# A single-cell atlas of breast cancer cell lines to study tumour heterogeneity and drug response.

#### 5 AUTHORS

Gambardella G<sup>1,2,\*</sup>, Viscido G<sup>1,2,\*</sup>, Tumaini B<sup>1</sup>, Isacchi A<sup>3</sup>, Bosotti R<sup>3</sup>, di Bernardo D<sup>1,2</sup>.

#### AFFILIATIONS

9 <sup>1</sup>Telethon Institute of Genetics and Medicine, Naples, Italy

<sup>2</sup>University of Naples Federico II, Department of Chemical Materials and Industrial Engineering, Naples, Italy

11 <sup>3</sup>NMSsrl, Nerviano Medical Sciences, 20014, Nerviano, Milan, Italy

12

3 4

6 7 8

# 1314 CORRESPONDENCE:

### 15 <u>dibernardo@tigem.it</u>

16

#### 17 \*These authors contributed equally to this work.

18

#### **ABSTRACT**

20 Breast cancer patient stratification is mainly driven by tumour receptor status and 21 histological grading and subtyping, with about twenty percent of patients for which absence 22 of any actionable biomarkers results in no clear therapeutic intervention. Cancer cells 23 within the same tumour have heterogeneous phenotypes and exhibit dynamic plasticity. 24 However, how to evaluate such heterogeneity and its impact on outcome and drug response 25 is still unclear. Here, we transcriptionally profiled 35,276 individual cells from 32 breast 26 cancer cell lines covering all main breast cancer subtypes to yield a breast cancer cell line 27 atlas. We found high degree of heterogeneity in the expression of clinically relevant 28 biomarkers across individual cells within the same cell line; such heterogeneity is non-29 genetic and dynamic. We computationally mapped single cell transcriptional profiles of 30 patients' tumour biopsies to the atlas to determine their composition in terms of cell lines. 31 Each tumour was found to be heterogenous and composed of multiple cell lines mostly, 32 but not exclusively, of the same subtype. We then trained an algorithm on the atlas to 33 determine cell line composition from bulk gene expression profiles of tumour biopsies, 34 thus providing a novel approach to patient stratification. Finally, we linked results from large-scale in vitro drug screening<sup>1,2</sup> to the single cell data to computationally predict 35 responses to more than 450 anticancer agents starting from single-cell transcriptional 36 37 profiles. We thus found that transcriptional heterogeneity enables cells with differential 38 drug sensitivity to co-exist in the same population. Our work provides a unique resource 39 and a novel framework to determine tumour heterogeneity and drug response in breast 40 cancer patients.

41

#### 42 MAIN TEXT

43

#### 44 Introduction

45 One of the main roadblocks to personalized medicine of cancer is the lack of biomarkers 46 to predict outcome and drug sensitivity from a tumour biopsy. Multigene assays such as

MammaPrint<sup>3</sup>, Oncotype DX<sup>4,5</sup> and PAM50<sup>6</sup> can classify Breast Cancer (BC) tumour types 47 48 and risk of relapse<sup>7</sup> but with limited clinical utility<sup>7,8</sup>. Genomic and transcriptional biomarkers of drug sensitivity are available only for a restricted number of drugs<sup>1,2,9</sup>. As a 49 50 consequence, BC patient stratification is still mainly driven by receptor status and histological grading and subtyping<sup>7</sup>, with about twenty percent<sup>10</sup> of patients for which 51 52 paucity of actionable biomarkers limits personalized therapies. Moreover, even when a 53 targeted treatment option is available, drug resistance may arise<sup>7</sup> partly because of rare 54 drug tolerant cells characterized by distinct transcriptional or mutational states<sup>11–17</sup>.

55 Determining tumour heterogeneity and its impact on drug response is essential to better 56 stratify patients and aid in the development of personalized therapies. Expression-based 57 biomarkers measured from bulk RNA-sequencing of a tumour biopsy are powerful 58 predictors of drug response in vitro<sup>1,2,18</sup>, but average out tumour heterogeneity. Single-cell transcriptomics yields a molecular profile of each cell<sup>19,20</sup>, however, it is still unclear if and 59 how it can inform clinical decision making. Here, we focused on tumour-derived breast 60 cancer cell lines. We hypothesized that despite being simplistic models of tumours, cancer 61 62 cell lines may exhibit themselves heterogeneous phenotypes, and serve as cell-state 63 "primitives" to deconvolve tumour cell composition from patients' biopsies for patient 64 stratification and prediction of drug response.

65

66

#### 67 **RESULTS**

68

#### 69 **1. Single-cell Transcriptome Profiling of Breast cancer cell lines.**

We performed single cell RNA-sequencing (scRNA-seq) of 31 breast cancer cell lines (Supplementary Table 01) and one non-cancer cell line, MCF12A<sup>21</sup>, by means of the Dropseq technology<sup>20</sup>. Following pre-processing (Methods), we retained a total of 35,276 cells, with an average of 1,069 cells per cell line and 3,248 genes captured per cell (Supplementary Figure 01 and Supplementary Table 01).

We next generated an atlas (<u>http://bcatlas.tigem.it</u>) encompassing the 32 BC cell lines, as shown in Figure 1A. In the atlas, luminal BC cell lines form a big "island" with multiple "peninsulas" with intermixing of cells from distinct cell lines; on the contrary, triple-negative breast cancer (TNBC) cell lines give rise to an "archipelago", where cells tend to separate into distinct islands according to the cell line of origin, thus suggesting that TNBC cell lines represent instances of distinct diseases.

Single-cell expression of clinically relevant biomarkers (Figure 1B,C) including
oestrogen receptor 1 (ESR1), progesterone receptor (PGR), Erb-B2 Receptor Tyrosine
Kinase 2 (ERBB2 a.k.a. HER2) and the epithelial growth factor receptor (EGFR) across
the different cell lines are in agreement with their reported status<sup>21-23</sup>.

To gain further insights into each cancer cell line, we analysed the expression of 48 literature-based biomarkers of clinical relevance<sup>24</sup>, as reported in Figure 1D. Luminal cell lines highly express luminal epithelium genes, but neither basal epithelial nor stromal markers; on the contrary, triple-negative BC cell lines (11 out of 15) show a basal-like phenotype with the expression of at least one of keratin 5, 14 or 17<sup>25,26</sup>, with triple-negative subtype B (TNB) cell lines also expressing vimentin (VIM) and Collagen Type VI Alpha Chains (COL6A1, COL6A2, COL6A3)<sup>21</sup>. Interestingly, two out of five HER2 overexpressing (HER2<sup>+</sup>) cell lines (JIMT1 and HCC1954) in the atlas are in the triplenegative "archipelago" and express keratin 5 (KRT5) (Figure 1A,D), which has been linked
to poor prognosis and trastuzumab resistance<sup>27</sup>. Indeed, both cell lines are resistant to antiHER2 treatments<sup>28</sup>. Finally, the non-tumorigenic MCF12A cell line lacks expression of
ESR1, PGR and HER2 and displays a basal-like phenotype characterized by the expression
of all basal-like marker genes including keratin 5, 14, 17 and TP63, in agreement with the
literature<sup>29</sup>.

99 Overall, these results show that single cell transcriptomics can be successfully used100 to capture the overall expression of clinically relevant markers.

101 102

#### 103 2. The BC single-cell atlas identifies clinically relevant transcriptional signatures.

By clustering the 35,276 single-cells in the atlas, we identified 22 clusters, as shown in Figure 1E. Within the luminal island, cells did not cluster according to their cell line of origin, indeed four out of the five luminal clusters contain cells from distinct cell lines (Figure 1F and Supplementary Figure 02). On the contrary, triple-negative cell lines clustered according to their cell line of origin, with each cluster containing mostly cells from the same cell line (Figure 1F).

110 We identified genes specifically expressed among cells in the same cluster for a 111 total of 22 biomarkers, one for each cluster (Figure 1G,H and Supplementary Figure 03). 112 Interestingly, neither ESR1 nor ERRB2 were part of this set. Literature mining confirmed 113 the significance of some of these markers: clusters in the luminal island (Figure 1G) were associated to genes involved in cancer progression (BCAS3<sup>30,31</sup> cluster 2), dissemination 114 (SCGB2A2<sup>32,33</sup> cluster 6), proliferation (DRAIC<sup>34,35</sup> cluster 1), migration and invasion 115 (CLCA2<sup>36,37</sup> cluster 8 and PIP<sup>38</sup> cluster 18). Interestingly, whereas DRAIC is correlated 116 with poorer survival of luminal BC patients<sup>35</sup>, both CLCA2 and PIP are significantly 117 associated with a favourable prognosis<sup>36,37,39,40</sup>. 118

119 To examine the clinical relevance of these 22 biomarkers, we analysed their 120 expression across 937 breast cancer patients from the TGCA collection encompassing all 121 four BC types. Out of the 22 biomarkers, two (MAGEA4 and XAGE2) could not be 122 mapped to the TGCA dataset. As shown in Figure 1H, there is a marked difference in the 123 expression of the 20 cluster-derived biomarkers across Luminal A, Luminal B, Her2 124 positive and Triple Negative patients. Moreover, it is possible to distinguish subtypes 125 within each category, which may lead to novel diagnostic/prognostic biomarkers (Figure 126 1H and Supplementary Figure 04). For example, one subset of triple-negative patients 127 strongly expresses the protease kallikrein-10 (KLK10), which has been associated with 128 poor prognosis, poor response to tamoxifen treatment<sup>41</sup> and identified as potential target to 129 reverse trastuzumab resistance<sup>42</sup>. Whereas a second subset is characterised by actin gamma 2 expression (ACTG2), which has been linked in BC to cell proliferation<sup>43</sup> and platinum-130 based chemotherapy sensitivity<sup>44–47</sup>. 131

Finally, we compared the performance of the 20 biomarker genes in classifying BC subtypes from bulk RNA-seq data (Methods) against the PAM50 gene signature (50 genes)<sup>6</sup> used in clinics to identify breast cancer subtypes (Figure 1I). The performances were overall comparable, with the obvious exceptions of HER2-overexpressing cancers. Indeed, when adding *ERBB2* to the list of 20 cluster-based biomarkers, classification of this subtypes markedly improved (Figure 1I). Altogether, these analyses confirm that the single cell BC cell line atlas allows
identifying clinically relevant gene signatures useful for patient stratification and tumour
type classification.

141

#### 142 **3.** The BC atlas as a reference for automated cancer diagnosis

143 The BC atlas can be used as a reference against which to compare single cell 144 transcriptomics data from a patient's tissue biopsy and to perform cancer subtype 145 classification and assessment of tumour heterogeneity. To this end, we developed an 146 algorithm able to map single-cell transcriptional profiles from a patient onto the BC atlas 147 and to assign a specific cell line to each of the patient's cells (Methods). We first tested the 148 ability of the algorithm in correctly classifying the very cells in the atlas starting from their 149 single-cell transcriptional profiles and correctly classified 92% of the cells (Supplementary 150 Figure 05). We then turned to single-cell transcriptional profiles obtained from five triplenegative breast cancer patients<sup>48</sup>. As shown in Figure 2A, most, but not all the patients' 151 152 cells mapped to the triple-negative "archipelago", except for the TNBC5 sample, for which 153 most cells mapped to the luminal island. As the algorithm assigns a specific cell line to 154 each tumour cell, it is also possible to look at the cell line composition of each patient, as 155 reported in Figure 2B. These results demonstrates that heterogeneity varies across patients 156 but is present in all the samples, as no patient's biopsy mapped to a single cell line. 157 Moreover, information on the drug sensitivity of the individual cell lines composing the 158 tumour may prove useful in guiding therapeutic choices.

159 We next tested the algorithm on spatial transcriptomics dataset obtained from the 160 tissue biopsy of two patients, one diagnosed with ESR1<sup>+</sup>/ERBB2<sup>+</sup> lobular oestrogen 161 positive carcinoma (Figure 2C-E and Supplementary Figure 06A) and the other with ESR1<sup>+</sup>/ERBB2<sup>+</sup> ductal carcinoma (Supplementary Figure 06C,D)<sup>49</sup>. The dataset consists 162 of 3,808 transcriptional profiles for patient 1 (Figure 2C) and 3,615 profiles for patient 2 163 (Supplementary Figure 06C), each obtained from a different tissue "tile" of size 100um x 164 165 100um x 100 um. The algorithm projected each of the spatial tiles onto the BC atlas and 166 assigned a cell line to each tile. We coloured the tiles according to the cell line and the BC 167 subtype of the cell line (Figure 2C) to yield an automatic cancer subtype classification of 168 tiles. Most of the tiles for both patients were assigned to just two cell lines and correctly 169 classified as luminal (A or B); the remaining 13% of the tiles for patient 1 and 20% for 170 patient 2 were instead classified either as HER2-overexpressing or Triple Negative, which 171 could be an important information to guide therapeutic choice and to predict the occurrence 172 of drug resistance.

As bulk gene expression profiles are more clinically relevant than single-cell gene 173 174 expression profiles, we next trained a deconvolution algorithm **Bisque**<sup>50</sup> (Methods and 175 Supplementary Figure 07) by leveraging our single-cell atlas to predict the cell line 176 composition of a tumour sample. To test the effectiveness of this algorithm, we collected 177 937 bulk gene expression profiles from breast cancer patients in TGCA whose BC subtypes 178 were annotated, and then assigned to each patient the corresponding cell line composition, 179 as shown in Figure 2D,E. Reassuringly, patients diagnosed with a specific breast cancer 180 subtype tend to have a tumour cell line composition consisting of cell lines of the same subtype. We quantified this observation in Figure 2F and observed some interesting 181 182 exceptions: JIMT-1 is an HER2-overexpressing cell line with an amplified ERBB2 locus, 183 but no HER2+ patient was mapped to this cell line. Interestingly, JIMT-1 cells are resistant

to anti-HER2 treatments<sup>51</sup>; another example is the HS578T cell line, which is reported to
be triple-negative, however the majority of patients who map to it are luminal; surprisingly,
this cell line has been reported to be sensitive to fulvestrant<sup>1,2</sup>, an anti-ESR1 drug.

187 These results show that this single cell atlas of cancer cell can be used to 188 automatically assign cell line composition and cancer subtypes both from single-cell 189 expression profiles and bulk gene expression profile.

4. Clinically relevant biomarkers exhibit heterogenous and dynamic expression in BCcell lines.

- 193 Clinically relevant receptors are heterogeneously expressed across cells belonging to the 194 same cell line, as assessed by computing the percentage of cells in a cell line expressing 195 the receptor as in Figure 3A. Consider the seven Luminal B and HER2<sup>+</sup> cell lines present 196 in the BC atlas, which by definition overexpress HER2: whereas more than 90% of cells 197 in AU565, BT574 and HCC1954 cell lines express *ERBB2*, in the remaining four cell lines 198 *ERBB2* expression ranged from 31% of EVSAT cells to 46% of JIMT1 cells and up to 64% 199 of MDA-MB-361 cells. This happens despite both JIMT1 and MDA-MB-361 harbour a 200 copy number gain of the locus containing the *ERBB2<sup>52</sup>*. We first excluded the possibility 201 that these results were artifacts of single-cell RNA-sequencing technology (Supplementary 202 Figure 08). We then assessed HER2 protein levels by flow cytometry in three 203 representative cell lines: AU565 (high HER2 expression), MDA-MB-361 (heterogeneous 204 HER2 expression) and HCC38 cell lines (low HER2 expression). As shown in Figure 3B, 205 single-cell transcriptional data agree with the cytometric analysis; however, the origin of 206 this heterogeneity is unclear. To exclude hereditable genetic differences as a source of 207 heterogeneity, we sorted MDA-MB-361 cells into HER2<sup>+</sup> and HER2<sup>-</sup> subpopulations 208 (Methods) and checked whether these homogenous subpopulations were stable over time, 209 or rather spontaneously gave rise to heterogeneous populations. As shown in Figure 3C, 210 after 18 days in culture, both subpopulations re-established the original heterogeneity, 211 demonstrating that HER2 expression in these cells is dynamic and driven by a yet 212 undiscovered mechanism.
- 213 Interestingly, HER2<sup>+</sup> circulating tumour cells (CTCs) isolated from an ER<sup>+</sup>/HER2<sup>-</sup> 214 breast cancer patient were shown to spontaneously interconvert from HER2<sup>-</sup> and HER2<sup>+</sup>, 215 with cells harbouring a phenotype producing daughters of the opposite one<sup>53</sup>. To check if cell-cycle phase could explain the observed heterogeneity in the MDA-MB-361 cell line, 216 217 we computationally predicted (Methods) the cell cycle phase of each cell in both the HER2<sup>-</sup> and HER2<sup>+</sup> subpopulations from single cell transcriptomics data<sup>54</sup>. A higher proportion of 218 219 HER2<sup>-</sup> cells was predicted to be in S/G2/M phases when compared to HER2+ cells (Figure 220 3D). This result is consistent with previous observations that report cell cycle arrest in 221 G2/M phase following HER2 inhibition<sup>55</sup>.
- 222 We next set to identify biological processes differing between the two 223 subpopulations by computing differentially expressed genes (DEGs) from the single-cell transcriptional profiles of HER2<sup>+</sup> cells against HER2<sup>-</sup> cells (Supplementary Table 02). 224 Gene Set Enrichment Analyses (GSEA) <sup>56</sup> against the ranked list of DEGs, reported in 225 226 Figure 3E, revealed seven significantly enriched pathways (FDR<10%): four of which 227 were upregulated in HER2<sup>+</sup> cells, but downregulated in HER2<sup>-</sup> cells, and included 228 adipogenesis, myogenesis and OXPHOS, all indicative of EMT engagement, which has been reported in HER2<sup>+</sup> cells<sup>57-59</sup>; the remaining three pathways were upregulated in 229

HER2<sup>-</sup> cells and related to cell-cycle and specifically to G2/M phase, in agreement with our previous analysis, suggesting that cell cycle may play a role in HER2 expression in this cell line.

These results show that heterogeneity in the expression of clinically relevant biomarkers is present even in cell lines and that it can also be dynamic and of a non-genetic nature.

236

#### 237 5. Heterogeneity in gene expression affects drug response.

238 To investigate the role of heterogeneity in gene expression within a cell line on drug 239 response, we collected large-scale in vitro drug screening data<sup>1,2</sup> reporting the effect of 450 drugs on 658 cancer cell lines from solid tumours. As show in Figure 3F and 240 241 Supplementary Figure 09, sensitivity of the BC cell lines to HER2 inhibitors was 242 significantly correlated with the percentage of cells in the cell line expressing *ERBB2* 243 (Supplementary Table 03). Receptor expression level is substantially the same across cells 244 expressing it, irrespective of the cell line they belong to (Supplementary Figure 10), except 245 for cell lines harbouring CNVs of the ERBB2 locus. Furthermore, we found that the 246 correlation between drug target expression and drug sensitivity holds true also for several 247 other targets (Figure 3G), thus suggesting that variability in gene expression within cells 248 of the same tumour may cause some cells to respond poorly to the drug treatment.

249 Starting from these observations, we developed DREEP (DRug Estimation from 250 single-cell Expression Profiles), a novel bioinformatics tool that, starting from single-cell 251 transcriptional profiles, allows to predict drug response at the single cell level. To this end, 252 we first detected expression-based biomarkers of drug sensitivity for 450  $drugs^2$ , as 253 schematised in Figure 4H,I (Methods). Briefly, we crossed data from the Cancer Cell Line 254 Encyclopaedia (CCLE) on the response to 450 drugs across 658 cancer cell lines from solid 255 tumours with their gene expression profiles from bulk RNA-seq. In the CCLE, drug 256 potency is evaluated as the inverse of the Area Under the Curve (AUC) of the dose-257 response graph, with low values of the AUC indicating drug sensitivity, while high values 258 implying drug resistance (Figure 3H). For each gene and for each drug, we computed the 259 correlation between the expression of the gene across the 658 cell lines with the drug 260 potency in the same cell lines. Hence, genes positively correlated with the AUC are 261 potential markers of resistance, vice-versa, negatively correlated genes are markers of 262 sensitivity (Figure 3H). In this way, we generated a ranked list of expression-based 263 biomarkers of drug sensitivity and resistance for each of the 450 drugs. We then used these 264 biomarkers to predict drug sensitivity at the single-cell level (Figure 3I). To this end, the 250 genes most expressed of each cell in the atlas were compared against the ranked list of 265 266 biomarkers for each one of 450 drugs by means of GSEA<sup>56</sup> and thus associated to the drug 267 it is most sensitive to, or to no drug, if no significant enrichment score from GSEA is found 268 (Figure 3I).

To assess the algorithm's performance, we applied it to the single-cell BC atlas and estimated its performance by checking how well we could predict sensitivity of the 32 BC cell lines to 86 drugs for which this information was publicly available<sup>60</sup> (Figure 3J). To convert single-cell predictions to predictions at the cell line level, we simply used the percentage of cells in the cell line deemed to be sensitive to the drug by the algorithm. To experimentally validate DREEP, we turned to the MDA-MB-361 cell line for which we found coexistence of two distinct and dynamic cell subpopulations (HER2<sup>+</sup> and HER2<sup>-</sup>). 276 We applied DREEP to each subpopulation to identify drugs able to selectively inhibit 277 growth of either the HER2<sup>-</sup> subpopulation or the HER2<sup>+</sup> subpopulation: 42 drugs (FDR <278 1%, Supplementary Table 04) were predicted to preferentially inhibit growth of HER2<sup>-</sup> 279 cells; the most overrepresented class among these drugs was that of inhibitors of DNA 280 topoisomerases (TOP1/TOP2A) (Supplementary Figure 11) such as Etoposide. 281 Surprisingly, no drug was found to specifically inhibit growth of HER<sup>+</sup> cells, whereas 44 282 drugs (FDR <1%) were predicted to be equally effective on both subpopulations and 283 unexpectedly included HER2 inhibitors, such as afatinib (Supplementary Table 03 and 284 Supplementary Figure 12).

285 We selected etoposide and afatinib for further experimental validation. MDA-MB-286 361 cells were first sorted by FACS into HER2<sup>+</sup> and HER2<sup>-</sup> subpopulations and then cell 287 viability was measured following 72h drug treatment at five different concentrations as 288 shown in Figure 3K (and Supplementary Table 05). In agreement with DREEP predictions, 289 HER2<sup>-</sup> cells were much more sensitive to etoposide than HER2<sup>+</sup> cells, while afatinib was 290 equally effective on both subpopulations. This counterintuitive result was similar to that 291 observed by Jordan et al<sup>53</sup> using circulating tumour cells from a BC patient sorted into 292 HER2<sup>-</sup>and HER2<sup>+</sup> subpopulations, which were found to be equally sensitive to Lapatinib 293 (another HER2 inhibitor), but no mechanism of action was put forward.

294 We hypothesise that the dynamic interconversion of MDA-MB-361 cells between 295 the HER2<sup>-</sup> and the HER2<sup>+</sup> state may explain this surprising result: when the starting 296 population consists of HER2<sup>-</sup> cells only, some of these cells will nevertheless interconvert 297 to HER2<sup>+</sup> cells during afatinib treatment, and they will thus become sensitive to HER2 298 inhibition, explaining the observed results. We mathematically formalised this hypothesis 299 with a simple mathematical model (Supplementary Figure 13 and in the Supplementary 300 Material) where two species (HER2<sup>+</sup> and HER2<sup>-</sup> cells) can replicate and interconvert, but 301 only one (HER2<sup>+</sup>) is affected by afatinib treatment. The model shows that if the 302 interconversion time between the two cell states is comparable to that of the cell cycle, then 303 afatinib treatment will have the same effect on both subpopulations. If instead the 304 interconversion time is much longer than the cell cycle, then afatinib will have little effect 305 on HER2<sup>-</sup> sorted cells, but maximal effects on HER2<sup>+</sup> sorted cells, and vice-versa, if the 306 interconversion time is much shorter than the cell cycle, then afatinib's effect would be 307 minimal on both HER2<sup>-</sup> and HER2<sup>+</sup> sorted cells.

Comparison of the modelling results with the experimental results thus suggests that the interconversion rate should be of the same order of the cell cycle (about 72h for MDAM361 cells). The model further predicts that treating the unsorted population of MDA-MB-361 cells with afatinib reduces the percentage of HER2<sup>+</sup> cells, since only HER2<sup>+</sup> will be affected, but that this percentage quickly recovers once Afatinib treatment is interrupted (Supplementary Figure 14 and 15 and Supplementary Material).

314 To test modelling predictions, we treated the MDAM361 cell line (without sorting) 315 with a fatinib and etoposide and then assessed by cytofluorimetry the percentage of HER2+ 316 and HER2<sup>-</sup> cells before and after the treatment. As shown in Figure 3L,M (Supplementary 317 Table 06 and Supplementary Table 07) etoposide increased the percentage of HER2<sup>+</sup> cells, 318 in agreement with the increased sensitivity of HER2<sup>-</sup> cells to this treatment, whereas 319 afatinib strongly decreased the percentage of HER2<sup>+</sup> cells, confirming that its effect is 320 specific for HER2<sup>+</sup> cells only. We next measured the percentage of HER2<sup>+</sup> cells following 321 removal of afatinib from the medium; as shown in Figure 3N,O the percentage of HER2+

322 cells quickly increased confirming the modelling results (Supplementary Figure 15 and323 Supplementary Material).

All together our results show that DREEP can predict drug sensitivity from singlecell transcriptional profiles and that dynamic heterogeneity in gene expression does play a significant role in how the cell population will respond to the drug treatment.

#### 327 229 **Di**so

#### 328 Discussion

329 In this study we provide the first transcriptional characterization at single cell level of a 330 panel of 32 breast cell lines. We show that single cell transcriptomics can be used to capture 331 the expression of clinically relevant markers. We show that breast-cancer cell lines express 332 clinically relevant BC receptors heterogeneously among cells within the same cell line. 333 Moreover, we observed dynamic plasticity in the regulation of HER2 expression in the 334 MDA-MB-361 cell line with striking consequences on drug response. This phenomenon 335 has been recently observed also in circulating tumour cells of a BC patient<sup>53</sup> and in other cell lines<sup>17,61</sup>. 336

We determined cell line composition of patients' biopsies both from both singlecell and bulk gene expression profiles. Estimation of cancer cell line composition provides an alternative and more information-rich framework to link bulk gene expression measurement of patient's biopsies to preclinical cancer models. Knowledge of drugs to which cancer cell lines are sensitive to may also inform drug treatment for patients for which bulk gene expression profiles have been measured.

343 Single cell transcriptomics is still not clinically ready because of the costs and time 344 needed, however this work shows the importance of performing single-cell sequencing on 345 the available cancer models, including cell lines and organoids to build a set of cell cancer 346 states with known phenotypes and drug response to which patients' tumour can be mapped 347 to make a leap in personalised diagnosis, prognosis and treatment of cancer patients.

#### References

- Iorio, F. *et al.* A Landscape of Pharmacogenomic Interactions in Cancer. *Cell* 166, 740–754 (2016).
- Rees, M. G. *et al.* Correlating chemical sensitivity and basal gene expression reveals mechanism of action. *Nat Chem Biol* 12, 109–116 (2016).
- Cardoso, F. *et al.* 70-Gene Signature as an Aid to Treatment Decisions in Early-Stage Breast Cancer. *N. Engl. J. Med.* 375, 717–729 (2016).
- 4. Sparano, J. A. *et al.* Prospective Validation of a 21-Gene Expression Assay in Breast Cancer. *N. Engl. J. Med.* **373**, 2005–2014 (2015).
- Sparano, J. A. *et al.* Adjuvant Chemotherapy Guided by a 21-Gene Expression Assay in Breast Cancer. *N. Engl. J. Med.* 379, 111–121 (2018).
- Cheang, M. C. U. *et al.* Defining Breast Cancer Intrinsic Subtypes by Quantitative Receptor Expression. *Oncologist* 20, 474–482 (2015).
- 7. Harbeck, N. et al. Breast cancer. Nature Reviews Disease Primers 5, (2019).
- Andre, F. *et al.* Use of Biomarkers to Guide Decisions on Adjuvant Systemic Therapy for Women With Early-Stage Invasive Breast Cancer: ASCO Clinical Practice Guideline Update—Integration of Results From TAILORx. *J. Clin. Oncol.* 37, 1956–1964 (2019).
- Garnett, M. J. *et al.* Systematic identification of genomic markers of drug sensitivity in cancer cells. *Nature* 483, 570– 575 (2012).
- Foulkes, W. D., Smith, I. E. & Reis-Filho, J. S. Triple-Negative Breast Cancer. N. Engl. J. Med. 363, 1938–1948 (2010).
- Sharma, S. V. *et al.* A Chromatin-Mediated Reversible Drug-Tolerant State in Cancer Cell Subpopulations. *Cell* 141, 69–80 (2010).
- Shaffer, S. M. *et al.* Rare cell variability and drug-induced reprogramming as a mode of cancer drug resistance. *Nature* 546, 431–435 (2017).
- Ebinger, S. *et al.* Characterization of Rare, Dormant, and Therapy-Resistant Cells in Acute Lymphoblastic Leukemia. *Cancer Cell* 30, 849–862 (2016).
- Meyer, A. S. & Heiser, L. M. Systems biology approaches to measure and model phenotypic heterogeneity in cancer. *Curr. Opin. Syst. Biol.* 17, 35–40 (2019).
- Marusyk, A., Janiszewska, M. & Polyak, K. Intratumor Heterogeneity: The Rosetta Stone of Therapy Resistance. *Cancer Cell* 37, 471–484 (2020).
- Shaffer, S. M. *et al.* Memory Sequencing Reveals Heritable Single-Cell Gene Expression Programs Associated with Distinct Cellular Behaviors. *Cell* 182, 947-959.e17 (2020).
- Schuh, L. *et al.* Gene Networks with Transcriptional Bursting Recapitulate Rare Transient Coordinated High Expression States in Cancer. *Cell Syst.* 10, 363-378.e12 (2020).
- Costello, J. C. *et al.* A community effort to assess and improve drug sensitivity prediction algorithms. *Nat Biotechnol* 32, 1202–1212 (2014).

- Zheng, G. X. Y. *et al.* Massively parallel digital transcriptional profiling of single cells. *Nat. Commun.* 8, 1–12 (2017).
- Macosko, E. Z. *et al.* Highly Parallel Genome-wide Expression Profiling of Individual Cells Using Nanoliter Droplets. *Cell* 161, 1202–1214 (2015).
- Dai, X., Cheng, H., Bai, Z. & Li, J. Breast cancer cell line classification and Its relevance with breast tumor subtyping. *J. Cancer* 8, 3131–3141 (2017).
- Soliman, N. A. & Yussif, S. M. Ki-67 as a prognostic marker according to breast cancer molecular subtype. *Cancer Biol. Med.* 13, 496–504 (2016).
- 23. Tajadura-Ortega, V. *et al.* O-linked mucin-type glycosylation regulates the transcriptional programme downstream of EGFR. *Glycobiology* (2020). doi:10.1093/glycob/cwaa075
- Karaayvaz, M. *et al.* Unravelling subclonal heterogeneity and aggressive disease states in TNBC through single-cell RNAseq. *Nat. Commun.* 9, (2018).
- Badve, S. *et al.* Basal-like and triple-negative breast cancers: A critical review with an emphasis on the implications for pathologists and oncologists. *Mod. Pathol.* 24, 157–167 (2011).
- 26. Gusterson, B. Do 'basal-like' breast cancers really exist? *Nat. Rev. Cancer* **9**, 128–134 (2009).
- Martin-Castillo, B. *et al.* Cytokeratin 5/6 fingerprinting in HER2-positive tumors identifies a poor prognosis and trastuzumab-resistant Basal-HER2 subtype of breast cancer. *Oncotarget* 6, 7104–7122 (2015).
- Jernström, S. *et al.* Drug-screening and genomic analyses of HER2-positive breast cancer cell lines reveal predictors for treatment response. *Breast Cancer Targets Ther.* 9, 185–198 (2017).
- Sweeney, M. F., Sonnenschein, C. & Soto, A. M. Characterization of MCF-12A cell phenotype, response to estrogens, and growth in 3D. *Cancer Cell Int.* 18, 1–12 (2018).
- Gururaj, A. E. *et al.* MTA1, a transcriptional activator of breast cancer amplified sequence 3. *Proc. Natl. Acad. Sci. U. S. A.* 103, 6670–6675 (2006).
- Bärlund, M. *et al.* Cloning of BCAS3 (17q23) and BCAS4 (20q13) genes that undergo amplification, overexpression, and fusion in breast cancer. *Genes Chromosom. Cancer* 35, 311–317 (2002).
- Zehentner, B. K. & Carter, D. Mammaglobin: A candidate diagnostic marker for breast cancer. *Clin. Biochem.* 37, 249– 257 (2004).
- Al Joudi, F. S. Human mammaglobin in breast cancer: A brief review of its clinical utility. *Indian J. Med. Res.* 139, 675–685 (2014).
- Sun, M., Gadad, S. S., Kim, D. S. & Kraus, W. L. Discovery, Annotation, and Functional Analysis of Long Noncoding RNAs Controlling Cell-Cycle Gene Expression and Proliferation in Breast Cancer Cells. *Mol. Cell* 59, 698–711 (2015).
- Zhao, D. & Dong, J. T. Upregulation of long non-coding RNA DRAIC correlates with adverse features of breast cancer. *Non-coding RNA* 4, 1–9 (2018).
- 36. Qiang, Y. Y. *et al.* Along with its favorable prognostic role, CLCA2 inhibits growth and metastasis of nasopharyngeal

carcinoma cells via inhibition of FAK/ERK signaling. J. Exp. Clin. Cancer Res. **37**, 1–14 (2018).

- Li, X., Cowell, J. K. & Sossey-Alaoui, K. CLCA2 tumour suppressor gene in 1p31 is epigenetically regulated in breast cancer. *Oncogene* 23, 1474–1480 (2004).
- Urbaniak, A., Jablonska, K., Podhorska-Okolow, M., Ugorski, M. & Dziegiel, P. Prolactin-induced protein (PIP)characterization and role in breast cancer progression. *Am. J. Cancer Res.* 8, 2150–2164 (2018).
- Debily, M. A. *et al.* A functional and regulatory network associated with PIP expression in human breast cancer. *PLoS One* 4, (2009).
- 40. Gruber, A. D. & Pauli, B. U. Tumorigenicity of Human Breast Cancer Is Associated with Loss of the Ca<sup&gt;2+&lt;/sup&gt;&lt;em&gt;-</em&gt;activated Chloride Channel CLCA2. *Cancer Res.* 59, 5488 LP – 5491 (1999).
- Wang, Z. *et al.* Identification of KLK10 as a therapeutic target to reverse trastuzumab resistance in breast cancer. *Oncotarget* 7, 79494–79502 (2016).
- Luo, L.-Y., Diamandis, E. P., Look, M. P., Soosaipillai, A. P. & Foekens, J. A. Higher expression of human kallikrein 10 in breast cancer tissue predicts tamoxifen resistance. *Br. J. Cancer* 86, 1790–1796 (2002).
- Dugina, V., Shagieva, G., Khromova, N. & Kopnin, P. Divergent impact of actin isoforms on cell cycle regulation. *Cell Cycle* 17, 2610–2621 (2018).
- Lu, X. *et al.* Establishment of a Predictive Genetic Model for Estimating Chemotherapy Sensitivity of Colorectal Cancer with Synchronous Liver Metastasis. *Cancer Biother. Radiopharm.* 28, 552–558 (2013).
- Edfeldt, K., Hellman, P., Westin, G. & Stalberg, P. A plausible role for actin gamma smooth muscle 2 (ACTG2) in small intestinal neuroendocrine tumorigenesis. *BMC Endocr. Disord.* 16, 19 (2016).
- Xu, C.-Z. *et al.* Gene and microRNA expression reveals sensitivity to paclitaxel in laryngeal cancer cell line. *Int. J. Clin. Exp. Pathol.* 6, 1351–1361 (2013).
- Verrills, N. M. *et al.* Alterations in γ-Actin and Tubulin-Targeted Drug Resistance in Childhood Leukemia. *JNCI J. Natl. Cancer Inst.* 98, 1363–1374 (2006).
- Gao, R. *et al.* Delineating copy number and clonal substructure in human tumors from single-cell transcriptomes. *Nat. Biotechnol.* (2021). doi:10.1038/s41587-020-00795-2
- Genomics, 10x. 10X Genomics datasets. Available at: https://wp.10xgenomics.com/resources/datasets.
- Jew, B. *et al.* Accurate estimation of cell composition in bulk expression through robust integration of single-cell information. *Nat. Commun.* 11, 1971 (2020).
- Tanner, M. *et al.* Characterization of a novel cell line established from a patient with Herceptin-resistant breast cancer. *Mol. Cancer Ther.* 3, 1585 LP – 1592 (2004).
- Ghandi, M. *et al.* Next-generation characterization of the Cancer Cell Line Encyclopedia. *Nature* 569, 503–508 (2019).
- 53. Jordan, N. V. et al. HER2 expression identifies dynamic

functional states within circulating breast cancer cells. *Nature* **537**, 102–106 (2016).

- Butler, A., Hoffman, P., Smibert, P., Papalexi, E. & Satija, R. Integrating single-cell transcriptomic data across different conditions, technologies, and species. *Nat. Biotechnol.* 36, 411 (2018).
- Yan, Y. *et al.* A novel function of HER2/Neu in the activation of G2/M checkpoint in response to γ-irradiation. *Oncogene* 34, 2215–2226 (2015).
- Subramanian, A. *et al.* Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proc. Natl. Acad. Sci. U. S. A.* **102**, 15545– 50 (2005).
- Ishay-Ronen, D. *et al.* Gain Fat—Lose Metastasis: Converting Invasive Breast Cancer Cells into Adipocytes Inhibits Cancer Metastasis. *Cancer Cell* 35, 17-32.e6 (2019).
- Ingthorsson, S. *et al.* HER2 induced EMT and tumorigenicity in breast epithelial progenitor cells is inhibited by coexpression of EGFR. *Oncogene* 35, 4244–4255 (2016).
- Savci-Heijink, C. D. *et al.* Epithelial-to-mesenchymal transition status of primary breast carcinomas and its correlation with metastatic behavior. *Breast Cancer Res. Treat.* 174, 649–659 (2019).
- Yang, W. *et al.* Genomics of Drug Sensitivity in Cancer (GDSC): a resource for therapeutic biomarker discovery in cancer cells. *Nucleic Acids Res.* 41, D955–D961 (2013).
- Gupta, P. B. *et al.* Stochastic State Transitions Give Rise to Phenotypic Equilibrium in Populations of Cancer Cells. *Cell* 146, 633–644 (2011).
- 62. Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15–21 (2012).
- Harrow, J. *et al.* GENCODE: The reference human genome annotation for The ENCODE Project. *Genome Res.* 22, 1760– 1774 (2012).
- Smith, T., Heger, A. & Sudbery, I. UMI-tools: modeling sequencing errors in Unique Molecular Identifiers to improve quantification accuracy. *Genome Res.* 27, 491–499 (2017).
- Gambardella, G. & di Bernardo, D. A Tool for Visualization and Analysis of Single-Cell RNA-Seq Data Based on Text Mining. *Front. Genet.* 10, (2019).
- Slovin, S. *et al.* Single-Cell RNA Sequencing Analysis: A Stepby-Step Overview. in (ed. Picardi, E.) 343–365 (Springer US, 2021). doi:10.1007/978-1-0716-1307-8\_19
- Robinson, M. D., McCarthy, D. J. & Smyth, G. K. edgeR: a Bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics* 26, (2010).
- McInnes, L., Healy, J. & Melville, J. UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction. *Prepr. https://arxiv.org/abs/1802.03426* (2018).
- Levine, J. H. *et al.* Data-Driven Phenotypic Dissection of AML Reveals Progenitor-like Cells that Correlate with Prognosis. *Cell* 162, 184–197 (2015).
- Grossman, R. L. *et al.* Toward a Shared Vision for Cancer Genomic Data. *N. Engl. J. Med.* 375, 1109–1112 (2016).

- Colaprico, A. *et al.* TCGAbiolinks: an R/Bioconductor package for integrative analysis of TCGA data. *Nucleic Acids Res.* 44, e71–e71 (2015).
- 72. Barretina, J. *et al.* The Cancer Cell Line Encyclopedia enables predictive modelling of anticancer drug sensitivity. *Nature* **483**, 603–607 (2012).

## 1 Figures

2

3 Figure 1 – The Breast Cancer Single Cell Atlas. (A) Representation of single-cell expression profiles of 35,276 cells from 32 cell lines color-coded according to cancer 4 5 subtype (LA=Luminal A, LB=Luminal B, H=Her2 positive, TNA = Triple Negative A, 6 TNB = Triple Negative B). (B) Expression levels of the indicated genes in the atlas, 7 with red indicating expression, together with their (C) distribution within the cell 8 lines, shown as a violin plot. (**D**) Dotplot of literature-based biomarker genes along 9 the columns for each of the 32 sequenced cell lines along the rows. Biomarker genes 10 are grouped by type (Basal Epith. = Basal Epithelial, Luminal Epith. = Luminal Epithelial, L.P. = Luminal Progenitor, EMT = Epithelial to Mesenchymal Transition). 11 12 (E) Graphical representation of 35,276 cells color-coded according to their cluster of 13 origin. Clusters are numbered from 1 to 22. (F) For the indicated cluster, the 14 corresponding pie-chart represents the cluster composition in terms of cell lines. Cell 15 lines in the same pie-chart are distinguished by colour. Only the top 10 most heterogenous clusters are shown. Cluster 2 is the most heterogeneous while cluster 16 17 19 is the most homogeneous. (G) Expression levels in the atlas of the five luminal 18 biomarkers identified as the most differentially expressed in each of the five luminal 19 clusters (1, 2, 6, 8 and 18). (H) Expression of 20 out of 22 atlas-derived biomarkers in 20 the biopsies of 937 breast cancer patient from TCGA. (I) Accuracy in classifying 21 tumour subtype for 937 patients from TCGA by using either PAM50 or the 20 atlas 22 derived biomarker genes (scCCL) alone or augmented with HER2 gene (scCCL + 23 HER2).

24

25 Figure 2 -Automatic classification of patients' tumour cells (A) Cancer cells from 26 triple negative breast cancer (TNBC) biopsies of 5 patients are embedded in the BC 27 atlas to predict their tumour type. (B) For each patient, the pie chart shows cell line 28 composition obtained by mapping patient's cells onto the atlas. (C) Tissue-slide of an 29 oestrogen positive breast tumour biopsy sequenced using 10x Visium spatial 30 transcriptomics (top-left) and the position of the mapped tissue tiles onto the atlas (top-left). Tiles are colour-coded according to the cell line (bottom-left) and to tumor 31 32 subtype (bottom-right) as predicted by the mapping algorithm. (D) Cell line 33 composition for each patient as estimated by the algorithm from bulk RNA-seq of 937 BC patients. For ease of interpretation, in the heatmap patients are clustered 34 35 according to their cell line composition. The bottom row reports the annotated cancer 36 subtype in TGCA. (E) Predicted cell-line composition for four representative patients. 37 (F) The distribution of the 937 BC patients across the 32 cell lines. For each cell line, 38 the stacked bars report the percentage of patients of a given cancer subtype assigned 39 by the algorithm to that cell line.

40

Figure 3 - Transcriptional heterogeneity in breast cancer cell lines and its
 impact on drug response. (A) Percentage of cells expressing the indicated genes in
 each of the sequenced 32 cell lines. (B) Fluorescence cytometry of HCC38, MDA-MB-

44 361 and AU565 cell lines stained with a fluorescent antibody against Her2. (C) 45 Expression of HER2 protein in MDA-MB-361 cells is dynamic and re-established in less than 3 weeks. (D) Analysis of the cell cycle phase for the HER2+ and HER2-46 47 subpopulations of MDA-MB-361 cells. The cell cycle of each cell is estimated from its 48 single-cell transcriptomics profile. (E) Enriched pathways (GSEA, FDR<10%) across 49 the genes differentially expressed between the HER2+ and HER2- subpopulations of 50 MDA-MB-361 cells. Orange refers to HER2+ subpopulation and blue to the HER2-51 ones. (F) Relationship between gene expression and drug potency for four anti-HER2 52 drugs. Each dot corresponds to a cell line reporting the percentage of cells expressing 53 ERBB2 or EGFR in the cell line [y-axis] and the drug potency [x-axis]. PCC (pearson 54 correlation coefficient) and p-value are also shown. (G) Box-plot reporting the 55 distribution of PCCs between percentage of cells expressing the cognate drug target 56 and the potency of the drug across cell lines for 66 drugs for two different drug 57 potency databases. For comparison, the PCC distribution when choosing a random 58 gene in place of the cognate drug target is also shown. (H) Bioinformatics pipeline for 59 the identification of drug sensitivity biomarkers for 450 drugs. For each drug, the 60 expression of a gene across 658 cell lines is correlated with drug potency in the same 61 cell lines; genes are then ranked from most positively correlated to the most 62 negatively correlated. (I) The top 250 most expressed genes in a single cell are used 63 as input for a Gene Set Enrichment Analysis (GSEA) against the ranked list of genes 64 for each one of the 450 drugs to predict its drug sensitivity. At the end of the process, each cell in the sample is associated to the drug it is most sensitive to, or to no drug, 65 66 if no significant enrichment score is found. Finally, for each of the 450 drugs, the 67 number of cells predicted to be either sensitive, resistant, or not classified in the 68 considered sample is estimated. (J) Validation of DREEP on the Breast Cancer Single 69 Cell atlas data to predict drug sensitivity to 86 drugs. The PPV (Positive Predicted 70 Value) is shown as a function of the percentage of cells in a cell line predicted to be 71 sensitive to the same drug. Dashed line represents the performance of a random 72 algorithm. (K) Dose-response curve for afatinib and etoposide on sorted MDA-MB-73 361 cell populations (triplicate experiment). (L) Percentage of HER2+ cells in MDA-74 MB-361 after 72h treatment with either afatinib (statistic: two-sided t-test. \*P  $\leq$  0.05: 75 \*\*P  $\leq$  0.01; \*\*\*P  $\leq$  0.001) or etoposide and (M) measured cell viability after the treatment. (N) Percentage of HER2 positive cells in MDA-MB-361 cell-line at the 76 77 indicated time-points either after 48h of afatinib pre-treatment (red bars) or without 78 any afatinib pre-treatment (black bars) and (**0**) the relative number of cells rescaled 79 for the number of cells at the beginning of the experiment.

- 80
- 81

#### 82 Methods

83

Cell culture: The 32 cell lines used in this study were obtained from commercial providers and cultured in
 ATCC recommended complete media at 37°C and 5% CO2.

86

87 **DROP-seq platform set-up:** Single cell transcriptomic of the 32 cell lines was performed by implementing 88 in-house the DROP-seq technology<sup>20</sup>. The microfluidics device for the generation of droplet was fabricated 89 using a bio-compatible, silicon-based polymer, polydimethylsiloxane (PDMS) that was rendered 90 hydrophobic with Aquapel® treatment as per protocol<sup>20</sup>. In each sequencing experiment, cell suspension, 91 bead suspension and carrier oil (QX200 droplet generation oil, Bio-Rad) were first loaded in syringes and 92 then placed in syringe pumps (Leafluid). Flow rates of syringe pumps were set at 4,000  $\mu$ L/hr for both cell 93 and barcoded bead suspensions while carrier oil syringe pump was set at 15,000  $\mu$ L/hr. In each sequencing 94 experiment, cells and barcoded beads were respectively diluted at the concentration of 200 cell/µL in PBS 95 with BSA 0.01% (Merck) and 120 bead/uL in lysis buffer. A self-built magnetic stirrer system was used to 96 keep in suspension barcoded beads. To count the occurrence of a single cell together with a barcoded bead 97 several tests were performed without lyses buffer in the bead suspension. In these tests, we observed about 98 5% of generated droplets filled with just one bead and one cell.

99 Single cell RNA library preparation and sequencing: For each sequencing experiment, the targeted 100 number of cells to sequence was set to 2,000. Droplets were collected in a 50 mL falcon and broke by adding 101 1 mL of Perfluoro-1-octanol. Captured RNA was reverse transcribed in a single reaction following the 102 original protocol <sup>20</sup> and then digested with exonuclease 1 to degrade unbound primers. Next, cDNA was first 103 amplified with a total of 12 PCR cycles and then purified using AMPure XP beads at 0.6X ratio. Finally, the 104 quality of the resulting cDNA library was quantified with the BioAnalyzer High Sensitivity DNA Chip and 105 its concentration measured using the Qubit Fluorometer. The Illumina Nextera XT v2 kit was used to produce 106 the next generation sequencing (NGS) libraries using four aliquots of 600pg of each cDNA library. Quality 107 and concentration of NGS libraries were respectively quantified on the BioAnalyzer High Sensitivity DNA 108 Chip and Qubit Fluorometer. Finally, either Illumina NextSeq 500/550 or NovaSeq 6000 machines were used 109 to sequence the produced NGS libraries (Supplementary Table 01). Samples processed with NextSeq500/550 110 NGS library were diluted at the final concentration of 3 nM and sequenced using the 75-cycle high output 111 flow cell while samples processed with NovaSeq 6000 machine were diluted at the final concentration of 250 112 pM and sequenced using the S1 100 cycles flow cell.

113 Read alignment and gene expression quantification: Raw data processing was performed using the Drop-114 seq tools package version 1.13 and following the Drop-seq Core Computational Protocol 115 (http://mccarrolllab.org/dropseq). Briefly, raw sequence data was filtered to remove all read pairs with at 116 least one base in their barcode or UMI with a quality score less than 10. Then read 2 was trimmed at the 5' 117 end to remove any TSO adapter sequence, and at the 3' end to remove polyA tails. Reads were then aligned 118 using STAR <sup>62</sup> on hg38 human genome (primary assembly, version 28) downloaded from GENCODE <sup>63</sup>. 119 After reads alignment, UMI tool <sup>64</sup> was used to perform UMI deduplication and quantify the number of gene 120 transcripts in each cell. The initial number of sequenced cells was identified using a simple (knee-like) 121 filtering rule as implemented by CellRanger 2.2.x. After this, only high depth cells with at least 2,500 UMI, 122 more than 1,000 captured genes and with less than 50% of reads aligned on mitochondrial gene were retained. 123 Putative multiples among the sequenced cells of each BC cell line were simply discarded identifying outliers 124 in the count depth distribution by using Tukey's method based on lower and upper quartiles with k equal to 125 3.

BC Atlas Construction: Single cells expression profiles were normalized using GF-ICF (Gene Frequency
 Inverse Cell Frequency) normalization using the *gficf* package<sup>65,66</sup> for R statistical environment
 (<u>https://github.com/dibbelab/gficf</u>). GF-ICF is based on a data transformation model called term frequency inverse document frequency (TF-IDF) that has been extensively used in the field of text mining. GF-ICF
 transformation was applied on CPM (count per million) after *EdgeR* normalization <sup>67</sup> and discarding genes

expressed in less than 5% of the total number of sequenced cells. Finally, each cell was summarized with its

first 10 Principal Components (PCs) and projected with UMAP <sup>68</sup> into a two dimensional embedded space.
The number of principal components was chosen as the "elbow" point on the plot of the first 50 PCs. UMAP

134 projection was performed by using the *uwot* package in the R statistical environment 3.6.

135 Cell clustering and identification of marker genes: Transcriptionally similar subpopulations of cells were 136 found using a Phenograph like approach<sup>69</sup> as implemented in the *clustcells* function of *gficf* package<sup>65</sup>. 137 Briefly, we initially built a graph of cells by using the K-Nearest Neighbours (KNN) algorithm applied to the 138 PC-reduced space where each cell was connected to its 50 most similar cells using the manhattan distance. 139 Then, to build the final graph of cells, the edge weight between any two cells was computed as the Jaccard 140 similarity, i.e. the proportion of neighbours they share. The Louvain algorithm with resolution parameter 141 equal to 0.25 was used to find communities of cells in this graph. Differentially expressed genes in each 142 cluster were identified by the *findClusterMarkers* function of *gficf* package, which compares the expression 143 of a gene in each cluster versus all the other by using the Wilcoxon rank-sum test<sup>65</sup>.

144 TGCA bulk expression dataset and cell-line deconvolution: Raw bulk expression data and relative patient 145 clinical information were collected from the Genomic Data Commons (GDC) portal<sup>70</sup> by using the 146 TCGAbiolinks package<sup>71</sup>. Then, raw counts were normalized using the EdgeR package<sup>67</sup> into R statistical 147 environment 3.6. Bisque tool<sup>50</sup> (available at https://github.com/cozygene/bisque) was used to estimate the 148 cell-line composition from the patient's bulk gene expression profile. Specifically, we applied the 149 ReferenceBasedDecomposition function with parameters: bulk.eset set to the bulk gene expression dataset in 150 log2 scale; sc.eset set to our single-cell BC atlas with normalized raw counts rescaled in log2; use overlap 151 set to FALSE and markers set to the marker genes across the 32 BC cell-lines estimated by using the function 152 *findClusterMarkers* of *gficf* package. As in the original manuscript describing the Bisque tool<sup>50</sup>, only marker 153 genes with an FDR<0.5 and Log2 fold change greaten then 0.25 were used for deconvolution purpose.

154 Spatial sequencing data: Spatial transcriptomic data of two BC patients were download from 10x Genomic
 155 website (<u>https://www.10xgenomics.com/resources/datasets</u>). Only tiles reported to be "in tissue" according
 156 to the related metadata of each patient slide were used.

Mapping new cells into the BC atlas and estimation of the cancer subtype: New points were mapped to the UMAP space via *embedNewCells* function of *gficf* package<sup>65</sup>. Briefly, tiles from 10x spatial transcriptomics were normalized with *gficf* package using the ICF weight estimated on the BC atlas. Then tiles were projected to the existing PC space using gene loadings from the BC atlas. After this transformation, tiles were mapped to the BC atlas via *umap\_transform* function of *uwot* package. Finally the cancer subtype of each mapped tile was predicted with the function *classify.cells* of the package *gficf* with the k nearestneighbour parameter set to 7.

164 Single-cell drug sensitivity prediction: The naïve gene expression profile (RNA-seq) of about 1,000 cancer cell line was obtained from the Cancer Cell Line Encyclopaedia (CCLE) portal<sup>72</sup>. Cell lines belonging to 165 166 liquid tumour were discarded and only 658 cell lines belonging to solid tumours were retained and used for 167 further analysis. The raw counts of each gene were normalized with edgeR package <sup>67</sup> and transformed in 168 log10(CPM+1). Poorly expressed genes and genes whose entropy was in the fifth percentile were excluded 169 from the analysis. Expression profiles of the 658 CCLs were then crossed with drug sensitivity data<sup>2</sup>. This 170 dataset was originally composed of 481 small molecules, but, after removing drugs for which the in vitro 171 response was available for less than 25 CCLs, only 450 small molecules were retained for further analysis. 172 For each gene and for each of the 450 drugs, we computed the Pearson correlation coefficient (PCC) between 173 the expression of the gene across the 658 cell lines and the effect of the drug expressed in terms of Area 174 Under the Curve (AUC). Since the AUC reflects the in vitro response of a cell line to different concertation 175 of a drug in a timeframe of 72 hours, lower values of AUC are associated with sensitivity whereas higher 176 values with resistance to the drug. Hence, genes positively correlated with the AUC are potential markers of 177 resistance (the more expressed the gene, the higher the concentration needed to inhibit growth), vice-versa, 178 negatively correlated genes are markers of sensitivity. We this approach, we generated a ranked list of 179 expression-based biomarkers of drug sensitivity and resistance for each of the 450 drugs where genes

positively correlated with the AUC are at the top, and those negatively correlated at the bottom. Finally, to predict drug sensitivity at the single-cell level, we used the top 250 expressed genes of each cell as input of Gene Set Enrichment Analysis (GSEA) <sup>56</sup> against the ranked list of biomarkers for each one of 450 drugs built as described above. Hence, while a negative enrichment score implies that genes associated to drug sensitivity are highly expressed by the cell, a positive one indicates the cell express genes conferring drug resistance. GSEA and associated p-values were estimating using the *fgsea* package in the R statistical environment version 3.6.

187 Drug sensitivity of the HER2+ and HER2- subpopulations in the MDA-MB-361 cell line: For each 188 sequenced cell of the MDA-MB-361 cell line, the enrichment score of 450 anticancer drugs was predicted as 189 described above. Then, to identify drugs exhibiting differential sensitivity for the two subpopulations, we 190 used the Mann-Whitney test was to assess if there was a difference between the enrichment scores of HER2+ 191 and HER2- subpopulations. P-values were corrected for false discovery rate using Benjamini-Hochberg 192 correction. A drug was considered specific for HER2- cell population if and only if its FDR was less than 193 0.05 and the median enrichment score across HER2- cells less than zero while its median enrichment score 194 across HER2+ cells greater than zero. Conversely, a drug was considered specific for HER2+ cell population 195 if and only if FDR was less than 0.05 and the median enrichment score across HER2+ cells less than zero 196 while its median enrichment score across HER2- cells greater than zero.

197 Validation of drug sensitivity prediction: Precision of the DREEP method in predicting drug sensitivity 198 from single cell transcriptional profiles was evaluated using an independent publicly available drug screening 199 dataset<sup>9</sup> composed by 1,001 CCLs and their maximal inhibitory concentration (IC50) values for 265 small 200 molecules. Hence, we applied DREEP to the single-cell profiles of the 32 BC cell lines to predict the 201 percentage of sensitive cells in each cell line for the 86 drugs. The "golden standard" was built by assigning 202 to each of  $32 \times 86$  (=2,752) cell line/drug pair the value 1 if the cell line was sensitive to the drug and 0 203 otherwise. To determine if a cell line was sensitive or not to a specific drug from the experimental data, we 204 converted for each drug its IC50 distribution in Z-scores using all the 1,001 available cell lines and then 205 defined a cell line sensitive to the drug if and only if its Z-score was in the 5% percentile. Finally, Positive 206 Predicted Values (PPV) were defined as TP/(TP+FP) where TP represents the number of true positives and 207 FP the number of false positives predicted cell lines/drug pairs.

Prediction of cell cycle phase from scRNA-seq: The cell cycle phase of each sequenced cell was predicted
 using the function *CellCycleScoring* of the *Seurat* tool with default parameter and following what was
 suggested in the corresponding vignette (<u>https://satijalab.org/seurat</u>).

211 HER2 antibody staining procedure for flow cytometry analysis: Cells were first washed with phosphate-212 buffered saline (PBS) 1x, detached with 0.05% trypsin-EDTA, resuspended and harvested with the 213 appropriate medium in single-cell suspension. Then, cells were counted, washed with PBS-FBS 1%, and 214 finally incubated for 15 min at 4° in the dark at the concentration of  $1.0 \times 10^6$  cell/µL with staining buffer. 215 The staining buffer was prepared diluting the mouse anti-human HER2 antibody (BD BB700) at the final 216 concentration of 0.00114 ng/µL. Then, to remove unbound antibody, cells were washed three times with 217 PBS-FBS 1%. Flow cytometry measurements were performed on either BD Accuri C6 or BD FACSAria III 218 instruments. To define antibody positive and negative cells, the unstained samples were used to set the gate. 219 To record data, at least  $1.0 \times 10^4$  events were collected for each sample. Data analysis was performed using 220 the either BD FACSDiva 8.0.1 or BD Accuri C6 software.

221

222HER2 expression dynamics experiment: Sorting of MDA-MB-361 HER2-positive and HER2-negative223cells was performed following the antibody staining procedure described above with the only exception that224before sorting, each sample was resuspended in sorting buffer (PBS 1x, FBS 1%, trypsin 0.1%, EDTA 2mM).225Then,  $4.0 \times 10^5$  cells were collected for each cell subpopulation (*i.e.* HER2-positive and HER2-negative),226plated in their appropriate medium, and incubated at 37°. After 18 days, the percentage of cells expressing227HER2 protein was checked by performing the antibody staining procedure described above.

229 Drug sensitivity assay: Cells were seeded in 96-well microplates (PerkinElmer); the seeding cell confluency 230 was specifically optimized for each cancer cell line to have cells in growth phase at the end of the assay. 231 After overnight incubation at 37°, cells were treated with DMSO (Merck) for the negative control and with 232 five concentrations of selected drugs in triplicate. Cells were then incubated at 37° for 72hr. Cell viability 233 was assessed by measuring either luminescence with GloMax<sup>®</sup> Discover instrument from Promega or by 234 nuclei count using the Operetta instrument from PerkinElmer. Luminescence measurements were normalized 235 using background wells as manufacturer protocol. For luminescence measurement, cells were treated with 236 Promega CellTiter-Glo® Luminescent Cell Viability Assay according to the manufacturer protocol. For 237 nuclei count, cells were washed with PBS 1x, fixed with paraformaldehyde (PFA) 4% for 10 min at room 238 temperature, washed with PBS 1x, incubated at room temperature in the dark with HOECHST 33342 239 (Thermo Fisher Scientific) diluted 1:1000 in PBS 1x for 10 min and finally washed with PBS 1x. Nuclei 240 count was performed using Columbus image analysis software (PerkinElmer). All drug used in this study 241 were purchased from Selleckchem.

242

#### 243 Data availability:

Raw sequence data of BC single cell atlas are available on Gene Expression Omnibus(GEO) repository under the accession number.

246

Code availability: The code to reproduce main results in the manuscript is available on
github at the following address <u>https://github.com/dibbelab/singlecell\_bcatlas</u>. Moreover,
the single cell atlas can be explored at <u>http://bcatlas.tigem.it</u>.

250

Acknowledgments: This work was supported by the STAR (Sostegno Territoriale alle
Attività di Ricerca) grant of University of Naples Federico II and the AIRC (Associazione
Italiana Ricerca sul Cancro) GRANT MFAG 23162 to GG and by the AIRC (Associazione
Italiana Ricerca sul Cancro) Grant IG 2016-18479 to DB and by iPC project H2020 826121
for both GG and DB.

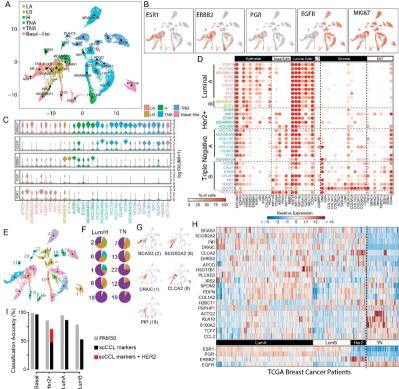
256

Author Contribution: GG performed all computational analysis, conceived the method for single-cell drug sensitivity prediction and contributed to the writing of the manuscript. GV implemented the dropseq platform, performed single-cell RNA sequencing and drug response validations. BT performed cytometric analyses, helped with cell culture and RNA-seq library preparation. AI and RB contributed to data discussion and writing of the manuscript. DdB supervised the work, contributed to the writing of the manuscript, and conceived the original idea.

264

#### 265 **Conflicts of interest**

- 266 The Authors declare no conflict of interests.
- 267



Her2+ -LumA umB.

