THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

# Machine Learning Enables Live Label-Free Phenotypic Screening in Three Dimensions

**Citation for published version:**
O'Duibhir, E, Paris, J, Lawson, H, Pires Sepulveda, C, Doughty Shenton, D, Carragher, N & Kranc, K 2018, 'Machine Learning Enables Live Label-Free Phenotypic Screening in Three Dimensions' Assay and Drug Development Technologies, vol. 16, no. 1. DOI: 10.1089/adt.2017.819

**Digital Object Identifier (DOI):**
10.1089/adt.2017.819

**Link:**
Link to publication record in Edinburgh Research Explorer

**Document Version:**
Peer reviewed version

**Published In:**
Assay and Drug Development Technologies

# Machine Learning Enables Live Label-Free Phenotypic Screening in 3D

| Journal: | *ASSAY and Drug Development Technologies* |
|---|---|
| Manuscript ID | ADT-2017-819.R1 |
| Manuscript Type: | SBI2 Special Issue |
| Date Submitted by the Author: | 01-Dec-2017 |
| Complete List of Authors: | O'Duibhir, Eoghan; University of Edinburgh , Centre for Regenerative Medicine<br>Paris, Jasmin; University of Edinburgh , Centre for Regenerative Medicine<br>Lawson, Hannah; University of Edinburgh , Centre for Regenerative Medicine<br>Sepulveda, Catarina; University of Edinburgh , Centre for Regenerative Medicine<br>Doughty Shenton, Dahlia; University of Edinburgh, Edinburgh Phenotypic Assay Centre, The Queen's Medical Research Institute<br>Carragher, Neil; University of edinburgh, Edinburgh Cancer Research UK Centre<br>Kranc, Kamil; University of Edinburgh , Centre for Regenerative Medicine; University of edinburgh, Edinburgh Cancer Research UK Centre |
| Keyword: | Computational, Imaging, Screening, Cell-based |
| Manuscript Keywords (Search Terms): | Machine Learning, Leukaemia, 3D, Epigenetic, Phenotypic, High Content |
| | |

**SCHOLARONE™**
Manuscripts

**Machine Learning Enables Live Label-Free Phenotypic Screening in 3D**

Eoghan O'Duibhir[1], Jasmin Paris[1], Hannah Lawson[1], Catarina Sepulveda[1], Dahlia Doughty Shenton[2], Neil O. Carragher[3] and Kamil R. Kranc[1,3]

1. Centre for Regenerative Medicine, University of Edinburgh.

2. Edinburgh Phenotypic Assay Centre, The Queen's Medical Research Institute, University of Edinburgh.

3. Cancer Research UK Edinburgh Centre, Institute of Genetics and Molecular Medicine, University of Edinburgh.

Address correspondence to: Eoghan O'Duibhir

Kamil R Kranc and Neil Carragher contributed equally to this work.

Email: eoduibh@staffmail.ed.ac.uk, Jasmin.Paris@ed.ac.uk, Hannah.Lawson@ed.ac.uk, C.Sepulveda@sms.ed.ac.uk, D.Shenton@ed.ac.uk, N.Carragher@ed.ac.uk, Kamil.Kranc@ed.ac.uk

**Keywords**

Machine Learning, Leukaemia, 3D, Epigenetic, Phenotypic, High Content

**Abstract**

There is a large amount of information in brightfield images that was previously inaccessible using traditional microscopy techniques. This information can now be exploited using machine learning approaches for both image segmentation and the classification of objects. We have combined these approaches with a label-free assay for growth and differentiation of leukemic colonies, to generate a novel platform for phenotypic drug discovery. Initially a supervised machine learning algorithm was used to identify in-focus colonies growing in a 3D methylcellulose gel. Once identified, unsupervised clustering and principle component analysis of texture based phenotypic profiles were applied to ~~identify novel~~group similar phenotypes. In a proof of concept study we successfully identified a novel phenotype induced by a compound that is currently in clinical trials for the treatment of leukaemia. We believe that our platform will be of great benefit for the utilization of patient-derived 3D cell culture systems for both drug discovery and diagnostic applications.

**Disclosure Statement**

No competing financial interests exist.

**Abbreviations**

3D          Three dimensional

AML         Acute myeloid leukaemia

BET         Bromodomain and extraterminal domain

BF          Brightfield

CFC         Colony forming cell

DMSO        Dimethyl sulfoxide

GFP         Green Fluorescent Protein

H3          Histone three

IMDM        Iscove's Modified Dulbecco's Medium

LSC         Leukemic stem cell

MLL         Mixed lineage leukaemia

PCA         Principle component analysis

**Introduction**

As a model disease for understanding cancer biology, leukaemia has been exceptionally revealing [1].

Leukemic stem cells (LSCs) driving acute myeloid leukaemia (AML) were the first described cancer

stem cell [2], ultimately leading to the more generalized 'cancer-stem-cell hypothesis'. Various

translocations involving the mixed lineage leukaemia (*MLL*) gene lead to multiple haematological

malignancies, including AML, and are often associated with a poor prognosis. MLL is a DNA-binding

protein and epigenetic regulator that methylates histone H3 lysine 4 [3]. When present as a

leukaemogenic fusion protein MLL has been shown to bind to the promoters of the *Hoxa9* and *Meis1*

genes and ~~promote~~ be associated with histone modification [4]. ~~–~~ When grown *in vitro*, LSC colonies

display graded phenotypes depending on the initiating mutation~~–~~ [5,6]. Looser colonies are surrounded

by a spectrum of more differentiated blast-like cells, while denser colonies contain more

undifferentiated cells [7]. These phenotypes are potentially clinically relevant as it has been shown that

colony morphology is correlated with the disease prognosis in mice [6]. Because the phenotype is easily

visualized, it is possible to use image based screening to identify agents that can drive leukaemic cells

towards a more benign, differentiated phenotype. We have developed a method for high-throughput,

high-content screening of live colonies cultured and imaged in 3D. To validate the sensitivity of our

approach to variations in genetic background we performed a pilot screen in three different cell lines.

This allowed comparison of effects between human and mouse species and, in mouse, between

primary cells transformed by different oncogenes.

Colony formation assays are typically performed in 6-well plates and scored manually by a researcher.

After initial isolation, cells are mixed with cytokine-containing semi-solid methylcellulose-based media

formulated to promote leukaemic colony growth in three dimensions through proliferation and

differentiation [8]. The methylcellulose colony forming cell (CFC) assay [9], is a preferred *in vitro* assay

used in the study of primitive hematopoietic cells, and cells can readily be recovered from

methylcellulose for further phenotypic and molecular characterization. Due to observed auto-

fluorescence of the growth gel (the methylcellulose scaffold and growth media mix), direct fluorescent imaging of GFP expressing cell colonies *in situ* could notcannot be utilized for our growth conditions. These colony forming assays are therefore low throughput, susceptible to bias due to manual scoring and generally unsuitable for arrayed chemical or genetic screening. Being able to employ these 3D assays for automated high throughput screening of peturbagens would clearly be advantageous, in both probing for mechanistic insights relating to disease biology and unearthing new therapeutic agents. In addition, the ability to perform high content screening for agents that are not simply preventing colony growth toxic but could drive colonies from a dense to loose phenotype would have added utility for drug discovery [10].

Brightfield (BF) images contain rich texture information which, until recently, was inaccessible to automated image analysis [11–13]. BF imaging of live cells also has several advantages over fluorescent imaging. Being label-free, there is no need to modify the cells with either a fluorescent protein expression cassette or the addition of dyes that could perturb normal cell function. Quantification of label-free BF images of colonies *in situ* would also support both short- and long-term live cell kinetic studies. We have previously been successful in developing a simple machine learning based analysis pipeline that could determine colony number and size from BF images [14]. Here, we investigate whether a similar approach could be employed in a screening campaign, not only to count and size colonies, but additionally to use the texture information to phenotypically profile colonies and potentially identify compounds that can induce novel phenotypes.

**Materials and Methods**

See also **table 1** for a summary of the screen protocol

**Colony Culture**

THP-1 cells were cultured at 500,000 cells/ml in RPMI-1640 GlutaMAX containing 10% FBS, 100 U/ml penicillin, and 100 μg/ml streptomycin.

Formatted: Font: Not Bold
Formatted: Font: Not Bold

MMA (*MLL-AF9*[KI/+] cells): foetal liver haematopoietic cells were extracted from a E14.5 *MLL-AF9*[KI/+]

embryo (*MLL-AF9*[KI/+] mice [15] were obtained from The Jackson Laboratory). After c-Kit enrichment using

MACS LS columns (Miltenyi Biotec), cells were serially replated every 6 d in MethoCult M3231

(STEMCELL Technologies) supplemented with 20ng/ml SCF, 10 ng/ml IL-3, 10 ng/ml IL-6 and 10 ng/ml

GM-CSF. After 3 rounds of plating, cells were cultured at 300,000 cells/ml in IMDM containing 10%

FBS, 100 U/ml penicillin, and 100 µg/ml streptomycin, supplemented with SCF, IL-3, and IL-6.

MMH (*Meis1/Hoxa9* cells): foetal liver haematopoietic cells were extracted from a E14.5 C57Bl/6

embryo. Following c-Kit enrichment using MACS LS columns (Miltenyi Biotec), cells were transduced

with MSCV-Meis1a-puro and MSCV-Hoxa9-neo retroviruses as per [14]. Following selection for

puromycin/neomycin co-resistance, cells were serially replated every 6 days in MethoCult M3231

(STEMCELL Technologies) supplemented with 20 ng/ml SCF, 10 ng/ml IL-3, 10 ng/ml IL-6 and 10 ng/ml

GM-CSF. After 3 rounds of plating, cells were cultured at 200,000 cells/ml in IMDM containing 10%

FBS, 100 U/ml penicillin, and 100 µg/ml streptomycin, supplemented with SCF, IL-3, and IL-6.

Animal experimentation complied with local and national requirements (UK Animals Act 1986)

For methylcellulose medium, 20 ml IMDM (Life Technologies) was added to 80 ml MethoCult 3231

(STEMCELL Technologies, Catalog #03231), vortexed, and allowed to settle. For primary murine cell

lines, the methylcellulose was supplemented with cytokines 20 ng/ml SCF, 10 ng/ml IL-3, 10 ng/ml IL-6

and 10 ng/ml GM-CSF. No antibiotics were added. Cells (THP-1 cells, *MLL-AF9*[KI/+] foetal liver cells, or

murine foetal liver transformed with *Meis1* and *Hoxa9* retroviruses) were suspended in IMDM and

added to the prepared methylcellulose at a ratio of 1:9. The mixture was vortexed and allowed to

settle. Compounds were added as a single dose. 5 µl of 2.1% test drug compound was pipetted into

the centre of each well of a 96-well non-tissue culture treated edge plate (Thermo Scientific, Cat. #

267313) with a CyBio FeLix. Subsequently, 100 µl of pre-mixed methylcellulose containing 400 cells

(THP-1) or 600 cells (*MLL-AF9*[KI/+] foetal liver cells, or murine foetal liver transformed with *Meis1* and

*Hoxa9* retroviruses) was syringed into each well (using BD Microlance 3 18 Gauge 1.5" needles,

resultant ~~drug~~ compound concentration 0.1%). The plate was vortexed, and the side troughs and unused wells were half filled with PBS (Sigma) to prevent edge effects due to uneven evaporation. Plates were incubated at 37°C 5% $CO_2$ (day 0), and ~~then~~ scanned on day 6 (murine cells), or day 9 (THP-1 cells).

**Imaging,~~ image and data analysis~~**

Images were acquired at 37°C 5% $CO_2$ on an Operetta high content microscope (Perkin Elmer) equipped with a live cell chamber. The imaging pattern for plates consisted of a snaking pattern across columns beginning with the top left gel containing well (B2), down to B7, across to C7 up to C2 and so on. In each well the imaging pattern began with the middle field and followed a snaking pattern beginning at the top left field, across rows and avoiding imaging of the central field twice. We choose 9 fields of view to maximise well coverage at 10 X magnification while avoiding the well edges. The edge of each of the wells had a texture that the algorithm sometimes identified as a colony and was therefore best to avoid. After testing various z-stack options during assay development focal planes separated by 150 µm were chosen to avoid repeated counting of the same colonies. Above a height of 600 µm there were no colonies found and plate scan times were unnecessarily increased.

**Image and numerical data analysis**

Image and subsequent numerical analysis was performed using a variety of software tools:

~~Image analysis was performed in~~ Columbus 2.7.1 (Perkin Elmer) was used for the initial image analysis step by manually training the "Find texture region" PhenoLogic machine learning ~~module~~ module to find two classes of texture regions in brightfield images. One class contained in-focus colonies (texture A) and the other class contained background and out of focus colonies (texture B). Texture A was ~~then~~ split into discrete objects, the outer border was shrunk by 6 pixels and any holes were filled. Objects greater than 2000 µm² were ~~then~~ considered as colonies and morphology and texture properties were calculated using the "Calculate morphology properties" and "Calculate texture properties" modules.

[Formatted: Underline]

Well level aggregated data and ~~Colony data~~ data for individual colonies including morphology and texture features ~~was~~ exported as separate text files.

~~subsequently analysed in~~ Spotfire HCP 7.5.0 (Perkin Elmer informatics) http://www.cambridgesoft.com was used for rapid initial visualization of the colony count data as plate heatmaps at colony and well level and scatterplots at well level for quality control purposes. Wells were tagged for positive and negative controls, compounds and concentration added. Hierarchical clustering of aggregated well level data and ~~(P~~principle ~~C~~component ~~A~~analysis [16]~~)~~ was performed using the built in HCP tools in the software. Principal components and tagged data at the well level were exported as text files for further plotting in Python.

**Formatted:** Underline

~~,~~ HC StratoMineR~~,~~ (Core Life Analytics) www.corelifeanalytics.com was used for ~~(for~~ hit calling of well level data based solely on colony number [17]~~)~~. All p-values were calculated using the z-test based on negative controls with a median estimator with a p-value of <0.0001 considered significant.

**Formatted:** Underline

~~and~~ ~~2~~Python ~~www.python.org~~ www.python.org was used for plotting of all data, except dose response curves. Although not necessarily required for the analysis Python was used so as to maintain consistent formatting of figures across the manuscript figures. ~~(all plotting,~~ Python was also used to calculate the Z-score normalization and perform the hierarchal clustering shown in figure 4 with~~:~~ sns.clustermap, method='average', metric='cosine'~~)~~.

**Formatted:** Underline

~~All p-values were calculated using the z-test with a p-value of <0.0001 considered significant.~~

**Results**

**Supervised machine learning-based segmentation of colonies in three dimensions.**

The following automated image acquisition parameters were developed to enable optimal label-free imaging of colonies grown in a 96-well plate while avoiding common pitfalls of assay miniaturization. The imaging pattern avoided issues with both imaging the well wall (**figure 1a**) and identifying the

same colony in more than one focal plane (**figure 1b**). Due to their relatively larger size, the number of

objects per well of a 96-well plate is limited when measuring colonies rather than cells. To maximise

image coverage while minimising the time taken for imaging each plate, we employed a 10 X

objective. This resulted in flatter illumination across fields than the 2 X lens but did result in more

colonies that were clipped by the edge of the field (**figure 1c**). Nine fields of view were imaged in each

well of the 96-well assay plate (**figure 1a**) covering approximately 50% of the well, with each field

acquired at five focal planes each separated by 150 μm (**figure 1b**). All images were subsequently

segmented using an algorithm (supervised texture segmentation module in the Columbus image

analysis software) that had previously been trained on an independent training set [14]. We tested the

algorithm on three independent cell lines: a human AML (M5) cell line harboring a MLL-AF9

translocation (THP-1 cells); cells obtained from a mouse (*MLL-AF9$^{KI/+}$*) with a genomic rearrangement

leading to expression of the MLL-AF9 fusion protein (further referred to as MMA cells); and a primary

mouse cell line containing retroviral constructs that overexpress *Meis1* and *Hoxa9* (further referred to

as MMH cells), each of which display differences in size and number of colonies. Upon visual

inspection the segmentation algorithm performed equally well in identifying colonies grown from each

cell line (**figure 1d**-**f**). As a positive control for compound addition to each plate we used iBET [18], a

known inhibitor of leukemic cell growth and colony formation [19]. In our assay, iBET proved effective at

inhibiting the growth of all three cell lines (**figure 1g**-**i**).

**Epigenetic tool compound library**

Abnormal epigenetic regulation of gene expression has been implicated as potentially causative in

several types of myeloid malignancies [20]. We therefore employed the high quality epigenetic tool

compound library from the Structural Genomics Consortium (SGC) [21] to map which epigenetic

regulators are involved in colony growth and differentiation across the three different leukaemic cell

lines. The compounds used are listed in **table 2**, along with their plate location and known targets. A

six point dose response was performed starting at 10 μM with a 1 in 5 dilution at each step (giving: 10

µM; 2 µM; 400 nM; 80 nM; 16 nM; and 3.2 nM). Although SGC do not recommend using their

compounds at concentrations higher than 1 µM we had previously observed that in semi-solid

methylcellulose medium our positive control iBET was only effective at concentrations approximately

10 fold higher than in liquid culture (unpublished data). We therefore began the dose response at 10

µM. A ~~simple visual schematic~~ summary of the screening ~~experimental design is provided~~ protocol is

shown in **table 1**, with more detailed procedures in~~in~~ the materials and methods section.

**Digitized colonies: size, number and location**

There was almost complete ablation of colonies in the positive control wells for each cell line (example

plates shown in **figure 2a-c**, with iBET added to first 4 wells of rows 2 and last 3 wells of row 11).

Compounds ~~displaying toxicity~~ ablating colony formation in all three cell lines are also plainly visible

(**figure 2a-c**) at the highest concentration used (10 µM). At this concentration the lack of colonies is

most likely due to toxicity due to the complete lack of cells found after manual inspection of the full

resolution images. Colony location and size are clearly recapitulated by the segmentation algorithm

(**figure 2d-f**). Visualizing the performance of the algorithm as an entire digital plate gave added

confidence of accurate measurement of colony number and size.

Quantification of total number of colonies across all plates in the screen shows several compounds to

~~be toxic~~ reduce CFC number at lower concentrations (**figure 3a**). There are no obvious edge effects on

colony size or number in the outer wells of the plate. There appears to be a general reduction in CFC

numbers, possibly due to a general~~ly~~ toxic effect of the compounds at the highest~~r~~ concentrations~~,~~

most apparent in the MMH cell line at 10 µM (**figure 3b**). Surprisingly there is also a single compound

(GSK-LSD1) that increases colony number across a range of concentrations (**figure 3a** and effect size

shown in **3b**). Z-prime (Z') scores based on colony number are excellent for THP-1 (0.57) and for MMH

(0.54) cell lines but only -0.52 for the MMA cell line (calculated on 42 positive and 60 negative wells

spread across 6 plates for each cell line). The reduced Z' for this primary cell line is due to increased

overall noise in the measurements because of 1) the lower colony numbers leading to reduced

number of colonies quantified, and 2) the greatly increased colony size which results in more frequent

colony clipping. This is also reflected in called hits based on a reduction in colony number. THP-1 and

MMH cell lines have almost perfect ~~toxic~~ hit overlap for reduction of colony numbers (**table ~~2~~3**, all

with a p-value < 0.0001 and dose response curves for overlapping compounds in **supplemental figure**

**1**). Most of the hits are at the 10 µM concentration. If compounds with a potency below 10 µM are

considered, only LAQ824 and JQ1 remain. JQ1 is clearly potent down to 2 µM with an IC50 of 1.6 µM

for THP-1 derived colonies and 0.9 µM for MMH derived colonies. ~~and~~JQ1 has a similar chemical

structure to iBET [22], also targeting bromodomains. Far more potent however is LAQ824, ~~killing colonies~~

~~down to 80 nM in both species~~with IC50s of 65 nM for THP-1 and 20 nM for MMH derived colonies.

The MMA cell line displayed no statistically significant hits at any concentration.

**Unsupervised clustering and PCA analysis identify novel colony phenotypes**

Although we had discovered clear ~~toxic~~ hits based on a reduction in colony number, ultimately our

goal was to find compounds which induce differentiation within the leukemic colonies, ideally

resulting in a less aggressive clinical phenotype and potentially having more specificity (with fewer side

effects than a toxic compound that indiscriminately kills proliferating stem cells). To this end we

performed morphology and texture analysis to give 21 further parameters describing each colony

(examples in **figure 4a**). Well level data for the entire screen was ~~then~~ further analysed using~~with a~~

hierarchical~~n~~ ~~unsupervised~~ clustering ~~algorithm~~ (**figure 4b**). Wells containing colonies from the same

cell line largely cluster together, demonstrating a specific morphology profile for colonies derived from

each cell type. Where there is intermingling of profiles from different cell lines, most of these wells

had been treated either with the iBET positive control (green) or a compound that ~~had a toxic effect~~

~~at~~reduced colony number at a particular dose (red). After treatment with a ~~toxic~~ compound that

affects colony number, wells containing affected colonies cluster together, rather than with their own

genotype. This indicates that the phenotypic effect elicited by the compound is stronger than the

original phenotypic similarity due to the genetics of each cell line.

To investigate the presence of potentially novel phenotypes, colony morphology and texture was

further analysed by principle component analysis (PCA). PCA was applied to the entire dataset,

containing all cell lines and compound concentrations. The first three principal components (PC1, 2

and 3) respectively capture 48%, 16% and 12% of the variance in the data. In this PCA space a clear

separation of positive (green) and negative (blue) controls can be seen, particularlyespecially for the

THP-1– and MMH cell lines (**figure 5 a** and **c**). This separation is not as clear for the MMA derived

colonies (**figure 5b**). In all cases the majority of compounds (yellow) are found clustering together with

the DMSO controls having no effect. Many compounds are found in the same space as the positive

controls (group **i** in **figure 5 a-c**). These compounds overlap exactly with the toxic hits based on a

reduction in colony number (LAQ824, PFI-1, JQ1, GSK J4, NVS-1, OLAPARIB, Bromosporine and CL994

in both THP-1 and MMH cell lines). As was the case for colony number, when only considering

compounds at concentrations less than 10 µM, we are again left with JQ1 and LAQ824 and in the case

of the THP-1 cell line also PFI-1. Most interestingly a single compound, GSK-LSD1 (at concentrations

ranging from 10 µM to 16 nM) occupies PCA space orthogonal to the positive and negative controls

(**figure 5 a** and **c**, group **ii**), and was not previously called as a hit based on a reduction in colony

number. Visual inspection of this phenotype shows colonies that have differentiated into single cells.

**Discussion**

Due to the high failure rate in target based drug discovery approaches [23] there is a need for renewed emphasis on phenotypic based approaches [24] that recognise the complexity of the biology involved [10]. Recent advances in imaging, cell culture and genetic engineering technologies [25], combined with advances in machine learning [26,27] are converging to facilitate a high throughput renaissance in empirical drug discovery using more complex and relevant cell-based models of disease. Here, we present a simple image based screening methodology that relies on a complex but commercially available analysis pipeline. Our objective was not to come as close as possible to ground truth measurements or improve the error rate of manual counting, but to Our aim was be able to increase assay throughput while readily quantifying a phenotypic difference. In this study we have used a machine learning approach to automate the quantification of, using a label-free 3D methylcellulose colony formation assay, to allow classification of compound activity identifying a novel based phenotype based on their induced morphological profiles.

BF is less perturbing and faster than fluorescent imaging in multiple channels and thus particularly well suited to complex live-cell kinetic and/or 3D assays. Combined with machine learning facilitated analysis, BF images provide a rich source of texture and morphology information that can be mined for novel phenotypes. Because our segmentation algorithm was texture rather than intensity based and trained specifically to only find in-focus colonies this meant we could screen in 3D and overcome the issues of uneven illumination across a well due to the gel meniscus. Furthermore, because BF imaging is label-free and permits live imaging with minimal genetic or chemical perturbation, the methods described here may be beneficial for personalised diagnostic applications using primary patient-derived cells. We have also used this approach to identify BF imaged liver organoids and in-focus cystic embryoid bodies grown in matrigel and stained with DAPI, followed by further nuclear segmentation (based on standard methods), estimation of relative cell numbers per cyst and classification of cells

based on fluorescent immunohistochemistry labelled markers (data unpublished). Thus, combining BF and fluorescent imaging can lead to even richer phenotypes in multiple tissue types and systems.

In order to identify and segment colonies in a brightfield image it is critical that the colonies do not overlap. Typical image analysis strategies for segmenting touching objects in fluorescent images include peak intensity and shape or the more recently developed approach by the Horvath lab [28] that includes assumptions about nuclear shape and additive pixel intensities of overlapping nuclei. These approaches cannot be employed here as the method for identifying the colonies is texture based. This is a limitation of our approach and necessitates a lower object density to avoid overlap.

During initial assay development we found it necessary to use non-tissue culture treated edge plates (Nunc Cat. # 267313) both to prevent colonies in contact with the bottom of the plate spreading over the plastic and to avoid what was obvious growth retardation in the outer wells, probably due to evaporation. As the number of compounds tested in this pilot screen allowed for only the inner 60 wells of each plate to be used this further avoided any edge effects. However, for scale up compound numbers it would be desirable to use all 96-wells in a plate. In this case use of the edge plates would be necessary.

MMA colonies did not display an orthogonal phenotype in PCA space when treated with GSK-LSD1. However, manual examination of GSK-LSD1 treated wells in this cell line reveals a similar differentiation effect but with a greatly reduced numbers of cells. These cells however had a curious elongated morphology (example seen in **figure 5b**, GSK-LSD1 at 400 nM). Because the cells were sparse they were not grouped as colonies by the algorithm and were lost during the size exclusion step after image segmentation. This compound has promise as a therapeutic agent, being potent down to 16 nM and producing the desired differentiation phenotype without an obvious toxic effect based on the continued presence of cells (and depending on genotype). Indeed GSK-LSD1 has been through phase I clinical trials to assess safety and activity in patients with relapsed AML (under the generic name GSK2879552, https://www.gsk-clinicalstudyregister.com/study/200200#ps). Other lysine

demethylase targeting inhibitors in the SGC set did not show same phenotype. These inhibitors target proteins other than LSD1 (see **table 2**), which has been identified as the target of GSK-LSD1 [29]. Another lysine demethylase identified as a ~~toxic~~ hit reducing colony number is GSK-J4. This compound targets the JMJD3, UTX and JARID1B proteins [30] and displays effects only at the highest concentration (10 µM) in our assay. This difference between compounds targeting separate lysine demethylases could be mechanistically informative, pointing to a specific differentiating effect upon LSD1 inhibition. Although only showing ~~toxic~~ a reduction in colony formation rather than purely differentiation effects in this assay, LAQ824 has also been used in a phase I clinical trial for patients with advanced solid tumours [31] and has shown activity against myeloma [32] and human acute leukaemia [33].

Future scale up of this screening method would require development of a pipetting head and automation platform capable of dispensing large amounts of methylcellulose gel containing cells. The current analysis pipeline holds enormous potential for repurposing to a variety of other 3D assay formats. We expect that future use of machine learning to analyse label-free images will aid in the identification of novel leads to treat a variety of diseases and in their initial diagnosis.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

**References**

1. Greaves, M. Leukaemia 'firsts' in cancer research and treatment. *Nat. Rev. Cancer* **16,** 163–172 (2016).

2. Huntly, B. J. P. & Gilliland, D. G. Leukaemia stem cells and the evolution of cancer-stem-cell research. *Nat. Rev. Cancer* **5,** 311–321 (2005).

3. Patel, A., Dharmarajan, V., Vought, V. E. & Cosgrove, M. S. On the mechanism of multiple lysine methylation by the human mixed lineage leukemia protein-1 (MLL1) core complex. *J. Biol. Chem.* **284,** 24242–24256 (2009).

4. Milne, T. A., Martin, M. E., Brock, H. W., Slany, R. K. & Hess, J. L. Leukemogenic MLL fusion proteins bind across a broad region of the Hox a9 locus, promoting transcription and multiple histone modifications. *Cancer Res.* **65,** 11367–11374 (2005).

5. Giustacchini, A. *et al.* Single-cell transcriptomics uncovers distinct molecular signatures of stem cells in chronic myeloid leukemia. *Nat. Med.* **23,** 692–702 (2017).

6. Somervaille, T. C. P. *et al.* Hierarchical maintenance of MLL myeloid leukemia stem cells employs a transcriptional program shared with embryonic rather than adult stem cells. *Cell Stem Cell* **4,** 129–140 (2009).

7. Lavau, C., Szilvassy, S. J., Slany, R. & Cleary, M. L. Immortalization and leukemic transformation of a myelomonocytic precursor by retrovirally transduced HRX-ENL. *EMBO J.* **16,** 4226–4237 (1997).

8. Borowicz, S. *et al.* The soft agar colony formation assay. *J. Vis. Exp. JoVE* e51998 (2014). doi:10.3791/51998

9. Sarma, N. J., Takeda, A. & Yaseen, N. R. Colony forming cell (CFC) assay for human hematopoietic cells. *J. Vis. Exp. JoVE* (2010). doi:10.3791/2195

10. Horvath, P. *et al.* Screening out irrelevant cell-based models of disease. *Nat. Rev. Drug Discov.* **15,** 751–769 (2016).

11. Blasi, T. *et al.* Label-free cell cycle analysis for high-throughput imaging flow cytometry. *Nat. Commun.* **7,** ncomms10256 (2016).

12. Buggenthin, F. *et al.* An automatic method for robust and fast cell detection in bright field images from high-throughput microscopy. *BMC Bioinformatics* **14,** 297 (2013).

13. Kraus, O. Z. *et al.* Automated analysis of high-content microscopy data with deep learning. *Mol. Syst. Biol.* **13,** 924 (2017).

14. Vukovic, M. *et al.* Hif-1α and Hif-2α synergize to suppress AML development but are dispensable for disease maintenance. *J. Exp. Med.* **212,** 2223–2234 (2015).

15. Chen, W. *et al.* Malignant transformation initiated by Mll-AF9: Gene dosage and critical target cells. *Cancer Cell* **13,** 432–440 (2008).

16. Bro, R., Acar, E. & Kolda, T. G. Resolving the sign ambiguity in the singular value decomposition. *J. Chemom.* **22,** 135–140 (2008).

17. Omta, W. A. *et al.* HC StratoMineR: A Web-Based Tool for the Rapid Analysis of High-Content Datasets. *Assay Drug Dev. Technol.* **14,** 439–452 (2016).

18. Nicodeme, E. *et al.* Suppression of inflammation by a synthetic histone mimic. *Nature* **468,** 1119–1123 (2010).

19. Dawson, M. A. *et al.* Inhibition of BET recruitment to chromatin as an effective treatment for MLL-fusion leukaemia. *Nature* **478,** 529–533 (2011).

20. Fong, C. Y., Morison, J. & Dawson, M. A. Epigenetics in the hematologic malignancies. *Haematologica* **99,** 1772–1783 (2014).

21. Brown, P. J. & Muller, S. Open access chemical probes for epigenetic targets. *Future Med. Chem.* **7,** 1901–1917 (2015).

22. Delmore, J. E. *et al.* BET Bromodomain Inhibition as a Therapeutic Strategy to Target c-Myc. *Cell* **146,** 904–917 (2011).

23. Hay, M., Thomas, D. W., Craighead, J. L., Economides, C. & Rosenthal, J. Clinical development success rates for investigational drugs. *Nat. Biotechnol.* **32,** 40–51 (2014).

24. Swinney, D. C. & Anthony, J. How were new medicines discovered? *Nat. Rev. Drug Discov.* **10,** 507–519 (2011).

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

25. O'Duibhir, E., Carragher, N. O. & Pollard, S. M. Accelerating glioblastoma drug discovery: Convergence of patient-derived models, genome editing and phenotypic screening. *Mol. Cell. Neurosci.* **80,** 198–207 (2017).

26. Grys, B. T. *et al.* Machine learning and computer vision approaches for phenotypic profiling. *J Cell Biol* **216,** 65–71 (2017).

27. Sommer, C. & Gerlich, D. W. Machine learning in cell biology - teaching computers to recognize phenotypes. *J. Cell Sci.* **126,** 5529–5539 (2013).

28. Molnar, C. *et al.* Accurate Morphology Preserving Segmentation of Overlapping Cells based on Active Contours. *Sci. Rep.* **6,** srep32412 (2016).

29. Mohammad, H. P. *et al.* A DNA Hypomethylation Signature Predicts Antitumor Activity of LSD1 Inhibitors in SCLC. *Cancer Cell* **28,** 57–69 (2015).

30. Kruidenier, L. *et al.* A selective jumonji H3K27 demethylase inhibitor modulates the proinflammatory macrophage response. *Nature* **488,** 404–408 (2012).

31. de Bono, J. S. *et al.* Phase I pharmacokinetic and pharmacodynamic study of LAQ824, a hydroxamate histone deacetylase inhibitor with a heat shock protein-90 inhibitory profile, in patients with advanced solid tumors. *Clin. Cancer Res. Off. J. Am. Assoc. Cancer Res.* **14,** 6663–6673 (2008).

32. Catley, L. *et al.* NVP-LAQ824 is a potent novel histone deacetylase inhibitor with significant activity against multiple myeloma. *Blood* **102,** 2615–2622 (2003).

33. Guo, F. *et al.* Cotreatment with Histone Deacetylase Inhibitor LAQ824 Enhances Apo-2L/Tumor Necrosis Factor-Related Apoptosis Inducing Ligand-Induced Death Inducing Signaling Complex Activity and Apoptosis of Human Acute Leukemia Cells. *Cancer Res.* **64,** 2580–2589 (2004).

**Figure Legends**

*Figure 1. Imaging strategy*. *Example of brightfield (BF) images showing: (**a**) approximate well coverage of nine tiled BF images avoiding well wall; (**b**) an example of a single stack both pre- and post-image processing; (**c**) even illumination and varied colony morphology; (**d**-**f**) performance of the algorithm throughout the gel for each of the cell lines (all images taken from top left field of DMSO negative control at the same plate location, well F2); (**g**-**i**) action of positive control on colony growth of each genotype (9 tiled images shown per cell line, all images taken from plane 1 in either well F2 for DMSO or C2 for 10 µM iBET positive control).*

*Figure 2. Digitisation of colonies*. *Tiled BF images showing plane 1 of an entire plate at the highest compound concentration for each cell line (**a**-**c**) and the performance of the algorithm across the entire plate shown as scatterplots (**d**-**f**). Row and column numbers are relative to well position in a 96-well plate.*

*Figure 3. Colony numbers across entire screen*. *Heatmaps showing effect of compounds while maintaining positional information for each plate (**a**) and the same data displayed as scatterplots (**b**) more clearly displaying the effect size. Data were normalized to the median DMSO value for each cell line.*

*Figure 4. Hierarchical Clustering of morphological phenotypes*. *An example brightfield image with segmentation and representations of the spot, edge and ridge texture features (**a**). Clustered heatmap of Z-score normalized profiling data (**b**). Wells containing each cell line are marked pink, dark grey or yellow. Compounds are marked in green for iBET positive control, red for ~~toxic~~ hits reducing colony number (as per table 3) and the remaining compounds are white. Empty attribute values (coming from wells with no colonies to profile) are light grey.*

*Figure 5. Orthogonal phenotype in PCA space*. *Three-dimensional scatter plots of first three principle components, plotted for each genotype (**a**-**c**) with example brightfield images directly below each plot.*

*Supplemental Figure 1. Dose response curves.*

*Dose response curves are shown for all overlapping compounds that significantly reduce colony number (as per table 3). A line shows a logistic regression curve was fitted to data for each compound and cell line. Single data points for each concentration without replicates are shown as circles with inflection points, corresponding to the IC50, shown as triangles.*

**Formatted:** Font: Not Bold

**Formatted:** Font: Not Bold

**Formatted:** Font: Not Bold

**Formatted:** Font: Not Bold

**Formatted:** Font: Not Bold

**Figure 1**



Figure 1. Imaging strategy. Example of brightfield (BF) images showing: (a) approximate well coverage of nine tiled BF images avoiding well wall; (b) an example of a single stack both pre- and post-image processing; (c) even illumination and varied colony morphology; (d-f) performance of the algorithm throughout the gel for each of the cell lines (all images taken from top left field of DMSO negative control at the same plate location, well F2); (g-i) action of positive control on colony growth of each genotype (9 tiled images shown per cell line, all images taken from plane 1 in either well F2 for DMSO or C2 for iBET positive control).

171x166mm (300 x 300 DPI)

Figure 2. Digitisation of colonies. Tiled BF images showing plane 1 of an entire plate at the highest compound concentration for each cell line (a-c) and the performance of the algorithm across the entire plate shown as scatterplots (d-f). Row and column numbers are relative to well position in a 96-well plate.

147x120mm (300 x 300 DPI)

**Figure 3**

a.



b.



Figure 3. Colony numbers across entire screen. Heatmaps showing effect of compounds while maintaining positional information for each plate (a) and the same data displayed as scatterplots (b) more clearly displaying the effect size. Data were normalized to the median DMSO value for each cell line.

160x136mm (300 x 300 DPI)

**Figure 4**

a.



b.



Figure 4. Hierarchical Clustering of morphological phenotypes. An example brightfield image with segmentation and representations of the spot, edge and ridge texture features (a). Clustered heatmap of Z-score normalized profiling data (b). Wells containing each cell line are marked pink, dark grey or yellow. Compounds are marked in green for iBET positive control, red for hits reducing colony number (as per table 3) and the remaining compounds are white. Empty attribute values (coming from wells with no colonies to profile) are light grey.

214x252mm (300 x 300 DPI)

**Figure 5**



Figure 5. Orthogonal phenotype in PCA space. Three-dimensional scatter plots of first three principle components, plotted for each genotype (a-c) with example brightfield images directly below each plot.

193x210mm (300 x 300 DPI)

**Supplemental figure 1**



Supplemental Figure 1. Dose response curves are shown for all overlapping compounds that significantly reduce colony number (as per table 3). A line shows a logistic regression curve was fitted to data for each compound and cell line. Single data points for each concentration without replicates are shown as circles with inflection points, corresponding to the IC50, shown as triangles.

199x249mm (300 x 300 DPI)

| Table 1: Protocol Table | | | |
|---|---|---|---|
| **Step** | **Parameter** | **Value** | **Description** |
| 1 | Compound addition | 5 µl /well | To empty 96 well plate, 2.1% DMSO |
| 2 | Mix cells and semi-solid media | 20 ml/cell line | 4000 cells/ml for human, 6000 cells/ml for mouse |
| 3 | Add cell mix to plates | 100 µl/well | Manually with syringe |
| 4 | Vortex | 5 seconds | |
| 5 | Incubation | 6 - 9 days | 6 days for mouse, 9 days for human |
| 6 | Imaging | 30 ms/field | BF, 37°C and 5% $CO_2$ |
| 7 | Image analysis | PhenoLogic module | Columbus image analysis server |
| 8 | Data analysis | Well level | Hierarchical clustering and PCA |
| **Step** | **Notes** | | |
| 1 | CyBio FeliX, non-tissue culture treated edge plate | | |
| 2 | Media pre-warmed to 37°C | | |
| 3 | Side trough and unused wells half filled with PBS | | |
| 4 | Ensures mixing of compound with media | | |
| 6 | Operetta microscope | | |
| 8 | With Spotfire HCP or HC Stratominer | | |

**Table 2**. Compounds used in this study.

| Compound | Row | Column | Protein Family | Specific Targets |
|---|---|---|---|---|
| iBET (positive control) | 2,3,4,5,5,6,7 | 2,2,2,2,11,11,11 | Bromodomains | BRD2, BRD3, BRD4, BRDT |
| GSK2801 | 2 | 10 | Bromodomains | BAZ2A, BAZ2B |
| BAZ2-ICR | 2 | 9 | Bromodomains | BAZ2A, BAZ2B |
| PFI-4 | 2 | 8 | Bromodomains | BRPF1B |
| JQ1 | 2 | 7 | Bromodomains | BRD2, BRD3, BRD4, BRDT |
| PFI-1 | 2 | 6 | Bromodomains | BRD2, BRD3, BRD4, BRDT |
| LP99 | 2 | 5 | Bromodomains | BRD9, BRD7 |
| BI-9564 | 2 | 4 | Bromodomains | BRD9, BRD7 |
| OF-1 | 2 | 3 | Bromodomains | BRPF1, BRPF2, BRPF3 |
| NI-57 | 3 | 10 | Bromodomains | BRPF1, BRPF2, BRPF3 |
| SGC-CBP30 | 3 | 9 | Bromodomains | CREBBP, EP300 |
| I-CBP112 | 3 | 8 | Bromodomains | CREBBP, EP300 |
| NVS-CECR2-1 | 3 | 7 | Bromodomains | CECR2 |
| IOX1 | 3 | 6 | Lysine demethylase | pan-2-OG |
| KDOAM25 | 3 | 5 | Lysine demethylase | KDM5 |
| SGC0946 | 3, 7 | 4, 9 | Methyltransferase | DOT1L |
| UNC1999 | 3 | 3 | Methyltransferase | EZH2 |
| GSK343 | 4 | 10 | Methyltransferase | EZH2 |
| UNC0638 | 4 | 9 | Methyltransferase | G9a, GLP |
| UNC0642 | 4 | 8 | Methyltransferase | G9a, GLP |
| A-366 | 4 | 7 | Methyltransferase | G9a, GLP |
| GSK-J4 | 4 | 6 | Lysine demethylase | JMJD3, UTX, JARID1B |
| UNC1215 | 4 | 5 | Methyl Lysine Binder | L3MBTL3 |
| GSK-LSD1 | 4 | 4 | Lysine demethylase | LSD1 |
| GSK484 | 4 | 3 | Arginine deiminases | PAD-4 |
| Bromosporine | 5 | 10 | Bromodomains | pan-Bromodomain |
| IOX2 | 5 | 9 | 2-oxoglutarate dependent oxygenases | PHD2 |
| SGC707 | 5 | 8 | Methyltransferase | PRMT3 |
| PFI-2 | 5 | 7 | Methyltransferase | SETD7 |
| PFI-3 | 5 | 6 | Bromodomains | SMARCA,PB1 |
| LLY-507 | 5 | 5 | Methyltransferase | SMYD2 |
| BAY-598 | 5 | 4 | Methyltransferase | SMYD2 |
| A-196 | 5 | 3 | Methyltransferase | SUV420H1/H2 |
| OICR-9429 | 6 | 10 | WD40 repeat | WDR5 |
| LAQ-824 | 6 | 9 | Histone deacetylases | - |
| OLAPARIB | 6 | 8 | DNA repair | PARP |
| C-646 | 6 | 7 | Histone acetyltransferases | p300/CBP |
| CL-994 | 6 | 6 | Histone deacetylases | HDAC 1, 2, 3, and 8 |

| IOX-2 | 6 | 5 | 2-oxoglutarate dependent oxygenases | PHD2 |
| I-BRD9 | 6 | 3 | Bromodomains | BRD9 |
| GSK-J1 | 7 | 3 | Lysine demethylase | JMJD3, UTX, JARID1B |
| DMSO (negative control) | 2,3,4, 6,7,7,7,7,7,7 | 11,11,11,2,2,4,5, 6,8,10 | - | - |

| **Table 3**. Hits based on a reduction in colony numbers (p<0.0001) | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **Cell line** | **Species** | **Onco-gene** | **10 uM** | **2 uM** | **400 nM** | **80 nM** | **16 nM** | **3.2 nM** |
| THP-1 | human | *MLL-AF9* | LAQ824, PFI-1, JQ1, GSK J4, NVS-CECR2-1, OLAPARIB, Bromosporine | LAQ824, JQ1 | LAQ824 | LAQ824 | - | - |
| MMA | mouse | *MLL-AF9* | - | - | - | - | - | - |
| MMH | mouse | *Meis1/ Hoxa9* | LAQ824, PFI-1, JQ1, GSK J4, NVS-CECR2-1, OLAPARIB, Bromosporine, CL994 | LAQ824, JQ1 | LAQ824 | LAQ824 | - | - |