

CellPhe: a toolkit for cell phenotyping using time-lapse imaging and pattern recognition

Julie C Wilson (✉ julie.wilson@york.ac.uk)

University of York

Laura Wiggins

University of York

Alice Lord

University of York

Peter O'Toole

University of York

William Brackenbury

University of York

Article

Keywords: cell phenotyping, pattern recognition, ensemble classification

Posted Date: October 28th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-971415/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

CellPhe: a toolkit for cell phenotyping using time-lapse imaging and pattern recognition

Laura Wiggins^{1,2}, Alice Lord², Peter J. O’Toole^{1,2}, William J. Brackenbury^{1,2} and Julie Wilson^{3*}

¹ York Biomedical Research Institute, University of York, York, UK

² Department of Biology, University of York, York, UK

³ Department of Mathematics, University of York, York, UK

* Corresponding author

email: julie.wilson@york.ac.uk

Abstract. With phenotypic heterogeneity in whole cell populations widely recognised, the demand for quantitative and temporal analysis approaches to characterise single cell morphology and dynamics has increased. We present CellPhe, a pattern recognition toolkit for the characterisation of cellular phenotypes within time-lapse videos. To maximise data quality for downstream analysis, our toolkit includes automated recognition and removal of erroneous cell boundaries induced by inaccurate tracking and segmentation. We provide an extensive list of features extracted from individual cell time series, with custom feature selection to identify variables that provide greatest discrimination for the analysis in question. We demonstrate the use of ensemble classification for accurate prediction of cellular phenotype and clustering algorithms for the characterisation of heterogeneous subsets. We validate and prove adaptability using different cell types and experimental conditions. Our methods could be extended to other imaging modalities, such as fluorescence, and would be suitable for all time-lapse studies including clinical applications.

1 Introduction

Heterogeneity in whole cell populations is a long-standing area of interest^{1,2,3} and previous studies have identified cell-to-cell phenotypic and genotypic diversity even within clonally derived populations.⁴ The emergence of methods such as single-cell RNA sequencing has enabled characterisation of subsets within a population from gene expression profiles,⁵ yet these methods involve collection of data at discrete time points, missing the subtle temporal changes in gene expression on a continuous scale. Such methods exclude information on single-cell morphology and dynamics, yet cellular phenotype plays a crucial role in determining cell function,^{6,7} disease progression,⁸ and response to treatment.⁹ There remains a demand for quantitative and temporal analysis approaches to describe the subtleties of single-cell heterogeneity and the complexities of cell behaviour.

Modern microscopy advancements facilitate the ability to produce information-rich images of cells and tissue, at high-throughput and of high quality. Temporal changes in cell behaviour can be observed through time-lapse imaging, but the task of identifying individual cells and following them over time is an ongoing computer vision challenge.^{10,11} Processing of cell time-lapse images can be broken down into two stages: segmentation, the ability to detect cells as regions of interest (ROIs) and distinguish them from background, and tracking, providing each identified cell with a unique identifier which is retained in consecutive frames. Imaging artefacts vary between experiments and finding solutions to issues such as background noise, inhomogeneity of cell size and texture and overlapping cells is still a challenge for biomedical research.¹² Reliable cell segmentation protocols are non-deterministic and experiment-specific.¹³

Recent years have seen the emergence of user-friendly software that uses machine learning algorithms to assist in objective, high-throughput cell segmentation and tracking, allowing users to focus their attention on data analysis.^{14,15} However, to ensure precise quantification of cell morphology and motility, and to monitor major cellular events such as mitosis and apoptosis, it is vital that instances of erroneous segmentation and tracking are removed from data sets prior to downstream analysis methods.¹⁶ Manual removal of such errors is heavily labour-intensive, particularly when time-lapses take place over several days.

49 Here we present CellPhe, a pattern recognition toolkit that uses the output of segmentation
50 and tracking software such as CellProfiler¹⁷ and Ilastik.¹⁸ To maximise data quality for down-
51 stream analysis, CellPhe includes the recognition and removal of erroneous cell boundaries induced
52 by inaccurate segmentation and tracking. Customised feature selection is used to identify the most
53 discriminatory variables for a particular objective from an extensive list of features extracted from
54 the time-course images. These extracted variables quantify cell morphology, texture and dynam-
55 ics and describe temporal changes and can be used to reliably characterise and classify individual
56 cells as well as cell populations. We demonstrate the use of ensemble classification for accurate
57 prediction of cellular phenotype and clustering algorithms for identification of heterogeneous sub-
58 sets. We exemplify CellPhe by characterising the behaviour of untreated and chemotherapy treated
59 breast cancer cells from ptychographic time-lapse videos. Quantitative phase images (QPI)^{19,20,21}
60 do not make use of fluorescent labels, though their usage enhances the difference in intensity of
61 cells and background which often improves cell segmentation accuracy. We show that our methods
62 successfully recognise and remove a population of erroneously segmented cells, improving data set
63 quality without fluorescence-induced perturbation.²² Morphological and dynamical changes induced
64 by chemotherapeutics, particularly at low drug concentration, are often more subtle than those that
65 discriminate distinct cell types and we demonstrate the ability of CellPhe to automatically identify
66 time series differences induced by chemotherapy treatment, with the chosen variables proving sta-
67 tistically significant even when not observable by eye.

68 The complexities of heterogeneous drug response and the problem of drug resistance further mo-
69 tivate our chosen application. The ability to identify discriminatory features between treated and
70 untreated cells can allow automated detection of "non-conforming" cells such as those that possess
71 cellular drug resistance. Further investigation of such features could elucidate the underlying bio-
72 logical mechanisms responsible for chemotherapy resistance and cancer recurrence. We validate the
73 adaptability of CellPhe with both a different cell type and a different drug treatment and show that
74 variables are selected according to experimental conditions, tailored to properties of the cell type
75 and drug mechanism of action.

76 CellPhe would extend to other imaging modalities such as fluorescence images, where further
77 extracted variables such as fluorescence intensity would complement our existing list of metrics. A
78 comprehensive manual with a working example is provided to guide the user through the complete
79 workflow.

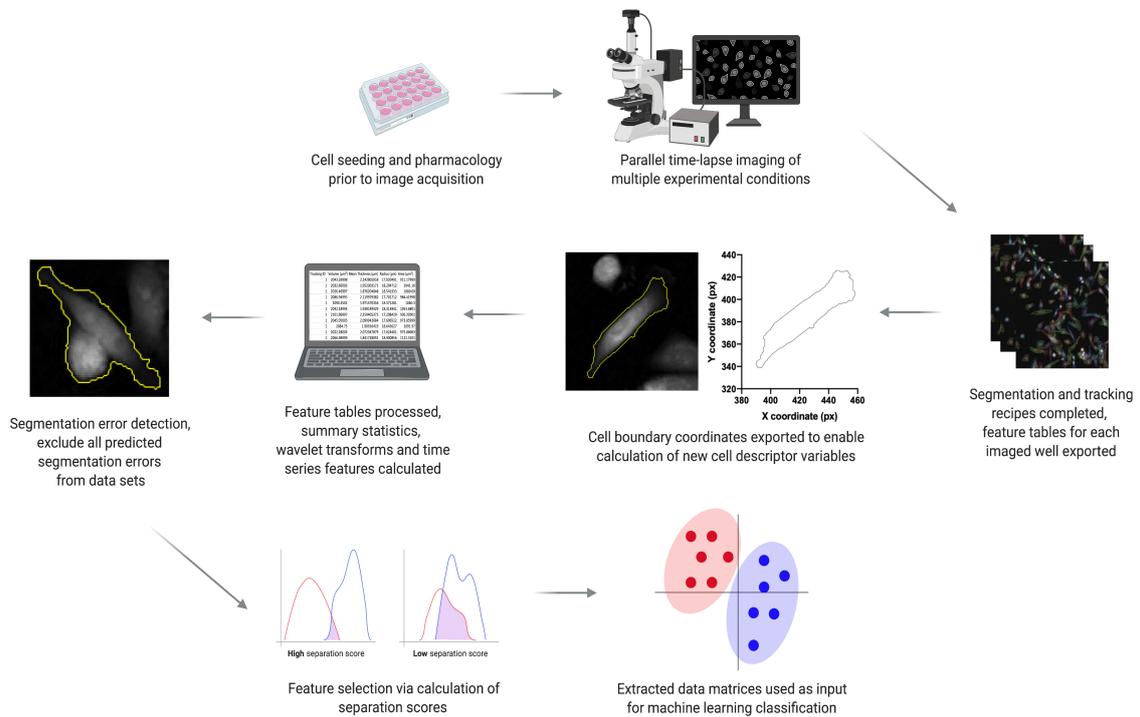
80 2 Results

81 Overview of CellPhe

82 CellPhe is a toolkit for the characterisation and classification of cellular phenotypes from time-
83 lapse videos, a diagrammatic summary of CellPhe is provided in **Figure 1**. Experimental design
84 is determined by the user prior to image acquisition where seeded cell types and pharmacology are
85 specific to the user's own analysis. Example uses are discrimination of cell types (e.g. neurons vs.
86 astrocytes), characterisation of disease (e.g. healthy vs. cancer) , or assessment of drug response
87 (e.g. untreated vs. treated). The user can then time-lapse image cells for the desired amount of time,
88 using an imaging modality of their choice. Once images are acquired and segmentation and tracking
89 of cells are complete, cell boundary coordinates are exported and used for calculation of an extensive
90 list of morphology and texture features. These together with dynamical features and extracted time
91 series variables are used to aid removal of erroneous segmentation by recognition of error-induced
92 interruption to cell time series. Once all predicted segmentation errors have been removed from
93 data sets, feature selection is performed and only features providing separation above an optimised
94 threshold are retained. This identifies a list of most discriminatory features and allows the user
95 to explore biological interpretation of these findings. The extracted data matrices are then used
96 as input for ensemble classification, where the phenotype of new cells can be accurately predicted.
97 Furthermore, clustering algorithms can be used to identify heterogeneous subsets of cells within the
98 user's data, both inter- and intra-class.

99 The remaining results exemplify the use of CellPhe with a biological application, characterisation

100 and classification of chemotherapeutic drug response. We look at each of the CellPhe stages in
 101 detail (segmentation error removal, feature selection, ensemble classification and cluster analysis) and
 102 demonstrate that each step provides interpretable, biologically relevant results to answer experiment
 103 specific questions and aid further research.



104

Figure 1: Summary of the CellPhe toolkit. Following time-lapse imaging, acquired images are processed and segmentation and tracking recipes implemented. Cell boundary coordinates are exported, features extracted for each tracked cell and the time series summarised by characteristic variables. Predicted segmentation errors are excluded and optimised feature selection performed using a threshold on the class separation achieved. Finally, multiple machine learning algorithms are combined for classification of cell phenotype and clustering algorithms utilised for identification of heterogeneous cell subsets.

105

106 **CellPhe application: characterising chemotherapeutic drug response**

107 The 231Docetaxel data set, obtained from multiple experiments involving MDA-MB-231 cells,
 108 both untreated and treated with 30 μ M docetaxel, is the main data set used to demonstrate our
 109 method. We show that the same analysis pipeline can be applied to other data sets by considering
 110 both a different cell line, MCF-7, in the MCF7Docetaxel data set, and a different drug, doxorubicin,
 111 with the 231Doxorubicin data set. In each case, we remove segmentation errors, as described in
 112 Section 2.5, before using feature selection (Section 2.6) to identify discriminatory variables tailored
 113 to the particular data set. We show that different variables are chosen depending on the inherent
 114 nature of the cell line and the effect of the drug in question. By using these features in classification
 115 algorithms, we aimed to characterise the behaviour over time of untreated and treated cells.

116

117 **Segmentation Error Removal**

118 The purpose of this analysis was to improve the quality of our data sets prior to untreated
 119 vs. treated cell classification by automating detection of segmentation errors and optimising the
 120 exclusion criteria of predicted errors.

121 Comparison of time series for cells with and without segmentation errors showed many of our
 122 features to be sensitive to such errors, motivating the need to remove these cells prior to treatment
 123 classification. Size metrics, such as volume, were particularly affected by segmentation errors as
 124 under- or over-segmentation could result in halving or doubling of cell volume respectively (**Figure**

125 **2a)**. This noticeable disruption to the time series of several features suggested that reliable detection
126 of segmentation errors would be possible.

127 After excluding 62 instances identified as tracked cell debris, a training data set for MDA-MB-231
128 cells (from the 231Docetaxel data set), was obtained, consisting of 1185 correctly segmented cells
129 and 278 cells with segmentation errors. The number of cells in the segmentation error class was
130 doubled using SMOTE and the resulting data set with 1741 observations used for the classification
131 of segmentation errors as described in Section 2.5. The MDA-MB-231 cells (from 231Docetaxel and
132 231Doxorubicin, both untreated and treated) that were not used for training formed independent
133 test sets (Table 1).

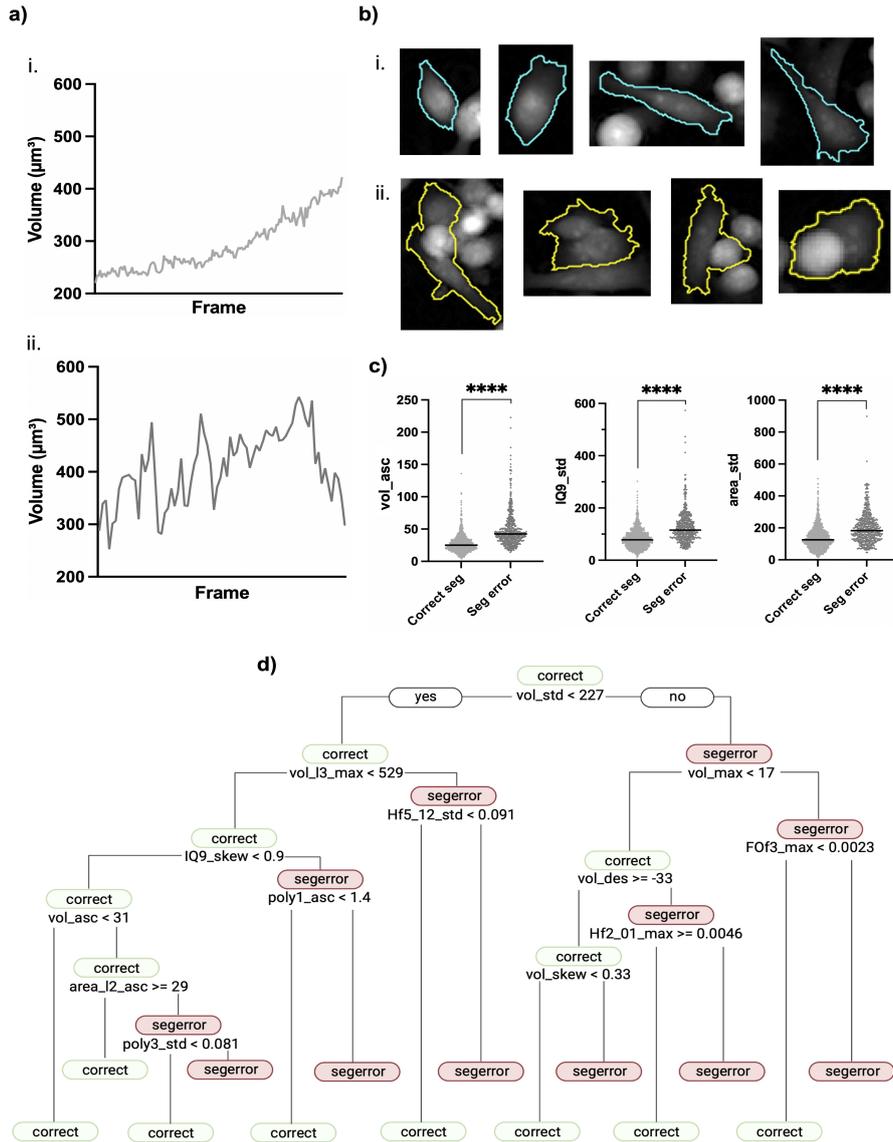
134 A total of 225 of the 1478 cells in the 231Docetaxel test set were predicted to be segmentation
135 errors. Of these, 219 were confirmed by eye to be true segmentation errors, most of which were due
136 to under- or over-segmentation throughout their time series. Other segmentation issues observed
137 included background pickup, cells swapping cell ID, and cells repeatedly entering and exiting the
138 field of view, all of which result in problem time series (**Figure 2b**). Although the 6 cells that were
139 misclassified as segmentation errors were all large cells, further investigation showed that removal of
140 these cells did not exclude an important subset from the data. Comparison of the area and volume
141 of these cells with correctly segmented cells, confirmed that a population of very large cells were
142 correctly classified as having no segmentation errors so that such cells were still represented in the
143 test set (**Figure S1**).

144 This classifier was also used to identify a further 83 segmentation errors from the 1005 cells in
145 the 231Doxorubicin data set, all 83 were confirmed by eye to be true segmentation errors (Table 1).
146 It was necessary to train a new classifier for MCF-7 segmentation error detection due to differences
147 between the cell lines. In this case 308 correctly segmented cells and 192 segmentation errors were
148 identified by eye. After applying SMOTE to double the number of segmentation error observations,
149 a classifier was trained with the resulting 692 observations as described in section 2.5. Just 4 of the
150 287 cells in the MCF7Docetaxel test set were classified as segmentation errors and were confirmed
151 by eye to be true segmentation errors.

Data set	TP	FP
231Docetaxel (1478)	219	6
231Doxorubicin (1005)	83	0
MCF7Docetaxel (287)	4	0

Table 1: Segmentation error prediction on the test data. The number of correctly classified segmentation errors (True Positives, TP) and the number of correctly segmented time series incorrectly classified as segmentation errors (False Positives, FP) are shown. The number of cells in each test data before segmentation error removal is shown in parentheses.

152



153

Figure 2: (a) Volume time series for **i.** a correctly segmented cell and **ii.** a cell experiencing segmentation errors, demonstrating greater fluctuation in volume when a cell experiences segmentation errors. (b) Examples of test set cells classified as **i.** correct segmentation and **ii.** segmentation error. (c) Beeswarm plots of features that are significant for identifying segmentation errors in the 231Docetaxel training set (****: $p < 0.0001$). Interruptions to time series induced by segmentation errors are identified by increased ascent and standard deviation for affected features. (d) A representative 231Docetaxel trained decision tree, demonstrating how volume, texture and polygon (shape) variables are used in combination to make classifications.

154

155 Although feature selection was not used in the identification of segmentation errors, we did calcu-
 156 lulate separation scores for the MDA-MB-231 training data to investigate the effect of such errors.
 157 As might be expected, volume was most affected, with segmentation errors resulting in larger stan-
 158 dard deviation, ascent and maximum value. Other features with high separation scores included
 159 area as well as spatial distribution descriptors with the highest thresholds, features that detect the
 160 clustering of high intensity pixels, characteristic of cell overlap and over-segmentation (**Figure 2c**).
 161 Analysis of the trained decision trees showed that a combination of size, texture and shape variables

162 frequently formed the most important features for detecting segmentation errors with MDA-MB-231
163 cells, see **Figure 2d** for an example.

164 For the MCF7Docetaxel data set, velocity was found to be important in determining whether
165 or not a cell experienced segmentation errors in addition to texture and shape variables. The cell
166 centroid, used to determine position and hence velocity, is affected by boundary errors and so high
167 velocity, uncharacteristic of MCF-7 cells, is a good indication of segmentation error for these cells.

168 **Feature Selection**

169 For the 231Docetaxel data set, the calculation of separation scores identified variables that pro-
170 vided good discrimination between untreated MDA-MB-231 cells and those treated with 30 μ M doc-
171 etaxel. As separation scores do not provide information on how these variables work in combination,
172 we performed Principal Component Analysis (PCA) to explore relationships between discriminatory
173 variables.

174
175 Differences in the appearance of MDA-MB-231 cells induced by docetaxel treatment were ob-
176 served by eye from cell timelapses. Untreated cells displayed a spindle-shaped morphology (a circular
177 cross-section with tapering at both ends), with contractions and protrusions facilitating migration.
178 Cells that received treatment were generally dense and spherical, and increased in size following a
179 failed attempt at cytokinesis (**Figure 3a**). Discriminatory features identified by calculation of sepa-
180 ration scores were consistent with differences observed by eye, the 20 variables that achieved greatest
181 separation are shown in **Figure 3b**. The most discriminatory feature was the ratio of pixels within
182 the cell boundary to the number of pixels within the minimal bounding box (A2B). Low values
183 indicate a small cell area within a large bounding box, as is the case for spindle-shaped cells. Other
184 morphological features that successfully characterised untreated MDA-MB-231 cell morphology were
185 length, radius and polygonal representation variables, all of which discriminate between irregular
186 untreated cell shape and the more rounded morphology of treated cells. Furthermore, separation
187 scores highlighted differences in the texture of cells following treatment, with both first order and
188 Haralick features providing good discrimination between untreated and treated cells.

189 Principal Component Analysis (PCA) demonstrated that the main variance within the data
190 arises due to class differences, with separation of classes observed across PC1 which explains 54.2%
191 of the total variance (**Figure 3c**). The dispersion of points within the scores plot illustrates hetero-
192 geneity of cells both inter- and intra-class. The non-conformity of some cells, for example treated
193 cells behaving as untreated cells, is demonstrated by points clustering within the opposite class.
194 The PC1 loadings with greatest absolute values are given in **Figure 3c**. Notably, only average
195 area to minimal box ratio (A2B_mean) and variance in edge length of the polygonal approximation
196 (poly2_mean) had positive PC1 loadings, indicating that treated cells had greater values for these
197 variables in comparison to untreated cells. The remaining variables had negative PC1 loadings, in-
198 cluding standard deviation, ascent and descent of morphological features such as length and minimal
199 box, and several texture features. As the majority of untreated cells had negative PC1 scores we
200 deduced that greater standard deviation, ascent and descent of features for untreated cells indicates
201 that these cells experience increased fluctuation throughout their time series. As treated cells mainly
202 had positive PC1 scores, they experience less fluctuation throughout their time series and instead
203 display greater stability. Identified differences in feature time series are visualised in **Figure 3d**.

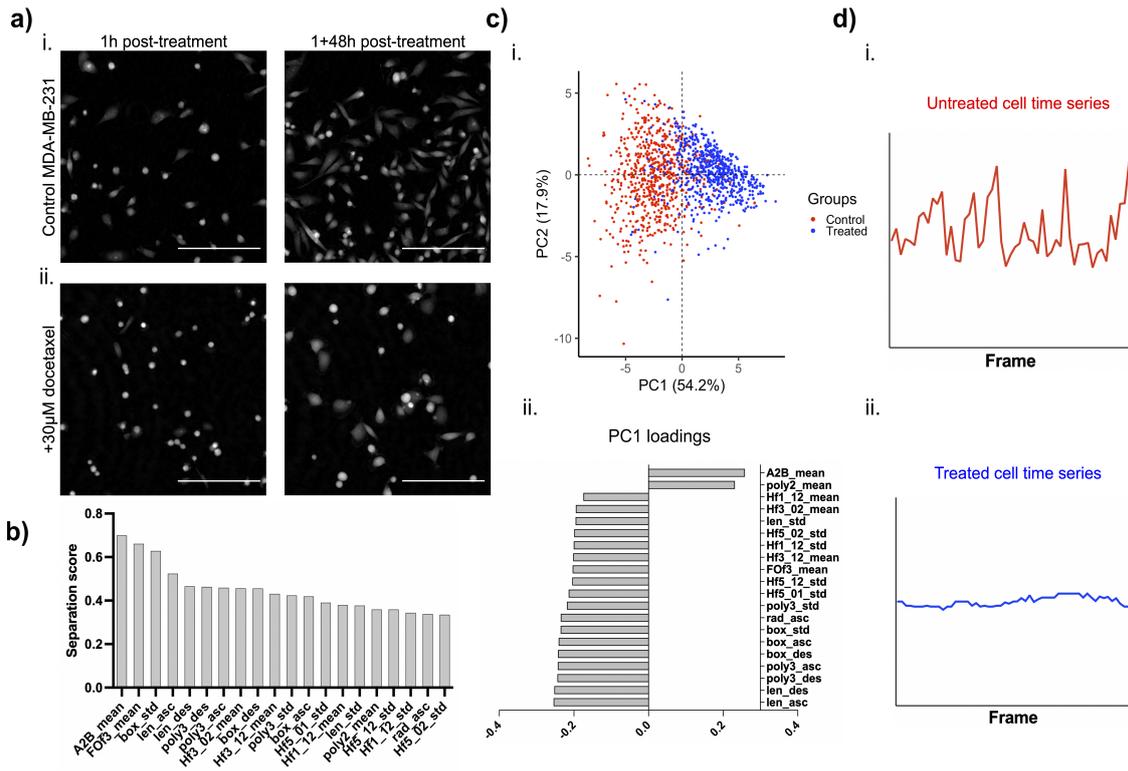


Figure 3: a) Images taken from cell timelapses of **i.** untreated MDA-MB-231 cells and **ii.** 30 μ M docetaxel treated MDA-MB-231 cells. Scale bar = 200 μ m. Increased cell count at 1+48h post treatment demonstrates healthy proliferation of untreated cells. Static cell count at 1+48h for treated cells is a result of cell cycle arrest and failed cytokinesis, leading to enlarged cell phenotype. **b)** The 20 most discriminatory features identified through calculation of 231Docetaxel separation scores. **c) i.** PCA scores plot with points coloured according to true class label. Observable separation of classes along PC1 demonstrates that the greatest source of variance within the data arises due to class differences. **ii.** PC1 loadings suggest increased activity within untreated cell time series, with standard deviation, ascent and descent variables obtaining negative PC1 loadings and scores. **d)** Representative feature time series plots for **i.** untreated MDA-MB-231 cells and **ii.** 30 μ M docetaxel treated MDA-MB-231 cells. Untreated cells experience greater fluctuation within their time series in comparison to treated cells where activity is more stabilised.

We assessed the adaptability of our feature selection method by calculating separation scores for both a different cell line and a different treatment, using PCA to evaluate the main sources of variance. We compared MCF-7 cells treated with 1 μ M docetaxel with untreated MCF-7 cells, and MDA-MB-231 cells that were treated with 1 μ M doxorubicin with untreated MDA-MB-231 cells and found that changes in the morphology and motility of cells upon treatment were both drug and cell-line specific with different variables selected. (**Figure 4**).

As was observed within the 231Docetaxel timelapses, cells increased in size due to failed cytokinesis. However, MCF-7 cells maintained a polygonal, epithelial-like morphology following treatment similar to that of the untreated population. Furthermore, no differences in movement were observed within the MCF7Docetaxel data set due to the poorly aggressive, non-invasive nature of MCF-7 cells described previously.²³ Conversely, remarkable differences in cellular dynamics were observed within the 231Doxorubicin data set, with motility of cells being severely hindered following treatment, particularly after the 24-hour time point. Only subtle differences in size and morphology of cells were observed by eye, with doxorubicin treated cells appearing slightly enlarged as a result of cell cycle arrest. Both untreated and treated sets contained examples of cells in G1 and G2, hence varied cell morphology can be observed within both (elongated and adherent cells in G1, round and dense morphology of cells in G2.)

223 The 20 variables that achieved greatest separation for each of the MCF7Docetaxel and 231Dox-
 224 orubicin data sets are shown in **Figure 4b**. Primarily size and texture variables were identified as
 225 most discriminatory for MCF7Docetaxel with variables such as length, width and area characterising
 226 the enlarged cell shape of treated cells. Spatial distribution variables were chosen for several inten-
 227 sity thresholds, demonstrating differences in the clustering of pixels, following docetaxel treatment.
 228 Furthermore, mean cell density was also identified as a discriminatory variable with the untreated
 229 cell population having greater mean cell density than the treated population, likely a result of de-
 230 creased cell proliferation and cell-cell adhesion for treated cells. As was observed by eye, movement
 231 features formed the majority of discriminatory variables for the 231Doxorubicin data set, with un-
 232 treated cells having greater velocity, tracklength and displacement than treated cells. Differences in
 233 movement were also described through density ascent and descent, as cell density fluctuated more
 234 for untreated cells due to the increased likelihood of passing neighbouring cells when migrating.
 235 Subtle differences in cell shape and size observed by eye upon doxorubicin treatment were described
 236 by changes in the area to minimal box ratio and width and radius variables. Notably both data
 237 sets received lower separation scores than the 231Docetaxel data set, with MCF7Docetaxel having
 238 the lowest. This effectively provides a measure of class similarity, with high separation scores for
 239 231Docetaxel indicative of significant changes to cells upon treatment and low separation scores for
 240 MCF7Docetaxel suggesting these changes are more subtle.

241 PCA scores plots obtained with the selected features are shown in **Figure 4c**. Differences be-
 242 tween classes can be observed for the MCF7Docetaxel data set although the separation involves
 243 both PC1 and PC2 whilst the greatest source of variance, along PC1 (55% of the total variance),
 244 due to heterogeneity within treated cells. We found that the subset of treated cells with lowest
 245 PC1 scores were primarily cells tracked during a failed attempt at mitosis. This subset of treated
 246 cells is characterised by fluctuations in the times series of spatial distribution features indication
 247 changes in texture as cells enter and exit failed mitosis. On the other hand, the PCA scores plot for
 248 231Doxorubicin shows the greatest source of variance to be due to class differences, with separation
 249 of classes along PC1 (69% of the total variance). All PCA scores plots demonstrated the potential
 250 to characterise untreated and treated cell behaviour, with feature selected variables providing good
 251 distinction of classes which was improved by using variables in combination.

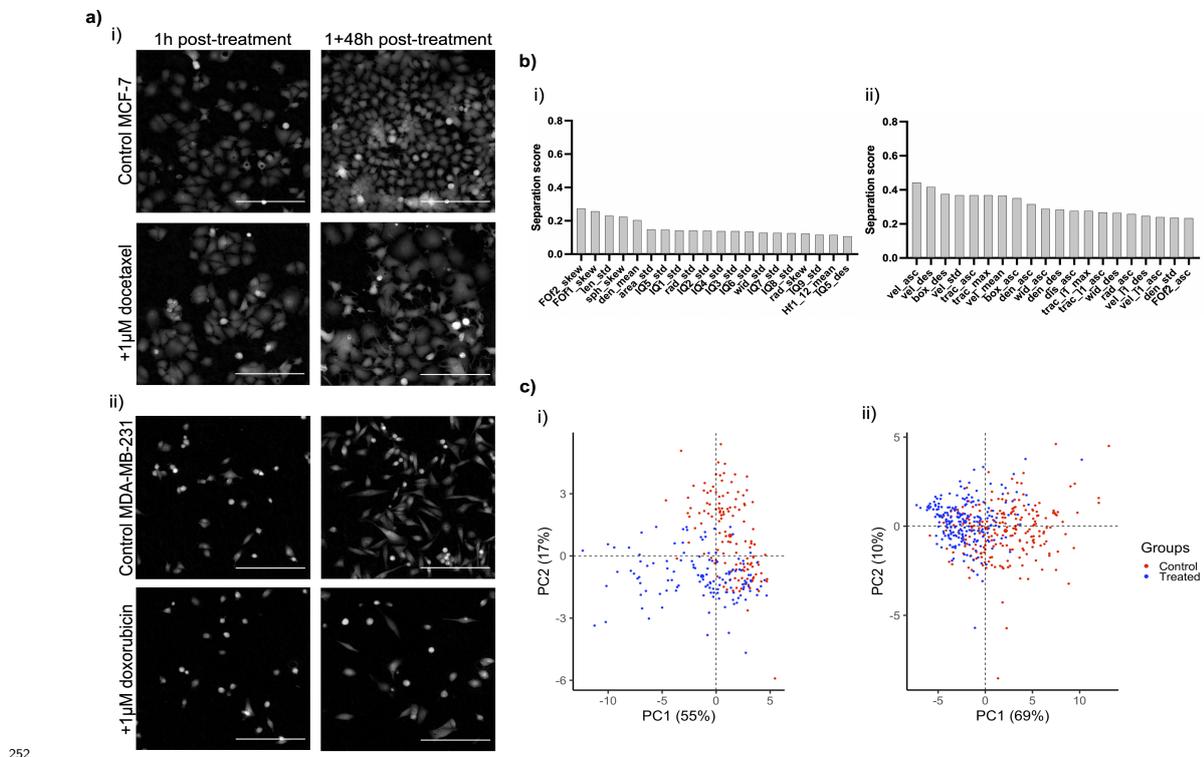


Figure 4: **a)** Images taken from cell timelapses of **i.** untreated and $1\mu\text{M}$ docetaxel treated MCF-7 cells and **ii.** untreated and $1\mu\text{M}$ doxorubicin treated MDA-MB-231 cells. Scale bar = $200\mu\text{m}$. Differences in cell count following treatment can be observed for both due to cell cycle arrest induced by docetaxel or doxorubicin respectively. Docetaxel treated MCF-7 cells display enlarged cell phenotype at the 1+48h time point due to failed cytokinesis. In comparison, differences in morphology are more subtle for doxorubicin treated MDA-MB-231 cells at the 1+48h time point. **b)** The 20 most discriminatory features identified through calculation of separation scores for **i.** MCF7Docetaxel and **ii.** 231Doxorubicin. Low separation scores for MCF7Docetaxel indicate less discrimination between untreated and treated cell populations within this data set. Several movement features provide good separation for 231Doxorubicin demonstrating differences in cell motility following doxorubicin treatment. **c)** PCA scores plots for **i.** MCF7Docetaxel and **ii.** 231Doxorubicin, colour-coded according to true class label. Notably the greatest source of variance within MCF7Docetaxel is a result of treated cell variability, with separation of classes increased when PC1 and PC2 are used in combination. The greatest source of variance within the 231Doxorubicin data set arises due to class differences with separation of classes observable along PC1.

253

254 **Classification of Treated and Untreated Cells**

255

256

257

258

259

260

We found that the distribution of separation scores differed for each data set, with the 231Docetaxel set having the greatest number of variables achieving high separation, followed by 231Doxorubicin and with MCF7Docetaxel generally having much lower separation scores (**Figure 5a**). Optimal separation thresholds of 0.2, 0.05 and 0.025 were obtained for 231Docetaxel, 231Doxorubicin and MCF7Docetaxel respectively, resulting in 82, 162 and 195 variables (of a possible 702) being selected for classifier training.

261

262

263

264

265

266

267

268

269

270

271

272

Having chosen an optimal separation threshold, we trained an ensemble classifier for each data set as described in Section 2.6. Classification accuracy scores for training and test sets obtained using our ensemble classifier are provided in **Table 2**, whilst classification accuracy scores for the individual classification algorithms before their combination can be found in **Table S3**. Through visual inspection, we found that misclassifications formed subsets of cells whose behaviour deviated from the behaviour of the main population, we call this subset "non-conforming". (**Figure 5b**). For untreated cells, we found that healthy, proliferating cells were correctly classified whereas less motile cells, cell debris or large, non-motile mutant cells were instead classified as treated. For treated cells, we found that cells experiencing the drug-induced phenotypic differences identified through feature selection were classified as treated. However, treated cells displaying behaviour similar to that of an untreated cell, such as increased migration or fluctuation and elongation in cell shape, and were classified as untreated (**Figure 5c**).

273

274

275

276

277

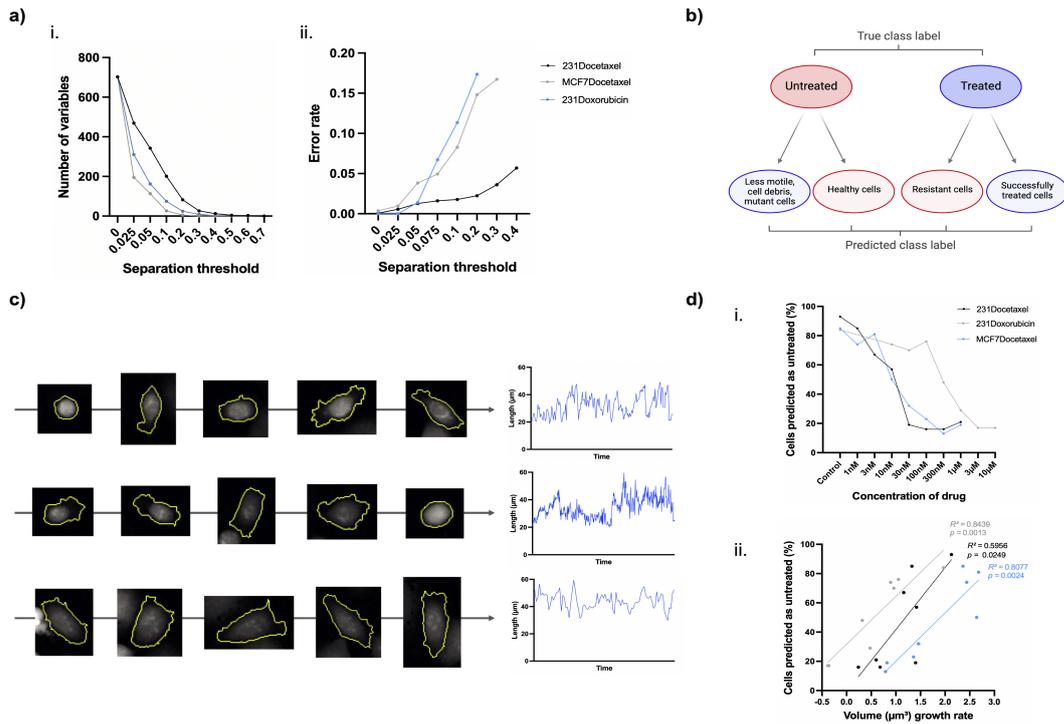
278

279

280

281

We found that the proportion of non-conforming treated cells, those classified as untreated, decreased as drug concentration increased for all three data sets (**Figure 5d**). To explore the connection between the proportion of non-conforming treated cells and the population drug response of each treated set, we considered the total volume growth rate at each drug concentration in relation to the percentage of cells predicted as untreated (**Figure 5d**). We found that the overall growth rate decreased with increased drug concentration due to more cells responding at higher concentrations. This correlated positively with the percentage of cells predicted as untreated, with a greater percentage of cells predicted as untreated for high volume growth rate with proliferation still occurring.



282

Figure 5: **a) i.** The number of variables with separation scores above different thresholds. A greater number of variables achieve high separation for 231Docetaxel in comparison to 231Doxorubicin and MCF7Docetaxel. **ii.** Optimisation of separation threshold for each data set. Thresholds of 0.2, 0.05 and 0.025 were selected for 231Docetaxel, 231Doxorubicin and MCF7Docetaxel respectively resulting in 82, 162 and 195 variables being used for classifier training. **b)** Sub-populations within each class, colour-coded according to the ideal final classification of each sub-population. Non-conforming cells for each class form a subset of misclassified cells. **c)** Examples of docetaxel treated MDA-MB-231 cells misclassified as untreated. Timelapse images demonstrate how these cells exhibit an elongated morphology characteristic of migratory untreated cells. Time series plots for cell length demonstrate the fluctuation in shape of these cells, typical of untreated cells. **d) i.** The percentage of cells predicted as untreated for a range of drug concentrations. For all three data sets, this percentage decreases as drug concentration increases due to a greater number of cells responding at higher concentrations. **ii.** Positive correlation between the total volume rate of growth and the percentage of cells predicted as untreated, with higher volume growth rates associated with a higher number of cells being predicted as untreated. Linear regression slopes were found to be significant (p values shown). R^2 correlation coefficients are also provided, demonstrating positive correlation for each data set.

283

	231Docetaxel	MCF7Docetaxel	231Doxorubicin
Train	Untreated: 96%	Untreated: 100%	Untreated: 97%
	Treated: 99%	Treated: 99%	Treated: 100%
	Overall: 98%	Overall: 100%	Overall: 99%
Test	Untreated: 97%	Untreated: 85%	Untreated: 84%
	Treated: 80%	Treated: 81%	Treated: 71%
	Overall: 93%	Overall: 85%	Overall: 81%

Table 2: Ensemble classification accuracy scores for each data set. All percentages have been rounded to the nearest whole number.

284 Subset Identification

285 Classification accuracy scores for the untreated cell population were notably greater than those
 286 of the treated population across all three of the data sets (**Table 2**), suggesting a greater proportion
 287 of non-conforming treated cells in comparison to non-conforming untreated cells. Imbalance of

288 classification accuracy scores in binary classification is often a result of hidden stratification,²⁴ where
289 poor performance of one class is a result of misclassifications of important, unlabeled subsets. To
290 investigate this phenomenon we performed hierarchical clustering on 231Docetaxel treated cells and
291 the obtained dendrogram is provided in **Figure 6a**, with examples of cells from each cluster in
292 **Figure 6b**.

293 **Figure 6c** shows the distribution of mean volumes for each cluster in comparison to the untreated
294 MDA-MB-231 population. Clusters 1, 2 and 4 span a similar range of volumes to the untreated set,
295 whereas clusters 3 and 5 have greater mean volumes. Cluster 6 is formed primarily of cell debris as
296 a result of cell death with mean volumes much lower than those of the untreated set.

297 Cells in the same cluster share similar properties (**Figure 6b**), and morphological differences
298 between clusters of different cell cycle states can be observed. For example cells in clusters 1, 2
299 and 4 are much smaller and brighter than cells in clusters 3 and 5 as the cells are heading towards
300 attempted mitosis, confirmed by visual inspection of cell timelapses, and hence resemble untreated
301 mitotic cells. However, discrimination between clusters at similar stages of the cell cycle were not
302 as readily identifiable by eye and we therefore calculated separation scores between these clusters to
303 identify discriminatory variables (**Figure 6d**). Differences between cluster 2 and clusters 1 and 4
304 were primarily based on cell shape with cells in cluster 2 having greater fluctuation in variables such
305 as length and width as they experienced repeated elongation and contraction during their tracking.
306 A range of textural variables provided highest separation between cluster 1 and cluster 4, including
307 many Haralick features, indicating differences in distribution of pixel intensities for cells within
308 these clusters. Highest separation between clusters 3 and 5 was achieved by both shape and texture
309 features. Cells in cluster 5 showed greater fluctuation in shape whereas the shape of cells in cluster
310 3 remained relatively spherical for the duration of tracking. Textural variables also achieved high
311 separation between clusters 3 and 5, with cells in cluster 5 having greater changes in the distribution
312 of interior pixel intensities.

313 Clusters also spanned a range of mean cell volumes beyond those of the untreated set when
314 hierarchical clustering was repeated for MCF7Docetaxel treated cells. However, this was not the
315 case for 231Doxorubicin treated cells and therefore k -means clustering was used to explore the
316 connection between misclassifications and hidden subsets in the 231Doxorubicin treated cell test
317 set. Two distinct clusters were obtained (**Figure 6ei**), where several movement variables achieved
318 greatest separation between the two clusters. Cells in cluster 2 experienced increased migration in
319 comparison to cells in cluster 1, with greater velocity, tracklength and displacement. We calculated
320 classification accuracy scores for the two clusters individually and found that 94% of cells in cluster 1
321 were correctly classified as treated, but only 30% in cluster 2 (**Figure 6eii**). The increased migration
322 of cells in cluster 2 mean these cells have greater similarity to the untreated population. These non-
323 conforming treated cells form the majority of treated cell misclassifications in the 231Doxorubicin
324 test set.

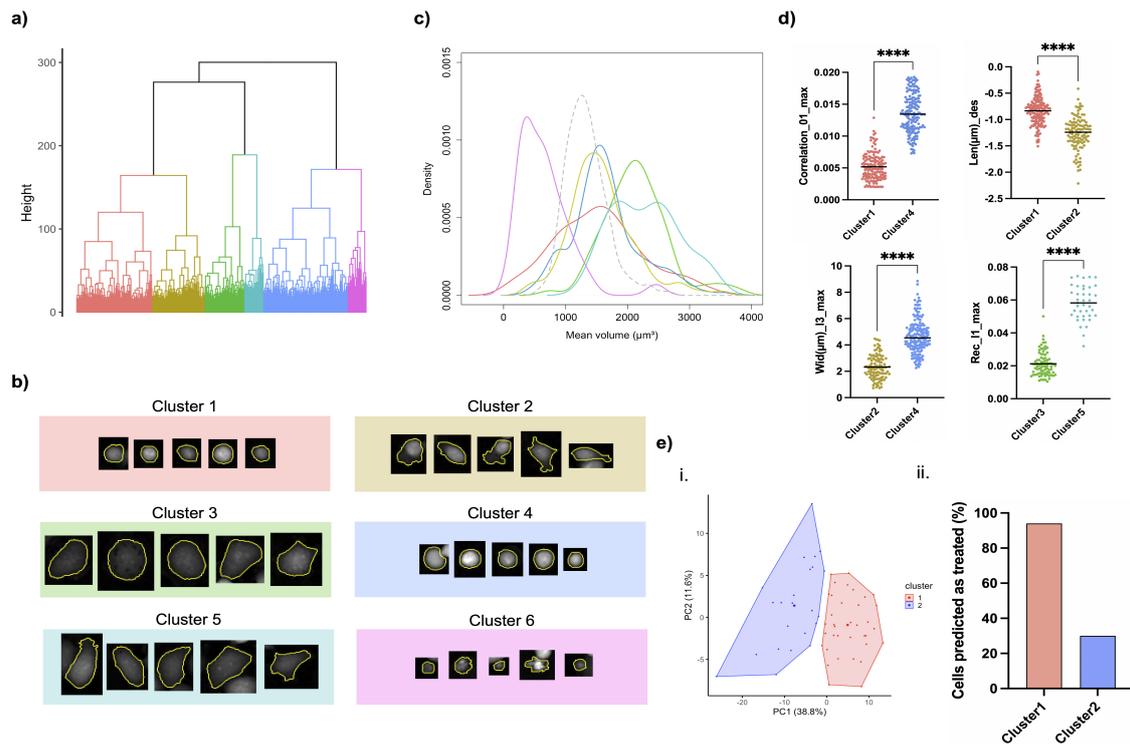


Figure 6: a) Dendrogram obtained from hierarchical clustering of 231Docetaxel treated cells, with 6 clusters coloured. b) Examples of cells from each cluster with background colours identifying the cluster. Cells within a cluster share similar properties but differ to cells in other clusters. c) Density plots of mean cell volume, colour-coded according to cluster. The gray, dashed density plot represents 231Docetaxel untreated cells. Cluster 6 (cell debris cluster) has the greatest leftward shift due to cells losing volume upon cell death. Clusters 1, 2 and 4 primarily span the same range of volumes as the untreated set as cells in these clusters have not yet attempted cytokinesis. Clusters 3 and 5 have mean volumes greater than the untreated set as cells in these clusters have continued to grow following failed cytokinesis. d) Beeswarm plots of example discriminatory variables, identified from separation scores, for clusters in similar cell cycle stages (****: $p < 0.0001$). e) i. k -means clustering of 231Doxorubicin test set treated cells. Cells are colour-coded according to which cluster they were assigned. ii. The number of cells predicted as treated for each of the clusters. Cluster 1 was formed of successfully treated cells with 94% (34/36) of cells correctly classified as treated, whereas cluster 2 formed a subset of non-conforming treated cells, with only 30% (7/23) correctly classified as treated.

3 Discussion

The CellPhe toolkit complements existing software for automated cell segmentation and tracking, using their output as a starting point for time series feature extraction, cell classification and cluster analysis. Erroneous cell segmentation and tracking can significantly reduce data quality but such errors often go undetected and can negatively influence the results of automated pattern recognition. CellPhe's extensive feature extraction followed by bespoke feature selection not only allows the characterisation and classification of cellular phenotypes from time-lapse videos but provides a method for the identification and removal of erroneous cell tracks prior to these analyses. Attribute analysis showed that different features were chosen to identify segmentation errors for different cell lines. For example, sudden increases in movement resulting from large boundary changes can indicate segmentation errors for MCF-7 cells, contrasting with their innate low motility. On the other hand, size and texture variables provide better characterisation of the unexpected fluctuations in cell size and clusters of high intensity pixels induced by segmentation errors for MDA-MB-231 cells. Current approaches for removal of segmentation errors are subjective and labour-intensive, requiring

341 manual input of parameters such as expected cell size that need to be fine-tuned for different data
342 sets. CellPhe provides an objective, automated approach to segmentation error removal with the
343 ability to adapt to new data sets.

344 For cell characterisation, we have shown that CellPhe’s feature selection method is able to adapt
345 to different experimental conditions, providing discrimination between untreated and treated groups
346 of two different breast cancer cell lines (MDA-MB-231 and MCF-7) and two different chemotherapy
347 treatments (docetaxel and doxorubicin). The discriminatory variables identified here coincide with
348 previously reported effects of docetaxel or doxorubicin treatment and can be interpreted in terms
349 of the mechanism of action of each drug. Previous studies have identified a subset of polyploid,
350 multinucleated cells following docetaxel treatment due to cell cycle arrest and occasionally cell cycle
351 slippage.²⁵ Our findings support this with shape and size variables providing the greatest separation
352 for docetaxel treatment in both MDA-MB-231 and MCF-7 cells. Many texture variables were also
353 identified as discriminatory following docetaxel treatment, providing label-free identification of the
354 multiple clusters of high intensity pixels in treated cells, likely a result of docetaxel-induced multin-
355 ucleation. We found that at a higher, sub-lethal concentration of 1 μ M, migration of MDA-MB-231
356 cells was reduced with variables associated with movement providing greatest discrimination be-
357 tween untreated and doxorubicin treated cells. This is supported by studies that have identified
358 changes in migration of doxorubicin treated cells, noting that low drug concentrations in fact facili-
359 tate increased invasion.^{26,27}

360 We found an imbalance in untreated and treated classification accuracy scores, with a greater
361 proportion of treated cells misclassified for all three data sets. This consistent imbalance suggests the
362 misclassifications are in fact representative of a subset of non-conforming, and potentially chemore-
363 sistant, cells. The concept of hidden stratification, where an unlabelled subset performs poorly
364 during classification, has been described previously²⁸ and poses a challenge in medical research as
365 important subsets (such as rare forms of disease) could be overlooked. Here, the misclassified cells
366 could be of most interest and the ability to identify non-conforming behaviour is precisely what
367 is required from a classifier as treated cells that display behaviour similar to untreated cells could
368 indicate a reduced response to drug treatment. The classification of cells treated with a range of
369 concentrations supported this hypothesis as a greater proportion of cells were classified as untreated
370 at lower drug concentrations, demonstrating that our trained ensemble classifier can be used to
371 quantify drug response, at both single-cell and populational level.

372 Cluster analysis revealed cell subsets that appear to represent different responses to drug treat-
373 ment. Heterogeneity of cellular drug response is a commonly reported phenomenon in cancer treat-
374 ment, yet mechanisms underlying this are not well understood.²⁹ Analysis of cell volumes showed
375 the mean volume of treated and untreated cells to be comparable for doxorubicin reflecting the fact
376 that this treatment can induce G1, S or G2 cell cycle arrest.³⁰ However, for docetaxel treated cells,
377 we found that clusters spanned a range of mean cell volumes beyond those of the untreated set for
378 both cell lines. Clustering allowed identification of three general responses to docetaxel treatment:
379 pre-”cytokinesis attempt”, with cells having similar volumes to the untreated MDA-MB-231 popula-
380 tion; post-”cytokinesis attempt”, where cells were tracked following failed cytokinesis and therefore
381 continued to grow to volumes beyond those of the late stages of the untreated cell cycle; and cell
382 death, with a final cluster, composed primarily of cell debris. Furthermore, giant cell morphology
383 has been linked with docetaxel resistance, a potential cause of relapse in breast cancer patients⁹ and
384 through cluster analysis we were able to identify a potentially resistant subset of very large, treated
385 cells that could be isolated for further investigation.

386 Our chosen application demonstrated the breadth of quantification and biological insight that
387 can be made by following our workflow, with characterisation of drug response and detection of
388 potentially resistant cells just two of many potential applications for CellPhe. CellPhe offers several
389 benefits for the quantification of cell behaviour from time-lapse images. First, errors in cell segmen-
390 tation and tracking can be identified and removed, improving the quality of input for downstream
391 data analysis. This is particularly important with machine learning where automation means that
392 such errors can easily be missed, and algorithms consequently trained with poor data.

393 Second, cell behaviour is characterised over time by extracting variables from the time series
394 of various features whereas many studies explore temporal changes by collecting data at discrete

395 time points (for example, 0 and 24 hours post-treatment) and using metrics from each static image,
396 missing behavioural changes experienced by cells on a continuous level. With CellPhe, changes over
397 time in features that provide information on morphology, movement and texture are quantified not
398 just by summary statistics but by variables extracted from wavelet transformation of the time series
399 allowing changes on different scales to be identified.

400 Third, whilst most studies use a limited number of metrics, assessed individually for discrimi-
401 nation between groups,^{31,22} CellPhe provides an extensive list of novel metrics and automatically
402 determines the combination that offers greatest discrimination. The bespoke feature selection fre-
403 quently found the most discriminatory variables to be those with the ability to detect changes in
404 cell behaviour over time. Previous research in this field has focused on identification of cell types
405 from co-cultures³² for use in automated diagnosis of disease such as cancer. Analysis methods for
406 these studies are often cell line specific whereas CellPhe’s feature selection method is successful in
407 identifying discriminatory variables tailored to different experimental conditions.

408 Finally, CellPhe uses an ensemble of classifiers to predict cell status with high accuracy and
409 we show that separation scores can be used to identify the variables associated with different cell
410 subsets identified in cluster analysis to explore cell heterogeneity within a population, even when
411 subtle differences are not readily visible by eye.

412 The interactive, interpretable, high-throughput nature of CellPhe deems it suitable for all cell
413 time-lapse applications, including drug screening or prediction of disease prognosis. We provide a
414 comprehensive manual with a working example and real data to guide users through the workflow
415 step-by-step, where users can interact with each stage of the workflow and customise to suit their
416 own experiments. Here we demonstrated the abundance of information and insight that can be
417 made by following the CellPhe workflow to quantify cell behaviour from QPI images. CellPhe can
418 also be extended to other imaging modalities and work is underway to determine further variables,
419 such as fluorescence intensity, that would complement our existing list of metrics.

420

421 4 Acknowledgements

422 We would like to thank Dr. Jon Pitchford for his ongoing valuable advice and Dr. Fiona Frame
423 for her useful contributions to the project. We would also like to thank the University of York
424 Bioscience Technology Facility - Imaging and Cytometry Team for the helpful technical assistance
425 they provided throughout the project. We express gratitude to Phasefocus UK for the Livecyte and
426 CATbox systems that were used to acquire and export all time-lapse data presented here, and for
427 their technical support throughout. Furthermore, we would like to sincerely thank BBSRC for their
428 generosity in funding the project, grant number: BB/S507416/1.

429 5 Contributions

430 Conceptualisation: W.B., P.O’T., J.W., and L.W.; cell culture, pharmacology and imaging: L.W.
431 and A.L.; data analysis and validation: L.W. and J.W.; software development J.W. and L.W; super-
432 vision: J.W., W.B. and P.O’T.; writing-original draft preparation, L.W. and J.W.; writing-review
433 and editing, W.B. and P.O’T.

434

435 6 Competing interests

436 The authors declare no competing interests.

References

- 437
- 438 ¹ S. Turajlic, A. Sottoriva, T. Graham, et al. Resolving genetic heterogeneity in cancer. *Nature*
439 *Reviews Genetics*, 20:404–416, 2019.
- 440 ² S. Goldman, M. MacKay, E. Afshinnekoo, et al. The impact of heterogeneity on single-cell
441 sequencing. *Frontiers in Genetics*, 10:8, 2019.
- 442 ³ S.J. Altschuler and L.F. Wu. Cellular heterogeneity: do differences make a difference? *National*
443 *Institute of Health, Cell*, 141:559–563, 2010.
- 444 ⁴ B. Carter and K. Zhao. The epigenetic basis of cellular heterogeneity. *Nature Reviews Genetics*,
445 22:235–250, 2021.
- 446 ⁵ F. Buettner, K. Natarajan, F. Casale, et al. Computational analysis of cell-to-cell heterogeneity
447 in single-cell rna-sequencing data reveals hidden subpopulations of cells. *Nature Biotechnology*,
448 33:155–160, 2015.
- 449 ⁶ B.M. Davis, M. Salinas-Navarro, M.F. Cordeiro, et al. Characterizing microglia activation: a
450 spatial statistics approach to maximize information extraction. *Scientific Reports*, 7, 2017.
- 451 ⁷ T. Henser-Brownhill, R.J. Ju, N.K. Haass, et al. Estimation of cell cycle states of human melanoma
452 cells with quantitative phase imaging and deep learning. *IEEE 17th International Symposium on*
453 *Biomedical Imaging (ISBI)*, pages 1617–1621, 2020.
- 454 ⁸ L.A. Tashireva, M.V. Zavyalova, O.E. Savelieva, et al. Single tumor cells with epithelial-like
455 morphology are associated with breast cancer metastasis. *Frontiers in Oncology*, 10:50, 2020.
- 456 ⁹ R. Mirzayans, B. Andrais, and D. Murray. Roles of polyploid/multinucleated giant cancer cells in
457 metastasis and disease relapse following anticancer treatment. *Cancers (Basel)*, 10(4):118, 2018.
- 458 ¹⁰ A. Voulodimos, N. Doulamis, A. Doulamis, et al. Deep learning for computer vision: A brief
459 review. *Hindawi*, 2018(Computational Intelligence and Neuroscience), 2018.
- 460 ¹¹ W. Chen, W. Li, X. Dong, and J. Pei. A review of biological image analysis. *Current Bioinform-*
461 *atics*, 12, 2017.
- 462 ¹² H.E. Munim and A.A. Farag. A shape-based segmentation approach: an improved technique using
463 level sets. *Tenth IEEE International Conference on Computer Vision, Volume 1*, 2:930–935, 2005.
- 464 ¹³ Z. Wang and H. Li. Generalizing cell segmentation and quantification. *BMC Bioinformatics*, 18,
465 2017.
- 466 ¹⁴ E. Gómez-de Mariscal, C. García-López-de Haro, W. Ouyang, et al. Deepimagej: A user-friendly
467 environment to run deep learning models in imagej. *Nature Methods*, 18:1192–1195, 2021.
- 468 ¹⁵ D. Ershov, Minh-Son Phan, J. Pylvänäinen, et al. Bringing trackmate in the era of machine-
469 learning and deep-learning. *bioRxiv*, 2021.
- 470 ¹⁶ R.F. Laine, I. Arganda-Carreras, R. Henriques, et al. Avoiding a replication crisis in deep-learning-
471 based bioimage analysis. *Nature Methods*, 18:1136–1144, 2021.
- 472 ¹⁷ A. E. Carpenter, T. R. Jones, M. R. Lamprecht, et al. Cellprofiler: image analysis software for
473 identifying and quantifying cell phenotypes. *Genome Biology*, 7(10), 2006.
- 474 ¹⁸ S. Berg, D. Kutra, T. Kroeger, et al. ilastik: interactive machine learning for (bio)image analysis.
475 *Nature Methods*, 2019.
- 476 ¹⁹ J. Marrison, L. Rätty, P. Marriott, et al. Ptychography - a label-free, high contrast imaging
477 technique for live cells using quantitative phase information. *Scientific Reports*, 3(2369), 2013.

- 478 ²⁰ Y. Rivenson, Y. Zhang, H. Günaydin, et al. Phase recovery and holographic image reconstruction
479 using deep learning in neural networks. *Nature Light Sci Appl.*, 7(17141), 2018.
- 480 ²¹ Y. Park, C Depeursinge, and G. Popescu. Quantitative phase imaging in biomedicine. *Nature*
481 *Photon*, 12:578–589, 2018.
- 482 ²² R. Suman, G. Smith, and K. E.A. Hazel et al. Label free imaging to study phenotypic behavioural
483 traits of cells in complex co-cultures. *Scientific Reports*, 6(1):22–32, 2016.
- 484 ²³ S. Comsa, A.M. Cimpean, and M. Raica. The story of mcf-7 breast cancer cell line: 40 years of
485 experience in research. *BioMed Central Cancer, Anticancer Research*, 35:3147–3154, 2015.
- 486 ²⁴ L. Oakden-Rayner, J. Dunnmon, G. Carneiro, et al. Hidden stratification causes clinically mean-
487 ingful failures in machine learning for medical imaging. *arXiv*, 2019.
- 488 ²⁵ H. Hernandez-Vargas, J. Palacios, and G. Moreno-Bueno. Molecular profiling of docetaxel cytotox-
489 icity in breast cancer cells: uncoupling of aberrant mitosis and apoptosis. *Oncogene*, 26:2902–2913,
490 2007.
- 491 ²⁶ J. Liu, L. Qu, L. Meng, et al. Topoisomerase inhibitors promote cancer cell motility via ros-
492 mediated activation of jak2-stat1-cxcl1 pathway. *Journal of Experimental and Clinical Cancer*
493 *Research*, 38:370, 2019.
- 494 ²⁷ C.L. Liu, M.J. Chen, J.C. Lin, et al. Migration and invasion of breast cancer cells through the
495 upregulation of the rhoa/mlc pathway. *J breast cancer*. *Journal of Breast Cancer*, 22:185–195,
496 2019.
- 497 ²⁸ N.S. Sohoni, J.A. Dunnmon, G. Angus, et al. No subclass left behind: Fine-grained robustness in
498 coarse-grained classification problems. *CoRR*, 2020.
- 499 ²⁹ R. Wang, C. Jin, and X. Hu. Evidence of drug-response heterogeneity rapidly generated from a
500 single cancer cell. *Oncotarget*, 8:25, 2017.
- 501 ³⁰ X. Wang, Z. Chen, A.K. Mishra, et al. Chemotherapy-induced differential cell cycle arrest in b-cell
502 lymphomas affects their sensitivity to weel inhibition. *Haematologica.*, 103(3):466–476, 2018.
- 503 ³¹ F. M. Frame, A. R. Noble, S. Klein, et al. Tumor heterogeneity and therapy resistance - impli-
504 cations for future treatments of prostate cancer. *Journal of Cancer Metastasis and Treatment*,
505 3:302–314, 2017.
- 506 ³² Y. Ozaki, H. Yamada, H. Kikuchi, et al. Label-free classification of cells based on supervised
507 machine learning of subcellular structures. *PLoS One*, 14(1), 2019.
- 508 ³³ M. Yang, D.J. Kozminski, L. Wold, et al. Therapeutic potential for phenytoin: targeting nav1.5
509 sodium channels to reduce migration and invasion in metastatic breast cancer. *Breast Cancer*
510 *Research and Treatment*, 134(2):603–615, 2012.
- 511 ³⁴ C. Uphoff, S. Gignac, and H. Drexler. Mycoplasma contamination in human leukemia cell lines.
512 i. comparison of various detection methods. *Journal of Immunological Methods*, 149:43–53, 1992.
- 513 ³⁵ L.K. Soh and C. Tsatsoulis. Texture analysis of sar sea ice imagery using gray level co-occurrence
514 matrices. *IEEE Transactions on geoscience and remote sensing*, 37(2):780–795, 1999.
- 515 ³⁶ S. Mallat. A theory for multiresolution signal decomposition: the wavelet representation. *IEEE*
516 *Transactions on Pattern Analysis & Machine Intelligence*, 7:674–693, 1989.
- 517 ³⁷ R.M. Haralick, K. Shanmugam, and I. Dinstein. Textural features for image classification. *IEEE*
518 *Transactions on systems, man, and cybernetics*, 6:610–621, 1973.
- 519 ³⁸ A. Baddeley and R. Turner. spatstat: An R package for analyzing spatial point patterns. *Journal*
520 *of Statistical Software*, 12(6):1–42, 2005.

- 521 ³⁹ A. Haar. Zur theorie der orthogonalen funktionensysteme. *Mathematische Annalen*, 69(3):331–
522 371, 1910.
- 523 ⁴⁰ N.V. Chawla, K.W. Bowyer, L.O. Hall, et al. Smote: Synthetic minority over-sampling technique.
524 *Journal of Artificial Intelligence Research*, 16:321–357, 2002.
- 525 ⁴¹ R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for
526 Statistical Computing, Vienna, Austria, 2019.
- 527 ⁴² W. N. Venables and B. D. Ripley. *Modern Applied Statistics with S*. Springer, New York, fourth
528 edition, 2002.
- 529 ⁴³ D. Meyer, E. Dimitriadou, et al. *e1071: Misc Functions of the Department of Statistics, Probability
530 Theory Group (Formerly: E1071)*, TU Wien, 2019. R package version 1.7-3.
- 531 ⁴⁴ A. Liaw and M. Wiener. Classification and regression by randomforest. *R News*, 2(3):18–22, 2002.
- 532 ⁴⁵ A. Kassambara and F. Mundt. *factoextra: Extract and Visualize the Results of Multivariate Data
533 Analyses*, 2020. R package version 1.0.7.

534 Methods

535 **Cell Culture.** MDA-MB-231 cells and MCF-7 cells were cultured separately in Dulbecco’s modi-
536 fied eagle medium supplemented with 5% fetal bovine serum and 4mM L-glutamine.³³ Fetal bovine
537 serum was filtered using a 0.22 μ m syringe filter prior to use to reduce artefacts when imaging. Cells
538 were incubated at 37°C in plastic filter-cap T-25 flasks and were split at a 1:6 ratio when passaged.
539 No antibiotics were added to cell culture medium. Cells were confirmed to be mycoplasma-free by
540 4’,6-diamidino-2-phenylindole (DAPI) method.³⁴ To image the following day, cells were counted and
541 then seeded in a Corning Costar plastic, flat bottom 24-well plate. Cells were seeded at a density of
542 8000 cells per well with a final volume of 500 μ L in each of the 24 wells.

543 **Pharmacology.** Docetaxel (Cayman Chemical Company) was prepared as 5mg/mL of DMSO and
544 doxorubicin (AdooQ Bioscience) as 25mg/mL of DMSO, both were then frozen into aliquots. Once
545 thawed, docetaxel and doxorubicin stock solutions were diluted in culture medium to give final
546 working concentrations. Docetaxel dose response analysis for both MDA-MB-231 and MCF-7 cells
547 involved imaging eight wells treated with the following concentrations of docetaxel: 0nM, 1nM, 3nM,
548 10nM, 30nM, 100nM, 300nM, 1 μ M, with additional concentrations 3 μ M, 10 μ M and 30 μ M imaged for
549 MDA-MB-231 cells. Doxorubicin dose response analysis for MDA-MB-231 cells involved imaging
550 eight wells treated with the following concentrations of doxorubicin: 0nM, 10nM, 30nM, 100nM,
551 300nM, 1 μ M, 3 μ M, 10 μ M.

552 Medium was removed from wells selected to receive treatment 30 minutes prior to image acqui-
553 sition, and 500 μ L of desired drug concentration was added to each well. Control wells received a
554 medium change and were treated with DMSO vehicle on the day of imaging to maintain consistent
555 DMSO concentration throughout.

556 **Image Acquisition and Exportation.** Cells were placed onto the Phasefocus Livecyte 2 (Phase-
557 focus Limited, Sheffield, UK) to incubate for 30 minutes prior to image acquisition to allow for
558 temperature equilibration. One 500 μ m x 500 μ m field of view per well was imaged to capture as
559 many cells, and therefore data observations, as possible. Selected wells were imaged in parallel for
560 48 hours at 20x magnification with 6 minute intervals between frames, resulting in full time-lapses
561 of 481 frames per imaged well.

562 Phasefocus’ Cell Analysis Toolbox[®] software was utilised for cell segmentation, cell tracking
563 and data exportation. Segmentation thresholds were optimised for a range of image processing tech-
564 niques such as rolling ball algorithm to remove background noise, image smoothing for cell edge
565 detection and local pixel maxima detection to identify seed points for final consolidation.

568 The Phasefocus software outputs a feature table for each imaged well. This table consists of
569 variables that describe the size and orientation of each cell on each frame as well as variables that
570 describe the cell's movement up to the current frame. As some of these variables are highly cor-
571 related, for example dry mass and volume, we only include a subset of Phasefocus' features in our
572 analyses to reduce the number of redundant variables. The first nine features in **Supplementary**
573 **table S1** are those that were retained.

574

575 **Implementation of CellPhe.** The CellPhe toolkit for the calculation of variables from time series
576 data and bespoke feature extraction is implemented in the C programming language and is freely
577 available from the corresponding author. Further analysis, including classification and cluster anal-
578 ysis, were carried out using freely available R packages as detailed and R scripts can be supplied
579 upon request.

580

581 *New Features.* Using the Regions of Interest (ROIs) produced by the Phasefocus software to identify
582 each cell's boundary pixels, a range of additional morphological and texture features were extracted.
583 In addition to shape descriptors calculated from the cell boundaries, a filling algorithm was used
584 to determine the interior pixels from which texture and spatial features were extracted. The local
585 density was also calculated as the proportion of the area, A , around the cell containing pixels from
586 other cells. Here A is the annulus around the cell from its radius r to $r + a$ and a is the average
587 radius of cells in the data set. A complete list of features together with their definitions is provided
588 in **Supplementary table S1**.

589

590 *Texture descriptors.* Gray-level co-occurrence matrices (GLCMs) are widely used in texture analysis
591 to investigate spatial relationships between the pixels with similar intensities.³⁵ Here, rather than
592 consider the positions of pixels within a cell, we calculated GLCMs between the image of the cell
593 at different resolutions to differentiate textures that are sharp and would be lost at lower resolu-
594 tion from those that are smooth and would remain. This was achieved by performing a two-level
595 2-D wavelet transform³⁶ on the pixels within the axis-aligned minimum rectangle containing a cell.
596 GLCMs were then calculated between the original interior pixels and the corresponding values from
597 the first and second levels of the transform as well as between the two sets of transformed pixels.
598 Statistics first described by Haralick³⁷ were then calculated from each GLCM as follows:

599

$$\begin{aligned}
\text{Energy} &= \sum_{i,j=1}^N M_{ij}^2 \\
\text{Contrast} &= \sum_{n=0}^{N-1} \sum_{i,j=1}^N M_{ij} \\
\text{Homogeneity} &= \sum_{i,j=1}^N \frac{M_{ij}}{1 + (i + j)^2} \\
\text{Correlation} &= \sum_{i,j=1}^N M_{ij} \left(\frac{(i - \bar{i})(j - \bar{j})}{\sigma_i \sigma_j} \right) \\
\text{Entropy} &= - \sum_{i,j=1}^N M_{ij} \log(M_{ij})
\end{aligned}$$

600 where M_{ij} is the (i, j) th entry in the normalised GLCM and $\bar{i}, \bar{j}, \sigma_i$ and σ_j are the means and stan-
601 dard deviations of the marginal densities $M_i = \sum_j^N M_{ij}$ and $M_j = \sum_i^N M_{ij}$ respectively. With three
602 co-occurrence matrices, this added 15 texture features.

603

604 *Spatial distribution descriptors.* We calculated spatial distribution descriptors to quantify the uni-
605 formity or clustering of cell interior pixels with intensities above 9 different thresholds. For IQn,

interior pixels with intensities greater than or equal to the $(n \times 10)$ th quantile are modelled as a spatial point process using the `ppp` function in the R *spatstat* package.³⁸ The cell boundary coordinates form the boundary window of each spatial point process and we make use of Ripley’s reduced second moment K function to compare the empirical distribution of pixels with a theoretical Poisson distribution. Each spatial distribution descriptor is then defined as follows:

$$IQn = K_{emp}(maxr) - K_{theo}(maxr)$$

where $maxr$ is one quarter of the smallest side of the enclosing rectangle within the cell boundary. The K functions, K_{emp} and K_{theo} , are evaluated for our empirical data and for a theoretical Poisson distribution respectively, such that:

$$K_{emp}(r) = \frac{a}{(n \times (n - 1))} \times \sum_{i,j} I(d(i, j) \leq r),$$

and

$$K_{theo}(r) = \pi r^2,$$

where a is the area of the cell boundary, n is the number of pixels being considered, and the sum is taken over all ordered pairs of pixels i and j . $I(d(i, j) \leq r)$ is an indicator variable that equals 1 if the distance $d(i, j)$ between pixels i and j is $\leq r$.

A negative IQn value indicates that the pixels are uniformly distributed, whereas positive indicates clustering of points, with large absolute values of IQn indicative of greater uniformity or clustering.

Summarising Cell Time Series. Cell tracking provides a time series for each of the 47 features extracted for a cell. The length of the time series depends on how many frames the cell has been tracked for and so differs between cells. In order to apply pattern recognition methods, we extracted a fixed number of characteristic variables for each cell from the time series for each feature. Statistical measures (mean, standard deviation and skewness) summarise time series of varying length, but may not be representative of changes throughout the time series. Therefore, in addition to summary statistics, we calculated variables inspired by elevation profiles in walking guides, that is, the sum of any increases between consecutive frames (total ascent), the sum of any decreases (total descent) and the maximum value of the time series (maximum altitude gain). These “elevation” variables were calculated for different levels of the wavelet transform of the time series to allow changes at different scales to be considered. The wavelet transform decomposes a time series to give a lower resolution approximation together with different levels of detail that need to be added to the approximation to restore the original time series. Using the Haar wavelet basis³⁹ with the multiresolution analysis of Mallat³⁶ allows increases and decreases in the values of the variables to be determined over different time scales.

Occasionally the automated cell tracking misses a frame or even several frames, for example when a cell temporarily leaves the field of view. To prevent jumps in the time series, we interpolated values for the missing frames, although these values were not used to calculate statistics. After interpolation, the three elevation variables were calculated from the original time series and three wavelet levels which, together with the summary statistics, provided 15 variables for each feature **Supplementary table S2**. This would have given $47 \times 15 = 705$ variables in total, but, as one feature, the tracklength or total distance travelled up to the current frame, is monotonically increasing, the total descent is always zero and therefore the 4 descent variables were not used.

One further variable was introduced to summarise cell movement as the area of the minimal bounding box around a cell’s full trajectory. This area will be large for migratory cells and small for cells whose movement remains local for the duration of the time series. If, within a cell’s trajectory, $minX$ and $minY$ are the minimal X and Y positions respectively with $maxX$ and $maxY$ the corresponding maximal positions, then the trajectory area is defined as

$$\text{trajectory area} = (maxX - minX) \times (maxY - minY).$$

656 Thus, a total of 702 characteristic variables were available for analysis and classification.

657

658 *Segmentation Error Removal.* To improve characterisation of cellular phenotype, we only included
659 cells that were tracked for at least 50 frames in our analyses. Whilst the majority of these cells
660 were correctly tracked, others had segmentation errors, with confusion between neighbouring cells,
661 missing parts of a cell or multiple cells included.

662 In order to increase the reliability of our results, we developed a classification process to identify
663 and remove such cells prior to further analysis. Cells (both treated and untreated) were classified by
664 eye to provide a training data set. Due to class imbalance, with the number of segmentation errors
665 far less than the number of correct segmentations, the Synthetic Minority Oversampling Technique
666 (SMOTE)⁴⁰ was performed using the *smotefamily* package in R, with the number of neighbours K
667 set to 3, to double the number of instances representing segmentation errors.

668 The resulting data set with all 702 variables was used to train a set of 50 decision trees using
669 the *tree* package in R with default parameters. For each tree, the observations from cells with seg-
670 mentation errors were used together with the same number of observations randomly selected from
671 the correctly segmented cells to further address class imbalance. For each cell, a voting procedure
672 was used to provide a classification from the predictions of the 50 decision trees. To minimise the
673 number of correctly tracked cells being falsely classified as segmentation errors, this class was only
674 assigned when it received at least 70% of the votes (i.e. 35). To add further stringency, the training
675 of 50 decision trees was repeated ten times and a cell only given a final classification of segmentation
676 error if predicted this label in at least five of the ten runs. MDA-MB-231 cells that were not used
677 for training formed an independent test set. All cells either manually labelled as segmentation error
678 or predicted as such were excluded from further analyses.

679

680 *Classification of Untreated and Treated Cells.* After removing segmentation errors, the remaining
681 data were used to form training and test sets for the classification of untreated and treated cells.
682 Training sets were balanced prior to classifier training to mitigate bias and data from cells in the
683 independent test sets were never used during training.

684 A separate classifier was trained for each cell line - treatment combination, as shown in **Table**
685 **3** and feature selection performed to determine the most appropriate variables in each case. Each
686 variable was assessed using the group separation, $S = V_B/V_W$, where V_B is the between-group
687 variance:

$$688 \quad V_B = \frac{n_1(\bar{x}_1 - \bar{x})^2 + n_2(\bar{x}_2 - \bar{x})^2}{(n_1 + n_2 - 2)}$$

689 and V_W is the within-group variance:

$$690 \quad V_W = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{(n_1 + n_2 - 2)}.$$

691 Here n_1 and n_2 denote the sample size of group 1 and group 2 respectively, \bar{x}_1 and \bar{x}_2 are the sample
692 means, \bar{x} the overall mean, and s_1^2 and s_2^2 are the sample variances. The most discriminatory vari-
693 ables were chosen for a particular data set by assessing the classification error on the training data to
694 optimise the threshold on separation. Starting with a threshold of zero, the n^{th} separation threshold
695 was maximised such that the classification error rate did not increase by more than 2% from that
696 obtained for the $(n - 1)^{\text{th}}$ threshold. The aim here was to reduce the risk of overfitting by only retain-
697 ing variables achieving greater than or equal to this threshold for the next stage of classifier training.

698

699 Data were scaled to prevent large variables dominating the analysis and ensemble classification
700 used to take advantage of different classifier properties. The predictions from three classification
701 algorithms, Linear Discriminant Analysis (LDA), Random Forest (RF) and Support Vector Ma-
702 chine (SVM) with radial basis kernel were combined using the majority vote. Model performance
703 was evaluated by classification accuracy, taking into account the number of false positives and false
704 negatives. All classification was performed in RStudio⁴¹ using open-source packages. LDA was per-
705 formed using the *lda* function from the *MASS* library,⁴² SVM classification used the *svm* function

706 from the *e1071* package⁴³ with a radial basis kernel and the *randomForest* package⁴⁴ was used to
 707 train random forest classifiers with 200 trees and 5 features randomly sampled as candidates at each
 708 split.

709

710 *Cluster analysis.* Both hierarchical clustering and *k*-means clustering were used to investigate sub-
 711 groups within single-class data sets (i.e. treated and untreated cells separately). Data were scaled
 712 prior to clustering and analyses performed in R. Hierarchical clustering was implemented with the
 713 *factoextra* package⁴⁵ using the *hcut* function to cut the dendrogram into *k* clusters. Agglomerative
 714 nesting (AGNES) was used with Ward’s minimum variance as the agglomeration method and the
 715 Euclidean distance metric to quantify similarity between cells. *k*-means clustering was performed
 716 using the R *stats* package, with the number of random initial configurations set to 50. The number of
 717 clusters *k* was chosen to obtain clusters with meaningful interpretation. Similarities and differences
 718 between clusters were identified through evaluation of separation scores to determine discriminatory
 719 features, as well as through observation of cells within each cluster by eye.

720

721 *Data.* Three data sets are used to demonstrate our pipeline for the classification of untreated and
 722 treated cells. For brevity we use abbreviations throughout to refer to each data set, for example
 723 "231Docetaxel" is a data set consisting of MDA-MB-231 cells, both untreated and treated with
 724 30 μ M docetaxel. This is the main data set used to develop the methods, with a training data set
 725 compiled from 6 experiments performed on different days and an independent test data set compiled
 726 from a further 3 experiments, also performed on separate days and by a different individual.

727 We validate our methods using two further datasets, the 231Doxorubicin and MCF7Docetaxel
 728 data sets, details of which are given in **Table 3**. This table also includes details of the number of
 729 cells within each training and test set. We show that the classification pipeline can be successfully
 730 reproduced using fewer experimental repeats for the 231Doxorubicin and MCF7Docetaxel data sets.
 731 The 231Doxorubicin training set consists of data from one experiment with a further, independent
 732 experiment performed on a separate day used as a test set. Training and test sets for MCF7Docetaxel
 733 are from the same experiment, with random sampling used to produce independent training and test
 734 sets. Each training data set contains a balanced number of untreated and treated cells, treated with
 735 a single drug concentration. We selected 30 μ M docetaxel and 1 μ M doxorubicin for the experiments
 736 with MDA-MB-231 cells as the optimal doses with which to induce changes in cell morphology and
 737 migration without inducing cell death. However, a lower concentration (1 μ M) of docetaxel was used
 738 for MCF-7 cells as we found that this induced similar morphological and dynamical changes to those
 739 induced by higher concentrations but with reduced cell death (**Table 3**).

740

Data set	Cell line	Treatment	Training set	Test set
231Docetaxel	MDA-MB-231	30 μ M Docetaxel	Untreated: 646 Treated: 600	Untreated: 939 Treated: 314
231Doxorubicin	MDA-MB-231	1 μ M Doxorubicin	Untreated: 216 Treated: 216	Untreated: 185 Treated: 59
MCF7Docetaxel	MCF-7	1 μ M Docetaxel	Untreated: 130 Treated: 130	Untreated: 304 Treated: 27

741 **Table 3:** The three data sets used in this study with the number of cells in training and test sets used for untreated vs treated classification.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [CellPh>manual.pdf](#)
- [CellPhesupplementary.pdf](#)
- [WilsonCodeFlat.pdf](#)
- [WilsonEPCflat.pdf](#)
- [WilsonSupplInfo.pdf](#)