

# Utilizing Electronic Health Records (EHR) and Tumor Panel Sequencing to Demystify Prognosis of Cancer of Unknown Primary (CUP) patients

Intae Moon<sup>1,5,\*</sup>, Jaclyn LoPiccolo<sup>2</sup>, Sylvan C. Baca<sup>2, 3</sup>, Lynette M. Sholl<sup>4</sup>, Kenneth L. Kehl<sup>5</sup>, Michael J. Hassett<sup>5</sup>, David Liu<sup>2, 6</sup>, Deborah Schrag<sup>7</sup>, and Alexander Gusev<sup>5, 6, 8, \*</sup>

<sup>1</sup>Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, Cambridge, MA, USA

<sup>2</sup>Department of Medical Oncology, Dana-Farber Cancer Institute, Boston, MA, USA

<sup>3</sup>Center for Functional Cancer Epigenetics, Dana-Farber Cancer Institute, Boston, Massachusetts.

<sup>4</sup>Department of Pathology, Brigham and Women's Hospital, Harvard Medical School, Boston, MA, USA

<sup>5</sup>Division of Population Sciences, Dana-Farber Cancer Institute and Harvard Medical School, Boston, MA, USA

<sup>6</sup>The Broad Institute of MIT & Harvard, Cambridge, MA, USA

<sup>7</sup>Memorial Sloan Kettering Cancer Center, New York, NY, USA

<sup>8</sup>Division of Genetics, Brigham and Women's Hospital and Harvard Medical School, Boston, MA, USA

\*Corresponding author, [alexander\\_gusev@dfci.harvard.edu](mailto:alexander_gusev@dfci.harvard.edu)

## Abstract

When a standardized diagnostic test fails to locate the primary site of a metastatic cancer, it is diagnosed as a cancer of unknown primary (CUP). CUPs account for 3-5% of all cancers but do not have established targeted therapies, leading to typically dismal outcomes. Here, we develop OncoNPC, a machine learning classifier of CUP, trained on targeted next generation sequencing data from 34,567 tumors across 22 primary cancer types collected as part of routine clinical care at three institutions under AACR Project GENIE initiative [1]. OncoNPC achieved a weighted F1 score of 0.94 for high confidence predictions on known cancer types (65% of held-out samples). To evaluate its clinical utility, we applied OncoNPC to 971 CUP tumor samples from patients treated at the Dana-Farber Cancer Institute (DFCI). OncoNPC CUP

31 subtypes exhibited significantly different survival outcomes, and identified potentially actionable  
32 molecular alterations in 23% of tumors. Importantly, patients with CUP, who received first  
33 palliative intent treatments concordant with their OncoNPC predicted sites, showed significantly  
34 better outcomes (Hazard Ratio 0.348, 95% C.I. 0.210 - 0.570, p-value  $2.32 \times 10^{-5}$ ) after accounting  
35 for potential measured confounders. As validation, we showed that OncoNPC CUP subtypes  
36 exhibited significantly higher polygenic germline risk for the predicted cancer type. OncoNPC  
37 thus provides evidence of distinct CUP subtypes and offers the potential for clinical decision  
38 support for managing patients with CUP.

## 39 Introduction

40 When a standardized diagnostic work-up, including radiology and pathology review, fails to locate  
41 the primary site of a metastatic cancer, it is diagnosed as a cancer of unknown primary (CUP). CUP  
42 represents about 3-5% of all cancers worldwide [2] and is characterized by aggressive progression  
43 and poor prognosis (survival of 6 to 16 months [3]). The hidden nature of the primary cancer  
44 types for a CUP limits treatment options since clinical responses to some treatments are known to  
45 vary based on patients' tumor types (e.g., identical BRAF V600 mutations targetable in melanoma  
46 but no colorectal cancer[4]). Emerging cancer treatments targeting actionable molecular alterations  
47 are typically developed for specific cancer types: HER2 in breast cancer and EGFR mutation or  
48 ALK/ROS1 rearrangement in Non-small cell lung cancer (NSCLC) [5], and are thus inaccessible to  
49 CUP patients. Accurately identifying the latent primary site for CUPs and demonstrating clinical  
50 benefit from site-specific therapies may thus open many existing treatment options for patients with  
51 CUP.

52 Pathology review plays a key role in determining primary cancer types of malignant tumors based  
53 on immunohistochemistry (IHC) results as well as tumor morphology and clinical findings; however,  
54 pathological diagnosis can be challenging for highly metastatic or poorly differentiated tumors. For  
55 known cancer types, prior studies showed that an IHC-based diagnostic work-up correctly identified  
56 77 - 86% of primary tumors, which further decreased to 60 - 71% for metastatic tumors [6]. For  
57 patients with CUP, IHC results suggestive of a single primary diagnosis account for only 25% of  
58 tumors [3]. The subjective nature of pathological interpretation and guidelines, as well as the  
59 variability in IHC staining techniques across institutions thus makes it challenging to establish  
60 consistent protocols for CUP diagnosis [7].

61 Molecular tumor profiling has been proposed as an alternative for CUP primary classification  
62 due to its quantitative nature and high accuracy on tumors with known cancer types [8–12]. Such  
63 tools rely on microarray DNA methylation [8], whole genome sequencing (WGS) [9, 12], or RNA-  
64 seq data [11] to train machine learning classifiers using reference data from known-primary tumors.  
65 However, molecular sequencing remains prohibitive and not integrated into the existing standard  
66 of care, limiting the translational potential of such methods. Recently, key work by Penson et  
67 al. [10] demonstrated that accurate primary cancer type classifications could be made from next  
68 generation sequencing (NGS) of targeted panels, now routinely collected at many cancer centers and  
69 applicable to hundreds of thousands of tumors [1]. However, its clinical utility in diagnosing and

70 aiding treatment for patients with CUP was not systematically investigated.

71 Several recent studies have investigated the potential clinical benefit of molecular CUP clas-  
72 sification, in non-randomized prospective studies [13–15] as well as the randomized clinical trials  
73 [16]. These trials have often struggled to recruit sufficient numbers of representative patients and  
74 explore the full range of available therapies. A recent randomized phase II trial [16] did not find  
75 significant improvement in 1-year survival for the treatment group receiving site-specific therapy  
76 guided by molecular profiling. However, this study was limited by a small number of patients ( $n =$   
77 101) recruited over 7 years, with few common solid tumor types and well-established therapies [17].  
78 Assessing the clinical benefits of molecular CUP classification thus poses both an opportunity for  
79 precision medicine and a major challenge for conventional randomized studies.

80 In contrast to prospective trials, retrospective Electronic Health Records (EHR) data can cap-  
81 ture a larger and more heterogeneous patient population, despite potential biases due to informative  
82 missingness and unobserved heterogeneity. Coupling EHR data with tumor sequencing can offer  
83 insights into the molecular mechanisms of CUPs and their relationship to clinical outcomes. As  
84 panel sequencing is often part of the standard of care, such insights also have the potential to assist  
85 diagnostic efforts and clinical management within existing molecular workflows. Here, we utilized  
86 multi-center, Next Generation Sequencing (NGS) targeted panel sequencing data from 36,445 tumor  
87 samples with known primary cancers to train and evaluate a machine learning classifier predicting a  
88 primary cancer type of a given tumor sample. We applied this classifier, named *OncoNPC* (**O**ncology  
89 **N**GS-based **P**rimarily cancer type **C**lassifier), to 971 patients with CUP with clinical follow up at the  
90 Dana-Farber Cancer Institute (DFCI). Using the *OncoNPC* cancer type predictions, we identified  
91 CUP subtypes that shared specific characteristics with their corresponding predicted primaries in-  
92 cluding: significant differences in clinical outcomes, elevated germline risk, and prognostic somatic  
93 alterations. 23% of *OncoNPC* classified CUP tumors had actionable somatic variants enabled by  
94 their corresponding *OncoNPC* cancer type predictions. Finally, using EHR-based treatment and  
95 survival data, we showed that site-specific treatments concordant with the *OncoNPC* cancer type  
96 predictions led to longer survival than those discordant with the cancer type predictions. Our find-  
97 ings suggest that many CUPs can be classified into meaningful subtypes with the potential to aid  
98 clinical decision making.

## 99 Results

### 100 *OncoNPC* accurately classifies 22 known cancer types

101 We developed *OncoNPC* (**O**ncology **N**GS-based **P**rimarily cancer type **C**lassifier), a molecular can-  
102 cer type classifier trained on multicenter targeted panel sequencing data (Fig. 1). *OncoNPC* utilized  
103 all somatic alterations including mutations (single nucleotide variants and indels), mutational signa-  
104 tures, copy number alterations, as well as patient age at sequencing and sex to jointly predict cancer  
105 type using a XGBoost algorithm (see Methods). Importantly, no other aspects of tumor morphol-  
106 ogy, pathology, or patient demographics were used so as not to bias the classifier towards known  
107 cancers. *OncoNPC* was trained and validated on the processed data consisting of 29,176 primary

108 and metastasis tumor samples from 22 known cancer types collected at the DFCI, MSK, and VICC  
109 (see Table 1 for details). Across all 22 cancer types, OncoNPC achieved a weighted F1 score of 0.784  
110 on the held-out test tumor samples consisting of 7,289 tumor samples (weighted precision and recall  
111 : 0.789 and 0.791, respectively). Across 10 cancer groups (grouped by sites and treatment options  
112 (Table 1), OncoNPC achieved an overall weighted F1 score of 0.824 (weighted precision and recall :  
113 0.829 and 0.826, respectively). Despite the evident class imbalance across cancer types, OncoNPC  
114 showed well-balanced precision across the cancer types (Fig. 2a) and cancer groups (Fig. 2b).  
115 Thresholding on prediction confidence ( $p_{max}$ , the maximum posterior probability across all labels)  
116 further increased the performance: weighted F1 score of 0.830 with 91.6 % remaining samples at  
117  $p_{max} \geq 0.5$  and 0.942 with 65.2 % remaining samples at  $p_{max} \geq 0.9$  (Fig. 2c, 2d). While rarer cancer  
118 types had generally lower overall performance, increasing the  $p_{max}$  threshold reduced this difference  
119 between common/rare cancer types (Fig. 2c, 2d). At  $p_{max} \geq 0$ , common cancer types in the upper  
120 quartile in terms of the number of tumor samples (NSCLC, BRCA, COADREAD, DIFG, PRAD,  
121 and PAAD) had a mean F1 of 0.84 while rare cancer types in the lower quartile (WDTC, MNGT,  
122 GINET, PANET, AML, and NHL) had a mean F1 of 0.58, whereas at  $p_{max} \geq 0.9$  common and rare  
123 cancer had a mean F1 of 0.95 and 0.86, respectively. This demonstrates that the OncoNPC was  
124 still able to do high-quality predictions for a subset of tumor samples in rare cancer types, for which  
125 training data was limited.

126 OncoNPC achieved robust performance against potential dataset shifts due to the factors includ-  
127 ing cancer center, biopsy site type, sequence panel version, and patient ethnicity (Fig. 2e). OncoNPC  
128 showed comparable performance for tumor samples from DFCI (AUC-PR, area under the precision  
129 recall curve = 0.89, n = 3,690) and those from MSK (AUC-PR = 0.85, n = 3,331). OncoNPC  
130 performance for those from VICC was slightly lower (AUC-PR = 0.76, n = 268). OncoNPC showed  
131 comparable performance for primary tumor samples (AUC-PR = 0.87, n = 4,525) and metastatic  
132 tumor samples (AUC-PR = 0.87, n = 2,605), demonstrating its capability to predict the primary  
133 cancer site of metastatic cancers without loss of performance. To assess the OncoNPC performance  
134 over time, we investigated its performance across sequence panel versions utilized at DFCI, as the  
135 panel version is a proxy for sequence dates of tumor samples (see Table 1). The OncoNPC perfor-  
136 mance on tumor samples from earlier versions of DFCI sequence panels (OncoPanel v1 : AUC-PR  
137 = 0.82, n = 414 and OncoPnael v2 : AUC-PR = 0.89, n = 1,050) was slightly lower than the per-  
138 formance on the tumor samples from the most recent panel (OncoPanel v3 : AUC-PR = 0.91, n =  
139 2,226) which also contained the largest number of genes. As all tumor samples have been collected  
140 from OncoPanel v3 since October 2016, we expect our model to make high-quality predictions in  
141 a prospective setting. Finally, OncoNPC demonstrated consistent performance across patient eth-  
142 nicity, an important consideration to avoid introducing algorithmic disparities. See Supplementary  
143 Fig. S1a for more detailed center-specific OncoNPC performance.

## 144 **Applying OncoNPC to CUP tumor samples**

145 We applied OncoNPC to classify 971 CUP tumors from patients who were admitted to DFCI and  
146 sequenced as part of routine clinical care. Compared to the held-out cohort of Cancer with Known

147 Primary (CKP;  $n = 7,289$ ), OncoNPC classifications for CUPs had prediction probabilities lower  
148 than those of the DFCI held-out cohort of Cancer with Known Primary (CKP;  $n = 3,690$ ), but  
149 comparable to those of the DFCI held-out cohort of CKPs including other cancer types ( $n = 8,025$ ),  
150 indicating that CUPs may contain other hard-to-classify cancer types: mean prediction probabil-  
151 ity 0.764 (95% C.I. 0.750 - 778) for CUPs versus 0.881 (95% C.I. 0.875 - 0.887) for the held-out  
152 CKPs at DFCI and 0.769 (95% C.I. 0.764 - 0.774) for all held-out CKPs at DFCI (Fig. S1 and  
153 Supplementary Fig. S1b). However, more than half of the CUP tumors (518/971) could still be  
154 classified with high confidence (i.e., prediction probability  $> 0.8$ ), and multiple classified types had  
155 distributions of posterior probabilities comparable to their corresponding CKPs: Non-small Cell  
156 Lung Cancer (NSCLC), Invasive Breast Carcinoma (BRCA), Pancreatic Adenocarcinoma (PAAD),  
157 Prostate Adenocarcinoma (PRAD), and Gastrointestinal Neuroendocrine Tumors (GINET). Inter-  
158 estingly, CUPs with predicted GINET were highly confident, despite their small number of tumor  
159 samples in the training cohort ( $n = 359$ ; 0.99% of the training cohort), suggesting some rarer cancer  
160 types may nevertheless be confidently identifiable. As shown in Fig. 3b, the most common CUP can-  
161 cer types were Non-small Cell Lung Cancer (NSCLC), Pancreatic Adenocarcinoma (PAAD), Invasive  
162 Breast Carcinoma (BRCA), Esophagogastric Adenocarcinoma (EGC), and Colorectal Adenocarci-  
163 noma (COADREAD); of which NSCLC, BRCA, and COADREAD were also the most common  
164 CKP types. These rates are broadly consistent with prior findings that the most frequently revealed  
165 underlying primary cancers for CUPs by autopsy include lung, large bowel, and pancreas cancers  
166 [18]. Finally, comparable rates were observed upon applying OncoNPC to 581 CUP tumors at MSK  
167 (Supplementary Fig. S4)

## 168 **Explaining OncoNPC cancer type predictions**

169 OncoNPC learns complex non-linear relationships between input somatic variants and clinical fea-  
170 tures and provides interpretable primary cancer type predictions, where impact of each input feature  
171 on a prediction is quantified as a SHAP value [19]. We investigated the most impactful features in  
172 predicting each cancer type across the CKP and CUP cohorts to evaluate face validity of OncoNPC  
173 (see Fig. 3d for the top 3 most frequent cancer types in the cohort: NSCLC, BRCA, and PAAD, and  
174 Supplementary Fig. S2 and S3 for other cancer types). For NSCLC, the most important features  
175 were EGFR mutation and SBS4, a tobacco smoking-associated mutation signature [20], for CKP  
176 tumor samples and CUP with NSCLC predicted tumor samples, respectively; both consistent with  
177 the known etiology of lung cancer. Somatic mutation in the EGFR gene is frequently observed in  
178 NSCLC tumors and the gene itself is a well-known therapeutic target for patients with NSCLC [21,  
179 22]. Carcinogens in tobacco smoke have been known to cause lung cancer [23]. For BRCA, the  
180 most important feature for both CKP and CUP was sex, as expected, followed by CNA events in  
181 GATA3 and CCND1 genes, known drivers and prognostic indicators in breast cancer [24, 25]. For  
182 PAAD, KRAS mutation was significantly more common than the population averages and by far  
183 the most important somatic feature. Mutations in the KRAS gene occur frequently among patients  
184 with colorectal cancer and are known to have prognostic significance [26, 27].

185 OncoNPC provides intuitive illustrations of an explanation for individual-level predictions (Fig.

186 3e). As an example, we show the explained classification for a tumor sample biopsied from the  
187 liver of the 76 year-old male patient and subsequently diagnosed with CUP. From the chart review,  
188 we found that the patient reported a 60-pack year smoking history, as well as having lived near a  
189 tar and chemical factory as a child. Despite the CUP diagnosis, OncoNPC confidently classified  
190 the primary site as NSCLC with posterior probability of 0.98. SBS4, a tobacco smoking-associated  
191 mutation signature, was significantly enriched in the patient’s tumor sample, which has, by far,  
192 the most impact on the prediction; followed by SBS24 mutation signature associated with known  
193 exposures to aflatoxin [20]; and KRAS mutation. Note that inhalation of aflatoxin has been linked  
194 to cause primary lung cancer [28–30], and KRAS mutation is one of the most common drivers  
195 of NSCLC [31, 32]. The feature interpretation analysis demonstrated that OncoNPC was able to  
196 capture biologically consistent, cancer-type specific signals from interpretable somatic mutation and  
197 clinical features at an individual tumor level as well as a cohort level.

### 198 Germline PRS-based validation on CUP tumor samples

199 We hypothesized that, if OncoNPC was accurately identifying latent primary cancers, the classified  
200 CUP cancer types would exhibit increased germline risk for the corresponding cancers. To that end,  
201 we imputed common germline variation for each CUP patient and quantified their polygenic risk  
202 scores (PRS) across 8 common cancers using external cancer GWAS data (see Methods). PRSs are  
203 a continuous estimate of the underlying germline liability for a given cancer and orthogonal from the  
204 somatic data used to train OncoNPC. As hypothesized, patients with CUP had a significantly higher  
205 mean germline PRS for the OncoNPC predicted cancers (Fig. 3c and see Supplementary Fig. S6  
206 for cancer type-specific analysis) compared to other cancer types. The magnitude of the difference  
207 (i.e.,  $\hat{\Delta}_{\text{PRS}}$ ) increased for more confident OncoNPC predictions ( $\hat{\Delta}_{\text{PRS}} = 0.142$ , 95% C.I. 0.0494 –  
208 0.235, Wald test p-value:  $2.66 \times 10^{-3}$  and  $\hat{\Delta}_{\text{PRS}} = 0.204$ , 95% C.I. 0.0655 – 0.344, Wald test p-value:  
209  $3.98 \times 10^{-3}$  at  $p_{\text{max}}$  threshold = 0.0 and  $p_{\text{max}}$  threshold = 0.9, respectively). As a negative control,  
210 the same analysis conducted with randomly shuffled OncoNPC labels showed no enrichment. As a  
211 positive control, the same analysis conducted on CKPs, with available imputed PRS ( $n = 11,332$ ),  
212 also demonstrated a highly significant germline enrichment, as expected. Notably, the enrichment for  
213 CUPs was in between that of CKPs and random tumors, suggesting that while OncoNPC classified  
214 CUPs are genetically correlated with CKPs, they still exhibit additional heterogeneity.

### 215 OncoNPC-based risk stratification among patients with CUP

216 To demonstrate clinical utility of OncoNPC, we examined if OncoNPC cancer type predictions can  
217 stratify risk among patients with CUP. Using overall survival, we identified subtypes which had  
218 significant prognostic differences in median survival based on the OncoNPC classifications (Figure  
219 4a, Chi-squared test, p-value:  $4.90 \times 10^{-14}$ ). Overall, the poorest prognosis was observed in patients  
220 with CUP predicted to be Esophagogastric Adenocarcinoma (EGC) and Pancreatic Adenocarci-  
221 noma (PAAD): median survival 8.44 months for the combined cohort (95% C.I. 5.39 - 10.5,  $n =$   
222 107). The most favorable prognosis was observed in patients with CUP predicted to be Head and  
223 Neck Squamous Cell Carcinoma (HNSCC), Gastrointestinal Neuroendocrine Tumors (GINET), and

224 Pancreatic Neuroendocrine Tumors (PANET): median survival 48.2 months for HNSCC (95% C.I.  
225 19.6 - not estimable, n = 41) and not estimable median survival (i.e. the estimated survival curve  
226 never reached the median) for the combined GINET and PANET cohort (n = 57), respectively.  
227 Our identified favorable subtypes are consistent with established favorable CUP subtypes such as  
228 poorly or well differentiated neuroendocrine carcinomas of unknown primary and squamous cell car-  
229 cinoma of non-supraclavicular cervical lymph nodes [33]. OncoNPC subtypes can thus be leveraged  
230 to meaningfully stratify patients by expected median survival.

### 231 **CUP-CKP metastatic survival comparison**

232 We investigated if cancer-specific prognosis is shared between CUP predicted cancer and their cor-  
233 responding CKP metastatic cancers. Utilizing overall survival data linked to the National Death  
234 Index and in-house follow-up data (see Methods), we found that median survival times of CUP-  
235 metastatic CKP pairs were significantly correlated across the cancer types (Spearman's  $\rho$ : 0.964,  
236 p-value:  $4.54 \times 10^{-4}$ ; Fig. 4b). This significant relationship provides evidence that genetics-based  
237 OncoNPC predictions capture prognostic signals specific to each predicted cancer type. While corre-  
238 lated, median survival times were significantly lower for patients with CUP compared to those with  
239 metastatic CKP: CUP median survival 14.0 months (95% C.I. 11.9 - 15.8, n = 685) vs. metastatic  
240 CKP median survival 23.1 months (95% C.I. 21.8 - 24.2, n = 7,797). This is expected as CUPs  
241 are an advanced metastatic cancer with limited treatment options [33]. The absolute difference in  
242 median survival was significant across all predicted CUP - metastatic CKP pairs with the exception  
243 of Pancreatic Adenocarcinoma (CUP PAAD median survival 8.61 months 95% C.I. 5.09 - 10.8 vs.  
244 metastatic CKP PAAD median survival 6.73 months 95% C.I. 5.98 - 8.02), known to be a particularly  
245 deadly cancer type.

### 246 **Shared prognostic somatic variants in CUP-metastatic CKP pairs**

247 We aimed to identify prognostic somatic variants shared between OncoNPC CUP subtypes and their  
248 corresponding metastatic CKP cancers. Three out of 14 tested CUP-metastatic CKP pairs (NSCLC,  
249 PAAD, and COADREAD) exhibited shared prognostic somatic variants significantly associated with  
250 overall survival with nominal p-value cut-off at 0.05 (Fig. 4c and 4d). In patients with known  
251 or classified NSCLC, three somatic mutations were associated with poor survival in both groups:  
252 SMARCA4 (CUP: H.R. 1.86, 95% C.I. 1.19 - 2.89, p-value  $6.23 \times 10^{-3}$ , CKP mets: H.R. 1.73,  
253 95% C.I. 1.44 - 2.09, p-value  $9.30 \times 10^{-9}$ ), STK11 (CUP: H.R. 1.76, 95% C.I. 1.14 - 2.71, p-value  
254  $1.05 \times 10^{-2}$ , CKP mets: H.R. 1.43, 95% C.I. 1.22 - 1.68, p-value  $1.00 \times 10^{-5}$ ), and KEAP1 (CUP:  
255 H.R. 1.83, 95% C.I. 1.18 - 2.85, p-value  $6.82 \times 10^{-3}$ , CKP mets: H.R. 1.40, 95% C.I. 1.18 - 1.66,  
256 p-value  $1.27 \times 10^{-4}$ ). These associations of somatic mutations in SMARCA4, STK11, and KEAP1  
257 genes with overall survival are well established for NSCLC [34–36]. Interestingly, a CNA event in  
258 NKX2-1 was associated with improved survival in the patients from the NSCLC pair (CUP: H.R.  
259 0.542, 95% C.I. 0.326 - 0.901, p-value  $1.83 \times 10^{-2}$ , CKP mets: H.R. 0.770, 95% C.I. 0.662 - 0.894,  
260 p-value  $6.28 \times 10^{-4}$ ), consistent with prior meta-analyses [37]. In patients with known or classified  
261 COADREAD tumors, SBS10b mutation signature, linked to polymerase epsilon exonuclease domain

262 mutations [20], was associated with longer overall survival (CUP: H.R. 0.371, 95% C.I. 0.148 - 0.928,  
263 p-value  $3.41 \times 10^{-2}$ , CKP mets: H.R. 0.495, 95% C.I. 0.255 - 0.958, p-value  $3.68 \times 10^{-2}$ ). Finally, in  
264 patients with known or classified PAAD tumors, the SBS29 mutation signature (commonly found in  
265 tumor samples from individuals with a tobacco chewing habit [20]) was associated with poor survival  
266 in CUPs but nominally protective in metastatic CKPs (CUP: H.R. 2.66, 95% C.I. 1.02 - 6.93, p-value  
267  $4.46 \times 10^{-2}$ , CKP mets: H.R. 0.657, 95% C.I. 0.438 - 0.986, p-value  $4.28 \times 10^{-2}$ ). Although these  
268 somatic associations remain to be validated in independent cohorts, by categorizing patients with  
269 CUP based on their OncoNPC predictions, we were able to identify prognostic somatic variants,  
270 consistent with recent research findings.

## 271 Identifying actionable somatic variants in CUP tumors based on OncoNPC 272 predictions

273 We investigated if OncoNPC classifications could identify genetically driven, site-specific treatment  
274 options that are typically available for cancers with known primaries. We utilized OncoKB [38] as  
275 a knowledge base and considered three different categories of actionable somatic variants: onco-  
276 genic mutation, amplification, and fusion (see Methods). OncoNPC cancer type predictions enabled  
277 identification of actionable somatic variants across CUP tumor samples (total 22.8% of the eligible  
278 CUP tumor samples; see Fig. 5a and Fig. 5b). The majority of actionable somatic variants for  
279 patients with CUP were oncogenic mutations (183 counts; 87.1%), followed by amplifications (22  
280 counts; 9.52%) and fusions (7 counts; 3.33%) as shown in Fig. 5a. The four most frequent oncogenic  
281 mutations were in PIK3CA, KRAS, ALK, and ERBB2 genes, occurring in CUP tumor samples  
282 classified as BRCA (PIK3CA and ERBB2 genes) and NSCLC (KRAS, ALK, and ERBB2 genes).  
283 Overall, among the eligible CUPs whose prediction confidences are greater than 0.5 ( $N = 794$ ; see  
284 Supplementary Fig. S5 for more details on the exclusion criteria), OncoNPC predictions identified  
285 actionable somatic variants for 11.5% of the CUP tumor samples for Level 1 therapeutic level (FDA-  
286 approved drugs), 3.63% for Level 2 (Standard care), 6.64% for Level 3 (Clinical evidence), and 1.00%  
287 for Level 4 (Biological evidence), summing up to the total 22.8% of the eligible CUP tumor samples  
288 (Fig. 5b).

## 289 Survival benefit of treatment concordance with OncoNPC predictions

290 We performed retrospective survival analysis to investigate whether patients with CUP achieved  
291 clinical benefit when treated in concordance with their OncoNPC classifications. We restricted to  
292 a cohort of 158 patients with CUP, received first treatment at DFCI with a palliative intent (see  
293 the exclusion criteria in Supplementary Fig. S5). Each case was then manually chart reviewed  
294 by a certified oncologist to determine whether the treatment administered was concordant with the  
295 OncoNPC prediction per National Comprehensive Cancer Network (NCCN) guidelines or standard of  
296 care (see Methods, Fig. 5c, and Fig. 5d). Strikingly, patients with CUP who received first palliative  
297 treatments concordant with their OncoNPC predicted cancer types exhibited significantly better  
298 survival than those who received discordant treatments as shown in Fig. 5e and 5f (*multivariable Cox*  
299 *regression*: H.R. 0.348, 95% C.I. 0.210 - 0.570, p-value  $2.32 \times 10^{-5}$ , Proportional Hazard assumption



300 test [39]: Chi-squared test with 17 degrees of freedom p-value 0.156, *IPTW Kaplan-Meier estimator*:  
301 weighted log-rank test p-value  $4.25 \times 10^{-10}$ ). Finally, after stratifying by OncoNPC predicted cancers  
302 and repeating the IPTW Kaplan-Meier analysis, we found that the treatment concordant group had  
303 improved survival across cancer cohorts (breast, GI, and others), with the exception of the lung  
304 cancer cohort (Supplementary Fig. S7).

305 We note that as this was not a randomized analysis, a potential concern may be systematic  
306 differences between the concordant and discordant groups leading to a significant prognostic but  
307 not predictive difference [40]. For example, treatment discordant patients may have systematically  
308 more advanced/de-differentiated tumors and thus exhibit poorer survival regardless of their treat-  
309 ment regimen. (see Table 2 for comparison of the two groups across the measured covariates). To  
310 minimize biases from potential confounders and move towards a predictive estimate of treatment  
311 concordance on patient survival, we adopted two estimation strategies: multivariable Cox regression  
312 [41] (i.e., covariate adjustment) and Inverse Probability of Treatment Weighted (IPTW) Kaplan-  
313 Meier estimator [42] (see Methods), which have recently been employed to emulate estimates from  
314 randomized trials [43, 44]. In both multivariable Cox regression and IPTW Kaplan-Meier estimator  
315 strategies, patients treated like their OncoNPC predicted cancer types (i.e. those in the concor-  
316 dant treatment group) consistently showed significantly better survival compared to those in the  
317 discordant treatment group. The multivariable Cox regression (Fig. 5e) additionally identified sig-  
318 nificant hazardous effects of age, gastrointestinal (GI) cancer types predicted by OncoNPC, and  
319 bone metastasis (H.R. 1.27, 95% C.I. 1.02 – 1.58, p-value  $3.10 \times 10^{-2}$ , H.R. 4.20, 95% C.I. 2.06 –  
320 8.55, p-value  $7.78 \times 10^{-5}$ , and H.R. 3.73, 95% C.I. 1.84 – 7.59, p-value  $2.71 \times 10^{-4}$ , respectively), and  
321 significantly protective effects of tumor mutational burden (TMB), as well as adenocarcinoma and  
322 neuroendocrine tumor group determined by the histopathology results (H.R. 0.537, 95% C.I. 0.388  
323 - 0.742, p-value  $1.64 \times 10^{-4}$ , H.R. 0.439, 95% C.I. 0.272 - 0.710, p-value  $7.85 \times 10^{-4}$  and H.R. 0.0854,  
324 95% C.I. 0.0298 - 0.245, p-value  $4.79 \times 10^{-6}$ , respectively). In the IPTW Kaplan-Meier analysis, we  
325 found that treatment concordance with the OncoNPC prediction was associated with Gastrointesti-  
326 nal (GI) cancer types (coefficient 1.916, 95% C.I. 0.627 - 3.205, p-value  $3.57 \times 10^{-3}$ ), whereas male  
327 sex and OncoNPC prediction uncertainty (i.e., entropy of predicted probability distribution over the  
328 considered cancer types) were inversely associated with receiving concordant treatment (coefficient  
329 -1.259, 95% C.I. -2.283 - -0.234, p-value  $1.61 \times 10^{-2}$ , and coefficient -1.693, 95% C.I. -2.458 - -0.927,  
330 p-value  $1.46 \times 10^{-5}$ ) (see Supplementary Fig. S8). These associations with treatment concordance are  
331 consistent with likely GI CUPs being more clinically identifiable and low OncoNPC confidence CUPs  
332 being less clinically identifiable. We note, however, that the IPTW approach specifically adjusts for  
333 these systematic differences when estimating the effect of treatment concordance on survival.

## 334 Discussion

335 Our work provides unique insights into the genetic and prognostic landscapes of CUP tumor samples  
336 by utilizing routinely collected EHR and multicenter NGS tumor panel sequencing data. We have  
337 developed OncoNPC, a machine learning model for molecular classification of tumor samples based  
338 on the NGS panel data. When evaluated with the held-out multicenter test data, OncoNPC provided

339 robust and interpretable predictions. Applying OncoNPC to CUP tumor samples, we demonstrated  
340 that the OncoNPC CUP subtypes showed significantly higher germline PRS risk for their predicted  
341 cancer. To our knowledge, this is the first evidence of germline genetic correlation between CUPs  
342 and corresponding known primaries, and lends orthogonal support to the molecular classification of  
343 CUPs into subtypes. We demonstrated clinical utility of the OncoNPC CUP subtypes by showing  
344 significant survival differences across subtypes, and, within subtypes, potentially actionable somatic  
345 alterations in 11.5% (Level 1 therapeutic level) and 22.8% (all levels) of tumors. Finally, in a  
346 retrospective analysis, we showed that patients with CUP, that had been treated in a consistent  
347 manner with their OncoNPC classification, achieved significantly longer survival than those treated  
348 in an inconsistent manner (multivariable Cox regression: H.R. 0.348, 95% C.I. 0.210 - 0.570, p-value  
349  $2.32 \times 10^{-5}$ ). Our findings suggest that CUP tumors share a genetic and prognostic architecture with  
350 known cancer types, and may benefit from molecular classification with OncoNPC for prognosis as  
351 well as treatment decision-making.

352 The question of whether CUP tumors consist of heterogeneous latent primaries or are a unique  
353 cancer type in and of themselves has been actively investigated [18, 45, 46]. Prior studies have  
354 demonstrated accurate classification of known tumors using Whole-Genome Sequencing [12], NGS  
355 panels [10], RNA-seq [11], methylation [8], and other platforms [47, 48]. However, these algorithms  
356 typically applied classification to metastatic tumors of known types and did not investigate the  
357 clinical implications for CUPs at large scale. Moran et al., [8] observed a nominally significant  
358 difference in survival between patients with CUP who received site-specific treatments concordant  
359 with their molecular primary site predictions and those who received empiric treatments. While  
360 promising, it remains unknown whether this difference is due to accurate classification for the site-  
361 specific group or systematically worse outcomes for the empirically treated group, which is typically  
362 a more challenging patient population [49]. To explicitly distinguish these scenarios, our analysis  
363 instead restricted to a CUP cohort wherein all patients received site-specific treatments as the  
364 first palliative-intent therapy and estimated a significant survival benefit of concordant treatment  
365 vs. discordant treatment (excluding the empirically treated group). Our findings were obtained  
366 after adjusting for left-truncation for sequencing time and measured potential confounders through  
367 covariate adjustment as well as propensity score weighting, which have been recently employed to  
368 mimic clinical trials in Real World data [43, 44]. Although we cannot rule out potential biases from  
369 unmeasured confounders, our cohort includes more heterogeneous populations compared to recruited  
370 cohorts in randomized controlled trials (RCT), and the proposed intervention (concordant treatment  
371 vs. discordant treatment) is challenging to ethically evaluate through RCTs, necessitating the use  
372 of retrospective causal inference.

373 Our study has several limitations. Firstly, although we utilized multicenter NGS tumor panel  
374 sequencing data to train OncoNPC model for cancer type prediction, we utilized retrospective EHR  
375 data from a single institution for the downstream clinical analyses. As a result, these analyses may  
376 be susceptible to systematic ascertainment patterns or biases specific to a tertiary academic cancer  
377 center. Replication of our clinical findings in other institutions is thus necessary to generalize our  
378 results. Secondly, we considered only the 22 most common cancer types in the cohort as classification  
379 labels (68.1 % of all tumor samples at DFCI, and 69.9 % across all three centers). As a result, if a

380 CUP tumor sample harbors a distinct yet not modeled primary cancer type, then the tumor sample  
381 will likely have high uncertainty in the prediction (see Supplementary Fig. S1b). Nevertheless, prior  
382 work has shown that the majority of resolvable primary sites of CUP tumor samples were from  
383 common cancers (e.g., lung, pancreas, and GI) [18], consistent with our findings. As more diverse  
384 tumor samples are collected across multiple institutions, our model can be augmented to robustly  
385 predict rare cancer types as well. Thirdly, our classifier and analyses relied on data from panel  
386 sequencing assays targeting 300-500 genes, which are inherently only sensitive to coding mutations  
387 and deep copy number alterations in the targeted genes. Other features captured by whole-genome  
388 sequencing or molecular assays may thus achieve better classification performance. Our focus in this  
389 work was on assays that are in routine clinical use as those are linked to Real World clinical data  
390 and offer the most immediate translational potential.

391 Our findings strongly suggest that routinely collected targeted tumor panel sequencing data have  
392 clinical utility in assisting diagnostic work-up and prognosis, and may additionally inform treatment  
393 decisions. To date, clinical sequencing is primarily used for identification of known biomarkers and  
394 corresponding clinical trial enrollment [50–53], and our findings additionally support use of panel  
395 sequencing for diagnosis. Conventional IHC-based pathology reviews are often unable to identify a  
396 primary diagnosis for advanced metastatic tumor samples [3, 6], particularly in community clinics  
397 where resources are limited. And in many cases, patients do not receive the complete diagnostic  
398 work-up that is recommended for CUPs [54]. As a result, oncologists resort to empiric treatment  
399 regimens to treat many patients with CUP [18] even when targeted therapies would otherwise be  
400 the standard of care for a corresponding known primary. In future work, we envision a multimodal  
401 framework that incorporates molecular sequencing together with patient pathology images [48],  
402 physiological data, and clinical notes to directly predict optimal treatment regimens rather than  
403 just cancer types. Our work thus paves a way for incorporating routine panel sequencing data into  
404 clinical decision support tools for clinically challenging cases.

## 405 **Methods**

### 406 **Patients and tumor samples**

407 We used the next generation sequencing (NGS) targeted panel sequencing data collected at three  
408 institutions in routine clinical care as part of the AACR project GENIE [1]: Dana-Farber Can-  
409 cer Institute (DFCI, n=18,816), Memorial Sloan Kettering Cancer (MSK, n=16,294) center, and  
410 Vanderbilt-Ingram Cancer Center (VICC, n=1,335). The collected tumor samples represented 22  
411 different cancer types and included 971 total samples from cancer of unknown primary (CUP). Na-  
412 tional Death Index (NDI) and clinical death and last clinical appointment records were available  
413 for 20,281 DFCI patients (n = 16,376 for CKP and n = 838 for CUP). Demographic details of the  
414 patients and tumor samples can be found in Table 1.

415 The cancer centers, DFCI, MSK, and VICC, were chosen because of similar genomic data char-  
416 acterization of their sequence panels in terms of coverage and alteration types [1]. DFCI samples  
417 were sequenced using a custom, hybridization-based panel called OncoPanel which targeted exons of

418 275-447 genes across three panel versions [1, 52]. MSK samples were sequenced using a custom panel  
419 called MSK-IMPACT which targeted 341-468 genes across 3 panel versions [1, 51]. VICC samples  
420 were sequenced using custom panels called VICC-01-T5A and VICC-01-T7, which targeted 322 and  
421 429 genes, respectively [1]. All panels were capable of detecting single nucleotide variants (SNVs),  
422 small indels, copy number alterations, and structural variants [1].

423 The DFCI CUP cohort consisted of 971 sequenced tumor samples (from 962 patients) with  
424 a cancer diagnosis of CUP and the following detailed cancer type: Adenocarcinoma, Not Other-  
425 wise Specified (NOS) (n = 345), Cancer of Unknown Primary, NOS (n = 194), Squamous Cell  
426 Carcinoma, NOS (n = 114), Poorly Differentiated Carcinoma, NOS (n = 118), Neuroendocrine  
427 Tumor/Carcinoma, NOS (n = 170), Small Cell Carcinoma of Unknown Primary (n = 16), Undiffer-  
428 entiated Malignant Neoplasm (n = 12), and Mixed Cancer Types (n = 2). For downstream clinical  
429 analyses, we applied additional exclusion criteria, described in Supplementary Fig. S5.

### 430 **Developing OncoNPC cancer type classifier**

431 We used a gradient tree boosting framework (XGBoost [55]) to develop OncoNPC for predicting  
432 cancer types from molecular features. In this framework, decision trees for the input features are  
433 sequentially added to an existing ensemble of the trees, such that the algorithm fits the new tree to  
434 the residuals from the ensembles with regularization on the tree structure. As the trees (a.k.a. weak  
435 learners) are added, the model learns optimal weights to combine their predictions and produces the  
436 improved outcome from the combined ensemble [55]. Owing to its high performance and scalability,  
437 the XGBoost method has been used across a wide range of applications in the healthcare space  
438 [56–58].

439 OncoNPC was trained and evaluated using tumors from 22 known cancer types split into 29,176  
440 training samples and 7,289 test samples. Hyper-parameter selection was conducted using random  
441 search [59] with 10-fold cross validation within the training set while utilizing weighted F1 score as an  
442 evaluation metric. The optimal hyper-parameters were then selected and the model was evaluated  
443 on the held-out test set (n = 7,289). To predict primary sites of CUP tumors, the model was  
444 then re-trained on all CKP tumor samples and applied to the CUP tumors to estimate posterior  
445 probabilities across the 22 different cancer labels. For each tumor sample, a cancer type with the  
446 highest probability was chosen as the predicted primary site.

### 447 **Feature selection and OncoNPC model interpretation**

448 The OncoNPC model was trained on somatic variant features from tumor sequencing data, as  
449 well as patient age at sequencing and sex. Other demographic/clinical features were intentionally  
450 not used so as not to bias the model toward cancer types with more available information. Somatic  
451 variant features included: mutations (i.e., single nucleotide variants (SNV) and indels), Copy Number  
452 Alteration (CNA) events, and mutational signatures [60]. For each gene, the total count of a somatic  
453 mutation (i.e., single nucleotide variants and indels) was encoded as a positive integer feature. The  
454 presence of a CNA event for each gene was encoded as a categorical variable with 5 levels: -2 (deep  
455 loss), -1 (single-copy loss), 0 (no event), 1 (low-level gain), and 2 (high-level amplification); note

456 that CNA events data for tumor samples from MSK and VICC were encoded as -2 (deep loss), 0  
457 (no event), and 2 (high-level amplification). Each of 60 different mutation signatures was inferred as  
458 the dot product of the weights derived from [60] and 96 single base substitutions in a trinucleotide  
459 context. The single base substitutions were computed using the `deconstructSigs` R library [61].  
460 See Supplementary Table S1 for the full set of features.

461 To identify important features in the OncoNPC’s predictions, we used the recently proposed  
462 feature interpretation tool for tree-based models, called TreeExplainer [19] (Python package `shap`).  
463 TreeExplainer uses an efficient polynomial time algorithm ( $O(TLD^2)$ ,  $T$  : number of trees,  $L$  :  
464 number of leaves,  $D$  : maximum depth) to approximate Shapley values which capture the impact  
465 of each feature on each individual model prediction. The Shapley value assigned to each feature  
466 is modeled as the average change in the model’s conditional expectation function over all possible  
467 feature orderings when introducing the corresponding feature into the model; it is formulated as  
468  $\mathbb{E}_S[f(X)|do(X_S = x_S)]$ , where  $S$  is the set of features,  $X$  is a random variable for the feature to  
469 perturb, and `do` notation [62] reflects the causal feature perturbation formulation. See [19] for more  
470 details on the algorithm and its properties.

471 Applying TreeExplainer on the model outcome at each fold across the 10-fold cross-fitting pro-  
472 cedure, we obtained out-of-sample local explanations for all individual model predictions of primary  
473 cancer types. By combining local explanations of correct predictions for each cancer type, we charac-  
474 terized the cancer type in terms of the most important or predictive features based on their Shapley  
475 values, which provided insights into the somatic variants and clinical features most relevant to the  
476 classification of each cancer type.

## 477 Germline PRS-based validation on CUP tumor samples

478 To validate the OncoNPC predictions for CUP tumor samples (which do not otherwise have a ground  
479 truth), we utilized germline Polygenic Risk Scores (PRS) which were never available to OncoNPC  
480 for training. Germline imputation from the off-target tumor sequencing data was conducted as  
481 previously described in [63]. Using weights from external GWAS data, we imputed PRS for Non-  
482 Small Cell Lung Cancer (NSCLC), Invasive Breast Carcinoma (BRCA), Colorectal Adenocarcinoma  
483 (COADREAD), Diffuse Glioma (DIFG), Melanoma (MEL), Ovarian Epithelial Tumor (OVT), Renal  
484 Cell Carcinoma (RCC), and Prostate Adenocarcinoma (PRAD). Pearson correlation between the  
485 PRS from off-target tumor data versus matched germline SNP array was previously shown to be  
486 higher than 0.9 without observable outliers [63].

487 We hypothesized that germline PRS specific to the underlying primary cancer type of a CUP  
488 tumor sample would be enriched in a manner similar to how the PRS specific to CKP tumor sample  
489 with the same primary cancer type is enriched. To that end, given the set of 8 different cancer types  
490  $\mathcal{C}$  we have the imputed PRS available for, we first restricted the cohort of CUP tumor samples to  
491 those with OncoNPC predictions in  $\mathcal{C}$  ( $N_{CUP,\mathcal{C}} = 505$ ). Then, we obtained standardized germline  
492 PRS values for the chosen CUP tumor samples over all the cancer types in  $\mathcal{C}$ . Finally, we defined  
493  $\hat{\Delta}_{PRS}$  as the estimated mean difference between the PRS specific to the predicted primary cancer  
494 type  $C$  (i.e. concordant PRS;  $PRS_C$ ) and average of PRSs corresponding to the rest of the cancer

495 types (i.e. discordant PRS;  $\text{PRS}_D$ , where  $D \in \mathcal{C} \setminus C$ ) as follows

$$\hat{\Delta}_{\text{PRS}} = \hat{\mathbb{E}}[\text{PRS}_C - \hat{\mathbb{E}}_D[\text{PRS}_D|C]] = \frac{1}{N_{\text{CUP},\mathcal{C}}} \sum_i^{N_{\text{CUP},\mathcal{C}}} \left( \text{PRS}_{c_i} - \frac{1}{|\mathcal{C} \setminus c_i|} \sum_{d_i \in \mathcal{C} \setminus c_i} \text{PRS}_{d_i} \right) \quad (1)$$

496 . As a true positive reference, we repeated the above procedure for the CKP tumor samples.  
497 Finally, as a true negative null, we estimated  $\hat{\Delta}_{\text{PRS-random}}$ , where the concordant cancer type was  
498 randomly assigned. We then repeated the random assignment 100 times to obtain estimated mean  
499 and standard errors.

## 500 Survival function estimation

501 National Death Index (NDI) and in-house clinical records were available for 20,281 DFCI patients  
502 ( $n = 16,376$  for CKP and  $n = 838$  for CUP). A patient's lost to follow-up date was determined  
503 at either the last NDI update date (12/31/2020) or their corresponding last contact date from the  
504 in-house records, whichever date is later. A patient's death date was determined from the in-house  
505 records, or the NDI data if the patient was lost to follow-up.

## 506 CUP-metastatic CKP survival comparison

507 We estimated median survival times of patients across CUP - metastatic CKP pairs using the  
508 Kaplan-Meier estimator [64] to account for patients lost to follow-up. For the CUP cohort, we  
509 excluded patients with CUP that were lost to follow up at the time of tumor sequencing and those  
510 whose primary cancer types were predicted with low probability (see Supplementary Fig. S5). The  
511 resulting CUP cohort ( $n = 685$ ), was then restricted to OncoNPCcancer types with more than  
512 35 CUP patients. For the CKP metastatic cohort, we excluded patients lost to follow up at the  
513 tumor sequencing time in the same manner and chose patients with one of the known cancers,  
514 where either the biopsy was metastatic or the patient had an ICD-10 code indicative of secondary  
515 malignant neoplasms within a year prior to sequencing dates. A total of 521 and 5,937 patients were  
516 thus retained from the CUP cohort and metastatic CKP cohort, respectively: Non-Small Cell Lung  
517 Cancer (NSCLC;  $n_{\text{CUP}} = 200$ ,  $n_{\text{met-CKP}} = 1,559$ ), Pancreatic Adenocarcinoma (PAAD;  $n_{\text{CUP}} = 80$ ,  
518  $n_{\text{met-CKP}} = 357$ ), Invasive Breast Carcinoma (BRCA;  $n_{\text{CUP}} = 67$ ,  $n_{\text{met-CKP}} = 1,656$ ), Colorectal  
519 Adenocarcinoma (COADREAD;  $n_{\text{CUP}} = 54$ ,  $n_{\text{met-CKP}} = 1,198$ ), Head and Neck Squamous Cell  
520 Carcinoma (HNSCC;  $n_{\text{CUP}} = 44$ ,  $n_{\text{met-CKP}} = 216$ ), Esophagogastric Adenocarcinoma (EGC;  $n_{\text{CUP}}$   
521  $= 40$ ,  $n_{\text{met-CKP}} = 336$ ), and Ovarian Epithelial Tumor (OVT;  $n_{\text{CUP}} = 36$ ,  $n_{\text{met-CKP}} = 615$ ). Note  
522 that patients with CUP, whose predicted cancer type is Gastrointestinal Neuroendocrine Tumors  
523 (GINET;  $n_{\text{CUP}} = 39$ ,  $n_{\text{CKP}} = 118$ ), were excluded due to the fact that the estimated survival function  
524 for the CUP cohort never reached 50 percent.

## 525 OncoNPC-based risk stratification among patients with CUP

526 To identify OncoNPC CUP subtypes with significant prognostic differences, we estimated survival  
527 functions for 7 common OncoNPC subtypes with more than 35 CUP patients: NSCLC, PAAD,  
528 BRCA, HNSCC, EGC, GINET, and Pancreatic Neuroendocrine Tumor (PANET). Patients that

529 were lost to follow up at time of sequencing were again excluded, as were CUPs with an OncoNPC  
530 prediction probability lower than 0.5 (i.e., same criteria as the CUP - metastatic CKP survival com-  
531 parison analysis). We merged subtypes with similar morphology and estimated survival functions:  
532 PAAD and EGC; GINET and PANET. To statistically test survival differences between these 5  
533 groups, we utilized Chi-squared test with 4 degrees of freedom.

## 534 **Identifying prognostic somatic variants shared in CUP-metastatic CKP** 535 **pairs**

536 To identify prognostic somatic variants shared between CUP/metastatic-CKP pairs, we again re-  
537 stricted to the 7 common OncoNPC subtypes with at least 35 CUP patients: NSCLC, PAAD,  
538 BRCA, COADREAD, HNSCC, EGC, GINET, and OVT. For somatic variants, we utilized the same  
539 processed features utilized in the OncoNPC model training (see Methods: Feature selection and  
540 OncoNPC model interpretation). To ensure sufficient statistical power, we restricted to candidate  
541 somatic variants (i.e., mutated genes and CNA genes) present in at least 15 samples in a given On-  
542 coNPC subtype and corresponding metastatic CKP cohort, as well as all 96 mutational signatures.

543 After selecting the cancer types to consider in the CUP-metastatic CKP pairs and candidate  
544 somatic variants for each pair, we iteratively tested each feature for association with survival in  
545 each OncoNPC subtype and in each corresponding metastatic CKP cohort. A multivariable Cox  
546 Proportional Hazard regression [41] model was used with time-to-death from sequencing as the  
547 outcome. To adjust for baseline effects, we included age at sequencing, sex, tumor sequencing panel  
548 version, mutational burden (i.e., sum of total somatic mutations in each tumor sample), and CNA  
549 burden (i.e., sum of total CNA events in each tumor sample) as covariates. Finally, to identify  
550 shared prognostic somatic variants for each CUP-metastatic CKP pair, we retained somatic variants  
551 which passed Schoenfeld residuals-based proportional hazard tests (`lifelines` Python library [65]:  
552 p-value threshold: 0.05) and were nominally significant ( $p < 0.05$ ) for both CUP and CKP cancer  
553 types in each pair.

## 554 **Actionable somatic variants in CUP tumors**

555 We estimated the frequency of known, actionable somatic alterations in each OncoNPC CUP subtype  
556 using the OncoKB knowledge base [38]. OncoNPC CUP predictions with a probability greater than  
557 0.5 were retained (see Supplementary Fig. S5). We considered 3 different types for somatic variants:  
558 oncogenic mutations such as indels, missense mutations, and splice site mutations, amplifications  
559 such as high-level amplifications, and finally fusions such as gene-gene and gene-intergenic fusions as  
560 specified in OncoKB. For each actionable somatic variant, we assigned one of the four therapeutic  
561 levels: level 1 for FDA-approved drugs, level 2 for standard care drugs, level 3 for drugs supported  
562 by clinical evidence, and level 4 for drugs supported by biological evidence.

## 563 Estimating impacts of treatment concordance on survival of patients with 564 CUP

565 We estimated the impact of the concordance between treatment and OncoNPC CUP predictions on  
566 a mortality outcome in a retrospective survival analysis. We utilized the in-house patient follow-up  
567 and treatment data to identify patients with CUP who received first treatment at DFCI with a  
568 palliative intent (Supplementary Fig. S5 for the exclusion criteria). Each patient was reviewed by  
569 a trained oncologist to determine whether the OncoNPC predicted cancer type was concordant or  
570 discordant with the first line of treatment received, per National Comprehensive Cancer Network  
571 (NCCN) guidelines or standard of care, in most reasonable situations, and within the clinical context  
572 delineated in the medical record. See Supplementary Section: *Determining treatment-OncoNPC*  
573 *concordance* for more details, and Supplementary Table S3 for clinical information, including primary  
574 cancer diagnosis, biopsy site, and first chemotherapy plan at DFCI, of patients with CUP in the  
575 analysis.

576 As we were interested in the counterfactual causal impact of the OncoNPC-treatment concor-  
577 dance, we utilized the principles of causal inference to account for potential patient heterogeneity  
578 and confounding. Specifically, we estimated the effect of treatment concordance specified by the  
579 indicator variable,  $A$ , which was 1 when the first palliative treatment for a patient with CUP was  
580 concordant with the corresponding OncoNPC prediction and 0 otherwise. Our analyses make the  
581 following identifiability assumptions:

- 582 • Conditional ignorability :  $A_i \perp\!\!\!\perp T_i^{a_i} | X_i$ , where  $A_i \in \{0, 1\}$ . It means that given patient  $i$ 's a set of  
583 covariates  $X_i$ , the patient's treatment concordance  $A_i$  is as good as random.
- 584 • Consistency :  $T_i^{a_i} = T_i$ , which means that a counterfactual outcome  $T_i^{a_i}$  for patient  $i$  is the  
585 observed outcome for the patient with a treatment concordance  $a_i$ .
- 586 • Overlap :  $P(0 < p(X_i) < 1) = 1$  where  $p(X_i) = P(A_i = 1 | X_i)$ , which means all patients have a  
587 strictly positive probability for receiving concordant treatment ( $A_i = 1$ ).

588 In addition to the above identifiability assumptions, we made independent censoring (i.e.  $C_i \perp\!\!\!\perp T_i | X_i$ )  
589 and independent entry assumption given the covariates (i.e.  $E_i \perp\!\!\!\perp T_i | X_i$ ).

590 We adopted two different estimation strategies to obtain the impact of treatment concordance:  
591 semi-parametric Cox Proportional Hazard estimator adjusted with a set of measured confounders  
592  $X$  [41] and non-parametric Kaplan Meier estimator adjusted with Inverse Probability Treatment  
593 Weighting (IPTW). We formulated an IPTW,  $w_i$  for each sample as  $w_i = \frac{P(A=a_i)}{P(A_i=a_i|X_i)}$  [42] and  
594 estimated  $P(A)$  non-parametrically and  $P(A|X)$  using a logistic regression model (R `glm` package  
595 [66]) in a 10-fold cross-fitting. A set of measured confounders (i.e.,  $X_i$ ) included patients' sex,  
596 age, OncoNPC prediction uncertainty (in entropy of posterior distribution over 22 cancer types),  
597 sequencing panel (i.e., OncoPanel) version, mutational burden, CNA burden, subsets of OncoNPC  
598 predicted cancer types and metastasis sites, and finally pathological histology (e.g., adenocarcinoma  
599 tumor or neuroendocrine tumor). Since patients with CUP who met the treatment criteria (i.e.,  
600 follow-up start time) but did not receive clinical panel sequencing (i.e., entry time) could not be



601 included in the analysis, we adjusted for the left truncation by defining the risk set  $\mathcal{R}(t)$  at time  $t$ ,  
602 which corresponds to the set of patients followed up in the analysis up to time  $t$  as follows

$$\mathcal{R}(t) = \{i | E_i \leq t \leq T_i\}$$

603 , where  $E_i$  is the entry time of patient  $i$ . With the independent entry assumption as stated before,  
604 we obtained survival function from Kaplan-Meier estimator as follows

$$\hat{S}(t) = \prod_{i:T_i \leq t} \left( 1 - \frac{\sum_{k:T_k=T_i} w_k}{\sum_{j:j \in \mathcal{R}(T_i)} w_j} \right)$$

605 . In this formulation, each individual is weighted by the corresponding IPTW,  $w_i$ , and we obtained  
606 two different survival functions for the treatment concordant and discordant groups. The adjusted  
607 Kaplan-Meier estimator provides a consistent estimate of impact of the treatment concordance under  
608 the assumptions stated above [42]. Once we obtained the survival estimates for the two groups, we  
609 used a weighted log-rank test [67] to test for a significant difference in survival.

610 In the Cox proportional hazard regression framework, we estimated the hazard function of patient  
611  $i$  as follows:  $\lambda(t|A_i, X_i) = \lambda_0(t) \exp(\alpha A_i + \beta^T X_i)$ , where  $\alpha, A_i \in \mathbb{R}$  and  $\beta, X_i \in \mathbb{R}^m$  ( $m$  is the number  
612 of measured confounders). Under the above identifiability assumptions and validity of the estimation  
613 model,  $e^\alpha$  is the hazard ratio capturing the causal effect of the treatment concordance  $A$ . Finally,  
614 under the assumption of no ties between event times across the patients, the parameters  $\alpha$  and  $\beta$   
615 are estimated by maximizing the following partial likelihood

$$L(\alpha, \beta) = \prod_{i:\delta_i=1} \frac{\exp(\alpha A_i + \beta X_i)}{\sum_{j:j \in \mathcal{R}(T_i)} \exp(\alpha A_j + \beta X_j)}$$

616 [41].

## 617 Acknowledgments

618 The participation of patients and the efforts of an institutional data collection system made this study  
619 possible, and we are grateful for their contributions. We would also like to express our appreciation  
620 to the DFCI Oncology Data Retrieval System (OncDRS) and AACR Project GENIE team for their  
621 role in aggregating, managing, and delivering the data used in this project.

## 622 Funding

623 IM and AG were supported by R01 CA227237, R01 CA244569, as well as grants from The Louis  
624 B. Mayer Foundation, The Doris Duke Charitable Foundation, The Phi Beta Psi Sorority, and The  
625 Emerson Collective.

## 626 **Supplementary Note**

### 627 **Determining treatment-OncoNPC concordance**

628 Concordance of OncoNPC predicted cancer type with a first palliative treatment assignments at  
629 DFCI was classified in one of five categories: 1) “TRUE”: the OncoNPC cancer type matched  
630 the clinically proven/suspected tumor type and the predicted treatment matched the treatment  
631 received, which was dictated by NCCN guidelines and/or standard of care, within the clinical context  
632 provided by the medical record; 2) “FALSE”: the OncoNPC cancer type did not match the clinically  
633 proven/suspected cancer type and the predicted treatment was not appropriate per NCCN guidelines  
634 or standard of care, in most reasonable situations, and within the context of the medical record; 3)  
635 “SOFT FALSE”: the OncoNPC cancer type did not match the clinically proven/suspected cancer  
636 type, but the treatment received was not chosen based on NCCN guidelines or standard of care, owing  
637 to the unique clinical context provided by the medical record, 4) “EMPIRIC”: treatment received was  
638 empiric treatment for cancer of unknown primary (e.g., carboplatin/taxol or gemcitabine/cisplatin)  
639 with the corresponding clinical rationale; in cases where patients received these regimens but not  
640 with the clinical intent of empiric CUP treatment (i.e., as regimens intended for treating other tumor  
641 types), the predicted treatment was not labeled as “EMPIRIC” and the case was instead evaluated  
642 in context of the proven/suspected tumor type. In our analysis, we considered the TRUE group  
643 as the concordant group, and FALSE and SOFT FALSE groups as the discordant group. We did  
644 not include the EMPIRIC group, which is typically a more challenging patient population with  
645 systematically worse outcomes [49].

### 646 **Code Availability**

647 Please see <https://github.com/itmoon7/onconpc> for the pre-processing script, the trained OncoNPC  
648 model, and other reference materials.

## 649 Figure and Table Legends

650 **Figure 1. Overview of model development and analysis workflow.** (a) OncoNPC, a  
651 XGBoost-based classifier, was trained and evaluated using 36,729 tumor samples across 22 cancer  
652 types from Cancers of Known Primary (CKP) collected from three different cancer centers. (b)  
653 OncoNPC performance was evaluated on the held-out tumor samples ( $n = 7,289$ ). (c) OncoNPC  
654 was applied to 971 CUPs at a single institution to predict primary cancer types. OncoNPC pre-  
655 dicted CUP subtypes were then investigated for association with: (d) elevated germline risk, (e)  
656 actionable molecular alterations, (f) overall survival, and (g) prognostic somatic features. (h) A  
657 subset of CUP patients with detailed treatment data were evaluated for treatment-specific outcomes.

658  
659 **Figure 2. Cancer type prediction performance of OncoNPC.** (a),(b) The normalized con-  
660 fusion matrix of OncoNPC classification performance on the held-out test set ( $n = 7,289$ ) for (a)  
661 22 detailed cancer types and (b) 10 broad cancer groups based on site and treatment (see Table 1).  
662 The sensitivity for each cancer type or cancer group is shown below each confusion matrix and the  
663 sample size is shown to the left of each confusion matrix. (c), (d) The performance (by F1 score)  
664 of OncoNPC on the test set across cancer types (c) and groups (d) at 4 different prediction confi-  
665 dences (i.e., minimum  $p_{max}$  thresholds). Each dot size is scaled by the proportion of tumor samples  
666 retained. (e) Precision-recall curves showing the performance of the OncoNPC across different co-  
667 horts in the test set by: cancer center, biopsy site type, sequence panel version, and ethnicity (color  
668 coded), with the yellow dotted curve corresponding to the baseline performance on the full test set.

669  
670 **Figure 3. Applying OncoNPC to CUP tumor samples and interpreting cancer type pre-  
671 dictions.** (a) Empirical distributions of prediction probabilities for correctly predicted, held-out  
672 CKP tumor samples ( $n = 3,429$ ) and CUP tumor samples ( $n = 934$ ) across CKP cancer types (blue)  
673 and their corresponding OncoNPC predicted cancer types for CUP tumors (green). Only OncoNPC  
674 classifications with at least 20 CUP tumor samples are shown. (b) Proportion of each CKP cancer  
675 type and the corresponding OncoNPC predicted CUP cancer type. All training CKP tumor samples  
676 ( $n = 36,445$ ) and all held-out CUP tumor samples ( $n = 971$ ) are shown. For both (a) and (b), the  
677 cancer types (x-axis) are ordered by the number of CKP tumor samples in each cancer type. (c)  
678 Germline Polygenic Risk Score (PRS) enrichment of the CKP tumor samples ( $n = 11,332$ ) and CUP  
679 tumor samples with available PRS data ( $n = 505$ ) averaged across 8 cancer types. The magnitude of  
680 the enrichment is quantified by  $\hat{\Delta}_{PRS}$ : the mean difference between the concordant (i.e. OncoNPC  
681 matching) cancer type PRS and mean of PRSs of discordant cancer types (see Methods).  $\hat{\Delta}_{PRS}$   
682 is shown for CKPs in blue (for reference) and CUPs in green. As a negative control,  $\hat{\Delta}_{PRS-random}$   
683 is also shown after permuting the OncoNPC labels. (d) Top 15 most important features based on  
684 mean absolute SHAP values (i.e.,  $\hat{\mu}(|SHAP|)$  [19]) for the top 3 most frequent cancer types in the  
685 cohort: Non-Small Cell Lung Cancer (NSCLC), Invasive Breast Carcinoma (BRCA), and Pancre-  
686 atic Adenocarcinoma (PAAD). The carrier rate for each feature in corresponding CKP and CUP  
687 cancer cohorts as well as the entire CKP and CUP cohorts are shown as bars going downwards  
688 and star-shaped markers, respectively. For mutation signature features that have continuous values,

689 individuals with feature values one standard deviation above the mean were treated as positives and  
690 the rest as negative. For age, individuals above the population mean were treated as positives and  
691 the rest as negatives. (e) Explanation of OncoNPC cancer type prediction for a sample patient with  
692 CUP. The patient is a 76 year-old male, with a tumor biopsy from the liver. The pie chart on the  
693 left shows the Top 10 important features across three different feature categories (i.e., CNA events,  
694 somatic mutation, and mutation signatures), and the scatter plot on the right shows their SHAP  
695 values and feature values. The size of each dot is scaled by corresponding absolute SHAP value.

696  
697 **Figure 4. Consistent survival and prognostic biomarkers between OncoNPC classifi-**  
698 **cations and known cancers.** (a) Survival stratification for patients with CUP based on their  
699 OncoNPC predicted cancer types. The Kaplan-Meier estimator [64] was used to estimate survival  
700 probability for each predicted cancer type over the follow-up time of 60 months from sequence date,  
701 with statistical significance assessed by Chi-square test. (b) Correspondence between median sur-  
702 vival time (in months) of CUP predicted cancer types (x-axis) and those of metastatic CKP cancer  
703 types (y-axis): Spearman’s rho 0.964 (p-value:  $4.54 \times 10^{-4}$ ). The size of each dot reflects the p-value  
704 of log-rank test for significant difference in median survival between CUP - metastatic CKP pairs.  
705 Only cancer types with at least 30 CUP tumor samples having OncoNPC probabilities greater than  
706 0.5 are shown. (c), (d) Prognostic somatic variants significantly associated with overall survival,  
707 shared between three different CUP (c)-metastatic CKP (d) pairs (NSCLC, PAAD, and COAD-  
708 READ; indicated by point shape). Variant types are indicated by colors: red for somatic mutations,  
709 green for CNAs, and blue for mutation signatures.

710  
711 **Figure 5. Potential for clinical decision support among OncoNPC classified CUPs.** (a)  
712 The number of CUP tumor samples with actionable targets, based on OncoKB [38], across actionable  
713 somatic variants (mutations, amplifications, and fusions). Each bar corresponds to an actionable  
714 target, color-coded by the number of each OncoNPC classified CUP carrier. Note that each tumor  
715 sample may contain more than one actionable somatic variant. (b) Proportions of CUP tumor  
716 samples with actionable somatic variants ( $N_{action}$ ) to the total number of patients ( $N_{total}$ ) across  
717 OncoNPC predicted cancer types. Proportions for 4 different therapeutic levels based on OncoKB  
718 [38], are shown in each bar: Level 1 - FDA-approved drugs, Level 2 - standard of care drugs, Level  
719 3 - drugs supported by clinical evidence, and Level 4 - drugs supported by biological evidence. (c),  
720 (d) Treatment diagrams for a group of patients with CUP who received treatments that were concor-  
721 dant with the OncoNPC classification (c) and the remaining CUP patients who received discordant  
722 treatments (d). OncoNPC classification is shown on the left and treatment groups are shown on  
723 the right, with each patient connected from left to right. (e) Forest plot of a multivariable Cox  
724 Proportional Hazards Regression on patients in the CUP cohort with first-line palliative treatment  
725 records at DFCI (n = 159; see Appendix Fig. S5 for the exclusion criteria). Treatment concordance  
726 (colored in blue), encoded as 1 when the first treatment a patient received at DFCI is *concordant*  
727 with their corresponding OncoNPC prediction and 0 otherwise, was significantly associated with  
728 mortality of patients in the cohort (H.R. 0.321, 95% C.I. 0.165 - 0.620, p-value:  $< 0.001$ ). (f)  
729 Estimated survival curves for patients with CUP in the concordant treatment group (shown in blue)

730 and discordant treatment group (shown in red), respectively. To estimate the survival function for  
731 each group, we utilized Inverse Probability of Treatment Weighted (IPTW) Kaplan-Meier estimator  
732 while adjusting for left truncation until time of sequencing (see Methods). Statistical significance of  
733 the survival difference between the two groups was estimated by a weighted log-rank test [68].

734

735 **Table 1.** Demographic details of the patients and tumor samples across DFCI, MSK, and VICC.

736

737 **Table 2.** Demographic details of patients with CUP in the concordant and discordant treatment  
738 groups.

739

740 **Supplementary Figure S1. OncoNPC prediction performances and confidences (i.e.,**  
741  $p_{\max}$ **) across centers. (a)** Center-specific OncoNPC performance (in weighted F1) on the test  
742 CKP tumor samples ( $n = 7,289$ ). The figure is a decomposed version of Fig. 2c by cancer center  
743 (DFCI:  $\circ$ , MSK:  $\square$ , VICC:  $\diamond$ ). The performance was evaluated at 4 different prediction confidences  
744 (i.e., minimum  $p_{\max}$  thresholds). Each dot size is scaled by the proportion of tumor samples re-  
745 tained. See Table S2 for the center-specific number of test CKP tumor samples across cancer types.  
746 **(b), (c)** Box plots of prediction confidences ( $p_{\max}$ ) across **(b)** DFCI CUP tumors, MSK CUP tu-  
747 mors, all DFCI CKP tumors, DFCI held-out CKP tumors, and DFCI excluded CKP tumors, and  
748 **(c)** DFCI held-out CKP tumors, MSK held-out CKP tumors, and VICC held-out CKP tumors.  
749 Medians and lower and upper quartiles are shown on the figures along with corresponding number  
750 of tumor samples as well as means and 95% confidence intervals.

751

752 **Supplementary Figure S2. Interpreting OncoNPC predictions.** Top 15 most important  
753 features based on mean absolute SHAP values (i.e.,  $\hat{\mu}(|\text{SHAP}|)$  [19]) for cancer types with at least  
754 20 CUP tumors samples were classified into.

755

756 **Supplementary Figure S3. SHAP summary plot** [19] for cancer types with at least 20 CUP  
757 tumors samples were classified into. SHAP values (i.e., impact on OncoNPC predictions) are shown  
758 on the x-axis, while features values are shown with a color map (from purple to yellow). In each  
759 plot, CUP and CKP tumor samples were combined into one cohort for the corresponding cancer.

760

761 **Supplementary Figure S4. Applying OncoNPC to MSK CUP tumor samples. (a)** Em-  
762 pirical distributions of prediction probabilities for correctly predicted, held-out CKP tumor samples  
763 ( $n = 3,429$ ) and MSK CUP tumor samples ( $n = 496$ ) across CKP cancer types (blue) and their  
764 corresponding OncoNPC predicted cancer types for CUP tumors (green). Only OncoNPC classifi-  
765 cations with at least 20 CUP tumor samples are shown. **(b)** Proportion of each CKP cancer type  
766 and the corresponding OncoNPC predicted CUP cancer type. All training CKP tumor samples ( $n$   
767  $= 36,445$ ) and all MSK CUP tumor samples ( $n = 581$ ) are shown. For both **(a)** and **(b)**, the cancer  
768 types (x-axis) are ordered by the number of CKP tumor samples in each cancer type.

769

770 **Supplementary Figure S5.** Exclusion criteria for downstream clinical analyses.

771

772 **Supplementary Figure S6.** Germline Polygenic Risk Score (PRS) enrichment of CKP tumor  
773 samples and CUP tumor samples across 8 different cancer types: **(a)** Colorectal Adenocarcinoma  
774 (COADREAD), **(b)** Diffuse Glioma (DIFG), **(c)** Invasive Breast Carcinoma (BRCA), **(d)** Melanoma  
775 (MEL), **(e)** Non-Small Cell Lung Cancer (NSCLC), **(f)** Ovarian Epithelial Tumor (OVT), **(g)**  
776 Prostate Adenocarcinoma (PRAD), and **(h)** Renal Cell Carcinoma (RCC). The magnitude of the  
777 enrichment is quantified by  $\hat{\Delta}_{\text{PRS}}$ : the mean difference between the concordant (i.e. OncoNPC  
778 matching) cancer type PRS and mean of PRSs of discordant cancer types (see Methods).  $\hat{\Delta}_{\text{PRS}}$   
779 is shown for CKPs in blue (for reference) and CUPs in green.

780

781 **Supplementary Figure S7.** Estimated survival curves for patients with CUP, broken down by  
782 OncoNPC predicted cancer types: **(a)** BRCA, **(b)** Gastrointestinal (GI) group (CHOL, COAD-  
783 READ, EGC, and PAAD), **(c)** Lung (NSCLC and PLMESO), and **(d)** other OncoNPC cancer  
784 types (BLCA, DIFG, GINET, HNSCC, MEL, OVT, PANET, PRAD, RCC, and UCEC). In each  
785 figure, the concordant treatment group and discordant treatment group are shown in blue and red,  
786 respectively. To estimate the survival function for each group, we utilized Inverse Probability of  
787 Treatment Weighted (IPTW) Kaplan-Meier estimator while adjusting for left truncation until time  
788 of sequencing (see Methods). Statistical significance of the survival difference between the two groups  
789 was estimated by a weighted log-rank test [68].

790

791 **Supplementary Figure S8. Summary of coefficients for estimating treatment-OncoNPC**  
792 **concordance.** Formally, we estimated out-of-sample  $P(A|X)$ , where  $A$  corresponds to the treatment-  
793 OncoNPC concordance, using a logistic regression model in a 10-fold cross-fitting. The coefficients  
794 were obtained from the first fold. See Methods: Estimating impacts of treatment concordance on  
795 survival of patients with CUP for more details.

796

797 **Supplementary Table S1.** A full set of 861 somatic input features for OncoNPC, abstracted from  
798 the next generation NGS targeted panel sequencing data. The features belong to three different  
799 categories (shown as columns of the table): somatic mutations (i.e., single nucleotide variants and  
800 indels: 316 features), Copy Number Alterations (CNA: 491 features), and mutational signatures  
801 (54 features). Note that we included patients' sex and age in addition to the somatic features. See  
802 Methods for more details on how the features were encoded.

803

804 **Supplementary Table S2.** Center-specific number of held-out CKP tumor samples across cancer  
805 types and prediction confidence (i.e.,  $p_{\text{max}}$ ) thresholds.

806

807 **Supplementary Table S3.** Clinical information of patient with CUP in the treatment concordance  
808 analysis ( $n = 158$ ).

809

## References

810

- 811 [1] A. P. G. Consortium *et al.*, “Aacr project genie: Powering precision medicine through an  
812 international consortium,” *Cancer discovery*, vol. 7, no. 8, pp. 818–831, 2017.
- 813 [2] N. Pavlidis, H. Khaled, and R. Gaafar, “A mini review on cancer of unknown primary site: A  
814 clinical puzzle for the oncologists,” *Journal of advanced research*, vol. 6, no. 3, pp. 375–382,  
815 2015.
- 816 [3] G. R. Varadhachary and M. N. Raber, “Cancer of unknown primary site,” *New England Journal*  
817 *of Medicine*, vol. 371, no. 8, pp. 757–765, 2014.
- 818 [4] D. M. Hyman *et al.*, “Vemurafenib in multiple nonmelanoma cancers with braf v600 muta-  
819 tions,” *New England Journal of Medicine*, vol. 373, no. 8, pp. 726–736, 2015.
- 820 [5] J. D. Hainsworth and F. A. Greco, “Cancer of unknown primary site: New treatment paradigms  
821 in the era of precision medicine,” *American Society of Clinical Oncology Educational Book*,  
822 vol. 38, pp. 20–25, 2018.
- 823 [6] G. G. Anderson and L. M. Weiss, “Determining tissue of origin for metastatic cancers: Meta-  
824 analysis and literature review of immunohistochemistry performance,” *Applied Immunohisto-*  
825 *chemistry & Molecular Morphology*, vol. 18, no. 1, pp. 3–8, 2010.
- 826 [7] K. A. Oien and J. L. Dennis, “Diagnostic work-up of carcinoma of unknown primary: from  
827 immunohistochemistry to molecular profiling,” *Ann Oncol*, vol. 23 Suppl 10, pp. x271–277,  
828 Sep. 2012.
- 829 [8] S. Moran *et al.*, “Epigenetic profiling to classify cancer of unknown primary: A multicentre,  
830 retrospective analysis,” *The Lancet Oncology*, vol. 17, no. 10, pp. 1386–1395, 2016.
- 831 [9] W. Jiao *et al.*, “A deep learning system can accurately classify primary and metastatic cancers  
832 based on patterns of passenger mutations,” *bioRxiv*, p. 214494, 2019.
- 833 [10] A. Penson *et al.*, “Development of genome-derived tumor type prediction to inform clinical  
834 cancer care,” *JAMA oncology*, vol. 6, no. 1, pp. 84–91, 2020.
- 835 [11] B. He *et al.*, “A neural network framework for predicting the tissue-of-origin of 15 common  
836 cancer types based on rna-seq data,” *Frontiers in Bioengineering and Biotechnology*, vol. 8,  
837 p. 737, 2020.
- 838 [12] L. Nguyen, A. Van Hoeck, and E. Cuppen, “Machine learning-based tissue of origin classifica-  
839 tion for cancer of unknown primary diagnostics using genome-wide mutation features,” *Nature*  
840 *communications*, vol. 13, 2022.
- 841 [13] J. D. Hainsworth *et al.*, “Molecular gene expression profiling to predict the tissue of origin and  
842 direct site-specific therapy in patients with carcinoma of unknown primary site: A prospective  
843 trial of the sarah cannon research institute,” *Journal of Clinical Oncology*, vol. 31, no. 2,  
844 pp. 217–223, 2013.
- 845 [14] H. Yoon *et al.*, “Gene expression profiling identifies responsive patients with cancer of unknown  
846 primary treated with carboplatin, paclitaxel, and everolimus: Ncctg n0871 (alliance),” *Annals*  
847 *of Oncology*, vol. 27, no. 2, pp. 339–344, 2016.

- 848 [15] H. Hayashi *et al.*, “Site-specific and targeted therapy based on molecular profiling by next-  
849 generation sequencing for cancer of unknown primary site: A nonrandomized phase 2 clinical  
850 trial,” *JAMA oncology*, vol. 6, no. 12, pp. 1931–1938, 2020.
- 851 [16] H. Hayashi *et al.*, “Randomized phase ii trial comparing site-specific treatment based on gene  
852 expression profiling with carboplatin and paclitaxel for patients with cancer of unknown pri-  
853 mary site,” *Journal of Clinical Oncology*, vol. 37, no. 7, pp. 570–579, 2019.
- 854 [17] A.-M. Conway, C. Mitchell, and N. Cook, “Challenge of the unknown: How can we improve  
855 clinical outcomes in cancer of unknown primary?” *Journal of clinical oncology: official journal  
856 of the American Society of Clinical Oncology*, vol. 37, no. 23, pp. 2089–2090, 2019.
- 857 [18] T. Bochtler and A. Krämer, “Does cancer of unknown primary (cup) truly exist as a distinct  
858 cancer entity?” *Frontiers in oncology*, vol. 9, p. 402, 2019.
- 859 [19] S. M. Lundberg *et al.*, “From local explanations to global understanding with explainable ai  
860 for trees,” *Nature machine intelligence*, vol. 2, no. 1, pp. 56–67, 2020.
- 861 [20] J. G. Tate *et al.*, “Cosmic: The catalogue of somatic mutations in cancer,” *Nucleic acids  
862 research*, vol. 47, no. D1, pp. D941–D947, 2019.
- 863 [21] G. da Cunha Santos, F. A. Shepherd, and M. S. Tsao, “Egfr mutations and lung cancer,”  
864 *Annual Review of Pathology: Mechanisms of Disease*, vol. 6, pp. 49–69, 2011.
- 865 [22] Y.-L. Zhang *et al.*, “The prevalence of egfr mutation in patients with non-small cell lung cancer:  
866 A systematic review and meta-analysis,” *Oncotarget*, vol. 7, no. 48, p. 78 985, 2016.
- 867 [23] S. S. Hecht, “Tobacco smoke carcinogens and lung cancer,” *JNCI: Journal of the National  
868 Cancer Institute*, vol. 91, no. 14, pp. 1194–1210, 1999.
- 869 [24] R. Mehra *et al.*, “Identification of gata3 as a breast cancer prognostic marker by global gene  
870 expression meta-analysis,” *Cancer research*, vol. 65, no. 24, pp. 11 259–11 264, 2005.
- 871 [25] S. Elsheikh *et al.*, “Ccnd1 amplification and cyclin d1 expression in breast cancer and their  
872 relation with proteomic subgroups and patient outcome,” *Breast cancer research and treatment*,  
873 vol. 109, no. 2, pp. 325–335, 2008.
- 874 [26] X. Liu, M. Jakubowski, and J. L. Hunt, “Kras gene mutation in colorectal cancer is corre-  
875 lated with increased proliferation and spontaneous apoptosis,” *American journal of clinical  
876 pathology*, vol. 135, no. 2, pp. 245–252, 2011.
- 877 [27] D. Dinu *et al.*, “Prognostic significance of kras gene mutations in colorectal cancer-preliminary  
878 study,” *Journal of medicine and life*, vol. 7, no. 4, p. 581, 2014.
- 879 [28] R. Hayes, J. Van Nieuwenhuize, J. Raatgever, and F. Ten Kate, “Aflatoxin exposures in the in-  
880 dustrial setting: An epidemiological study of mortality,” *Food and Chemical Toxicology*, vol. 22,  
881 no. 1, pp. 39–43, 1984.
- 882 [29] M. A. Ahmed Adam, Y. M. Tabana, K. B. Musa, and D. A. Sandai, “Effects of different  
883 mycotoxins on humans, cell genome and their involvement in cancer,” *Oncology Reports*, vol. 37,  
884 no. 3, pp. 1321–1336, 2017.



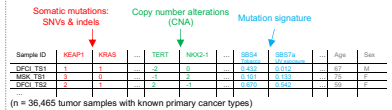
- 885 [30] S. Marchese, A. Polo, A. Ariano, S. Velotto, S. Costantini, and L. Severino, “Aflatoxin b1 and  
886 m1: Biological properties and their involvement in cancer development,” *Toxins*, vol. 10, no. 6,  
887 p. 214, 2018.
- 888 [31] P. M. Westcott and M. D. To, “The genetics and biology of kras in lung cancer,” *Chinese*  
889 *journal of cancer*, vol. 32, no. 2, p. 63, 2013.
- 890 [32] H. Adderley, F. H. Blackhall, and C. R. Lindsay, “Kras-mutant non-small cell lung cancer: Con-  
891 verging small molecules and immune checkpoint inhibition,” *EBioMedicine*, vol. 41, pp. 711–  
892 716, 2019.
- 893 [33] A. M. Conway, C. Mitchell, E. Kilgour, G. Brady, C. Dive, and N. Cook, “Br J CancerMolecular  
894 characterisation and liquid biomarkers in Carcinoma of Unknown Primary (CUP): taking the  
895 ‘U’ out of ‘CUP’,” *Br J Cancer*, vol. 120, no. 2, pp. 141–153, Jan. 2019.
- 896 [34] A. J. Schoenfeld *et al.*, “The genomic landscape of smarca4 alterations and associations with  
897 outcomes in patients with lung cancersmarca4 alterations in lung cancer,” *Clinical Cancer*  
898 *Research*, vol. 26, no. 21, pp. 5701–5708, 2020.
- 899 [35] S. Papillon-Cavanagh, P. Doshi, R. Dobrin, J. Szustakowski, and A. M. Walsh, “Stk11 and  
900 keap1 mutations as prognostic biomarkers in an observational real-world lung adenocarcinoma  
901 cohort,” *ESMO open*, vol. 5, no. 2, e000706, 2020.
- 902 [36] T. Takahashi *et al.*, “Mutations in keap1 are a potential prognostic factor in resected non-small  
903 cell lung cancer,” *Journal of surgical oncology*, vol. 101, no. 6, pp. 500–506, 2010.
- 904 [37] T. Berghmans *et al.*, “Thyroid transcription factor 1—a new prognostic factor in lung cancer:  
905 A meta-analysis,” *Annals of oncology*, vol. 17, no. 11, pp. 1673–1676, 2006.
- 906 [38] D. Chakravarty *et al.*, “Oncokb: A precision oncology knowledge base,” *JCO precision oncology*,  
907 vol. 1, pp. 1–16, 2017.
- 908 [39] P. M. Grambsch and T. M. Therneau, “Proportional hazards tests and diagnostics based on  
909 weighted residuals,” *Biometrika*, vol. 81, no. 3, pp. 515–526, 1994.
- 910 [40] L. Simms, H. Barraclough, and R. Govindan, “Biostatistics primer: What a clinician ought to  
911 know—prognostic and predictive factors,” *Journal of Thoracic Oncology*, vol. 8, no. 6, pp. 808–  
912 813, 2013.
- 913 [41] D. R. Cox, “Regression models and life-tables,” *Journal of the Royal Statistical Society: Series*  
914 *B (Methodological)*, vol. 34, no. 2, pp. 187–202, 1972.
- 915 [42] J. Xie and C. Liu, “Adjusted kaplan–meier estimator and log-rank test with inverse probability  
916 of treatment weighting for survival data,” *Statistics in medicine*, vol. 24, no. 20, pp. 3089–3110,  
917 2005.
- 918 [43] R. Liu *et al.*, “Systematic pan-cancer analysis of mutation–treatment interactions using large  
919 real-world clinicogenomics data,” *Nature Medicine*, vol. 28, no. 8, pp. 1656–1661, 2022.
- 920 [44] R. Liu *et al.*, “Evaluating eligibility criteria of oncology trials using real-world data and ai,”  
921 *Nature*, vol. 592, no. 7855, pp. 629–633, 2021.

- 922 [45] S. Kolling *et al.*, ““metastatic cancer of unknown primary” or “primary metastatic cancer?””  
923 *Frontiers in Oncology*, vol. 9, p. 1546, 2020.
- 924 [46] T. Olivier *et al.*, “Redefining cancer of unknown primary: Is precision medicine really shifting  
925 the paradigm?” *Cancer treatment reviews*, vol. 97, p. 102 204, 2021.
- 926 [47] E. Moiso *et al.*, “Developmental deconvolution for classification of cancer origin,” *medRxiv*,  
927 2021.
- 928 [48] M. Y. Lu *et al.*, “Ai-based pathology predicts origins for cancers of unknown primary,” *Nature*,  
929 vol. 594, no. 7861, pp. 106–110, 2021.
- 930 [49] K. Fizazi, F. Greco, N. Pavlidis, G. Daugaard, K. Oien, and G. Pentheroudakis, “Cancers of  
931 unknown primary site: Esmo clinical practice guidelines for diagnosis, treatment and follow-  
932 up,” *Annals of Oncology*, vol. 26, pp. v133–v138, 2015.
- 933 [50] T. L. Stockley *et al.*, “Molecular profiling of advanced solid tumors and patient outcomes  
934 with genotype-matched clinical trials: The princess margaret impact/compact trial,” *Genome  
935 medicine*, vol. 8, no. 1, pp. 1–12, 2016.
- 936 [51] D. T. Cheng *et al.*, “Memorial sloan kettering-integrated mutation profiling of actionable cancer  
937 targets (msk-impact): A hybridization capture-based next-generation sequencing clinical assay  
938 for solid tumor molecular oncology,” *The Journal of molecular diagnostics*, vol. 17, no. 3,  
939 pp. 251–264, 2015.
- 940 [52] E. P. Garcia *et al.*, “Validation of oncopanel: A targeted next-generation sequencing assay for  
941 the detection of somatic variants in cancer,” *Archives of Pathology and Laboratory Medicine*,  
942 vol. 141, no. 6, pp. 751–758, 2017.
- 943 [53] D. Sha, Z. Jin, J. Budczies, K. Kluck, A. Stenzinger, and F. A. Sinicrope, “Tumor mutational  
944 burden as a predictive biomarker in solid tumors,” *Cancer discovery*, vol. 10, no. 12, pp. 1808–  
945 1825, 2020.
- 946 [54] L. Mileshekin *et al.*, “Cancer-of-unknown-primary-origin: A seer–medicare study of patterns  
947 of care and outcomes among elderly patients in clinical practice,” *Cancers*, vol. 14, no. 12,  
948 p. 2905, 2022.
- 949 [55] T. Chen and C. Guestrin, “Xgboost: A scalable tree boosting system,” in *Proceedings of the  
950 22nd acm sigkdd international conference on knowledge discovery and data mining*, 2016,  
951 pp. 785–794.
- 952 [56] Y. Chen *et al.*, “Physiol MeasClassification of short single-lead electrocardiograms (ECGs) for  
953 atrial fibrillation detection using piecewise linear spline and XGBoost,” *Physiol Meas*, vol. 39,  
954 no. 10, p. 104 006, Oct. 2018.
- 955 [57] C. M. Hatton, L. W. Paton, D. McMillan, J. Cussens, S. Gilbody, and P. A. Tiffin, “Predicting  
956 persistent depressive symptoms in older adults: A machine learning approach to personalised  
957 mental healthcare,” *Journal of affective disorders*, vol. 246, pp. 857–860, 2019.
- 958 [58] A. Ogunleye and Q.-G. Wang, “Xgboost model for chronic kidney disease diagnosis,” *IEEE/ACM  
959 transactions on computational biology and bioinformatics*, vol. 17, no. 6, pp. 2131–2140, 2019.

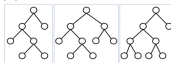
- 960 [59] J. Bergstra and Y. Bengio, “Random search for hyper-parameter optimization.,” *Journal of*  
961 *machine learning research*, vol. 13, no. 2, 2012.
- 962 [60] L. B. Alexandrov *et al.*, “NatureThe repertoire of mutational signatures in human cancer,”  
963 *Nature*, vol. 578, no. 7793, pp. 94–101, Feb. 2020.
- 964 [61] R. Rosenthal, N. McGranahan, J. Herrero, B. S. Taylor, and C. Swanton, “Genome BiolD-  
965 econstructSigs: delineating mutational processes in single tumors distinguishes DNA repair  
966 deficiencies and patterns of carcinoma evolution,” *Genome Biol*, vol. 17, p. 31, Feb. 2016.
- 967 [62] D. Janzing, L. Minorics, and P. Blöbaum, “Feature relevance quantification in explainable ai:  
968 A causal problem,” in *International Conference on artificial intelligence and statistics*, PMLR,  
969 2020, pp. 2907–2916.
- 970 [63] A. Gusev, S. Groha, K. Taraszka, Y. R. Semenov, and N. Zaitlen, “Constructing germline  
971 research cohorts from the discarded reads of clinical tumor sequences,” *Genome medicine*,  
972 vol. 13, no. 1, pp. 1–14, 2021.
- 973 [64] E. L. Kaplan and P. Meier, “Nonparametric estimation from incomplete observations,” *Journal*  
974 *of the American statistical association*, vol. 53, no. 282, pp. 457–481, 1958.
- 975 [65] C. Davidson-Pilon, “Lifelines: Survival analysis in python,” *Journal of Open Source Software*,  
976 vol. 4, no. 40, p. 1317, 2019.
- 977 [66] I. Marschner, M. W. Donoghoe, and M. M. W. Donoghoe, “Package ‘glm2’,” *Journal, Vol*,  
978 vol. 3, no. 2, pp. 12–15, 2018.
- 979 [67] A. Pezzi, M. Cavo, A. Biggeri, E. Zamagni, and O. Nanni, “Inverse probability weighting to  
980 estimate causal effect of a singular phase in a multiphase randomized clinical trial for multiple  
981 myeloma,” *BMC medical research methodology*, vol. 16, no. 1, pp. 1–10, 2016.
- 982 [68] S. Xu *et al.*, “Extension of kaplan-meier methods in observational studies with time-varying  
983 treatment,” *Value in Health*, vol. 15, no. 1, pp. 167–174, 2012.

Targeted clinical NGS assays : DFCI OncoPanel, MSK IMPACT, and VICC Panel

Somatic variants pre-processing



(a)



OncoNPC: XGBoost-based primary cancer type classifier

Posterior probabilities across 22 cancer types

Sample ID	NSCLC	PRAD	BRCA	BLCA
DFCI TS1	0.92	0.00	0.04	0.00
MSK TS1	0.11	0.85	0.00	0.01
DFCI TS2	0.01	0.01	0.96	0.00

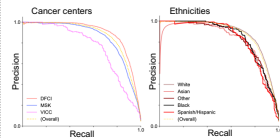
OncoNPC predictions for held-out CKP tumors

OncoNPC predictions for CUP tumors

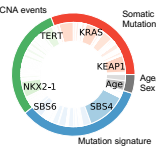
(c)

Model evaluation and interpretation

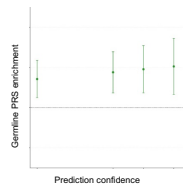
(b) Model performance evaluation



Model interpretation



(d) Germline PRS validation for CUP tumor samples



Clinical utility of OncoNPC classifications for patients with CUP

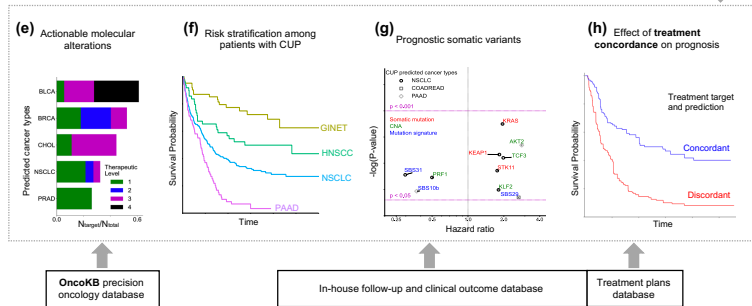
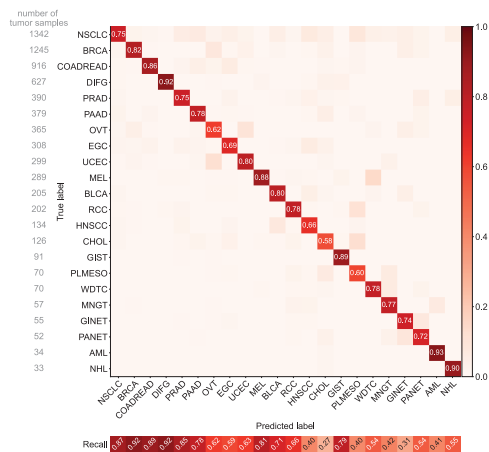


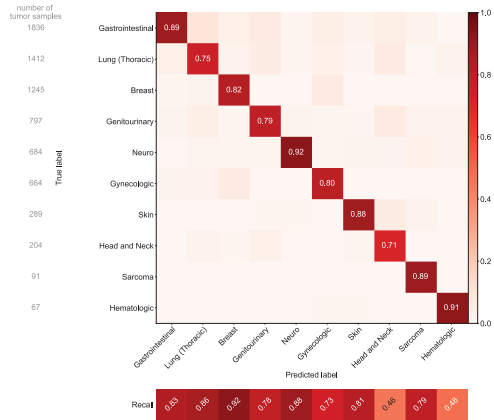
Figure 1

Table 1

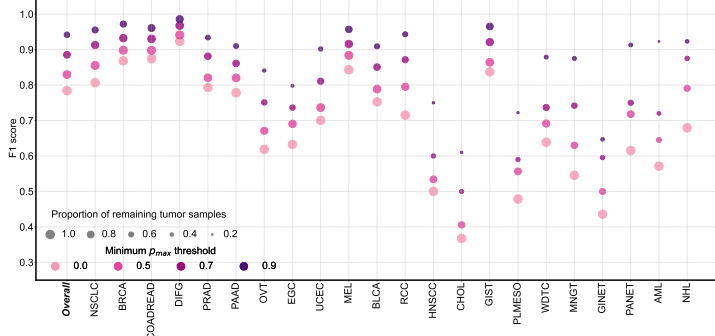
	DFCI	MSK	VICC	DFCI CUP	
Number of patients	18,106	15,151	1,310	962	
Patients age at sequence (95 % C.I.)	60.7 (60.5 - 60.9)	60.2 (60.0 - 60.4)	58.3 (57.6 - 59.0)	61.9 (61.1 - 62.7)	
Sex; male:female ratio	43.8 - 56.2	43.5 - 56.5	44.5 - 55.5	50.0 - 50.0	
Patients ethnicity (proportion %)					
White	16,105 (88.9 %)	11,575 (76.4 %)	1,089 (83.1 %)	853 (88.7 %)	
Black	538 (3.0 %)	866 (5.7 %)	72 (5.5 %)	38 (4.0 %)	
Asian	554 (3.1 %)	956 (6.3 %)	17 (1.3 %)	34 (3.5 %)	
Hispanic	379 (2.1 %)	744 (4.9 %)	14 (1.1 %)	15 (1.6 %)	
Others	530 (2.9 %)	1010 (6.7 %)	118 (9.0 %)	22 (2.2 %)	
Sequenced Tumor Samples					
Total number of samples	18,816	16,294	1,335	971	
Panel version (proportion %; 95% sequence date range)					
v1	OncoPanel v1	MSK-IMPACT341	VICC-01-T5A	OncoPanel v1	
	1,924 (10.2 %; 2013-8-20 - 2014-8-17)	1,803 (11.1 %; Not available)	307 (23.0 %; Not available)	47 (4.8 %; 2013-9-8 - 2014-8-12)	
v2	OncoPanel v2	MSK-IMPACT410	VICC-01-T7	OncoPanel v2	
	5,304 (28.2 %; 2014-9-28 - 2016-10-5)	6,917 (42.5 %; Not available)	1,028 (77.0 %; Not available)	203 (20.9 %; 2014-11-5 - 2016-10-5)	
v3	OncoPanel v3	MSK-IMPACT468		OncoPanel v3	
11,588 (61.6 %; 2016-11-11 - 2021-1-6)	7,574 (46.5 %; Not available)			721 (74.3 %; 2016-12-14 - 2020-12-23)	
Biopsy site type					
Primary	11,662 (62.0 %)	9,576 (58.8 %)	622 (46.6 %)	.	
Metastatic recurrence	5,737 (30.5 %)	6,718 (41.2 %)	637 (47.7 %)	.	
Local recurrence	673 (3.6 %)	Not available	64 (4.8 %)	.	
Unspecified/others	744 (4.0 %)	Not available	12 (0.9 %)	.	
Cancer group	OncoTree Cancer type	Predicted cancer type			
Lung (Thoracic)	Non-Small Cell Lung Cancer (NSCLC)	3,489 (18.5 %)	3,183 (19.5 %)	137 (10.3 %)	280 (28.8 %)
	Pleural Mesothelioma (PLMESO)	258 (1.4 %)	118 (0.7 %)	2 (0.1 %)	9 (0.9 %)
Gastrointestinal	Colorectal Adenocarcinoma (COADREAD)	2,525 (13.4 %)	1,919 (11.8 %)	232 (17.4 %)	63 (6.5 %)
	Esophagogastric Adenocarcinoma (EGC)	988 (5.3 %)	495 (3.0 %)	59 (4.4 %)	69 (7.1 %)
	Pancreatic Adenocarcinoma (PAAD)	772 (4.1 %)	980 (6.0 %)	53 (4.0 %)	85 (8.8 %)
	Cholangiocarcinoma (CHOL)	241 (1.3 %)	338 (2.1 %)	44 (3.3 %)	33 (3.4 %)
	Gastrointestinal Neuroendocrine Tumors (GINET)	219 (1.2 %)	76 (0.5 %)	18 (1.3 %)	46 (4.7 %)
	Pancreatic Neuroendocrine Tumor (PANET)	121 (0.6 %)	133 (0.8 %)	12 (0.9 %)	23 (2.4 %)
Sarcoma	Gastrointestinal Stromal Tumor (GIST)	273 (1.5 %)	217 (1.3 %)	5 (0.4 %)	3 (0.3 %)
Head and Neck	Head and Neck Squamous Cell Carcinoma (HNSCC)	473 (2.5 %)	285 (1.7 %)	20 (1.5 %)	52 (5.4 %)
	Well-Differentiated Thyroid Cancer (WDTC)	166 (0.9 %)	166 (1.0 %)	8 (0.6 %)	1 (0.1 %)
Skin	Melanoma (MEL)	729 (3.9 %)	619 (3.8 %)	187 (14.0 %)	43 (4.4 %)
Breast	Invasive Breast Carcinoma (BRCA)	2,558 (13.6 %)	3,113 (19.1 %)	274 (20.5 %)	85 (8.8 %)
Gynecologic	Ovarian Epithelial Tumor (OVT)	1,213 (6.4 %)	525 (3.2 %)	81 (6.1 %)	58 (6.0 %)
	Endometrial Carcinoma (UCEC)	703 (3.7 %)	703 (4.3 %)	34 (2.5 %)	18 (1.9 %)
Hematologic	Acute Myeloid Leukemia (AML)	150 (0.8 %)	1 (0.0 %)	0 (0.0 %)	1 (0.1 %)
	Non-Hodgkin Lymphoma (NHL)	110 (0.6 %)	88 (0.5 %)	0 (0.0 %)	1 (0.1 %)
Genitourinary	Prostate Adenocarcinoma (PRAD)	601 (3.2 %)	1,222 (7.5 %)	27 (2.0 %)	27 (2.8 %)
	Renal Cell Carcinoma (RCC)	457 (2.4 %)	497 (3.1 %)	39 (2.9 %)	24 (2.5 %)
	Bladder Urothelial Carcinoma (BLCA)	550 (2.9 %)	505 (3.1 %)	41 (3.1 %)	21 (2.2 %)
Neuro	Diffuse Glioma (DIFG)	2,041 (10.8 %)	1,069 (6.6 %)	47 (3.5 %)	25 (2.6 %)
	Meningothelial Tumor (MNGT)	179 (1.0 %)	42 (0.3 %)	15 (1.1 %)	4 (0.4 %)



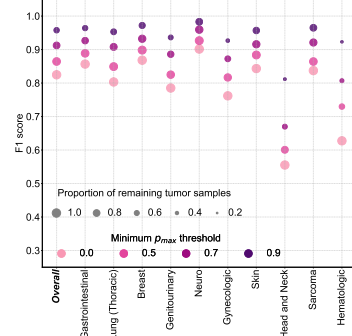
(a)



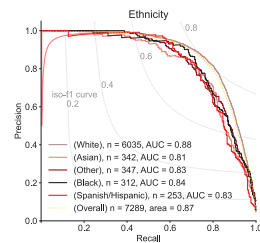
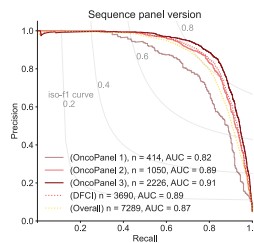
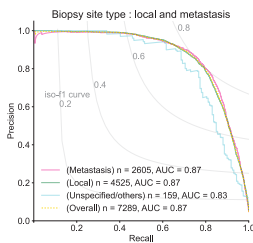
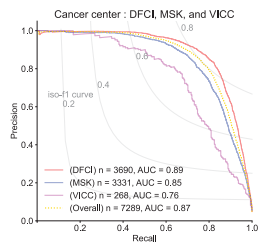
(b)



(c)

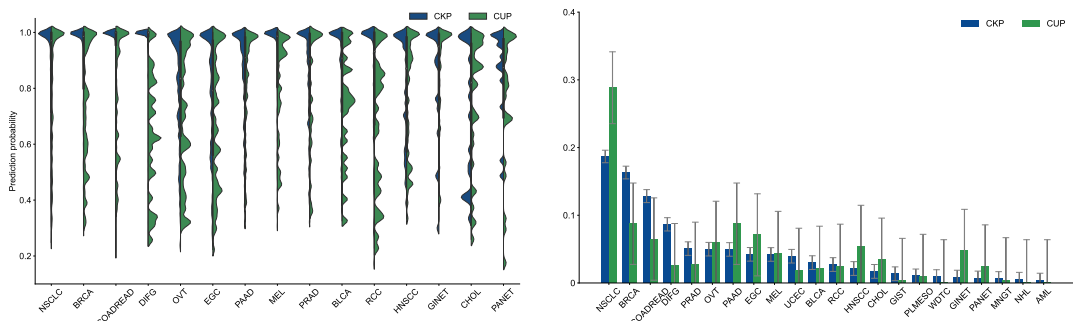


(d)



(e)

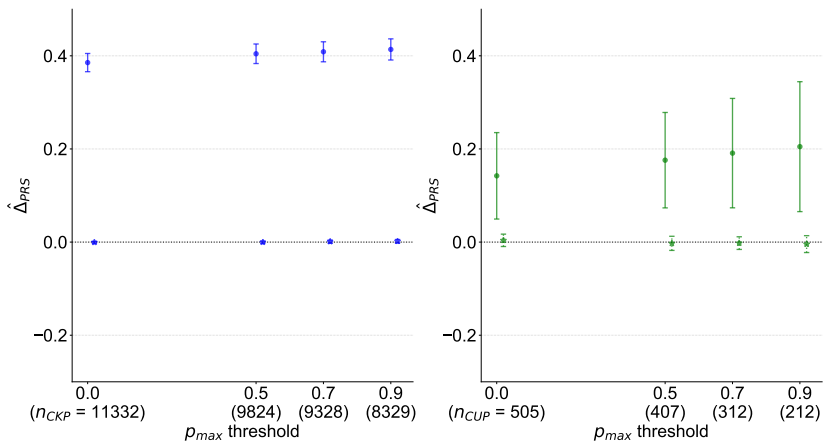
Figure 2



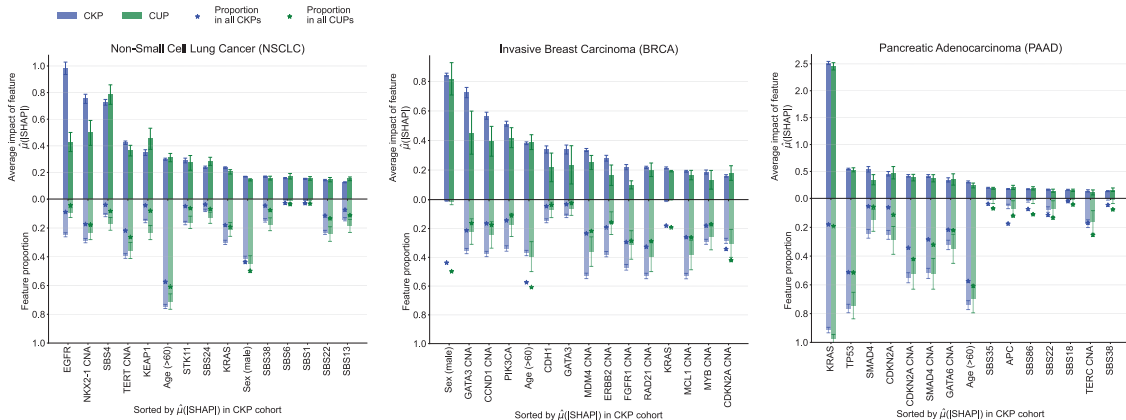
(a)

(b)

$\hat{\Delta}_{PRS}$  for CKP       $\hat{\Delta}_{PRS} - random$  for CKP       $\hat{\Delta}_{PRS}$  for CUP       $\hat{\Delta}_{PRS} - random$  for CUP



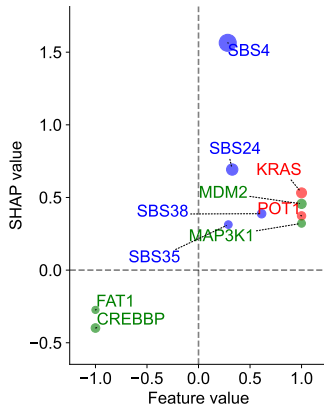
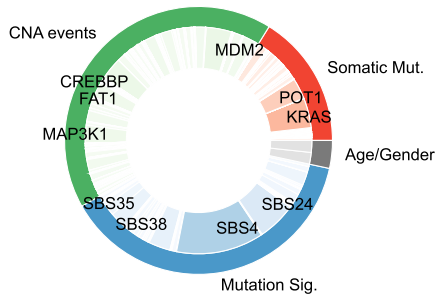
(c)



(d)

Figure 3

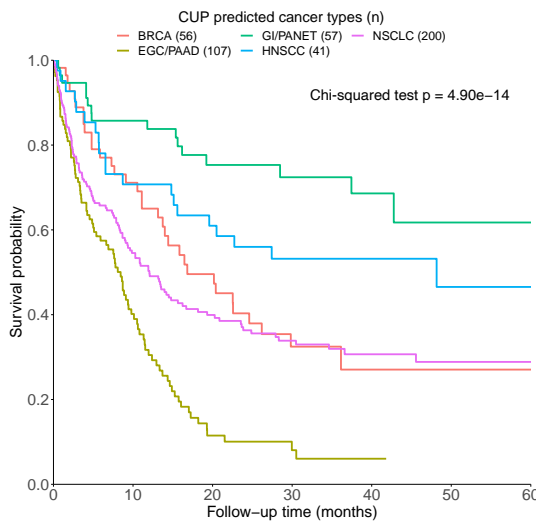
Predicted cancer type : Non-Small Cell Lung Cancer (NSCLC)  
Posterior probability : 0.98



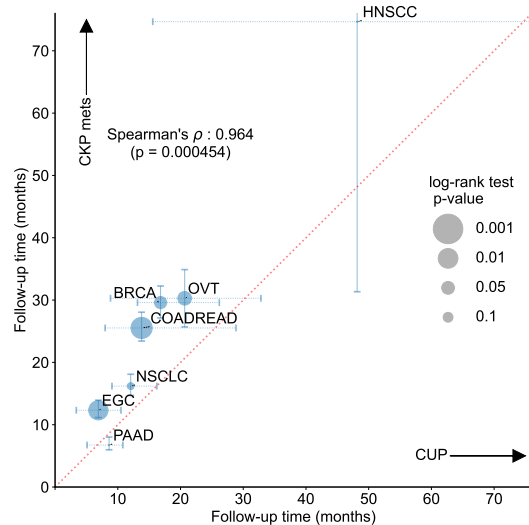
(e)

Figure 3

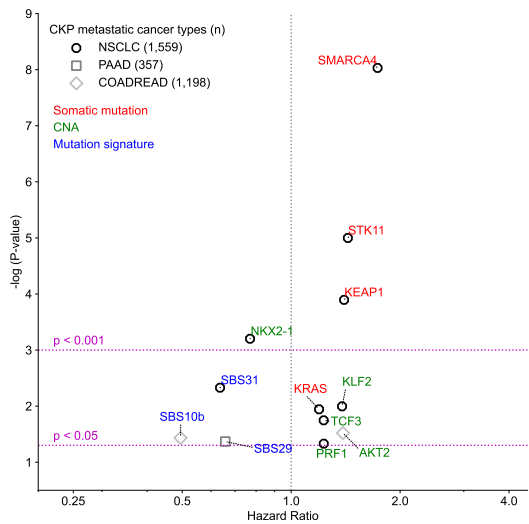




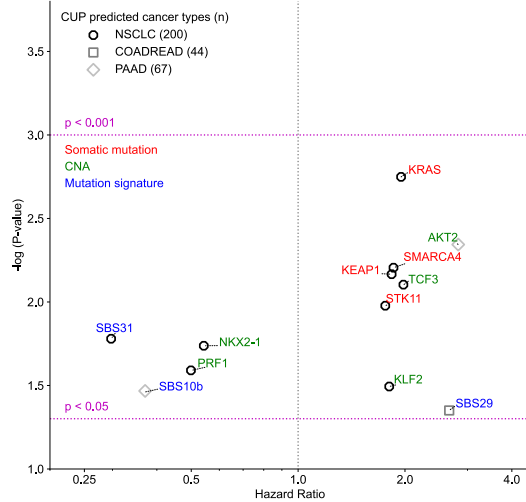
(a)



(b)



(c)

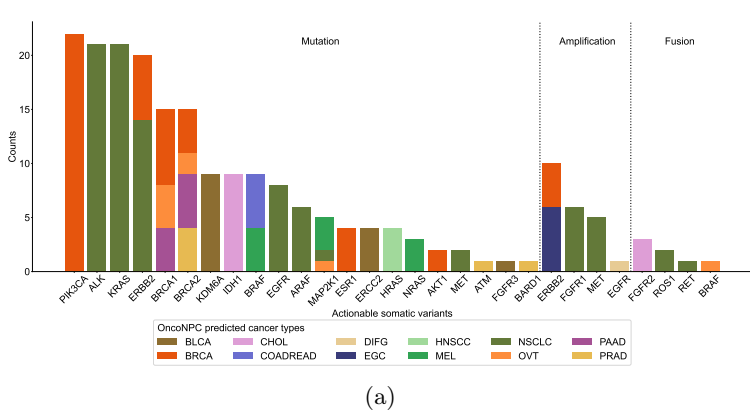


(d)

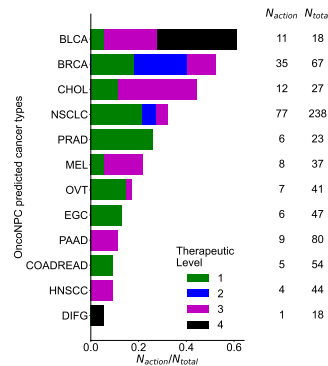
Figure 4

Table 2

	Concordant treatment group (n = 77)	Discordant treatment group (n = 81)
Sex; male-female ratio	0.442-0.558	0.556-0.444
Age at sequencing (95% C.I.)	64 (61.6 - 66.4)	62 (59.4 - 64.6)
Prediction uncertainty (in entropy; 95% C.I.)	0.550 (0.426 - 0.675)	0.988 (0.850 - 1.127)
OncoPanel version (proportion in %)		
v1	1 (1.30%)	1 (1.24%)
v2	9 (11.7%)	15 (18.5%)
v3	67 (87.0%)	65 (80.2%)
Mutational burden (95% C.I.)	0.027 (0.021 - 0.033)	0.033 (0.027 - 0.040)
CNA burden (95% C.I.)	0.201 (0.166 - 0.236)	0.186 (0.155 - 0.217)
Predicted primary cancer groups (proportion in %)		
Lung	15 (19.5%)	24 (29.6%)
Breast	5 (6.50%)	11 (13.6%)
GI	33 (42.9%)	16 (19.8%)
Gyn	9 (11.7%)	5 (6.17%)
Others	15 (19.5%)	25 (31.0%)
Metastatic sites (proportion in %)		
Brain	4 (5.20%)	8 (9.88%)
Bone	7 (9.10%)	10 (12.3%)
Soft tissue	6 (7.79%)	5 (6.17%)
Others	60 (77.9%)	58 (71.6%)
Histology (proportion in %)		
Adenocarcinoma	41 (53.2%)	32 (39.5%)
Neuroendocrine	9 (11.7%)	11 (13.6%)
Squamous cell	2 (2.60%)	6 (7.41%)
Others	25 (32.5%)	32 (39.5%)
Treatment start date (95% C.I.)	2018-4-30 (2017-12-24 - 2018-9-3)	2018-3-1 (2017-10-28 - 2018-7-3)

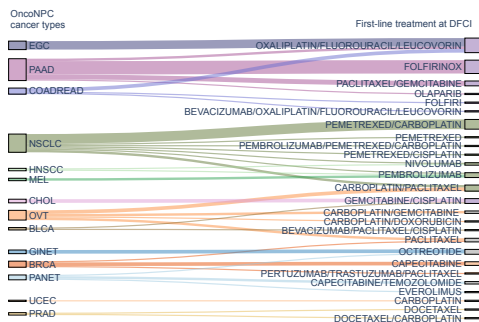


(a)



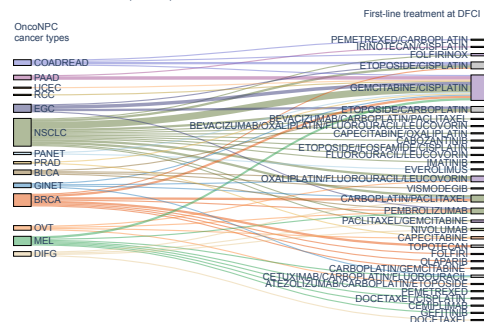
(b)

Concordant treatments (n = 77)

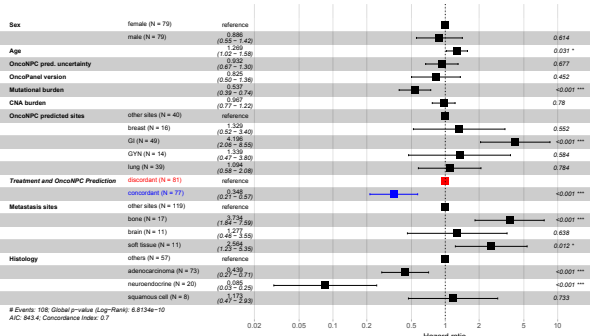


(c)

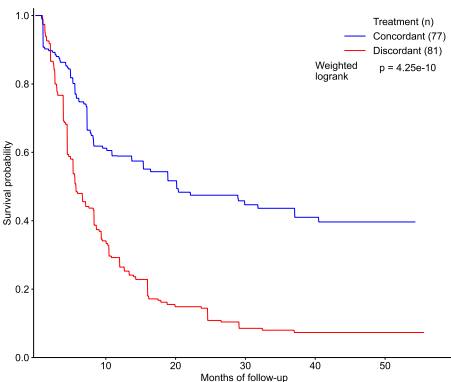
Discordant treatments (n = 81)



(d)



(e)



(f)

Figure 5