UNIVERSITY of York

This is a repository copy of *The CellPhe toolkit for cell phenotyping using time-lapse imaging and pattern recognition*.

White Rose Research Online URL for this paper: <u>https://eprints.whiterose.ac.uk/197666/</u>

Version: Accepted Version

Article:

Wiggins, Laura, Lord, Alice, Murphy, Killian L et al. (4 more authors) (2023) The CellPhe toolkit for cell phenotyping using time-lapse imaging and pattern recognition. Nature Communications. 1854. ISSN 2041-1723

https://doi.org/10.1038/s41467-023-37447-3

Reuse

Items deposited in White Rose Research Online are protected by copyright, with all rights reserved unless indicated otherwise. They may be downloaded and/or printed for private study, or other acts as permitted by national copyright laws. The publisher or other rights holders may allow further reproduction and re-use of the full text version. This is indicated by the licence information on the White Rose Research Online record for the item.

Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing eprints@whiterose.ac.uk including the URL of the record and the reason for the withdrawal request.



eprints@whiterose.ac.uk https://eprints.whiterose.ac.uk/

2	CellPhe: a toolkit for cell phenotyping using time-lapse imaging and
3	pattern recognition
4	Laura Wiggins ^{1,2} , Alice Lord ² , Killian L. Murphy ³ , Stuart E. Lacy ³ , Peter J. O'Toole ^{1,2} , William
5	J. Brackenbury ^{1,2} and Julie Wilson ^{4*}
6	¹ York Biomedical Research Institute, University of York, York, UK
7	² Department of Biology, University of York, York, UK
8	³ Wolfson Atmospheric Chemistry Laboratories, University of York, York, UK
9	⁴ Department of Mathematics, University of York, York, UK
10	* Corresponding author
11	email: julie.wilson@york.ac.uk
12	Abstract. With phenotypic heterogeneity in whole cell populations widely recognised, the de-
13	mand for quantitative and temporal analysis approaches to characterise single cell morphology and
14	dynamics has increased. We present CellPhe, a pattern recognition toolkit for the unbiased char-
15	acterisation of cellular phenotypes within time-lapse videos. CellPhe imports tracking information
16	from multiple segmentation and tracking algorithms to provide automated cell phenotyping from
17	different imaging modalities, including fluorescence. To maximise data quality for downstream anal-
18	ysis, our toolkit includes automated recognition and removal of erroneous cell boundaries induced
19	by inaccurate tracking and segmentation. We provide an extensive list of features extracted from
20	individual cell time series, with custom feature selection to identify variables that provide greatest

cellular phenotype and clustering algorithms for the characterisation of heterogeneous subsets, we validate and prove adaptability using different cell types and experimental conditions. 23

1 Introduction 24

1

20

21

22

Heterogeneity in whole cell populations is a long-standing area of $interest^{1,2,3}$ and previous 25 studies have identified cell-to-cell phenotypic and genotypic diversity even within clonally derived 26 populations.⁴ The emergence of methods such as single-cell RNA sequencing has enabled charac-27 terisation of subsets within a population from gene expression profiles, 5 yet these methods involve 28 collection of data at discrete time points, missing the subtle temporal changes in gene expression 29 on a continuous scale. Such methods exclude information on single-cell morphology and dynamics, 30 yet cellular phenotype plays a crucial role in determining cell function,^{6,7} disease progression,⁸ and 31 response to treatment.⁹ There remains a demand for quantitative and temporal analysis approaches 32 to describe the subtleties of single-cell heterogeneity and the complexities of cell behaviour. 33

discrimination for the analysis in question. Using ensemble classification for accurate prediction of

Modern microscopy advancements facilitate the ability to produce information-rich images of 34 cells and tissue, at high-throughput and of high quality. Temporal changes in cell behaviour can 35 be observed through time-lapse imaging and features describing the cells' behaviour over time can 36 be extracted for analysis. However, the task of identifying individual cells and following them 37 over time is an ongoing computer vision challenge.^{10,11} Initial processing requires segmentation, 38 the detection of cells as regions of interest (ROIs) distinguished from background, and tracking, 39 with each cell given a unique identifier that is retained over subsequent frames. Recent work using 40 the similarity between cell metrics on consecutive frames highlighted the importance of accurate 41 tracking to follow cell lineage.¹² Imaging artefacts vary between experiments and issues such as 42 background noise, inhomogeneity of cell size and overlapping cells are still challenges for biomedical 43 research.¹³ Reliable cell segmentation protocols are non-deterministic and experiment-specific¹⁴ but 44 user-friendly software systems that use machine learning algorithms are emerging to provide ob-45 jective, high-throughput cell segmentation and tracking.^{15,16} Recent developments to TrackMate¹⁷ 46 allow the results of various segmentation software to be integrated with flexible tracking algorithms 47 and provide visualisation tools to assess both segmentation and cell tracks. Although the time series 48

for certain cell properties, such as cell area and circularity, can be displayed, the extraction and analysis of descriptive time series is not within the scope of the TrackMate software. Comparison of the tracked cells behaviour is challenging as cells are tracked for different numbers of frames with frames missing where cells leave the field of view. This has meant that analysis of any extracted features has been limited to visualisation. CellPhe interpolates the time series and then calculates a fixed number of variables that characterise each feature's time series- the features of features!

Here we present CellPhe, a pattern recognition toolkit that uses the output of segmentation 55 and tracking software to provide an extensive list of features that characterise changes in the cells' 56 appearance and behaviour over time. Customised feature selection allows the most discriminatory 57 variables for a particular objective to be identified. These extracted variables quantify cell morphol-58 ogy, texture and dynamics and describe temporal changes and can be used to reliably characterise 59 and classify individual cells as well as cell populations. To ensure precise quantification of cell mor-60 phology and motility, and to monitor major cellular events such as mitosis and apoptosis, it is vital 61 that instances of erroneous segmentation and tracking are removed from data sets prior to down-62 stream analysis methods.¹⁸ Manual removal of such errors is heavily labour-intensive, particularly 63 when time-lapses take place over several days. To maximise data quality for downstream analysis, 64 CellPhe includes the recognition and removal of erroneous cell boundaries induced by inaccurate 65 segmentation and tracking. We demonstrate the use of ensemble classification for accurate predic-66 tion of cellular phenotype and clustering algorithms for identification of heterogeneous subsets. 67

We exemplify CellPhe by characterising the behaviour of untreated and chemotherapy treated 68 breast cancer cells from ptychographic time-lapse videos. Quantitative phase images (QPI)^{19, 20, 21} 69 avoid any fluorescence-induced perturbation of the cells but segmentation accuracy can be affected by 70 reduced differences in intensity between cells and background in comparison to fluorescent labelling. 71 We show that our methods successfully recognise and remove a population of erroneously segmented 72 cells, improving data set quality. Morphological and dynamical changes induced by chemothera-73 peutics, particularly at low drug concentration, are often more subtle than those that discriminate 74 distinct cell types and we demonstrate the ability of CellPhe to automatically identify time series 75 differences induced by chemotherapy treatment, with the chosen variables proving statistically sig-76 nificant even when not observable by eye. 77

The complexities of heterogeneous drug response and the problem of drug resistance further mo-78 tivate our chosen application. The ability to identify discriminatory features between treated and 79 untreated cells can allow automated detection of "non-conforming" cells such as those that possess 80 cellular drug resistance. Further investigation of such features could elucidate the underlying bio-81 logical mechanisms responsible for chemotherapy resistance and cancer recurrence. We validate the 82 adaptability of CellPhe with both a different cell type and a different drug treatment and show that 83 variables are selected according to experimental conditions, tailored to properties of the cell type 84 and drug mechanism of action. 85

⁸⁶ CellPhe is available on GitHub as an R package with a user-friendly interactive GUI that al ⁸⁷ lows completely unbiased cell phenotyping using time-lapse data from fluorescence imaging as well
 ⁸⁸ as ptychography. A working example guides the user through the complete workflow and a video
 ⁸⁹ demonstrating the GUI is also provided.

90 2 Results

91 Overview of CellPhe

CellPhe is a toolkit for the characterisation and classification of cellular phenotypes from time-92 lapse videos, a diagrammatic summary of CellPhe is provided in Figure 1. Experimental design 93 is determined by the user prior to image acquisition where seeded cell types and pharmacology are 94 specific to the user's own analysis. Example uses are discrimination of cell types (e.g. neurons vs. 95 astrocytes), characterisation of disease (e.g. healthy vs. cancer), or assessment of drug response 96 97 (e.g. untreated vs. treated). The user can then time-lapse image cells for the desired amount of time, using an imaging modality of their choice. Once images are acquired and segmentation and tracking 98 of cells are complete, cell boundary coordinates are exported and used for calculation of an extensive 99

list of morphology and texture features. These together with dynamical features and extracted time 100 series variables are used to aid removal of erroneous segmentation by recognition of error-induced 101 interruption to cell time series. Once all predicted segmentation errors have been removed from 102 data sets, feature selection is performed and only features providing separation above an optimised 103 threshold are retained. This identifies a list of most discriminatory features and allows the user 104 to explore biological interpretation of these findings. The extracted data matrices are then used 105 as input for ensemble classification, where the phenotype of new cells can be accurately predicted. 106 Furthermore, clustering algorithms can be used to identify heterogeneous subsets of cells within the 107 user's data, both inter- and intra-class. 108

The remaining results exemplify the use of CellPhe with a biological application, characterisation and classification of chemotherapeutic drug response. We look at each of the CellPhe stages in detail (segmentation error removal, feature selection, ensemble classification and cluster analysis) and demonstrate that each step provides interpretable, biologically relevant results to answer experiment

¹¹³ specific questions and aid further research.



114

115

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.

¹¹⁶ CellPhe application: characterising chemotherapeutic drug response

The 231Docetaxel data set, obtained from multiple experiments involving MDA-MB-231 cells, 117 both untreated and treated with 30μ M docetaxel, is the main data set used to demonstrate our 118 method. We show that the same analysis pipeline can be applied to other data sets by considering 119 both a different cell line, MCF-7, in the MCF7Docetaxel data set, and a different drug, doxorubicin, 120 with the 231Doxorubicin data set. In each case, we remove segmentation errors, as described in 121 Section 2.5, before using feature selection (Section 2.6) to identify discriminatory variables tailored 122 to the particular data set. We show that different variables are chosen depending on the inherent 123 nature of the cell line and the effect of the drug in question. Using these features in classification 124

¹²⁵ algorithms, we characterise and compare the behaviour over time of untreated and treated cells.

126 127

Segmentation Error Removal

We improve the quality of our data sets prior to untreated vs. treated cell classification by 128 automating detection of segmentation errors and optimising the exclusion criteria of predicted errors. 129 Comparison of time series for cells with and without segmentation errors showed many of our 130 features to be sensitive to such errors, motivating the need to remove these cells prior to treatment 131 classification. Size metrics, such as volume, were particularly affected by segmentation errors as 132 under- or over-segmentation could result in halving or doubling of cell volume respectively (Figure 133 **2a**). Such noticeable disruption to the time series of several features suggested that reliable detection 134 of segmentation errors would be possible. 135

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

A total of 223 of the 1478 cells in the 231Docetaxel test set were predicted to be segmentation 143 errors. Of these, 217 were confirmed by eye to be true segmentation errors, most of which were due 144 to under- or over-segmentation throughout their time series. Other segmentation issues observed 145 included background pickup, cells swapping cell ID, and cells repeatedly entering and exiting the 146 field of view, all of which result in problem time series (Figure 2b). Of the remaining six cells that 147 were misclassified as segmentation errors, one was a large cell and the other five were cells tracked 148 before, during and after attempted mitosis. Further investigation showed that removal of these cells 149 did not exclude an important subset from the data. 150

This classifier was also used to identify a further 78 segmentation errors from the 955 cells in 151 the 231Doxorubicin data set, all 78 were confirmed by eye to be true segmentation errors (Table 1). 152 It was necessary to train a new classifier for MCF-7 segmentation error detection due to differences 153 between the cell lines. In this case 308 correctly segmented cells and 192 segmentation errors were 154 identified by eye. After applying SMOTE to double the number of segmentation error observations, 155 a classifier was trained with the resulting 692 observations as described in section 2.5. 188 cells in 156 the MCF7Docetaxel data set (848 cells in total) were classified as segmentation errors. 185 of these 157 cells were confirmed by eye to be true segmentation errors, the remaining three were large cells or 158 cells tracked before, during and after attempted mitosis. 159

Data set	TP	FP
231 Docetaxel (1478)	217	6
231Doxorubicin (955)	78	0
MCF7Docetaxel (848)	185	3

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.



161

162

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) Box and whisker plots of features that are significant for identifying segmentation errors in the 231Docetaxel training set (****: p < 0.0001). The median value is shown by the line within the box representing the interquartile range (IQR) and the whiskers extend to the most extreme data points. (d) A representative 231Docetaxel trained decision tree, demonstrating how size, shape, texture and density are used in combination to make classifications.

As decision trees are used in the identification of segmentation errors, our feature selection is not required. However we still calculated separation scores for the MDA-MB-231 training data to investigate the effect of such errors. As might be expected, volume was most affected, with segmentation errors resulting in larger standard deviation, ascent and maximum value. Other features with high separation scores included area as well as spatial distribution descriptors with the highest thresholds, features that detect the clustering of high intensity pixels, characteristic of cell overlap and over-segmentation (**Figure 2c**). Analysis of the trained decision trees showed that a combination of size, shape, texture and density variables frequently formed the most important features for
detecting segmentation errors with MDA-MB-231 cells, see Figure 2d for an example.

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

177 Feature Selection

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

Differences in the appearance of MDA-MB-231 cells induced by docetaxel treatment were ob-183 served by eye from cell timelapses. Untreated cells displayed a spindle-shaped morphology (a circular 184 cross-section with tapering at both ends), with contractions and protrusions facilitating migration. 185 Cells that received treatment were generally dense and spherical, and increased in size following 186 a failed attempt at cytokinesis (Figure 3a). Discriminatory features identified by calculation of 187 separation scores were consistent with differences observed by eye, the 100 variables that achieved 188 greatest separation are shown in Figure 3b. Texture, shape and size variables provided great-189 est discrimination of untreated and treated cells. Untreated cells experienced increased elongation 190 throughout the time-lapse and displayed irregular, spindle-shaped morphology in comparison to the 191 generally spherical appearance of treated cells. Furthermore, separation scores highlighted differ-192 ences in the texture of cells, with intensity quantile metrics characterising changes in granularity of 193 cells induced by drug treatment. 194

Principal Component Analysis (PCA) demonstrated that the main variance within the data arises 195 due to class differences, with separation of classes observed across PC1 which explains 66% of the 196 total variance (Figure 3c). The dispersion of points within the scores plot illustrates heterogeneity 197 of cells both inter- and intra-class. The non-conformity of some cells, for example treated cells be-198 having as untreated cells, is demonstrated by points clustering within the opposite class. Analysis 199 of PCA loadings highlighted increased ascent, descent and standard deviation for untreated cells, 200 as can be observed from the PCA biplot in **Figure 3d**. Although descent variables appear to have 201 opposite loadings to all other variables, in fact, this is only due to their negative values. As the ma-202 jority of untreated cells had negative PC1 scores we deduced that greater standard deviation, ascent 203 and descent of features for untreated cells indicates that these cells experience increased fluctuation 204 throughout their time series. As treated cells mainly had positive PC1 scores, they experience less 205 fluctuation throughout their time series and instead display greater stability. Identified differences 206 in feature time series are visualised in Figure 3d. 207



209

210

Figure 3: a) Images taken from cell timelapses of i. untreated MDA-MB-231 cells and ii. 30μ M docetaxel treated MDA-MB-231 cells. Scale bar = 200μ m. Increased cell count at 49h post-treatment demonstrates healthy proliferation of untreated cells. Static cell count at 49h for treated cells is a result of cell cycle arrest and failed cytokinesis, leading to enlarged cell phenotype. b) Features with the top 100 highest separation scores, colour-coded according to feature type. Texture, shape and size features provide greatest separation. c) PCA scores plot with points colour-coded 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. Only features with the 100 highest separation scores were included in PCA. d) i. PCA biplot demonstrating how features with the 100 highest separation scores work in combination to discriminate between untreated and 30μ M docetaxel treated MDA-MB-231 cells. Greater ascent and descent can be observed for untreated cells, indicating greater activity across a range of features for untreated cells. ii. Representative feature time series plots for untreated and 30μ 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.

²¹¹ We assessed the adaptability of our feature selection method by calculating separation scores ²¹² for both a different cell line and a different treatment, using PCA to evaluate the main sources of ²¹³ variance. We compared MCF-7 cells treated with 1 μ M docetaxel with untreated MCF-7 cells, and ²¹⁴ MDA-MB-231 cells that were treated with 1 μ M doxorubicin with untreated MDA-MB-231 cells and ²¹⁵ found that changes in the morphology and motility of cells upon treatment were both drug and ²¹⁶ cell-line specific with different variables selected (**Figure 4**).

As was observed within the 231Docetaxel timelapses, cells increased in size due to failed cytoki-217 nesis. However, MCF-7 cells maintained a polygonal, epithelial-like morphology following treatment 218 similar to that of the untreated population. Conversely, remarkable differences in cellular dynamics 219 were observed within the 231Doxorubicin data set, with motility of cells being severely hindered 220 following treatment, particularly after the 24-hour time point. Only subtle differences in size and 221 morphology of cells were observed by eye, with doxorubicin treated cells appearing slightly enlarged 222 as a result of cell cycle arrest. Both untreated and treated sets contained examples of cells in G1 223 and G2, hence varied cell morphology can be observed within both (elongated and adherent cells in 224 G1, round and dense morphology of cells in G2.) 225

The 100 variables that achieved greatest separation for each of the MCF7Docetaxel and 231Dox-226 orubicin data sets are shown in Figure 4b. Density variables were highly discriminatory for un-227 treated and docetaxel treated MCF-7 cells, characterising decreased proliferation and cell-cell ad-228 hesion induced by drug treatment. Size, shape and texture variables were also identified as most 229 discriminatory with variables such as length, width and area characterising the enlarged cell shape 230 of treated cells. Spatial distribution variables were chosen for several intensity thresholds, demon-231 strating differences in the clustering of pixels, following docetaxel treatment. As was observed by 232 eye, movement features formed the majority of discriminatory variables for the 231Doxorubicin data 233 set, with untreated cells having greater velocity, tracklength and displacement than treated cells. 234 Differences in movement were also described through density ascent and descent, as cell density 235 fluctuated more for untreated cells due to the increased likelihood of passing neighbouring cells 236 when migrating. Subtle differences in cell shape and size observed by eye upon doxorubicin treat-237 ment were described by changes in rectangularity, width and radius variables. Notably both data 238 sets received lower separation scores than the 231Docetaxel data set, with 231Doxorubicin having 239 the lowest. This effectively provides a measure of class similarity, with high separation scores for 240 231Docetaxel indicative of significant changes to cells upon treatment and low separation scores for 241 231Doxorubicin suggesting these changes are more subtle. 242

PCA scores plots obtained with the selected features are shown in **Figure 4c**. Differences between classes can be observed for the MCF7Docetaxel data set, with separation of classes along PC1 (40% of the total variance) and PC2 (13% of the total variance). The PCA scores plot for 231Doxorubicin shows the greatest source of variance to be due to class differences, with separation of classes along PC1 (49% of the total variance). All PCA scores plots demonstrated the potential to characterise untreated and treated cell behaviour, with feature selected variables providing good distinction of classes which was improved by using variables in combination.



Figure 4: a) Images taken from cell timelapses of i. untreated and 1μ M docetaxel treated MCF-7 cells and ii. untreated and 1μ M doxorubicin treated MDA-MB-231 cells. Scale bar = 200μ 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 49h time point due to failed cytokinesis. In comparison, differences in morphology are more subtle for doxorubicin treated MDA-MB-231 cells at the 49h time point. b) Features with the top 100 highest separation scores, colour-coded according to feature type for i. MCF7Docetaxel, where cell density and texture provide greatest separation, and ii. 231Doxorubicin where shape and movement features provide greatest separation. c) PCA scores plot with points colour-coded according to true class label for i. MCF7Docetaxel and for ii. 231Doxorubicin. Only features with the 100 highest separation scores were included in PCA.

²⁵² Classification of Treated and Untreated Cells

251

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 MCF7Docetaxel and 231Doxorubicin generally having much lower separation scores (Figure 5a). Optimal separation thresholds of 0.075, 0.025 and 0.025 were obtained for 231Docetaxel, MCF7Docetaxel and 231Doxorubicin respectively, resulting in 437, 539 and 442 variables (of a possible 1111) being selected for classifier training.

Having chosen an optimal separation threshold, we trained an ensemble classifier for each data 259 set as described in Section 2.6. Classification accuracy scores for training and test sets obtained 260 using our ensemble classifier are provided in Table 2. Through visual inspection, we found that 261 misclassifications formed subsets of cells whose behaviour deviated from the behaviour of the main 262 population, we call this subset "non-conforming". (Figure 5b). For untreated cells, we found 263 that healthy, proliferating cells were correctly classified whereas less motile cells, cell debris or large, 264 non-motile mutant cells were instead classified as treated. For treated cells, we found that cells expe-265 riencing the drug-induced phenotypic differences identified through feature selection were classified 266 as treated. However, treated cells displaying behaviour similar to that of an untreated cell, such 267 as increased migration or fluctuation and elongation in cell shape, and were classified as untreated 268 (Figure 5c). 269

We found that the proportion of non-conforming treated cells, those classified as untreated, 270 decreased as drug concentration increased for all three data sets (Figure 5d). To explore the con-271 nection between the proportion of non-conforming treated cells and the population drug response 272 of each treated set, we considered the total volume growth rate at each drug concentration in re-273 lation to the percentage of cells predicted as untreated (Figure 5d). We found that the overall 274 growth rate decreased with increased drug concentration due to more cells responding at higher 275 concentrations. This correlated positively with the percentage of cells predicted as untreated, with 276 greater percentage of cells predicted as untreated for high volume growth rate with proliferation 277 still occurring. 278



279

280

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 0.075, 0.025 and 0.025 were selected for 231Docetaxel, MCF7Docetaxel and 231Doxorubicin respectively resulting in 437, 539 and 442 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. Time-lapse 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 (\log_{10} scale). For all three data sets, this percentage decreases as drug concentration increases due to a greater number of cells responding to treatment at higher concentrations. Lines were fitted using asymmetric, five parameter, non-linear regression. 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.

	231Docetaxel	MCF7Docetaxel	231Doxorubicin
	Untreated: 98%	Untreated: 100%	Untreated: 100%
Train	Treated: 100%	Treated: 99%	Treated: 100%
	Overall: 99%	Overall: 100%	Overall: 100%
	Untreated: 97%	Untreated: 83%	Untreated: 86%
\mathbf{Test}	Treated: 85%	Treated: 90%	Treated: 66%
	Overall: 94%	Overall: 85%	Overall: 81%

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

281 Subset Identification

²⁸² Classification accuracy scores for the untreated and treated cell populations were imbalanced ²⁸³ across all three of the data sets (**Table 2**). Imbalance of classification accuracy scores in binary classification is often a result of hidden stratification,²² where poor performance of one class is a result
 of misclassifications of important, unlabeled subsets. To investigate this phenomenon we performed
 hierarchical clustering on 231Docetaxel treated cells and the obtained dendogram is provided in

²⁸⁷ Figure 6a, with examples of cells from each cluster.

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

Cells in the same cluster share similar properties and morphological differences between clusters 292 of different cell cycle states can be observed. For example cells in clusters 1 and 2 are much smaller 293 and brighter than cells in clusters 3 and 5 as the cells are heading towards attempted mitosis, 294 confirmed by visual inspection of cell time-lapses, and hence resemble untreated mitotic cells. The 295 PCA biplot in **Figure 6c** shows how variables work in combination to determine cell clusters. 296 Clusters 1 and 2 are generally bright and spherical, similar to a mitotic treated cell, as these cells 297 are tracked prior to failed cytokinesis. Cells that have attempted to split, clusters 3 and 5, are 298 larger, longer, wider and display greater irregularity in shape. These cells become less dense and 299 are often multinucleated resulting in changes to texture features. Cell debris is best distinguished 300 by granularity, hence texture metrics are fundamental in identifying these instances. 301

Clusters also spanned a range of mean cell volumes beyond those of the untreated set when 302 hierarchical clustering was repeated for MCF7Docetaxel treated cells. However, this was not the 303 case for 231Doxorubicin treated cells and therefore k-means clustering was used to explore the 304 connection between misclassifications and hidden subsets in the 231Doxorubicin treated cell test set. 305 Two distinct clusters were obtained (Figure 6di), cluster 1 was formed of 33 cells and cluster 2 of 306 32 cells. We calculated classification accuracy scores for the two clusters individually and found that 307 91% of cells in cluster 1 were correctly classified as treated but only 31% in cluster 2 (Figure 6eii). 308 The increased migration and fluctuation in shape of cells in cluster 2 mean these cells have greater 309 similarity to the untreated population (Figure 6eiii). These non-conforming treated cells form the 310 majority of treated cell misclassifications in the 231Doxorubicin test set and highlight the presence 311 of heterogeneous subsets within a population. 312

Notably there was a greater number of misclassifications for untreated MCF-7 cells in comparison to the docetaxel treated set. Cluster analysis demonstrated the presence of heterogeneous subsets within the untreated population, with one cluster in particular consisting mainly of misclassified cells (**Figure S1**). Texture metrics discerned this cluster from other untreated cell clusters, containing several instances of cell debris that were understandably classified as "non-conforming". Other cells within this cluster shared similarities in texture to cell debris.



319

320

Figure 6: a) i. Dendogram obtained from hierarchical clustering of 231Docetaxel treated cells, with 5 clusters coloured. ii. 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. b) Density plots of mean cell volume, colour-coded according to cluster. The grey, dashed density plot represents 231Docetaxel untreated cells for reference. Cluster 4 (cell debris cluster) has the greatest leftward shift due to cells losing volume upon cell death. Clusters 1 and 2 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. c) d) 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 91% (30/33) of cells correctly classified as treated, whereas cluster 1 formed a subset of non-conforming treated cells, with only 31% (10/32) correctly classified as treated. iii. Increased velocity and ascent in cell elongation are characteristic of untreated cells. These metrics show extremely significant decrease for cells in cluster 1 but no significant difference for cells in cluster 2. Extremely significant differences are observed between cluster 1 and cluster 2, highlighting the presence of subsets within the treated cell population (ns: $p \ge 0.05$, ****: p < 0.0001, dashed lines in violin plots are representative of the lower quartile, median and upper quartile).

321 Compatibility with fluorescence images and TrackMate

TrackMate-Cellpose¹⁷ was used to demonstrate the compatibility of CellPhe with outputs ob-322 tained from alternative segmentation and tracking software and show that CellPhe extends to flu-323 orescence time-lapse imaging. Ptychographic and fluorescence time-lapse images of untreated and 324 docetaxel treated MDA-MB-231 cells stably expressing dsRed were acquired in parallel (Fig 7a). 325 Cell segmentation from the fluorescence images was performed using Cellpose and a representative 326 image is provided in Fig 7bi. Segmented cells were then tracked using TrackMate resulting in 123 327 cell tracks of greater than or equal to 50 frames (Fig 7bii). The resulting folders of cell ROIs and 328 TrackMate feature tables were used as input for CellPhe to extract single-cell phenotypic metrics 329 to describe cell behaviour over time. An optimal separation threshold of 0.3 was determined for 330 discrimination between untreated and treated cells, with 231 variables achieving separation scores 331 greater than the threshold (Fig 7c). As observed with the phase images, size, shape and texture 332 variables provide greatest separation, with cell density amongst the most discriminatory variables. 333

 $_{334}$ Good separation of untreated and treated cells can be observed within the PCA scores plot in Fig

³³⁵ 7d, supporting the use of CellPhe for cell phenotyping from fluorescence images.



337

338

Figure 7: a) Images taken from cell timelapses of untreated and 1μ M docetaxel treated MDA-MB-231 cells stably expressing dsRed. Phase and fluorescence images were acquired in parallel. Scale bar = 200 μ m. b) i. Representative image of Cellpose segmentation on a fluorescent image of MDA-MB-231 cells stably expressing dsRed. ii. Cell tracks obtained from TrackMate for untreated MDA-MB-231 cells stably expressing dsRed. Only cell tracks greater than or equal to 50 frames are displayed. c) Features with separation scores greater than or equal to 0.3, the optimal separation threshold, colour-coded according to feature type. Texture, density, shape and size features provide greatest separation. d) PCA scores plot with points colour-coded 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. Only features with separation score greater than or equal to 0.3 were included in PCA.

339 **3** Discussion

The CellPhe toolkit complements existing software for automated cell segmentation and track-340 ing, using their output as a starting point for bespoke time series feature extraction and selection, 341 cell classification and cluster analysis. Erroneous cell segmentation and tracking can significantly 342 reduce data quality but such errors often go undetected and can negatively influence the results 343 of automated pattern recognition. CellPhe's extensive feature extraction followed by customised 344 feature selection not only allows the characterisation and classification of cellular phenotypes from 345 time-lapse videos but provides a method for the identification and removal of erroneous cell tracks 346 prior to these analyses. Attribute analysis showed that different features were chosen to identify seg-347 mentation errors for different cell lines. For example, sudden increases in movement resulting from 348 large boundary changes can indicate segmentation errors for MCF-7 cells, contrasting with their 349 innate low motility. On the other hand, size and texture variables provide better characterisation of 350

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

For cell characterisation, we have shown that CellPhe's feature selection method is able to adapt 356 to different experimental conditions, providing discrimination between untreated and treated groups 357 of two different breast cancer cell lines (MDA-MB-231 and MCF-7) and two different chemotherapy 358 treatments (docetaxel and doxorubicin). The discriminatory variables identified here coincide with 359 previously reported effects of docetaxel or doxorubicin treatment and can be interpreted in terms 360 of the mechanism of action of each drug. Previous studies have identified a subset of polyploid, 361 multinucleated cells following docetaxel treatment due to cell cycle arrest and occasionally cell cycle 362 slippage.²³ Our findings support this with shape and size variables providing the greatest separation 363 for docetaxel treatment in both MDA-MB-231 and MCF-7 cells. Many texture variables were also 364 identified as discriminatory following docetaxel treatment, providing label-free identification of the 365 multiple clusters of high intensity pixels in treated cells, likely a result of docetaxel-induced multin-366 ucleation. We found that at a higher, sub-lethal concentration of 1µM, migration of MDA-MB-231 367 cells was reduced with variables associated with movement providing greatest discrimination be-368 tween untreated and doxorubicin treated cells. This is supported by studies that have identified 369 changes in migration of doxorubicin treated cells, noting that low drug concentrations in fact facili-370 tate increased invasion.^{24,25} 371

We found an imbalance in untreated and treated classification accuracy scores, with a greater 372 proportion of treated cells misclassified for all three data sets. This consistent imbalance suggests the 373 misclassifications are in fact representative of a subset of non-conforming, and potentially chemore-374 sistant, cells. The concept of hidden stratification, where an unlabelled subset performs poorly 375 during classification, has been described previously 26 and poses a challenge in medical research as 376 important subsets (such as rare forms of disease) could be overlooked. Here, the misclassified cells 377 could be of most interest and the ability to identify non-conforming behaviour is precisely what 378 is required from a classifier as treated cells that display behaviour similar to untreated cells could 379 indicate a reduced response to drug treatment. The classification of cells treated with a range of 380 concentrations supported this hypothesis as a greater proportion of cells were classified as untreated 381 at lower drug concentrations, demonstrating that our trained ensemble classifier can be used to 382 quantify drug response, at both single-cell and populational level. 383

Cluster analysis revealed cell subsets that appear to represent different responses to drug treat-384 ment. Heterogeneity of cellular drug response is a commonly reported phenomenon in cancer treat-385 ment, yet mechanisms underlying this are not well understood.²⁷ Analysis of cell volumes showed 386 the mean volume of treated and untreated cells to be comparable for doxorubicin reflecting the fact 387 that this treatment can induce G1, S or G2 cell cycle arrest.²⁸ However, for docetaxel treated cells, 388 we found that clusters spanned a range of mean cell volumes beyond those of the untreated set for 389 both cell lines. Clustering allowed identification of three general responses to docetaxel treatment: 390 pre-"cytokinesis attempt", with cells having similar volumes to the untreated MDA-MB-231 popula-391 tion; post-"cytokinesis attempt", where cells were tracked following failed cytokinesis and therefore 392 continued to grow to volumes beyond those of the late stages of the untreated cell cycle; and cell 393 death, with a final cluster, composed primarily of cell debris. Furthermore, giant cell morphology 394 has been linked with docetaxel resistance, a potential cause of relapse in breast cancer patients⁹ and 395 through cluster analysis we were able to identify a potentially resistant subset of very large, treated 396 cells that could be isolated for further investigation. 397

Our chosen application demonstrated the breadth of quantification and biological insight that can be made by following our workflow, with characterisation of drug response and detection of potentially resistant cells just two of many potential applications for CellPhe. CellPhe offers several benefits for the quantification of cell behaviour from time-lapse images. First, errors in cell segmentation and tracking can be identified and removed, improving the quality of input for downstream data analysis. This is particularly important with machine learning where automation means that such errors can easily be missed, and algorithms consequently trained with poor data. Although different cell lines have different properties that allow segmentation errors to be recognised, we have shown that ground truth data for a particular cell-line can be re-used for different experiments, in our case, different drug treatments.

Second, cell behaviour is characterised over time by extracting variables from the time series of various features whereas many studies explore temporal changes by collecting data at discrete time points (for example, 0 and 24 hours post-treatment) and using metrics from each static image, missing behavioural changes experienced by cells on a continuous level. With CellPhe, changes over time in features that provide information on morphology, movement and texture are quantified not just by summary statistics but by variables extracted from wavelet transformation of the time series allowing changes on different scales to be identified.

Third, whilst most studies use a limited number of metrics, assessed individually for discrimi-415 nation between groups,^{29,30} CellPhe provides an extensive list of novel metrics and automatically 416 determines the combination that offers greatest discrimination. The bespoke feature selection fre-417 quently found the most discriminatory variables to be those with the ability to detect changes in 418 cell behaviour over time. Previous research in this field has focused on identification of cell types 419 from co-cultures³¹ for use in automated diagnosis of disease such as cancer. Analysis methods for 420 these studies are often cell line specific whereas CellPhe's feature selection method is successful in 421 identifying discriminatory variables tailored to different experimental conditions. 422

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

The interactive, interpretable, high-throughput nature of CellPhe deems it suitable for all cell 427 time-lapse applications, including drug screening or prediction of disease prognosis. We provide a 428 comprehensive manual with a working example and real data to guide users through the workflow 429 step-by-step, where users can interact with each stage of the workflow and customise to suit their 430 own experiments. Here we demonstrated the abundance of information and insight that can be 431 made by following the CellPhe workflow to quantify cell behaviour from QPI images. CellPhe can 432 be used with tracking information from multiple segmentation and tracking algorithms and different 433 imaging modalities, including fluorescence, and would be suitable for all time-lapse studies including 434 clinical applications. 435

436

437 4 Acknowledgements

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

445 5 Contributions

⁴⁴⁶ Conceptualisation: W.B., P.O'T., J.W., and L.W.; cell culture, pharmacology and imaging: L.W.
⁴⁴⁷ and A.L.; data analysis and validation: L.W. and J.W.; software development J.W., L.W, S.L.
⁴⁴⁸ and K.M; supervision: J.W., W.B. and P.O'T.; writing-original draft preparation, L.W. and J.W.;
⁴⁴⁹ writing-review and editing, W.B. and P.O'T.

451 6 Competing interests

⁴⁵² The authors declare no competing interests.

453 7 Code availability

The source code for algorithms developed during this research has been deposited in GitHub: https:
 //github.com/uoy-research/CellPhe/. The interactive CellPhe GUI can be accessed here: https:
 //cellphegui.shinyapps.io/app_to_host/.

457 8 Data availability

All data used to produce the results in the manuscript, including separate data that will allow the user to follow the worked example in the CellPhe user guide, are available from https://doi.org/ 10.15124/936b6b09-a341-40ee-b08f-c049316ac247. Here, the file example_data.zip contains all the data required to follow the worked example and the file CellPhe_GUI_demo_vid.mov is a video that explains how to use the GUI.

463 **References**

- ⁴⁶⁴ ¹ S. Turajlic, A. Sottoriva, T. Graham, et al. Resolving genetic heterogeneity in cancer. Nature
 ⁴⁶⁵ Reviews Genetics, 20:404–416, 2019.
- ⁴⁶⁶ ² S. Goldman, M. MacKay, E. Afshinnekoo, et al. The impact of heterogeneity on single-cell
 ⁴⁶⁷ sequencing. *Frontiers in Genetics*, 10:8, 2019.
- ³ S.J. Altschuler and L.F. Wu. Cellular heterogeneity: do differences make a difference? National
 Institute of Health, Cell, 141:559–563, 2010.
- ⁴⁷⁰ ⁴ B. Carter and K. Zhao. The epigenetic basis of cellular heterogeneity. *Nature Reviews Genetics*, ⁴⁷¹ 22:235–250, 2021.
- ⁵ F. Buettner, K. Natarajan, F. Casale, et al. Computational analysis of cell-to-cell heterogeneity
 in single-cell rna-sequencing data reveals hidden subpopulations of cells. *Nature Biotechnology*,
 33:155–160, 2015.
- ⁴⁷⁵ ⁶ B.M. Davis, M. Salinas-Navarro, M.F. Cordeiro, et al. Characterizing microglia activation: a ⁴⁷⁶ spatial statistics approach to maximize information extraction. *Scientific Reports*, 7, 2017.
- ⁷ T. Henser-Brownhill, R.J. Ju, N.K. Haass, et al. Estimation of cell cycle states of human melanoma cells with quantitative phase imaging and deep learning. *IEEE 17th International Symposium on Biomedical Imaging (ISBI)*, pages 1617–1621, 2020.
- ⁴⁸⁰ ⁸ L.A. Tashireva, M.V. Zavyalova, O.E. Savelieva, et al. Single tumor cells with epithelial-like ⁴⁸¹ morphology are associated with breast cancer metastasis. *Frontiers in Oncology*, 10:50, 2020.
- ⁹ R. Mirzayans, B. Andrais, and D. Murray. Roles of polyploid/multinucleated giant cancer cells in
 metastasis and disease relapse following anticancer treatment. *Cancers (Basel)*, 10(4):118, 2018.
- ⁴⁸⁴ ¹⁰ A. Voulodimos, N. Doulamis, A. Doulamis, et al. Deep learning for computer vision: A brief
 ⁴⁸⁵ review. *Hindawi*, 2018(Computational Intelligence and Neuroscience), 2018.
- ⁴⁸⁶ ¹¹ W. Chen, W. Li, X. Dong, and J. Pei. A review of biological image analysis. *Current Bioinformatics*, 12, 2017.

- ⁴⁸⁸ ¹² Andreas P Cuny, Aaron Ponti, Tomas Kündig, Fabian Rudolf, and Jörg Stelling. Cell region
 ⁴⁸⁹ fingerprints enable highly precise single-cell tracking and lineage reconstruction. *Nature Methods*,
 ⁴⁹⁰ pages 1–10, 2022.
- ⁴⁹¹ ¹³ H.E. Munim and A.A. Farag. A shape-based segmentation approach: an improved technique using
 ⁴⁹² level sets. *Tenth IEEE International Conference on Computer Vision, Volume 1*, 2:930–935, 2005.
- ⁴⁹³ ¹⁴ Z. Wang and H. Li. Generalizing cell segmentation and quantification. BMC Bioinformatics, 18,
 ⁴⁹⁴ 2017.
- ⁴⁹⁵ ¹⁵ E. Gómez-de Mariscal, C. García-López-de Haro, W. Ouyang, et al. Deepimagej: A user-friendly
 ⁴⁹⁶ environment to run deep learning models in imagej. *Nature Methods*, 18:1192–1195, 2021.
- ⁴⁹⁷ ¹⁶ D. Ershov, Minh-Son Phan, J. Pylvänäinen, et al. Bringing trackmate in the era of machine ⁴⁹⁸ learning and deep-learning. *bioRxiv*, 2021.
- ⁴⁹⁹ ¹⁷ Dmitry Ershov, Minh-Son Phan, Joanna W Pylvänäinen, Stéphane U Rigaud, Laure Le Blanc,
 ⁵⁰⁰ Arthur Charles-Orszag, James RW Conway, Romain F Laine, Nathan H Roy, Daria Bonazzi, et al.
 ⁵⁰¹ Trackmate 7: integrating state-of-the-art segmentation algorithms into tracking pipelines. *Nature* ⁵⁰² *Methods*, pages 1–4, 2022.
- ¹⁸ R.F. Laine, I. Arganda-Carreras, R. Henriques, et al. Avoiding a replication crisis in deep-learning based bioimage analysis. *Nature Methods*, 18:1136–1144, 2021.
- ¹⁹ J. Marrison, L. Räty, P. Marriott, et al. Ptychography a label-free, high contrast imaging technique for live cells using quantitative phase information. *Scientific Reports*, 3(2369), 2013.
- ⁵⁰⁷ ²⁰ Y. Rivenson, Y. Zhang, H. Günaydın, et al. Phase recovery and holographic image reconstruction
 ⁵⁰⁸ using deep learning in neural networks. *Nature Light Sci Appl.*, 7(17141), 2018.
- ⁵⁰⁹ ²¹ Y. Park, C Depeursinge, and G. Popescu. Quantitative phase imaging in biomedicine. Nature Photon, 12:578–589, 2018.
- ⁵¹¹ ²² L. Oakden-Rayner, J. Dunnmon, G. Carneiro, et al. Hidden stratification causes clinically mean ⁵¹² ingful failures in machine learning for medical imaging. arXiv, 2019.
- ⁵¹³ ²³ H. Hernandez-Vargas, J. Palacios, and G. Moreno-Bueno. Molecular profiling of docetaxel cytotox ⁵¹⁴ icity in breast cancer cells: uncoupling of aberrant mitosis and apoptosis. *Oncogene*, 26:2902–2913,
 ⁵¹⁵ 2007.
- ⁵¹⁶ ²⁴ J. Liu, L. Qu, L. Meng, et al. Topoisomerase inhibitors promote cancer cell motility via ros ⁵¹⁷ mediated activation of jak2-stat1-cxcl1 pathway. *Journal of Experimental and Clinical Cancer* ⁵¹⁸ Research, 38:370, 2019.
- ²⁵ C.L. Liu, M.J. Chen, J.C. Lin, et al. Migration and invasion of breast cancer cells through the
 ²⁵ upregulation of the rhoa/mlc pathway. j breast cancer. Journal of Breast Cancer, 22:185–195,
 ²⁰ 2019.
- ⁵²² ²⁶ N.S. Sohoni, J.A. Dunnmon, G. Angus, et al. No subclass left behind: Fine-grained robustness in
 ⁵²³ coarse-grained classification problems. *CoRR*, 2020.
- ⁵²⁴ ²⁷ R. Wang, C. Jin, and X. Hu. Evidence of drug-response heterogeneity rapidly generated from a
 ⁵²⁵ single cancer cell. *Oncotarget*, 8:25, 2017.
- ²⁸ X. Wang, Z. Chen, A.K. Mishra, et al. Chemotherapy-induced differential cell cycle arrest in b-cell
 lymphomas affects their sensitivity to weel inhibition. *Haematologica.*, 103(3):466–476, 2018.
- ⁵²⁸ ²⁹ F. M. Frame, A. R. Noble, S. Klein, et al. Tumor heterogeneity and therapy resistance impli ⁵²⁹ cations for future treatments of prostate cancer. *Journal of Cancer Metastasis and Treatment*,
 ⁵³⁰ 3:302–314, 2017.

- ³⁰ R. Suman, G. Smith, and K. E.A. Hazel et al. Label free imaging to study phenotypic behavioural
 traits of cells in complex co-cultures. *Scientific Reports*, 6(1):22–32, 2016.
- ⁵³³ ³¹ Y. Ozaki, H. Yamada, H. Kikuchi, et al. Label-free classification of cells based on supervised
 ⁵³⁴ machine learning of subcellular structures. *PLoS One*, 14(1), 2019.
- ³² M. Yang, D.J. Kozminski, L. Wold, et al. Therapeutic potential for phenytoin: targeting nav1.5
 ^{sodium} channels to reduce migration and invasion in metastatic breast cancer. Breast Cancer
 Research and Treatment, 134(2):603–615, 2012.
- ³³ C. Uphoff, S. Gignac, and H. Drexler. Mycoplasma contamination in human leukemia cell lines.
 i. comparison of various detection methods. *Journal of Immunological Methods*, 149:43–53, 1992.
- ³⁴ R. Kasprowicz, R. Suman, and P. O'Toole. Characterising live cell behaviour: Traditional label free and quantitative phase imaging approaches. *The international journal of biochemistry & cell biology*, 84:89–95, 2017.
- ⁵⁴³ ³⁵ Godfried T Toussaint. Solving geometric problems with the rotating calipers. In *Proc. IEEE* ⁵⁴⁴ *Melecon*, volume 83, page A10, 1983.
- ³⁶ Julie Wilson. Towards the automated evaluation of crystallization trials. Acta Crystallographica
 Section D: Biological Crystallography, 58(11):1907–1914, 2002.
- ⁵⁴⁷ ³⁷ Julie Wilson, Karen Hardy, Richard Allen, Les Copeland, Richard Wrangham, and Matthew
 ⁵⁴⁸ Collins. Automated classification of starch granules using supervised pattern recognition of morphological properties. *Journal of Archaeological Science*, 37(3):594–604, 2010.
- ⁵⁵⁰ ³⁸ Urs Ramer. An iterative procedure for the polygonal approximation of plane curves. *Computer* ⁵⁵¹ graphics and image processing, 1(3):244–256, 1972.
- ³⁹ Namita Aggarwal and R.K. Agrawal. First and second order statistics features for classification of
 magnetic resonance brain images. *Journal of Signal and Information Processing*, 3:146–153, 2012.
- ⁴⁰ L.K. Soh and C. Tsatsoulis. Texture analysis of sar sea ice imagery using gray level co-occurrence
 matrices. *IEEE Transactions on geoscience and remote sensing*, 37(2):780–795, 1999.
- ⁵⁵⁶ ⁴¹ S. Mallat. A theory for multiresolution signal decomposition: the wavelet representation. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, 7:674–693, 1989.
- ⁴² R.M. Haralick, K. Shanmugam, and I. Dinstein. Textural features for image classification. *IEEE Transactions on systems, man, and cybernetics*, 6:610–621, 1973.
- ⁴³ Tommy Löfstedt, Patrik Brynolfsson, Thomas Asklund, Tufve Nyholm, and Anders Garpebring.
 Gray-level invariant haralick texture features. *PloS one*, 14(2):e0212110, 2019.
- ⁴⁴ A. Haar. Zur theorie der orthogonalen funktionensysteme. Mathematische Annalen, 69(3):331–
 371, 1910.
- ⁴⁵ N.V. Chawla, K.W. Bowyer, L.O. Hall, et al. Smote: Synthetic minority over-sampling technique.
 Journal of Artificial Intelligence Research, 16:321–357, 2002.
- ⁴⁶ R Core Team. R: A Language and Environment for Statistical Computing. R Foundation for
 Statistical Computing, Vienna, Austria, 2019.
- ⁴⁷ W. N. Venables and B. D. Ripley. *Modern Applied Statistics with S.* Springer, New York, fourth
 edition, 2002.
- ⁴⁸ D. Meyer, E. Dimitriadou, et al. e1071: Misc Functions of the Department of Statistics, Probability Theory Group (Formerly: E1071), TU Wien, 2019. R package version 1.7-3.
- ⁵⁷² ⁴⁹ A. Liaw and M. Wiener. Classification and regression by randomforest. *R News*, 2(3):18–22, 2002.

⁵⁷³ ⁵⁰ A. Kassambara and F. Mundt. factoextra: Extract and Visualize the Results of Multivariate Data
 ⁵⁷⁴ Analyses, 2020. R package version 1.0.7.

575 Methods

⁵⁷⁶ **Cell Culture.** MDA-MB-231 cells and MCF-7 cells were cultured separately in Dulbecco's modified ⁵⁷⁷ eagle medium supplemented with 5% fetal bovine serum and 4mM L-glutamine.³² Fetal bovine serum ⁵⁷⁸ was filtered using a 0.22μ m syringe filter prior to use to reduce artefacts when imaging. Cells were ⁵⁷⁹ incubated at 37°C in plastic filter-cap T-25 flasks and were split at a 1:6 ratio when passaged. No ⁵⁸⁰ antibiotics were added to cell culture medium. Cells were confirmed to be mycoplasma-free by 4',6-⁵⁸¹ diamidino-2-phenylindole (DAPI) method.³³ In cases where dsRed expressing MDA-MB-231 cells ⁵⁸² were used, cells were sorted via FACS prior to imaging to enrich for a transfected cell population.

To image the following day, cells were counted and then seeded in a Corning Costar plastic, flat bottom 24-well plate. Cells were seeded at a density of 8000 cells per well with a final volume of 500μ L in each of the 24 wells.

586

Pharmacology. Docetaxel (Cavman Chemical Company) was prepared as 5mg/mL of DMSO and 587 doxorubicin (AdooQ Bioscience) as 25mg/mL of DMSO, both were then frozen into aliquots. Once 588 thawed, docetaxel and doxorubicin stock solutions were diluted in culture medium to give final 589 working concentrations. Docetaxel dose response analysis for both MDA-MB-231 and MCF-7 cells 590 involved imaging eight wells treated with the following concentrations of docetaxel: 0nM, 1nM, 3nM, 591 10nM, 30nM, 100nM, 300nM, 1 μ M, with additional concentrations 3μ M, 10μ M and 30μ imaged for 592 MDA-MB-231 cells. Doxorubicin dose response analysis for MDA-MB-231 cells involved imaging 593 eight wells treated with the following concentrations of doxorubicin: 0nM, 10nM, 30nM, 100nM, 594 $300nM, 1\mu M, 3\mu M, 10\mu M.$ 595

⁵⁹⁶ Medium was removed from wells selected to receive treatment 30 minutes prior to image acqui-⁵⁹⁷ sition, and 500μ L of desired drug concentration was added to each well. Control wells received a ⁵⁹⁸ medium change and were treated with DMSO vehicle on the day of imaging to maintain consistent ⁵⁹⁹ DMSO concentration throughout.

600

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

For phase images, Phasefocus' Cell Analysis Toolbox[®] software was utilised for cell segmentation, cell tracking and data exportation. Segmentation thresholds were optimised for a range of image processing techniques such as rolling ball algorithm to remove background noise, image smoothing for cell edge detection and local pixel maxima detection to identify seed points for final consolidation.

The Phasefocus software outputs a feature table for each imaged well. Information on missing frames for tracked cells can be obtained from this table which also provides descriptive features. However, most features are calculated within CellPhe and we only utilise the Phasefocus' features that rely on phase information, these being the volume of the cell and sphericity.³⁴

For fluorescence images, the TrackMate-Cellpose ImageJ plugin was used for cell segmentation and tracking. Cells were segmented using Cellpose's pre-trained cytoplasm model and image contrast was enhanced prior to segmentation to improve detection of cell boundaries. Once complete, TrackMate feature tables and individual cell ROIs were exported from ImageJ. Prior to use with CellPhe, it was necessary to interpolate TrackMate-Cellpose ROIs to obtain a complete list of cell boundary coordinates. Interpolation of ROIs was performed using a custom ImageJ macro.

622

⁶²³ Implementation of CellPhe.

Feature extraction. Using cell boundary information from Regions of Interest (ROIs) produced by 624 the Phasefocus software or TrackMate, a range of morphological and texture features were extracted 625 for each cell that was tracked for at least 50 frames. In addition to size and shape descriptors cal-626 culated from the cell boundaries, a filling algorithm was used to determine the interior pixels from 627 which texture and spatial features were extracted. The local density was also calculated as the sum 628 of inverse distances from the cell centroid to those of neighbouring cells within three times the cells 629 diameter. A complete list of features together with their definitions is provided in **Supplementary** 630 table S1. 631

632

Movement descriptors. By considering the position of a cell's centroid on subsequent frames, vari-633 ables describing the cell's movement were extracted from the images. The current speed of the cell 634 estimated by considering its position in consecutive frames, taking into account any missing frames. 635 The measure provided is proportional to rather than equal to velocity as this would require the rate 636 at which frames were produced to be entered by the user for no gain in discriminatory power. The 637 displacement, or straight line distance between the cell centroid on the current frame and the frame 638 it was first detected in, and the tracklength or total path length travelled by the cell up to the 639 current frame, are also calculated. To see how these vary, the quotient current tracklength/current 640 displacement is also calculated. 641

642

Size descriptors. In addition to volume, calculated using phase information, the size variables determined are cell area, as the number of pixels within (or on) the cell boundary, the length and width of the cell, determined from the minimal rectangular box that the cell can be enclosed by,³⁵ and the radius, as the average distance of boundary pixels from the cell centroid.

647

Shape descriptors. We make use of an imported feature, sphericity, which requires phase information 648 for calculation, but extract a number of other shape features within CellPhe. As well as determining 649 the length and width from the arbitrarily oriented minimum bounding box, we use this to provide 650 a measure of 'rectangularity' as max(x,y)/(x+y) where x and y are the length and width of the 651 minimal bounding box.³⁶ We also consider the shape of the cell by calculating the fraction of the 652 minimal box area that the cell area covers and by comparing the number of pixels on the boundary 653 with the total pixels within the cell.³⁶ Here the number of boundary pixels is squared in the quotient 654 to avoid the effect of cell size. We also calculate the variance on the distance from the centroid to 655 the boundary pixels, with more circular cells having less variance 36 and an measure of boundary 656 curvature based of the triangle inequality.³⁷ Finally 4 shape descriptors are obtained from a poly-657 gon fitted to the cell boundary, being the mean and variance of both edge length and interior angle.³⁸ 658 659

Texture descriptors. Textural features of each cell are represented in terms of three first order 660 statistics calculated from the pixel intensities within the cell: mean, variance and skewness.³⁹ For 661 second order texture features, we used gray-level co-occurrence matrices (GLCMs)⁴⁰ but, rather than 662 consider the positions of pixels within a cell, we calculated GLCMs between the image of the cell at 663 different resolutions to differentiate textures that are sharp and would be lost at lower resolution from 664 those that are smooth and would remain. This was achieved by performing a two-level 2-D wavelet 665 transform⁴¹ on the pixels within the axis-aligned minimum rectangle containing a cell. GLCMs were 666 then calculated between the original interior pixels and the corresponding values from the first and 667 second levels of the transform as well as between the two sets of transformed pixels (levels 1 and 668 2). Statistics first described by $Haralick^{42}$ were then calculated from each GLCM. We use 14 of the 669 20 Haralick features described by Löfstedt et al.:⁴³ Angular Second Moment, Contrast, Correlation, 670 Variance, Homogeneity, Sum Average, Sum Variance, Entropy, Sum Entropy, Difference Variance, 671 Difference Entropy, Information Measure of Correlation 2, Cluster Shade, Cluster Prominence. With 672 three co-occurrence matrices, this gives 42 Haralick features. 673

We calculated spatial distribution descriptors to quantify the uniformity or clustering of cell interior pixels at different intensity levels. IQn is a measure of dispersion calculated for the subset of interior pixels with intensities greater than or equal to the $(n \times 10)$ th quantile. Based on a Poisson distribution, for which the mean is equal to the variance, the measure is calculated as the variance divided by the mean, calculated over the pairwise distances between pixels within the *n*th subset. IQn = 1 indicates a random distribution whereas a value of IQn less than 1 indicates that the pixels are more uniformly distributed and a value greater than 1 indicates clustering.

681

Characterising Cell Time Series. Cell tracking provides a time series for each of the 74 features 682 extracted for a cell. The length of the time series depends on how many frames the cell has been 683 tracked for and so differs between cells. In order to apply pattern recognition methods, we extracted a 684 fixed number of characteristic variables for each cell from the time series for each feature. Statistical 685 measures (mean, standard deviation and skewness) summarise time series of varying length, but 686 may not be representative of changes throughout the time series. Therefore, in addition to summary 687 statistics, we calculated variables inspired by elevation profiles in walking guides, that is, the sum of 688 any increases between consecutive frames (total ascent), the sum of any decreases (total descent) and 689 the maximum value of the time series (maximum altitude gain). Similar variables were calculated for 690 different levels of the wavelet transform of the time series to allow changes at different scales to be 691 considered. The wavelet transform decomposes a time series to give a lower resolution approximation 692 together with different levels of detail that need to be added to the approximation to restore the 693 original time series. Using the Haar wavelet $basis^{44}$ with the multiresolution analysis of Mallat⁴¹ 694 allows increases and decreases in the values of the variables to be determined over different time 695 scales. With Haar wavelets, a negative detail coefficient represents an increase from one point to 696 the next, and so we used the sum of the negative detail coefficients to provide the equivalent to 697 total ascent and the sum of the positive detail coefficients as total descent. Rather than an overall 698 maximum, we use the maximum detail coefficient for the transformed time series. 699

Occasionally the automated cell tracking misses a frame or even several frames, for example 700 when a cell temporarily leaves the field of view. To prevent jumps in the time series, we interpolated 701 values for the missing frames, although these values were not used to calculate statistics. After 702 interpolation, the three elevation variables were calculated from the original time series and three 703 wavelet levels which, together with the summary statistics, provided 15 variables for each feature 704 (Supplementary table S2). The 72 extracted features together with the 2 imported features would 705 have given $74 \ge 15 = 1110$ variables in total, but, as one feature, the tracklength or total distance 706 travelled up to the current frame, is monotonically increasing, the total descent is always zero and 707 therefore variables related to tracklength descent were not used. Similarly, as the tracklength and 708 displacement are the same for the first frame and the displacement can never be greater than the 709 tracklength, the maximum value for their quotient will always be 1 and this variable is also not used. 710 One further variable was introduced to summarise cell movement as the area of the minimal 711 bounding box around a cell's full trajectory. This area will be large for migratory cells and small for 712 cells whose movement remains local for the duration of the time series. If, within a cell's trajectory, 713 minX and minY are the minimal X and Y positions respectively with maxX and maxY the 714 corresponding maximal positions, then the trajectory area is defined as 715

$$trajectory area = (maxX - minX) \times (maxY - minY).$$
(1)

⁷¹⁶ Thus, a total of 1106 characteristic variables were available for analysis and classification.

717

Segmentation Error Removal. To improve characterisation of cellular phenotype, we only included cells that were tracked for at least 50 frames in our analyses. Whilst the majority of these cells were correctly tracked, others had segmentation errors, with confusion between neighbouring cells,

⁷²¹ missing parts of a cell or multiple cells included. ⁷²² In order to increase the reliability of our results, we developed a classification process to identify ⁷²³ and remove such cells prior to further analysis. Cells (both treated and untreated) were classified by ⁷²⁴ eye to provide a training data set. Due to class imbalance, with the number of segmentation errors ⁷²⁵ far less than the number of correct segmentations, the Synthetic Minority Oversampling Technique ⁷²⁶ (SMOTE)⁴⁵ was performed using the *smotefamily* package in R, with the number of neighbours K ⁷²⁷ set to 3, to double the number of instances representing segmentation errors.

The resulting data set with all 1111 variables was used to train a set of 50 decision trees using the *tree* package in R with default parameters. For each tree, the observations from cells with seg-

mentation errors were used together with the same number of observations randomly selected from 730 the correctly segmented cells to further address class imbalance. For each cell, a voting procedure 731 was used to provide a classification from the predictions of the 50 decision trees. To minimise the 732 number of correctly tracked cells being falsely classified as segmentation errors, this class was only 733 assigned when it received at least 70% of the votes (i.e. 35). To add further stringency, the training 734 of 50 decision trees was repeated ten times and a cell only given a final classification of segmentation 735 error if predicted this label in at least five of the ten runs. MDA-MB-231 cells that were not used 736 for training formed an independent test set. All cells either manually labelled as segmentation error 737 or predicted as such were excluded from further analyses. 738

739

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

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

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

⁷⁴⁸ and V_W is the within-group variance:

$$V_W = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{(n_1 + n_2 - 2)}.$$
(3)

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 means, \bar{x} the overall mean, and s_1^2 and s_2^2 are the sample variances. The most discriminatory variables were chosen for a particular data set by assessing the classification error on the training data to optimise the threshold on separation. Starting with a threshold of zero, the n^{th} separation threshold was minimised such that the classification error rate did not increase by more than 2% from that obtained for the $(n-1)^{\text{th}}$ threshold. The aim here was to reduce the risk of overfitting by only retaining variables achieving greater than or equal to this threshold for the next stage of classifier training.

Data were scaled to prevent large variables dominating the analysis and ensemble classification 757 used to take advantage of different classifier properties. The predictions from three classification 758 algorithms, Linear Discriminant Analysis (LDA), Random Forest (RF) and Support Vector Ma-759 chine (SVM) with radial basis kernel were combined using the majority vote. Model performance 760 was evaluated by classification accuracy, taking into account the number of false positives and false 761 negatives. All classification was performed in RStudio⁴⁶ using open-source packages. LDA was per-762 formed using the *lda* function from the MASS library,⁴⁷ SVM classification used the *svm* function 763 from the e1071 package⁴⁸ with a radial basis kernel and the randomForest package⁴⁹ was used to 764 train random forest classifiers with 200 trees and 5 features randomly sampled as candidates at each 765 split. 766

767

Cluster analysis. Both hierarchical clustering and k-means clustering were used to investigate sub-768 groups within single-class data sets (i.e. treated and untreated cells separately). Data were scaled 769 prior to clustering and analyses performed in R. Hierarchical clustering was implemented with the 770 factoextra package⁵⁰ using the *hcut* function to cut the dendrogram into k clusters. Agglomerative 771 nesting (AGNES) was used with Ward's minimum variance as the agglomeration method and the 772 Euclidean distance metric to quantify similarity between cells. k-means clustering was performed 773 using the R stats package, with the number of random initial configurations set to 50. The number of 774 clusters k was chosen to obtain clusters with meaningful interpretation. Similarities and differences 775 between clusters were identified through evaluation of separation scores to determine discriminatory 776 features, as well as through observation of cells within each cluster by eye. 777

Statistical tests. All tests of statistical significance within this study are two-tailed, non-parametric
 Mann-Whitney t-tests performed using Graphpad Prism 9.1.0 (GraphPad Software, San Diego, CA).

⁷⁸² Data. Three data sets were used to demonstrate our pipeline for the classification of untreated and ⁷⁸³ treated cells. For brevity we use abbreviations throughout to refer to each data set, for example ⁷⁸⁴ "231Docetaxel" is a data set consisting of MDA-MB-231 cells, both untreated and treated with ⁷⁸⁵ 30μ M docetaxel. This is the main data set used to develop the methods, with a training data set ⁷⁸⁶ compiled from 6 experiments performed on different days and an independent test data set compiled ⁷⁸⁷ from a further 3 experiments, also performed on separate days and by a different individual.

We validate our methods using two further datasets, the 231Doxorubicin and MCF7Docetaxel 788 data sets, details of which are given in Table 3. This table also includes details of the number of 789 cells within each training and test set. We show that the classification pipeline can be successfully 790 reproduced using fewer experimental repeats for the 231Doxorubicin and MCF7Docetaxel data sets. 791 The 231Doxorubicin training set consists of data from one experiment with a further, independent 792 experiment performed on a separate day used as a test set. Training and test sets for MCF7Docetaxel 793 are from the same two experiments, with random sampling used to produce independent training 794 and test sets. Each training data set contains a balanced number of untreated and treated cells, 795 treated with a single drug concentration. We selected $30\mu M$ docetaxel and $1\mu M$ doxorubicin for 796 the experiments with MDA-MB-231 cells as the optimal doses with which to induce changes in 797 cell morphology and migration without inducing cell death. However, a lower concentration $(1\mu M)$ 798 of docetaxel was used for MCF-7 cells as we found that this induced similar morphological and 799 dynamical changes to those induced by higher concentrations but with reduced cell death (Table 800 **3**). 801

	Data set	Cell line	Treatment	Training set	Test set
	231Docetaxel	MDA-MB-231	$30\mu M$ Docetaxel	Untreated: 646	Untreated: 913
				Treated: 600	Treated: 300
12	231Doxorubicin	MDA-MB-231	$1\mu M$ Doxorubicin	Untreated: 213	Untreated: 191
				Treated: 215	Treated: 60
	MCF7Docetaxel	MCF-7	$1\mu M$ Docetaxel	Untreated: 200	Untreated: 441
				Treated: 200	Treated: 128

80

Table 3: The three data sets used in this study with the number of cells in training and test sets used for203untreated vs treated classification.