



# Detection of isoforms and genomic alterations by high-throughput full-length single-cell RNA sequencing for personalized oncology

## Working Paper

### Author(s):

Dondi, Arthur; Lischetti, Ulrike; Jacob, Francis; Singer, Franziska; [Borgsmüller, Nico](#) ; Tumor Profiler Consortium; Heinzelmann-Schwarz, Viola; Beisel, Christian; [Beerenwinkel, Niko](#) 

### Publication date:

2022-12-15

### Permanent link:

<https://doi.org/10.3929/ethz-b-000594543>

### Rights / license:

[Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International](#)

### Originally published in:

bioRxiv, <https://doi.org/10.1101/2022.12.12.520051>

### Funding acknowledgement:

190413 - Single-cell transcript isoform sequencing for precision oncology diagnostics and treatment (SNF)  
766030 - Computational ONcology TRaining Alliance (EC)

1 Detection of isoforms and genomic alterations by high-  
2 throughput full-length single-cell RNA sequencing for  
3 personalized oncology

4

5

6 **Arthur Dondi<sup>1,2,\*</sup>, Ulrike Lischetti<sup>3,\*,\*\*</sup>, Francis Jacob<sup>3</sup>, Franziska Singer<sup>2,4</sup>, Nico**  
7 **Borgsmüller<sup>1,2</sup>, Tumor Profiler Consortium, Viola Heinzelmann-Schwarz<sup>3</sup>, Christian**  
8 **Beisel<sup>1,\*\*</sup>, Niko Beerenwinkel<sup>1,2,\*\*</sup>**

9

10 \* Equal contributions

11 \*\* Corresponding authors

12

13 <sup>1</sup> ETH Zurich, Department of Biosystems Science and Engineering, Mattenstrasse 26, 4058  
14 Basel, Switzerland

15 <sup>2</sup> SIB Swiss Institute of Bioinformatics, Mattenstrasse 26, 4058 Basel, Switzerland

16 <sup>3</sup> University Hospital Basel and University of Basel, Ovarian Cancer Research, Department of  
17 Biomedicine, Hebelstrasse 20, 4031 Basel, Switzerland

18 <sup>4</sup> ETH Zurich, NEXUS Personalized Health Technologies, Wagistrasse 18, 8952 Schlieren,  
19 Switzerland

20

21

22

23

24 **Keywords:** full-length single-cell RNA sequencing, long-read PacBio sequencing, transcript  
25 concatenation, isoforms, mutations, gene fusions, ovarian cancer

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

## 42 Abstract

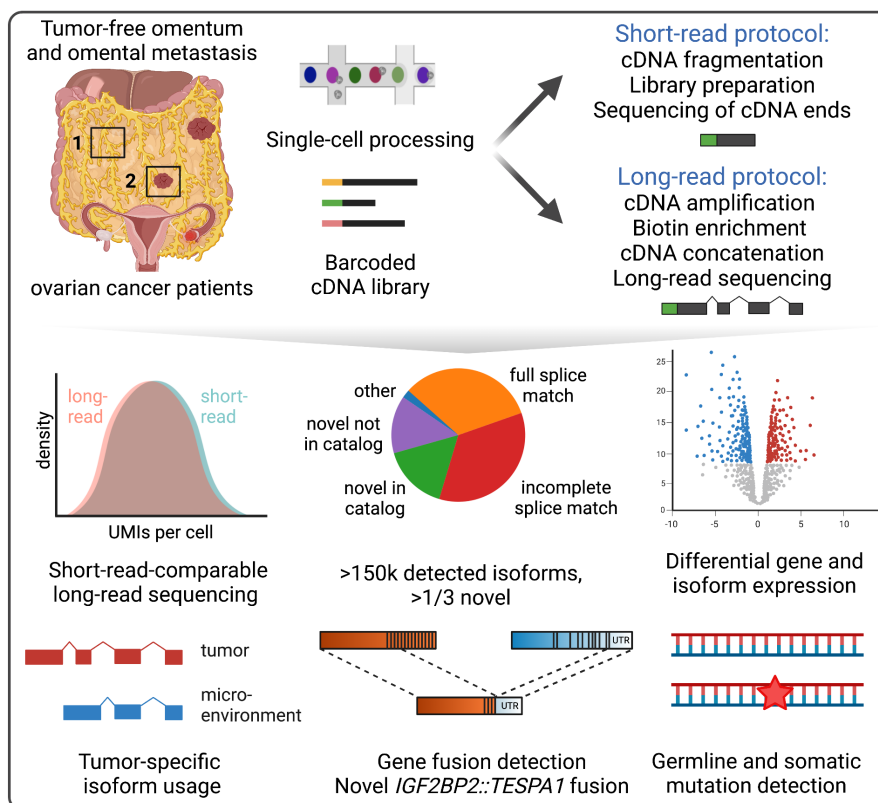
43

44 Understanding the complex background of cancer requires genotype-phenotype information  
45 in single-cell resolution. Long-read single-cell RNA sequencing (scRNA-seq), capturing full-  
46 length transcripts, lacked the depth to provide this information so far. Here, we increased the  
47 PacBio sequencing depth to 12,000 reads per cell, leveraging multiple strategies, including  
48 artifact removal and transcript concatenation, and applied the technology to samples from  
49 three human ovarian cancer patients. Our approach captured 152,000 isoforms, of which over  
50 52,000 were novel, detected cell type- and cell-specific isoform usage, and revealed  
51 differential isoform expression in tumor and mesothelial cells. Furthermore, we identified gene  
52 fusions, including a novel scDNA sequencing-validated *IGF2BP2::TESPA1* fusion, which was  
53 misclassified as high *TESPA1* expression in matched short-read data, and called somatic and  
54 germline mutations, confirming targeted NGS cancer gene panel results. With multiple new  
55 opportunities, especially for cancer biology, we envision long-read scRNA-seq to become  
56 increasingly relevant in oncology and personalized medicine.

57

58

59 Graphical Abstract:



60

61

62

63

64

## 65 Introduction

66 Cancer is a complex disease characterized by genomic and transcriptomic alterations<sup>1</sup> that  
67 drive multiple tumor-promoting capabilities or hallmarks<sup>2</sup>. Among others, these alterations  
68 include point mutations, insertions and deletions (indels), and gene fusions on the genomic  
69 level, and splice isoforms on the transcriptomic level. Their detection offers great potential for  
70 personalized oncology as they can serve as direct therapeutic targets<sup>3,4</sup> or potential  
71 neoantigens informing on the immunogenicity of the tumor<sup>5</sup>. Gene fusions arising from large-  
72 scale genomic rearrangements, for example, play an oncogenic role in a variety of tumor  
73 types<sup>6</sup>, and are successfully used as therapeutic targets<sup>7,8</sup>. Like mutations<sup>9</sup> and copy number  
74 variations<sup>10</sup>, fusion rates can vary widely across cancer types, and gene fusions are thought  
75 to be drivers in 16.5% of cancer cases, and even the only driver in more than 1%<sup>11</sup>.  
76 Furthermore, out-of-frame gene fusions are more immunogenic than mutations and indels,  
77 making them an ideal target for immunotherapies and cancer vaccines<sup>12,13</sup>. On the  
78 transcriptomic level, alternative splicing is a major mechanism for the diversification of a cell's  
79 transcriptome and proteome<sup>14</sup> and can impact all hallmarks of tumorigenesis. It also presents  
80 a fairly novel non-genomic source of potential neoantigens<sup>15</sup>. In breast and ovarian cancer,  
81 68% of samples had at least one isoform with novel exon-exon junction (neojunction) detected  
82 in proteomic data<sup>16</sup>.

83 The complexity of cancer further extends to intra-tumor heterogeneity<sup>17</sup> and its intricate  
84 interplay with the tumor microenvironment (TME)<sup>18</sup>. Ultimately, to fully decipher functional  
85 tumor heterogeneity and its effect on the TME, single-cell resolution providing both phenotype  
86 and genotype information is required. Single-cell RNA sequencing (scRNA-seq) is now widely  
87 used for the phenotypic dissection of heterogeneous tissues. It can be divided into short-read,  
88 high-throughput technologies allowing for gene expression quantification and long-read, low-  
89 throughput technologies that cover full-length transcripts<sup>19</sup>. Up to now, short- and long-read  
90 methods had to be used in parallel to combine the advantages of each technology. The long-  
91 read scRNA-seq field is rapidly expanding, with methods being constantly developed and  
92 improved on Nanopore<sup>20,21</sup> and PacBio<sup>22-26</sup> long-read platforms. So far, long-read RNA-seq  
93 has however only been applied on the bulk level in the field of oncology<sup>24,27,28</sup>. High-quality,  
94 high-throughput, long-read scRNA-seq has the potential to provide isoform-level cell type-  
95 specific readouts and capture tumor-specific genomic alterations. With near ubiquitous p53  
96 mutations and defective DNA repair pathways causing frequent non-recurrent gene fusions,  
97 high-grade serous ovarian cancer (HGSOC) is an ideal candidate to investigate these  
98 alterations<sup>10,29,30</sup>.

99 Here, for the first time, we used high-quality, high-throughput long-read scRNA-seq to capture  
100 cell type-specific genomic and transcriptomic alterations on clinical cancer patients. We  
101 applied both short-read and long-read scRNA-seq to five samples from three HGSOC patients,  
102 comprising 2,571 cells, and generated the largest PacBio scRNA-seq dataset to date. We  
103 were able to identify over 150,000 isoforms, of which a third were novel, as well as novel cell  
104 type- and cell-specific isoforms. We detected differential isoform usage in tumor cells and cells  
105 of the TME. Additionally, we discovered dysregulations in the insulin like growth factor (IGF)  
106 network in tumor cells on the genomic and transcriptomic level. Thereby, we demonstrated  
107 that scRNA-seq can capture genomic alterations accurately, including cancer- and patient-  
108 specific germline and somatic mutations in genes such as *TP53*, as well as gene fusions,  
109 including a novel *IGF2BP2::TESPA1* fusion.

## 110 Results

### 111 Long-read scRNA-seq creates a catalog of isoforms in ovarian cancer 112 patient-derived tissue samples

113 We generated short-read and long-read scRNA-seq data from five omentum biopsy samples  
114 (**Extended Data Table 1**) from three HGSOC patients. Three samples were derived from  
115 HGSOC omental metastases and two from matching distal tumor-free omental tissues  
116 (**Fig. 1a**). To generate long reads, we opted for the PacBio platform for its generation of high-  
117 fidelity (HiFi) reads through circular consensus sequencing (CCS). To overcome its limitations  
118 in sequencing output and optimize for longer library length, we 1) removed template-switch  
119 oligo artifacts that can account for up to 50% of reads through biotin enrichment, 2)  
120 concatenated transcripts to sequence multiple cDNA molecules per CCS read, and 3)  
121 sequenced on the PacBio Sequel II platform (2-4 SMRT 8M cells per sample, **Methods**). This  
122 allowed the generation of a total of 212 Mio HiFi reads in 2,571 cells, which, after  
123 demultiplexing, deduplication, and intrapriming removal, resulted in 30.7 Mio unique molecular  
124 identifiers (UMIs) (**Extended Data Table 1**). On average, 12k UMIs were detected per cell.

125  
126 The long-read dataset revealed 152,546 isoforms, each associated with at least three UMIs.  
127 We classified the isoforms according to the SQANTI classification<sup>31</sup> and calculated their  
128 proportions (**Methods, Fig. 1b,c**): full splice match (FSM) - isoforms already in the GENCODE  
129 database (32.8%), incomplete splice matches (ISM) - isoforms corresponding to shorter  
130 versions of the FSM (35.1%), novel in catalog (NIC) - isoforms presenting combinations of  
131 known splice donors and acceptors (15.9%), and novel not in catalog (NNC) - isoforms  
132 harboring at least one unknown splice site, or neojunction (14.4%). Novel isoforms (classes  
133 NIC and NNC) accounted for 30% of the isoforms, and 11% of the total reads in all samples,  
134 while FSM accounted for 33% of the isoforms and 80% of the reads (**Fig. 1c,d**), indicating that  
135 high coverage is required for the reliable detection of new, low abundant, transcripts.

136  
137 To evaluate the structural integrity of all isoforms, we compared their 5' end to the FANTOM5  
138 CAGE database<sup>32</sup> and their 3' end to the PolyASite database<sup>33</sup> (**Fig. 1e**). More than 82% of  
139 the NIC and 74% of NNC isoforms could be validated on 3' and 5' ends, similarly to FSM. As  
140 expected, fewer ISM isoforms were found to be complete (42%): they are either incompletely  
141 sequenced isoforms missing their 5' end (30%) or the result of early 3' termination (55%).

142  
143 FSM, NIC, and NNC had overall better 3' and 5' validation than the full-length tagged isoforms  
144 in the GENCODE database (**Fig. 1e**). Only the 'Matched Annotation from NCBI and EMBL-  
145 EBI' (MANE<sup>34</sup>) containing curated representative transcripts cross-validated between the  
146 GENCODE and RefSeq database had a better 3' and 5' validation of 95%. A total of 52,884  
147 novel isoforms were complete (NIC+NNC), of which 40,046 were confirmed as valid novel  
148 isoforms by GENCODE, corresponding to 17% of the current GENCODE v36 database.  
149 Isoforms that were not confirmed were mainly either "partially redundant with existing  
150 transcripts", or "overlapping with multiple loci". Finally, we assessed the biotypes of our newly  
151 discovered isoforms, indicative of their presumed functional categorization. We found that 42%  
152 are protein coding, more than the 36% of protein coding isoforms found in the GENCODE  
153 database (230k entries) (**Fig. 1f-g**). This demonstrates the ability of concatenated long-read

154 sequencing to generate high yield, high-quality data and discover novel isoforms with  
155 enhanced annotation.

## 156 Long-read sequencing allows for short-read-independent cell type 157 identification

158 Next, through comparison to short-read data, we assessed the ability of long-read sequencing  
159 to cluster cells and to identify cell types. We generated short- and long-read gene count  
160 matrices and removed non-protein-coding, ribosomal, and mitochondrial genes. After filtering,  
161 we obtained 16.5 Mio unique long reads associated with 12,757 genes, and 26.3 Mio unique  
162 short reads associated with 13,122 genes (**Extended Data Table 1**). The short- and long-read  
163 datasets were of similar sequencing depth with a median of 4,930 and 2,750 UMIs per cell,  
164 respectively (average 10,235 and 6,413 UMIs, **Extended Data Fig. 1a**). Also, the genes  
165 detected in both datasets overlapped by 86.4% (**Extended Data Fig. 1b,c**).

166  
167 We first identified cell types independently per cell, using cell type marker gene lists  
168 (**Methods**). We compared short- and long-read data and found that both data types identified  
169 cell types with similar percentages, namely HGSOc (13% in short-read vs 15% in long-read  
170 data), mesothelial cells (22 vs 23%), fibroblasts (9 vs 8%), T cells (38 vs 37%), myeloid cells  
171 (both 14%), B cells (3 vs 1%), and endothelial cells (both 1%). Those cell populations  
172 expressed cell type specific marker genes (**Extended Data Fig. 1c**). We then projected short-  
173 read gene, long-read gene, and long-read isoform expression onto 2-dimensional embeddings  
174 using UMAP<sup>35</sup> (**Fig. 2a**). We manually clustered cell types based on the embeddings and  
175 calculated the Jaccard distance between clusters. Cell clusters based on short- and long-  
176 reads were very similar, with a Jaccard distance >94% for all cell types except B-cells, where  
177 the Jaccard distance was >75% (**Fig. 2b**). Furthermore, Jaccard similarity analysis between  
178 cell type clusters and attributed cell type labels were analogous between short- and long-read  
179 data, with a better prediction of B cells and endothelial cells for long reads (**Extended Data**  
180 **Fig. 1b**). These findings show that long-read gene and isoform expression data can be used  
181 to identify cell types reliably and independently from short-read data.

## 182 Long-read sequencing captures germline and somatic mutations and 183 identifies increased neojunctions in tumor cells

184 Next, we assessed the potential of long-read data for mutation detection, and used somatic  
185 mutations to further validate the cell type annotation. Germline mutations are expected in all  
186 cell types, whereas somatic mutations should be present only in tumor cells. As reference, we  
187 used mutations called from a panel covering 324 genes on patient-matched bulk DNA samples  
188 (**Methods**). We identified germline variants in 48 cells belonging to all cell types from distal  
189 omentum and tumor sites (**Fig. 2c, Supplementary Table 1**). Somatic mutations were called  
190 in 34 cells, all in the cell cluster annotated as tumor cells (**Fig. 2d**). In 20 of those cells, *TP53*  
191 was found mutated (**Supplementary Table 1**). Thus, high-fidelity long-read data can be  
192 leveraged for both germline and somatic mutation calling.

193 We analyzed the expression of cell type-specific isoforms. HGSOc cells expressed more  
194 genes, transcript isoforms, and RNA molecules than other cell types (**Extended Data Fig. 3a-**  
195 **c**). This difference does however not translate into mean UMIs per isoform, as isoforms  
196 expressed in cancer cells harbor fewer UMIs than in mesothelial cells, for example. This



197 means that cancer cells express more low-abundant isoforms (**Extended Data Fig. 3d**)  
198 suggesting wider isoform diversity and broader cellular functions and controls. Isoform class  
199 distribution between cell types revealed a higher fraction of novel isoforms and neojunctions  
200 (NNC) in tumor cells (**Fig. 2e**).

201 We then looked into isoforms uniquely expressed in the different cell types. At the cell type  
202 level, cancer cells contained more than 8% (9,476) of cell type-specific isoforms, between 2.3-  
203 10.6 times more than the most frequent other cell types (myeloids, T/NK cells, fibroblasts and  
204 mesothelial cells) (**Methods, Extended Data Fig. 3e**). At the cellular level, 0.5% of the cancer-  
205 specific isoforms were also unique to a single cell, which is between 3-6 times the percentage  
206 of unique isoforms in other cell types (**Extended Data Fig. 3e**). In all cell types, cell type-  
207 specific isoforms (**Extended Data Fig. 3f**) had a higher percentage of novel isoforms than  
208 non-specific isoforms distributed across cells (**Fig. 2e**). This phenomenon was even stronger  
209 in cell-specific isoforms: in cancer, more than 75% of isoforms unique to cells were novel, and  
210 50% of these were neojunctions (NNC) (**Extended Data Fig. 3e**). Those rare isoforms were  
211 difficult to detect for previous methods, hence their novelty. Taken together, cancer cells  
212 expressed at least twice as many unique isoforms than other cell types, indicating an  
213 increased transcriptomic diversification and support previous findings of cancer-specific  
214 neojunction expression in bulk data<sup>16</sup>.

## 215 Differential isoform expression in the tumor microenvironment reveals 216 epithelial-to-mesenchymal transition

217 Comparing cells from metastatic and tumor-free samples, we found that mesothelial and  
218 fibroblast cells showed distinct clustering, in both short- and long-read embeddings (**Fig. 3a**).  
219 We observed a bridge between TME fibroblasts and mesothelial cells on the UMAPs,  
220 suggesting that TME cells might undergo a form of transdifferentiation. To understand this  
221 phenomenon, we analyzed differential isoform and gene expression in TME vs. distal  
222 mesothelial and fibroblast cells. For mesothelial cells, the gene with the highest change in  
223 relative isoform abundance amongst all its transcripts was the collagen type 1 alpha chain  
224 (*COL1A1*) ( $P_{\text{corr}}=6.34 \times 10^{-49}$ ,  $|\Delta\Pi|=0.86$ , **Methods**) (**Fig. 3b**). TME mesothelial cells used the  
225 canonical 3' transcription termination site, while distal cells had a premature transcription  
226 termination, resulting in a truncated protein (**Fig. 3c**). *COL1A1* was also the top differentially  
227 expressed gene ( $P = 2 \times 10^{-3}$ ) between TME and distal mesothelial cells, and the fifth most  
228 differentially expressed gene between TME and distal fibroblasts ( $P = 0.015$ ), with TME cells  
229 overexpressing it in both cases compared to their distal counterparts. *COL1A2*, was also found  
230 to be differentially spliced in TME mesothelial cells ( $P_{\text{corr}}=6.85 \times 10^{-91}$ ,  $|\Delta\Pi|=0.37$ ) and  
231 fibroblasts ( $P_{\text{corr}}=2.02 \times 10^{-77}$ ,  $|\Delta\Pi|=0.36$ ). HGSOC cells showed the same *COL1A2* splicing  
232 pattern as TME cells when compared to all non-tumor cells ( $P_{\text{corr}}=6.54 \times 10^{-79}$ ,  $|\Delta\Pi|=0.42$ ). Both  
233 expressed transcripts with a canonical 3'UTR, longer than the 3'UTR expressed in distal cells  
234 (**Fig. 3d**). Thus, in two cases, tumor-associated stromal cells overexpressed and used longer  
235 collagen matrix isoforms than their distal counterparts. Another top differentially expressed  
236 isoform in TME vs. distal mesothelial cells was gelsolin (*GSN*), which exists in two main protein  
237 variants: one residing in the cytoplasm (*cGSN*), the other in the extracellular (plasma)  
238 environments (*pGSN*)<sup>36</sup>. At the gene level, *GSN* was not significantly overexpressed in TME  
239 vs. distal or in HGSOC vs. non-HGSOC cells. However, TME mesothelial cells had a  
240 significantly higher *cGSN/pGSN* isoform ratio than distal ones ( $P_{\text{corr}}=2.49 \times 10^{-18}$ ,  $|\Delta\Pi|=0.34$ )  
241 (**Fig. 3e**). Similarly, cancer cells had a significantly higher *cGSN/pGSN* ratio than non-cancer

242 cells ( $P_{\text{corr}}=3.4 \times 10^{-127}$ ,  $|\Delta\Pi|=0.28$ ), and consistent with findings for *COL1A*, TME cells displayed  
243 a cancer-like isoform expression profile compared to cells from distal sites, suggesting tissue  
244 mimicry. To test if the differential expression of those structural isoforms in TME cells could  
245 be linked to epithelial-to-mesenchymal transition (EMT), we performed gene set enrichment,  
246 which revealed the EMT pathway as enriched in TME mesothelial and fibroblasts cells (**Fig.**  
247 **3f**) supporting the idea of a tumor-transformed stroma.

## 248 Differential isoform expression in cancer reveals isoform-specific *IGF1* 249 usage

250 HGSOC cells significantly expressed different isoforms in 17% of the genes, compared to  
251 all distal cells, but only 0.6% were switched with  $|\Delta\Pi|>0.5$  (6,841 genes tested, **Methods**)  
252 (**Extended Data Fig. 4a**). One of the most significant switches was found in the insulin-like  
253 growth factor gene *IGF1* ( $P_{\text{corr}}=1.1 \times 10^{-130}$ ,  $|\Delta\Pi|=0.68$ ), a gene coding for a hormone linked to  
254 the development, progression, survival, and chemoresistance of many cancer types including  
255 ovarian cancer<sup>37</sup>. Cancer cells from all patients almost exclusively used the second exon of  
256 the gene as their transcription start site (Class II isoform), whereas other cells mainly used the  
257 first exon (Class I isoform)<sup>38</sup> (**Fig. 4a,b**). The Class II isoform was highly expressed in HGSOC,  
258 with 95% of cancer cells expressing it (**Fig. 4c,d**). Reflecting the findings of the DIE analysis  
259 in mesothelial cells, fibroblasts and mesothelial cells in the TME also expressed a higher  
260 fraction of class II isoforms than cells derived from distal biopsies (**Fig. 4d**). *IGF1* was found  
261 to be significantly higher expressed in cancer cells ( $P_{\text{corr}}=4.8 \times 10^{-32}$ ) as well as in TME  
262 mesothelial cells and fibroblasts compared to distal mesothelial cells and fibroblasts  
263 ( $P_{\text{corr}}=4.05 \times 10^{-32}$ ).

264  
265 Similarly, cancer and TME cells differentially expressed multiple isoforms in the two actin-  
266 associated tropomyosin genes *TPM1* and *TPM2*. Cancer cells expressed terminal exon 9a  
267 and exon 6b of *TPM2* ( $P_{\text{corr}}<10^{-293}$ ,  $|\Delta\Pi|=0.28$ ), and TME cells also expressed those exons  
268 more than distal ones (**Extended Data Fig. 4a-d**). Cancer cells also preferentially expressed  
269 exon 1b and 6a of *TPM1* (**Extended Data Fig. 4e**). Another strongly switched gene in cancer  
270 cells is vesicle-associated *VAMP5* ( $P_{\text{corr}}=4.59 \times 10^{-17}$ ,  $|\Delta\Pi|=0.70$ ). Indeed, the overexpressed  
271 isoforms in HGSOC cells were a (predicted protein-coding) *VAMP8-VAMP5* read-through  
272 gene, i.e., a novel gene formed of two adjacent genes (**Extended Data Fig. 4f**). HGSOC cells  
273 expressed almost no wild-type (wt) *VAMP5* but had a significantly higher *VAMP8* expression  
274 than other cells ( $P_{\text{corr}}=1.0 \times 10^{-15t}$ ), indicating that this read-through gene was under  
275 transcriptional control of *VAMP8*. Amongst others, HGSOC cells also differentially expressed  
276 isoforms in the Golgi vesicle-associated *AP1S2* gene ( $P_{\text{corr}}=6.52 \times 10^{-97}$ ,  $|\Delta\Pi|=0.60$ ).  
277 Fibroblasts, mesothelial, and myeloid cells expressed the canonical isoform (Uniprot: P56377-  
278 1), whereas HGSOC cells used another terminal 3' exon (Uniprot: A0A5F9ZHW1) (**Extended**  
279 **Data Fig. 4g**). Last, patient 2 cancer cells highly expressed a novel shortened isoform of  
280 ceramide kinase gene *CERK*, ( $P_{\text{corr}}=1.38 \times 10^{-39}$ ,  $|\Delta\Pi|=0.78$ ) (**Extended Data Fig. 4h**). In  
281 summary, tumor cells showed differential isoform usage in genes associated with hormonal  
282 (*IGF1*), actin (*TPM1*, *TPM2*, *GSN*), vesicle (*VAMP8-VAMP5*, *APS1A*), and sphingolipid  
283 (*CERK*) functions.



## 284 Long-read sequencing captures gene fusions and identifies an 285 *IGF2BP2::TESPA1* fusion that was misidentified in short-read data

286 To detect fusion transcripts, we aligned long reads to the reference genome and filtered for  
287 reads split-aligned across multiple genes. We then ranked fusion transcripts with counts  
288 across all cells of more than 10 UMIs (**Supplementary Table 2**). Out of the 34 detected fusion  
289 entries, 21 were genes fused with mitochondrial ribosomal RNA (*mt-rRNA1-2*) and ubiquitous  
290 among all cell types, 11 isoforms were *IGF2BP2::TESPA1* fusions specific to patient 2, one  
291 was a cancer cell-specific *CBLC* (chr8:43.064.215) fusion to a long non-coding RNA (lncRNA)  
292 expressed in patient 3, and one was a cancer cell-specific fusion of *FNTA* with a lncRNA  
293 expressed in patient 1. The ubiquitous *mt-rRNA* fusions were likely template-switching artifacts  
294 from the library preparation, as *rRNA* makes up to 80% of RNA in cells<sup>39</sup>. *IGF2BP2::TESPA1*  
295 was a highly expressed fusion event in patient 2: 2,174 long-reads mapped to both *IGF2BP2*  
296 (Chr3) and *TESPA1* (Chr12). The gene fusion consisted of 5' located exons 1-4 of *IGF2BP2*,  
297 corresponding to 112 amino acids (aa) and including the RNA recognition motif 1 (RRM1) and  
298 half of the RRM2 domain, linked to the terminal *TESPA1* 3' untranslated region (UTR) exon,  
299 encoding 69 aa as in-frame fusion and including no known domains (**Fig. 5a**). In total, the  
300 gene fusion encoded 181 aa, compared to 599 aa of wt *IGF2BP2* and 521 aa of wt *TESPA1*  
301 (**Fig. 5b**). 98.9% of fusion reads were found in HGSOC cells and the fusion was detected in  
302 86.8% of patient 2's cancer cells, making it a highly cancer cell- and patient-specific fusion  
303 event (**Fig. 5c**). Cancer cells lacking the gene fusion had lower overall UMI counts, suggesting  
304 low coverage as a possible reason for the absence of the gene fusion (**Fig. 5d**).

305  
306 We next investigated the footprint of the gene fusion in the short-read data. The *TESPA1* gene  
307 was expressed uniquely in T cells and highly expressed only in patient 2, almost exclusively  
308 in HGSOC cells, and colocalized with *IGF2BP2* expression (**Fig. 5e,f**). In short-read data,  
309 *TESPA1* was the highest differentially expressed gene in cancer cells compared to non-cancer  
310 cells in patient 2 ( $P_{\text{corr}}=1.17 \times 10^{-14}$ ). Next, we designed a custom reference including the  
311 *IGF2BP2::TESPA1* transcriptomic breakpoint as well as wt *TESPA1* and wt *IGF2BP2*  
312 junctions and re-aligned Patient 2's short-reads (**Extended Data Fig. 5, Methods**). Out of the  
313 989 reads mapping to the custom reference, 94% preferentially aligned to *IGF2BP2::TESPA1*  
314 (99.8% of those in HGSOC cells). This implies that the reported overexpression of *TESPA1* in  
315 short-reads is false, as nearly all junction reads map to the fusion and not the wt gene. Reads  
316 covering the *TESPA1* 3' UTR region harbored three heterozygous single nucleotide  
317 polymorphisms (hSNPs): chr12:54.950.144 A>T (rs1047039), chr12:54.950.240 G>A  
318 (rs1801876), and chr12:54.950.349 C>G (rs2171497). In long reads, wt *TESPA1* was either  
319 triple-mutated or not mutated at all, indicating two different alleles. All fusion long reads,  
320 however, were triple-mutated, indicating a genomic origin and monoallelic expression of the  
321 fusion (**Fig. 5g**). In short reads, the three loci were mutated in nearly all reads, supporting the  
322 hypothesis that the observed *TESPA1* expression represents almost completely  
323 *IGF2BP2::TESPA1* expression and that it has a genomic origin.

## 324 Genomic breakpoint validation of the *IGF2BP2::TESPA1* fusion

325 To validate that the *IGF2BP2::TESPA1* gene fusion is the result of genomic rearrangements,  
326 we looked for a breakpoint in single-cell DNA sequencing (scDNA-seq) data from a patient  
327 2-matched metastatic sample. Two RNA fusion long reads mapped to intronic regions of  
328 *IGF2BP2* and *TESPA1* (**Extended Data Fig. 5**) indicating the location of the breakpoint at

329 chr3:185.604.020-chr12:54.960.603. We then estimated the scDNA-seq copy number  
330 profiles of all cells and identified two clones among the 162 cells of the scDNA sample: a  
331 cancer clone (Subclone 0) and a copy number-neutral non-cancer clone (Subclone 1)  
332 (**Fig. 6a**). We next aligned the scDNA data to a custom reference covering the breakpoint  
333 (**Methods, Supplementary Dataset 1**), including the wt *TESPA1*, wt *IGF2BP2*, and  
334 *IGF2BP2::TESPA1* fusion sequences. We found nine reads mapping to the breakpoint (nine  
335 in subclone 0 cancer cells, zero in subclone 1 cells,  $P=0.0321$ ) (**Fig. 6b**). We also found 14  
336 reads mapping to wt *IGF2BP2* (ten in subclone 0 cells, four in subclone 1 cells,  $P=0.78$ ) (**Fig.**  
337 **6c**), and eight reads mapping to wt *TESPA1* (five subclone 0 cells, three subclone 1 cells,  
338  $P=1.0$ ) (**Fig. 6d**). Thus, scDNA-seq data confirmed the breakpoint in the intronic region  
339 detected by the long-read scRNA-seq. The scDNA-seq data also confirmed that the  
340 *IGF2BP2::TESPA1* fusion was cancer-cell specific, as suggested by long-read scRNA-seq  
341 data. *IGF2* RNA, which is bound by the wt IGF2BP2 protein, is also largely overexpressed in  
342 patient 2 cancer cells compared to other patients ( $P_{\text{corr}} < 2.54 \times 10^{-15}$ ). The genomic region  
343 containing *IGF2BP2* has an increased copy number (**Fig. 6a**) in patient 2, so the fact that  
344 one allele is a fusion allele does not impair the wt *IGF2BP2* transcription.

## 345 Discussion

346 Detecting genomic alterations such as mutations<sup>40,41</sup> and gene fusions<sup>42,43</sup> in combination with  
347 isoform-level<sup>15</sup> transcriptomic readouts on the single-cell level can provide valuable  
348 information on cancer formation, progression, the role of the TME, drug targets, and therapy  
349 response<sup>44</sup>. Here, we applied PacBio HiFi high-throughput long-read RNA-seq on five omental  
350 metastases and tumor-free samples from chemo-naive HGSOc patients to detect and quantify  
351 all of these alterations.

352 Until now, a combination of single-cell short- and long-read sequencing was necessary to  
353 identify cell-specific isoforms: the higher depth of short-read sequencing allowed for cell typing  
354 based on gene expression, while long-read sequencing was used to identify isoforms<sup>22</sup>.  
355 Leveraging multiple strategies to generate high PacBio sequencing output, we achieved a 50-  
356 fold increased sequencing depth compared to the first long-read PacBio scRNA-seq study<sup>22</sup>  
357 allowing for short read-comparable cell type identification. Consequently, future studies with  
358 similar or increased long-read throughput will not have to rely on parallel short-read  
359 sequencing, thereby saving cost and labor.

360 Our analysis revealed a differential isoform usage between distal tumor-free and TME  
361 mesothelial cells in extracellular matrix associated genes (*COL1A1*, *COL1A2*, *GSN*). A  
362 geneset enrichment analysis between the two sites revealed higher EMT pathway enrichment  
363 in TME-derived mesothelial cells and fibroblasts. These findings are consistent with increasing  
364 evidence that EMT in the TME is induced by cancer cells, leading to cancer-associated  
365 phenotypes<sup>45</sup> including TGF $\beta$ 1-induced mesenchymal states of mesothelial cells in ovarian  
366 cancer<sup>46</sup>. Notably, in *IGF1*, *TPM2*, *GSN* and *COL1A2* genes, we found overlap in isoform  
367 usage between cancer and TME cells (fibroblasts and mesothelial cells). Whether this cancer  
368 mimicry of the TME is caused by signaling or the result of mRNA exchange via tumor-secreted  
369 extracellular vesicles<sup>47</sup>, as it was shown for *GSN*<sup>48</sup>, requires further investigation.

370 Additionally, we demonstrated the potential of the technology in terms of coverage and  
371 sequencing accuracy to detect mutations and gene fusions. In particular, in one patient, the

372 novel fusion *IGF2BP2::TESPA1* was highly overexpressed compared to wt *IGF2BP2* (~10x  
373 more) and *TESPA1* (~150x more). *IGF2BP2* is known to be regulated via 3'UTR miRNA  
374 silencing<sup>49</sup>, however the *IGF2BP2::TESPA1* fusion has the unregulated 3'UTR of *TESPA1*,  
375 which could explain its overexpression. *TESPA1* is normally expressed in T cells<sup>50</sup> and long-  
376 read data confirmed T cell-specific wt *TESPA1* expression. Short read data however  
377 erroneously reported *TESPA1* as the most differentially expressed gene in cancer cells,  
378 resulting from 3' end capture of the fusion transcripts. This highlights that short-read scRNA-  
379 seq data fails to distinguish between gene and fusion expression, potentially leading to wrong  
380 biological conclusions.

381

382 Overall, HGSOC cells revealed a profoundly modified IGF system in all patients, with a drastic  
383 switch from *IGF1* Class I to Class II isoform, *IGF2* overexpression, and a highly expressed  
384 *IGF2BP2* gene fusion in one patient. The *IGF* protein family promotes cancer growth, survival,  
385 proliferation, and drug resistance through signaling via *PI3K-AKT* or *MAPK*, and is a known  
386 clinical target in ovarian cancer<sup>37</sup>. Secreted (Class II) *IGF1* is associated with the progression  
387 of ovarian cancer<sup>51</sup> and the observed overexpression of Class II IGF1 in HGSOC cells could  
388 mediate uncontrolled cell proliferation in the tumor.

389 Although the achieved sequencing depth allowed for short-read independent cell typing and  
390 clustering, a further increased depth is needed to capture low abundance transcripts. For  
391 example, we did not obtain sufficient reads to retrieve and characterize the T cell receptor  
392 repertoire. This is consistent with a long-read scRNA-seq study in blood lymphocytes that  
393 reported a 3.6-fold lower pairing rate for T cell receptors than the higher abundant B cell  
394 receptors from plasmablasts<sup>52</sup>. With further technological advances and decreased  
395 sequencing costs, however, we expect that these limitations can and will be overcome.  
396 Enrichment for low abundant transcripts for long-read sequencing or depletion of mitochondrial  
397 and ribosomal RNA<sup>53</sup> represent interesting avenues forward.

398 Altogether, we demonstrate that long-read sequencing provides a more complete picture of  
399 cancer-specific changes. These findings highlight the manifold advantages and new  
400 opportunities that this technology provides to the field of precision oncology, opening the  
401 premise of personalized drug prediction and neoantigen detection for cancer vaccines<sup>54,55</sup>.

## 402 Materials and Methods

### 403 Omentum patient cohort

404 The use of material for research purposes was approved by the corresponding cantonal ethic  
405 commissions (EKNZ: 2017–01900, to V.H.S.) and informed consent was obtained for all  
406 human primary material. Tissue samples were immediately collected from the theater and  
407 transferred on ice to the department of biomedicine of the University Hospital Basel for tissue  
408 dissociation.

### 409 Sample processing

410 Fresh omentum and omental HGSOC tumor metastasis biopsy samples were cut into small  
411 pieces and dissociated in digestion solution (1 mg/mL collagenase/Dispase [Sigma cat. no.  
412 10269638001], 1 unit/mL DNase I [NEB, cat. no. M0303] and 10% FBS in DMEM [Sigma, cat.

413 no. D8437-500mL]) for 30 min at 37°C. To focus on the non-adipose cell fraction, adipocytes  
414 were separated by centrifugation and the cell pellet was collected. Red blood cell lysis (RBC)  
415 was performed using MACS red blood lysis solution (cat. no. 130-094-183). Then, the cell  
416 pellet was resuspended into MACS dead cell removal microbeads (cat. no. 130-090-101) and  
417 was loaded into the AutoMACS separator to remove dead cells. After counting cell number,  
418 cells were resuspended in PBS with 1% BSA and transferred to the Genomics Facility Basel.  
419 The cell suspension was again filtered and cell number and viability was assessed on a  
420 Cellometer K2 Image Cytometer (Nexcelom Bioscience, cat. no. Cellometer K2) using  
421 ViaStain AOP1 Staining Solution (Nexcelom Bioscience, cat. no. CS2-0106-5mL) and PD100  
422 cell counting slides (Nexcelom Bioscience, cat. no. CHT4-PD100-003). For samples with  
423 viability below 70% and when cell numbers allowed ( $>10^5$  cells total), apoptotic and dead cells  
424 were removed by immunomagnetic cell separation using the Annexin Dead Cell Removal Kit  
425 (StemCell Technologies, cat. no. 17899) and EasySep Magnet (StemCell Technologies, cat.  
426 no. 18000). If the cell pellet appeared still red, additional RBC lysis was performed. Cells were  
427 washed with a resuspension buffer (PBS with 0.04% BSA), spun down and resuspended in a  
428 resuspension buffer. Finally, cells were again counted and their viability determined. The cell  
429 concentration was set according to 10x Genomics protocols (700-1,200 cells/ $\mu$ L).

#### 430 10x Genomics single-cell capture and short-read sequencing

431 Cell suspensions were loaded and processed using the 10x Genomics Chromium platform  
432 with the 3P v3.1 kit on the 10x Genomics Chromium Single Cell Controller (10x Genomics,  
433 PN-120263) according to the manufacturer's instructions. 500 or 1,000 cells were targeted per  
434 lane. The quality of cDNA traces and GEX libraries were profiled on a 5200 Fragment Analyzer  
435 (Agilent Technologies).  
436 Paired-end sequencing was performed on the Illumina NovaSeq platform (100 cycles, 380pm  
437 loading concentration with 1% addition of PhiX) at recommended sequencing depth (20,000-  
438 50,000 reads/cell).

#### 439 Long-read library preparation and PacBio sequencing

440 To increase long-read PacBio sequencing throughput, we followed the strategy of cDNA  
441 concatenation of the HIT-sclSOseq protocol<sup>23</sup> with the modification of two rounds of biotin-  
442 PCR in order to further reduce template-switch oligo (TSO) artifacts from the data.  
443 Full protocol details:

#### 444 cDNA amplification and biotin-enrichment

445 15 ng of each patient's cDNA library were amplified using the KAPA HiFi HotStart Uracil+  
446 ReadyMix 2x (Kapa Biosystems, cat. no. KK2801) with 0.5  $\mu$ M final concentration of custom-  
447 primers (Integrated DNA Technologies, HPLC purified). Primers contained overhang  
448 sequences adapted from Hebelstrup *et al.*<sup>56</sup> with a single deoxyuridine (dU) residue at a 10  
449 nt distance from the 5' terminus enabling USER enzyme digestion and creating single-  
450 stranded overhangs. Generated PCR fragments thus contain a single dU residue per DNA  
451 strand. The forward primer was specific to the 10x Genomics partial Read 1 sequence and  
452 contained a biotin modification allowing for biotin enrichment of amplified full-length cDNA  
453 molecules. The reverse primer was specific to the 10x Genomics partial TSO sequence.  
454 Forward Primer: /5Biosg/AGGTCTTAA/ideoxyU/CTACACGACGCCTTCCGATCT  
455 Reverse Primer: ATTAAGACC/ideoxyU/AAGCAGTGGTATCAACGCAGAG



456 The PCR was run according to the manufacturer's instruction with two cycles at an annealing  
457 temperature of 63°C followed by 7 cycles at an annealing temperature of 67°C; annealing time  
458 was 30 seconds. Extension was performed at 72°C for 90 seconds. PCR products were  
459 purified at 0.6X SPRIselect bead cleanup (Beckman Coulter, cat. no. B23318) according to  
460 the manufacturer's instructions and eluted in 22 µL EB buffer (Qiagen, cat. no. 19086). DNA  
461 concentrations were measured using the Qubit dsDNA HS Assay Kit (Thermo Fisher  
462 Scientific, cat. no. Q32854), which were in the range of 1.5 µg per sample. cDNA traces were  
463 additionally evaluated on a 5200 Fragment Analyzer System (Agilent Technologies) using the  
464 HS NGS Fragment Kit, 1-6000 bp (Agilent, cat. no. DNF-474-0500). Full-length cDNAs were  
465 enriched through capture on 5 µL streptavidin-coated M-280 dynabeads using the  
466 Dynabeads™ kilobaseBINDER™ Kit (Invitrogen, cat. no. 60101), thus depleting TSO-TSO  
467 artifacts. Washed Dynabeads containing the DNA-complexes were directly resuspended in 20  
468 µL USER reaction buffer containing 10 µL StickTogether DNA Ligase Buffer 2x (NEB, cat. no.  
469 B0535S), 1.5 µL USER Enzyme (NEB, cat. no. M5505S) and 8.5 µL Nuclease-free water  
470 (Invitrogen, AM9939) and incubated in a thermocycler at 37°C for 20 min and held at 10°C (no  
471 annealing). This created a nick at the deoxyuracil site forming palindrome overhangs and  
472 releasing the biotin-bound DNA molecules from the beads. Beads were removed by magnetic  
473 separation and the supernatant with the biotin-released cleaved PCR products was subjected  
474 to a 0.6X SPRIselect cleanup step. Approximately 100 ng of purified product per sample were  
475 split into two aliquots and subjected to a second PCR amplification step with 6 cycles using  
476 an annealing temperature of 67°C. Reactions were pooled, purified by 0.6X SPRIselect  
477 cleanup and quality checked on both Qubit and Fragment Analyzer. Total DNA yield was  
478 between 5-8 µg, which were subjected to a second round of streptavidin-purification using 10  
479 µL of beads.

#### 480 Transcript ligation

481 Beads were incubated in 19 µL USER reaction buffer at 37°C for 20 min for USER digestion  
482 and 25°C for 17 min for overhang annealing. Beads were then removed by magnetic  
483 separation and the supernatant was transferred to a new PCR tube. 1 µL of T4 DNA ligase  
484 high-concentration (2,000,000, units/mL, NEB, cat. no. M0202T) was added, mixed and  
485 incubated at 10°C for >24hrs and heat inactivated at 65°C for 10 min. To efficiently deplete  
486 any non-ligated transcripts, 0.38X SPRIselect cleanup was performed, eluted in 20 µL EB  
487 buffer and traces were evaluated on the Fragment Analyzer using the HS Large Fragment kit  
488 (Agilent Technologies, cat. no. DNF-492-0500) at 1:5 dilutions. Ligation products were 8-11kb  
489 long; average yield was 100 ng per sample.

#### 490 End repair/dA tailing, adapter ligation and PCR amplification

491 To enable PCR-amplification of the ligated construct, the NEBNext Ultra II DNA Library Prep  
492 Kit for Illumina was followed (NEB, cat. no. E7645S) using total DNA yield as input material.  
493 2.5 µL of 5 µM dT overhang adapter (Roche, cat. no. KK8727) were used for the End Prep  
494 reaction. Adapter-ligated libraries were purified by 0.39X SPRIselect cleanup, eluted in 22 µL  
495 EB buffer and products were evaluated by HS Large Fragment kit. Total yield of around 40 ng  
496 was split in two and PCR amplified using 2X KAPA HiFi Hot-Start ReadyMix (Roche, cat. no.  
497 KK2602) and KAPA Library Amplification Primer Mix (10X concentration, Roche, cat. no.  
498 KK2623), 10 µL library input each with 11 cycles and 9 min extension time. Following a 0.38X  
499 SPRIselect cleanup and elution in 48 µL EB buffer, products were evaluated on a large  
500 fragment gel revealing an average fragment length of libraries of 4.6 kb and average total of



501 1.1 µg DNA. To increase total yield to 2 µg DNA required for SMRTbell library preparation of  
502 a product with 5 kb amplicon size, the PCR was repeated with three additional cycles and  
503 5 min extension time. After 0.4X SPRI cleanup and Fragment Analyzer inspection, the final  
504 yield was 2 µg per library.

#### 505 PacBio SMRTbell library preparation

506 The SMRTbell Express Template Kit (PacBio, cat. no. 100-938-900) was used following  
507 manufacturer's instructions for DNA damage repair, end repair/dA-tailing and ligation of a  
508 hairpin adapter (double amount used). Final purification of the SMRTbell template was  
509 performed by 0.42X SPRIselect cleanup and elution in 43 µL EB buffer. Exonuclease  
510 treatment was performed by addition of 5 µL of NEBuffer1 (NEB, cat. no. B7001S) and 1 µL  
511 of each Exonuclease I (NEB, cat. no. M0293S) and Exonuclease III (NEB, cat. no. M0206S)  
512 bringing the total volume to 50 µL per reaction. Enzyme treatment was performed at 37°C for  
513 60 min. After SPRIselect cleanup, products were quantified on a large fragment gel at 1:30  
514 dilution. Final yield was approximately 650 ng per sample, a sufficient amount for long-read  
515 sequencing.

#### 516 PacBio Sequel II sequencing

517 Libraries were sequenced on the PacBio Sequel II platform with the SMRT cell 8M. Omentum  
518 metastasis and tumor-free omentum were run on three and two 8M cells, respectively.

#### 519 Single-cell DNA-sequencing

520 Cell suspensions were loaded and processed using the 10x Genomics Chromium platform  
521 with the single-cell CNV kit on the 10x Genomics Chromium Single Cell Controller (10x  
522 Genomics, PN-120263) according to the manufacturer's instructions. Paired-end sequencing  
523 was performed on the Illumina NovaSeq platform (100 cycles, 380pm loading concentration  
524 with 1% addition of PhiX) at recommended sequencing depth.

#### 525 Data Analysis

##### 526 Short-read data analysis

##### 527 Preprocessing

528 Raw reads were mapped to the GRCh38 reference genome using 10x Genomics Cell Ranger  
529 3.1.0 to infer read counts per gene per cell. We performed index-hopping removal using a  
530 method developed by Griffiths *et al.*<sup>57</sup>.

##### 531 10x Genomics short-read analysis

532 GEX data of each sample was analyzed using the scAmp workflow<sup>58</sup>. In brief, UMI counts  
533 were quality controlled and cells and genes filtered to remove known contaminants. Cells  
534 where over 50% of the reads mapped to mitochondrial genes and cells with fewer than 400  
535 different expressed genes were removed, as well as non protein-coding genes and genes that  
536 were expressed in less than 20 cells. Doublet detection was performed using scDblFinder<sup>59</sup>.  
537 Subsequently, counts were normalized and corrected for cell cycle effects, library size, and  
538 sample effect using sctransform<sup>60</sup>. Similar cells were grouped based on unsupervised

539 clustering using Phenograph<sup>61</sup> and an automated cell type classification was performed  
540 independently for each cell<sup>62</sup> using gene lists defining highly expressed genes in different cell  
541 types from previous publications. Major cell type marker lists were developed in-house based  
542 on unpublished datasets (manuscripts in preparation) including the Tumor Profiler Study<sup>63</sup>  
543 using the Seurat FindMarkers method<sup>64</sup>. Immune subtype marker gene lists were obtained  
544 from Newman *et al.*<sup>65</sup>, enriched with T cell subtypes from Sade-Feldman *et al.*<sup>66</sup>

## 545 Long-read data analysis

### 546 Generating CCS

547 Using SMRT-Link (version 9.0.0.92188), we performed circular consensus sequencing (CCS)  
548 with the following modified parameters: maximum subread length 50,000 bp, minimum  
549 subread length 10 bp, and minimum number of passes 3.

### 550 Unconcatenating long reads

551 We used NCBI BLAST (version 2.5.0+) to map the 5' and 3' primers to CCS constructs, with  
552 parameters: "-outfmt 7 -word\_size 5" as described previously<sup>23</sup>. Sequences between two  
553 successive primers were used as input for primer trimming using IsoSeq3 Lima (parameters:  
554 --isoseq --dump-clips --min-passes 3). Cell barcodes and UMIs were then demultiplexed using  
555 IsoSeq3 tag with parameter --design T-12U-16B. Finally, we used IsoSeq3 refine with option  
556 --require-polya to remove concatemers and trim polyA tails. Only reads with a correct 5'-3'  
557 primer pair, a barcode also found in the short-read data, a UMI, and a polyA tail were retained.

### 558 Isoform classification

559 Demultiplexing UMIs with IsoSeq3 dedup and calling isoforms on the cohort level with  
560 collapse\_isoforms\_by\_sam.py resulted in unfeasible runtimes. Therefore, we called isoforms  
561 first on the cell level as a pre-filtering step. Long-reads were split according to their cell  
562 barcodes, and UMI deduplication was performed using IsoSeq3 dedup. Next, reads were  
563 mapped and aligned to the reference genome (hg38) using minimap2 with parameters: -ax  
564 splice -uf --secondary=no -C5. Identical isoforms were merged based on their aligned exonic  
565 structure using collapse\_isoforms\_by\_sam.py with parameters: -c 0.99 -i 0.95 --  
566 gen\_mol\_count. We then classified isoforms using SQANTI3<sup>31</sup> with arguments: --skipORF --  
567 fl\_count --skip\_report. We finally filtered artifacts including intrapriming (accidental priming of  
568 pre-mRNA 'A's), reverse-transcriptase template switching artifacts, and mismapping to non-  
569 canonical junctions. In order to have a unique isoform catalog for all our samples, we then  
570 retained only reads associated to isoforms passing the SQANTI3 filter, and we ran  
571 collapse\_isoforms\_by\_sam.py, SQANTI3 classification and filtering again on all cells together.  
572 The described pipeline is available [here](#) and was implemented in Snakemake, a reproducible  
573 and scalable workflow management system<sup>67</sup>.

### 574 3' and 5' isoform filtering

575 For SQANTI3-defined isoforms, incomplete splice match, novel in catalog and novel not in  
576 catalog, we only retained isoforms falling within 50 bp of a CAGE-validated transcription start  
577 site (FANTOM5 CAGE database), and 50 bp of a polyA site from the PolyASite database<sup>33</sup>.  
578 The GENCODE database was used as a comparison, all protein-coding isoforms were  
579 grouped under the GENCODE.full label, a subset including only full-length isoforms was

580 labeled as GENCODE.FL, and the Matched Annotation from NCBI and EMBL-EBI (MANE<sup>34</sup>)  
581 was named GENCODE.MANE.

## 582 Isoforms biotypes

583 Novel isoform biotypes were assessed internally by the GENCODE team with biotypes  
584 matching those described by Frankish *et al.*<sup>68</sup>.

## 585 Cell type-specific isoforms

586 Considering only the SQANTI3-defined 'full splice match', 'novel not in catalog' and 'novel in  
587 catalog' isoforms with at least 3 reads, we established the following classification: "Cell-  
588 specific" isoforms are present in only 1 cell and "cell type specific" isoforms are present in  $\geq 3$   
589 cells of a unique cell type.

590

## 591 Cell type annotation

592

593 Cells were annotated with long-reads the same way as short-reads, using scROSHI. The  
594 major cell types were modified according to gene expression in long-reads. Immune subtype  
595 marker gene lists were unchanged.

## 596 Mutation detection

597 Positions of mutations from Foundation Medicine's targeted NGS panel (Foundation One CDx)  
598 mutations described in Table 1 were used as reference. One mutation not present in the list,  
599 TP53\_P151H, was visually detected in Patient 1 and added to the list. If a position was  
600 mutated at least in one cell belonging to a distal biopsy sample, the mutation was classified  
601 as a germline variant. Cells with one mutated read in one of the positions were considered  
602 mutated.

## 603 Differential isoform tests

604 Differential isoform testing was performed using a  $\chi^2$  test as previously described in  
605 Scisorseqr<sup>25</sup>. Briefly, counts for each isoform ID were assigned to individual cell types, and  
606 genes were discarded if they did not reach sufficient depth per condition (25 reads per  
607 condition per gene). P-values from a  $\chi^2$  test for differential isoform usage were computed per  
608 gene where a sufficient depth was reached, and we corrected for multiple testing using  
609 Benjamini Hochberg correction with a 5% false discovery rate. If the corrected p-value was  
610  $\leq 0.05$  and the sum of change in the relative percent of isoform ( $\Delta\pi$ ) of the top two isoforms in  
611 either positive or negative direction was more than 10%, then the gene was called differentially  
612 spliced. To classify the top differentially spliced genes, we took the rank of genes by  $\Delta\pi$  and  
613 corrected p-values, and summed those two ranks. The smallest sum of ranks were considered  
614 as the top differentially expressed genes. Differentially used isoforms were visualized using  
615 ScisorWiz<sup>69</sup>.

## 616 Pathway enrichment analysis

617 We used GSVA to perform pathway enrichment analysis. Gene sets were obtained from the  
618 default scAmp workflow<sup>70</sup>, with the addition of the  
619 EPITHELIAL\_MESENCHYMAL\_TRANSITION pathway from GSEA.

## 620 Fusion Discovery

621 Mapped reads from isoform classification were pooled. We called reads mapping to two  
622 separate genes at a distance of more than 100,000 bp or to different chromosomes using  
623 fusion\_finder.py (cDNA\_Cupcake package, [https://github.com/Magdoll/cDNA\\_Cupcake](https://github.com/Magdoll/cDNA_Cupcake)) with  
624 parameters --min\_locus\_coverage\_bp 200 -d 1000000. Fusion isoforms with sufficient depth  
625 (min. 10 reads) were kept, and their breakpoint, expression per cell type and number of cells  
626 in which they are expressed was assessed.

## 627 Short-reads re-alignment to *IGF2BP2::TESPA1*

628 We designed a custom reference including *IGF2BP2::TESPA1* transcriptomic breakpoint as  
629 well as the wild-type *IGF2BP2* and *TESPA1* exon junction covering the breakpoint. The  
630 reference was composed of 5 sequences of 80 nucleotides (40 bases upstream and  
631 downstream of the breakpoint), sequences XXX\_1 and XXX\_2 represent the breakpoints of  
632 the two main isoforms seen in each gene:

```
633  
634 >TESPA1_wt_1  
635 TTCTGTCAGACCACATGCTGTTGTGGTGGTGGAGAAAGCAATTCTGGAGGCTGGCAAATCCAAG  
636 GTCAAAGCCTGCA  
637  
638 >TESPA1_wt_2  
639 TTCTGTCAGACCACATGCTGTTGTGGTGGTGGAGAAAGCTTCACGAGTCTTGCCAGCAAAAAGTC  
640 TGGTGGTGGTGGG  
641  
642 >IGF2BP2_wt_1  
643 ATGTGACGTTGACAACGGCGTTTTCTGTGTCTGTGTTGACTTGTTCCACATTCTCCACTGTCCCA  
644 TATTGAGCCAAAA  
645  
646 >IGF2BP2_wt_2  
647 ATCACTGGATTGTGTGTTCTTCTGAATTACTTCTTAGGCTTGTTCCACATTCTCCACTGTCCCAT  
648 ATTGAGCCAAAA  
649  
650 >TESPA1_IGF2BP2_fusion_1  
651 TTCTGTCAGACCACATGCTGTTGTGGTGGTGGAGAAAGCCTTGTTCCACATTCTCCACTGTCCCA  
652 TATTGAGCCAAAA  
653  
654 >TESPA1_IGF2BP2_fusion_2  
655 CAAATCCAAGGTCAAAGCCTGCATCTGGTGGAGGGCCTCCTTGTTCCACATTCTCCACTGTCCCA  
656 TATTGAGCCAAAA
```

657  
658 Patient 2 reads were aligned to this reference using minimap2 with parameters: -ax sr --  
659 secondary=no. Reads mapping unambiguously to one of those reference sequences were  
660 then attributed to the cell type to which their cell barcode belonged.

## 661 scDNA analysis

662 Cell Ranger DNA was used to demultiplex and align Chromium-prepared sequencing  
663 samples. We used the cellranger-dna mkfastq command to generate FASTQ files from the  
664 Illumina raw BCL files, and we ran the command cellranger-dna cnv to align FASTQ files to  
665 the hg38 reference genome, call cells, and estimate copy numbers. We obtained the copy  
666 number profiles and detected the main clonal structure of samples using SCICoNE<sup>71</sup>.

## 667 DNA breakpoint validation

668 To validate in scDNA data breakpoints found in scRNA data, we used the putative scRNA  
669 breakpoint reads as a reference to re-align scDNA reads using BWA with options: -pt8 -CH.  
670 For the *IGF2BP2::TESPA1* fusion, the reference was composed of 3 sequences of 184  
671 nucleotides (92 bases upstream and downstream of the breakpoint):

672  
673 >IGF2BP2\_WT  
674 CAAACTTGTAGAAATGTGAATTTTTCTTGTTATTTTACAAGATTTGCAAAGGGACCTGAGACCCCG  
675 AAAAGCTTAAGGACTACTGTTAAAAATACTGTTTGTAAATAACTTTAAAGCAGCTGCAGCCTTTAT  
676 GGGTTGCAGGGAGTTGTATGTAATGCTCAGAAAGAGCTGCCACTGAGAAT  
677

678 >TESPA1\_WT  
679 TTCAATGATGTGGGCTGATTAGAACATAGCTGAAAGCAGGTGTTGGGATATTGATTTCCATGGCT  
680 GGTCTCACCTGTTACAAAACCTTCTACTACAATGAGTTTCAAACCTTCAATATGCAATCAATTATCTA  
681 ACCTAAAGATCTTGGTAAAACCTGTGATTCATTAGGTCTGGGGTGGGGGCTG  
682

683 >IGF2BP2\_TESPA1\_Fusion  
684 TTCAATGATGTGGGCTGATTAGAACATAGCTGAAAGCAGGTGTTGGGATATTGATTTCCATGGCT  
685 GGTCTCACCTGTTACAAAACCTTCTACTACTGTTTGTAAATAACTTTAAAGCAGCTGCAGCCTTT  
686 ATGGGTTGCAGGGAGTTGTATGTAATGCTCAGAAAGAGCTGCCACTGAGAAT  
687

688 Reads mapping unambiguously to one of those reference sequences were then attributed to  
689 the clone to which their cell barcode belonged.  
690

## 691 Data and code availability

692 The raw sequencing files reported in this study have been deposited in the European  
693 Genome-phenome Archive (EGA) under the accession number EGAS00001006807. The  
694 software used to analyze the data of this study has been deposited at the GitHub repository:  
695 <https://github.com/cbg-ethz/scIsoPrep>

## 696 Acknowledgements

697 We thank Ina Nissen (Genomics Facility Basel) for technical support with PacBio sequencing,  
698 Ching-Yeu Liang (University Hospital Basel) for assistance with sample dissociation, and Anne  
699 Bertolini for her help with the scAmp R package. We also wish to thank Elisabeth Tseng  
700 (Pacific Bioscience) for her help with cDNA\_Cupcake, SQANTI3 and long-read sequencing  
701 analysis, Anoushka Joglekar (Weill Cornell Medicine) for her help with ScisorSeq and  
702 ScisorWiz as well as Lucia Csepregi for her help with immune repertoire analysis.  
703 Furthermore, we are grateful to Adam Frankish (Wellcome Sanger Institute) and the  
704 GENCODE Project team for their help with the manual annotation of transcripts. Graphical



705 illustrations were created with BioRender.com. Illumina sequencing was carried out in the  
706 Genomics Facility Basel of the University of Basel and the Department of Biosystems Science  
707 and Engineering, ETHZI. PacBio sequencing was done in the Functional Genomics Center  
708 Zurich of the University of Zurich and ETHZ, and the Lausanne Genomic Technologies  
709 Facility, University of Lausanne.

## 710 Funding Information

711 Part of this work was funded by the SNSF SPARK grant #190413, the grant #2017-510 of the  
712 Strategic Focal Area “Personalized Health and Related Technologies (PHRT)” of the ETH  
713 Domain, and the European Union’s Horizon 2020 research and innovation programme under  
714 the Marie Skłodowska-Curie grant agreement #766030.

## 715 Author contributions

716 UL and CB acquired funding and conceived and designed the experiments. VHS provided  
717 patient material. UL, FJ, AD, and CB selected the clinical cohort. UL performed 10x Genomics  
718 sample processing as well as short-read and long-read sequencing library preparation. The  
719 Tumor Profiler Consortium provided scDNA-seq data and NGS panel results. AD designed  
720 the analysis pipeline, and implemented it with the help of NBo. AD conducted all computational  
721 analyses. FS assisted in short-read scRNA-seq analysis. AD, UL, FJ, NBe, and CB interpreted  
722 the data. AD and UL wrote the manuscript with contributions of all authors. All authors read  
723 and approved the final manuscript.

## 724 Tumor Profiler Consortium authors list

725 Rudolf Aebersold<sup>2</sup>, Melike Ak<sup>28</sup>, Faisal S Al-Quaddoomi<sup>9,17</sup>, Silvana I Albert<sup>7</sup>, Jonas Albinus<sup>7</sup>,  
726 Ilaria Alborelli<sup>24</sup>, Sonali Andani<sup>6,17,26,31</sup>, Per-Olof Attinger<sup>11</sup>, Marina Bacac<sup>16</sup>, Daniel  
727 Baumhoer<sup>24</sup>, Beatrice Beck-Schimmer<sup>39</sup>, Niko Beerenwinkel<sup>4,17</sup>, Christian Beisel<sup>4</sup>, Lara  
728 Bernasconi<sup>27</sup>, Anne Bertolini<sup>9,17</sup>, Bernd Bodenmiller<sup>8,35</sup>, Ximena Bonilla<sup>6,17,26</sup>, Lars  
729 Bosshard<sup>9,17</sup>, Byron Calgua<sup>24</sup>, Ruben Casanova<sup>35</sup>, Stéphane Chevrier<sup>35</sup>, Natalia  
730 Chicherova<sup>9,17</sup>, Ricardo Coelho<sup>18</sup>, Maya D’Costa<sup>10</sup>, Esther Danenberg<sup>37</sup>, Natalie  
731 Davidson<sup>6,17,26</sup>, Monica-Andreea Drăgan<sup>4</sup>, Reinhard Dummer<sup>28</sup>, Stefanie Engler<sup>35</sup>, Martin  
732 Erkens<sup>14</sup>, Katja Eschbach<sup>4</sup>, Cinzia Esposito<sup>37</sup>, André Fedier<sup>18</sup>, Pedro Ferreira<sup>4</sup>, Joanna  
733 Ficek<sup>6,17,26</sup>, Anja L Frei<sup>31</sup>, Bruno Frey<sup>13</sup>, Sandra Goetze<sup>7</sup>, Linda Grob<sup>9,17</sup>, Gabriele Gut<sup>37</sup>,  
734 Detlef Günther<sup>5</sup>, Martina Haberecker<sup>31</sup>, Pirmin Haeuptle<sup>1</sup>, Viola Heinzelmann-Schwarz<sup>18,23</sup>,  
735 Sylvia Herter<sup>16</sup>, Rene Holtackers<sup>37</sup>, Tamara Huesser<sup>16</sup>, Alexander Immer<sup>6,12</sup>, Anja Irmisch<sup>28</sup>,  
736 Francis Jacob<sup>18</sup>, Andrea Jacobs<sup>35</sup>, Tim M Jaeger<sup>11</sup>, Katharina Jahn<sup>4</sup>, Alva R James<sup>6,17,26</sup>,  
737 Philip M Jermann<sup>24</sup>, André Kahles<sup>6,17,26</sup>, Abdullah Kahraman<sup>17,31</sup>, Viktor H Koelzer<sup>31</sup>, Werner  
738 Kuebler<sup>25</sup>, Jack Kuipers<sup>4,17</sup>, Christian P Kunze<sup>22</sup>, Christian Kurzeder<sup>21</sup>, Kjong-Van  
739 Lehmann<sup>6,17,26</sup>, Mitchell Levesque<sup>28</sup>, Ulrike Lischetti<sup>18</sup>, Sebastian Lugert<sup>10</sup>, Gerd Maass<sup>13</sup>,  
740 Markus G Manz<sup>30</sup>, Philipp Markolin<sup>6,17,26</sup>, Martin Mehnert<sup>7</sup>, Julien Mena<sup>2</sup>, Julian M Metzler<sup>29</sup>,  
741 Nicola Miglino<sup>1</sup>, Emanuela S Milani<sup>7</sup>, Holger Moch<sup>31</sup>, Simone Muenst<sup>24</sup>, Riccardo Murri<sup>38</sup>,  
742 Charlotte KY Ng<sup>24,34</sup>, Stefan Nicolet<sup>24</sup>, Marta Nowak<sup>31</sup>, Monica Nunez Lopez<sup>18</sup>, Patrick GA  
743 Pedrioli<sup>3</sup>, Lucas Pelkmans<sup>37</sup>, Salvatore Piscuoglio<sup>18,24</sup>, Michael Prummer<sup>9,17</sup>, Natalie  
744 Rimmer<sup>18</sup>, Mathilde Ritter<sup>18</sup>, Christian Rommel<sup>14</sup>, María L Rosano-González<sup>9,17</sup>, Gunnar  
745 Rättsch<sup>3,6,17,26</sup>, Natascha Santacroce<sup>4</sup>, Jacobo Sarabia del Castillo<sup>37</sup>, Ramona Schlenker<sup>15</sup>,

746 Petra C Schwalie<sup>14</sup>, Severin Schwan<sup>11</sup>, Tobias Schär<sup>4</sup>, Gabriela Senti<sup>27</sup>, Wenguang Shao<sup>7</sup>,  
747 Franziska Singer<sup>9,17</sup>, Sujana Sivapatham<sup>35</sup>, Berend Snijder<sup>2,17</sup>, Bettina Sobottka<sup>31</sup>, Vipin T  
748 Sreedharan<sup>9,17</sup>, Stefan Stark<sup>6,17,26</sup>, Daniel J Stekhoven<sup>9,17</sup>, Tanmay Tanna<sup>4,6</sup>, Alexandre PA  
749 Theocharides<sup>30</sup>, Tinu M Thomas<sup>6,17,26</sup>, Markus Tolnay<sup>24</sup>, Vinko Tosevski<sup>16</sup>, Nora C  
750 Toussaint<sup>9,17</sup>, Mustafa A Tuncel<sup>4,17</sup>, Marina Tusup<sup>28</sup>, Audrey Van Drogen<sup>7</sup>, Marcus Vetter<sup>20</sup>,  
751 Tatjana Vlajnic<sup>24</sup>, Sandra Weber<sup>27</sup>, Walter P Weber<sup>19</sup>, Rebekka Wegmann<sup>2</sup>, Michael  
752 Weller<sup>33</sup>, Fabian Wendt<sup>7</sup>, Norbert Wey<sup>31</sup>, Andreas Wicki<sup>30,36</sup>, Mattheus HE Wildschut<sup>2,30</sup>,  
753 Bernd Wollscheid<sup>7</sup>, Shuqing Yu<sup>9,17</sup>, Johanna Ziegler<sup>28</sup>, Marc Zimmermann<sup>6,17,26</sup>, Martin  
754 Zoche<sup>31</sup>, Gregor Zuend<sup>32</sup>

755 <sup>1</sup>Cantonal Hospital Baselland, Medical University Clinic, Rheinstrasse 26, 4410 Liestal,  
756 Switzerland, <sup>2</sup>ETH Zurich, Department of Biology, Institute of Molecular Systems Biology,  
757 Otto-Stern-Weg 3, 8093 Zurich, Switzerland, <sup>3</sup>ETH Zurich, Department of Biology, Wolfgang-  
758 Pauli-Strasse 27, 8093 Zurich, Switzerland, <sup>4</sup>ETH Zurich, Department of Biosystems Science  
759 and Engineering, Mattenstrasse 26, 4058 Basel, Switzerland, <sup>5</sup>ETH Zurich, Department of  
760 Chemistry and Applied Biosciences, Vladimir-Prelog-Weg 1-5/10, 8093 Zurich, Switzerland,  
761 <sup>6</sup>ETH Zurich, Department of Computer Science, Institute of Machine Learning,  
762 Universitätstrasse 6, 8092 Zurich, Switzerland, <sup>7</sup>ETH Zurich, Department of Health Sciences  
763 and Technology, Otto-Stern-Weg 3, 8093 Zurich, Switzerland, <sup>8</sup>ETH Zurich, Institute of  
764 Molecular Health Sciences, Otto-Stern-Weg 7, 8093 Zurich, Switzerland, <sup>9</sup>ETH Zurich,  
765 NEXUS Personalized Health Technologies, John-von-Neumann-Weg 9, 8093 Zurich,  
766 Switzerland, <sup>10</sup>F. Hoffmann-La Roche Ltd, Grenzacherstrasse 124, 4070 Basel, Switzerland,  
767 <sup>11</sup>F. Hoffmann-La Roche Ltd, Grenzacherstrasse 124, 4070 Basel, Switzerland, , <sup>12</sup>Max  
768 Planck ETH Center for Learning Systems, , <sup>13</sup>Roche Diagnostics GmbH, Nonnenwald 2,  
769 82377 Penzberg, Germany, <sup>14</sup>Roche Pharmaceutical Research and Early Development,  
770 Roche Innovation Center Basel, Grenzacherstrasse 124, 4070 Basel, Switzerland, <sup>15</sup>Roche  
771 Pharmaceutical Research and Early Development, Roche Innovation Center Munich, Roche  
772 Diagnostics GmbH, Nonnenwald 2, 82377 Penzberg, Germany, <sup>16</sup>Roche Pharmaceutical  
773 Research and Early Development, Roche Innovation Center Zurich, Wagistrasse 10, 8952  
774 Schlieren, Switzerland, <sup>17</sup>SIB Swiss Institute of Bioinformatics, Lausanne, Switzerland,  
775 <sup>18</sup>University Hospital Basel and University of Basel, Department of Biomedicine,  
776 Hebelstrasse 20, 4031 Basel, Switzerland, <sup>19</sup>University Hospital Basel and University of  
777 Basel, Department of Surgery, Brustzentrum, Spitalstrasse 21, 4031 Basel, Switzerland,  
778 <sup>20</sup>University Hospital Basel, Brustzentrum & Tumorzentrum, Petersgraben 4, 4031 Basel,  
779 Switzerland, <sup>21</sup>University Hospital Basel, Brustzentrum, Spitalstrasse 21, 4031 Basel,  
780 Switzerland, <sup>22</sup>University Hospital Basel, Department of Information- and Communication  
781 Technology, Spitalstrasse 26, 4031 Basel, Switzerland, <sup>23</sup>University Hospital Basel,  
782 Gynecological Cancer Center, Spitalstrasse 21, 4031 Basel, Switzerland, <sup>24</sup>University  
783 Hospital Basel, Institute of Medical Genetics and Pathology, Schönbeinstrasse 40, 4031  
784 Basel, Switzerland, <sup>25</sup>University Hospital Basel, Spitalstrasse 21/Petersgraben 4, 4031  
785 Basel, Switzerland, <sup>26</sup>University Hospital Zurich, Biomedical Informatics,  
786 Schmelzbergstrasse 26, 8006 Zurich, Switzerland, <sup>27</sup>University Hospital Zurich, Clinical  
787 Trials Center, Rämistrasse 100, 8091 Zurich, Switzerland, <sup>28</sup>University Hospital Zurich,  
788 Department of Dermatology, Gloriastrasse 31, 8091 Zurich, Switzerland, <sup>29</sup>University  
789 Hospital Zurich, Department of Gynecology, Frauenklinikstrasse 10, 8091 Zurich,  
790 Switzerland, <sup>30</sup>University Hospital Zurich, Department of Medical Oncology and Hematology,  
791 Rämistrasse 100, 8091 Zurich, Switzerland, <sup>31</sup>University Hospital Zurich, Department of  
792 Pathology and Molecular Pathology, Schmelzbergstrasse 12, 8091 Zurich, Switzerland,

793 <sup>32</sup>University Hospital Zurich, Rämistrasse 100, 8091 Zurich, Switzerland, <sup>33</sup>University  
794 Hospital and University of Zurich, Department of Neurology, Frauenklinikstrasse 26, 8091  
795 Zurich, Switzerland, <sup>34</sup>University of Bern, Department of BioMedical Research,  
796 Murtenstrasse 35, 3008 Bern, Switzerland, <sup>35</sup>University of Zurich, Department of Quantitative  
797 Biomedicine, Winterthurerstrasse 190, 8057 Zurich, Switzerland, <sup>36</sup>University of Zurich,  
798 Faculty of Medicine, Zurich, Switzerland, <sup>37</sup>University of Zurich, Institute of Molecular Life  
799 Sciences, Winterthurerstrasse 190, 8057 Zurich, Switzerland, <sup>38</sup>University of Zurich,  
800 Services and Support for Science IT, Winterthurerstrasse 190, 8057 Zurich, Switzerland,  
801 <sup>39</sup>University of Zurich, VP Medicine, Künstlergasse 15, 8001 Zurich, Switzerland

## 802 Conflict of interest

803 The authors declare no competing interests.

804

## 805 ORCiDs

806 Ulrike Lischetti: 0000-0002-9956-3043

807 Arthur Dondi: 0000-0003-3234-2550

808 Francis Jacob: 0000-0002-0446-1942

809 Christian Beisel: 0000-0001-5360-2193

810 Niko Beerenwinkel: 0000-0002-0573-6119

811 Franziska Singer: 0000-0002-6017-1595

812 Nico Borgsmüller: 0000-0003-4073-3877

813 Viola Heizelmann: 0000-0002-4056-3225

814

## 815 References

816 1. Garraway, L. A. & Lander, E. S. Lessons from the cancer genome. *Cell* **153**, 17–37  
817 (2013).

818 2. Hanahan, D. Hallmarks of cancer: new dimensions. *Cancer Discov.* **12**, 31–46 (2022).

819 3. Hertzman Johansson, C. & Egyhazi Brage, S. BRAF inhibitors in cancer therapy.  
820 *Pharmacol. Ther.* **142**, 176–182 (2014).

821 4. Li, J. *et al.* A functional genomic approach to actionable gene fusions for precision  
822 oncology. *Sci. Adv.* **8**, eabm2382 (2022).

823 5. Schumacher, T. N. & Schreiber, R. D. Neoantigens in cancer immunotherapy. *Science*  
824 **348**, 69–74 (2015).

825 6. Yu, Y.-P. *et al.* Identification of recurrent fusion genes across multiple cancer types. *Sci.*

- 826        *Rep.* **9**, 1074 (2019).
- 827    7. Bower, H. *et al.* Life expectancy of patients with chronic myeloid leukemia approaches  
828        the life expectancy of the general population. *J. Clin. Oncol.* **34**, 2851–2857 (2016).
- 829    8. Khan, M. *et al.* ALK inhibitors in the treatment of ALK positive NSCLC. *Front. Oncol.* **8**,  
830        557 (2018).
- 831    9. Alexandrov, L. B. *et al.* Signatures of mutational processes in human cancer. *Nature* **500**,  
832        415–421 (2013).
- 833    10. Zack, T. I. *et al.* Pan-cancer patterns of somatic copy number alteration. *Nat. Genet.* **45**,  
834        1134–1140 (2013).
- 835    11. Gao, Q. *et al.* Driver fusions and their implications in the development and treatment of  
836        human cancers. *Cell Rep.* **23**, 227–238.e3 (2018).
- 837    12. Wang, Y., Shi, T., Song, X., Liu, B. & Wei, J. Gene fusion neoantigens: Emerging targets  
838        for cancer immunotherapy. *Cancer Lett.* **506**, 45–54 (2021).
- 839    13. Wei, Z. *et al.* The Landscape of Tumor Fusion Neoantigens: A Pan-Cancer Analysis.  
840        *iScience* **21**, 249–260 (2019).
- 841    14. Blencowe, B. J. Alternative splicing: new insights from global analyses. *Cell* **126**, 37–47  
842        (2006).
- 843    15. Bonnal, S. C., López-Oreja, I. & Valcárcel, J. Roles and mechanisms of alternative  
844        splicing in cancer - implications for care. *Nat. Rev. Clin. Oncol.* **17**, 457–474 (2020).
- 845    16. Kahles, A. *et al.* Comprehensive Analysis of Alternative Splicing Across Tumors from  
846        8,705 Patients. *Cancer Cell* **34**, 211–224.e6 (2018).
- 847    17. Marusyk, A. & Polyak, K. Tumor heterogeneity: causes and consequences. *Biochim.*  
848        *Biophys. Acta* **1805**, 105–117 (2010).
- 849    18. Whiteside, T. L. The tumor microenvironment and its role in promoting tumor growth.  
850        *Oncogene* **27**, 5904–5912 (2008).
- 851    19. Hedlund, E. & Deng, Q. Single-cell RNA sequencing: Technical advancements and  
852        biological applications. *Mol. Aspects Med.* **59**, 36–46 (2018).
- 853    20. Glinos, D. A. *et al.* Transcriptome variation in human tissues revealed by long-read

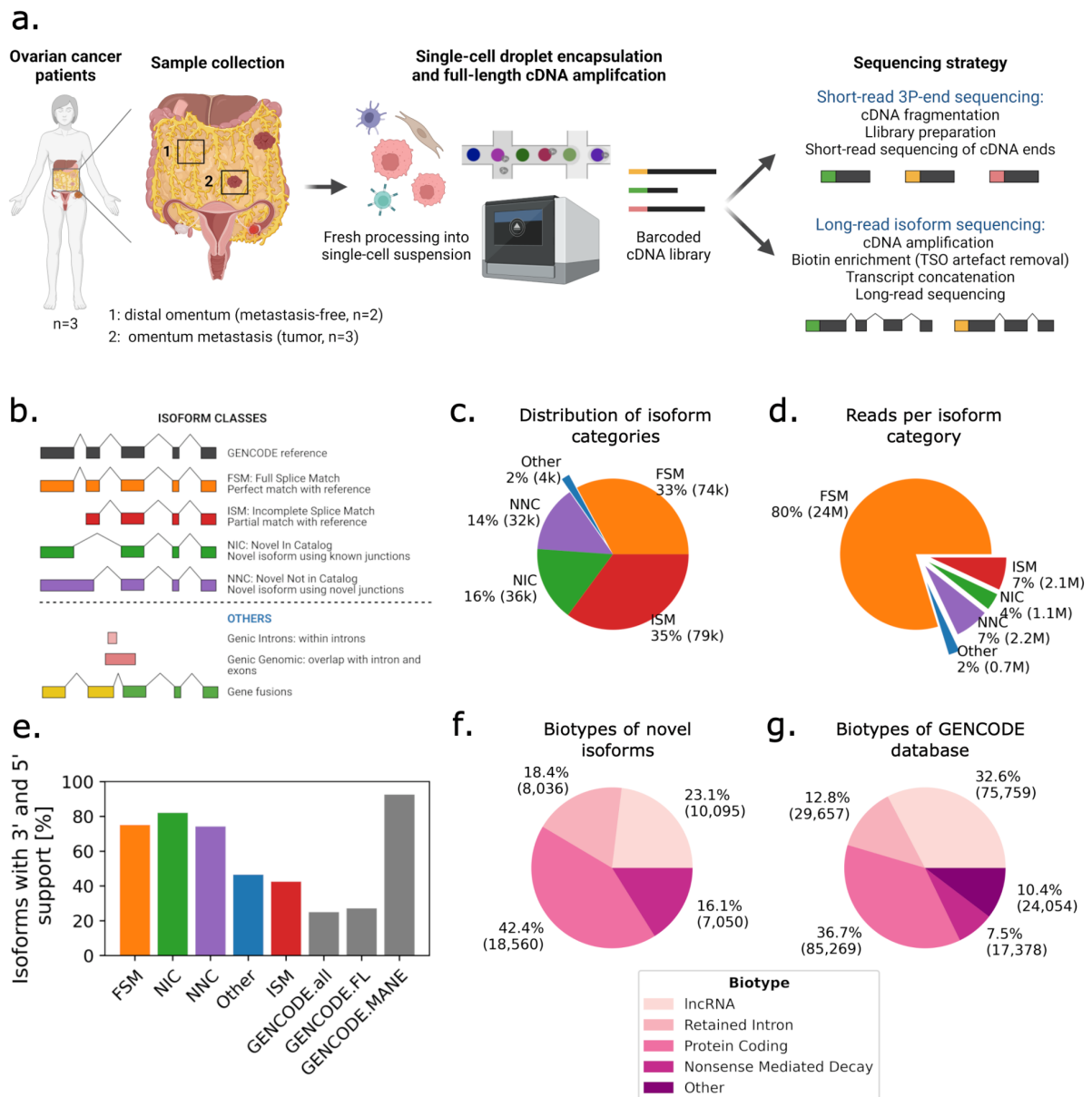
- 854 sequencing. *Nature* **608**, 353–359 (2022).
- 855 21. Philpott, M. *et al.* Nanopore sequencing of single-cell transcriptomes with scCOLOR-seq.  
856 *Nat. Biotechnol.* **39**, 1517–1520 (2021).
- 857 22. Gupta, I. *et al.* Single-cell isoform RNA sequencing characterizes isoforms in thousands  
858 of cerebellar cells. *Nat. Biotechnol.* **36**, 1197–1202 (2018).
- 859 23. Zheng, Y.-F. *et al.* HIT-sclSOseq: High-throughput and High-accuracy Single-cell Full-  
860 length Isoform Sequencing for Corneal Epithelium. *BioRxiv* (2020)  
861 doi:10.1101/2020.07.27.222349.
- 862 24. Al'Khafaji, A. M. *et al.* High-throughput RNA isoform sequencing using programmable  
863 cDNA concatenation. *BioRxiv* (2021) doi:10.1101/2021.10.01.462818.
- 864 25. Joglekar, A. *et al.* A spatially resolved brain region- and cell type-specific isoform atlas of  
865 the postnatal mouse brain. *Nat. Commun.* **12**, 463 (2021).
- 866 26. Hardwick, S. A. *et al.* Single-nuclei isoform RNA sequencing unlocks barcoded exon  
867 connectivity in frozen brain tissue. *Nat. Biotechnol.* **40**, 1082–1092 (2022).
- 868 27. Veiga, D. F. T. *et al.* A comprehensive long-read isoform analysis platform and  
869 sequencing resource for breast cancer. *Sci. Adv.* **8**, eabg6711 (2022).
- 870 28. Namba, S. *et al.* Transcript-targeted analysis reveals isoform alterations and double-hop  
871 fusions in breast cancer. *Commun. Biol.* **4**, 1320 (2021).
- 872 29. Macintyre, G. *et al.* Copy number signatures and mutational processes in ovarian  
873 carcinoma. *Nat. Genet.* **50**, 1262–1270 (2018).
- 874 30. Cancer Genome Atlas Research Network. Integrated genomic analyses of ovarian  
875 carcinoma. *Nature* **474**, 609–615 (2011).
- 876 31. Tardaguila, M. *et al.* SQANTI: extensive characterization of long-read transcript  
877 sequences for quality control in full-length transcriptome identification and quantification.  
878 *Genome Res.* (2018) doi:10.1101/gr.222976.117.
- 879 32. Abugessaisa, I. *et al.* FANTOM5 CAGE profiles of human and mouse reprocessed for  
880 GRCh38 and GRCm38 genome assemblies. *Sci. Data* **4**, 170107 (2017).
- 881 33. Herrmann, C. J. *et al.* PolyASite 2.0: a consolidated atlas of polyadenylation sites from 3'



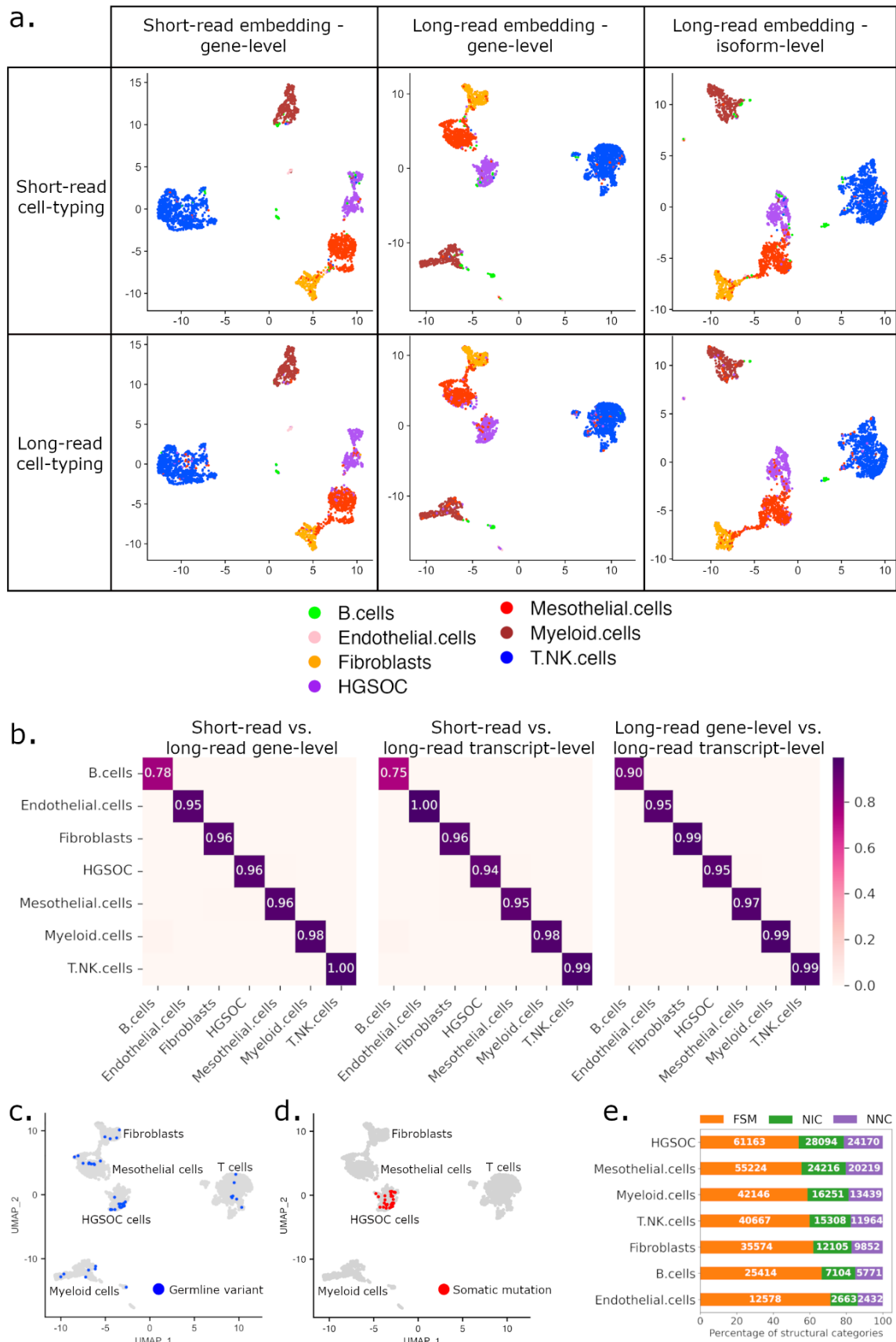
- 882 end sequencing. *Nucleic Acids Res.* **48**, D174–D179 (2020).
- 883 34. Morales, J. *et al.* A joint NCBI and EMBL-EBI transcript set for clinical genomics and  
884 research. *Nature* **604**, 310–315 (2022).
- 885 35. Becht, E. *et al.* Dimensionality reduction for visualizing single-cell data using UMAP. *Nat.*  
886 *Biotechnol.* **37**, 38–44 (2018).
- 887 36. García-Bartolomé, A. *et al.* Altered Expression Ratio of Actin-Binding Gelsolin Isoforms  
888 Is a Novel Hallmark of Mitochondrial OXPHOS Dysfunction. *Cells* **9**, (2020).
- 889 37. Liefers-Visser, J. A. L., Meijering, R. A. M., Reyners, A. K. L., van der Zee, A. G. J. & de  
890 Jong, S. IGF system targeted therapy: Therapeutic opportunities for ovarian cancer.  
891 *Cancer Treat. Rev.* **60**, 90–99 (2017).
- 892 38. Philippou, A., Maridaki, M., Pneumaticos, S. & Koutsilieris, M. The complexity of the IGF1  
893 gene splicing, posttranslational modification and bioactivity. *Mol. Med.* **20**, 202–214  
894 (2014).
- 895 39. Choy, J. Y. H., Boon, P. L. S., Bertin, N. & Fullwood, M. J. A resource of ribosomal RNA-  
896 depleted RNA-Seq data from different normal adult and fetal human tissues. *Sci. Data* **2**,  
897 150063 (2015).
- 898 40. Martínez-Jiménez, F. *et al.* A compendium of mutational cancer driver genes. *Nat. Rev.*  
899 *Cancer* **20**, 555–572 (2020).
- 900 41. Tate, J. G. *et al.* COSMIC: the catalogue of somatic mutations in cancer. *Nucleic Acids*  
901 *Res.* **47**, D941–D947 (2019).
- 902 42. Mertens, F., Johansson, B., Fioretos, T. & Mitelman, F. The emerging complexity of gene  
903 fusions in cancer. *Nat. Rev. Cancer* **15**, 371–381 (2015).
- 904 43. Yang, W. *et al.* Immunogenic neoantigens derived from gene fusions stimulate T cell  
905 responses. *Nat. Med.* **25**, 767–775 (2019).
- 906 44. Berger, M. F. & Mardis, E. R. The emerging clinical relevance of genomics in cancer  
907 medicine. *Nat. Rev. Clin. Oncol.* **15**, 353–365 (2018).
- 908 45. Jing, Y., Han, Z., Zhang, S., Liu, Y. & Wei, L. Epithelial-Mesenchymal Transition in tumor  
909 microenvironment. *Cell Biosci.* **1**, 29 (2011).

- 910 46. Kenny, H. A. *et al.* Mesothelial cells promote early ovarian cancer metastasis through  
911 fibronectin secretion. *J. Clin. Invest.* **124**, 4614–4628 (2014).
- 912 47. Hoshino, A. *et al.* Tumour exosome integrins determine organotropic metastasis. *Nature*  
913 **527**, 329–335 (2015).
- 914 48. Asare-Werehene, M. *et al.* The exosome-mediated autocrine and paracrine actions of  
915 plasma gelsolin in ovarian cancer chemoresistance. *Oncogene* **39**, 1600–1616 (2020).
- 916 49. Yang, Y. *et al.* Tumor Suppressor microRNA-138 Suppresses Low-Grade Glioma  
917 Development and Metastasis via Regulating IGF2BP2. *Onco Targets Ther* **13**, 2247–  
918 2260 (2020).
- 919 50. Wang, D. *et al.* Tespa1 is involved in late thymocyte development through the regulation  
920 of TCR-mediated signaling. *Nat. Immunol.* **13**, 560–568 (2012).
- 921 51. Brokaw, J. *et al.* IGF-I in epithelial ovarian cancer and its role in disease progression.  
922 *Growth Factors* **25**, 346–354 (2007).
- 923 52. Volden, R. & Vollmers, C. Single-cell isoform analysis in human immune cells. *Genome*  
924 *Biol.* **23**, 47 (2022).
- 925 53. Dutton, G. CRISPR-Cas9 Technology Cuts Clutter from Sequencing Libraries. *Genetic*  
926 *Engineering & Biotechnology News* **41**, 24–25 (2021).
- 927 54. Lang, F., Schrörs, B., Löwer, M., Türeci, Ö. & Sahin, U. Identification of neoantigens for  
928 individualized therapeutic cancer vaccines. *Nat. Rev. Drug Discov.* **21**, 261–282 (2022).
- 929 55. Lin, M. J. *et al.* Cancer vaccines: the next immunotherapy frontier. *Nat. Cancer* **3**, 911–  
930 926 (2022).
- 931 56. Hebelstrup, K. H. *et al.* UCE: A uracil excision (USER)-based toolbox for transformation  
932 of cereals. *Plant Methods* **6**, 15 (2010).
- 933 57. Griffiths, J. A., Richard, A. C., Bach, K., Lun, A. T. L. & Marioni, J. C. Detection and  
934 removal of barcode swapping in single-cell RNA-seq data. *Nat. Commun.* **9**, 2667 (2018).
- 935 58. Bertolini, A. *et al.* scAmpI-A versatile pipeline for single-cell RNA-seq analysis from basics  
936 to clinics. *PLoS Comput. Biol.* **18**, e1010097 (2022).
- 937 59. Germain, P.-L., Lun, A., Garcia Meixide, C., Macnair, W. & Robinson, M. D. Doublet

- 938 identification in single-cell sequencing data using scDbIFinder. *F1000Res.* **10**, 979  
939 (2021).
- 940 60. Hafemeister, C. & Satija, R. Normalization and variance stabilization of single-cell RNA-  
941 seq data using regularized negative binomial regression. *Genome Biol.* **20**, 296 (2019).
- 942 61. Levine, J. H. *et al.* Data-Driven Phenotypic Dissection of AML Reveals Progenitor-like  
943 Cells that Correlate with Prognosis. *Cell* **162**, 184–197 (2015).
- 944 62. Prummer, M. *et al.* scROSHI - robust supervised hierarchical identification of single cells.  
945 *BioRxiv* (2022) doi:10.1101/2022.04.05.487176.
- 946 63. Irmisch, A. *et al.* The Tumor Profiler Study: integrated, multi-omic, functional tumor  
947 profiling for clinical decision support. *Cancer Cell* **39**, 288–293 (2021).
- 948 64. Hao, Y. *et al.* Integrated analysis of multimodal single-cell data. *Cell* **184**, 3573-3587.e29  
949 (2021).
- 950 65. Newman, A. M. *et al.* Robust enumeration of cell subsets from tissue expression profiles.  
951 *Nat. Methods* **12**, 453–457 (2015).
- 952 66. Sade-Feldman, M. *et al.* Defining T Cell States Associated with Response to Checkpoint  
953 Immunotherapy in Melanoma. *Cell* **175**, 998-1013.e20 (2018).
- 954 67. Köster, J. & Rahmann, S. Snakemake--a scalable bioinformatics workflow engine.  
955 *Bioinformatics* **28**, 2520–2522 (2012).
- 956 68. Frankish, A. *et al.* GENCODE reference annotation for the human and mouse genomes.  
957 *Nucleic Acids Res.* **47**, D766–D773 (2019).
- 958 69. Stein, A. N., Joglekar, A., Poon, C.-L. & Tilgner, H. U. ScisorWiz: Visualizing Differential  
959 Isoform Expression in Single-Cell Long-Read Data. *BioRxiv* (2022)  
960 doi:10.1101/2022.04.14.488347.
- 961 70. Bertolini, A. *et al.* scAmpi - A versatile pipeline for single-cell RNA-seq analysis from  
962 basics to clinics. *BioRxiv* (2021) doi:10.1101/2021.03.25.437054.
- 963 71. Kuipers, J., Tuncel, M. A., Ferreira, P., Jahn, K. & Beerenwinkel, N. Single-cell copy  
964 number calling and event history reconstruction. *BioRxiv* (2020)  
965 doi:10.1101/2020.04.28.065755.



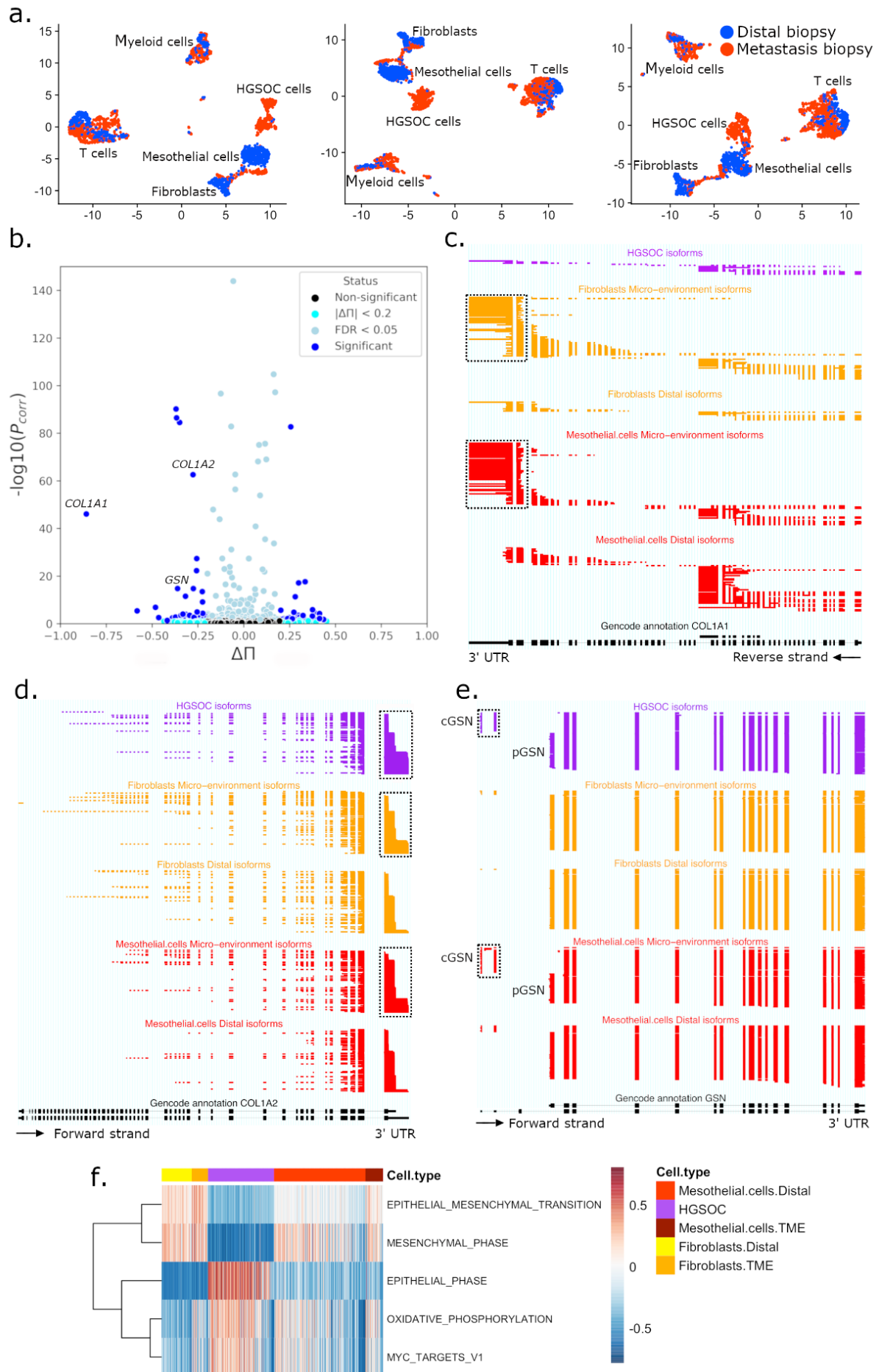
**Figure 1: Study design and long read data overview. (a)** Schematic of freshly processed HGSOc metastasis and patient-matched tumor-free omentum tissue biopsies, scRNA-seq. **(b)** Definition of SQANTI-defined isoform structural categories. **(c)** Proportions of isoform structural categories detected in merged metastasis and healthy omentum samples. Percentage and total number of isoforms per category are indicated. **(d)** Proportions of unique reads attributed to isoforms detected in **(c)**. Percentage and total number of UMIs per category are indicated. **(e)** Percentage of isoforms for which transcription start site is supported by CAGE (FANTOM5) data and transcription termination site is supported by polyA (PolyASite) data, per isoform structural categories. GENCODE.all indicates all protein-coding isoforms in the GENCODE database, GENCODE.FL is a subset of GENCODE.full containing only isoforms tagged as full-length, and GENCODE.MANE is a subset of canonical transcripts, one per human protein coding locus. **(f)** GENCODE defined biotypes composition of novel isoforms. **(g)** Biotypes composition of the GENCODE database.



**Figure 2: Clustering and cell type specific isoform distribution.**

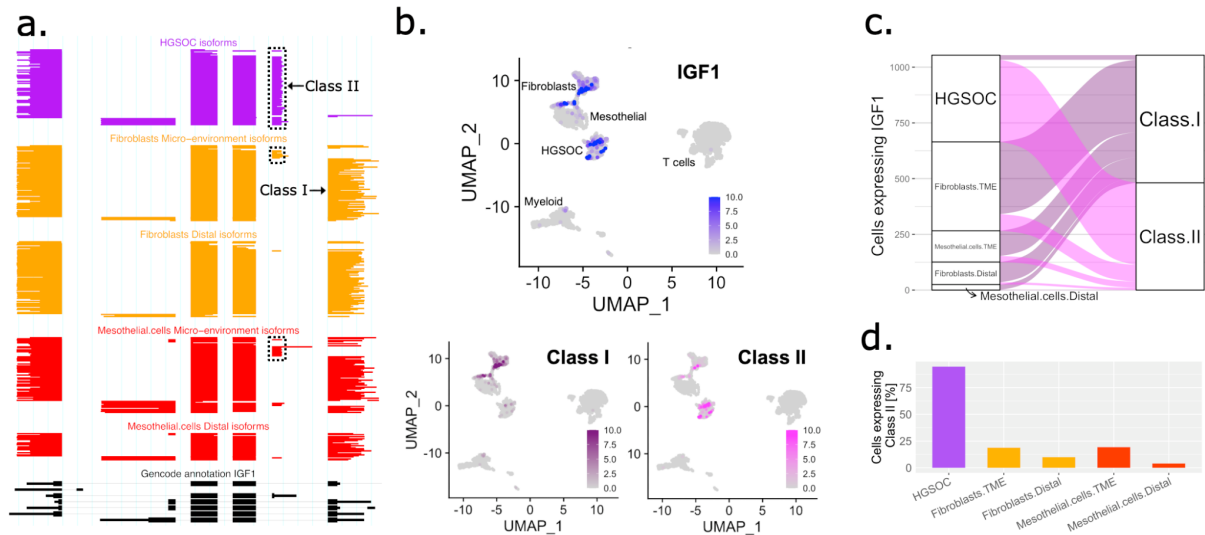


**(a)** Cohort UMAP embeddings by data types and automatic cell type annotation. Top and bottom rows: cell type labels based on short- and long-read data, respectively. Left column: embedding on short-read data - gene level, middle column: embedding on long-read data - gene level, right column: embedding on long-read data - isoform level. **(b)** Jaccard distance of cell populations in different UMAP embeddings: short-reads - gene level versus long-reads - gene level (left), short-reads - gene level versus long-reads - isoform level (middle), long-reads - gene level versus long-reads - isoform level (right). **(c)** Long-reads - gene level UMAP cohort visualizations of cells with at least one somatic mutation also found in bulk DNA. **(d)** Long-reads - gene level UMAP cohort visualization of cells with at least one germline variant. Germline variants are variants detected in healthy omentum distal samples. **(e)** SQANTI-defined structural category normalized distribution of isoforms detected per cell type (number of isoforms displayed in white).



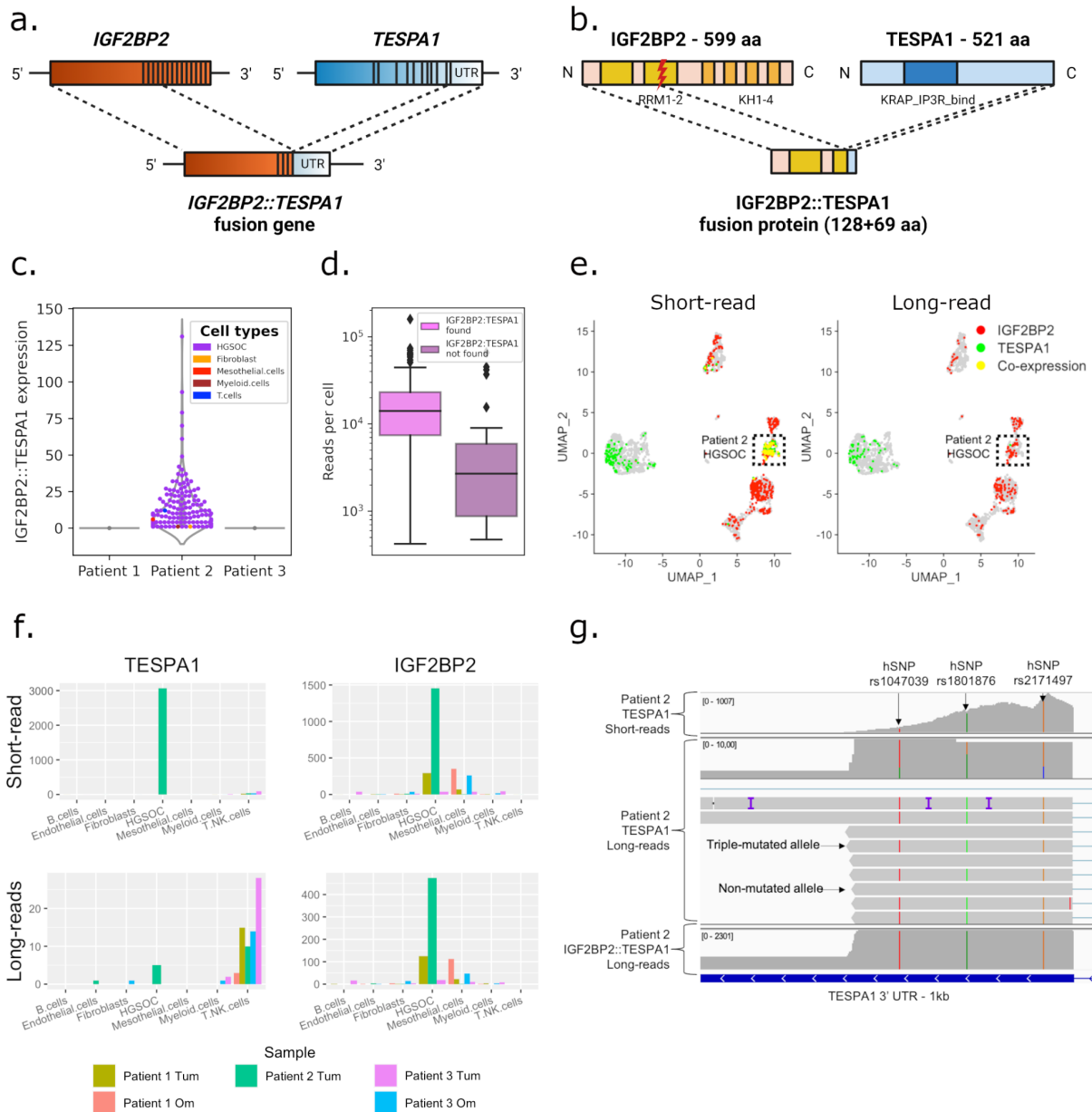
**Figure 3: Differential isoform expression in tumor microenvironment reveals epithelial-to-mesenchymal transition.**

**(a)** Cohort UMAPs embedding of short-read data - gene level (left), long-read data - gene level (middle), long-read data - isoform level (right), colored by tissue type. **(b)** Volcano plot of mesothelial TME vs. distal cells differential isoform usage. The X-axis represents the effect size in the gene, the Y-axis is the p-value derived from a  $\chi^2$  test corrected for multiple testing using the Benjamini–Hochberg method. **(c)** ScisorWiz representation of isoforms in *COL1A1*, each horizontal line represents a single read colored according to cell types. Dashed boxes highlight the use of the canonical 3' UTR in TME fibroblasts and mesothelial cells, while distal mesothelial cells use an earlier 3' exon termination. **(d)** ScisorWiz representation of isoforms in *COL1A2*. Dashed boxes highlight the 3'UTR, where TME and HGSOC cells differentially express a longer 3'UTR than distal cells. **(e)** ScisorWiz representation of isoforms in *GSN*. Dashed boxes highlight the TSS, where mesothelial TME and HGSOC cells differentially express the *cGSN* isoform, while mesothelial distal cells and fibroblasts use *pGSN*. **(f)** Gene set variation analysis (GSVA) scores for different cell types. Heatmap colors from blue to red represent low to high enrichment.



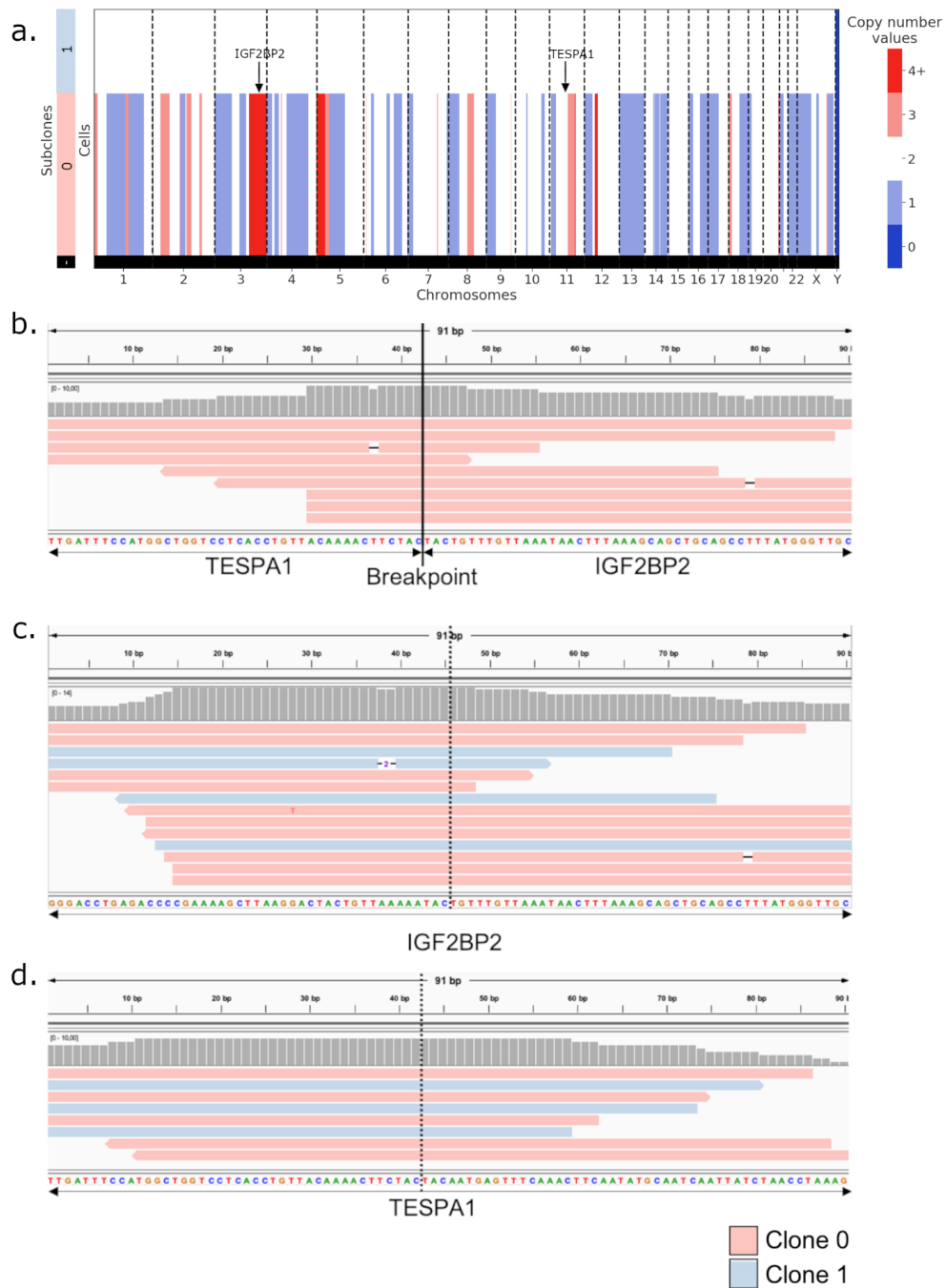
**Figure 4: Differential isoform expression of *IGF1* in tumor vs non-tumor cells.**

**(a)** ScisorWiz representation of isoforms in *IGF1*, each horizontal line represents a single isoform colored according to cell types. Colored areas are exons, and whitespace are intronic space, not drawn to scale. Exons are numbered according to the Gencode reference, Class I and II isoforms are isoforms with starting exons 1 and 2, respectively. Boxes highlight Class II expression in cancer and TME cells. **(b)** Projection of *IGF1* gene (top) and Class I/II isoform (bottom) expression on UMAP obtained from clustering on long-reads transcripts. **(c)** Alluvial plot of cells expressing *IGF1* in different cell types (left), divided between cells expressing Class I or II (right). **(d)** Barplot of percentage of cells expressing Class II isoform in different cell types and locations colored by cell type.



**Figure 5: Tumor and patient-specific detection of novel *IGF2BP2::TESPA1* gene fusion.** (a) Overview of wt *IGF2BP2*, wt *TESPA1* and gene fusions with exon structure. (b) Overview of wt *IGF2BP2*, wt *TESPA1* and fusion proteins and protein domains. RRM: RNA-recognition motif, KH: hnRNP K-homology domain, KRAP\_IP3R\_bind: Ki-ras-induced actin-interacting protein-IP3R-interacting domain. (c) Violin plot showing patient and tumor specific *IGF2BP2::TESPA1* fusion transcript detection in patient 2. (d) UMI count in fusion-containing vs -lacking patient 2 tumor cells. (e) scDNA copy-number profile clustering of the matched patient 2 sample. Subclone 0 (121 cells) exhibited multiple copy number alterations along its genome representing a single tumor clone, while subclone 1 (62 cells) had a diploid genome representing non-HGSOc cells. (f) patient 2 scDNA reads aligning to custom *IGF2BP2::TESPA1* gene fusion breakpoint reference. Only tumor subclone reads were found to align to it.





**Figure 6: *IGF2BP2::TESPA1* fusion breakpoint validation in scDNA.**

(a) Copy number values per subclone in Patient 2 scDNA. Subclone 0 has multiple copy number alterations, indicative of cancer, while Subclone 1 is copy-number neutral, non-cancer. (b) IGV view of scDNA reads aligning unambiguously to the *TESPA1::IGF2BP2* genomic breakpoint. In red, reads from Subclone 0 cells, in blue, reads from Subclone 1 cells. (c) IGV view of scDNA reads aligning unambiguously to wt *IGF2BP2*. The dashed line indicates the location of the (putative) breakpoint. (d) IGV view of scDNA reads aligning unambiguously to wt *TESPA1*. The dashed line indicates the location of the breakpoint.