40 ¹⁶Department of Human Genetics, Laboratory of Genetics of Malignant Diseases, KU Leuven,
41 Leuven, Belgium.

42 ¹⁷Department of Hematology, University Hospitals Leuven, Leuven, Belgium.

43 ¹⁸Center for Gynecological Oncology Amsterdam: Academic Medical Centre Amsterdam-
44 University of Amsterdam and The Netherlands Cancer Institute-Antoni van Leeuwenhoek
45 Hospital, Amsterdam, the Netherlands.

46

47 *Corresponding author: Joris Robert Vermeesch, Katholieke Universiteit Leuven (KU Leuven)
48 Universitaire Ziekenhuizen Leuven (University Hospitals Leuven), Human Genetics,
49 Herestraat 49, box 602, 3000 Leuven, Belgium.

50 Telephone: +32 16 34 5941

51 Email: joris.vermeesch@uzleuven.be

52 **Abstract**

53 **Background**

54 Cell-free DNA (cfDNA) analysis holds great promise for non-invasive cancer screening,
55 diagnosis and monitoring. We hypothesized that mining the patterns of big datasets of shallow
56 whole genome sequencing cfDNA from cancer patients could improve cancer detection.

57 **Methods**

58 By applying unsupervised clustering and supervised machine learning on large shallow whole-
59 genome sequencing cfDNA datasets from healthy individuals (n=367), patients with different
60 hematological (n=238) and solid malignancies (n=320), we identify cfDNA signatures that
61 enable cancer detection and typing.

62 **Results**

63 Unsupervised clustering revealed cancer-type-specific sub-grouping. Classification using
64 supervised machine learning model yielded an overall accuracy of 81.62% in discriminating
65 malignant from control samples. The accuracy of disease type prediction was 85% and 70%
66 for the hematological and solid cancers, respectively. We demonstrate the clinical utility of our
67 approach by classifying benign from invasive and borderline adnexal masses with an AUC of
68 0.8656 and 0.7388, respectively.

69 **Conclusions**

70 This approach provides a generic and cost-effective strategy for non-invasive pan-cancer
71 detection.

72

73 **key words:** liquid biopsy, cfDNA, ctDNA, hematological malignancies, solid tumors, ovarian
74 tumors, machine learning

75 **Introduction**

76 Cell-free DNA (cfDNA) is a promising non-invasive biomarker in liquid biopsy for cancer
77 management. Shallow whole-genome sequencing (sWGS) of cfDNA can identify cancer-
78 specific copy number aberrations (CNAs) in cancer patients (1,2). Using genome-wide cfDNA
79 sequencing data to profile genomic imbalances, we reported that CNAs in the asymptomatic
80 population can be indicative of incipient tumors and has potential as a cancer screening tool
81 (3).

82

83 In addition to CNAs, sequencing of cfDNA provides a unique view on the genome-wide
84 cfDNA fragmentation profile (4,5). CfDNA fragments carry tissue-associated nucleosome and
85 preferred end position information (6,7), reflecting tissue-specific degradation, chromatin
86 accessibility and nucleosome organization of its cellular origin (8,9). In healthy individuals,
87 plasma cfDNA comprises DNA fragments that are mainly resulting from apoptotic release of
88 DNA from the cells of hematopoietic origin (10). In plasma of cancer patients, circulating
89 tumor DNA (ctDNA) has decreased fragment sizes and signatures of the tissue of origin (8,11).
90 Consequently, fragmentomics is emerging as an approach to reveal cfDNA properties,
91 broadening the potential of cfDNA as a biomarker (4,12).

92

93 Increasing availability of cfDNA sWGS data from large-scale liquid biopsy projects offer
94 unique opportunities to explore the cfDNA profiles by machine learning. We hypothesized that
95 mining variation between sWGS profiles may uncover distinct patterns that can be associated
96 with different pathological or physiological states. Hence, we applied an unsupervised
97 clustering analysis and supervised machine learning workflow, which we term *GIPXplore*, on
98 a large number of genome-wide sWGS cfDNA profiles from cancer patients with different
99 hematological and solid tumors and unveiled cancer-type-specific and also shared tumor-

100 associated signatures that are absent in healthy individuals. This approach enables accurate
101 detection of different cancers and allows prediction of the cancer types.

102 **Materials and Methods**

103 **Patients and clinical data**

104 The study was approved by the ethical committee of the University Hospitals Leuven (S57999,
105 S62285, S62795, S50623, S56534, S63240, S51375, S59207, S64205 and S64035). Samples
106 and consents were obtained from healthy controls and cancer patients. Blood was collected
107 either into Streck Cell-Free DNA BCT or Roche Cell-Free DNA Collection Tubes. Plasma was
108 isolated via a standard, two-step centrifugation procedure and stored at -80°C. Previously
109 published sequencing data from 260 healthy subjects (3) and 177 patients with Hodgkin's
110 lymphoma (13) were included in the study.

111

112 **sWGS Analysis**

113 cfDNA was extracted from plasma using standard processing procedures and sWGS
114 sequencing (14) (details described in the Supplemental Materials). Each sample ended with
115 57509 autosome bin features from standard processing. Principal component analysis (PCA)
116 was used for dimension reduction to transform data from high dimension to low dimension.
117 We performed the supervised learning on both the original data space and PCA transformed
118 space and found marginal gains of performance in the majority of analyses with the original
119 data space. As the computational time was much higher using the original data space, we used
120 PCA features in the main analyses such that features being used in both unsupervised and
121 supervised learning were consistent.

122

123 **GIPXplore**

124 As illustrated in **Fig. 1**, we developed GIPXplore to mine sWGS cfDNA data for identification
125 of signatures. We utilized unsupervised clustering and supervised machine learning. For
126 unsupervised clustering, we evaluated the variance being explained from principal components

127 (PCs) in the tumor data. Overall, the top 30, 50 and 100 PCs explained above 80%, 85% and
128 90% of the variance in the data, respectively. While there is no absolute optimal number of
129 PCs to be used for further analysis, non-trivial components - 50 PCs (**Supplemental**
130 **Materials**) were determined as a default number for downstream analyses in the results. The
131 Euclidean metric was used to measure dissimilarity among samples for clustering analysis.
132 Proximity matrix based on dissimilarity of samples was generated. The t-distributed stochastic
133 neighbor embedding (tSNE) (15) was used to map high-dimensional data to two (or three)
134 dimensions and to visualize the clusters. Due to the random process of tSNE, we applied
135 Walktrap community (16) detection on the original proximity matrix for cluster assignments
136 regardless of the presentation of tSNE visualization. In running tSNE, we set parameters
137 perplexity of 15/30 and iteration of 10000 with exact tSNE for accuracy, and the process was
138 repeated for 10 times with different seeds. For Walktrap, we used the parameters of 8 initial
139 numbers of neighbors search and a walk step of 2. Clusters defined from the community
140 detection were used for annotation. In supervised learning, PCA transformed genome-wide
141 features were used in the machine learning model for training. PCA was performed on training
142 data, and test data was projected on PCA space of training data for classification tasks. We
143 measured performance by repeating the tenfold cross-validation 10 times and leave-one-out
144 (LOO) procedures. For cross fold (CV) validation, the ROC curve and performance was
145 calculated by averaging over 10 repeats. For classifiers, we used a support vector machine
146 (SVM) and hyperparameters were chosen based on the grid search with a subset of the data. A
147 separate model was trained to localize tissue of origin and LOO was used to evaluate
148 performance characteristics. Weighted sample size was accounted for in the model for
149 imbalanced classes.

150 **Results**

151 **GIPXplore detects and classifies hematological malignancies with high accuracy**

152 To assess the potential, we applied our method on a set of cfDNA samples from healthy
153 controls (n=260) and patients with hematological malignancies that included Hodgkin's
154 lymphoma (HL; n=179), diffuse large B-cell lymphoma (DLBCL; n=37), multiple myeloma
155 (MM; n=22) (**Table 1**). Walktrap community detection was performed on the dataset, and 15
156 clusters were defined. Visualization with the tSNE yielded separations between malignant and
157 healthy control profiles, and the tSNE representation was largely in agreement with the clusters
158 found by Walktrap (**Fig. 2, A**). Moreover, we observed cancer type-specific clusters. Cluster
159 1, 3 and 4 was exclusively composed of HL samples. Cluster 9 was enriched for DLBCL
160 samples, and cluster 13 was specific to MM samples (**Fig. 2, B** and **Supplemental Fig. 1**).

161
162 In parallel, we benchmarked our method against the ichorCNA (17) algorithm for copy number
163 profiling and tumor fraction (TF) estimation from sWGS data. IchorCNA utilizes the depth of
164 coverage to evaluate the presence of large-scale copy number aberrations and the probabilistic
165 model is used to infer copy number states and estimate fraction of tumor. Overall, only 52.95%
166 of hematological cancer samples had detectable tumor-derived cfDNA levels (**Fig. 2, C**,
167 **Supplemental Fig. 2** and **Supplemental Table 1**), using the 3% detection limit suggested in
168 the ichorCNA for detecting the presence of tumor. The above-mentioned clusters 1, 3, and 4
169 consisted of profiles characterized by large chromosomal aberrations and high tumor load.
170 Clusters 2 and 8 consisted of profiles from HL patients with both high and low tumor fractions,
171 implying that the clustering was not completely CNA-driven. In particular, ten out of 65
172 (15.38%) lymphoma samples in cluster 2 with normal-like profiles (without detectable CNAs)
173 grouped together with samples characterized by detectable CNAs. A less pronounced
174 separation could be observed between clusters containing healthy controls and cluster 8, in

175 which 76.47% (26 out of 34) malignant cases had normal-like profiles with less than 3% TFs.
176 Nine HL samples in cluster 10 showed higher bin-to-bin log₂ ratio variations and were more
177 likely to be noisy on a genome-wide scale (**Fig. 2, D** and **Supplemental Fig. 3**). The remaining
178 malignant cases without detectable CNAs co-localized with healthy controls. To further
179 explore whether clustering of malignant samples would be mainly CNA-driven, we performed
180 clustering analysis using the log₂ copy ratio values produced by ichorCNA. The analysis
181 revealed that genome-wide copy number ratios alone were less informative (**Supplemental**
182 **Fig. 4**). In addition, we tested whether our method could detect underlying genome-wide
183 changes irrespective of the presence of CNAs by restricting the clustering to the cancer samples
184 with low TF (< 3%). The separation between some malignant and healthy samples still
185 remained (**Supplemental Fig. 5**). Collectively, the clustering analysis on genome-wide
186 features showed separation between malignant and healthy profiles and grouping of similar
187 cancer type-specific profiles.

188
189 The unsupervised learning delineated cancer-associated profile changes, which suggested that
190 a more precise prediction can be made by learning representations within different tumor types
191 using supervised classification. Therefore, we evaluated the capability to detect cancer signals
192 and identify cancer types with supervised learning on the hematological cohort. Both leave-
193 one-out (LOO) and repeated 10-fold cross validation (CV) was used to assess the performance
194 of the classifier. Incorporating transformed genome-wide features, the SVM machine learning
195 model correctly classified 220 (out of 238) malignant cases in LOO analysis, at a sensitivity of
196 92.44% (95% CI: 88.31% - 95.46%) and a specificity of 98.46% (95% CI: 96.11% - 99.58%),
197 including 170 HL, 32 DLBCL and 18 MM cfDNA samples (**Supplemental Table 1**). The
198 remaining 18 misclassified malignant samples had normal-like profiles and clustered together
199 with healthy controls (**Supplemental Fig. 6**). The detection sensitivity was the highest for HL

200 **(Supplemental Table 1)**. The sensitivity did not differ substantially between early (I-II) and
201 advanced (III-IV) stages for these cancer types, though the distribution of the cases across
202 clinical stages was unequal (**Fig. 3, A**). ROC analysis had an AUC value of 0.989 (95% CI:
203 0.980 – 0.998) in distinguishing malignant from healthy samples, compared to ichorCNA TF-
204 based analysis which had an AUC of 0.929 (**Fig. 3, B**). Repeated 10-fold CV also revealed a
205 stable performance at an averaged AUC of 0.989 (**Supplemental Fig. 7**). As the clustering
206 analysis demonstrated the co-localization of samples originating from the same cancer type,
207 we then attempted to determine the accuracy of our *GIPXplore* in cancer type classification.
208 For this purpose, we trained the classification model using the 220 correctly predicted
209 malignant samples. The analysis showed an overall accuracy of 85.45% (95% CI: 80.09% -
210 89.83%), with the highest accuracy in HL prediction (**Fig. 3, C** and **Supplemental Table 2**).
211 Consistent with the exploratory clustering analysis, where some of the profiles from DLBCL
212 patients colocalized together with those from HL patients, DLBCL samples were more likely
213 to be misclassified.

214

215 ***GIPXplore* identifies and classifies different types of solid malignancies and allows disease** 216 **stratification**

217 Extending our analyses, we applied our method on a solid tumor dataset, consisting of 320
218 cfDNA profiles from cancer patients, and a set of 107 cfDNA profiles from healthy controls.
219 The malignant cohort was represented by five tumor types: breast (n=46), colorectal (n=70),
220 gastrointestinal stromal tumor (GIST; n=35), lung (n=44) and ovarian (n=125; **Table 1**). Using
221 *GIPXplore*, 19 clusters were identified in the solid tumor dataset (**Fig. 4, A** and **Supplemental**
222 **Fig. 8**). The separations between malignant and control cfDNA profiles were less distinct
223 compared to clustering results of the hematological cancer dataset. Clusters 4, 8, 10 and 12
224 were found to be cancer type-specific, in which cluster 4 was mainly enriched with ovarian

225 cancer samples, cluster 8 was primarily consisting of cfDNA profiles from lung cancer patients,
226 cluster 10 was GIST-specific and cluster 13 was mainly composed of colorectal samples (**Fig.**
227 **4, B**). Cluster 2, adjacent to clusters 4 and 8, was enriched with ovarian samples, although it
228 co-localized with other tumors. Clusters 9 (mostly ovarian cancer) and 15 (intermixed cancer
229 types) deviated from healthy and other malignant clusters. Majority of the cfDNA profiles from
230 breast cancer patients resembled profiles from healthy controls, while one advanced stage
231 breast cancer sample was found in cluster 8, and 2 samples from patients with advanced stage
232 primary metastatic disease were found in cluster 2 (**Fig. 4, A and B**).

233

234 Compared with the hematological cancer dataset, ctDNA levels estimated by ichorCNA were
235 generally lower in the solid malignant cohort (**Fig. 4, C**). Tumor fraction varied among
236 different types of cancer and increased with the stage (**Supplemental Fig. 9**). The malignant
237 cases with detectable CNAs and therefore higher TF were more likely to separate from the
238 healthy controls (**Supplemental Fig. 10**). Cluster 4 contained ovarian cancer samples with
239 detectable chromosome instability. Among lung cancer profiles in cluster 8, 64.29% (9 out of
240 14) had detectable CNAs. Clusters 16 to 19 included four ovarian samples with high
241 chromosomal instability that greatly deviated from other profiles. Overall, in clusters 9 and 15,
242 profiles tended to be noisy, without clear CNAs (**Fig. 4, D**), however they deviated from
243 healthy control and other malignant clusters (**Fig. 4, A**). When using the log₂ copy ratio profiles
244 from the CNA analysis to investigate whether the sub-grouping of cfDNA profiles was driven
245 by CNAs, cancer type-specific clustering patterns were diminished (**Supplemental Fig. 11**).
246 When restricting the clustering analysis to samples with TF lower than 3%, samples from
247 clusters 9 and 15 still showed deviations (**Fig. 4, A clusters 8 and 9**) from normal profiles
248 (**Supplemental Fig. 12**).

249

250 We next investigated whether supervised learning using genome-wide features can enhance
251 the detection of solid malignancy signals in sWGS cfDNA data. Classification of samples as
252 either healthy or malignant (107 healthy controls and 320 malignancies) was performed using
253 the SVM model, with performance estimated by LOO and repeated 10-fold CV. With an
254 overall accuracy of 65.34%, we correctly detected 177 out of 320 cancer profiles (55.31%
255 sensitivity, 95% CI: 49.68% - 60.84%), at a specificity of 95%. Performance in individual
256 tumor types ranged from 15.22% (95% CI: 6.34% - 28.87%) for classifying breast cancer to
257 80.00% (95% CI: 63.06% - 91.56%) for GIST (**Supplemental Table 3**). Stage of the disease
258 affected the detection, with a sensitivity of 26.17% (95% CI: 18.15% - 35.55%) in the early
259 stage (I-II) *versus* 69.95% (95% CI: 63.31% - 76.03%) in the advanced stage (III-IV). In
260 individual tumor types, it remained true that higher sensitivities were found for the advanced
261 stages than for the early-stage diseases (**Fig. 5, A**). Colorectal cancer was an exception as
262 sensitivities were almost the same for early and advanced cancer stages. Misclassified
263 malignant samples had low tumor fraction, which potentially restricted the detection of
264 underlying tumor-specific patterns (**Supplemental Fig. 13**). We could distinguish malignancy
265 from healthy samples with an AUC of 0.827 (95% CI: 0.787 – 0.867), which again was superior
266 to ichorCNA TF-based analysis (0.733 AUC, 95% CI: 0.687 – 0.780; **Fig. 5, B** and
267 **Supplemental Fig. 14**). Subsequently, we explored the potential of our GIPX*lore* method for
268 tumor classification. When performing tumor type-specific prediction with the 171 correctly
269 predicted primary tumor samples, the LOO validation resulted in a 69.01% (95% CI: 61.49%
270 - 75.84%) overall accuracy. Highest sensitivities (>70%) were obtained for cfDNA samples
271 from ovarian cancer and GIST patients. At the same time, ovarian and colorectal tumor cfDNA
272 profiles were more likely to be misassigned to each other (**Fig. 5, C** and **Supplemental Table**
273 **4**).

274

275 Moreover, among this solid malignant cohort, we had nine cfDNA samples from patients with
276 ovarian metastases, of which four patients had gastrointestinal primary site, one lymphoma,
277 one leiomyosarcoma, one uterine origin, and the remaining two had Krukenberg tumors.
278 Annotation of these nine cases on the tSNE plot showed that metastatic profiles could resemble
279 profiles of either the primary tumor or the distant site (**Supplemental Fig. 15, A**). Applying
280 the type-specific classifier to the six metastatic cases that were predicted as malignant cases by
281 the malignancy classifier, the case with gastrointestinal origin that was co-clustered with
282 colorectal samples was classified to be colorectal class. Two out of three metastatic cases that
283 were identified in intermixed clusters of lung and ovarian tumors were predicted to be lung
284 class and the other one was assigned to ovarian class. The additional two cases identified by
285 the classifier were classified to the ovarian and colorectal classes, respectively (**Supplemental**
286 **Fig. 15, B**).

287

288 **Accurate classification of benign from invasive and borderline adnexal masses may** 289 **improve clinical management**

290 In addition to the invasive ovarian tumor samples, our cohort contained 160 benign and 63
291 borderline ovarian samples. To assess the potential utility of the method for ovarian cancer
292 management, we analyzed the ovarian tumor cohort independently by performing clustering
293 analysis and building the ovarian-specific classifier to differentiate benign from malignant
294 adnexal masses. Benign and borderline samples were less likely to have detectable ctDNA
295 levels (**Supplemental Fig. 16**). In the clustering analysis, 35 invasive samples formed a distinct
296 group in cluster 1. In clusters 6-9, common patterns were found for invasive, benign and
297 borderline samples, although they remained distinct from controls (**Supplemental Fig. 17**).
298 The classification analysis exhibited an AUC of 0.8656 (95% CI: 0.7761 – 0.8689) in
299 discriminating benign from invasive samples, and an AUC of 0.7388 (95% CI: 0.6857 –

300 0.7920) in discriminating benign from borderline and invasive samples (**Supplemental Fig. 18**
301 and **19**).

302 **DISCUSSION**

303 We present a generic approach for cancer identification and classification by mapping genome-
304 wide cfDNA signatures, without prior knowledge of genetic alterations or predefined
305 signatures in the sequencing data. The unsupervised clustering allows the discovery of hidden
306 genome-wide patterns, and the supervised learning model can be trained to detect such
307 underlying signatures. This method can be used to classify cfDNA samples by matching to
308 existing datasets and has the potential to be used as a pan-cancer assay for detection and typing
309 of multiple cancers from one blood draw.

310

311 Current sWGS cfDNA analyses mainly focus on the detection of somatic CNAs (17–19). These
312 methods are blind to events that involve copy neutral abnormalities. Our approach also differs
313 from the previous method that classified tumor types based on selected CNAs, and in which
314 normal-like profiles were incapable of tumor classification (20). We demonstrate that even
315 profiles without detectable CNAs carry informative and discriminative patterns in sWGS data.
316 Different recent studies have utilized methylation, transcription factor binding, fragment
317 lengths, or cfChIP-seq for cancer detection (4,12,21–25). While these studies have important
318 implications and show cfDNA as a promising biomarker, they require more specific workup
319 and/or deeper sequencing. In contrast, analysis of sWGS data can be easily adapted in clinical
320 settings and complement CNA analysis. By mapping differences among the cfDNA profiles,
321 shared abnormality patterns are captured.

322

323 To date, Liu *et al.* have reported the largest population-level cfDNA methylation study for
324 multi-cancer detection, in which the targeted methylation analysis of cfDNA enabled detection
325 of more than 50 cancer types at a sensitivity of 54.9% and at a specificity of 99% (21). This
326 test was refined and validated in an independent follow-up study, with an overall sensitivity of

327 51.5% at 99.5% specificity was reported (26). In line with these findings, we estimate the
328 combined sensitivity of 68.03% at above 95% specificity for the hematological and solid cancer
329 cohorts. Performance for cancer signal detection varied among the different cancer types and
330 stages. The prediction accuracy was highest for hematological malignancies and lowest for
331 breast cancer. Shedding of the ctDNA from breast cancer is known to be low (27,28). Also,
332 our cohort had an over-representation of early-stage cancers, with 50% of the samples from
333 stage I. Apart from potential screening applications, we also demonstrated that *GIPXplore*
334 could be used for risk stratification and management of a specific cancer type. Discrimination
335 between malignant, borderline and benign masses at diagnosis is of critical importance to
336 improve patient management (29,30).

337

338 The accuracy of tumor type-specific prediction might depend on the intrinsic tumor
339 characteristics. For example, DLBCL being more heterogeneous on molecular level (31,32),
340 had lower classification accuracy than HL and MM. The subtype of colorectal and ovarian
341 tumors is of similar cellular origin, and histological subtypes can be hard to distinguish (33–
342 35), which might be a reason for misclassification amongst the two cancer types. The
343 identification of the origin of some metastases, suggests the method may allow the
344 identification of unknown primary cancers. The metastatic cases were classified into profiles
345 of its primary or distant sites, possibly reflecting changes during the metastatic progression or
346 dynamic tumor DNA shedding from tumor tissues (36–38).

347

348 Interestingly, besides tumor type- or aberration-specific subgroups, our analysis revealed the
349 presence of additional clusters that segregated from healthy controls (**Fig. 4, B** and
350 **Supplemental Fig. 17**). Though the origin of such segregations remains unknown, we
351 hypothesize the method provides a system-wide insight, potentially reflecting

352 (patho)physiological conditions of these individuals. Dynamic cellular responses and
353 malignant cell proliferation with active involvement of immune response during (early)
354 carcinogenesis might lead to the observed common changes in cfDNA composition across
355 different cancer types (39,40). Therefore, it is possible that our analysis detected tumor-driven
356 immune or other biological responses or states.

357

358 *GIPXplore* provides an unbiased genome-wide scan of cfDNA profiles. However, it also has
359 some limitations. Increasing the sequencing depth might improve detection of disease-specific
360 cfDNA patterns and improve the sensitivity of our methodology further. The data presented
361 here has a larger proportion of HL and ovarian cancer samples and is limited in the number of
362 different cancer types, which may affect the aggregated sensitivity and distort tumor typing
363 accuracies. We foresee that expanding the breadth of the evaluated cancer types may improve
364 prediction of tissue/cell origin and facilitate a deeper understanding of cfDNA in the context
365 of tumors. Increasing the range of physiological states and diseases that are relevant for these
366 tumor samples will be essential to fully interrogate the potential and limitations of our
367 approach. The approach may also be further broadened to project and embed new treatment or
368 follow-up data for cancer prognosis and monitoring.

369

370 In summary, we have extended the scope of cfDNA analysis, allowing cost-effective
371 identification of genome-wide cancer-(type-)specific signatures from shallow sequencing data,
372 allowing improved discrimination between profiles from cancer patients and healthy
373 individuals. This study lays the foundation for enhanced genomic characterization of cfDNA
374 that can be used for improved cancer management. We foresee that the method can be scaled
375 up for detection of multiple pathological conditions.

376 **Acknowledgements:** We would like to acknowledge the patients and blood donors. We would
377 like to thank Gitte Thirion and Annick Van den Broeck for the collection of samples and the
378 extraction of ctDNA, Kate Stanley for helpful suggestions for the manuscript.

379 **Funding:** This study was supported by the Research Foundation-Flanders (FWO-Vlaanderen)
380 (G080217N to FA and JRV, G0A1116N to PV), Agentschap Innoveren en Ondernemen
381 (VLAIO; Flanders Innovation & Entrepreneurship grant HBC.2018.2108 to JRV), Kom Op
382 Tegen Kanker (Stand Up to Cancer, the Flemish Cancer Society under grant 2016/10728/2603
383 to AC), Stichting tegen Kanker (FAF-C/2016/836 to PV, 2018-134 to JRV and FA) and KU
384 Leuven funding (no C1/018 to JRV and DL).

385 **Conflict of interest:** Patent application pending on ‘Method for analyzing cell-free nucleic
386 acids’ (JRV and LD).

387 **Author contributions:** HC, TJ, LL, LD, FA and JRV conceptualized and designed the study.
388 KP, AW, EW, AC, PN, ST, DT, PV, HW and FA provided clinical samples and patient data.
389 LL, LV, NB, IP, KVDB, CD, DF, RH, SH, CL, L.Liekens, VP, ACT, AV and AW carried out
390 clinical sample procurement and processing. LV, NB, IP and KVDB coordinated sequencing
391 of cell-free DNA. TJ and LL conducted project coordination and administration. HC and LD
392 performed bioinformatics analysis of sWGS data. HC, TJ, LL, LD, KP, AW, EW, AC, DL, ST,
393 DT, PV, FA and JRV contributed to the interpretation of results. HC, TJ and JRV wrote the
394 manuscript; all co-authors reviewed the manuscript.

395 **Data and materials availability:** Processed alignments of sequencing data are archived to
396 ArrayExpress (<https://www.ebi.ac.uk/arrayexpress/>) with unrestricted access under accession
397 number E-MTAB-10934. Code will be available upon request. All other materials associated
398 with this study are present in the paper or the Supplementary Materials.

399 Reference

- 400 1. Vandenberghe P, Wlodarska I, Tousseyn T, Dehaspe L, Dierickx D, Verheecke M, et al.
401 Non-invasive detection of genomic imbalances in Hodgkin/Reed-Sternberg cells in
402 early and advanced stage Hodgkin's lymphoma by sequencing of circulating cell-free
403 DNA: a technical proof-of-principle study. *Lancet Haematol.* 2015;2:e55–65.
- 404 2. Lenaerts L, Che H, Brison N, Neofytou M, Jatsenko T, Lefrère H, et al. Breast Cancer
405 Detection and Treatment Monitoring Using a Noninvasive Prenatal Testing Platform:
406 Utility in Pregnant and Nonpregnant Populations. *Clin Chem.* 2020;66:1414–23.
- 407 3. Lenaerts L, Vandenberghe P, Brison N, Che H, Neofytou M, Verheecke M, et al.
408 Genomewide copy number alteration screening of circulating plasma DNA: potential
409 for the detection of incipient tumors. *Ann Oncol.* 2019;30:85–95.
- 410 4. Cristiano S, Leal A, Phallen J, Fiksel J, Adleff V, Bruhm DC, et al. Genome-wide cell-free
411 DNA fragmentation in patients with cancer. *Nature.* Nature Publishing Group;
412 2019;570:385–9.
- 413 5. Lo YMD, Han DSC, Jiang P, Chiu RWK. Epigenetics, fragmentomics, and topology of
414 cell-free DNA in liquid biopsies. *Science.* American Association for the Advancement
415 of Science; 2021;372. Available from:
416 <https://science.sciencemag.org/content/372/6538/eaaw3616>
- 417 6. Jiang P, Sun K, Tong YK, Cheng SH, Cheng THT, Heung MMS, et al. Preferred end
418 coordinates and somatic variants as signatures of circulating tumor DNA associated
419 with hepatocellular carcinoma. *Proc Natl Acad Sci. National Academy of Sciences;*
420 2018;115:E10925–33.
- 421 7. Chan KCA, Jiang P, Sun K, Cheng YKY, Tong YK, Cheng SH, et al. Second generation
422 noninvasive fetal genome analysis reveals de novo mutations, single-base parental
423 inheritance, and preferred DNA ends. *Proc Natl Acad Sci U S A.* 2016;113:E8159–68.
- 424 8. Mouliere F, Robert B, Peyrotte EA, Rio MD, Ychou M, Molina F, et al. High
425 Fragmentation Characterizes Tumour-Derived Circulating DNA. *PLOS ONE. Public*
426 *Library of Science;* 2011;6:e23418.
- 427 9. Snyder MW, Kircher M, Hill AJ, Daza RM, Shendure J. Cell-free DNA Comprises an In
428 Vivo Nucleosome Footprint that Informs Its Tissues-Of-Origin. *Cell.* 2016;164:57–68.
- 429 10. Jiang P, Lo YMD. The Long and Short of Circulating Cell-Free DNA and the Ins and
430 Outs of Molecular Diagnostics. *Trends Genet. Elsevier;* 2016;32:360–71.
- 431 11. Jiang P, Chan CWM, Chan KCA, Cheng SH, Wong J, Wong VW-S, et al. Lengthening
432 and shortening of plasma DNA in hepatocellular carcinoma patients. *Proc Natl Acad*
433 *Sci. National Academy of Sciences;* 2015;112:E1317–25.
- 434 12. Mouliere F, Chandrananda D, Piskorz AM, Moore EK, Morris J, Ahlborn LB, et al.
435 Enhanced detection of circulating tumor DNA by fragment size analysis. *Sci Transl*
436 *Med. American Association for the Advancement of Science;* 2018;10. Available from:
437 <https://stm.sciencemag.org/content/10/466/eaat4921>

- 438 13. Buedts L, Wlodarska I, Finalet-Ferreiro J, Gheysens O, Dehaspe L, Tousseyn T, et al.
439 The landscape of copy number variations in classical Hodgkin lymphoma: a joint KU
440 Leuven and LYSA study on cell-free DNA. *Blood Adv.* 2021;5:1991–2002.
- 441 14. Bayindir B, Dehaspe L, Brison N, Brady P, Ardui S, Kammoun M, et al. Noninvasive
442 prenatal testing using a novel analysis pipeline to screen for all autosomal fetal
443 aneuploidies improves pregnancy management. *Eur J Hum Genet.* 2015;23:1286–93.
- 444 15. Maaten L van der, Hinton G. Visualizing Data using t-SNE. *J Mach Learn Res.*
445 2008;9:2579–605.
- 446 16. Pons P, Latapy M. Computing Communities in Large Networks Using Random Walks.
447 In: Yolum pInar, Güngör T, Gürgeç F, Özturan C, editors. *Comput Inf Sci - ISCIS*
448 2005. Berlin, Heidelberg: Springer; 2005. page 284–93.
- 449 17. Adalsteinsson VA, Ha G, Freeman SS, Choudhury AD, Stover DG, Parsons HA, et al.
450 Scalable whole-exome sequencing of cell-free DNA reveals high concordance with
451 metastatic tumors. *Nat Commun.* Nature Publishing Group; 2017;8:1324.
- 452 18. Leary RJ, Sausen M, Kinde I, Papadopoulos N, Carpten JD, Craig D, et al. Detection of
453 Chromosomal Alterations in the Circulation of Cancer Patients with Whole-Genome
454 Sequencing. *Sci Transl Med.* American Association for the Advancement of Science;
455 2012;4:162ra154-162ra154.
- 456 19. Ulz P, Belic J, Graf R, Auer M, Lafer I, Fischereder K, et al. Whole-genome plasma
457 sequencing reveals focal amplifications as a driving force in metastatic prostate cancer.
458 *Nat Commun.* 2016;7:12008.
- 459 20. Molparia B, Nichani E, Torkamani A. Assessment of circulating copy number variant
460 detection for cancer screening. *PLOS ONE.* Public Library of Science;
461 2017;12:e0180647.
- 462 21. Liu MC, Oxnard GR, Klein EA, Swanton C, Seiden MV, Liu MC, et al. Sensitive and
463 specific multi-cancer detection and localization using methylation signatures in cell-free
464 DNA. *Ann Oncol.* 2020;31:745–59.
- 465 22. Sun K, Jiang P, Chan KCA, Wong J, Cheng YKY, Liang RHS, et al. Plasma DNA tissue
466 mapping by genome-wide methylation sequencing for noninvasive prenatal, cancer, and
467 transplantation assessments. *Proc Natl Acad Sci.* 2015;112:E5503–12.
- 468 23. Ulz P, Perakis S, Zhou Q, Moser T, Belic J, Lazzeri I, et al. Inference of transcription
469 factor binding from cell-free DNA enables tumor subtype prediction and early
470 detection. *Nat Commun.* Nature Publishing Group; 2019;10:4666.
- 471 24. Moss J, Magenheimer J, Neiman D, Zemmour H, Loyfer N, Korach A, et al.
472 Comprehensive human cell-type methylation atlas reveals origins of circulating cell-free
473 DNA in health and disease. *Nat Commun.* Nature Publishing Group; 2018;9:5068.
- 474 25. Sadeh R, Sharkia I, Fialkoff G, Rahat A, Gutin J, Chappleboim A, et al. ChIP-seq of
475 plasma cell-free nucleosomes identifies gene expression programs of the cells of origin.
476 *Nat Biotechnol.* Nature Publishing Group; 2021;1–13.

- 477 26. Klein EA, Richards D, Cohn A, Tummala M, Lapham R, Cosgrove D, et al. Clinical
478 validation of a targeted methylation-based multi-cancer early detection test using an
479 independent validation set. *Ann Oncol. Elsevier*; 2021;0. Available from:
480 [https://www.annalsofoncology.org/article/S0923-7534\(21\)02046-9/abstract](https://www.annalsofoncology.org/article/S0923-7534(21)02046-9/abstract)
- 481 27. Wan JCM, Heider K, Gale D, Murphy S, Fisher E, Mouliere F, et al. ctDNA monitoring
482 using patient-specific sequencing and integration of variant reads. *Sci Transl Med.*
483 *American Association for the Advancement of Science*; 2020;12. Available from:
484 <http://stm.sciencemag.org/content/12/548/eaaz8084>
- 485 28. Moss J, Zick A, Grinshpun A, Carmon E, Maoz M, Ochana BL, et al. Circulating breast-
486 derived DNA allows universal detection and monitoring of localized breast cancer. *Ann*
487 *Oncol. Elsevier*; 2020;31:395–403.
- 488 29. Vanderstichele A, Busschaert P, Smeets D, Landolfo C, Nieuwenhuysen EV, Leunen K,
489 et al. Chromosomal Instability in Cell-Free DNA as a Highly Specific Biomarker for
490 Detection of Ovarian Cancer in Women with Adnexal Masses. *Clin Cancer Res.*
491 *American Association for Cancer Research*; 2017;23:2223–31.
- 492 30. Kaijser J. Towards an evidence-based approach for diagnosis and management of adnexal
493 masses: findings of the International Ovarian Tumour Analysis (IOTA) studies. *Facts*
494 *Views Vis ObGyn.* 2015;7:42–59.
- 495 31. Cascione L, Aresu L, Baudis M, Bertoni F. DNA Copy Number Changes in Diffuse
496 Large B Cell Lymphomas. *Front Oncol. Frontiers*; 2020;10. Available from:
497 <https://www.frontiersin.org/articles/10.3389/fonc.2020.584095/full>
- 498 32. Zhang J, Grubor V, Love CL, Banerjee A, Richards KL, Mieczkowski PA, et al. Genetic
499 heterogeneity of diffuse large B-cell lymphoma. *Proc Natl Acad Sci.* 2013;110:1398–
500 403.
- 501 33. Kelemen LE, Köbel M. Mucinous carcinomas of the ovary and colorectum: different
502 organ, same dilemma. *Lancet Oncol.* 2011;12:1071–80.
- 503 34. Cheasley D, Wakefield MJ, Ryland GL, Allan PE, Alsop K, Amarasinghe KC, et al. The
504 molecular origin and taxonomy of mucinous ovarian carcinoma. *Nat Commun. Nature*
505 *Publishing Group*; 2019;10:3935.
- 506 35. Nishizuka S, Chen S-T, Gwadry FG, Alexander J, Major SM, Scherf U, et al. Diagnostic
507 Markers That Distinguish Colon and Ovarian Adenocarcinomas. *Cancer Res.*
508 2003;63:5243–50.
- 509 36. Wyatt AW, Annala M, Aggarwal R, Beja K, Feng F, Youngren J, et al. Concordance of
510 Circulating Tumor DNA and Matched Metastatic Tissue Biopsy in Prostate Cancer.
511 *JNCI J Natl Cancer Inst.* 2017;109. Available from: <https://doi.org/10.1093/jnci/djx118>
- 512 37. Wei T, Zhang J, Li J, Chen Q, Zhi X, Tao W, et al. Genome-wide profiling of circulating
513 tumor DNA depicts landscape of copy number alterations in pancreatic cancer with liver
514 metastasis. *Mol Oncol.* 2020;14:1966–77.

- 515 38. Cresswell GD, Nichol D, Spiteri I, Tari H, Zapata L, Heide T, et al. Mapping the breast
516 cancer metastatic cascade onto ctDNA using genetic and epigenetic clonal tracking. *Nat*
517 *Commun.* Nature Publishing Group; 2020;11:1446.
- 518 39. Abbosh C, Birkbak NJ, Wilson GA, Jamal-Hanjani M, Constantin T, Salari R, et al.
519 Phylogenetic ctDNA analysis depicts early-stage lung cancer evolution. *Nature.* Nature
520 Publishing Group; 2017;545:446–51.
- 521 40. Kustanovich A, Schwartz R, Peretz T, Grinshpun A. Life and death of circulating cell-
522 free DNA. *Cancer Biol Ther.* Taylor & Francis; 2019;20:1057–67.
- 523

524 **Table 1. Participant and characteristics**

	Stage*	Age, Mean ± SD	Female, n (%)	Total Samples
Hematological cancer dataset				
	Healthy	69 (±3)	164 (63%)	260
	Hodgkin's lymphoma (HL)	32 (±14)	98 (55%)	179
	I			10
	II			145
	III			9
	IV			15
	Diffuse large B-cell lymphoma (DLBCL)	59 (±13)	22 (60%)	37
	I			1
	II			5
	III			7
	IV			8
	unknown			16
	Multiple myeloma (MM)	67 (±9)	8 (36%)	22
	I			3
	II			7
	III			7
	unknown			5
Solid tumor dataset				
	Healthy	49 (±12)	107 (91%)	107
	Breast	56 (±12)	46 (100%)	46
	I			23
	II			12
	III			5
	IV			6
	Colorectal	66 (±12)	29 (41%)	70
	I			19
	II			17
	III			25
	IV			9
	Gastrointestinal stromal tumor (GIST)	64 (±11)	N.A.	35
	Advanced			35
	Lung		N.A.	44
	Advanced	N.A.	N.A.	44
	Ovarian invasive tumors	61 (±14)	125 (100%)	125
	I			25
	II			11
	III			49
	IV			31
	Metastatic			9
	Ovarian Benign	49 (±16)	160 (100%)	160
	Ovarian Borderline	51 (±17)	63 (100%)	63

525 * Multiple myeloma stratification refers to Revised International Staging System.

526 **Figure Captions**

527 **Fig. 1. Schematic illustration of GIPXplore.** Plasma cfDNA in healthy individuals (blue box)
528 comprises short nucleosome-protected DNA fragments mainly released from the cells of
529 hematopoietic origin. In patients with cancer (green box), cfDNA is also released from the
530 tumor. Since cfDNA fragmentation pattern is cell- or tissue-specific, sequencing and mapping
531 of cfDNA from a patient with cancer may have differential genome-wide distribution of DNA
532 fragments along the genome compared to a healthy one (green and blue profiles respectively).
533 The workflow of GIPXplore combines two tasks. First, explorative analysis of the high-
534 dimensional data is performed via unsupervised clustering. Data complexity is reduced by
535 using the first 50 linearly transformed genome-wide coverage features (non-trivial principal
536 components, PCs) from a large number of cfDNA profiles, which are used for dataset
537 exploration to unveil the potential biological signals or technical confounding factors based on
538 the sub-grouping of underlying patterns that facilitate the design of the supervised models.
539 Concurrently, classifiers are constructed to predict disease status and identify disease type to
540 assess the use of such transformed genome-wide features as a marker for diagnostic
541 application.

542

543 **Fig. 2. Genome-wide cfDNA profiles carry cancer type-specific patterns. A,** Two-
544 dimensional tSNE visualization of the clustering result. Sample type is annotated by point color
545 and community detection resulted clusters are annotated by point shape. Cluster numbers are
546 labeled in the center of the defined cluster. **B,** Sample distribution in each community detection
547 defined cluster is shown. The upper bar plot shows the total number of samples grouped in
548 each cluster and the lower bar plot depicts the proportion of each class of samples. **C,** Tumor
549 fraction estimated using ichorCNA. Red horizontal line indicates a detection limit of 3% tumor
550 fraction level. **D,** Examples of copy number profiles generated from ichorCNA for selected

551 clusters. In each copy number profile, color red represents copy number gains and green
552 represents copy number losses. The color is supposed to be interpreted together with the log
553 ratio values to pinpoint copy number gains or losses.

554

555 **Fig. 3. Plasma cfDNA genome-wide signatures enable hematological malignancies**

556 **detection and subtype prediction. A,** Sensitivities for detection of subtypes of hematological

557 malignancies. Performance for early and advanced stages for DLBCL and HL are shown. Three

558 (R-ISS) stages of MM are shown. 95% confidence interval is shown as an error bar. **B,** ROC

559 curves for performance comparison between the genome-wide feature analysis and ichorCNA

560 tumor fraction analysis. For the genome-wide feature analysis, decision value from SVM

561 prediction is used to build a dynamic threshold of true and false positives. Tumor fraction

562 values were used to construct ROC for ichorCNA analysis. **C,** Confusion matrix for tissue of

563 origin detection in hematological tumor. The color shading represents the proportion of

564 samples being correctly localized. The labeled numbers indicate the number of samples being

565 classified into the class.

566

567 **Fig. 4. Clustering analysis elucidates profile representations in solid tumors. A,** Two-

568 dimensional tSNE visualization of solid tumor dataset clustering result. Sample type is

569 annotated by point color and community detection resulted clusters are annotated by point

570 shape. Cluster numbers are labeled in the center of the defined cluster. **B,** Sample distribution

571 in each community detection defined cluster is shown. The upper bar plot shows the total

572 number of samples grouped in each cluster and the lower bar plot depicts the proportion of

573 each class of samples. **C,** Tumor fraction estimation for indicated types of solid tumors. Red

574 horizontal line indicates a detection limit of 3% tumor fraction level. **D,** Examples of copy

575 number profiles generated from ichorCNA for selected clusters.

576

577 **Fig. 5. Malignancy detection and typing in solid tumors.** **A**, Sensitivities for detection of
578 different types of solid malignancies detection. Performance for detection of early and
579 advanced stages of disease is shown. **B**, ROC curves for performance comparison between the
580 genome-wide feature analysis and ichorCNA tumor fraction analysis. **C**, Confusion matrix for
581 tissue of origin detection in solid malignancies. The color shading represents the proportion of
582 samples being correctly localized. The labeled numbers indicate the number of samples being
583 classified into the class.

Figure 1.

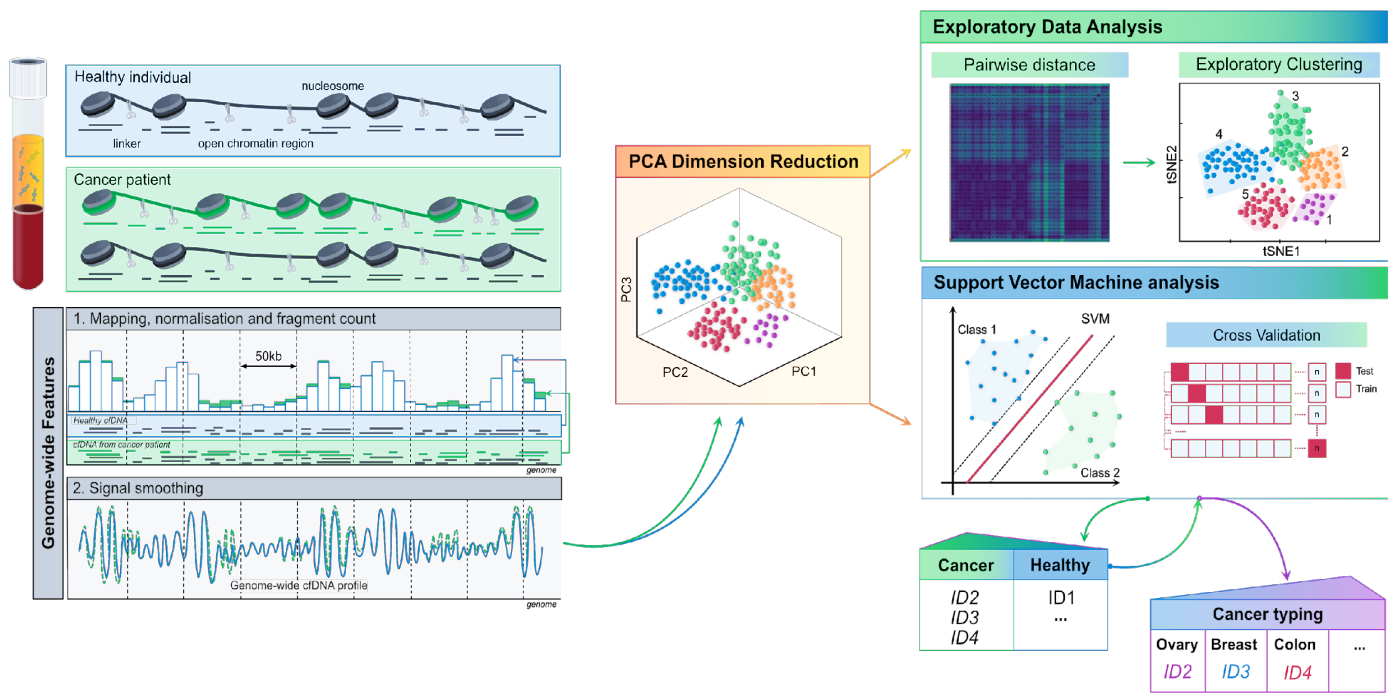


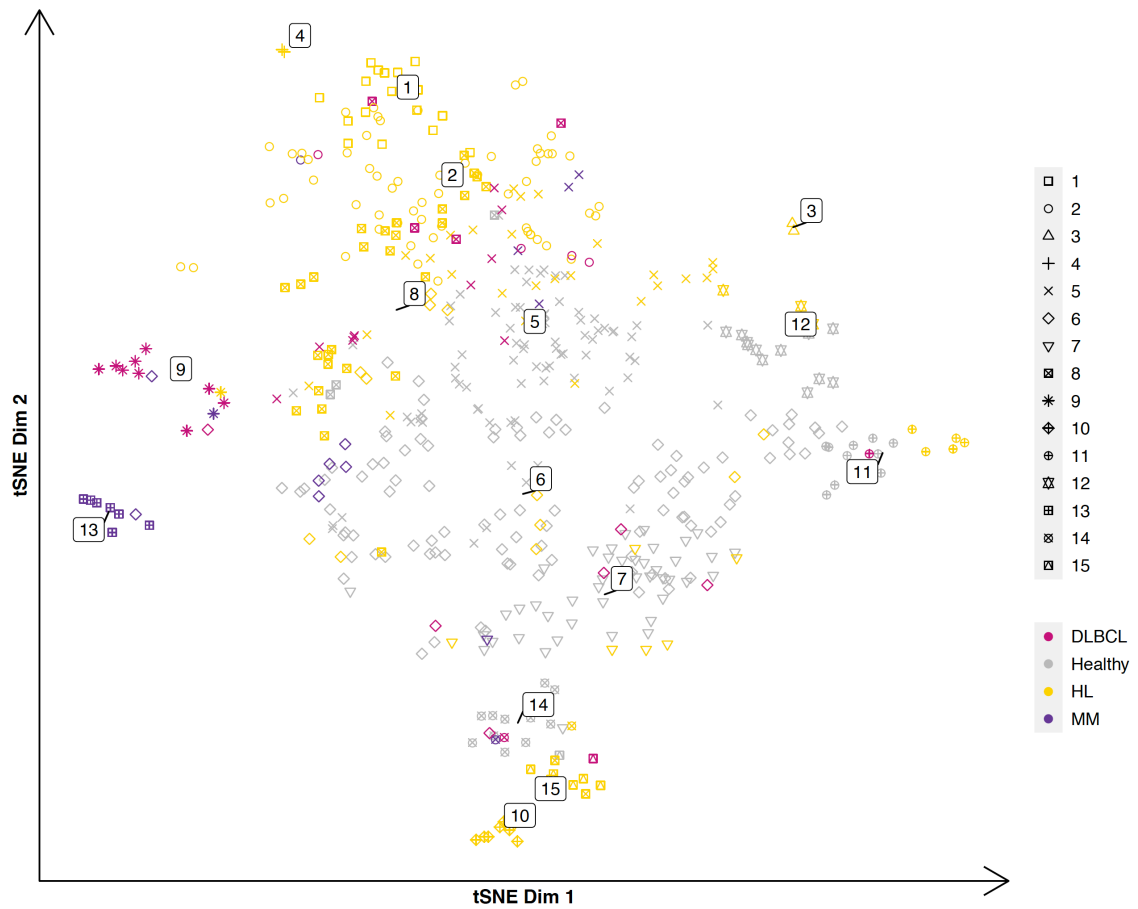
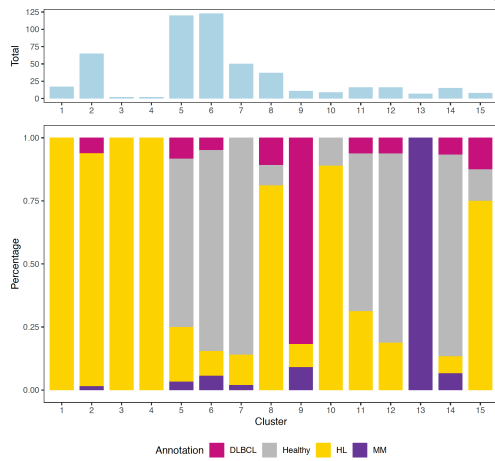
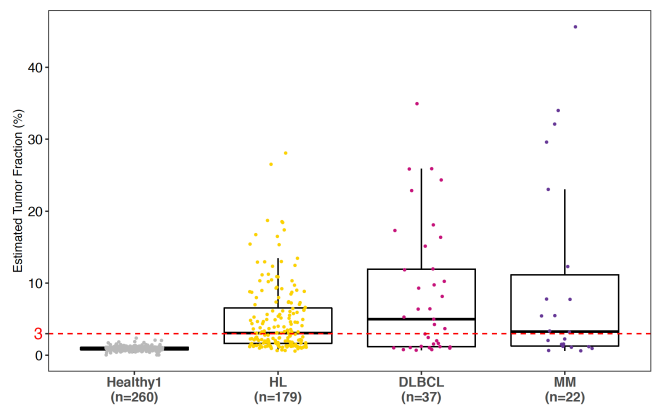
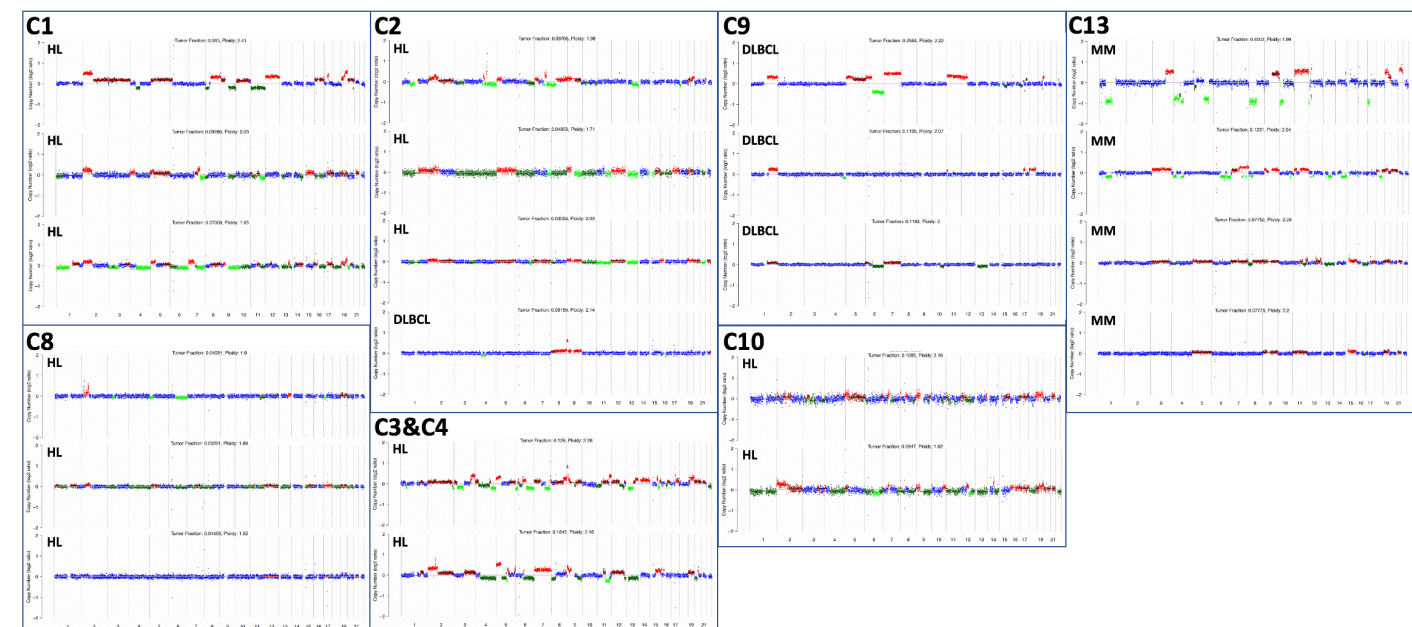
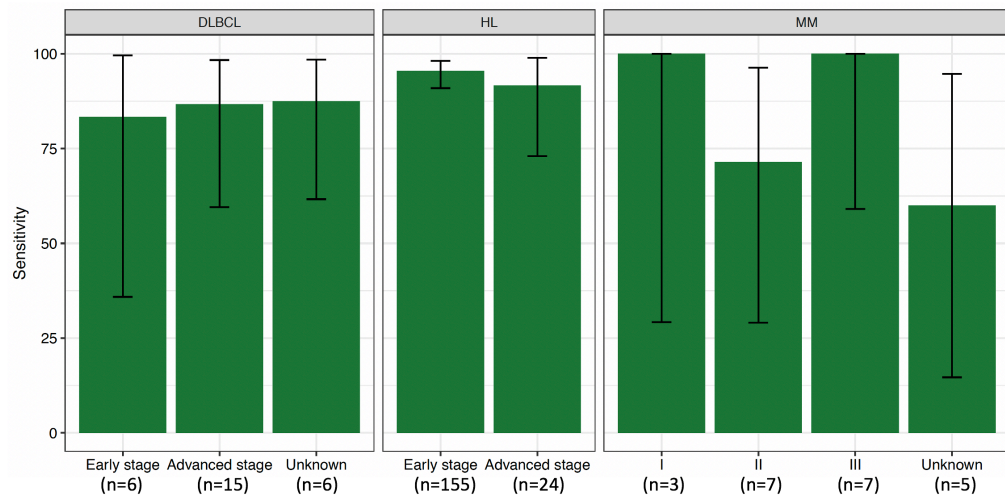
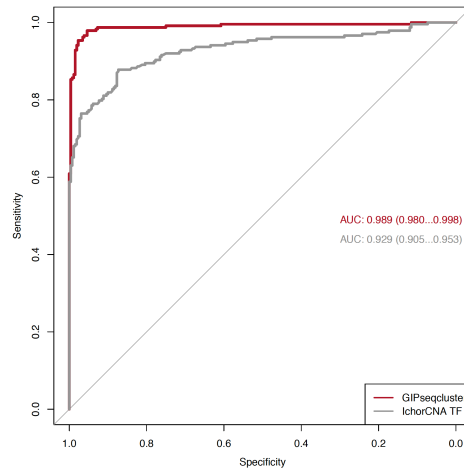
Figure 2.**A****B****C****D**

Figure 3.

A



B



C

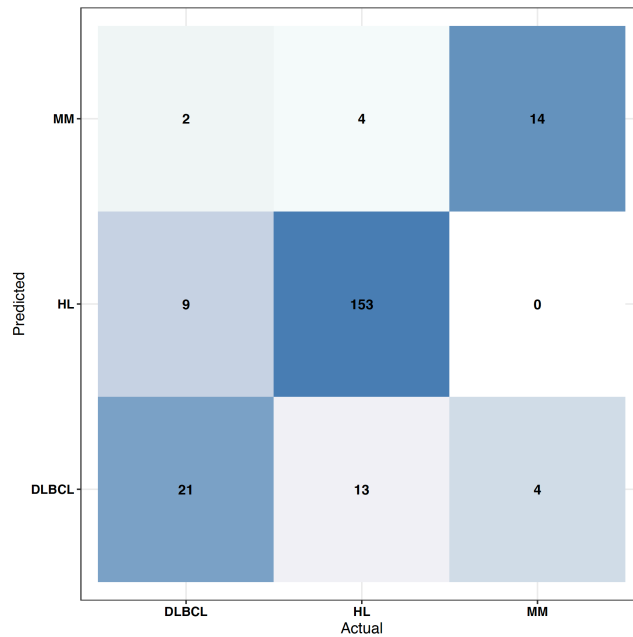
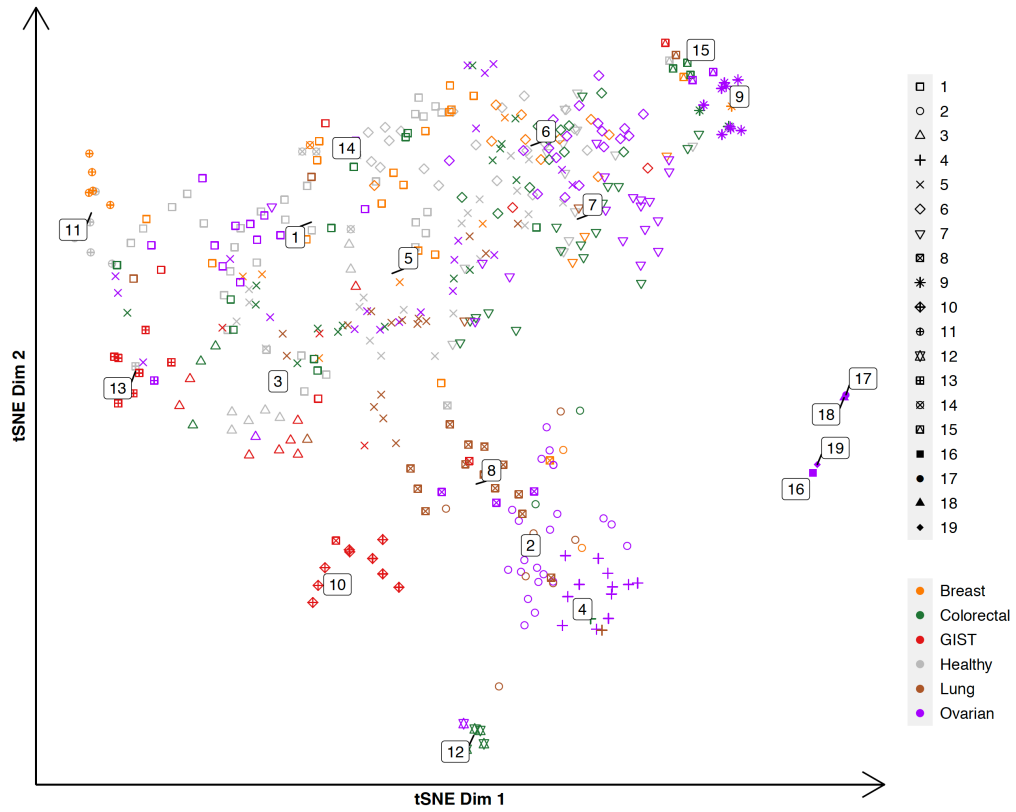
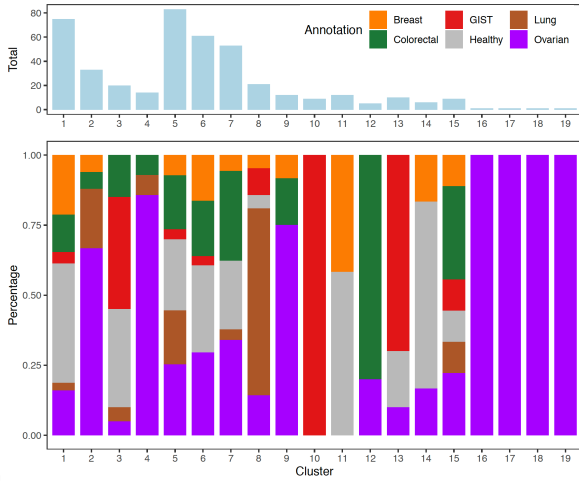


Figure 4.

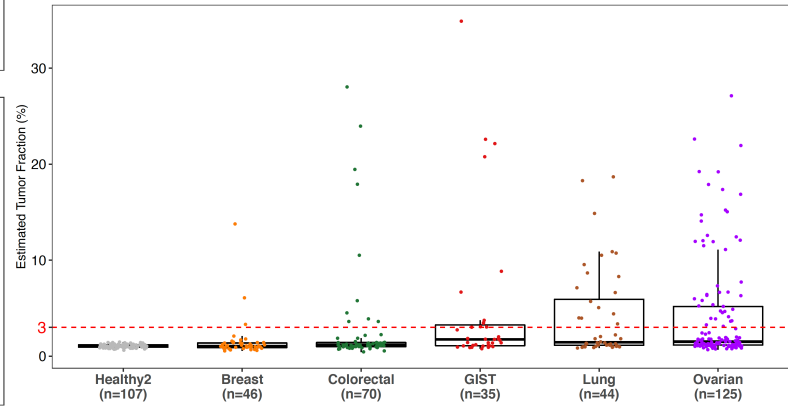
A



B



C



D

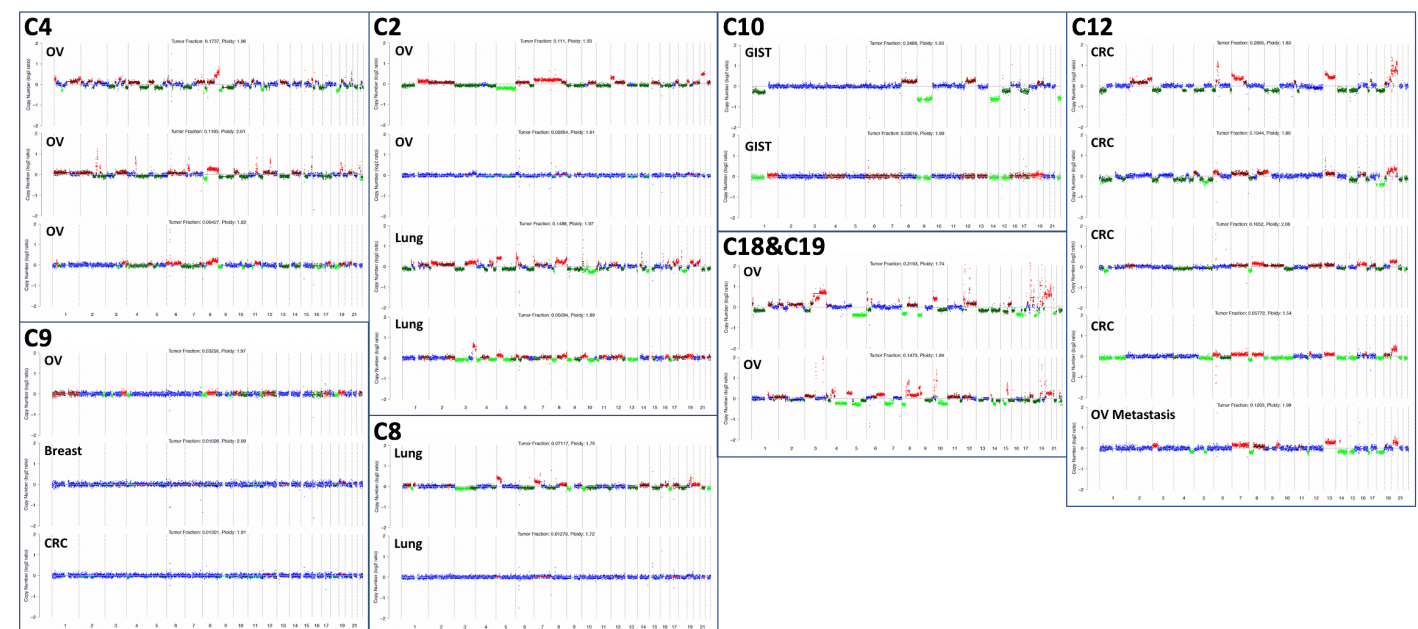
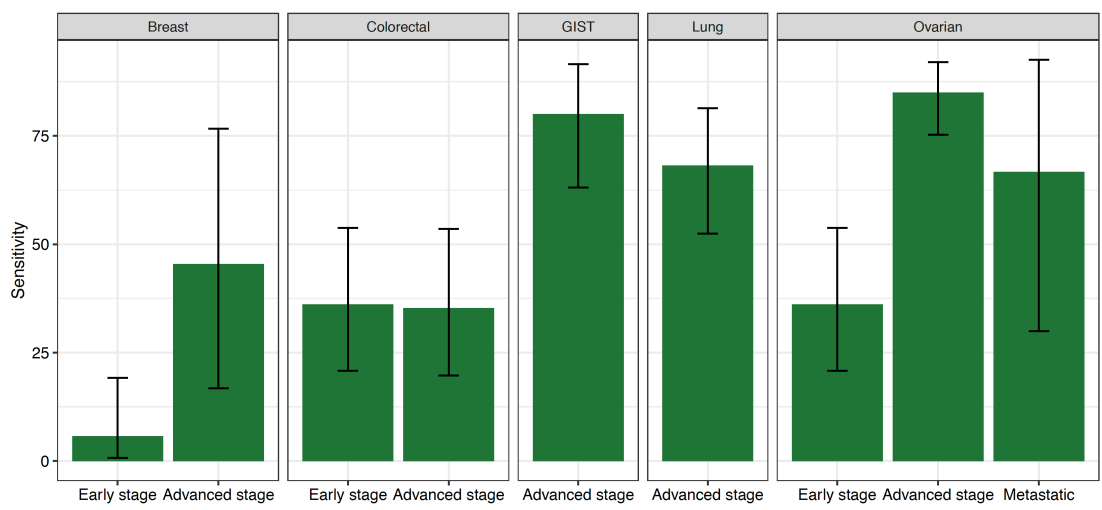
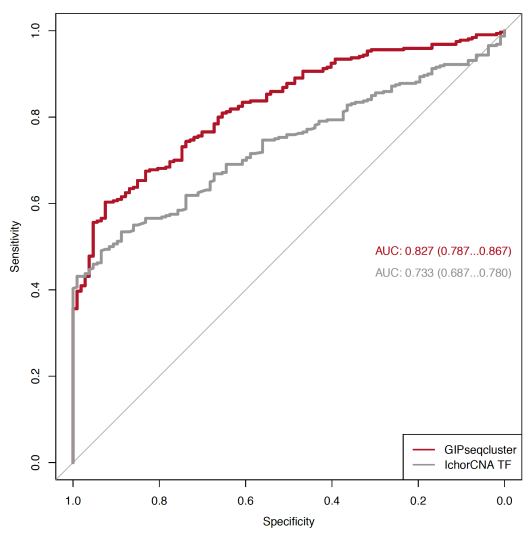


Figure 5.

A



B



C

