

1 **Convergent Usage of Amino Acids in Human Cancers as a Reversed**  
2 **Process of Tissue Development**

3

4

Yikai Luo<sup>1,2,a</sup>, and Han Liang<sup>2,3,1,\*,b</sup>

5

6 <sup>1</sup> *Graduate Program in Quantitative and Computational Biosciences, Baylor College of*  
7 *Medicine, Houston TX 77030, USA*

8 <sup>2</sup> *Department of Bioinformatics and Computational Biology, The University of Texas MD*  
9 *Anderson Cancer Center, Houston, TX 77030, USA*

10 <sup>3</sup> *Department of Systems Biology, The University of Texas MD Anderson Cancer Center,*  
11 *Houston TX 77030, USA*

12

13 \* Corresponding author.

14 E-mail: [hliang1@mdanderson.org](mailto:hliang1@mdanderson.org) (Liang H)

15

16 **Running title:** *Luo Y / Convergent usage of amino acids in cancer*

17

18 <sup>a</sup>ORCID: 0000-0001-7589-7981.

19 <sup>b</sup>ORCID: 0000-0001-7633-286X.

20

21 Total word counts: 4053

22 Total figures: 6

23 Total tables:0

24 Total supplementary figures: 7

25 Total supplementary tables: 0

26 Total supplementary files: 0

27 **Abstract**

28 Genome and transcriptome-wide amino acid usage preference across different species is a  
29 well-studied phenomenon in molecular evolution, but its characteristics and implication in  
30 cancer evolution and therapy remain largely unexplored. Here, we analyzed large-scale  
31 transcriptome/proteome profiles such as TCGA, GTEx, and CPTAC and found that compared  
32 to normal tissues, different cancer types showed a convergent pattern towards using  
33 biosynthetically low-cost amino acids. Such a pattern can be accurately captured by a single  
34 index based on the average biosynthetic energy cost of amino acids, termed Energy Cost Per  
35 Amino Acid (ECPA). With this index, we further compared the trends of amino acid usage  
36 and the contributing genes in cancer and tissue development and revealed their reversed  
37 patterns. Finally, focusing on the liver, a tissue with a dramatic increase in ECPA during  
38 development, we found that ECPA represented a powerful biomarker that could distinguish  
39 liver tumors from normal liver samples consistently across 11 independent patient cohorts  
40 (AUROC = ~0.99) and outperformed any index based on single genes. Our study reveals an  
41 important principle underlying cancer evolution and suggests the global amino acid usage as  
42 a system-level biomarker for cancer diagnosis.

43

44 **KEYWORDS:** Amino acid usage; Tissue development; Biosynthetic energy; Diagnostic  
45 biomarker

## 46 **Introduction**

47 Amino acids are the basic building blocks of a cell. Coding sequences and gene expression  
48 profiles are two key factors determining the overall amino acid usage of a cell. Through  
49 analysis of the genomes or transcriptomes of many species, preferred amino acid usage is a  
50 well-studied topic in macroevolution. The universal trend of “Cost-Usage anticorrelation”  
51 suggests that the relative abundance of amino acids, quantified as the number of codons  
52 encoding a specific amino acid in the genome of a species, is mainly driven by their  
53 biosynthetic energy costs [1–5]. However, it remains unclear how amino acid usage of cancer  
54 cells deviates from normal tissues and evolve in different tumor contexts.

55 From an evolutionary point of view, cancer cells are characterized by a low degree of  
56 divergence from its tissue of origin, measured by the limited amount of somatic changes,  
57 which is in contrast to the macroevolution that happens across different taxa or even the  
58 microevolution existing between within-species individuals [6]. However, such trifling  
59 transformation does yield a wide range of phenotypic commonalities shared by distinct  
60 cancer types, including activated proliferative signaling, resistance to programmed cell death,  
61 induction of angiogenesis, and metastatic capability [7]. Among many theories proposed to  
62 understand such convergence, one appealing concept is that cancer cells bear a set of  
63 genomic, transcriptomic, and epigenomic features that can be summed up as “stemness,” [8–  
64 11] which in the context of ontogeny, defines the level of reprogramming/dedifferentiation of  
65 adult tissue cells. The underlying mechanistic links between cancer evolution and tissue  
66 development have been hinted at by the observations of frequent mutations leading to  
67 reactivation of stem cell-related pathways in cancer [12,13]. However, little effort has been  
68 made to examine a potential association between these two seemingly non-overlapping  
69 processes with respect to amino acid usage.

70 Characterizing the amino acid usage of cancer cells not only helps us understand the  
71 evolutionary constraints in the tumor microenvironment but may also have clinical utility. In  
72 recent years, tremendous efforts have been made to identify gene expression-based  
73 biomarkers for cancer diagnosis, outcome prediction, and treatment selection, but successful  
74 cases with proven clinical values are still limited [14–16]. One factor that determines the  
75 feasibility of such biomarkers in clinical practice, the robustness, is rarely satisfied, meaning  
76 that a threshold chosen based on limited data is usually not generalizable to unseen scenarios.  
77 In contrast to conventional biomarkers based on individual genes, the amino acid usage  
78 represents a holistic property of a cellular state. Therefore, there is a possibility that its related

79 indices represent more robust biomarkers for clinical applications. To fill these knowledge  
80 gaps, here we performed a systematic analysis of the amino acid usage profiles across many  
81 cohorts of tumor and normal tissue samples.

82

## 83 **Results**

### 84 **A convergence of amino acid usage across cancer types**

85 Since gene expression levels are largely associated with amino acid usage in a cell, we first  
86 examined the gene expression patterns of 30 tissue types in the Genotype-Tissue Expression  
87 (GTEx) cohort [17] (Figure S1A) and 31 cancer types in The Cancer Genome Atlas (TCGA)  
88 cohort [18] (Figure S1B). Using the t-distributed stochastic neighborhood embedding (t-  
89 SNE)[19] projection, we found that samples of a common tissue origin largely formed a  
90 single cluster regardless of being normal or cancerous. In addition, cancer types with the  
91 same tissue origin, such as brain cancers (glioblastoma multiforme [GBM] and lower grade  
92 glioma [LGG]), kidney cancers (kidney renal clear cell carcinoma [KIRC] and kidney renal  
93 papillary cell carcinoma [KIRP]), lung cancers (lung adenocarcinoma [LUAD] and lung  
94 squamous cell carcinoma [LUSC]), and liver cancers (hepatocellular carcinoma [LIHC] and  
95 cholangiocarcinoma [CHOL]), tended to be mingled or closer to each other than to other  
96 cancer types. We observed similar patterns in two other large, pan-cancer cohorts, PCAWG  
97 [20], and MET500 [21] (Figure S1C and D). Consistent with previous studies [18,22], these  
98 results indicate that cancer cells largely retain their tissue-specific gene expression profiles.

99 To study whether this tissue-specific pattern holds for amino acid usage, we calculated the  
100 similarity of transcriptome-based amino acid usage by integrating the gene expression  
101 profiles and the amino acid frequencies of protein-coding genes (**Figure 1A**) and visualized  
102 their patterns in the same way. Similar to the strong tissue specificity observed in the gene  
103 expression analysis, we found that normal tissues of the GTEx cohort still had distinct amino  
104 acid usage patterns (Figure 1B). We further confirmed this result by co-clustering amino acid  
105 usage profiles of the Human Protein Atlas (HPA) cohort [23] with corresponding GTEx  
106 tissue types (Figure S2A). More intriguingly, samples of a multi-species multi-tissue cohort  
107 [24] were principally separated by tissue type rather than by species, suggesting that tissue-  
108 specific amino acid usage is highly conserved across mammals (Figure 1C).

109 In sharp contrast to normal tissues, when clustered by amino acid usage, samples of  
110 different cancer types were much less separated and did not segregate on the basis of tissue  
111 origins (Figure 1D). To further confirm this observation, we clustered amino acid usage

112 profiles of two other cancer cohorts, PCAWG and MET500, and observed a dramatic loss of  
113 tissue-specificity relative to the patterns observed in the gene expression-based analysis  
114 (Figure S1C and D and Figure S2B and C). To ensure that the detected pattern was not due to  
115 a disparity in sample size or unmatched tissue types, we leveraged a conservative GTEx-  
116 TCGA mapping to only include normal and tumor samples whose tissue origins are matched  
117 without ambiguity, then performed down-sampling within individual tissue-specific cohorts,  
118 and finally, applied t-SNE to redo a supervised clustering. The results remained the same for  
119 the comparison between down-sampled GTEx and TCGA samples (Figure S2D and E) as  
120 well as for that between TCGA tumor samples and the normal adjacent to tumor (NAT)  
121 (Figure S2F and G). This observation is important since, evaluating tumor purity and gene  
122 signatures, recent studies have shown that NAT samples reside in an intermediate state  
123 between healthy and tumor [25,26].

124 The observation that amino acid usage for cancer cells failed to preserve their distinct tissue  
125 origins raised two possibilities: (i) cancer cells evolved to possess highly stochastic amino  
126 acid usage profiles both within and between cancer types; or (ii) they went through  
127 convergence of amino acid usage, thereby losing the constraint of the original tissue  
128 specificity. To identify the correct hypothesis, we simply asked whether, in the 20-  
129 dimensional space (each dimension representing the frequency of specific amino acid), the  
130 distances between samples of different cancer types were shorter than those among samples  
131 of different normal tissues. Based on Pearson's distance, for each sample, we defined an  
132 amino acid usage convergence index that measured its distance to all other samples of  
133 different tissue or cancer types. Through a comparative analysis of GTEx normal vs. TCGA  
134 tumor and TCGA NAT vs. tumor, we found that tumor samples showed significantly  
135 increased convergence than normal samples, a pattern consistently observed across all  
136 surveyed cancer types (Figure 1E and F). Furthermore, we compared the variations of amino  
137 acid frequencies across NAT samples and tumor samples of different cancer types based on  
138 the same set of standard deviations. Indeed, the extent to which amino acids are differentially  
139 used in tumors was markedly reduced than that in NATs (Figure S4A and B). Collectively,  
140 these results indicated a strong convergence rather than a stochastic transformation of amino  
141 acid usage across cancer types, supporting our second hypothesis.

142

### 143 **Cancer cells tend to use biosynthetically low-cost amino acids**

144 To understand how such a convergent pattern occurs, we quantified the differential usage of  
145 each amino acid in tumors vs. normal tissues and found no highly consistent trend across

146 cancer types in terms of increased or decreased usage (Figure S4C). However, when taking a  
147 higher view of the heatmap, structurally complex amino acids, such as tryptophan and  
148 cysteine, tended to be significantly depleted in most cancer types, whereas those with  
149 relatively simpler structures tended to be significantly enriched in a majority of cancers.  
150 Because the structural complexity of the amino acids correlates well with the energy cost of  
151 their biosynthesis [1], we hypothesized an association between the biosynthetic energy cost  
152 of amino acid and its usage tendency in cancers. Indeed, we observed a strong negative  
153 correlation between the biosynthetic energy cost and the net number of cancer types in which  
154 the usage of an amino acid was significantly increased (**Figure 2A**,  $R_s = -0.56$ ,  $p = 0.01$ ),  
155 suggesting that cancer cells prefer amino acids with a lower biosynthetic energy cost. We  
156 previously introduced two indices,  $ECPA_{gene}$ , and  $ECPA_{cell}$ , which quantify the average  
157 biosynthetic energy cost per amino acid for a gene and a cell (or a sample), respectively [27]  
158 (Figure 2B).  $ECPA_{gene}$  is based on the amino acid frequency encoded in a gene, and  $ECPA_{cell}$   
159 considers the expression levels and amino acid frequencies of all the genes in a cell. A high  
160 ECPA value indicates that the gene or the cell tends to use biosynthetically expensive amino  
161 acids. We found that compared to NAT samples,  $ECPA_{cell}$  of the tumor samples became  
162 significantly lower for 9 out of the 15 tested cancer types, while no significantly opposite  
163 patterns were observed (Figure 2C). To confirm this pattern at the proteomic level, we  
164 extended these analyses to six cancer proteomics datasets from the Clinical Proteomic Tumor  
165 Analysis Consortium (CPTAC) [28] and others [29,30], covering five cancer types.  
166 Strikingly, in all the cases, proteins that were significantly up-regulated in tumor samples  
167 ( $\log_2FC > 0$ ,  $FDR < 0.05$ ) had significantly lower  $ECPA_{gene}$  than the proteins that were  
168 significantly down-regulated ( $\log_2FC < 0$ ,  $FDR < 0.05$ ) (Figure 2D). These results indicate  
169 that cancer cells reshaped their gene/protein expression programs to use biosynthetically  
170 inexpensive (or structurally simpler) amino acids, thereby losing their original tissue-specific  
171 amino acid usage profiles. Finally, we sought to test if our ECPA index is insensitive to the  
172 expression of genes with extremely high abundance, including those encoding certain  
173 housekeeping proteins as well as tissue-specific proteins. After removal of all genes that  
174 either encode cytoplasmic and mitochondrial ribosome proteins or rank within top 200 in  
175 median TPM of the same cancer type, we recalculated the ECPA index for each sample and  
176 found that the pattern of consistent decrease of  $ECPA_{cell}$  in tumor samples across multiple  
177 cancer types was almost perfectly reproduced (Figure S3).

178 We next tested whether the amino acid usage convergence level of a tumor was correlated  
179 with its  $ECPA_{cell}$ . Indeed, we found a strong inverse relationship for seven out of the nine

180 cancer types where  $ECPA_{cell}$  was significantly lower in tumors (Figure 2E). Thus, the more a  
181 tumor follows a convergent path to a common state of amino acid usage, the higher the bias it  
182 has toward using biosynthetically low-cost amino acids. These results also suggest that  
183  $ECPA_{cell}$  is a simple, informative, interpretable index that effectively captures the overall  
184 preference of amino acid usage for a specific sample. Therefore, we focused on this index in  
185 subsequent analyses.

186

### 187 **Biosynthetically expensive amino acids are increasingly used during tissue development**

188 To elucidate the underlying cause for the convergence of amino acid usage in cancer, we first  
189 sought to understand how tissue-specific amino acid usage patterns are established during  
190 development. Using the  $ECPA_{cell}$  index, we quantified the overall amino acid usage of liver  
191 and kidney tissues across different development stages in mammals, including humans, mice,  
192 rats, rabbits, and opossums. Intriguingly, both tissues showed an increasing trend of  $ECPA_{cell}$   
193 along their developmental trajectories in all five mammals (**Figure 3A and B**). A closer  
194 inspection of the  $ECPA_{cell}$  trend lines led to two observations: i) key turning points of  
195  $ECPA_{cell}$  in different species tend to happen at corresponding developmental stages; and ii)  
196 the rise of  $ECPA_{cell}$  in the liver takes concave trajectories while that in the kidney takes  
197 convex trajectories, suggesting that the establishment of high  $ECPA_{cell}$  status is driven by  
198 evolutionarily conserved synchronous molecular events that possess strong tissue specificity.  
199 To confirm this pattern, we collected another three independent RNA-seq datasets on mouse  
200 liver development and found a consistent  $ECPA_{cell}$  increase along the developmental paths in  
201 all three cases (Figure 3C-E).

202 To pinpoint which gene modules are responsible for the tissue-specific build-up of a high  
203  $ECPA_{cell}$  status, we first defined a “ $\Delta ECPA_{cell}$  contribution index” for each gene, which  
204 quantified the contribution of the gene to the global shift of  $ECPA_{cell}$  (see Materials and  
205 methods). We then divided all genes into 15 equal bins based on their index values and  
206 employed a mutual information-based enrichment identification algorithm called iPAGE [31]  
207 to detect the enrichment of these gene groups with well-established functional gene modules.  
208 We noted that genes contributing to the  $ECPA_{cell}$  increase were conserved among mammals  
209 but were tissue-specific. For the liver, the enriched modules included glucuronosyltransferase  
210 activity and complement activation (Figure 3F, Figure S5A, C and E); and for the kidney, the  
211 enriched modules included sphingolipid biosynthetic process and zinc/calcium ion  
212 homeostasis (Figure 3G, Figure S5B, D, and F).

213 Development-related cellular states that are instituted in adulthood can be prone to  
214 significant transformation or even complete collapse during aging [32]. To further understand  
215 how tissue-specific amino acid usage patterns alter when the tissue undergoes senescence, we  
216 gathered independent transcriptome profiles of aging livers and kidneys in humans, mice, and  
217 rats, and characterized the  $ECPA_{cell}$  patterns. Both tissues showed a stable pattern of high  
218  $ECPA_{cell}$  status with reasonable fluctuations (Figure S6A-C). We concluded that tissue-  
219 specific, preferred usage of biosynthetically expensive (or structurally complex) amino acids,  
220 characterized by a high- $ECPA_{cell}$  status, was gradually formed during development and  
221 remained largely unchanged in aging.

222

### 223 **Amino acid usage convergence of tumor follows a reversed path of tissue development**

224 The strong convergence of amino acid usage across different cancer types is reminiscent of  
225 the “reverse-evolution” concept for tumorigenesis. As demonstrated above, this idea is well  
226 illustrated by the observation that there is a consistent decline of  $ECPA_{cell}$  in tumors, whereas  
227 there is a gradual increase of  $ECPA_{cell}$  during tissue development. To test the hypothesis that  
228 cancer evolution and tissue development move in opposite directions with respect to amino  
229 acid usage, we assessed whether the genes that boosted  $ECPA_{cell}$  in tissue development were  
230 overlapped with those that reduced  $ECPA_{cell}$  in tumors of the corresponding tissue origin and  
231 vice versa. Following the same method of computing  $\Delta ECPA_{cell}$  contribution index for tissue  
232 development, we measured the contribution of individual genes to  $\Delta ECPA_{cell}$  in cancer  
233 evolution for three cancer types for which gene expression profiles of normal developing  
234 tissues are available, namely LIHC, KIRC, and KIRP. Based on their contributions to  
235  $\Delta ECPA_{cell}$  in either development or tumorigenesis, we divided individual genes into four  
236 quadrants with zero as the cutoff. We then used Fisher’s exact test to analyze the overlap of  
237 developmental  $\Delta ECPA_{cell}$ -positive-contributing genes with tumorigenic  $\Delta ECPA_{cell}$ -negative-  
238 contributing genes and vice versa. We observed that genes indeed tended to make opposite  
239 contributions to  $\Delta ECPA_{cell}$  in tumorigenesis and tissue development (**Figure 4A-C**, Fisher’s  
240 exact test, LIHC,  $p = 1.6 \times 10^{-156}$ ; KIRC,  $p = 1.9 \times 10^{-39}$ ; KIRP,  $p = 8.9 \times 10^{-30}$ ). Furthermore, for  
241 the genes reducing  $ECPA_{cell}$  in tumorigenesis and increasing  $ECPA_{cell}$  in development, their  
242  $\Delta ECPA_{cell}$  contribution index in these two processes were significantly negatively correlated  
243 (Figure 4D-F).

244 While the gene-level analyses above were possibly hindered by the fact that cancer  
245 progression is highly heterogeneous even within the same cancer type [33,34], we can expect  
246 that a sample-level analysis would be more efficient to detect potential reverse relationships



247 between cancer evolution and tissue development regarding amino acid usage. To this end,  
248 we defined the “developmental reversal index” for each tumor sample, which quantifies how  
249 strongly its gene expression pattern reversed what was instituted in tissue development.  
250 Specifically, we first calculated the gene-expression fold change of each tumor sample in  
251 terms of that averaged over the adjacent normal samples in order to measure the  
252 transcriptomic shift during tumorigenesis. We then measured the strength of anti-correlation  
253 between such a shift and the expression changes of the same gene set along the  
254 developmental trajectories of matched tissues (see Materials and methods). Interestingly,  
255 using this index to stratify cancer patients in terms of overall survival time, we found that a  
256 higher developmental reversal value was consistently associated with a worse prognosis  
257 (Figure 4G-I), suggesting that more aggressive tumors tend to have gene expression profiles  
258 more reversed in the tissue development trajectory.

259 Finally, we employed a multivariate linear regression model to clarify the associations  
260 between how biased a tumor sample tends to be in using biosynthetically inexpensive amino  
261 acids (represented by  $ECPA_{cell}$ ), how far it travels on the path of amino acid usage  
262 convergence relative to other cancer types (represented by amino acid usage convergence  
263 index), and how strongly its gene expression pattern reversed from what was instituted in  
264 tissue development (represented by the developmental reverse index). Remarkably, both the  
265 convergence level and the developmental reversal level were strongly anti-correlated with  
266  $ECPA_{cell}$  across cancer types (Figure 4J-L). We, therefore, put forward an integrated model in  
267 which cancer cells initiated from distinct tissue origins converge into a common state  
268 favoring the use of biosynthetically inexpensive amino acids through reversed paths of tissue  
269 development (Figure 4M).

270

### 271 **The amino acid usage index, $ECPA_{cell}$ , is a robust biomarker for liver cancer diagnosis**

272 Among different cancer types in our  $ECPA_{cell}$  analysis, the difference between liver normal  
273 and liver tumor samples was striking, making this tissue stand out from others (Figure 2C).  
274 Indeed, by quantifying the downward shift of  $ECPA_{cell}$  ( $\Delta ECPA_{cell}$ ) between tumor and the  
275 matched NAT pairs, the top two cancers were CHOL and LIHC, both of which originate  
276 from the liver (**Figure 5A**). We suspected that such a striking pattern could be attributed to  
277 liver-specific gene expression. To test this, we calculated  $ECPA_{cell}$  of both GTEx normal  
278 samples and TCGA NAT samples based only on tissue-specific genes [35] and ranked the  
279 tissues by their average  $ECPA_{cell}$ . Indeed, the liver  $ECPA_{cell}$  level was higher than almost all

280 other tissues (Figure 5B and C) (although the pancreas showed an even higher  $ECPA_{cell}$   
281 according to the GTEx samples, the pattern did not hold for TCGA NAT samples). Of note,  
282 while the sample size of LIHC-NAT was as large as 50, the variation of their  $ECPA_{cell}$  based  
283 on tissue-specific genes was low. Furthermore, a comparison of the developmental  $ECPA_{cell}$   
284 trend lines for different human tissues revealed that a fast and early build-up of a high-  
285  $ECPA_{cell}$  status only existed for the liver (Figure 5D). We observed similar patterns in other  
286 mammals as well (Figure S7A-D). These results suggest that during development, the liver  
287 acquires a very high  $ECPA_{cell}$  state, and the liver-specific genes are the underlying  
288 contributing factor.

289 Given (i) the extremely high  $ECPA_{cell}$  level of liver tissue, and (ii) the dramatic difference  
290 between liver tumor and matched normal samples, we speculated whether  $ECPA_{cell}$  could be  
291 utilized as a novel biomarker for detecting liver cancer. To this end, we first collected 11  
292 independent liver-cancer RNA-seq datasets (including TCGA LIHC and CHOL) where  
293 matched tumor and adjacent normal biopsies were simultaneously collected, thereby enabling  
294 a direct comparison of  $ECPA_{cell}$  between these conditions. In all cases, the tumor samples  
295 showed significantly reduced  $ECPA_{cell}$  with large effect sizes (**Figure 6A**).

296 To evaluate more rigorously the capacity of  $ECPA_{cell}$  to serve as a diagnostic marker in  
297 discriminating liver tumors from normal tissues, we employed the area under the receiver  
298 operating characteristic curve (AUROC) as a performance metric. To ensure the robustness  
299 of our analyses, we only included six datasets with sample size  $\geq 12$ . The  $ECPA_{cell}$  index was  
300 able to separate tumor vs. normal samples with very high ROC scores (median value = 0.993,  
301 range = 0.982 - 1.00, Figure 6B). To compare the predictive power of  $ECPA_{cell}$  relative to  
302 individual gene-based biomarkers, we calculated the average AUROC of all detectable genes  
303 across the six datasets and assessed their performance in the same way. Among 9,559 genes  
304 assessed, only three genes (*CCT3*, *DDX39A*, and *FLAD1*) showed slightly better performance  
305 than  $ECPA_{cell}$  (0.992), but none of them had statistically significant superiority (Figure 6C  
306 and D). In addition,  $ECPA_{cell}$  showed significantly higher discriminating power than the  
307 usage of any single amino acid (Figure 6E). Along with accuracy, a key feature of a  
308 successful biomarker is its robustness. To assess this feature, we computed the coefficient of  
309 variation (CV) for the optimal thresholds of  $ECPA_{cell}$  and individual genes across different  
310 datasets as an indicator of robustness.  $ECPA_{cell}$  showed exceptionally high robustness with its  
311 CV as low as  $7.9 \times 10^{-3}$ , about  $5 \times$  smaller than the lowest CV of any single gene-based  
312 biomarker (Figure 6F). Notably, the three genes that had a statistically insignificant  
313 advantage over  $ECPA_{cell}$  by AUROC had extremely unstable optimal cutoffs among different

314 datasets, suggesting their limited power in detecting liver cancer across diverse clinical  
315 scenarios. Collectively, these results suggest that, as a system-level feature capturing the  
316 global usage of amino acids in a sample,  $ECPA_{cell}$  represents a promising biomarker for liver  
317 cancer diagnosis, and possesses both high accuracy and exceptional robustness.

## 318 **Discussion**

319 Here we performed a systematic analysis on transcriptome and proteome-based amino acid  
320 usage across a broad range of cancer types. Using a previously introduced index,  $ECPA_{cell}$ ,  
321 our results revealed, for different tumors, a convergent pattern toward a cellular state of using  
322 more biosynthetically low-cost amino acids. In parallel, we studied the amino acid usage in  
323 the developmental trajectories of multiple organs and uncovered diverse paths into a tissue-  
324 specific high- $ECPA_{cell}$  status that were evolutionarily conserved across mammals. Thus, a  
325 reverse relationship existed between cancer evolution and tissue development, which can be  
326 viewed as reminiscent of the widely accepted concept of the cancer cell “stemness.”  
327 Furthermore, given the long-standing parallels between phylogeny and ontogeny [36],  
328 supported by recent evidence [24,37,38], it would be reasonable to interpret cancer evolution  
329 as a reversed process of not only the development of an organism or its tissues but also the  
330 evolution of species. It has been argued that one key mechanism adopted by cancer cells to  
331 obtain fitness in spite of the diversity of the microenvironments is to unleash the force that is  
332 suppressed in multicellular organisms but is borne by unicellular organisms that are at the  
333 very bottom of the evolutionary hierarchy [39–43]. Thus, amino acid usage, a key aspect of  
334 cellular metabolism, may provide a unique perspective to understand the fundamental  
335 principles governing cancer progression, tissue development, and macroevolution, three  
336 evolutionary processes on different scales.

337 With the advances in transcriptome profiling technology, gene expression-based  
338 biomarkers have attracted wide attention for tumor detection and patient stratification.  
339 However, due to the high heterogeneity of cancer and intrinsically stochastic nature of gene  
340 expression, biomarkers based on either a single gene or a set of genes tend to suffer from  
341 numerical instability, thereby performing poorly. As demonstrated for liver cancer diagnosis,  
342 our  $ECPA_{cell}$  index represents a system-level biomarker that has at least three remarkable  
343 advantages. First,  $ECPA_{cell}$  captures a global cellular state by retaining the entire  
344 transcriptome as its information source, thereby conferring unparalleled robustness. Second,  
345  $ECPA_{cell}$  was derived *de novo* from the gene expression profile of a sample, thus independent  
346 of external reference, which might introduce large noise predominantly attributable to batch

347 effect. Third, in contrast to data-driven metrics,  $ECPA_{\text{cell}}$  has a well-defined biological  
348 meaning, the biosynthetic energy cost of amino acids. Because of these properties,  $ECPA_{\text{cell}}$   
349 is an extremely robust diagnostic biomarker for liver cancer with a nearly constant threshold  
350 for tumor-normal segregation. Further efforts are warranted to assess the utility of this index  
351 in other cancer types and clinical applications.

352

## 353 **Materials and methods**

### 354 **Data acquisition and processing**

355 We obtained the gene-level expression values (e.g., fragments per kilobase per million  
356 [FPKM] or transcripts per million [TPM]) of the TCGA cancer sample cohorts, the GTEx  
357 normal tissue cohort, and the MET500 metastatic tumor cohort, from the Xena data portal  
358 (<https://xenabrowser.net/datapages/>); the HPA cohort from the HPA data portal  
359 (<http://www.proteinatlas.org/>); and the PCAWG cohort from the ICGC data portal  
360 (<https://dcc.icgc.org/releases/PCAWG/transcriptome/>). We also obtained the gene expression  
361 values of the mammalian tissue development cohorts  
362 from ArrayExpress (<https://www.ebi.ac.uk/arrayexpress/>) under the accession IDs E-MTAB-  
363 6769 (chicken), E-MTAB-6782 (rabbit), E-MTAB-6798 (mouse), E-MTAB-6811 (rat), E-  
364 MTAB-6813 (rhesus macaque), E-MTAB-6814 (human), and E-MTAB-6833 (opossum);  
365 and two independent RNA-seq datasets of mouse liver development from the Gene  
366 Expression Omnibus (GEO) under the accession IDs GSE58733 and GSE58827, as well as  
367 from ArrayExpress under the accession ID E-MTAB-2328. Finally, we obtained RNA-seq  
368 datasets of aging mouse liver and kidney from GEO under the accession ID GSE132040.

369 To convert gene-level FPKM values to TPM [44] values for a gene  $g_i$  in a sample  $s_k$ , we  
370 used the formula:

$$TPM_{g_i, s_k} = \frac{FPKM_{g_i, s_k}}{\sum_{j=1}^n FPKM_{g_j, s_k}} \times 10^6$$

371 where the denominator on the right side is the sum of FPKM values of all the genes for an  
372 individual sample.

373 We downloaded raw RNA-seq fastq files of human liver cancer from GEO under the  
374 accession IDs GSE65485, GSE119336, GSE77314, GSE77509, GSE63863, GSE94660,  
375 GSE25599, GSE124535, and GSE55758; files of aging rat liver from the Sequence Read  
376 Archive (SRA) under the accession ID PRJNA516151, and files of TCGA LIHC and CHOL  
377 cohorts from the GDC Data Portal (<https://portal.gdc.cancer.gov/>). MultiQC [45] was used to

378 assess the quality of the sequencing files and the performance of the preprocessing steps.  
379 Transcript-level abundances were quantified by Salmon [46] using the GRCh38  
380 transcriptome as the reference. Gene-level TPM values were aggregated from transcript-level  
381 TPM values by tximport [47].

382 We obtained the proteomics datasets of KIRC, COAD, LUAD, and OV patient cohorts  
383 from the CPTAC data portal (<https://cptac-data-portal.georgetown.edu/>). We obtained two  
384 proteomics datasets of liver cancer from the NODE data portal  
385 (<https://www.biosino.org/node/index/>) and the CNHPP data portal  
386 (<http://liver.cnhpp.ncpsb.org/>), respectively.

387

### 388 **Calculation of transcriptome-based amino acid usage**

389 We used the following formula to compute the amino acid frequency matrix given an RNA-  
390 seq dataset (see also Fig. 1a):

$$F_{m \times 20} = E_{m \times n} A_{n \times 20}^T$$

391 where  $E$  is a matrix of genes  $g_1, g_2, \dots, g_n$  by samples  $s_1, s_2, \dots, s_m$  with entries as TPM  
392 values, and  $A$  is a matrix of genes  $g_1, g_2, \dots, g_n$  by amino acids  $a_1, a_2, \dots, a_{20}$  with entries as  
393 relative frequencies of amino acids computed using the protein sequences annotated in the  
394 Swiss-Prot and TrEMBL databases hosted by the UniProt website (<https://www.uniprot.org/>).  
395 When a gene has multiple isoforms, we used its canonical sequence, as defined by UniProt  
396 based on criteria such as transcript length, relative abundance, and evolutionary conservation,  
397 in our analyses. We also repeated our analyses using transcript-level TPM data, where all  
398 isoforms annotated by ENSEMBL were included and had nearly identical results.

399

### 400 **Variation analysis of amino acid usage for TCGA samples**

401 To illustrate the variation of amino acid usage of NAT samples from different tissues, we  
402 computed z-scores based on the average frequencies for individual amino acids across tissues.  
403 To compare these with the variations in amino acid usage of tumor samples across cancer  
404 types, instead of using *de novo* standard deviations to compute z-scores, we used the set of  
405 standard deviations derived for the NAT samples to obtain z-scores for the tumor samples.  
406 We used hierarchically clustered heatmaps with Euclidean distance as the distance metric to  
407 visualize the tissue-specificity of amino acid usage. To identify differential amino acid usage  
408 between tumor and NAT samples, we performed the Wilcoxon rank-sum test for frequencies  
409 of individual amino acids using paired tumor and NAT samples and used an FDR-adjusted p-

410 value of 0.05 as the threshold for significance. Similarly, a hierarchically clustered heatmap  
411 was used to display amino acid de-regulation patterns across cancer types.

412

### 413 **Calculation of $ECPA_{gene}$ and $ECPA_{cell}$**

414 We calculated two indices of amino acid usage,  $ECPA_{gene}$ , and  $ECPA_{cell}$ , representing the  
415 average biosynthetic energy cost per amino acid of a gene and a cell, respectively, as  
416 described previously [27]. Briefly, the biosynthetic costs of amino acids are based on the  
417 amount of high-energy phosphate bond equivalents required for amino acid biosynthesis in  
418 yeast and are normalized by amino acid decay rates (the biosynthetic costs of amino acids are  
419 highly correlated between different species). We then calculated  $ECPA_{gene}$  and  $ECPA_{cell}$  by  
420 multiplying the biosynthetic energy costs with the relative amino acid frequency of a gene or  
421 a cell (sample).

422

### 423 **Quantification of amino acid usage convergence for TCGA samples**

424 To quantify the similarity of NAT or tumor samples in the TCGA cohort in terms of their  
425 amino acid usage patterns, we applied the Pearson's distance metric to the amino acid  
426 frequency profiles, derived as described above. We also employed the Spearman rank  
427 correlation coefficient as an alternative metric and obtained the same results. Specifically, to  
428 capture the convergent pattern of amino acid usage across cancer types, we defined, for a  
429 sample  $s_i$  of cancer type  $X$ , the amino acid usage convergence index as:

$$1 - \frac{\sum_{j=1}^N d_{s_i, s_j}}{N} (s_j \notin X)$$

430 where  $d_{s_i, s_j}$  is the Pearson's distance between sample  $s_i$  from cancer type  $X$  and sample  $s_j$   
431 not from cancer type  $X$ .

432

### 433 **Calculation of $\Delta ECPA_{cell}$ contribution index**

434 To estimate the contribution of individual genes to the alteration of  $ECPA_{cell}$  in a specific  
435 biological process, we considered both how different the  $ECPA$  of a gene is from the baseline  
436  $ECPA_{cell}$ , as well as how much its expression level has changed. Formally, we defined the  
437  $\Delta ECPA_{cell}$  contribution index of a gene  $g_i$  as:

$$(ECPA_{g_i} - ECPA_{baseline}) \times I_{g_i}$$

438 where  $I_{g_i}$  is an importance score that describes the extent of deregulation of  $g_i$ . In  
439 tumorigenesis, we employed the  $\log_2$  fold-change of average expression level between tumor

440 and NAT samples as the importance score. In tissue development, we employed a different  
441 importance score that was not based on binary comparison as in tumorigenesis since the  
442 nature of the dataset is time-course measurements. Specifically, we applied an R package  
443 designed for transcriptomic time courses, maSigPro [48], to build a polynomial regression  
444 model (degree = 3) for each gene using its expression level as the response variable and the  
445 log-transformed post-conception days as the independent variable. Such models yielded the  
446 goodness-of-fit ( $R^2$ ) values that were then signed by the corresponding Spearman correlation  
447 coefficients and were finally used as the importance score.

448

#### 449 **Pathway analysis of $\Delta ECPA_{\text{cell}}$ contribution in mammalian tissue development**

450 We employed an information-theoretic framework [31] to reveal gene modules or regulatory  
451 pathways that were enriched in genes with a significant contribution to the increase of  
452  $ECPA_{\text{cell}}$  during tissue development. First, we focused on down-regulated genes with lower-  
453 than-baseline  $ECPA_{\text{gene}}$  and up-regulated genes with higher-than-baseline  $ECPA_{\text{gene}}$ , both of  
454 which could contribute to the increase of developmental  $ECPA_{\text{cell}}$ . Second, we distinguished  
455 these two groups of genes by signing the index of down-regulated genes as negative,  
456 followed by rank-transforming all retained genes, and dividing the genes into equal bins.  
457 Third, we used the iPAGE algorithm that calculated the mutual information between the gene  
458 ranks and the pathway memberships (the number of genes belonging to a pathway in each bin)  
459 for every Gene Ontology term. A random-permutation test was used to estimate the  
460 significance of these mutual information (MI) values so that significantly informative  
461 pathways were identified with high MI values and low p values. Finally, the hypergeometric  
462 test was used to determine whether a specific pathway was over- or under-represented in each  
463 bin. For visualization, heatmaps of pathways by bins were drawn using log-transformed p  
464 values.

465

#### 466 **Calculation of developmental reversal index of tumor samples**

467 To assess the level of developmental reversal for tumor samples of TCGA LIHC, KIRC, and  
468 KIRP cohorts, we asked how greatly the shift of a tumor transcriptome from a mega NAT  
469 reference (averaging gene expressions over all NAT samples of a certain cancer type) had  
470 reversed the shift of the transcriptome along the developmental trajectory of a corresponding  
471 tissue. Formally, we defined, for a sample  $s_i$ , the developmental reversal index as:

$$\rho(\log_2(\vec{e}_{s_i} \oslash (E\vec{m}^{-1})), \vec{r})$$

472 where  $\oslash$  is element-wise division,  $\rho$  is the Spearman correlation coefficient,  $e_{s_i}$  is a vector of  
473  $n$  gene expressions for sample  $s_i$ ,  $E$  is a matrix of genes  $g_1, g_2, \dots, g_n$  by NAT samples  
474  $s_1, s_2, \dots, s_m$  of a certain cancer type with entries as expression level,  $\overrightarrow{m^{-1}}$  is a normalization  
475 vector of constant  $m^{-1}$ , and  $\vec{r}$  is a vector of signed goodness-of-fit values of genes  
476  $g_1, g_2, \dots, g_n$  derived from the developmental RNA-seq data of a matched tissue type. We  
477 examined the association of this index with patients' overall survival times in TCGA LIHC,  
478 KIRC, and KIRP cohorts using log-rank tests, where patients were split into two equal groups  
479 based on the median value of developmental reversal index.

480

### 481 **Evaluation of the utility of ECPA<sub>cell</sub> as a diagnostic biomarker**

482 To quantify the performance of ECPA<sub>cell</sub> in differentiating tumors from related normal  
483 samples, we used the AUROC metric to compare it with those of all detectable individual  
484 genes (TPM  $\geq 1$  in  $\geq 50\%$  of samples in the cohort). To determine the optimal threshold of  
485 ECPA<sub>cell</sub> or gene expression level for tumor-normal separation, we chose the value that  
486 maximizes Youden's J statistic, which equals to (sensitivity + specificity - 1). If multiple  
487 optimal cutoffs existed for a biomarker whose average level was higher in NAT than in  
488 tumors, the one with the highest value was picked and vice versa.

489

### 490 **CRedit author statement**

491 **Yikai Luo:** Conceptualization, Formal Analysis, Visualization, Writing – Original Draft.

492 **Han Liang:** Conceptualization, Supervision, Writing - Review & Editing, Funding  
493 acquisition. All authors read and approved the final manuscript.

494

### 495 **Competing interests**

496 H.L. is a shareholder and scientific advisor to Precision Scientific Ltd.

497

### 498 **Acknowledgements**

499 We thank H. Zhang, H. Chen, and other members of the Liang lab for helpful  
500 discussions. We thank H. Goodarzi for critical review of the manuscript. We also thank K.  
501 Mojumdar for editorial assistance. This work was supported by the US National Institutes of  
502 Health (U24CA209851 and Cancer Center Support Grant P30CA016672 to H.L.), an MD  
503 Anderson Faculty Scholar Award (to H.L.), and the Lorraine Dell Program in Bioinformatics  
504 for Personalization of Cancer Medicine (to H.L.).



505

506 **Authors' ORCID IDs**

507 0000-0001-7589-7981 (Yikai Luo)

508 0000-0001-7633-286X (Han Liang)

509

510

511 **References**

512 [1] Seligmann H. Cost-minimization of amino acid usage. *J Mol Evol* 2003;56:151–61.

513 <https://doi.org/10.1007/s00239-002-2388-z>.

514 [2] Heizer EM, Raiford DW, Raymer ML, Doom TE, Miller R V., Krane DE. Amino acid  
515 cost and codon-usage biases in 6 prokaryotic genomes: A whole-genome analysis. *Mol*  
516 *Biol Evol* 2006;23:1670–80. <https://doi.org/10.1093/molbev/msl029>.

517 [3] Raiford DW, Heizer EM, Miller R V., Akashi H, Raymer ML, Krane DE. Do amino  
518 acid biosynthetic costs constrain protein evolution in *Saccharomyces cerevisiae*? *J Mol*  
519 *Evol* 2008;67:621–30. <https://doi.org/10.1007/s00239-008-9162-9>.

520 [4] Harrison RJ, Charlesworth B. Biased gene conversion affects patterns of codon usage  
521 and amino acid usage in the *saccharomyces sensu stricto* group of yeasts. *Mol Biol*  
522 *Evol* 2011;28:117–29. <https://doi.org/10.1093/molbev/msq191>.

523 [5] Krick T, Verstraete N, Alonso LG, Shub DA, Ferreiro DU, Shub M, et al. Amino acid  
524 metabolism conflicts with protein diversity. *Mol Biol Evol* 2014.  
525 <https://doi.org/10.1093/molbev/msu228>.

526 [6] Wu C-I, Wang H-Y, Ling S, Lu X. The Ecology and Evolution of Cancer: The Ultra-  
527 Microevolutionary Process. *Annu Rev Genet* 2016;50:347–69.  
528 <https://doi.org/10.1146/annurev-genet-112414-054842>.

529 [7] Hanahan D, Weinberg RA. Hallmarks of cancer: The next generation. *Cell*  
530 2011;144:646–74. <https://doi.org/10.1016/j.cell.2011.02.013>.

531 [8] Miranda A, Hamilton PT, Zhang AW, Pattnaik S, Becht E, Mezheyski A, et al.  
532 Cancer stemness, intratumoral heterogeneity, and immune response across cancers.  
533 *Proc Natl Acad Sci U S A* 2019;116:9020–9. <https://doi.org/10.1073/pnas.1818210116>.

534 [9] Saygin C, Matei D, Majeti R, Reizes O, Lathia JD. Targeting Cancer Stemness in the  
535 Clinic: From Hype to Hope. *Cell Stem Cell* 2019;24:25–40.  
536 <https://doi.org/10.1016/j.stem.2018.11.017>.

537 [10] Milanovic M, Fan DNY, Belenki D, Däbritz JHM, Zhao Z, Yu Y, et al. Senescence-

- 538 associated reprogramming promotes cancer stemness. *Nature* 2018.  
539 <https://doi.org/10.1038/nature25167>.
- 540 [11] Peiris-Pagès M, Martinez-Outschoorn UE, Pestell RG, Sotgia F, Lisanti MP. Cancer  
541 stem cell metabolism. *Breast Cancer Res* 2016. [https://doi.org/10.1186/s13058-016-](https://doi.org/10.1186/s13058-016-0712-6)  
542 [0712-6](https://doi.org/10.1186/s13058-016-0712-6).
- 543 [12] Bellacosa A. Developmental disease and cancer: Biological and clinical overlaps. *Am*  
544 *J Med Genet Part A* 2013;161:2788–96. <https://doi.org/10.1002/ajmg.a.36267>.
- 545 [13] Aiello NM, Stanger BZ. Echoes of the embryo: Using the developmental biology  
546 toolkit to study cancer. *DMM Dis Model Mech* 2016;9:105–14.  
547 <https://doi.org/10.1242/dmm.023184>.
- 548 [14] Kamel HFM, Al-Amodi HSAB. Exploitation of Gene Expression and Cancer  
549 Biomarkers in Paving the Path to Era of Personalized Medicine. *Genomics, Proteomics*  
550 *Bioinforma* 2017;15:220–35. <https://doi.org/10.1016/j.gpb.2016.11.005>.
- 551 [15] Alkhateeb A, Rezaeian I, Singireddy S, Cavallo-Medved D, Porter LA, Rueda L.  
552 Transcriptomics Signature from Next-Generation Sequencing Data Reveals New  
553 Transcriptomic Biomarkers Related to Prostate Cancer. *Cancer Inform* 2019;18.  
554 <https://doi.org/10.1177/1176935119835522>.
- 555 [16] Rodon J, Soria JC, Berger R, Miller WH, Rubin E, Kugel A, et al. Genomic and  
556 transcriptomic profiling expands precision cancer medicine: the WINTHER trial. *Nat*  
557 *Med* 2019;25:751–8. <https://doi.org/10.1038/s41591-019-0424-4>.
- 558 [17] Aguet F, Brown AA, Castel SE, Davis JR, He Y, Jo B, et al. Genetic effects on gene  
559 expression across human tissues. *Nature* 2017;550:204–13.  
560 <https://doi.org/10.1038/nature24277>.
- 561 [18] Hoadley KA, Yau C, Hinoue T, Wolf DM, Lazar AJ, Drill E, et al. Cell-of-Origin  
562 Patterns Dominate the Molecular Classification of 10,000 Tumors from 33 Types of  
563 Cancer. *Cell* 2018;173:291-304.e6. <https://doi.org/10.1016/j.cell.2018.03.022>.
- 564 [19] Van Der Maaten L, Hinton G. Visualizing data using t-SNE. *J Mach Learn Res*  
565 2008;9:2579–625.
- 566 [20] Campbell PJ, Getz G, Korbel JO, Stuart JM, Jennings JL, Stein LD, et al. Pan-cancer  
567 analysis of whole genomes. *Nature* 2020;578:82–93. [https://doi.org/10.1038/s41586-](https://doi.org/10.1038/s41586-020-1969-6)  
568 [020-1969-6](https://doi.org/10.1038/s41586-020-1969-6).
- 569 [21] Robinson DR, Wu YM, Lonigro RJ, Vats P, Cobain E, Everett J, et al. Integrative  
570 clinical genomics of metastatic cancer. *Nature* 2017;548:297–303.  
571 <https://doi.org/10.1038/nature23306>.

- 572 [22] Hoadley KA, Yau C, Wolf DM, Cherniack AD, Tamborero D, Ng S, et al.  
573 Multiplatform analysis of 12 cancer types reveals molecular classification within and  
574 across tissues of origin. *Cell* 2014;158:929–44.  
575 <https://doi.org/10.1016/j.cell.2014.06.049>.
- 576 [23] Uhlen M, Zhang C, Lee S, Sjöstedt E, Fagerberg L, Bidkhorji G, et al. A pathology  
577 atlas of the human cancer transcriptome. *Science* (80- ) 2017;357.  
578 <https://doi.org/10.1126/science.aan2507>.
- 579 [24] Cardoso-Moreira M, Halbert J, Valloton D, Velten B, Chen C, Shao Y, et al. Gene  
580 expression across mammalian organ development. *Nature* 2019;571:505–9.  
581 <https://doi.org/10.1038/s41586-019-1338-5>.
- 582 [25] Aran D, Sirota M, Butte AJ. Systematic pan-cancer analysis of tumour purity. *Nat*  
583 *Commun* 2015;6:1–11. <https://doi.org/10.1038/ncomms9971>.
- 584 [26] Aran D, Camarda R, Odegaard J, Paik H, Oskotsky B, Krings G, et al. Comprehensive  
585 analysis of normal adjacent to tumor transcriptomes. *Nat Commun* 2017;8:1–13.  
586 <https://doi.org/10.1038/s41467-017-01027-z>.
- 587 [27] Zhang H, Wang Y, Li J, Chen H, He X, Zhang H, et al. Biosynthetic energy cost for  
588 amino acids decreases in cancer evolution. *Nat Commun* 2018;9:1–15.  
589 <https://doi.org/10.1038/s41467-018-06461-1>.
- 590 [28] Ellis MJ, Gillette M, Carr SA, Paulovich AG, Smith RD, Rodland KK, et al.  
591 Connecting genomic alterations to cancer biology with proteomics: The NCI clinical  
592 proteomic tumor analysis consortium. *Cancer Discov* 2013;3:1108–12.  
593 <https://doi.org/10.1158/2159-8290.CD-13-0219>.
- 594 [29] Jiang Y, Sun A, Zhao Y, Ying W, Sun H, Yang X, et al. Proteomics identifies new  
595 therapeutic targets of early-stage hepatocellular carcinoma. *Nature* 2019;567:257–61.  
596 <https://doi.org/10.1038/s41586-019-0987-8>.
- 597 [30] Gao Q, Zhu H, Dong L, Shi W, Chen R, Song Z, et al. Integrated Proteogenomic  
598 Characterization of HBV-Related Hepatocellular Carcinoma. *Cell* 2019;179:561-  
599 577.e22. <https://doi.org/10.1016/j.cell.2019.08.052>.
- 600 [31] Goodarzi H, Elemento O, Tavazoie S. Revealing Global Regulatory Perturbations  
601 across Human Cancers. *Mol Cell* 2009;36:900–11.  
602 <https://doi.org/10.1016/j.molcel.2009.11.016>.
- 603 [32] López-Otín C, Blasco MA, Partridge L, Serrano M, Kroemer G. The hallmarks of  
604 aging. *Cell* 2013;153:1194. <https://doi.org/10.1016/j.cell.2013.05.039>.
- 605 [33] Anglani R, Creanza TM, Liuzzi VC, Piepoli A, Panza A, Andriulli A, et al. Loss of

- 606 connectivity in cancer co-expression networks. *PLoS One* 2014;9.  
607 <https://doi.org/10.1371/journal.pone.0087075>.
- 608 [34] Han R, Huang G, Wang Y, Xu Y, Hu Y, Jiang W, et al. Increased gene expression  
609 noise in human cancers is correlated with low p53 and immune activities as well as  
610 late stage cancer. *Oncotarget* 2016;7:72011–20.  
611 <https://doi.org/10.18632/oncotarget.12457>.
- 612 [35] Yang RY, Quan J, Sodaei R, Aguet F, Segrè A V., Allen JA, et al. A systematic survey  
613 of human tissue-specific gene expression and splicing reveals new opportunities for  
614 therapeutic target identification and evaluation. *BioRxiv* 2018:311563.  
615 <https://doi.org/10.1101/311563>.
- 616 [36] Gould SJ. Ontogeny and phylogeny--revisited and reunited. *Bioessays* 1992;14:275–9.  
617 <https://doi.org/10.1002/bies.950140413>.
- 618 [37] Kalinka AT, Varga KM, Gerrard DT, Preibisch S, Corcoran DL, Jarrells J, et al. Gene  
619 expression divergence recapitulates the developmental hourglass model. *Nature*  
620 2010;468:811–6. <https://doi.org/10.1038/nature09634>.
- 621 [38] Domazet-Lošo T, Tautz D. A phylogenetically based transcriptome age index mirrors  
622 ontogenetic divergence patterns. *Nature* 2010;468:815–9.  
623 <https://doi.org/10.1038/nature09632>.
- 624 [39] Davies PCW, Lineweaver CH. Cancer tumors as Metazoa 1.0: Tapping genes of  
625 ancient ancestors. *Phys Biol* 2011;8. <https://doi.org/10.1088/1478-3975/8/1/015001>.
- 626 [40] Lineweaver CH, Davies PCW, Vincent MD. Targeting cancer's weaknesses (not its  
627 strengths): Therapeutic strategies suggested by the atavistic model. *BioEssays*  
628 2014;36:827–35. <https://doi.org/10.1002/bies.201400070>.
- 629 [41] Chen H, Lin F, Xing K, He X. The reverse evolution from multicellularity to  
630 unicellularity during carcinogenesis. *Nat Commun* 2015;6:1–9.  
631 <https://doi.org/10.1038/ncomms7367>.
- 632 [42] Trigos AS, Pearson RB, Papenfuss AT, Goode DL. Altered interactions between  
633 unicellular and multicellular genes drive hallmarks of transformation in a diverse range  
634 of solid tumors. *Proc Natl Acad Sci U S A* 2017;114:6406–11.  
635 <https://doi.org/10.1073/pnas.1617743114>.
- 636 [43] Trigos AS, Pearson RB, Papenfuss AT, Goode DL. Somatic mutations in early  
637 metazoan genes disrupt regulatory links between unicellular and multicellular genes in  
638 cancer. *Elife* 2019;8:1–28. <https://doi.org/10.7554/eLife.40947>.
- 639 [44] Li B, Ruotti V, Stewart RM, Thomson JA, Dewey CN. RNA-Seq gene expression

- 640 estimation with read mapping uncertainty. *Bioinformatics* 2009;26:493–500.  
641 <https://doi.org/10.1093/bioinformatics/btp692>.
- 642 [45] Ewels P, Magnusson M, Lundin S, Källér M. MultiQC: Summarize analysis results for  
643 multiple tools and samples in a single report. *Bioinformatics* 2016;32:3047–8.  
644 <https://doi.org/10.1093/bioinformatics/btw354>.
- 645 [46] Patro R, Duggal G, Love MI, Irizarry RA, Kingsford C. Salmon provides fast and bias-  
646 aware quantification of transcript expression. *Nat Methods* 2017;14:417–9.  
647 <https://doi.org/10.1038/nmeth.4197>.
- 648 [47] Sonesson C, Love MI, Robinson MD. Differential analyses for RNA-seq: Transcript-  
649 level estimates improve gene-level inferences [version 2; referees: 2 approved].  
650 *F1000Research* 2016;4:1521. <https://doi.org/10.12688/F1000RESEARCH.7563.2>.
- 651 [48] Nueda MJ, Tarazona S, Conesa A. Next maSigPro: Updating maSigPro bioconductor  
652 package for RNA-seq time series. *Bioinformatics* 2014;30:2598–602.  
653 <https://doi.org/10.1093/bioinformatics/btu333>.

654

## 655 **Figure legends**

### 656 **Figure 1 Pan-cancer convergence of transcriptome-based amino acid usage**

657 **A.** Schematic diagram showing the computation of amino acid usage frequency based on the  
658 gene expression profile derived from an RNA-seq sample. t-SNE projection of the GTE<sub>x</sub> (**B**),  
659 developing mammalian tissue (**C**), and TCGA tumor samples (**D**) based on their amino acid  
660 frequency profiles. Samples are color-coded based on tissue or cancer types. Marker shapes  
661 correspond to species. Developmental stages were classified into three categories and  
662 indicated by marker size. All t-SNE projections were generated using sklearn TSNE, with  
663 perplexity as 30, learning rate as 200, and the number of iterations as 1,000. Comparison of  
664 amino acid usage convergence index between tumor samples and either matched down-  
665 sampled normal samples (**E**) or adjacent normal samples (**F**) across multiple cancer types.  
666 Box plots show the quartiles, and the whiskers indicate quartile  $\pm 1.5 \times$  interquartile range. A  
667 two-sided Mann-Whitney U-test was used to calculate the p-value. \* $p < 0.05$ , \*\* $p < 0.01$ ,  
668 \*\*\* $p < 0.001$ .

669

### 670 **Figure 2 Amino acid usage preference in tumor evolution as quantified by ECPA<sub>cell</sub>**

671 **A.** Correlation between the biosynthetic energy cost of an amino acid and the net number of  
672 cancer types with significantly increased usage across 20 amino acids. The net number is

673 defined as the number of cancer types with significantly increased usage of the amino acid  
674 minus the number with significantly decreased usage. The colored region around the  
675 regression lines indicates a 95% confidence interval. **B.** Schematic diagram showing the  
676 computation of  $ECPA_{gene}$  and  $ECPA_{cell}$  based on the gene expression profile derived from  
677 RNA-seq data. **C.**  $ECPA_{cell}$  of tumor samples and matched normal tissue samples across  
678 TCGA cancer types. A paired two-sided Wilcoxon signed-rank test was used to calculate the  
679 p values. **D.** Bar plots showing  $ECPA_{gene}$  values of significantly down- and up-regulated  
680 proteins in several cancer proteomics datasets. Error bars denote 95% confidence intervals. A  
681 two-sided Mann–Whitney U-test was used to calculate the p values. **E.** Correlation between  
682  $ECPA_{cell}$  and amino acid usage convergence index across samples in nine cancer types. The  
683 colored regions around the regression lines indicate 95% confidence intervals. \* $p < 0.05$ , \*\* $p$   
684  $< 0.01$ , \*\*\* $p < 0.001$ .

685

### 686 **Figure 3 The increasing trend of $ECPA_{cell}$ throughout mammalian organogenesis**

687 Trend lines of  $ECPA_{cell}$  during the development of the liver (**A**), and the kidney (**B**) across  
688 five mammalian species. Developmental stages of non-mouse species correspond to the  
689 mouse stages shown in brackets. Error bars denote 95% confidence intervals. The trend line  
690 of  $ECPA_{cell}$  along the developmental trajectory of the mouse liver across three independent  
691 datasets (**C-E**). Error bars denote 95% confidence intervals. Heatmaps showing enrichment  
692 patterns of gene modules that contribute to  $\Delta ECPA_{cell}$  during the development of the human  
693 liver (**F**) and the human kidney (**G**). The red stripes embedded in the black background on  
694 top of each heatmap designate the range of  $\Delta ECPA_{cell}$  contribution index within every bin.

695

### 696 **Figure 4 A proposed model unifying developmental reversal, amino acid usage** 697 **convergence, and $ECPA_{cell}$ decline of cancer samples**

698 Stacked bar plots showing the proportion of genes that positively or negatively contribute to  
699  $\Delta ECPA_{cell}$  in either tumorigenesis or development for LIHC-liver (**A**), KIRC-kidney (**B**), and  
700 KIRP-kidney (**C**). Scatter plots showing, for genes with negative  $\Delta ECPA_{cell}$  contribution  
701 index in tumorigenesis and positive  $\Delta ECPA_{cell}$  contribution index in tissue development,  
702 scaled  $\Delta ECPA_{cell}$  contribution index in tumorigenesis versus scaled  $\Delta ECPA_{cell}$  contribution  
703 index in tissue development for LIHC-liver (**D**), KIRC-kidney (**E**), and KIRP-kidney (**F**).  
704 Colored regions around the regression lines indicate 95% confidence intervals. Kaplan-Meier  
705 plots show the overall survival for patients with LIHC (**G**), KIRC (**H**), or KIRP (**I**) stratified  
706 by developmental reversal index into two equal groups, respectively. The p values were

707 calculated from two-sided log-rank tests. Multivariate linear regression of  $ECPA_{cell}$  with  
708 developmental reversal index and amino acid usage convergence index as dependent  
709 variables for LIHC-liver (**J**), KIRC-kidney (**K**), and KIRP-kidney (**L**). **M**. Cartoon depicting  
710 a conceptual model in which cancer evolution is accompanied by the convergence of amino  
711 acid usage and decrease of  $ECPA_{cell}$ , which is a reversal of the tissue development process.

712

### 713 **Figure 5 The liver shows the most dramatic $ECPA_{cell}$ reduction in tumorigenesis**

714 Distributions of  $\Delta ECPA_{cell}$  between tumor samples and paired NAT samples across multiple  
715 cancer types (**A**), tissue-specific genes-based  $ECPA_{cell}$  of normal samples across multiple  
716 tissues (**B**), tissue-specific genes-based  $ECPA_{cell}$  of adjacent normal samples across multiple  
717 cancer types (**C**), ranked by the median values. The box plots show the quartiles. The  
718 whiskers indicate quartile  $\pm 1.5 \times$  interquartile range. The horizontal dashed line indicates the  
719 level of  $\Delta ECPA_{cell} = 0$ . **D**. Trend lines of  $ECPA_{cell}$  of multiple tissues across human  
720 developmental stages. Error bars denote 95% confidence interval. wpc, weeks post  
721 conception.

722

### 723 **Figure 6 $ECPA_{cell}$ is a robust diagnostic biomarker for liver cancer**

724 **A**.  $ECPA_{cell}$  of tumor samples and matched normal tissue samples in 11 independent RNA-  
725 seq datasets of liver cancer and their matched normal samples. A paired two-sided Wilcoxon  
726 signed-rank test was used to calculate the p values. **B**. ROC curves of  $ECPA_{cell}$  as a  
727 diagnostic biomarker in six independent liver cancer cohorts with sample size  $\geq 12$ . Dashed  
728 lines indicate the lines of identity. ROC, receiver operating characteristic; AUC, area under  
729 the ROC curve. **C**. Histogram showing the distribution of the average AUC across the six  
730 cohorts for tumor-normal segregation using the mRNA expression level of each of the 9,559  
731 detectable genes. The vertical dashed line corresponds to the average AUC of  $ECPA_{cell}$ . **D**.  
732 Box plots showing the AUC of the top four metrics, including three genes and  $ECPA_{cell}$ , in  
733 discriminating tumor samples from normal samples across the six cohorts. A paired two-  
734 sided Wilcoxon signed-rank test was used to calculate the p values. **E**. Box plots showing the  
735 AUC of  $ECPA_{cell}$  and the frequency of each amino acid in detecting tumors across the six  
736 cohorts. The box plots show the quartiles. The whiskers indicate quartile  $\pm 1.5 \times$  interquartile  
737 range. A paired two-sided Wilcoxon signed-rank test was used to calculate the p values. **F**.  
738 Histogram showing the distribution of coefficients of variation (CV) of the optimal  
739 thresholds in using individual genes for tumor-normal segregation. The vertical red dashed

740 line indicates the CV of  $ECPA_{cell}$ . Vertical lines in three other colors indicate the CV of three  
741 genes whose average AUCs are higher than  $ECPA_{cell}$ . \* $p < 0.05$ , \*\* $p < 0.01$ , \*\*\* $p < 0.001$ .  
742  
743  
744



745 **Supplementary materials**

746 **Figure S1 t-SNE projection of samples based on gene expression**

747 t-SNE projection of the GTEx (A), TCGA (B), PCAWG (C), and MET500 (D) samples  
748 based on their gene expression profiles.

749

750 **Figure S2 t-SNE projection of samples based on amino acid frequency**

751 t-SNE projection of the GTEx & HPA (A), PCAWG (B), MET500 (C), down-sampled GTEx  
752 (D), down-sampled TCGA tumor (E), matched TCGA NAT (F), and matched TCGA tumor  
753 (G) samples based on their amino acid frequency.

754

755 **Figure S3 Re-calculation of  $ECPA_{cell}$  in TCGA samples without highly-expressed genes**

756 A.  $ECPA_{cell}$  of tumor samples and matched normal tissue samples across TCGA cancer types  
757 after removal of genes encoding high-abundance housekeeping or tissue-specific proteins

758

759 **Figure S4 Differential amino acid usage within and between tumor and NAT samples**

760 Heatmaps showing the average frequency of individual amino acids for NAT samples (A)  
761 and tumor samples (B), normalized as z-scores, across 15 cancer types. C. Heatmap showing  
762 significantly increased or decreased usage of individual amino acids between tumor and NAT  
763 samples across 15 cancer types.

764

765 **Figure S5 Functional enrichment of genes with a positive contribution to  $\Delta ECPA_{cell}$  in  
766 non-human tissue development**

767 Heatmaps showing enrichment patterns of well-defined gene modules that contribute to  
768  $ECPA_{cell}$  increase during the development of the mouse liver (A) and kidney (B), the rat liver  
769 (C) and kidney (D), and the rabbit liver (E) and kidney (F). Red stripes embedded in the  
770 black background on top of each heatmap designate the range of  $\Delta ECPA_{cell}$  contribution  
771 index within every bin.

772

773 **Figure S6 Variations of  $ECPA_{cell}$  during tissue aging**

774 Trend lines of  $ECPA_{cell}$  during aging of the liver and the kidney in humans (A), mice (B), and  
775 rats (C). Blue and orange dashed lines indicate the average levels of  $ECPA_{cell}$  across age  
776 groups for the liver and the kidney, respectively. Error bars denote 95% confidence intervals.

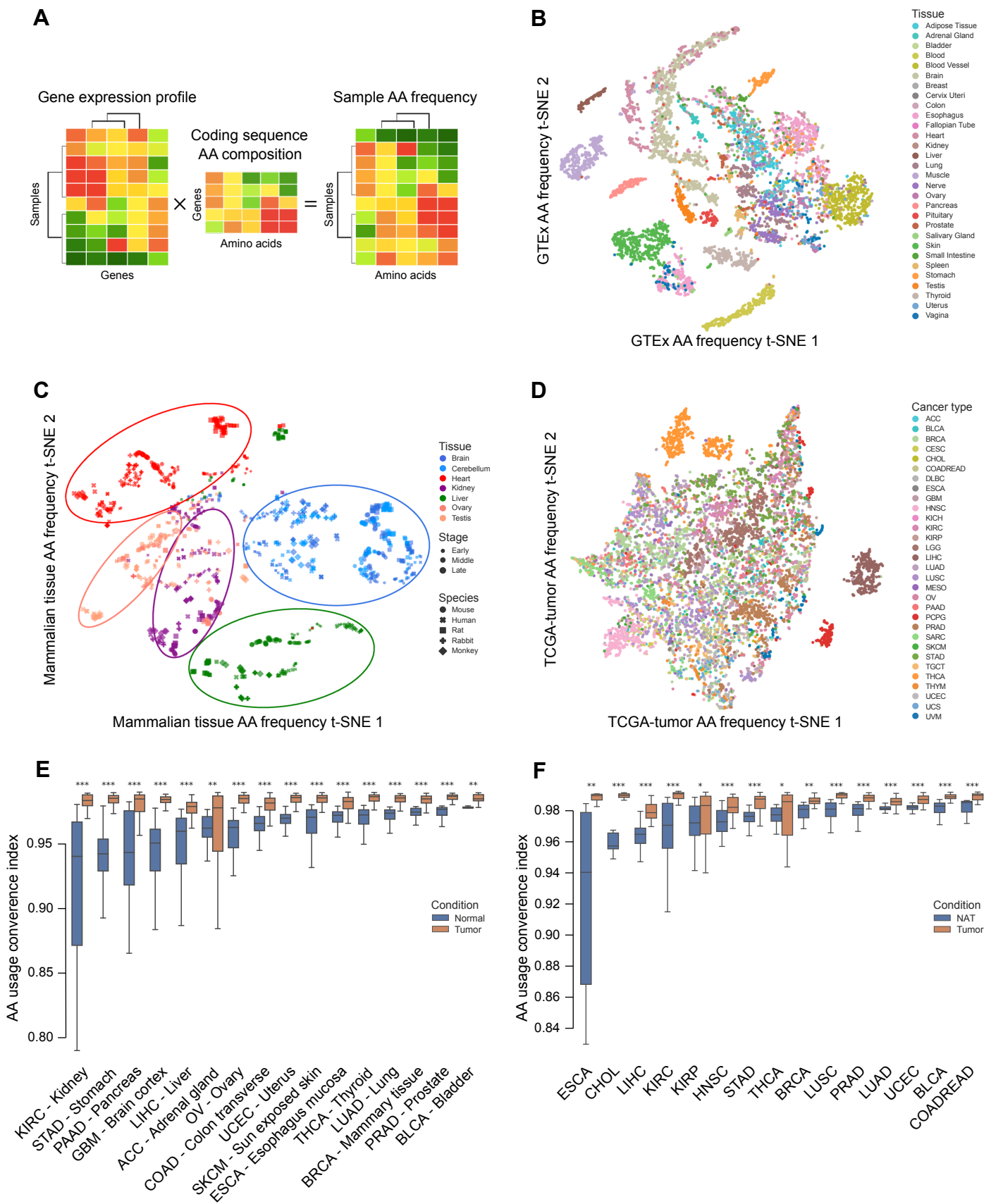
777

778 **Figure S7 Variations of  $ECPA_{cell}$  during the development of multiple tissues in non-**  
779 **human mammals**

780 Trend lines of  $ECPA_{cell}$  in seven tissues across four mammals, including mice (**A**), rats (**B**),  
781 rabbits (**C**), and opossums (**D**). Error bars denote 95% confidence intervals.

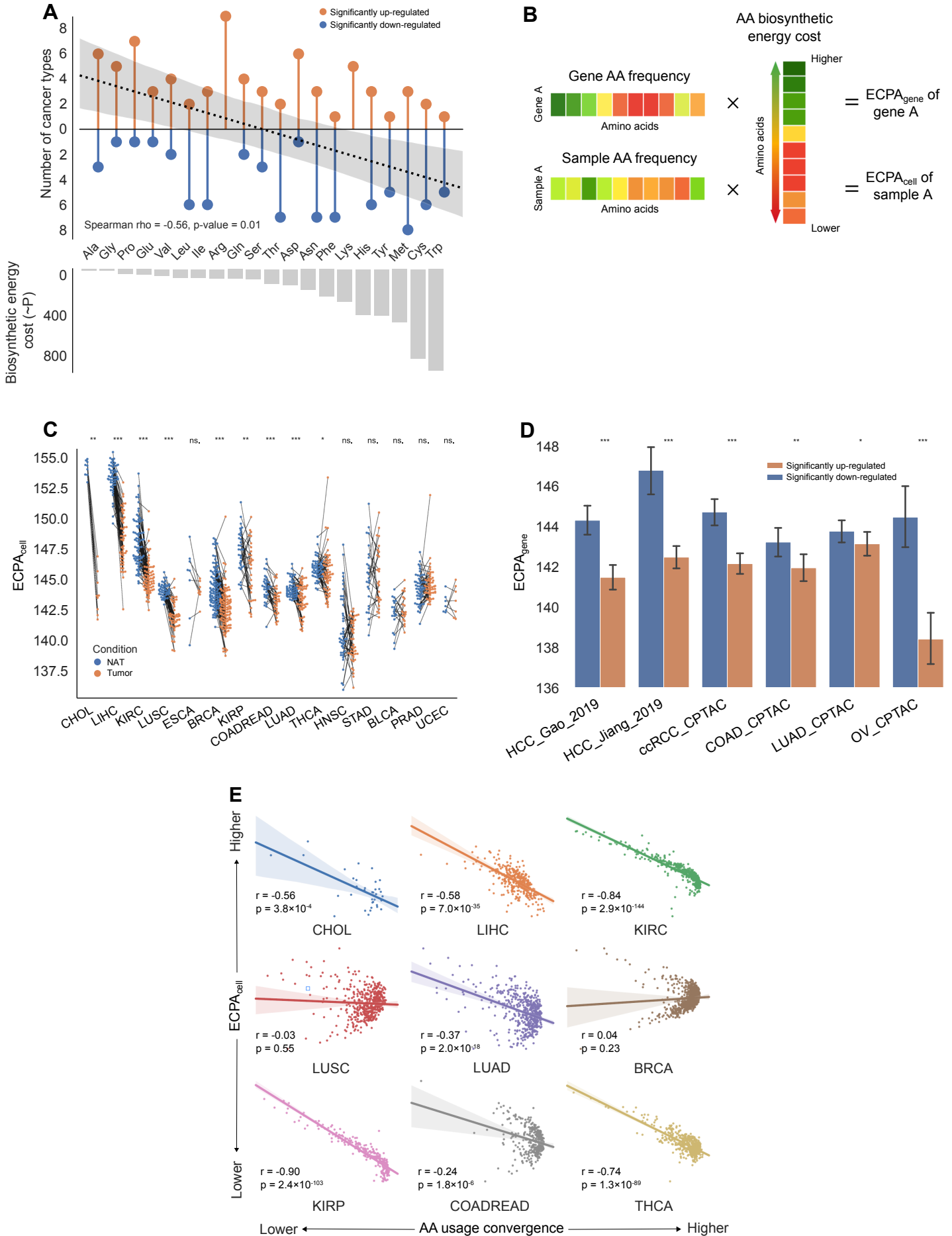
# Figure 1

bioRxiv preprint doi: <https://doi.org/10.1101/2021.03.18.436083>; this version posted March 20, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.

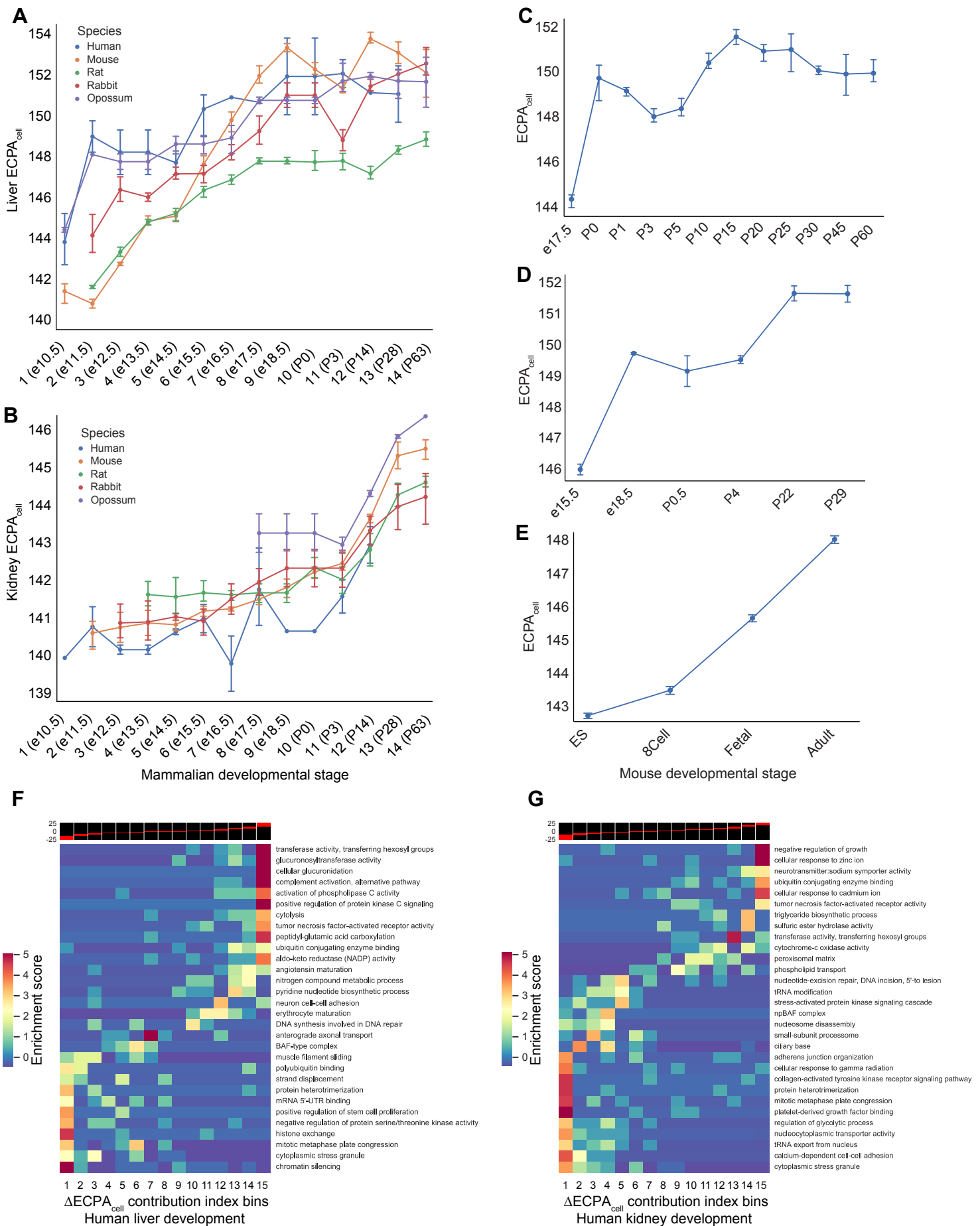


# Figure 2

bioRxiv preprint doi: <https://doi.org/10.1101/2021.03.18.436083>; this version posted March 20, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.

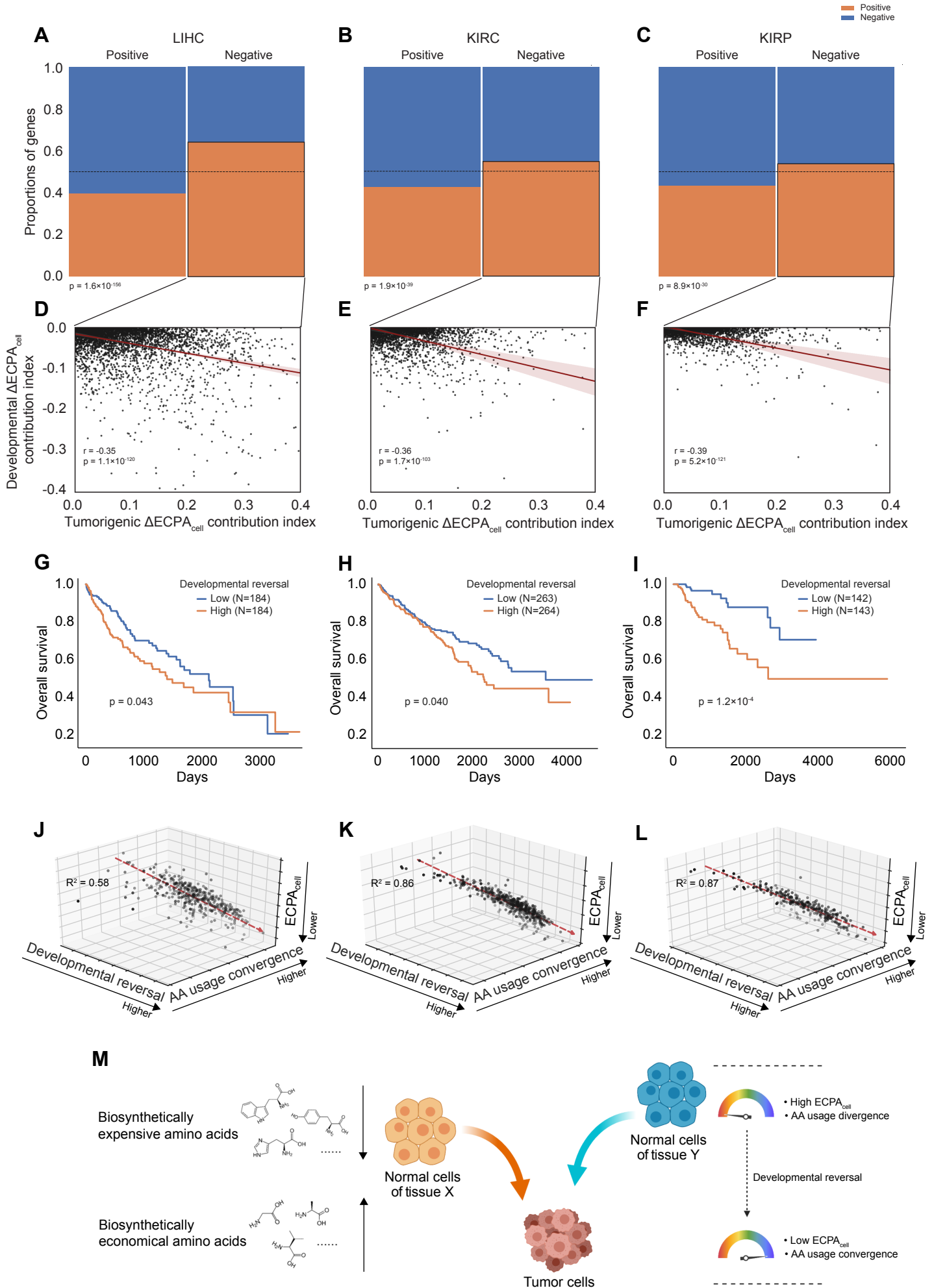


**Figure 3** | bioRxiv preprint doi: <https://doi.org/10.1101/2021.03.18.436083>; this version posted March 20, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.



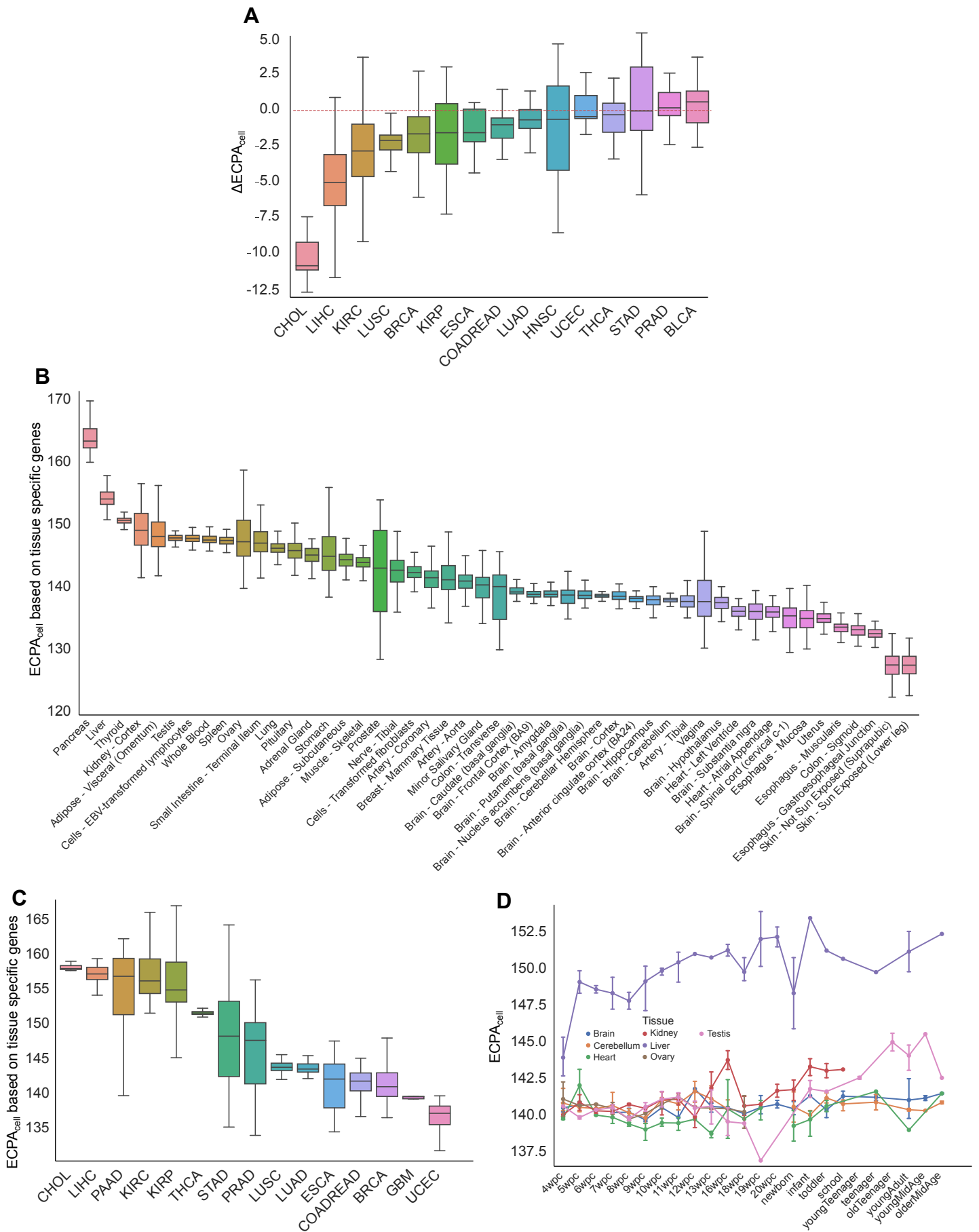
# Figure 4

bioRxiv preprint doi: <https://doi.org/10.1101/2021.03.18.436083>; this version posted March 20, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.

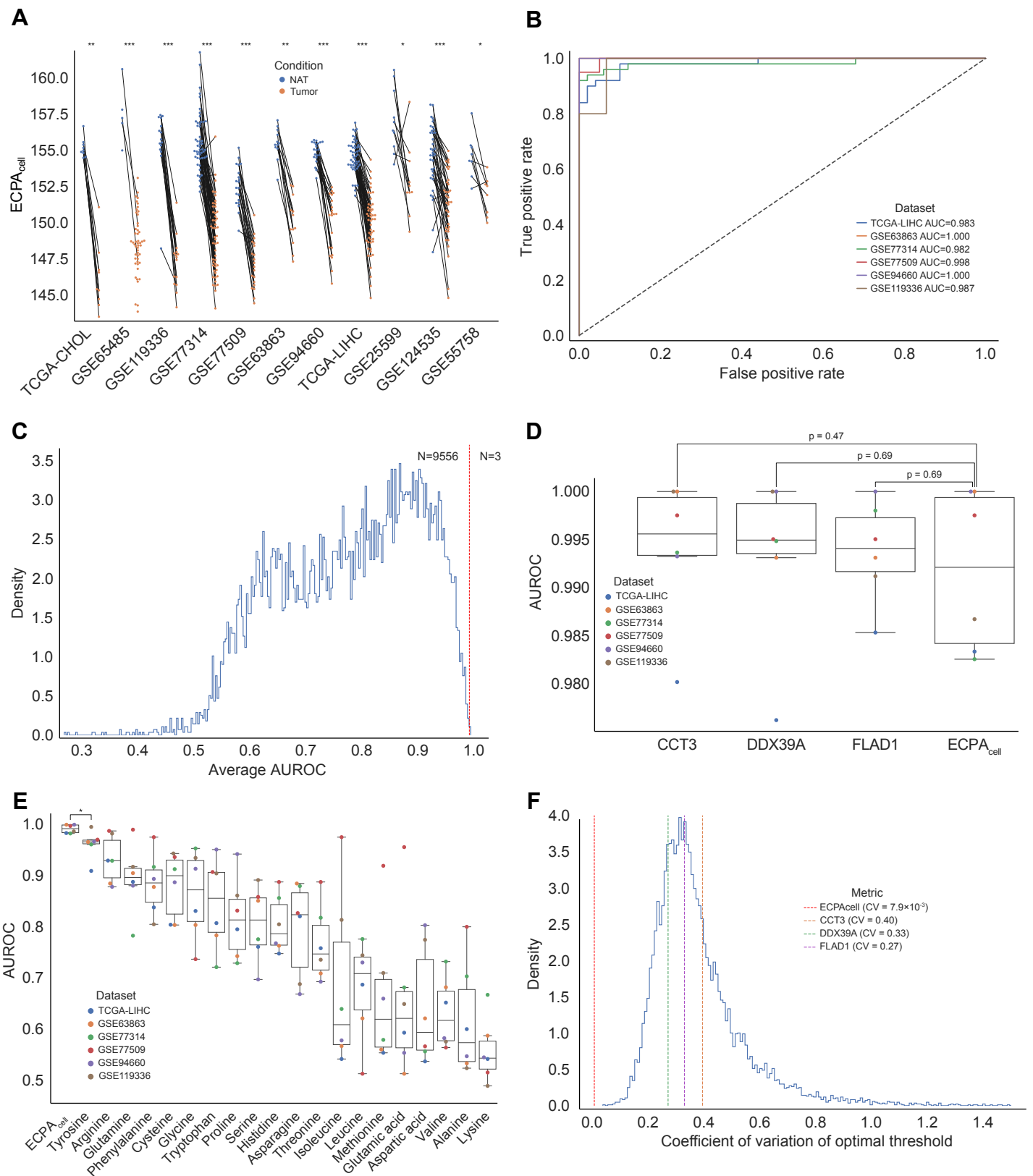


# Figure 5

bioRxiv preprint doi: <https://doi.org/10.1101/2021.03.18.436083>; this version posted March 20, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.

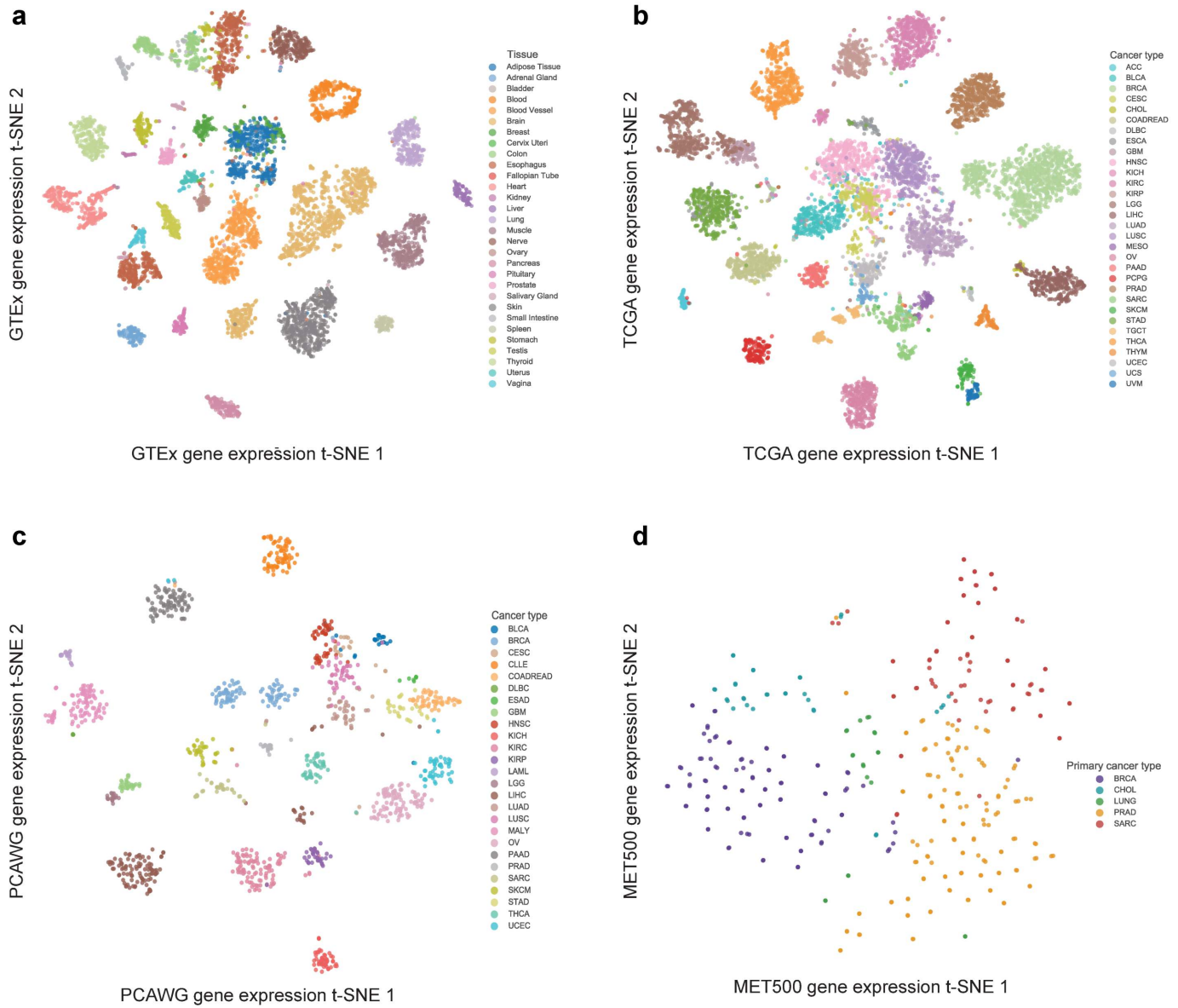


**Figure 6** bioRxiv preprint doi: <https://doi.org/10.1101/2021.03.18.436083>; this version posted March 20, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.

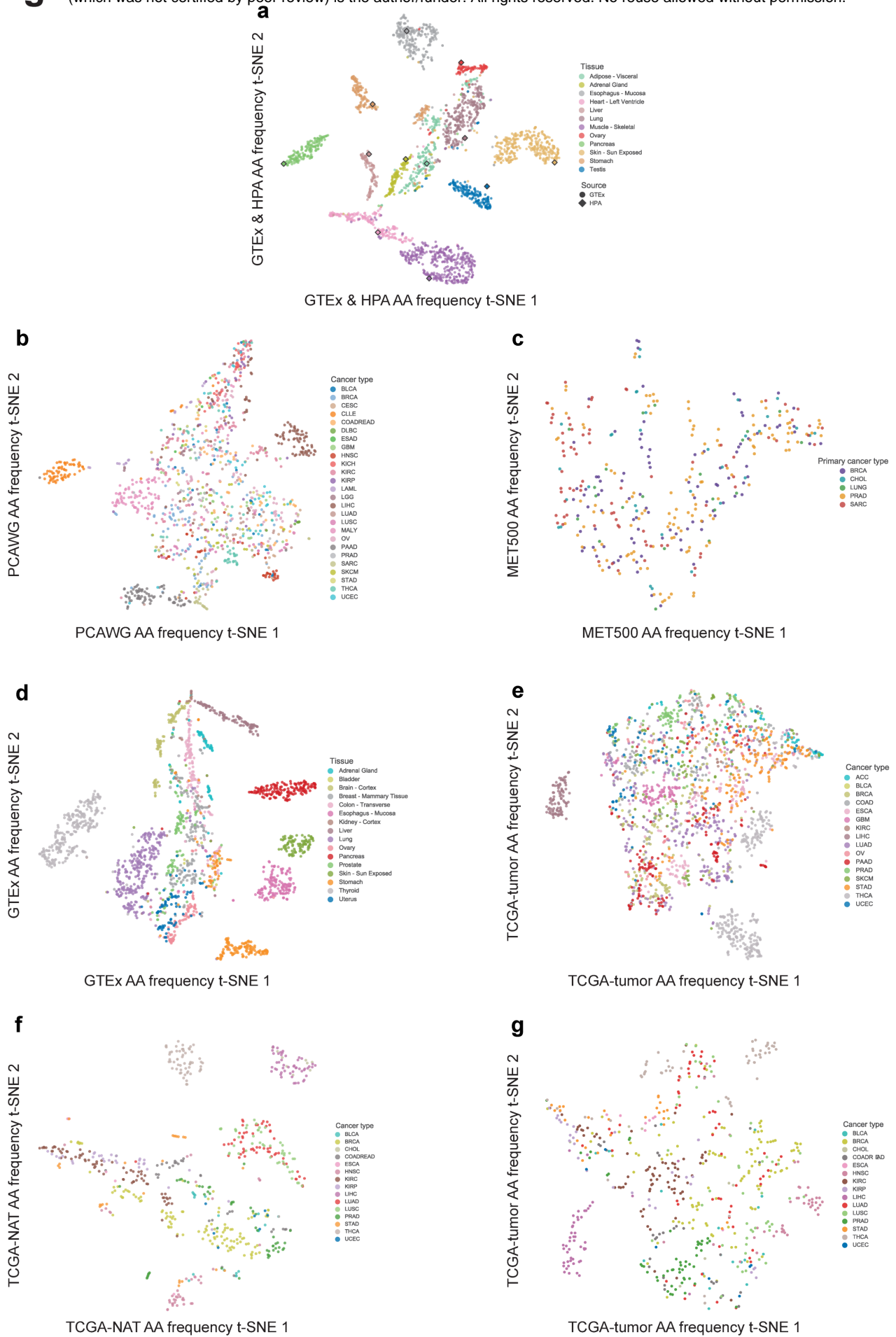




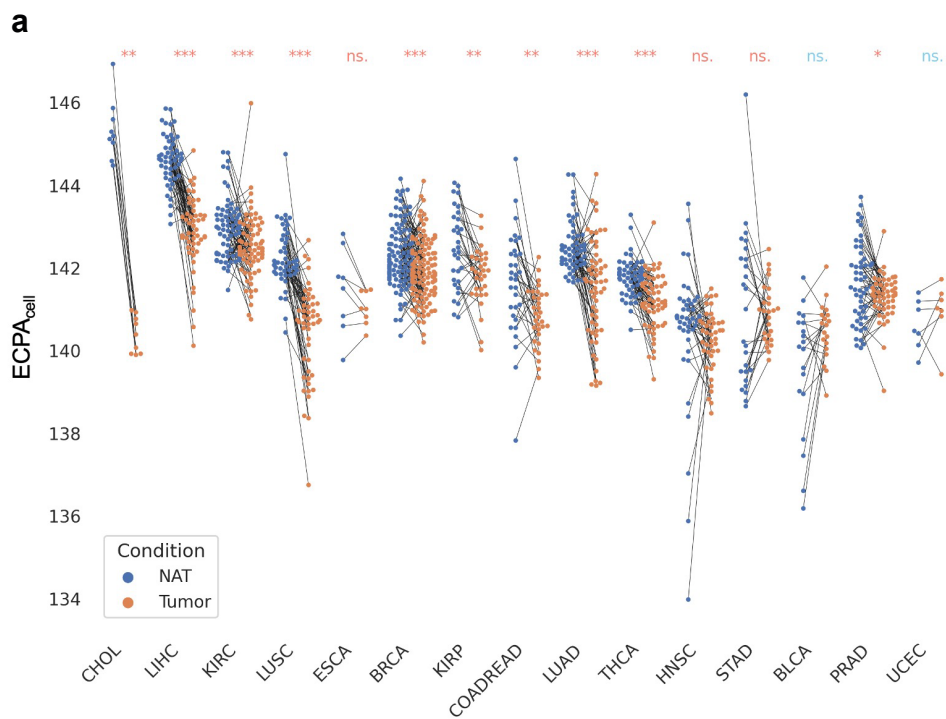
# Figure S1



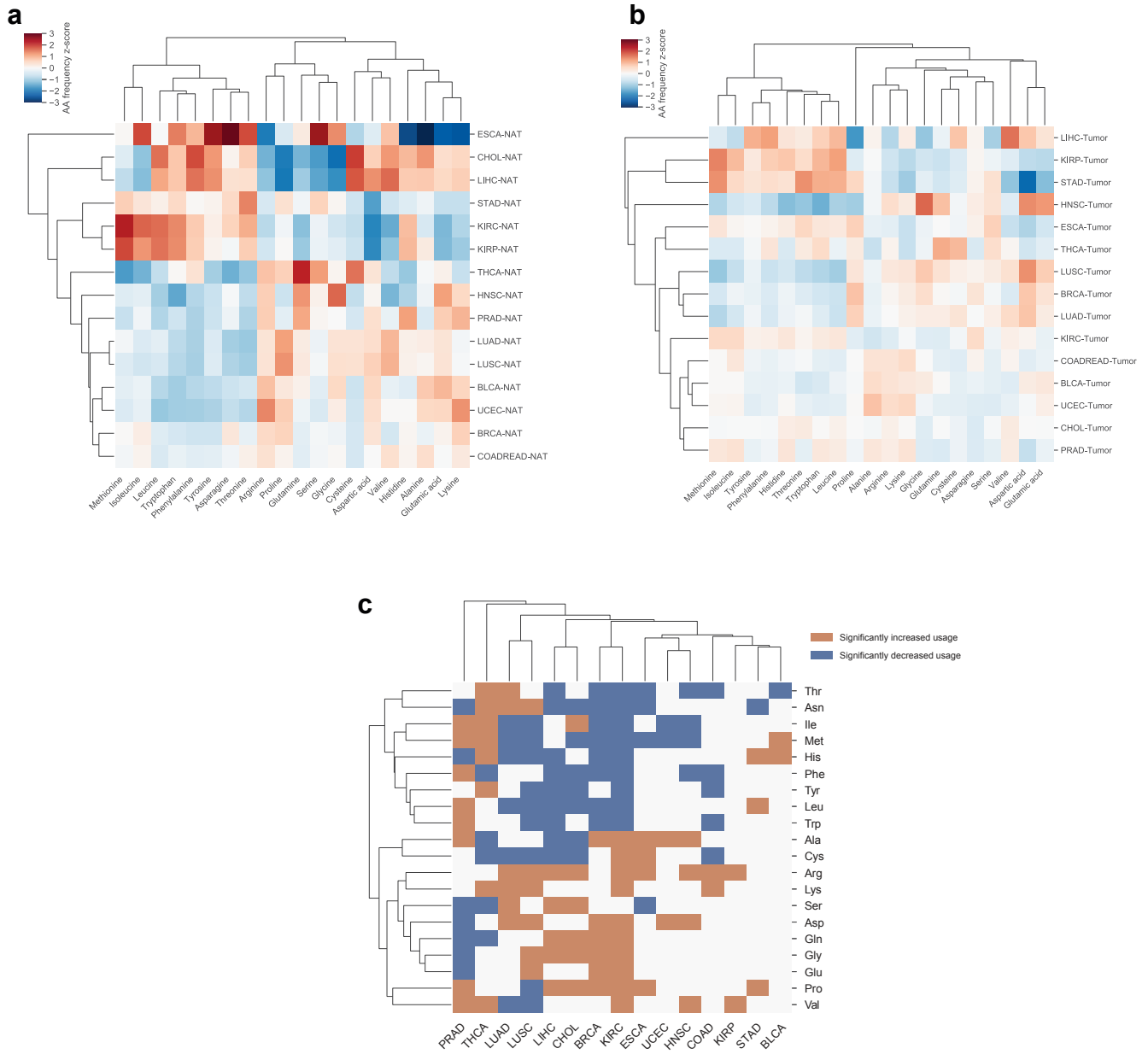
# Figure S2



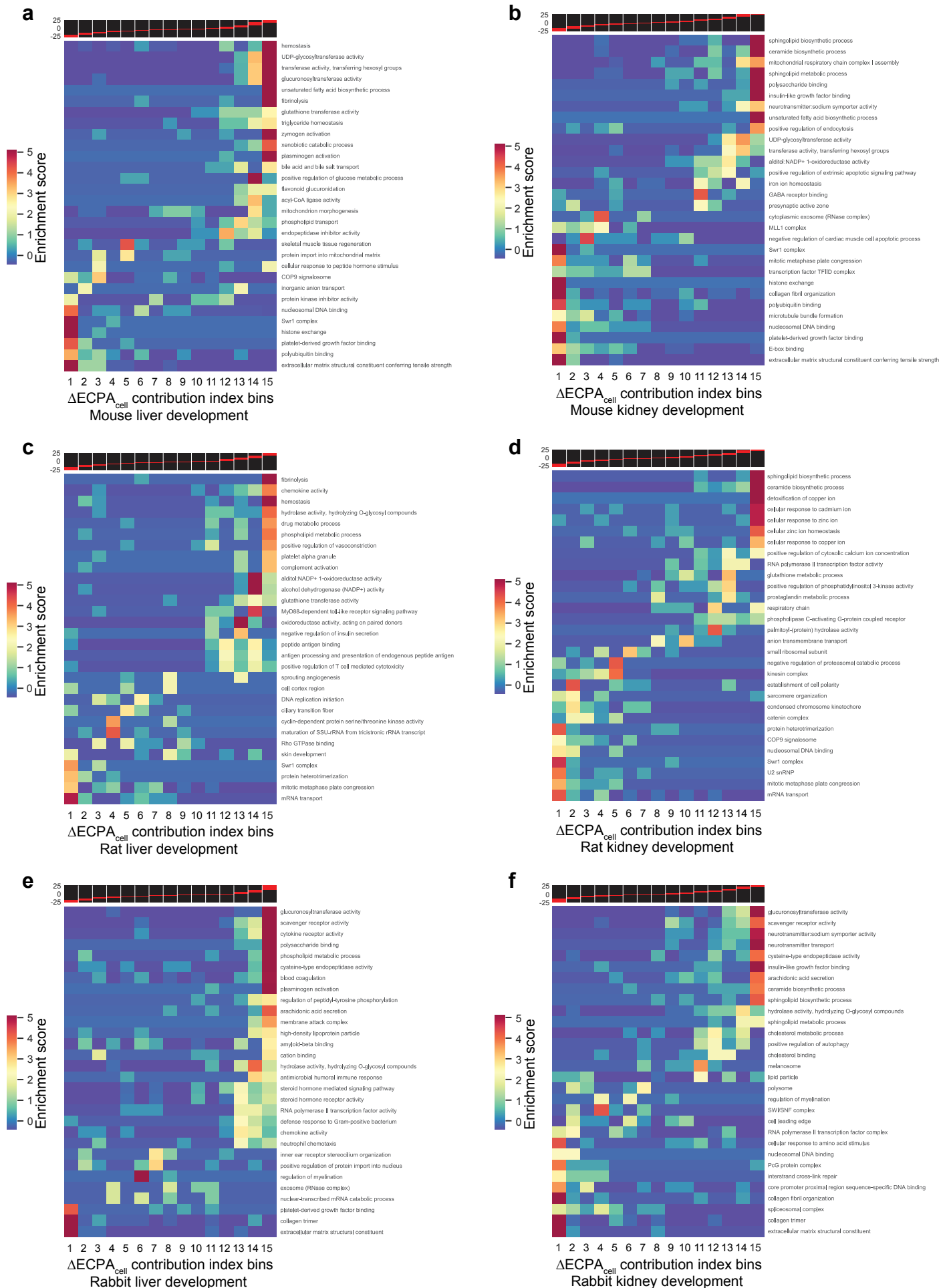
# Figure S3



# Figure S4



# Figure S5



# Figure S6

