1 **TITLE PAGE**

2 Inter-laboratory automation of the *in vitro* micronucleus assay using imaging flow cytometry and

3 deep learning

4

5 **Author information**:

6 John W. Wills*[1,2], Jatin R. Verma[3], Benjamin J. Rees[3], Danielle S. G. Harte[3], Qiellor Haxhiraj[3],

7 Claire M. Barnes[1], Rachel Barnes[3], Matthew A. Rodrigues[4], Minh Doan[5], Andrew Filby[6], Rachel E.

8 Hewitt[2], Catherine A. Thornton[3], James G. Cronin[3], Julia D. Kenny[7], Ruby Buckley[7], Anthony M.

9 Lynch[3,7], Anne E. Carpenter[8], Huw D. Summers[1], George Johnson*[#3], Paul Rees*[#1,8].

10

11 [#]These authors contributed equally.

12

13 *Joint-corresponding authors:

14 *John W. Wills       Email: jw2020@cam.ac.uk    +44 (0)1223 337701

15 *George E. Johnson      Email: g.johnson@swan.ac.uk   +44 (0)1792 295158

16 *Paul Rees             Email: p.rees@swan.ac.uk     +44 (0)1792 295197

17

18 **Affiliations:**

19 [1]College of Engineering, Swansea University, Swansea, UK.

20 [2]Department of Veterinary Medicine, Cambridge University, Cambridge, UK.

21 [3]Swansea University Medical School, Swansea University, Swansea, UK.

22 [4]Amnis Flow Cytometry, Luminex Corporation, Seattle, Washington, USA.

23 [5]Bioimaging Analytics, GlaxoSmithKline, Collegeville, USA.

24 [6]Faculty of Medical Sciences, Newcastle University, Newcastle upon Tyne, UK.

25 [7]GlaxoSmithKline Research and Development Platform, Ware, UK.

26 [8]Imaging Platform, Broad Institute of MIT and Harvard, Cambridge, Massachusetts, USA.

**DECLARATIONS**

**Acknowledgements**

**Funding**

**Conflicts of interest/competing interests**

M. A. R., is an employee of Luminex Corporation which manufactures the Amnis ImageStream imaging flow cytometers used in this research study.

**Availability of data and materials**

Imaging flow cytometry test images alongside the final DeepFlow neural network are provided for download from the BioStudies database (http://www.ebi.ac.uk/biostudies) under accession number S-BSST641.

**Code availability**

The presented deep learning image analysis pipeline is available for download at the BioStudies database (http://www.ebi.ac.uk/biostudies) in MATLAB and Python programming languages under accession number S-BSST641.

2

53

## **Ethics approval**

55    This study uses *in vitro* cell lines only. No ethical approval was required.

56

## **Consent to participate**

58    Not applicable

59

## **Consent for publication**

61    Not applicable

**ABSTRACT**

The *in vitro* micronucleus assay is a globally significant method for DNA damage quantification used for regulatory compound safety testing in addition to inter-individual monitoring of environmental, lifestyle and occupational factors. However it relies on time-consuming and user-subjective manual scoring. Here we show that imaging flow cytometry and deep learning image classification represents a capable platform for automated, inter-laboratory operation. Images were captured for the cytokinesis-block micronucleus (CBMN) assay across three laboratories using methyl methanesulphonate (1.25 – 5.0 μg/mL) and/or carbendazim (0.8 – 1.6 μg/mL) exposures to TK6 cells. Human-scored image sets were assembled and used to train and test the classification abilities of the "DeepFlow" neural network in both intra- and inter-laboratory contexts. Harnessing image diversity across laboratories yielded a network able to score unseen data from an entirely new laboratory without any user configuration. Image classification accuracies of 98%, 95%, 82% and 85% were achieved for 'mononucleates', 'binucleates', 'mononucleates with MN' and 'binucleates with MN', respectively. Successful classifications of 'trinucleates' (90%) and 'tetranucleates' (88%) in addition to 'other or unscorable' phenotypes (96%) were also achieved. Attempts to classify extremely rare, tri- and tetranucleated cells with micronuclei into their own categories were less successful (≤ 57%). Benchmark dose analyses of human or automatically scored micronucleus frequency data yielded quantitation of the same equipotent dose regardless of scoring method. We conclude that this automated approach offers significant potential to broaden the practical utility of the CBMN method across industry, research and clinical domains. We share our strategy using openly-accessible frameworks.

**Keywords**

Micronucleus test, genetic toxicology, compound screening, machine learning, high throughput, image analysis.

4

**INTRODUCTION**

Across industry, government and academic research institutions the *in vitro* micronucleus test is one of the most widely used bioassays for the identification and quantification of chromosomal damage (Decordier and Kirsch-Volders 2006; Fenech 2000; Fenech 2020; Kirsch-Volders et al. 2011). Because DNA damage at the chromosome level is recognised as a key event in the initiation of carcinogenesis, the assay has become an essential component of genetic toxicity screening programmes worldwide (Fenech 2000). Harmonised assay protocols and scoring approaches have been detailed by Organisation for Economic Cooperation and Development (OECD)-Test Guideline 487 (OECD 2016). In addition to regulatory compound screening, the assay is also widely used for more specific research and clinical purposes including compound mode-of-action determinations, tumour radiosensitivity prediction and inter-individual monitoring of lifestyle, occupational and environmental factors including radiation biodosimetry assessments (Decordier and Kirsch-Volders 2006; Fenech 2000; Fenech 2020; Kirsch-Volders et al. 2011; Wang et al. 2019).

The micronucleus assay operates through the detection of whole chromosomes or chromosome fragments expressed by cells after nuclear division as satellite 'micronucleus' (MN) events. Because complete nuclear division is required to enable expression of these events, the 'cytokinesis-block' version of the assay was developed. This method inhibits cell division into daughter entities (cytokinesis) using the microfilament assembly inhibitor cytochalasin-B. This yields cells that have successfully undergone division easily identifiable by their binucleated appearance. In this way, the cytokinesis-block micronucleus (CBMN) assay allows scoring of micronucleus events in cells known to have undergone division during the treatment period. This avoids misleading results otherwise present due to pre-existing damage, sub-optimal cell culture conditions or from the selection of overly cytotoxic compound doses that retard or inhibit cell division and concomitant micronucleus expression (Decordier and Kirsch-Volders 2006; Fenech 2000; Kirsch-Volders et al. 2011).

113

114    Despite almost global utilisation, CBMN assay scoring still often relies upon manual observation and

115    recording using light microscopy. Whilst manual scoring is the 'gold standard', it suffers from user

116    subjectivity and scorer variability in addition to being extremely time and labour-intensive

117    (Rodrigues et al. 2014a; Rodrigues et al. 2014b; Rodrigues et al. 2018). For these reasons, over the

118    last two decades significant efforts have been directed towards automated approaches for both image

119    collection and subsequent scoring. As recently reviewed (Rodrigues et al. 2018), these largely

120    involve slide and laser scanning microscopy systems that automate image collection in conjunction

121    with traditional, threshold-based image classification techniques (Darzynkiewicz et al. 2011;

122    Decordier et al. 2009; Decordier et al. 2011; François et al. 2014; Maertens and White 2015;

123    Rossnerova et al. 2011; Schunck et al. 2004; Seager et al. 2014; Smolewski et al. 2001; Varga et al.

124    2004; Verhaegen et al. 1994; Willems et al. 2010). Conventional flow cytometry methods have also

125    been developed that aim to identify isolated micronuclei using fluorescence intensity measurements

126    in the absence of image-based validation (Avlasevich et al. 2006; Bryce et al. 2008; Bryce et al.

127    2010; Bryce et al. 2013; Bryce et al. 2007).

128

129    More recently, imaging flow cytometry unites the acquisition approach of flow cytometry with

130    microscopical observation (Allemang et al. 2021; Rodrigues 2018; Rodrigues 2019; Rodrigues et al.

131    2014a; Rodrigues et al. 2014b; Rodrigues et al. 2016a; Rodrigues et al. 2018; Rodrigues et al. 2016b;

132    Wang et al. 2019; Wilkins et al. 2017). This fluidics-based approach is well suited for processing cell

133    suspension cultures (*e.g.,* TK6 B-lymphocytes commonly used for the CBMN assay) enabling rapid

134    collection of transmitted light brightfield, darkfield laser scatter and fluorescence images for

135    populations of tens of thousands of single cells. Simple inclusion of a single nuclear fluorescent stain

136    (*e.g.,* Hoechst 33342, propidium iodide or DRAQ5 *etc.*) allows detection of parent nuclei and

137    micronucleus events (Rodrigues 2018; Rodrigues 2019; Rodrigues et al. 2018; Rodrigues et al.

138    2016b). Without need of further labels, the brightfield images provide essential context for detecting

6

139   micronuclei associated with parent cells (Rodrigues et al. 2014a; Verma et al. 2018). The 'Amnis

140   ImageStream$^{X}$' series cytometers (Luminex Corporation) further support unassisted data acquisition

141   for multiple samples via a 96-well plate sampling attachment. Images are stored to sample-specific

142   data files enabling archiving should human validation or reevaluation be required (Rodrigues et al.

143   2018). Traditional image classification approaches deployed within the manufacturer-supplied

144   analysis software have shown utility for CBMN scoring automation (Rodrigues 2018; Rodrigues

145   2019; Rodrigues et al. 2014a; Rodrigues et al. 2014b; Rodrigues et al. 2016a; Rodrigues et al. 2018;

146   Rodrigues et al. 2016b; Wang et al. 2019; Wilkins et al. 2017). However, in our experience these

147   strategies require significant expertise to set up, in addition to frequent tuning to maintain acceptable

148   performance, even within a single laboratory (Verma et al. 2018). Deviations of around 30% from

149   the results obtained by manual microscopy scoring have also been reported in experiments utilising

150   this approach to study irradiated peripheral blood lymphocytes (Rodrigues et al. 2016b). This

151   outcome was in part attributed to the lack of flexibility of the implemented image analysis algorithms

152   relative to the expertise of human judgement (Rodrigues et al. 2018; Rodrigues et al. 2016b).

153

154   Building image classification strategies that *generalise* well enough to permit robust, entirely

155   automated image classifications without need of human intervention or configuration is a difficult

156   task. This is because, even when protocols are harmonised, there will always be variability (*e.g.,*

157   illumination, focus and fluorescence staining heterogeneity *etc.*) in the input image data. This

158   variation is even more extreme across laboratories due to the inevitable use of different imaging

159   equipment, calibration settings, personnel, cell culture and bioassay regimens. Recently, artificial

160   intelligence approaches have been achieving increasing success in providing generalised automation

161   of image classification tasks (Caicedo et al. 2019; Moen et al. 2019). These approaches can use

162   handcrafted features extracted from images in conjunction with machine learning algorithms, but

163   increasingly, the availability of computational power is enabling the application of deep learning on

164   image pixel data (Blasi et al. 2016; Eulenberg et al. 2017). This approach uses so-called deep

7

165    convolutional neural networks in a manner inspired by neural connectivity in the brain. A typical

166    image classification workflow involves assigning 'ground truth' class annotations to a large set of

167    images before subdividing them into 'train' and 'test' datasets. The weights connecting the nodes of

168    the neural network are then optimised during a training phase that attempts to match the input images

169    to the annotated classifications. A potential issue due to the flexibility of neural networks as non-

170    linear function approximators is that 'memorisation' due to over-fitting of training data can emerge

171    (Zhang et al. 2017). For this reason, final network accuracy is assessed by cross validation against a

172    test set that importantly was entirely 'unseen' during the training phase. Subsequently, the trained

173    neural net can be deployed for the classification of new images.

174

175    In the context of the CBMN assay, deep learning approaches were recently used on imaging flow

176    cytometry data using the cytometer manufacturer's 'Amnis Artificial Intelligence' software to

177    identify binucleated cells in the 3-D reconstructed skin micronucleus assay. This binucleated cell

178    population was then used as a refined start point from which to expedite manual identification of

179    micronucleus events (Allemang et al. 2021). However, there would be considerable value in openly

180    accessible frameworks for accessibility and for adaptability: the modular nature of modern, open

181    source deep learning interfaces allows new network architectures to be easily switched or

182    specifically tailored as they emerge. This flexibility provides complete ability to build bespoke

183    solutions using the latest tools to pursue maximal accuracy and the accommodation of diverse

184    research objectives.

185

186    Here, we used imaging flow cytometry to automate image capture for the CBMN assay across three

187    laboratories using differing local protocols for cell culture, bioassay procedure, DNA staining,

188    cytometer calibration and image collection. Given the inherent variability in the captured images, we

189    investigate the ability of deep learning to enable robust, inter-laboratory scoring automation. To do

190    this, we provide an open framework that utilises the powerful, yet lightweight DeepFlow neural

8

191 network architecture that has been previously optimised to achieve rapid training and classification

192 of imaging flow cytometry data (Eulenberg et al. 2017).

193

194

195 **MATERIALS & METHODS**

196 **Multi-centre image collection**

197 Image data was collected using three different Amnis ImageStream[X] imaging flow cytometers

198 (Luminex Corporation, USA) across three locations: Central Biotechnology Services, Cardiff

199 University School of Medicine (hereafter, Cardiff), the Department of Veterinary Medicine's

200 Imaging Facility, University of Cambridge, UK (Cambridge) and at GlaxoSmithKline Research and

201 Development, Stevenage, UK (GSK).

202

203 **Chemicals**

204 Methyl methanesulphonate (MMS) (#129925) (CAS registry number 66-27-3) and carbendazim

205 (#378674) (CAS no. 10605-21-7) were purchased from Sigma-Aldrich (Merck), UK.

206

207 **Cardiff and Cambridge: Cell culture and cytokinesis-block micronucleus assay**

208 P53 competent, virally transformed human B lymphoblastoid (TK6) cells were purchased from the

209 Health Protection Agency Culture Collections (Wiltshire, UK). The cells were cultured in RPMI

210 1640 media (#A1049101, ThermoFisher) supplemented with 100 U/mL penicillin and 100 μg/mL

211 streptomycin and containing 10% (v/v) heat-inactivated horse serum (#26050088, ThermoFisher).

212 Cells were seeded at 2 x $10^5$ cells/mL in 25 cm$^2$ flasks (ThermoFisher) and incubated at 37 $^o$C for ~

213 1.5 cell cycles (24-30 h) in the presence of MMS (0 / 1.25 / 2.5 / 5.0 μg/mL doses) or carbendazim (0

214 / 0.8 / 1.0 / 1.6 μg/mL doses) with co-exposed cytochalasin-B (#C6762, Sigma) added to a final

215 concentration of 3 μg/mL as a cytokinesis-block. Following exposure, cells were pelleted by

216 centrifugation (200xg, 10 min) and washed once with 10 mL phosphate buffered saline (PBS). Cells

9

217  were then pelleted and resuspended in 2 mL 1X BD FACS lysing solution (#349202, BD) for 12 min

218  to achieve fixation and permeabilisation.

219

220  **GSK: Cell culture and cytokinesis-block micronucleus assay**

221  TK6 (IVGT) cells (#13051501) purchased from ECACC, operated by Public Health England

222  (Wiltshire, UK). The cells were cultured in RPMI 1640 media with 2 mM glutamine (#52400-025,

223  ThermoFisher) supplemented with 100 U/mL penicillin and 100 μg/mL streptomycin (#15140-122,

224  ThermoFisher), 1.8 mM sodium pyruvate (#11360-039, ThermoFisher) and containing 10% (v/v)

225  heat-inactivated horse serum (#26050-088, BioSera, Labtech, UK). Cells were seeded at $2 \times 10^5$

226  cells/mL in 25 $cm^2$ flasks (ThermoFisher) and incubated at 37 $^o$C for 24 h in the presence of

227  carbendazim (0 / 0.8 / 1.2 / 1.6 μg/mL doses) with co-exposed cytochalasin-B (#C6762, Sigma)

228  added to a final concentration of 6 μg/mL as a cytokinesis-block. Following exposure, cells were

229  pelleted by centrifugation (200xg, 10 min) and washed once with 10 mL PBS (#10010-015,

230  ThermoFisher). Cells were then pelleted and resuspended in 2 mL 1X BD FACS lysing solution

231  (#349202, BD) for 12 min to achieve fixation and permeabilisation.

232

233  **Nuclear labelling**

234  Fixed, permeabilised cells were incubated with nuclear stains in PBS at room temperature. Nuclei

235  and micronuclei were stained at the Cardiff and GSK laboratories by 30 min incubation with 0.05

236  mM DRAQ5 (peak excitation: 647 nm, peak emission: 681 nm) (#564902, BD). Samples at the

237  Cambridge laboratory were stained with a 1:2500 dilution (8 μM) of Hoechst 33342 (peak excitation:

238  351 nm, peak emission: 461 nm) (#62249, ThermoFisher) for 30 mins. After labelling, cells were

239  pelleted, resuspended and final cell concentrations adjusted through addition of PBS towards an

240  optimal cell concentration for imaging flow cytometry (typically ~100 μL sample volumes at ~$10^7$

241  cells/mL).

242

10

243    **Imaging flow cytometry**

244    Brightfield and nuclear fluorescence images (20,000 images / sample) were collected using Amnis

245    ImageStream$^X$ (Luminex) flow cytometers using the 40X objective lens via the manufacturer's

246    INSPIRE software at the Cardiff, Cambridge and GSK laboratories (described above). At Cardiff

247    and GSK, DRAQ5-labelled cells were excited using 488 nm or 642 nm lasers (respectively) with the

248    brightfield collected in channel 1 and DRAQ5 in channel 11. At Cambridge, Hoechst 33342-labelled

249    cells were excited using a 405 nm laser with brightfield collection in channel 4 and nuclear

250    fluorescence collection in channel 1. At all locations, a brightfield area range of 100-900 $\mu m^2$ was

251    used to avoid debris, speed bead and large aggregate image collection. Full details of image

252    acquisition settings including the laser excitation powers the exact cytometer models utilised at each

253    location are provided in **Supp. Table S1.**

254

255    **Compensated image file generation using IDEAS**

256    Prior to image extraction, raw image files (.rif) acquired by the INSPIRE software were converted to

257    compensated image files (.cif) using identical settings via batch processing with a template using the

258    IDEAS (version 6.2) software (Luminex). During the process, populations of cell images suitable for

259    scoring were refined by gating out (brightfield area, 200 – 500 $\mu m^2$ versus aspect ratio, 0.75 – 1.0)

260    debris and identifying a single cell population that was also suitably in focus. This was achieved by

261    linescan gradient via the root mean square of the brightfield images ranging from 55 – 80.

262

263    **Image data pre-processing: CIF to TIF extraction**

264    Single, in-focus cell populations were exported from the IDEAS software in compensated image file

265    format (.cif). The individual cell images within these files were then extracted to 16-bit grayscale,

266    two-channel (nuclear fluorescence / brightfield) multipage TIF files using a custom script (code and

267    example available for download from the BioStudies database (http://www.ebi.ac.uk/biostudies) in

268    MATLAB and Python programming languages under accession number S-BSST641). During this

11

269    TIF extraction process, each channel image was also max/min rescaled to normalise illumination.

270    Images were also cropped and zero-padded to a standard 64x64 pixel-square size for input into the

271    DeepFlow network.

272

273    **Deep learning image classification**

274    Automated scoring was achieved using a nine-class, feed-forward, image classification deep neural

275    network built using our previously described "DeepFlow" architecture (Eulenberg et al. 2017). This

276    network is optimised for the relatively small input dimensions of imaging flow cytometry data, and

277    in itself utilises dual-path convolution / batch normalisation / nonlinearity subunits interspersed by

278    max pooling from the popular "Inception" architecture (Szegedy et al. 2015). These subunit layers

279    process and aggregate visual information at increasing scale before average pooling, the fully

280    connected layer and softmax classification (full network architecture shown, **Supp. Figure 1)**.

281    Images were passed to the network with an input size of 64x64x2 (x, y, channels), with augmentation

282    by random x/y reflection, rotation, translation, 90%-110% image scaling and zero-center batch

283    normalisation. Training lasted for 30 epochs using a batch size of 88 with optimisation under ADAM

284    using cross-entropy loss. The initial learn rate was $5x10^{-3}$, dropping every five epochs by 0.9, with

285    L2 regularisation $1x10^{-4}$ and epsilon $1x10^{-8}$. Images were shuffled every epoch. The final pre-trained

286    network alongside test images and all code detailing training hyper-parameters and final layer

287    weightings are available for download in MATLAB (using the Deep Learning Toolbox) or Python

288    (using TensorFlow / keras) languages at the BioStudies database (http://www.ebi.ac.uk/biostudies)

289    under accession number S-BSST641.

290

291    **Ground truth curation by human scoring**

292    For the Cardiff / Cambridge analyses, cell image data across compounds (carbendazim and MMS)

293    and doses (0 – 5 μg/mL) were merged to create diverse ground truth training sets that contained the

294    wide representation of different cell phenotypes essential for effective network training. Ground truth

295    classifications for each image were assigned by biologists with extensive experience manually

296    scoring the *in vitro* micronucleus assay, with phenotypes assigned through consideration of both the

297    nuclear fluorescence and the brightfield image (*i.e.,* ensuring nuclear events belonged to one cell

298    *etc.*). As per micronucleus assay test guidance, the aim was to only score cells positive for

299    micronucleus events where the micronuclei were fluorescently-labelled, were circular/oval in shape,

300    were within the size range of $1/3 - 1/16^{th}$ that of the parent nuclei, and that were clearly inside the

301    cell boundary of the parent cell (Fenech 2000; OECD 2016). At the GSK laboratory, TK6 cells were

302    exposed to just the carbendazim compound (0 / 0.8 / 1.2 / 1.6 μg/mL doses) with the experiment

303    conducted in triplicate. For the initial network cross validation with the GSK data, five thousand

304    human-scored cell images were used with these events equally accumulated from across all

305    carbendazim exposures. For the dose-response analysis, cell populations of two thousand events

306    were scored per dose in triplicate by either human-scoring or by the neural network.

307

308    **Statistical significance of micronucleus responses relative to control**

309    Assessment of micronucleus response significance was conducted according to the framework

310    described in Johnson et al., (Johnson et al. 2014). Response data was $log_{10}$ transformed and assessed

311    for normality and variance homogeneity by Shapiro-Wilk and Bartlett tests respectively. Where the

312    transformed data passed these tests (p > 0.05), comparisons of micronucleus responses relative to

313    untreated negative controls employed one sided *post hoc* Dunnett's test with alpha 0.05. Datasets

314    that failed these tests (p < 0.05) were analysed using the non-parametric *post hoc* Dunn's test.

315

316    **Benchmark dose analysis**

317    To compare the dose-response relationships obtained from human expert scoring relative to those

318    obtained from automatic scoring using the trained neural network, nonlinear regression analysis

319    using the Benchmark Dose (BMD) framework was used. Using the freely available PROAST

320    software, dose-response data were analysed using both the exponential and the Hill model family

13

321    recommended for the assessment of continuous toxicity data by the European Food Safety Authority

322    (EFSA) (Hardy et al. 2017). In each analysis, combined datasets (*i.e.,* across scoring methods) were

323    analysed together with 'scoring method' specified as a potential covariate (Wills et al. 2016). More

324    complex models with additional parameters were accepted if the fit significantly ($p < 0.05$; log-

325    likelihood) improved. Here, as in previous work, we found that the log-steepness (*parameter d*) and

326    maximum response (*parameter c*) could reasonably be held equal across dose-response curves,

327    whereas the parameters for background response (*parameter a*), potency (*parameter b*), and within-

328    group variance (*var*) were found to be covariate-dependent (Slob and Setzer 2014). The BMD output

329    describes the 'equipotent dose' of the modelled dose-response relationships in addition to the

330    bounding, two-sided 90% confidence interval for each level of the covariate. The benchmark

331    response (BMR) size (also termed the critical effect size) used was 50%, which represents a 50%

332    increase in response relative to the background established in the vehicle (zero-dose) control.

333

334

335    **RESULTS**

336    Here, we investigate the ability of deep learning to provide generalised automation of CBMN assay

337    scoring using imaging flow cytometry data acquired according to local protocols across three

338    different laboratories (Cardiff, Cambridge and GSK). **Fig. 1a** demonstrates our workflow. At the end

339    of the assay, cells were fixed and permeabilised before fluorescent nuclear staining. The choice of

340    nuclear stain varied across the different laboratories according to compatibility with the laser

341    configuration of the local imaging cytometer. At Cambridge, cells were labelled with the blue-

342    fluorescent dye Hoechst 33342 which was stimulated by a 405 nm laser with image capture using a

343    ImageStream[X] cytometer. At Cardiff and GSK, ImageStream[X] MKII cytometers were used in

344    conjunction with the red-emitting DRAQ5 nuclear stain and excitation by either a 488 nm or 642 nm

345    laser (respectively). Full details of image acquisition settings at each laboratory are shown in **Supp.**

346    **Table 1**. Image acquisition speeds depended on cell concentrations, in addition to the time taken to

14

347    purge the flow stream and load each new sample; approximately ~ 2000 – 5000 cell-images / minute

348    was typical.

349

350    After image collection, a template file created in the cytometer manufacturer's IDEAS software was

351    used to automatically batch-save populations of single cells that additionally met acceptable focus

352    criteria (see Methods). These cell populations served as the input into the deep learning scoring

353    pipeline. This workflow is provided for download in both MATLAB and Python programming

354    languages at the Biostudies database (accession no. S-BSST641). In brief – the download

355    demonstrates initial image pre-processing to normalise image illumination across cytometers in

356    addition to how to build and train the DeepFlow neural network using a human-scored training

357    image set. After successful training, the saved network can subsequently be used to automate the

358    scoring of new images. For example, **Fig. 1b-j** shows typical events classified by a pretrained, nine-

359    class network with cell classes for mononucleates, binucleates, trinucleates and quadranucleates with

360    or without micronucleus events in addition to a final class for 'other or unscorable' phenotypes.

361

362    As introduced above, an essential component of network testing involves cross validation with

363    human-scored test images unseen during the training phase. We display this evaluation as a

364    confusion matrix, which compares network outputs to the human scores for every image in the test

365    set (explained, **Fig. 1k**). In the subsequently presented results, we use this strategy to rigorously test

366    the ability of a range of trained networks to enable automated CBMN assay scoring in both intra- and

367    inter-laboratory contexts. In each instance, human-scored image sets were built from cell events

368    pooled across the available compounds and exposures. This strategy was chosen to maximise the

369    diversity of cellular phenotypes present, as well as to ensure that the rarer, micronucleated

370    phenotypes that predominately manifested at higher exposures were well represented.

371

15

**a** Tissue culture   Imaging flow cytometry   Deep learning

- Dose cells
- Cytochalasin-B

Harvest cells

Laser excitation
Bright-field illumination
Camera

- Fix / permeabilise
- Nuclear stain

Image cells

Brightfield
Nuclear fluorescence
Input

Convolution   Classification

- Train network
- Cross-validate against unseen images

Automatic scoring using pre-trained neural network

Mononucleate
Mononucleate +MN
Binucleate
Binucleate +MN
Trinucleate
Trinucleate +MN
Tetranucleate
Tetranucleate +MN
Other / unscorable

**b** Mononucleate

**c** Mononucleate +MN

**d** Binucleate

**e** Binucleate +MN

**f** Tetranucleate

**g** Tetranucleate +MN

**h** Other / unscorable

**i** Trinucleate

**j** Trinucleate +MN

**k**

Example cross-validation "confusion matrix":
*Assess neural network accuracy / generalisation across laboratories*

Deep learning classification

| | Binucleate | Binucleate +MN | Mononucleate | Mononucleate +MN | Other / unscorable | Tetranucleate | Tetranucleate +MN | Trinucleate | Trinucleate +MN | |
|---|---|---|---|---|---|---|---|---|---|---|
| Binucleate | **4000** 39.6% | 43 0.4% | 0 0.0% | 0 0.0% | 127 1.3% | 0 0.0% | 0 0.0% | 21 0.2% | 0 0.0% | 95.4% 4.6% |
| Binucleate +MN | 0 0.0% | **139** 1.4% | 0 0.0% | 0 0.0% | 1 0.0% | 0 0.0% | 0 0.0% | 0 0.0% | 7 0.1% | 94.6% 5.4% |
| Mononucleate | 8 0.1% | 0 0.0% | **2135** 21.1% | 10 0.1% | 89 0.9% | 0 0.0% | 0 0.0% | 0 0.0% | 0 0.0% | 95.2% 4.8% |
| Mononucleate +MN | 1 0.0% | 3 0.0% | 1 0.0% | **91** 0.9% | 4 0.0% | 0 0.0% | 0 0.0% | 0 0.0% | 0 0.0% | 91.0% 9.0% |
| Other / unscorable | 47 0.5% | 3 0.0% | 41 0.4% | 20 0.2% | **2467** 24.4% | 3 0.0% | 1 0.0% | 11 0.1% | 0 0.0% | 95.1% 4.9% |
| Tetranucleate | 0 0.0% | 0 0.0% | 0 0.0% | 0 0.0% | 4 0.0% | **258** 2.6% | 4 0.0% | 0 0.0% | 0 0.0% | 97.0% 3.0% |
| Tetranucleate +MN | 0 0.0% | 0 0.0% | 0 0.0% | 0 0.0% | 0 0.0% | 0 0.0% | **15** 0.1% | 0 0.0% | 0 0.0% | 100% 0.0% |
| Trinucleate | 0 0.0% | 6 0.1% | 0 0.0% | 0 0.0% | 3 0.0% | 87 0.9% | 5 0.0% | **410** 4.1% | 16 0.2% | 77.8% 22.2% |
| Trinucleate +MN | 0 0.0% | 0 0.0% | 0 0.0% | 0 0.0% | 0 0.0% | 0 0.0% | 2 0.0% | 0 0.0% | **24** 0.2% | 92.3% 7.7% |
| | 98.6% 1.4% | 71.6% 28.4% | 98.1% 1.9% | 75.2% 24.8% | 91.5% 8.5% | 74.1% 25.9% | 55.6% 44.4% | 92.8% 7.2% | 51.1% 48.9% | **94.4% 5.6%** |

Human scorer classification

**Fig. 1 Automating the *in vitro* micronucleus assay using imaging flow cytometry and deep learning image classification. a** Workflow: harvested cells were fixed and permeabilised before counterstaining the nuclei with a fluorescent DNA stain. Transmitted light brightfield (grey) and nuclear fluorescence (red) images were then automatically captured by high-throughput imaging flow cytometry. After initial training using a human-annotated image set, single cell images from the cytometer can be automatically classified using the neural network image classification algorithm. **b-j** Example image classifications according to a nine-class network developed to score the cytokinesis-block *in vitro* micronucleus assay in human lymphoblastoid TK6 cells. **k** An example cross-validation 'confusion matrix' obtained during preliminary network optimisations and presented here to demonstrate confusion matrix interpretation. The matrix represents an image set scored by humans that is 'unseen' during network training. The horizontal direction represents the human scorer classifications, whilst the vertical direction shows the automated output classifications from the network. The green diagonal represents correct, matching classifications: for example (indicated, red box) 4,000 'binucleate' images, representing 39.6% of the total test image set, were classified correctly. Away from this diagonal, misclassifications are shown *e.g.,* (yellow box) 21 images (0.2%) labelled as 'trinucleates' by human scoring were incorrectly classified as 'binucleates' by the network. In the bottom-right corner (green box) the overall network accuracy and overall misclassification rate are shown for all nine classes (94.4% and 5.6%, respectively). In the white squares down the right-hand side of the matrix, the network precision *i.e.,* true positive / (true positive plus false positive) (green percentages) and the false discovery rate *i.e.,* 100-precision (red percentages) are shown for each classification. The horizontal bottom white row shows the network sensitivity *i.e.,* true positive / (true positive plus false negatives) (green percentages) and false negative rates (red percentages), respectively. Therefore – by example – 95.4% of the images classified as binucleates by the network were binucleates by human-scoring (blue box) whereas the trained model can be expected to correctly assign the binucleate class 98.6% of the time (magenta box). *Scale bars equal 5 microns*
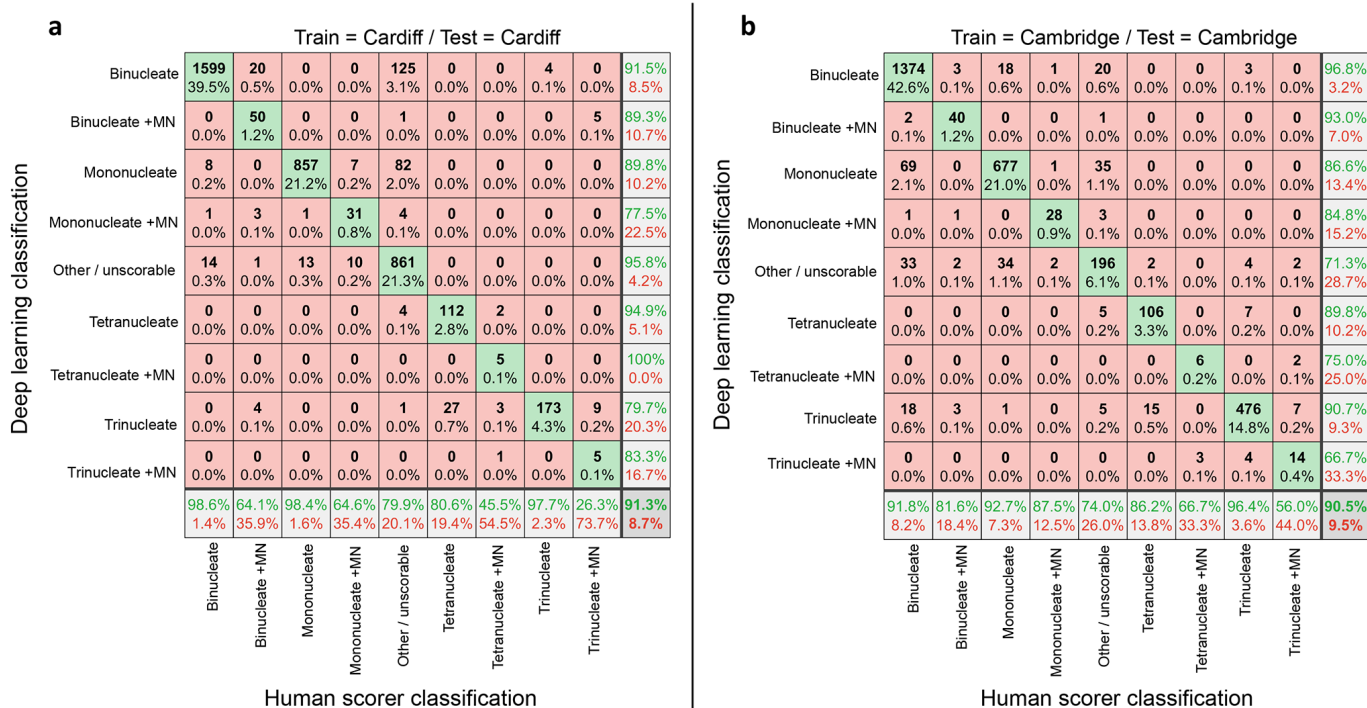
372   First, we tested the ability of a network trained on one laboratory's data to work well for unseen data

373   from that same laboratory (*i.e.,* 'single-laboratory testing') using imaging flow cytometry data

374   collected at either Cardiff or Cambridge (**Fig. 2**). In this single laboratory context, images were

375   randomly assigned to training (60%) and unseen testing (40%) groups. In both instances, the overall

376   accuracies within this single-laboratory context were very high (91.3% and 90.5% for Cardiff and

377   Cambridge, respectively). However, the compiled test sets were quite imbalanced in terms of the

378   numbers of images per class, with network performance with some of the sparser classifications less

379   well represented by the metric of overall accuracy.

380

381   For Cardiff (**Fig. 2a**), whereas accuracy in classification of the common parent nuclei classes (*i.e.,*

382   mononucleates, binucleates, trinucleates) was generally very good (> 97 %), 20 out of a total of 78

383   events (~ 25%) human-scored as 'binucleate + MN' were misclassified as 'binucleates' by the

384   network. Similarly, around 35% of the human-scored 'mononucleate + MN' events were outputted

385   into the 'mononucleate' or 'other/unscorable' classes, with a further ~ 20% of 'tetranucleated' test

386   images misclassified as 'trinucleates'. Despite scoring ~10,000 total events from the Cardiff

387   cytometer, the very rarest cell phenotypes represented by the 'tetranucleate with MN' and

388   'trinucleate with MN' classes presented at very low frequency (~ 0.27 % and 0.47 %, respectively).

389   This led to sparsity in the training set which appeared associated with the network missing

390   micronucleus events, as the 'trinucleate + MN' images were often misclassified into the 'trinucleate'

391   or 'tetranucleate' classes. In a similar manner, 'tetranucleate + MN' images were often misclassified

392   into the 'trinucleate' or 'binucleate + MN' categories.

393

394   Similar results were observed within the Cambridge laboratory (**Fig. 2b**). Whereas accuracies with

395   the 'mononucleate plus MN' and 'binucleate plus MN' classes showed slight improvement when

396   compared against Cardiff, accuracies with the sparser, micronucleated tri- and tetranucleated cells

397   again suffered (~ 44 and ~ 33% error rates, respectively).

16

**a** Train = Cardiff / Test = Cardiff

Deep learning classification (rows) vs Human scorer classification (columns)

| Deep learning \ Human | Binucleate | Binucleate +MN | Mononucleate | Mononucleate +MN | Other / unscorable | Tetranucleate | Tetranucleate +MN | Trinucleate | Trinucleate +MN | Row % |
|---|---|---|---|---|---|---|---|---|---|---|
| Binucleate | 1599 / 39.5% | 20 / 0.5% | 0 / 0.0% | 0 / 0.0% | 125 / 3.1% | 0 / 0.0% | 0 / 0.0% | 4 / 0.1% | 0 / 0.0% | 91.5% / 8.5% |
| Binucleate +MN | 0 / 0.0% | 50 / 1.2% | 0 / 0.0% | 0 / 0.0% | 1 / 0.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 5 / 0.1% | 89.3% / 10.7% |
| Mononucleate | 8 / 0.2% | 0 / 0.0% | 857 / 21.2% | 7 / 0.2% | 82 / 2.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 89.8% / 10.2% |
| Mononucleate +MN | 1 / 0.0% | 3 / 0.1% | 1 / 0.0% | 31 / 0.8% | 4 / 0.1% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 77.5% / 22.5% |
| Other / unscorable | 14 / 0.3% | 1 / 0.0% | 13 / 0.3% | 10 / 0.2% | 861 / 21.3% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 95.8% / 4.2% |
| Tetranucleate | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 4 / 0.1% | 112 / 2.8% | 2 / 0.0% | 0 / 0.0% | 0 / 0.0% | 94.9% / 5.1% |
| Tetranucleate +MN | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 5 / 0.1% | 0 / 0.0% | 0 / 0.0% | 100% / 0.0% |
| Trinucleate | 0 / 0.0% | 4 / 0.1% | 0 / 0.0% | 0 / 0.0% | 1 / 0.0% | 27 / 0.7% | 3 / 0.1% | 173 / 4.3% | 9 / 0.2% | 79.7% / 20.3% |
| Trinucleate +MN | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 1 / 0.0% | 0 / 0.0% | 5 / 0.1% | 83.3% / 16.7% |
| Col % | 98.6% / 1.4% | 64.1% / 35.9% | 98.4% / 1.6% | 64.6% / 35.4% | 79.9% / 20.1% | 80.6% / 19.4% | 45.5% / 54.5% | 97.7% / 2.3% | 26.3% / 73.7% | 91.3% / 8.7% |

**b** Train = Cambridge / Test = Cambridge

| Deep learning \ Human | Binucleate | Binucleate +MN | Mononucleate | Mononucleate +MN | Other / unscorable | Tetranucleate | Tetranucleate +MN | Trinucleate | Trinucleate +MN | Row % |
|---|---|---|---|---|---|---|---|---|---|---|
| Binucleate | 1374 / 42.6% | 3 / 0.1% | 18 / 0.6% | 1 / 0.0% | 20 / 0.6% | 0 / 0.0% | 0 / 0.0% | 3 / 0.1% | 0 / 0.0% | 96.8% / 3.2% |
| Binucleate +MN | 2 / 0.1% | 40 / 1.2% | 0 / 0.0% | 0 / 0.0% | 1 / 0.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 93.0% / 7.0% |
| Mononucleate | 69 / 2.1% | 0 / 0.0% | 677 / 21.0% | 1 / 0.0% | 35 / 1.1% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 86.6% / 13.4% |
| Mononucleate +MN | 1 / 0.0% | 1 / 0.0% | 0 / 0.0% | 28 / 0.9% | 3 / 0.1% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 84.8% / 15.2% |
| Other / unscorable | 33 / 1.0% | 2 / 0.1% | 34 / 1.1% | 2 / 0.1% | 196 / 6.1% | 2 / 0.1% | 0 / 0.0% | 4 / 0.1% | 2 / 0.1% | 71.3% / 28.7% |
| Tetranucleate | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 5 / 0.2% | 106 / 3.3% | 0 / 0.0% | 7 / 0.2% | 0 / 0.0% | 89.8% / 10.2% |
| Tetranucleate +MN | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 6 / 0.2% | 0 / 0.0% | 2 / 0.1% | 75.0% / 25.0% |
| Trinucleate | 18 / 0.6% | 3 / 0.1% | 1 / 0.0% | 0 / 0.0% | 5 / 0.2% | 15 / 0.5% | 0 / 0.0% | 476 / 14.8% | 7 / 0.2% | 90.7% / 9.3% |
| Trinucleate +MN | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 3 / 0.1% | 4 / 0.1% | 14 / 0.4% | 66.7% / 33.3% |
| Col % | 91.8% / 8.2% | 81.6% / 18.4% | 92.7% / 7.3% | 87.5% / 12.5% | 74.0% / 26.0% | 86.2% / 13.8% | 66.7% / 33.3% | 96.4% / 3.6% | 56.0% / 44.0% | 90.5% / 9.5% |

**Fig. 2 Assessing automated scoring accuracies using intra-laboratory train and test data. a/b** Confusion matrices comparing human scoring versus deep learning image classifications for test image sets of approximately four thousand unseen images. In each instance, the results reflect the outputs from nine-class networks trained and tested exclusively on image-data from one imaging cytometer at either the **a** Cardiff or **b** Cambridge laboratories
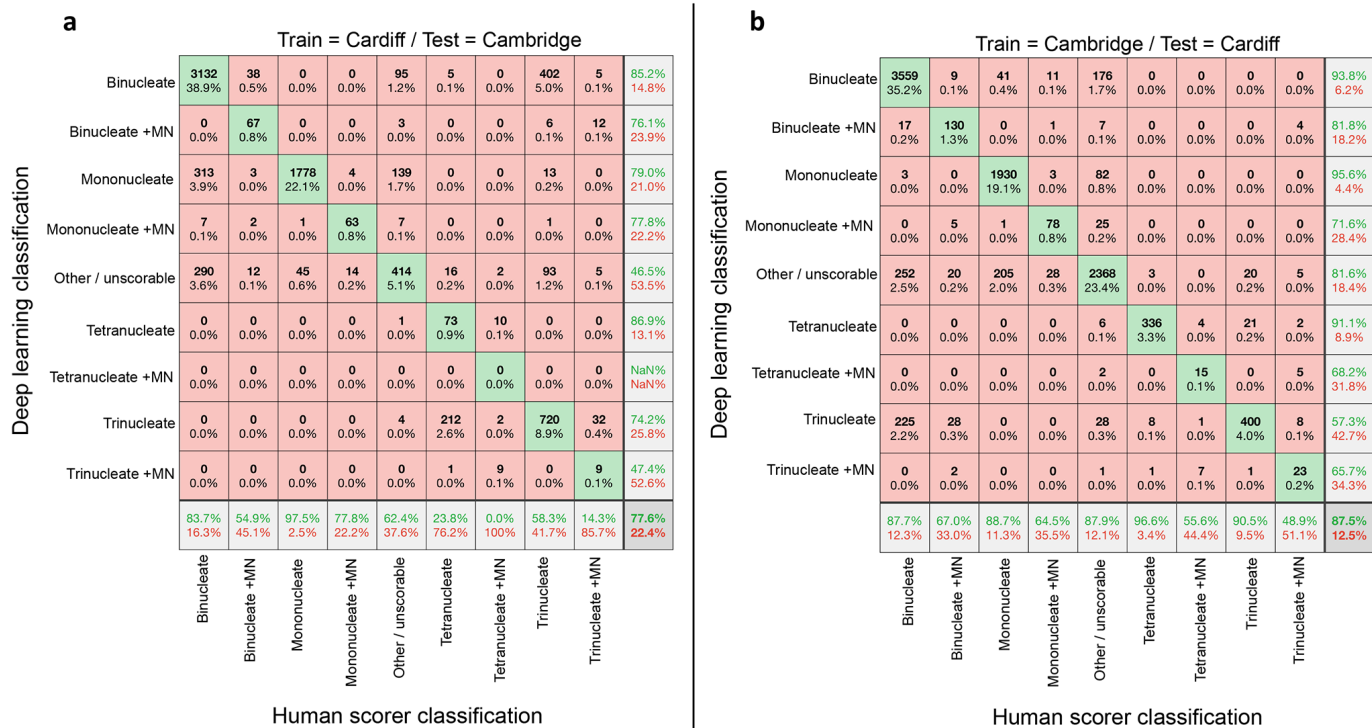
398

399    We next considered the ability of the networks trained on single-laboratory data to generalise to the

400    task of scoring the image data collected from the opposite Centre (**Fig. 3**). This was expected to be a

401    difficult task given that the networks had been trained initially with fairly small numbers of images

402    and because the two laboratories had utilised different cytometer models ($IS^X$ vs. $IS^X$ Mk II) and

403    nuclear stains (Hoechst at Cambridge or DRAQ5 at Cardiff). This presented the likelihood of

404    overfitting during training – yielding networks highly adapted to the task of scoring data from that

405    particular laboratory.

406

407    Despite these factors, at first-glance the overall accuracies appeared quite encouraging at 77.6% for

408    the Cardiff-trained network classifying the Cambridge images (**Fig. 3a**) and 87.5% for the

409    Cambridge network classifying Cardiff images (**Fig. 3b**). Comparing across the individual classes, it

410    was apparent that the Cambridge-trained model generalised slightly better to the task of scoring the

411    Cardiff data than was observed *vice-versa*. Closer examination however showed that the metric of

412    overall accuracy was weighted by the prevalence of the easily identified 'mononucleate' and

413    'binucleate' phenotypes, which masked assessment of the ability of the networks to identify the

414    micronucleated classes representing DNA-damage events (**Fig. 3a/b**). In this regard, in almost all

415    instances, the accuracy of micronucleated event detection suffered considerably compared to the

416    results achieved with laboratory-matched test data (**Fig. 2**).

417

418    With these single-laboratory results established, the images from Cambridge and Cardiff were

419    combined together. This increased the diversity of training exemplifications considerably given the

420    use of two different nuclear stains, two compounds, different imaging cytometers and no 'hold out'

421    requirement for cross validation testing. Training a new DeepFlow neural network on this combined

422    training set (~ 19,000 images) took approximately one hour using modest hardware (single RTX

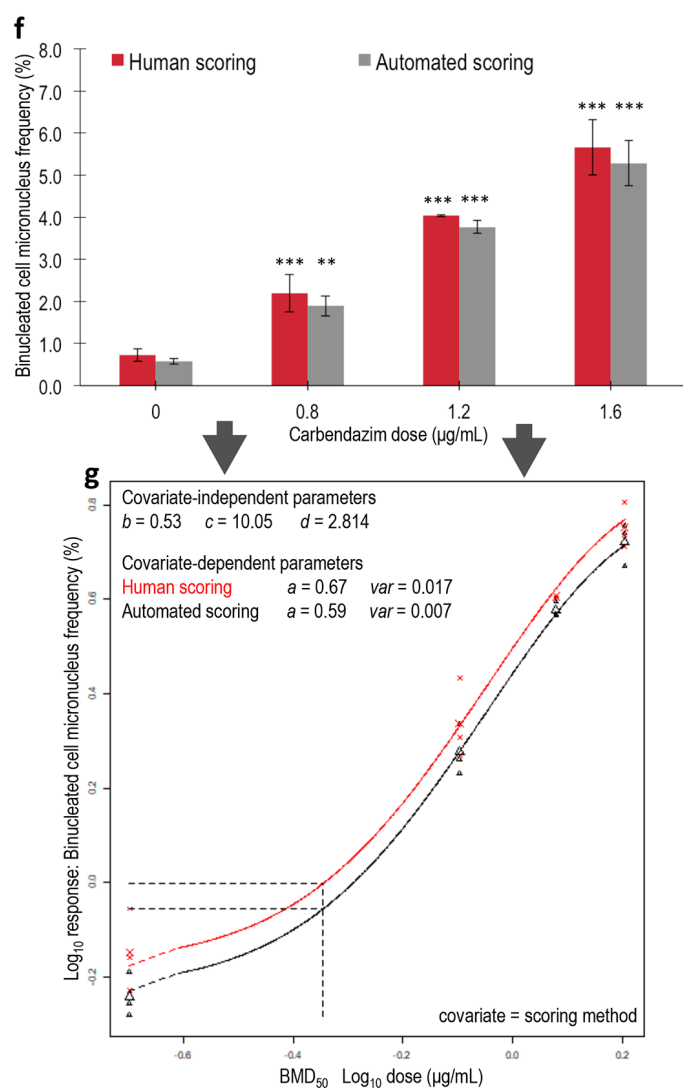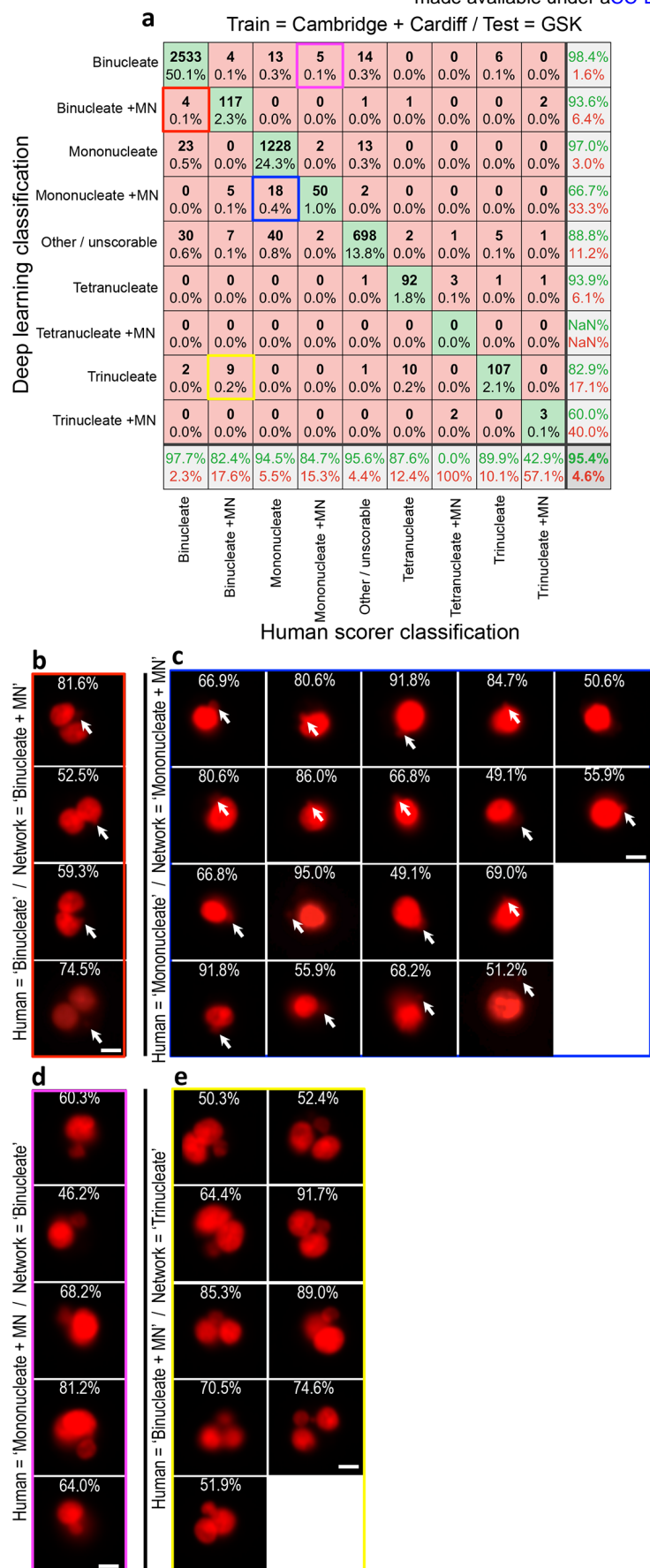423    2080 GPU). The resulting network was then cross validated using a test set where both the bioassay

17

**a**  Train = Cardiff / Test = Cambridge

Deep learning classification (rows) vs Human scorer classification (columns)

| DL \ Human | Binucleate | Binucleate +MN | Mononucleate | Mononucleate +MN | Other / unscorable | Tetranucleate | Tetranucleate +MN | Trinucleate | Trinucleate +MN | Accuracy |
|---|---|---|---|---|---|---|---|---|---|---|
| Binucleate | 3132 / 38.9% | 38 / 0.5% | 0 / 0.0% | 0 / 0.0% | 95 / 1.2% | 5 / 0.1% | 0 / 0.0% | 402 / 5.0% | 5 / 0.1% | 85.2% / 14.8% |
| Binucleate +MN | 0 / 0.0% | 67 / 0.8% | 0 / 0.0% | 0 / 0.0% | 3 / 0.0% | 0 / 0.0% | 0 / 0.0% | 6 / 0.1% | 12 / 0.1% | 76.1% / 23.9% |
| Mononucleate | 313 / 3.9% | 3 / 0.0% | 1778 / 22.1% | 4 / 0.0% | 139 / 1.7% | 0 / 0.0% | 0 / 0.0% | 13 / 0.2% | 0 / 0.0% | 79.0% / 21.0% |
| Mononucleate +MN | 7 / 0.1% | 2 / 0.0% | 1 / 0.0% | 63 / 0.8% | 7 / 0.1% | 0 / 0.0% | 0 / 0.0% | 1 / 0.0% | 0 / 0.0% | 77.8% / 22.2% |
| Other / unscorable | 290 / 3.6% | 12 / 0.1% | 45 / 0.6% | 14 / 0.2% | 414 / 5.1% | 16 / 0.2% | 2 / 0.0% | 93 / 1.2% | 5 / 0.1% | 46.5% / 53.5% |
| Tetranucleate | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 1 / 0.0% | 73 / 0.9% | 10 / 0.1% | 0 / 0.0% | 0 / 0.0% | 86.9% / 13.1% |
| Tetranucleate +MN | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | NaN% / NaN% |
| Trinucleate | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 4 / 0.0% | 212 / 2.6% | 2 / 0.0% | 720 / 8.9% | 32 / 0.4% | 74.2% / 25.8% |
| Trinucleate +MN | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 1 / 0.0% | 9 / 0.1% | 0 / 0.0% | 9 / 0.1% | 47.4% / 52.6% |
| Accuracy | 83.7% / 16.3% | 54.9% / 45.1% | 97.5% / 2.5% | 77.8% / 22.2% | 62.4% / 37.6% | 23.8% / 76.2% | 0.0% / 100% | 58.3% / 41.7% | 14.3% / 85.7% | 77.6% / 22.4% |

**b**  Train = Cambridge / Test = Cardiff

Deep learning classification (rows) vs Human scorer classification (columns)

| DL \ Human | Binucleate | Binucleate +MN | Mononucleate | Mononucleate +MN | Other / unscorable | Tetranucleate | Tetranucleate +MN | Trinucleate | Trinucleate +MN | Accuracy |
|---|---|---|---|---|---|---|---|---|---|---|
| Binucleate | 3559 / 35.2% | 9 / 0.1% | 41 / 0.4% | 11 / 0.1% | 176 / 1.7% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 93.8% / 6.2% |
| Binucleate +MN | 17 / 0.2% | 130 / 1.3% | 0 / 0.0% | 1 / 0.0% | 7 / 0.1% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 4 / 0.0% | 81.8% / 18.2% |
| Mononucleate | 3 / 0.0% | 0 / 0.0% | 1930 / 19.1% | 3 / 0.0% | 82 / 0.8% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 95.6% / 4.4% |
| Mononucleate +MN | 0 / 0.0% | 5 / 0.0% | 1 / 0.0% | 78 / 0.8% | 25 / 0.2% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 71.6% / 28.4% |
| Other / unscorable | 252 / 2.5% | 20 / 0.2% | 205 / 2.0% | 28 / 0.3% | 2368 / 23.4% | 3 / 0.0% | 0 / 0.0% | 20 / 0.2% | 5 / 0.0% | 81.6% / 18.4% |
| Tetranucleate | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 6 / 0.1% | 336 / 3.3% | 4 / 0.0% | 21 / 0.2% | 2 / 0.0% | 91.1% / 8.9% |
| Tetranucleate +MN | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 0 / 0.0% | 2 / 0.0% | 0 / 0.0% | 15 / 0.1% | 0 / 0.0% | 5 / 0.0% | 68.2% / 31.8% |
| Trinucleate | 225 / 2.2% | 28 / 0.3% | 0 / 0.0% | 0 / 0.0% | 28 / 0.3% | 8 / 0.1% | 1 / 0.0% | 400 / 4.0% | 8 / 0.1% | 57.3% / 42.7% |
| Trinucleate +MN | 0 / 0.0% | 2 / 0.0% | 0 / 0.0% | 0 / 0.0% | 1 / 0.0% | 0 / 0.0% | 7 / 0.1% | 1 / 0.0% | 23 / 0.2% | 65.7% / 34.3% |
| Accuracy | 87.7% / 12.3% | 67.0% / 33.0% | 88.7% / 11.3% | 64.5% / 35.5% | 87.9% / 12.1% | 96.6% / 3.4% | 55.6% / 44.4% | 90.5% / 9.5% | 48.9% / 51.1% | 87.5% / 12.5% |

**Fig. 3 Assessment of automated network scoring accuracies using inter-laboratory test data**. **a/b** Confusion matrices comparing human scoring versus deep learning image classifications for test image sets of approximately ten thousand unseen images. In each instance, the results reflect the outputs from nine-class networks trained exclusively on image data from one laboratory's imaging cytometer before cross-validation testing against image data collected at a different laboratory. **a** Network accuracies after training using Cardiff data before testing on unseen Cambridge data. **b** Network accuracies after training on Cambridge data then testing on unseen Cardiff data

424    and imaging cytometry were conducted at an entirely new, third laboratory (GSK). Scoring ~ 5,000

425    test-images took around six seconds on the RTX 2080 hardware or ~ 82 seconds on a single CPU.

426    This time, the network showed much better ability to generalise to the task of successfully scoring

427    the images from the new laboratory (**Fig. 4a**). Across the four core classes central to utilisation of

428    CBMN assay (*i.e.,* 'mononucleate', 'mononucleate plus MN', 'binucleate' and 'binucleate plus

429    MN'), and with no user input or configuration required, the network achieved 98%, 82%, 94%, and

430    85% accuracies, respectively.

431

432    We then examined failure cases, starting with 22 instances where the network detected micronucleus

433    events in cells scored by humans as just mono- or binucleated (**Fig. 4a**). Surprisingly, many did, in

434    fact, appear to have faint or partially occluded potential micronucleus or nuclear bud events that

435    would have been extremely difficult for the human scorer to detect (**Fig. 4b/c**). Similarly,

436    visualisation of cell events scored by humans as either 'mononucleate with MN' or 'binucleate with

437    MN', but outputted by the network as 'binucleate' or 'trinucleate' showed that these images often

438    contained very large micronucleus events (**Fig. 4d/e**). Indeed, some of these likely exceeded the

439    upper size limitation typically imposed on micronucleus classifications (*i.e.,* $\leq 1/3$ diameter of the

440    parent-nuclei) suggesting additional validity to the network's outputs.

441

442    Progressing towards the less frequent cell phenotypes, the accuracies achieved with the 'trinucleate'

443    and 'tetranucleate' cell classes were also good at 90% and 88% respectively. However, detection of

444    these cell types with micronucleus events was either quite poor or failed entirely. Again, this

445    outcome was likely related to extreme sparsity in occurrence (< 0.25 % frequency in the training

446    data). In an attempt to improve accuracies with these classes, we tried both class weighting the

447    classification layer and combining tri- and tetranucleated events with and without micronucleus

448    events into a single, 'polynucleated' class (**Supp. Figure 2**). Whereas both strategies somewhat

449    improved the classification accuracies with these rare events, they were also found to compromise

18

**Fig. 4 Network accuracy and dose-response assessment using unseen test data from a new laboratory. a** Confusion matrix showing human versus deep learning image classifications for a test image set of approximately five thousand unseen images. Here, the neural network was trained using image data from both the Cambridge and Cardiff laboratories before testing on new, unseen imaging cytometry data acquired at a third laboratory (GSK). **b** Cell events human scored as 'binucleates' but classified as 'binucleate plus MN' by the neural network (*i.e.,* red square in **A**). **c** Cell events human scored as 'mononucleates' but classified as 'mononucleate with MN' by the neural network (*i.e.,* blue square in **a**). **b/c,** Close examination of the purportedly misclassified cells shows that many display indistinct events that might be micronucleus or nuclear buds missed by the human scorer (indicated, white arrows). **d** Cell events human scored as 'mononucleate with MN' but classified as 'binucleate' by the neural network (*i.e.,* magenta square in **a**). **e** Events human scored as 'binucleate with MN' classified as 'trinucleate' by the neural network (*i.e.,* yellow square in **a**). **d/e** In both instances, some of the human-scored micronucleus events encroach upon the 1/3 parent nuclei upper-size limitation typically imposed on micronucleus classifications. **b-e** For each event, the white percentages represent neural network confidence in the outputted classification. **f** Binucleated-cell micronucleus frequencies for a three dose plus control dose-response experiment performed in triplicate for carbendazim exposure to TK6 cells. Scores were established from image sets of 2,000 events per replicate by human scoring or by the cross-validated network established in (**a**). (*) (**) (***) indicate statistical significance relative to control at p < 0.05, p < 0.01 and p < 0.001 respectively. **g** Covariate benchmark (BMD) dose modelling using dose-response data from either the human (black) or automated neural network (red) scores established in (**f**). The horizontal and vertical dashed lines represent interpolation to determine the equipotent, benchmark dose for a benchmark response size of 50%. Regardless of human or automated scoring, the model predicts the same benchmark dose. *Scale bars equal 5 microns*

450    the accuracies achieved with one or more of the four core phenotypes more central to successful

451    CBMN assay scoring.

452

453    Given that the frequency of binucleated cells with or without micronucleus events represents the core

454    readout for successful DNA damage assessment by the CBMN assay, after validating the network we

455    proceeded to assess the binucleated-cell micronucleus frequency for a three dose plus control

456    experiment conducted in triplicate with carbendazim at the GSK laboratory. For each dose and

457    replicate, 2000 cell images were scored both manually and automatically. Visually, the resultant

458    dose-response relationships appeared similar across the human and neural network scoring

459    approaches, with the human scores consistently fractionally higher for each dose-group (**Fig. 4f**). To

460    better understand the consequences of this using a recognised, quantitative framework for genotoxic

461    potency estimation, the dose-response relationships were fitted using both the exponential and the

462    Hill model families recommended for the assessment of continuous toxicity data using Benchmark

463    Dose (BMD) analysis (Hardy et al. 2017). With scoring method specified as a potential covariate,

464    model fitting with the PROAST package resulted in covariate-dependent parameterisation for the

465    background response (parameter *a*) and for within-group variation (*var*). For both model families,

466    this parameterisation subsequently allowed rejection of scoring method as covariate, yielding the

467    *same* estimation for the equipotent, benchmark dose from both manual and automated methods (**Fig.**

468    **4g**). Model fits to the data are presented in **Supp. Figure 3**.

469

470

471    **DISCUSSION**

472    The CBMN assay represents a globally significant method for the identification and quantification of

473    chromosomal damage (Fenech 2000; Fenech 2020; OECD 2016). Its utility reaches beyond

474    regulatory compound screening to encompass inter-individual monitoring of wide-ranging lifestyle,

475    occupational and environmental factors (Fenech 2020; Kirsch-Volders et al. 2011; Wang et al.

19

476   2019). Despite this, continued reliance upon time-consuming and user-subjective manual scoring

477   represents a bottleneck to broadening practical utilisation (Seager et al. 2014; Verma et al. 2018;

478   Verma et al. 2017). In this pilot study, we show that rapid image acquisition by imaging flow

479   cytometry in conjunction with deep learning image classification represents a capable platform for

480   automated, inter-laboratory operation. We share our strategy via openly accessible frameworks.

481

482   As an image acquisition method, imaging flow cytometry is now well established as a means for

483   high-throughput CBMN data capture with concomitant image archiving potential (Rodrigues et al.

484   2014a; Rodrigues et al. 2016a; Rodrigues et al. 2018). Moreover, this is achieved with simple sample

485   preparation involving a single nuclear stain and brightfield to provide the context that events lie

486   inside parent cells (Rodrigues et al. 2018). Comparison studies have shown that the captured images

487   contain dose-response information that aligns to results obtained from 'gold standard' manual

488   microscopy scoring (Verma et al. 2018). Whereas conventional flow cytometry offers faster

489   throughput, it lacks this image-based validation whilst additionally requiring cell lysis. This prevents

490   utilisation of the cytokinesis-block version of the assay with its associated advantages such as robust

491   utilisation of primary human cell lines, knowledge that cells have divided during the test period and

492   quantitation of mononucleated, binucleated and different classes of multinucleated cells. This

493   information is useful in the avoidance of misleading negative results and additionally enables

494   calculation of division and replication indexes that contribute to assessments of mitogen response

495   and cytostatic impact (Rodrigues et al. 2018).

496

497   Beyond image collection, automated scoring of imaging flow cytometry data – as with other

498   automated microscopy strategies – has thus far largely relied upon traditional, threshold-based image

499   classification techniques. These require image analysis expertise to implement, alongside user-

500   configuration and tuning to maintain performance (Rodrigues et al. 2018; Seager et al. 2014; Verma

20

501 et al. 2017). Unfortunately, much as with traditional manual scoring, this is time-consuming and

502 subjective.

503

504 In contrast, once successfully trained, the results achieved here suggest that deep learning image

505 classification has the potential to eliminate these expertise and user-input requirements, dramatically

506 reducing the time to results. This comes from encompassing image diversity during network training

507 and harnessing it to improve the consistency and robustness of subsequent classifications. To this

508 end, here we show that utilisation of diverse training data curated across two laboratories utilising

509 different nuclear stains, multiple compounds and two different cytometer models yielded a capable

510 neural network for scoring automation. Without user configuration, the network was able to classify

511 data collected from an entirely new laboratory with > 82% accuracy for each of the four cell

512 phenotypes central to CBMN performance (*i.e.,* mononucleate and binucleate cells with or without

513 micronucleus events) in addition to successfully classifying tri- and tetranucleated cells (> 88%

514 accuracy) and unscorable events (96% accuracy). Importantly, these seven classes encompassed

515 virtually all of the cell images encountered (>99%). Success at micronucleus detection in both

516 mononucleate and binucleate cell classes further suggests that this single network could be used to

517 automate scoring of both mononuclear and cytokinesis-block versions of the assay.

518

519 Despite this success with the assay classes central to CBMN scoring, the scarce, tri- and

520 tetranucleated phenotypes with micronucleus events proved more challenging. Commonly employed

521 methods such as class weighting or class combination offered little in the way of accuracy

522 improvements, and often compromised accuracy with the other classes. These findings suggest that

523 significant increases in the representation of these sparse events during training will likely be

524 required to improve success. In this context, imaging flow cytometry is well suited to examine

525 whether an improved image bank leads to enhanced accuracy in scoring given the high rates of

526 image capture achievable. Our results also suggests that class reduction does not necessarily simplify

21

527    the classification problem and may instead cause ambiguities. In this way, future expansions to the

528    number of classes to encompass all distinctive cellular phenotypes may represent a route to

529    improving overall network performance.

530

531    In this regard, we identified additional, potentially-scorable cell phenotypes (**Fig. 5**). In particular,

532    cell death events (*i.e.,* due to apoptosis and necrosis) were visually apparent, but we were unable to

533    determine apoptotic from necrotic events using just the brightfield and nuclear fluorescence images

534    alone. Cells caught during mitosis also represented distinctive events. At the same time, we were less

535    convinced that more subtle phenotypes relevant to the expanded, CBMN cytome assay such as

536    nuclear buds and bridges could reliably and consistently be detected – given the relatively low

537    resolution of the image data (Fenech 2007). However, it is important to note that previous studies

538    demonstrating capture of these phenotypes by imaging flow cytometry have utilised both the 60X

539    ImageStream objective lens in addition to hypotonic treatments to swell cell volumes prior to

540    imaging (Rodrigues 2019; Rodrigues et al. 2018). Hypotonic treatments were not used here but may

541    improve image capture of these more subtle phenotypes. With regards to network class expansion to

542    encompass these events – or, indeed for simultaneous measurement of other endpoints – the

543    ImageStream platform is capable of multiplexed imaging. Additional channels might therefore be

544    used to simultaneously measure other DNA-damage pathways (*e.g.,* ϒH2AX for DNA double-strand

545    breaks (Smart et al. 2011)), or to improve the reliability of ground truth image curations through use

546    of additional fluorescent markers to differentiate events such as apoptotic from necrotic cells.

547

548    Manual scoring of the images for this experiment was more challenging than the exemplar images

549    shown might suggest. Fundamentally, the acquired images are relatively low resolution (*i.e.,* cells

550    occupy ~ 64x64 pixels) and further image degradation is always present as a result of the capture of

551    moving objects by time delay integration. The acquired images also represent a central, 2-D

552    projection of a 3-D cell-object. This means that nuclei and micronucleus events may overlap each

22

**Fig. 5 Other scorable cell phenotypes captured by imaging flow cytometry**. **a** Cells undergoing mitosis were visually apparent according to metaphase spread-type nuclear fluorescence imagery (red) alongside large, brightfield-delineated cell sizes (grey). **b/c** Cell death events displayed shrunken cell sizes in conjunction with granular brightfield and fluorescence imagery. In the case of cell death, two distinctive cell phenotypes appeared visually separable according to cell size and the number, size and extent of nuclear foci formation (**b** versus **c**). Whether these observations represented distinct apoptotic versus necrotic events was unclear from the nuclear fluorescence and brightfield information alone. *Scale bars equal 5 microns*

553    other, or they may lie outside of the plane of optimal focus (Rodrigues et al. 2018). These factors all

554    served to make ground truth assignments more complicated, even for experienced CBMN scorers.

555    Whereas network accuracy assessments by confusion matrix provided a more representative

556    breakdown of outputs when compared to simplistic overall accuracy measures, it is a relatively

557    stringent success measure because any ambiguity in human score assignment is not captured. A

558    potential advantage of automated network classification approach is therefore likely greater

559    consistency – even in error – than arises from manual scoring.

560

561    Regarding image focussing, the ImageStream platform offers 'extended depth of field' (EDF)

562    technology, whereby image deconvolution is used to improve the utility of out of focus events

563    through projection onto a single plane (Ortyn et al.). Whereas previous studies have shown this

564    technique can improve accuracy in 'spot counting' applications, the strategy has been reported less

565    helpful for the provision of improved CBMN data (Parris et al. ; Rodrigues 2018; Rodrigues et al.

566    2014a). This was attributed to a slight degradation in overall image resolution, compromising

567    differentiation of micronucleus events from parent nuclei (Rodrigues 2018). On a similar theme, the

568    ImageStream platform is also configurable with 20X, 40X or 60X objective lenses. Here, image

569    collection was via the 'standard', 40X objective across all laboratories. This approach was chosen as

570    previous work has shown that whilst greater resolution is achievable with the 60X objective, focus

571    depth also decreases, reinforcing the out of plane difficulties described above (Rodrigues et al.

572    2018).

573

574    Whilst considering the nature and utility of imaging flow cytometry data, a relevant comparison is to

575    that provided by other automated imaging methods such as slide scanning platforms. In addition to

576    the potential for higher resolution imaging, here an overlooked advantage comes from the ability to

577    use slide-based preparations created by cytocentrifugation. This technique causes the flattening and

578    spreading of cellular content, presenting nuclear objects on a more two dimensional plane (Fitzgerald

23

579  and Hosking 1982; Shanholtzer et al. 1982). From a practical perspective however, this also

580  necessitates the consistent preparation of high-quality slides with optimal cell densities (Rodrigues et

581  al. 2018). Meanwhile, a major advantage of the imaging flow cytometry approach is that single cell

582  image data is inherently acquired by the fluidics-based processing of individualised cells.

583

584  **CONCLUSIONS**

585  As a platform for the CBMN assay, imaging flow cytometry combines the high throughput and

586  multiplexing potential of flow cytometry with the image-based validation and archiving attributes of

587  automated microscopy. Here we demonstrate accurate, automated assay scoring using a neural

588  network for data collected in a laboratory wholly separate to that in which the algorithm was trained.

589  This proves that without any human configuration, the machine is able to correctly anticipate the

590  decisions of the expert human on unseen images in a new setting. For the first time, this suggests the

591  possibility for generalised scoring automation through dissemination of a pretrained network for the

592  ImageStream platform established from ground truth agreed by a single, expert group. Such an

593  approach would provide the ultimate in terms of standardisation and result reliability, but more

594  importantly could enable adoption of the assay beyond current practitioners as local expertise in

595  scoring and/or image analysis would no longer be required. For these reasons, we believe that full

596  development of this automated, accessible, inter-laboratory approach would represent a truly twenty-

597  first century method with significant potential to transform CBMN utility across industry, research

598  and clinical domains.

599

24

600 **References**

601  Allemang A, Thacker R, DeMarco RA, Rodrigues MA, Pfuhler S (2021) The 3D reconstructed skin

602   micronucleus assay using imaging flow cytometry and deep learning: A proof-of-principle

603   investigation. Mutat Res Genet Toxicol Environ Mutagen 865:503314.

604   https://doi.org/10.1016/j.mrgentox.2021.503314

605

606 Avlasevich SL, Bryce SM, Cairns SE, Dertinger SD (2006) In vitro micronucleus scoring by flow

607   cytometry: differential staining of micronuclei versus apoptotic and necrotic chromatin

608   enhances assay reliability. Environ Mol Mutagen 47(1):56-66. https://doi:10.1002/em.20170

609

610 Blasi T, Hennig H, Summers HD, et al. (2016) Label-free cell cycle analysis for high-throughput

611   imaging flow cytometry. Nat Comms 7(1):10256. https://doi:10.1038/ncomms10256

612

613 Bryce SM, Avlasevich SL, Bemis JC, et al. (2008) Interlaboratory evaluation of a flow cytometric,

614   high content in vitro micronucleus assay. Mutat Res 650(2):181-95.

615   https://doi:10.1016/j.mrgentox.2007.11.006

616

617 Bryce SM, Avlasevich SL, Bemis JC, Phonethepswath S, Dertinger SD (2010) Miniaturized flow

618   cytometric in vitro micronucleus assay represents an efficient tool for comprehensively

619   characterizing genotoxicity dose-response relationships. Mutat Res 703(2):191-9.

620   https://doi:10.1016/j.mrgentox.2010.08.020

621

622  Bryce SM, Avlasevich SL, Bemis JC, et al. (2013) Flow cytometric 96-well microplate-based in

623   vitro micronucleus assay with human TK6 cells: protocol optimization and transferability

624   assessment. Environ Mol Mutagen 54(3):180-94. https://doi:10.1002/em.21760

625

626     Bryce SM, Bemis JC, Avlasevich SL, Dertinger SD (2007) In vitro micronucleus assay scored by

627         flow cytometry provides a comprehensive evaluation of cytogenetic damage and cytotoxicity.

628         Mutat Res 630(1-2):78-91. https://doi:10.1016/j.mrgentox.2007.03.002

629

630     Caicedo JC, Goodman A, Karhohs KW, et al. (2019) Nucleus segmentation across imaging

631         experiments: the 2018 Data Science Bowl. Nat Methods 16(12):1247-1253.

632         https://doi:10.1038/s41592-019-0612-7

633

634     Darzynkiewicz Z, Smolewski P, Holden E, et al. (2011) Laser scanning cytometry for automation of

635         the micronucleus assay. Mutagenesis 26(1):153-61. https://doi:10.1093/mutage/geq069

636

637     Decordier I, Kirsch-Volders M (2006) The in vitro micronucleus test: from past to future. Mutat Res

638         607(1):2-4. https://doi:10.1016/j.mrgentox.2006.04.008

639

640     Decordier I, Papine A, Plas G, et al. (2009) Automated image analysis of cytokinesis-blocked

641         micronuclei: an adapted protocol and a validated scoring procedure for biomonitoring.

642         Mutagenesis 24(1):85-93. https://doi:10.1093/mutage/gen057

643

644     Decordier I, Papine A, Vande Loock K, Plas G, Soussaline F, Kirsch-Volders M (2011) Automated

645         image analysis of micronuclei by IMSTAR for biomonitoring. Mutagenesis 26(1):163-8.

646         https://doi:10.1093/mutage/geq063

647

648     Eulenberg P, Köhler N, Blasi T, et al. (2017) Reconstructing cell cycle and disease progression using

649         deep learning. Nat Comms 8(1):463. https://doi:10.1038/s41467-017-00623-3

650

651    Fenech M (2000) The in vitro micronucleus technique. Mutat Res 455(1-2):81-95.

652        https://doi:10.1016/s0027-5107(00)00065-8

653

654    Fenech M (2007) Cytokinesis-block micronucleus cytome assay. Nat Protoc 2(5):1084-1104.

655        https://doi:10.1038/nprot.2007.77

656

657    Fenech M (2020) Cytokinesis-block micronucleus cytome assay evolution into a more

658        comprehensive method to measure chromosomal instability. Genes 11(10):1203.

659        https://doi:10.3390/genes11101203

660

661    Fitzgerald MG, Hosking CS (1982) Cell structure and percent viability by a slide centrifuge

662        technique. J Clin Pathol 35(2):191-194. https://doi:10.1136/jcp.35.2.191

663

664    François M, Hochstenbach K, Leifert W, Fenech MF (2014) Automation of the cytokinesis-block

665        micronucleus cytome assay by laser scanning cytometry and its potential application in

666        radiation biodosimetry. BioTechniques 57(6):309-12. https://doi:10.2144/000114239

667

668    Hardy A, Benford D, Halldorsson T, et al. (2017) Update: use of the benchmark dose approach in

669        risk assessment. EFSA J 15(1):e04658. https://doi:10.2903/j.efsa.2017.4658

670

671    Johnson GE, Soeteman-Hernández LG, Gollapudi BB, et al. (2014) Derivation of point of departure

672        (PoD) estimates in genetic toxicology studies and their potential applications in risk

673        assessment. Environ Mol Mutagen 55(8):609-23. https://doi:10.1002/em.21870

674

675 Kirsch-Volders M, Plas G, Elhajouji A, et al. (2011) The in vitro MN assay in 2011: origin and fate,

676     biological significance, protocols, high throughput methodologies and toxicological

677     relevance. Arch Toxicol 85(8):873-99. https://doi:10.1007/s00204-011-0691-4

678

679 Maertens RM, White PA (2015) RE: Recommendations, evaluation and validation of a semi-

680     automated, fluorescent-based scoring protocol for micronucleus testing in human cells.

681     Mutagenesis 30(2):311-2. https://doi:10.1093/mutage/geu066

682

683 Moen E, Bannon D, Kudo T, Graf W, Covert M, Van Valen D (2019) Deep learning for cellular

684     image analysis. Nat Methods 16(12):1233-1246. https://doi:10.1038/s41592-019-0403-1

685

686 OECD (2016) Test Guidline 487 Guideline for the Testing of Chemicals, In Vitro Mammalian Cell

687     Micronucleus Test. Organisation for Economic Cooperation.

688     https://doi.org/10.1787/9789264264861-en

689

690 Ortyn WE, Perry DJ, Venkatachalam V, Liang L, Hall BE, Frost K, Basiji DA (2007) Extended

691     depth of field imaging for high speed cell analysis. Cytometry A 71(4):215-31. https://doi:

692     10.1002/cyto.a.20370.

693

694 Parris CN, Adam Zahir S, Al-Ali H, Bourton EC, Plowman C, Plowman PN (2015) Enhanced γ-

695     H2AX DNA damage foci detection using multimagnification and extended depth of field in

696     imaging flow cytometry. Cytometry A 87(8):717-723. https://doi:10.1002/cyto.a.22697

697

698 Rodrigues MA (2018) Automation of the in vitro micronucleus assay using the Imagestream imaging

699     flow cytometer. Cytometry A 93(7):706-726. https://doi:10.1002/cyto.a.23493

700

701    Rodrigues MA (2019) An Automated Method to Perform The In Vitro Micronucleus Assay using

702        Multispectral Imaging Flow Cytometry. JoVE (147),e59324. https://doi:10.3791/59324

703

704    Rodrigues MA, Beaton-Green LA, Kutzner BC, Wilkins RC (2014a) Automated analysis of the

705        cytokinesis-block micronucleus assay for radiation biodosimetry using imaging flow

706        cytometry. Radiat Environ Biophys 53(2):273-82. https://doi:10.1007/s00411-014-0525-x

707

708    Rodrigues MA, Beaton-Green LA, Kutzner BC, Wilkins RC (2014b) Multi-parameter dose

709        estimations in radiation biodosimetry using the automated cytokinesis-block micronucleus

710        assay with imaging flow cytometry. Cytometry A 85(10):883-93.

711        https://doi:10.1002/cyto.a.22511

712

713    Rodrigues MA, Beaton-Green LA, Wilkins RC (2016a) Validation of the Cytokinesis-block

714        Micronucleus Assay Using Imaging Flow Cytometry for High Throughput Radiation

715        Biodosimetry. Health Phys 110(1):29-36. https://doi:10.1097/hp.0000000000000371

716

717    Rodrigues MA, Beaton-Green LA, Wilkins RC, Fenech MF (2018) The potential for complete

718        automated scoring of the cytokinesis block micronucleus cytome assay using imaging flow

719        cytometry. Mutat Res Genet Toxicol Environ Mutagen 836:53-64.

720        https://doi:10.1016/j.mrgentox.2018.05.003

721

722    Rodrigues MA, Probst CE, Beaton-Green LA, Wilkins RC (2016b) Optimized automated data

723        analysis for the cytokinesis-block micronucleus assay using imaging flow cytometry for high

724        throughput radiation biodosimetry. Cytometry A 89(7):653-62.

725        https://doi:10.1002/cyto.a.22887

726

727  Rossnerova A, Spatova M, Schunck C, Sram RJ (2011) Automated scoring of lymphocyte

728      micronuclei by the MetaSystems Metafer image cytometry system and its application in

729      studies of human mutagen sensitivity and biodosimetry of genotoxin exposure. Mutagenesis

730      26(1):169-75. https://doi:10.1093/mutage/geq057

731

732  Schunck C, Johannes T, Varga D, Lörch T, Plesch A (2004) New developments in automated

733      cytogenetic imaging: unattended scoring of dicentric chromosomes, micronuclei, single cell

734      gel electrophoresis, and fluorescence signals. Cytogenet. Genome Res. 104:383-9.

735      https://doi:10.1159/000077520

736

737  Seager AL, Shah UK, Brüsehafer K, et al. (2014) Recommendations, evaluation and validation of a

738      semi-automated, fluorescent-based scoring protocol for micronucleus testing in human cells.

739      Mutagenesis 29(3):155-64. https://doi:10.1093/mutage/geu008

740

741  Shanholtzer CJ, Schaper PJ, Peterson LR (1982) Concentrated gram stain smears prepared with a

742      cytospin centrifuge. J Clin Microbiol 16(6):1052.

743

744  Slob W, Setzer RW (2014) Shape and steepness of toxicological dose-response relationships of

745      continuous endpoints. Crit Rev Toxicol. 44(3):270-97.

746      https://doi:10.3109/10408444.2013.853726

747

748  Smart DJ, Ahmedi KP, Harvey JS, Lynch AM (2011) Genotoxicity screening via the γH2AX by

749      flow assay. Mutat Res Genet Toxicol Environ Mutagen 715(1):25-31.

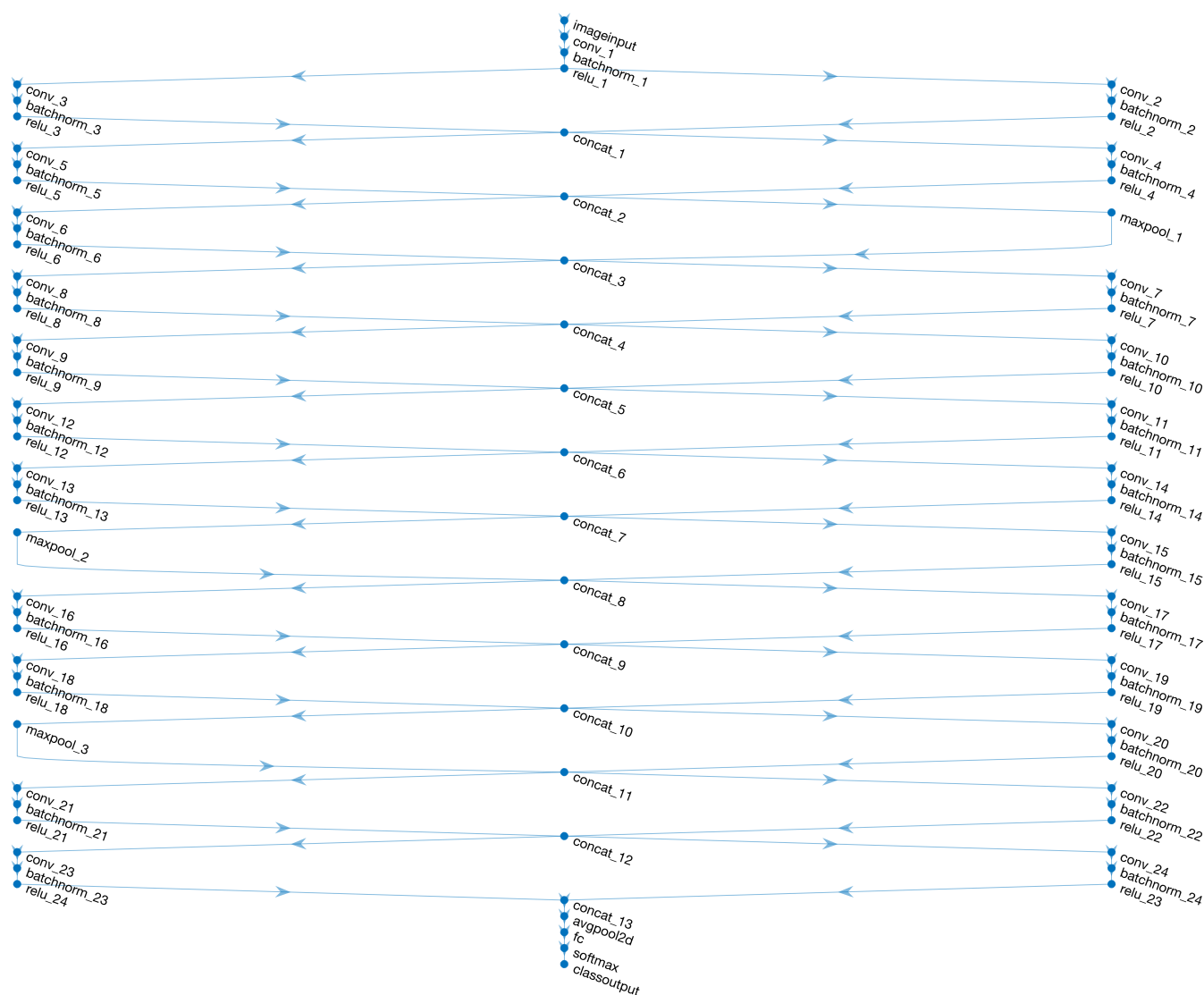750      https://doi.org/10.1016/j.mrfmmm.2011.07.001

751

752    Smolewski P, Ruan Q, Vellon L, Darzynkiewicz Z (2001) Micronuclei assay by laser scanning

753        cytometry. Cytometry 45(1):19-26. https://doi:10.1002/1097-0320(20010901)45

754

755    Szegedy C, Wei L, Yangqing J, et al. (2015) Going deeper with convolutions. 2015 IEEE

756        Conference on Computer Vision and Pattern Recognition (CVPR) arXiv:1409.4842.

757

758    Varga D, Johannes T, Jainta S, et al. (2004) An automated scoring procedure for the micronucleus

759        test by image analysis. Mutagenesis 19(5):391-7. https://doi:10.1093/mutage/geh047

760

761    Verhaegen F, Vral A, Seuntjens J, Schipper NW, de Ridder L, Thierens H (1994) Scoring of

762        radiation-induced micronuclei in cytokinesis-blocked human lymphocytes by automated

763        image analysis. Cytometry 17(2):119-27. https://doi:10.1002/cyto.990170203

764

765    Verma JR, Harte DSG, Shah UK, et al. (2018) Investigating FlowSight imaging flow cytometry as a

766        platform to assess chemically induced micronuclei using human lymphoblastoid cells in

767        vitro. Mutagenesis 33(4):283-289. https://doi:10.1093/mutage/gey021

768

769    Verma JR, Rees BJ, Wilde EC, et al. (2017) Evaluation of the automated MicroFlow and Metafer

770        platforms for high-throughput micronucleus scoring and dose response analysis in human

771        lymphoblastoid TK6 cells. Arch Toxicol 91(7):2689-2698. https://doi:10.1007/s00204-016-

772        1903-8

773

774    Wang Q, Rodrigues MA, Repin M, et al. (2019) Automated Triage Radiation Biodosimetry:

775        Integrating Imaging Flow Cytometry with High-Throughput Robotics to Perform the

776        Cytokinesis-Block Micronucleus Assay. Radiat Res 191(4):342-351.

777        https://doi:10.1667/rr15243.1

778

779   Wilkins RC, Rodrigues MA, Beaton-Green LA (2017) The application of imaging flow cytometry to

780       high-throughput biodosimetry. Genome Integr 8:7. https://doi:10.4103/2041-9414.198912

781

782   Willems P, August L, Slabbert J, et al. (2010) Automated micronucleus (MN) scoring for population

783       triage in case of large scale radiation events. Int J Radiat Biol. 86(1):2-11.

784       https://doi:10.3109/09553000903264481

785

786   Wills JW, Johnson GE, Doak SH, Soeteman-Hernández LG, Slob W, White PA (2016) Empirical

787       analysis of BMD metrics in genetic toxicology part I: in vitro analyses to provide robust

788       potency rankings and support MOA determinations. Mutagenesis 31(3):255-63.

789       https://doi:10.1093/mutage/gev085

790

791   Zhang C, Bengio S, Hardt M, Recht B, Vinyals O (2017) Understanding deep learning requires

792       rethinking generalization. ICLR. https://arxiv.org/abs/1611.03530
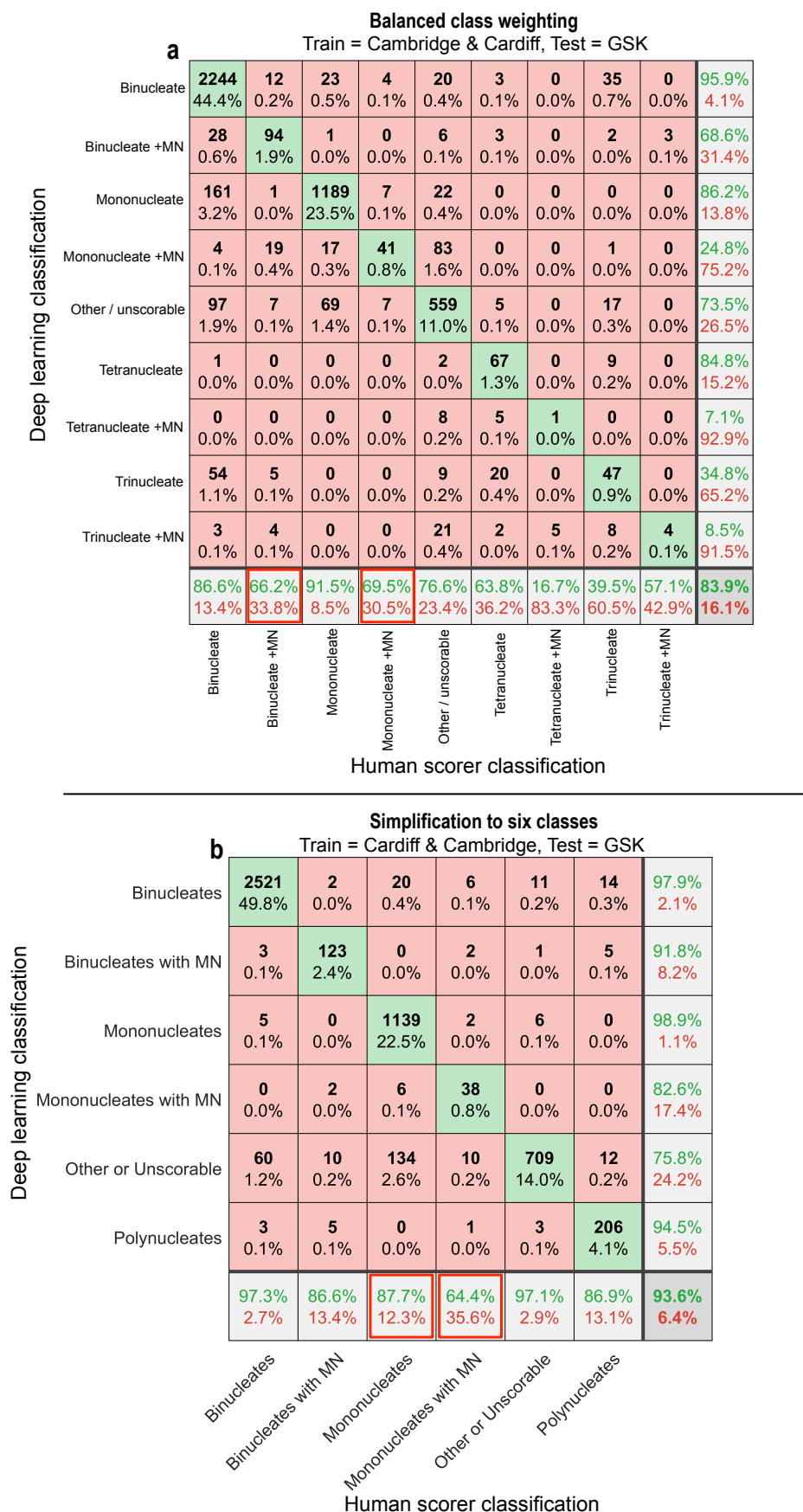
793
794

**Supp. Table 1 – Imaging flow cytometry data acquisition information**

| Centre | Excitation laser (nm) | Intensity (mW) | Brightfield channel | Nuclear fluorescence channel | Nuclear stain | Objective lens | Cytometer Model |
|--------|----------------------|----------------|---------------------|------------------------------|---------------|----------------|-----------------|
| Cambridge | 405 | 50 | Ch04 | Ch01 | Hoechst 33342 | 40X | Amnis ImageStream$^X$ |
| Cardiff | 488 | 100 | Ch01 | Ch11 | DRAQ5 | 40X | Amnis ImageStream$^X$ MkII |
| GSK | 642 | 55 | Ch01 | Ch11 | DRAQ5 | 40X | Amnis ImageStream$^X$ MkII |

Image data were collected using three different imaging flow cytometers located across three laboratories (Cambridge, Cardiff and GSK). At each laboratory, the choice of florescent nuclear stain depended upon local protocols and compatibility with the cytometer's laser configuration.
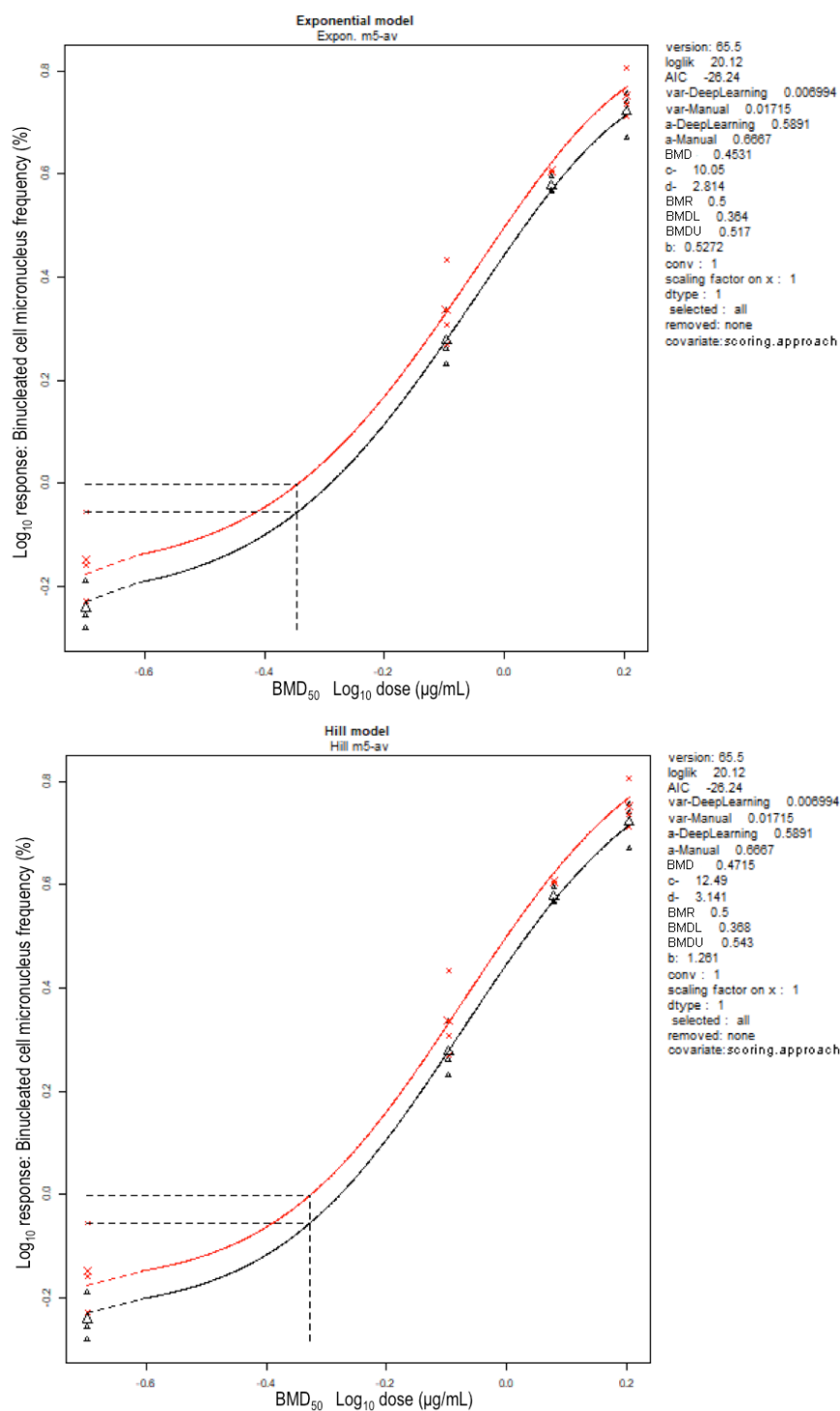
**Supp. Figure 1** DeepFlow neural network architecture schematic. The DeepFlow network utilises a 64x64x2 input layer (x, y, channels) followed by repeating dual-path subunits from the "Inception" architecture to aggregate visual information over increasing scales. The number of kernels used increases at each layer, yielding 336 features maps with size 8 x 8 before average pooling, the fully connected (fc) layer and softmax classification using cross-entropy loss.

**Supp. Figure 2** Cross validation testing using class weighting or class simplification strategies. **a/b** Confusion matrices comparing human scoring versus deep learning image classifications for a test set of ~ 5000 unseen images. In each instance, the results reflect the outputs after training using image data from both the Cambridge and Cardiff laboratories before cross validation on new imaging cytometry data acquired at a third laboratory (GSK). In **a** class weighted cross entropy loss was used at the classification layer in an attempt to improve performance with the sparsely-represented phenotypes (*i.e.*, tri and tetranucleates with or without micronucleus (MN) events). In **b** these sparse, multinucleated categories were combined together into a single 'polynucleated' class. Whilst some improvements were realised using these strategies, they both reduced achieved accuracies (indicated, red squares) with one or more of the four, core phenotypes central to successful CBMN scoring (*i.e.*, mono or binucleated cells with or without MN events).

**Supp. Figure 3** Benchmark dose (BMD) analysis using exponential and Hill model families. The curves represent fits to micronucleus dose-response data obtained either by human (red) or neural network (black) scoring using either the exponential (top) or the Hill (bottom) model families. Both models were fitted with covariate (scoring method) dependent parameters for the background (parameter *a*) and within-group variance (*var*), whilst constant parameters could be used for potency, shape and steepness (parameters *b*, *c* and *d*). Horizontal and vertical dashed lines represent interpolation at a benchmark response (BMR) size of 50% to determine the $BMD_{50}$ (respectively).