1 # A semi-automated technique for adenoma quantification in the
2 # *Apc^Min* mouse using *FeatureCounter*

3 Amy L. Shepherd[1], A. Alexander T. Smith[1], Kirsty A. Wakelin[1], Sabine Kuhn[1], Jianping
4 Yang[1], David A. Eccles[1] and Franca Ronchese[1*]
5
6 [1]Malaghan Institute of Medical Research, Wellington, New Zealand.
7
8
9
10 *To whom correspondence should be addressed at:
11 fronchese@malaghan.org.nz
12 Malaghan Institute of Medical Research
13 PO Box 7060
14 Newtown, Wellington 6242
15 New Zealand
16

19
20

21   **ABSTRACT**

22   Colorectal cancer is a major contributor to death and disease worldwide. The $Apc^{Min}$ mouse is

23   a widely used model of intestinal neoplasia, as it carries a mutation also found in human

24   colorectal cancers. However, the method most commonly used to quantify tumour burden in

25   these mice is manual adenoma counting, which is time consuming and poorly suited to

26   standardization across different laboratories. We describe a method to produce suitable

27   photographs of the small intestine, process them with an ImageJ macro, *FeatureCounter*,

28   which automatically locates image features potentially corresponding to adenomas, and a

29   machine learning pipeline to identify and quantify them. Compared to a manual method, the

30   specificity (or True Negative Rate, TNR) and sensitivity (or True Positive Rate, TPR) of this

31   method in detecting adenomas are similarly high at about 80% and 87%, respectively.

32   Importantly, total adenoma area measures derived from the automatically-called tumours were

33   just as capable of distinguishing high-burden from low-burden mice as those established

34   manually. Overall, our strategy is quicker, helps control experimenter bias and yields a greater

35   wealth of information about each tumour, thus providing a convenient route to getting

36   consistent and reliable results from a study.

37

38

39    **INTRODUCTION**

40    Human colorectal cancer is a major contributor to both disease and death in the Western

41    world, with approximately 1.36 million cases diagnosed in 2012 [1]. Due to the massive impact

42    of colorectal cancer worldwide, many animal models have been created to understand this

43    disease and test potential treatments. Mutations in the Wingless/Int-1 (Wnt) pathway are

44    commonplace in human colorectal cancer [2]. The Adenomatous polyposis coli (APC) protein is

45    part of the canonical Wnt pathway, which is strongly conserved across many species,

46    including humans and mice. APC promotes the destruction of ß-catenin and prevents Wnt

47    signalling. Interestingly, the *Apc* gene is mutated in over 80% of colorectal cancer cases, as

48    well as in some breast cancers [3]. One of the *Apc* mutations is particularly noteworthy, as it

49    causes Familial Adenomatous Polyposis [4]. This hereditary genetic disease causes thousands of

50    polyps to form in the colon of the patient, which will invariably lead to colorectal cancer if

51    that patient is not screened and treated.

52    The *Apc^{Min}* mouse is a widely used model of spontaneously occurring intestinal tumours that

53    closely model human Familial Adenomatous Polyposis[5]. *Apc^{Min}* mice have been highly

54    valuable in demonstrating key mechanisms in colorectal cancer, for example, the importance

55    of Vascular Endothelial Growth Factor in the initial growth of intestinal tumours [6], the role of

56    COX-2 in adenoma formation [7], and the role of IL-33 in promoting tumorigenesis by

57    modifying the tumor immune environment [8]. *Apc^{Min}* mice produce an inactive, truncated APC

58    protein due to a mutation leading to a premature stop codon in the *Apc* gene [9]. This functional

59    loss in *Apc^{Min}* mice favours aberrant cell growth and, ultimately, spontaneous adenoma

60    generation in the mouse intestinal tract. Adenomas continue to grow throughout the mouse's

61    life, eventually causing bleeding, anaemia, and death, suggesting that tumour size, rather than

3

62    tumour count, may be a relevant metric.

63    Despite the wide use of the $Apc^{Min}$ model, there is no standardized technique to quantify

64    adenoma burden in these mice. Most papers rely on complex protocols and report only on

65    manually-counted adenoma numbers, or numbers and areas in selected areas of the intestinal

66    tract, although some also include information on adenoma location and size. However, high

67    quality semi-automated methods are now becoming available to facilitate the identification of

68    tumour lesions in histological images[10], or guide the visual classification of macroscopic

69    tumour lesions including melanomas in patients[11]. Therefore, these methods can offer rapid

70    and objective tumour identification in a broad range of situations.

71    In this paper, we describe a protocol for preparing standardised, photography-based images of

72    mouse small intestine (SI), large intestine (LI) and caecum; a new ImageJ[12] software macro

73    called *FeatureCounter* that automatically identifies tumour-like features in the SI images and

74    extracts measures such as area; and a machine learning pipeline for classifying these features

75    as true adenomas or not. We illustrate this strategy's performance on 120 mice of different

76    genotypes, age and sex. On the whole, our approach extracts a more detailed picture of the

77    adenoma burden in mice in a standardized and reliable manner, enabling a rapid and more

78    sophisticated analysis of the experimental results.

79

80    **RESULTS**

81    **Adenoma enumeration approaches**

82    Unbiased and reliable evaluation of tumour burden is essential to the interpretation of the

83    results of any preclinical study addressing tumour biology and potential therapy. This is

84   normally achieved by blinding investigators to the treatment group and performing lengthy

85   manual quantification under a microscope. Nonetheless, individual variations in measurement

86   techniques make the standardization of results across different investigators difficult to

87   achieve. To overcome these limitations, we designed three new techniques to evaluate tumour

88   count and area in the SI of tumour-prone $Apc^{Min}$ mice. The three techniques differed in degree

89   of automation, in how "features" of interest were identified, and in how those features were

90   classified or "called" as true Adenomas or not. A diagrammatic representation of the steps and

91   approximate time taken to perform a traditional method and these three new techniques are

92   shown in **Fig. 1**. A summary of these approaches is provided below:

93   1.  The TRAD (Traditional) method involved dissecting the intestinal tract, longitudinally

94       opening the gut, spreading the tissue onto a petri dish or glass plate, and manually

95       enumerating tumours on fresh tissue using a stereomicroscope (**Supplementary Fig.**

96       **1**). The nature of these visually-identified tumours can be confirmed by standard

97       histological techniques as shown in **Supplementary Fig. 1**.

98   2.  The DRAW approach involved dissecting the SI and removing all fat tissue, opening it

99       longitudinally taking care to leave any visible tumours intact, and carefully spreading

100      the tissue flat on a suitable cardboard as detailed in the Methods section. This was then

101      photographed close-up with a white ruler in shot for scale, and the photo stitched

102      together and opened using the Java-based image processing programme ImageJ [12].

103      The image was scaled using the ruler, and the ImageJ 'freehand selection' function

104      was used to manually draw the margin of each of the visually-identified Ad. These

105      features were then measured and added up using the ImageJ's 'analyze particles'

106      function to generate adenoma numbers and area. The same approach was used also to

107     quantify tumours in the LI and caecum, after they were prepared similarly to the SI. In

108     this study, the DRAW approach identified no adenomas in the SI, caecum or LI from

109     control mice, indicating high researcher reliability when identifying tumours.

110     3.  The CALL approach followed the DRAW approach up to the full SI image opening in

111     ImageJ. At this point, the *FeatureCounter* macro was run in ImageJ to automatically

112     set the scale and outline the contour of interesting features that might be adenomas.

113     From here, a researcher manually located each feature and "called" (assigned) them as

114     'true Adenomas' (Ad) or 'not Adenomas' (nAd). The resulting information is used by

115     ImageJ 'analyse particles' function to calculate adenoma number and areas. Thus, the

116     CALL approach automatically identifies adenoma-like features that are then verified

117     by eye, providing a gold-standard training set for machine learning if required.

118     4.  Finally, the LDA approach used the *FeatureCounter* macro-identified features

119     generated using the CALL approach and Linear Discriminant Analysis (LDA, a

120     simple machine learning technique) to determine how to discriminate between Ad and

121     nAd features based on the feature measures. Once trained on a CALL dataset, this

122     method is fully automatic, and features can be delineated by *FeatureCounter* and then

123     classified as an Ad or nAd by the LDA.

124

**Photography and *FeatureCounter* can be faster than manual quantification**

126     We compared the time required to quantify SI tumours using the various approaches

127     described in **Fig. 1**. Preparing the SI for analysis using the TRAD approach took about 30

128     minutes. In contrast, the time to prepare and photograph one SI sample for all other

129     approaches took in total about 40 minutes, including sample dissection, washing,

6

130    photographing, and image stitching time. Similar quality images were obtained using either

131    fresh SI tissue or tissue that had been stored frozen and thawed before sample processing and

132    analysis. Use of frozen tissues added about 5-10 minutes to the total tissue preparation time,

133    but introduced a very useful experimental breakpoint option when immediate analysis was not

134    possible or highly inconvenient, as is often the case in survival studies.

135    The quantification of tumours using the TRAD approach, by visually quantifying tumours

136    under the dissecting microscope, took up to 60 minutes per sample depending on tumour

137    burden. Measurement of individual tumour sizes would add considerably to this time,

138    especially when the tumour burden is high. In the DRAW approach, tracing features by hand

139    in ImageJ took about 1 to 10 minutes per sample, again depending on tumour burden.

140    Running the *FeatureCounter* macro to automatically identify image features of interest took

141    about 15-30 seconds. Manually calling tumour features from the *FeatureCounter* macro's

142    features in the CALL approach took 1 to 5 minutes per sample, while the LDA approach

143    (assuming a streamlined processing pipeline) took only one minute to complete the analysis

144    across all 3188 features from 117 animals. It is immediately apparent that the main time gain

145    is in the ability to automatically identify and call features, which is highest on heavily tumour-

146    burdened mice. For low-burden mice, the extra preparation time would offset this gain;

147    however, the consistency and depth of data generated using the DRAW, CALL or LDA

148    methods may make the extra time investment beneficial compared to the TRAD approach.

149    Overall, the TRAD approach takes approximately 90 minutes per sample, the DRAW

150    approach 60 minutes, the CALL approach 50 minutes and the LDA method 45 minutes per

151    sample. **Figure 1** schematizes these four approaches along with time costs for each step of

152    each method.

7

153

154     **Tissue preparation and *FeatureCounter* True Positive Rate**

155     High quality tissue preparation is essential to tumour identification using the *FeatureCounter*

156     macro. **Figure 2A** shows a SI laid out on cardboard, before being bisected into two long

157     pieces which were then cut longitudinally and, using tweezers, opened out, spread flat with

158     smoothed edges, and cleaned with PBS to expose any adenomas present. A representative

159     image is presented in **Fig. 2B**. Tumours are visible as denser white areas on the blue

160     cardboard background. From these images, tumours were manually delineated by an

161     experienced researcher to generate the DRAW mask in **Fig. 2C**. Alternatively, the

162     *FeatureCounter* macro was used to automatically flag adenoma-like areas and generate a

163     mask as shown in **Fig. 2D**. *FeatureCounter* identified very few features from a good

164     preparation of control SI with no adenomas. Representative image and mask are shown in

165     **Fig. 2E** and **2F,** respectively. Common issues with tissue preparation and image analysis

166     include rolled edges, excess fat, patches of dried tissue, and light reflections which can all be

167     picked up as non-tumour features by the *FeatureCounter* macro (**Supplementary Fig. 2**).

168     These "false positive" image features can be largely avoided by first removing excess fat at

169     sample collection and then, during preparation, ensuring that the tissue edges are flat by

170     smoothing with tweezers, regularly moistening the samples once mounted, and finally

171     ensuring consistent camera and light placement during photography. Once the protocol is

172     learnt, it is relatively simple to avoid all these artifacts.

173

174     **Validation of tumour identification in the small intestine**

8

175    To ensure that our premise of identifying image features as actual adenomas was correct, we

176    carried out experiments where fresh SI tissue was spread on blue cardboard, analysed using

177    the DRAW method, and then used as a source of tissue for microscopic analysis. As shown in

178    **Fig. 3C** and **3D**, two putative adenomas were selected due to their relatively isolated location

179    away from other tumours in the same sample, removed using a scalpel, then formalin fixed,

180    embedded in paraffin, and stained with haematoxylin and eosin. **Figure 3A and 3E** show a

181    magnification of these adenomas. Microscopic images in **Fig. 3B and 3F** revealed a typical

182    morphology with thickened mucosa, glandular appearance and a sessile structure. This

183    appearance is characteristic of adenomas as described in $Apc^{Min}$ mice [5] and very similar to that

184    of $Apc^{Min}$ adenomas imaged in our Lab using standard methods such as Swiss rolling of

185    intestinal tissue (**Supplementary Fig. 1**).

186    As a further validation of the tumour-bearing status of $Apc^{Min}$ mice as determined using the

187    DRAW method, we compared spleen and body weight between groups of $Apc^{Min}$ mice and

188    their adenoma-free WT littermates, which were sacrificed at the same time or shortly after

189    euthanasia of the last surviving $Apc^{Min}$ mouse in the same litter. A total of 49 mice, 27 $Apc^{Min}$

190    and 22 WT, were assessed. The average age of the $Apc^{Min}$ mice was 149 days with SD of 37,

191    while the average age of the WT controls was $177 \pm 21$ days. The results in **Fig. 4** show that

192    spleen weight was significantly higher in $Apc^{Min}$ mice compared to WT controls, while body

193    weight was lower. This is consistent with the reported anemia that develops in $Apc^{Min}$ mice

194    with increasing tumour burden, which in turn leads to splenomegaly [5]. All $Apc^{Min}$ mice

195    harboured numerous adenomas in the SI and a considerable tumour burden measured as total

196    tumour surface throughout the SI. No tumours were detected in the WT littermates.

197

9

198     **Linear Discriminant Analysis setup and feasibility**

199     We postulated that it would be possible to identify the true adenomas amongst the SI image

200     features delineated by *FeatureCounter* using data from the 22 shape and colour feature

201     measures provided by ImageJ. For example, one might expect adenomas to have rounder

202     shapes and slightly different colour than fat deposits and other non-tumour features. We thus

203     investigated the use of machine learning techniques for separating the true adenomas, "Ad",

204     from not true adenomas, "nAd". To provide a full training data for such a classifier, all the

205     image features from 120 mice with complete measures were called as Ad or nAd by a blinded,

206     experienced researcher using the CALL method. The dataset ultimately contained 3447image

207     features (1286 Ad, 1919 nAd, rest unclassified).

208     As a first analysis, we performed a PCA of the of the image feature data generated using

209     *FeatureCounter*. It was quickly apparent that there was segregation – though imperfect –

210     between the Ad and nAd classes (**see Supplementary Fig. 3**), suggesting that it was likely

211     that the LDA would be able to identify true Ad from nAd. We thus pursued the LDA to try

212     and automatically separate the feature classes based on the measure data.

213     Non-independence of observations can be a major problem in any statistical methodology not

214     designed to take it into account, as is the case for LDA. Here, observations (image features)

215     are nested within mice, in other words, many features may be found in the same mouse,

216     potentially causing non-independence of observations. This may be an issue if, for example, a

217     generally low-quality gut preparation led to bias in one or more image feature measurements

218     across all features from that mouse: the LDA learning would include this bias and thus fail to

219     generalize properly to all features. We thus used the PCA in **Supplementary Fig. 3** to

220     highlight potential mouse-level biases. As shown in **Supplementary Fig. 3**, the barycentres of

10

221    most of the 120 mice clustered at the center of the PCA, indicating no major mouse-level bias.

222    For animals with barycentres not clustering within this central area, SI photographs were

223    retrieved and scrutinized for signs of substandard preparation. We concluded that 3 mice had

224    photography of insufficient quality due to either poor sample preparation or inappropriate

225    camera settings. After excluding these, no such bias was observed. This result emphasises the

226    importance of standardising the tissue preparation and photography protocols to minimise

227    sample batch effects. After this step, 3188 features with proper CALL classifications (1279

228    Ad (40.1%) and 1909 nAd (59.9%)) from 117 mice were retained for training the classifier.

229

230    **Linear Discriminant Analysis performance**

231    As with any classification strategy, it is good practice to perform a validation experiment to

232    assess the classifier's stability and performance when faced with novel data; in other words,

233    we wanted to check that the LDA classification strategy would perform well when applied to

234    real-world experimental numbers. Using a "bootstrapping" random sampling with

235    replacement strategy (see LDA validation in Methods), we generated a total of 4000

236    validation datasets, computationally representing the equivalent number of 'experiments' of

237    normal $Apc^{Min}$ and WT animals, and each was used to train a separate LDA. We chose a

238    bootstrapping approach due to the relatively smaller size of our dataset, and selected with

239    replacement to ensure that population distribution was maintained for selections within each

240    validation dataset. For each validation set, feature-level performance indicators including

241    accuracy, TPR and Positive Predictive Value (PPV, or precision) and dataset-wide

242    performance indicators (such as the ratio of positive adenoma calls over true adenomas) were

243    derived for Ad and nAd on the full dataset, and compared to those obtained using LDA on the

11

244    full dataset, as described in the Methods.

245    The distributions of the feature-level performance indicators are presented in **Fig. 5A and 5B**.

246    The accuracy achieved for the full dataset was of 87%. The TPR (or the percent of the true Ad

247    / nAd correctly identified as such by the LDA) for the full LDA of Ad and nAd were close to

248    80% and 90%, respectively, indicating that the LDA was identifying correctly the majority of

249    both real Ad and nAd features. The PPV (or the proportion of the features identified as Ad /

250    nAd by the LDA that were correctly identified) of Ad or nAd were approximately 85% and

251    87% respectively, again showing good performance of the full LDA to classify Ad and nAd

252    features. Unsurprisingly, the LDA done on the whole dataset outperformed the majority of

253    bootstrapping datasets, perhaps indicating a slight overfitting when using the full dataset.

254    Nonetheless, all the indicators obtained on the validated datasets remained strong (indeed, the

255    worst performing indicator was Ad.TPR, with only 75% of values above 75%).

256    Importantly, the LDA performed very well when considering mouse-level performance

257    indicators. The Ad.ratio represents the ratio of the LDA-derived Ad count over the CALL-

258    provided Ad-count; the nAd.ratio is a similar indicator for nAd features. If the LDA was, in

259    practice, perfect, these ratios would be of exactly 1 (although it should be noted that the

260    converse is not true, and a ratio of 1 does not correspond to perfect performance). We

261    observed that the majority (the "most average 50%", as indicated by the gray boxes in **Fig.**

262    **5B**) of validation dataset Ad.ratios were between 0.919 and 1.051, with a median of 0.984,

263    while the whole dataset achieved an Ad.ratio of 0.941. The nAd.ratio performed arguably

264    even better, with the majority of validation dataset nAd.ratios being between 0.965 and 1.054

265    with a median of 1.010, compared to an overall dataset performance of 1.04. Full indicator

266    quantiles are given in **Tables S3 & S4**, with the 0% and 100% quantiles indicating the

267    minimum and maximum values, that is, 0% and 100% of datasets below the indicated values,

268    respectively. Taken together, these results indicate that despite the presence of a low

269    frequency of inaccurate tumour callings, the estimated mouse-level tumour count is highly

270    accurate.

271    Both the adenoma numbers and the total adenoma areas calculated by LDA showed high

272    correlation to the values obtained using DRAW or CALL. Concordance between LDA and

273    CALL was very good, in general with LDA obtaining only slightly less Ad counts than

274    CALL, as shown by the regression line & confidence region thereof in **Fig. 5C**. Quite

275    interestingly, the total Ad area was a much more accurate and consistent mouse-level measure

276    compared to the number of Ad, as evidenced by the tight regression line in **Fig. 5E**.

277    Unsurprisingly, the LDA approach yields mouse-level measures closer to those of CALL

278    rather than that of DRAW, as it was trained and used on adenoma callings from the CALL

279    approach (**Fig. 5D for counts and 5F for area**); however, all three approaches generate

280    similar tumour number and total tumour area measures, indicating a good predictive value

281    across the three methods.

282

283    **Adenoma area is a valuable measure of tumour burden**

284    Many previous papers have used total tumour number as the only measurement of tumour

285    burden to assess the effects of various treatments on $Apc^{Min}$ mice (for example, [13-15]).

286    However, this does not take into account the size of the tumours, which can also be highly

287    variable.

288    The automated method described here greatly facilitates the measurement of total adenoma

13

289    area. We investigated how appropriate total adenoma area is as a measure of tumour burden.

290    **Figure 6A** illustrates why total area should be measured and recorded: it presents two

291    samples with identical tumour counts, but largely different tumour sizes. Biologically, larger

292    tumours in the colon have been shown to be associated with shorter patient survival, showing

293    the importance of considering tumour size as well as number in response to treatments [16,17].

294    Furthermore, **Fig. 6B** illustrates that, in a sample of 63 mice evaluated using the DRAW

295    method, the average area of each tumour varied between different sections of the intestinal

296    tract, with tumours in the LI being significantly larger on average than SI tumours. We also

297    examined the correlation between total adenoma area and adenoma count in the SI. As shown

298    in **Fig. 6C**, the correlation between adenoma area and count was high, but the spread

299    increased with tumour number, thus reinforcing the utility of both measurements in evaluating

300    tumour status. Finally, we correlated the number and total area of tumours in the SI to spleen

301    weight, which represents a good surrogate measure of health status in $Apc^{Min}$ mice. Total

302    tumour area in the SI was a better correlate of spleen weight than tumour number (**Fig. 6D**),

303    even when excluding a potential outlier ($R^2$= 0.36 vs. 0.43). We argue that these observations,

304    taken together, demonstrate the need to evaluate tumour area in addition to tumour count.

305

306    **Utility of the total adenoma area measurements as assessed by LDA**

307    To evaluate the usefulness and comparability of the tumour burden measures established by

308    the DRAW and LDA approaches, we compared their power to discriminate between tumour

309    burdens in mice of different ages (147 days or younger versus older than 147 days at the time

310    of sacrifice, which are expected to have different tumour burdens) as a proof of principle.

311    These comparisons are illustrated in **Fig. 7A-D**. Younger mice show a significantly lower

14

312   number of Ad and total Ad area than older mice, in both the DRAW and LDA method, thus

313   validating that both manual and automatic classification of SI features can distinguish

314   between lower numbers and area of adenomas. Unsurprisingly, differences were much more

315   pronounced for the area measures than the counts, further illustrating the utility of area as a

316   measure of tumour load.

317

## DISCUSSION

318   **DISCUSSION**

319   We have developed a standardized protocol for first preparing and photographing mouse SI

320   samples, then for the manual (using the DRAW approach) or automatic (using an ImageJ

321   macro, *FeatureCounter*) identification of interesting image features (CALL approach), and

322   finally an LDA-based method for the automatic classification of said features as true

323   Adenomas or not Adenomas. Taken as a whole, these strategies allow for the consistent, rapid

324   and robust derivation of mouse-level tumour burden measures (both as adenoma count and

325   total adenoma area) for subsequent analysis.

326   Each of the steps in this standardized protocol works towards reducing technical, mouse-,

327   experimenter- and even institute-level bias and variability, thus increasing result

328   comparability and reproducibility. Additionally, the benefits are synergistic: as already

329   pointed out, more controlled sample preparation allows for more consistent feature

330   identification; and more consistently-defined features make feature classification easier. To

331   note, best results for training the LDA classifier would be expected by using training sets

332   called manually by either a single experimenter (as in this study), allowing the LDA to

333   "learn" the same cues as that experimenter, or by as many different experimenters as possible

15

334    (preferably across the same mice) allowing the LDA to "learn" the common cues to all.

335    Even with the best practice, however, the correct classification of image features by our LDA

336    step was not perfect. Most certainly, each step of our proposed method can be further

337    improved in future research. The use of diffuse lighting (such as a photography tent) at the

338    photography stage would minimise reflections that can be picked up as image features. The

339    *FeatureCounter* may be adjusted to detect less features in tumour-less images (for example,

340    by increasing the threshold size to ignore small features), while the automatic classification

341    may be adjusted or replaced with another machine learning methodology. For example, a

342    GLMNET algorithm [18] would allow the simultaneous selection and estimation of input

343    variable coefficients, at the very least leading to more consistent, if not more accurate, results.

344    More advanced machine learning algorithms, such as neural networks, are now being used in

345    the analysis of images from pathological samples, with new quantification approaches

346    becoming available (reviewed in [10]). In some cases, deep neural networks have been shown to

347    deliver classifications that are as accurate as those of a specialist, as in the case of skin lesions

348    [11]. Therefore, neural networks, of which LDA is a simple, single node example, have the

349    potential to provide better classification of images such as those generated in this study. In

350    any case, manual verification by an experienced researcher can be rapidly and easily

351    associated with any of the protocols described here, and would be most conveniently carried

352    out after LDA corrects the most evident misclassifications, such as those resulting from

353    imperfect sample preparation or photography– although these are relatively rare once the

354    technique is learned (3/120 in this study).

355    Regardless, our semi-automated strategy is faster, more reliable and also more flexible than

356    previously used methods. Samples can be processed and analysed while fresh, or can be

16

357   frozen and analysed later at a convenient time. Through the sample freezing step, "break

358   points" are introduced into the experimental workflow, *i.e.* points at which the

359   experimentation for a single sample can be suspended temporarily, while in traditional

360   methods each sample is often prepared and counted the same day. The reduced time cost in

361   tumour quantification can be another major benefit in the DRAW, CALL and LDA

362   approaches. It is thus immediately apparent that, beyond the added flexibility, our automated

363   strategy may earn a considerable sample preparation and counting time gain when many mice

364   – especially heavily tumour-burdened ones – are being assessed. Furthermore, the preparation

365   techniques are accelerated further when processing multiple samples at a time. Additionally,

366   the wealth of data is higher using these approaches compared to the TRAD count method,

367   where just tumour number, or cumulative tumour area in a small section of the intestine, is

368   assessed. We also note that once digitized, the photographic information can be stored almost

369   indefinitely, allowing the data to be revisited if need be, for example, after a *FeatureCounter*

370   update, or after the implementation of a new classification methodology, or for meta-analysis.

371   Finally, if the effort of generating a large LDA training set was not justified, the CALL and

372   DRAW methods can be rapidly implemented, and are still quicker, more reliable, and

373   producing more detailed data than the traditional method.

374   Several previous papers (for example, [13-15]) have only reported on total adenoma number,

375   using this as the lone tumour burden measure to assess the effects of various treatments on

376   *Apc*[Min] mice. However, this does not take into account the size or aspect of the tumours, which

377   can be highly variable. For our part, we believe that adenoma count certainly cannot be used

378   alone, as area can differ for identical adenoma counts, and its distribution changes between

379   different segments of the mouse intestinal tract. The reasons for these similarities and

17

380    differences are multiple. For example, early studies of the $Apc^{Min}$ mouse strain reported that

381    adenomas develop mostly during early life and up to puberty, and their numbers did not

382    increase after 100 days of age [19]. After this stabilisation in numbers, the adenomas have been

383    observed to instead grow in size [20], thus increasing tumour burden in a way not captured by

384    adenoma count alone. Additionally, significant size differences have been found in some

385    cases, demonstrating that area measures can provide additional information about treatments

386    or exacerbating conditions [21]. For example, therapies may be effective at controlling adenoma

387    growth without fully eradicating tumours, an effect that would be detected as decreased

388    burden with little or no change in tumour number. We thus conclude that adenoma area, and

389    potentially other measures, are of sufficient importance and value to warrant the use of new

390    methods to facilitate collection of such information. As adenoma number is still generated

391    using our approach, comparisons to previous studies remain possible. Of note, with our

392    ImageJ feature-based approach, it is possible to derive several aggregate measures (for

393    example, average adenoma greyscale value per mouse, as listed in parameters in

394    **Supplementary Table 2**) that might relate back to tumour burden or other biological

395    indicators of interest. Further research in this direction may yield interesting insights.

396    In conclusion, we propose a semi-automated method to rapidly quantify tumour number and

397    associated tumour burden measures that will help alleviate biases, along with reproducibility

398    and consistency problems, which currently hamper efforts to interpret results across the

399    $Apc^{Min}$ mouse literature. Our method is convenient, can be adapted to provide measurements

400    of several tumour characteristics, and will facilitate the use of $Apc^{Min}$ mouse intestinal

401    adenoma model in a variety of applications.

402

18

403 **METHODS**

404 **Animals**

405 C57BL/6J-$Apc^{Min}$ ($Apc^{Min}$) mice were purchased from The Jackson Laboratory (Bar Harbor,

406 ME) and bred in SPF conditions at the Malaghan Institute of Medical Research by mating

407 C57BL/6J-$Apc^{Min/+}$ males with Wild-type (WT) C57BL/6J ($Apc^{+/+}$) females. $Apc^{Min/+}$ and WT

408 offspring were identified by PCR and were both used in experimental conditions and pipeline

409 development. Water and standard laboratory chow were available *ad libitum.* All mice were

410 checked regularly for signs of anaemia and sickness, and were euthanized for tissue collection

411 if they developed pallor, low haematocrit ($< 20\%$), weight loss, slow movement and/or

412 hunched posture.

413 All experimental protocols were approved by the Victoria University of Wellington Animal

414 Ethics Committee, and were carried out in accordance with the Victoria University of

415 Wellington Code of Ethical Conduct.

416

417 **Tissue preparation**

418 Mice were euthanized and the entire intestinal tract was extracted and sectioned into the SI,

419 caecum and LI. Special care was taken to remove as much mesenteric fat as possible. Sections

420 were washed thoroughly using PBS, drained, and analysed immediately or frozen in 6 well

421 plates at $-80°C$ until further use.

422

423 **Photography**

19

424    For image analysis, SI tract sections were thawed (if frozen) and spread out in a thin

425    horseshoe shape on pieces of Steel Blue Germination paper (Anchor Paper Company, St

426    Pauls, MN, USA) approximately 25x10 cm in size. This colour was selected to enhance the

427    contrast between adenomas and the rest of the intestine. Once laid out on the paper, the SI was

428    cut into 2 equal pieces. Each piece was then cut longitudinally along the tube, opened and

429    edges spread flat using the edge of curved tweezers. Mucus and intestinal contents were

430    removed by spraying PBS on the tissue preparation, revealing any adenomas present. The

431    preparation was then photographed with a Panasonic Lumix G Vario DMX-G5W and a 45-

432    150 mm lens with additional 4x filter (Marumi, Japan), with a white ruler in shot. Multiple

433    pictures were taken and stitched together to reconstruct an image of the entire SI using the

434    software *Hugin* 2013.0.0 [22].

435    For the LI and caecum, a similar strategy was undertaken, where the tissue was placed on the

436    same type of Steel Blue Germination Paper as the SI, cut longitudinally (with multiple cuts

437    needed for the caecum), spread as flat as possible with special care taken to flatten tissue near

438    an adenoma in the caecum, and photographed with the white ruler in shot. Both the LI and the

439    caecum are small enough that they could be captured in one photograph.

440

441    **Manual delineation of tumours in images (DRAW approach)**

442    In order to enumerate and measure the area of tumours in the stitched images of SI, LI and

443    caecum, we used the Java-based image processing programme "ImageJ"

444    (https://imagej.net,[12]), which is freely available and able to analyse images in a variety of

445    formats. Full detail on tissue preparation, photography and analysis can be found at

446    https://gitlab.com/gringer/featurecounter/blob/master/Sample_Photography.pdf.

20

447    Images were scaled using a small macro and the white ruler in shot as a reference. ImageJ's

448    'freehand selection' function was then used to manually delineate visually-identified image

449    regions corresponding to adenomas. A scaled mask image was created using ImageJ's 'create

450    mask' function, and was analysed with the 'analyze particles' function to generate adenoma

451    numbers and measurements such as area. This is referred to as the "DRAW" approach.

452

453    *FeatureCounter*, **an ImageJ macro for the automatic identification of image features**

454    In order to automate the identification of image regions potentially corresponding to

455    adenomas from the photographs of intestinal sections as described in the DRAW approach,

456    we developed a more extensive ImageJ macro, called "*FeatureCounter*", focusing on SI

457    sections as these contain the large majority of the tumours that develop in $Apc^{Min}$ mice. First,

458    *FeatureCounter* subtracts the blue background, leaving a grey scale image. It subsequently

459    performs automatic thresholding, before despeckling the image according to the parameters

460    listed in **Table S1**. This leaves areas of over 0.2mmsq in size, or "image features" that are

461    potentially tumours. The "analyse particles" function within ImageJ measures 22 variables for

462    each feature: Area, Perimeter, Mean, StdDev, Mode, Min, Max, Median, Skew, Kurt, Major,

463    Minor, Angle, Circularity, AR, Round, Solidity, Feret, FeretAngle, MinFeret, IntDen, and

464    RawIntDen. The details of these measures and their processing can be found in **Table S2.**

465    *FeatureCounter* was optimised to work on the SI due to its smooth and regular surface. It

466    does not perform as well at quantifying tumours in the LI, where the surface of the intestinal

467    wall is ridged, or in the caecum, where the tissue does not spread out flat particularly well. As

468    the number of tumours in the caecum and LI rarely exceeds 3 (mean and SD of LI and

469    caecum is $1.81 \pm 2.00$ and $0.41 \pm 0.75$ respectively), these tumours can be quickly and

21

470    accurately quantified manually from photos using the DRAW approach. Therefore, further

471    work to optimize *FeatureCounter* performance on the LI and caecum did not seem warranted,

472    and was not pursued.

473

474    **Manual validation of tumour features (CALL approach)**

475    Image features identified by *FeatureCounter* can be manually validated. After running

476    the macro, a user can manually assign or "call" which features are tumours, referring to

477    them as "Adenoma" (Ad) or "not-an-Adenoma" (nAd) or, for unclear features, 'Not

478    Assigned' (NA). In our study, there were relatively few NA features, and they were

479    consequently excluded from further analyses. We refer to this approach as the "CALL"

480    approach.

481    Of further interest, the image feature measures obtained from *FeatureCounter* can be

482    leveraged in a machine learning algorithm to automatically determine which features are

483    tumours and which are false positives. Such a machine learning algorithm would require

484    a gold-standard "training dataset", i.e., a dataset of image features, their measurements,

485    and a prior validation of which features are indeed tumours or not, to learn tumour-

486    specific patterns. The CALL approach can be used to generate such a training set.

487

488    **Linear Discriminant Analysis (LDA) for automatic classification of image features**

489    Using the image feature measurements from *FeatureCounter* and a training dataset as

490    prepared using the CALL approach above, a machine learning technique can be used to

491    attempt to automatically separate tumour features from non-tumour features using the feature

22

492    measurements. LDA is one such supervised classification technique. It determines

493    discriminant functions – or the optimal linear combinations of the various input variables

494    (here: the 22 feature measures) – that can be used to classify statistical observations (here:

495    image features) into different classes (here: Ad or nAd). In our implementation, the *squares* of

496    the input variables were included as further input variables, as this allows quadratic

497    separations within the original variable space. All data were analysed within the R statistical

498    programming framework [23].

499    LDA is sensitive to several influences, including 1) extreme non-normality in input variable

500    distributions and 2) extreme outliers in input variables. For these reasons, it is recommended

501    to pre-process the input variables. We manually examined the distributions of the feature

502    measure variables per class, and applied log10 transformations, shifted log10 transformations,

503    and imposed certain filters, as described in **Table S2**.

504    The applicability of LDA to the transformed feature data was first evaluated by performing a

505    Principal Components Analysis (PCA) with package *FactoMineR* [24], the assumption being

506    that if the major axes of variability in the measurement data cannot segregate the classes even

507    partially, there is no point in performing an LDA and more advanced machine learning

508    techniques need to be used. The LDA was then performed using the *lda* function in the R

509    package *MASS* [25] for features with no missing values. A link to the R script used to run the

510    LDA can be found in the Supplementary materials. We then proceeded to investigate the

511    performance of our LDA at two levels, described below: at the feature level (checking

512    whether the classification performed well) and at the mouse level (checking whether, in

513    practice, the methodology allowed for accurate tumour counting and area quantification).

514

23

**LDA feature-level and dataset-level performance**

515

516     We compared the LDA's feature-level predictions to the adenomas selected using the CALL

517     method, which were considered "true" adenomas in this instance. We considered as indicators

518     of the LDA's performance the True Positive Rate (TPR, or Sensitivity, here defined as the

519     proportion of all true Ad that were also identified as adenomas using LDA), the Positive

520     Predictive Value (PPV, or the proportion of the LDA-identified adenomas that were indeed

521     Ad), and the Accuracy (the proportion of all features correctly identified as Ad or nAd).

522     Similar calculations were done for the nAd classes.

523     As indicators of dataset-level performance of the CALL and LDA adenoma callings, we

524     counted the number of Ad and nAd calls, and calculated the ratios of the number of LDA-

525     predicted Ad and nAd over the number of CALL-provided Ad and nAd (Ad.ratio and

526     nAd.ratio, respectively). An LDA with perfect performance would generate ratios of exactly

527     1, although a value of 1 is not necessarily indicative of perfect performance.

528

**LDA validation**

529

530     To assess the robustness of the LDA's results, we performed a large validation experiment

531     with a complex re-sampling scheme inspired by those of mixed modelling/multi-level models.

532     We chose to randomly sample mice (with replacements, *i.e.* a same mouse can be sampled

533     more than once) from the 117 with appropriate data, including all their image features in each

534     validation dataset. Mice continued to be sampled until a) at least 12 mice (about 10.3% of the

535     total) had been sampled, and until b) at least 750 features (23.5% of total) had been sampled.

536     Indeed, as the choice of the feature number parameter in the re-sampling scheme strongly

24

537    influences the performance indicators, we empirically determined that a minimal feature

538    count of 750 presented the best trade-off between sample size and indicator performance

539    (**Supplementary Fig. 4**). Additionally, to ensure some measure of class balance, only datasets

540    with a composition containing at least 30% Ad features and 30% nAd features were retained.

541    A total of 4000 validation datasets were generated (computationally representing the

542    equivalent number of 'experiments' of normal $Apc^{Min}$ and WT animals), and each was used to

543    train a separate LDA. For each validation LDA model, feature-level performance indicators

544    (Accuracy, TPR, PPV) and dataset-level performance indicators (Ad.ratio and nAd.ratio)

545    described above were derived using the whole dataset. For all indicators, we established their

546    quantiles of interest (0, 5, 25, 50, 75, 95, 100%) to compare to the values obtained on the full

547    dataset LDA.

548

549    **Statistics used to compare mouse-level results**

550    Comparisons of mouse data (weight, tumor numbers etc) used the Mann-Whitney U test or a

551    Kruskal-Wallis test followed by Dunn's multiple comparison test, and were performed using

552    Prism 8.0 GraphPad software.

553    To compare adenoma results at the mouse level (counts, total areas) obtained using different

554    methods (CALL and LDA), we used Deming regression, a statistical technique used for

555    comparing two measurement methods for a same quantity, where *both* measurements are

556    assumed to have measurement error (typical linear regression only assumes error in the

557    outcome variable). We used the *mcreg* function implemented in package *mcr* [26] assuming a

558    variance ratio of 1, and using bootstrapping (n=999, 'Bias-corrected and accelerated' method)

559    to obtain a regression curve confidence area.

25

560

**Availability of Data and Materials**

562    The *FeatureCounter* ImageJ macro is freely available to download from

563    https://gitlab.com/gringer/featurecounter/ together with instructions for photography, and

564    macro installation, some examples of tumour images, and the R code for running the

565    LDA. The datasets generated during the current study are available from the

566    corresponding author on reasonable request. Tumour images are available from Zenodo

567    repository. doi:10.5281/zenodo.3365777.

568

569

570   **REFERENCES**

571   1     GLOBOCAN. *Estimated cancer incidence, mortality and prevalence worldwide*

572         *in 2012*, 2012).

573   2     Half, E., Bercovich, D. & Rozen, P. Familial adenomatous polyposis. *Orphanet J*

574         *Rare Dis* **4**, 22, doi:10.1186/1750-1172-4-22 (2009).

575   3     Kinzler, K. W. & Vogelstein, B. Lessons from hereditary colorectal cancer. *Cell*

576         **87**, 159-170 (1996).

577   4     Nishisho, I. *et al.* Mutations of chromosome 5q21 genes in FAP and colorectal

578         cancer patients. *Science* **253**, 665-669 (1991).

579   5     Moser, A. R., Pitot, H. C. & Dove, W. F. A dominant mutation that predisposes

580         to multiple intestinal neoplasia in the mouse. *Science* **247**, 322-324 (1990).

581   6     Korsisaari, N. *et al.* Inhibition of VEGF-A prevents the angiogenic switch and

582         results in increased survival of Apc+/min mice. *Proc Natl Acad Sci U S A* **104**,

583         10625-10630, doi:10.1073/pnas.0704213104 (2007).

584   7     Zhang, M. Z. *et al.* Inhibition of 11beta-hydroxysteroid dehydrogenase type II

585         selectively blocks the tumor COX-2 pathway and suppresses colon

586         carcinogenesis in mice and humans. *J Clin Invest* **119**, 876-885,

587         doi:10.1172/JCI37398 (2009).

588   8     He, Z. *et al.* Epithelial-derived IL-33 promotes intestinal tumorigenesis in Apc

589         (Min/+) mice. *Sci Rep* **7**, 5520, doi:10.1038/s41598-017-05716-z (2017).

590   9     Su, L. K. *et al.* Multiple intestinal neoplasia caused by a mutation in the

591         murine homolog of the APC gene. *Science* **256**, 668-670 (1992).

592    10    Wang, S., Yang, D. M., Rong, R., Zhan, X. & Xiao, G. Pathology Image Analysis

593           Using Segmentation Deep Learning Algorithms. *Am J Pathol*,

594           doi:10.1016/j.ajpath.2019.05.007 (2019).

595    11    Esteva, A. *et al.* Dermatologist-level classification of skin cancer with deep

596           neural networks. *Nature* **542**, 115-118, doi:10.1038/nature21056 (2017).

597    12    Schneider, C. A., Rasband, W. S. & Eliceiri, K. W. NIH Image to ImageJ: 25 years

598           of image analysis. *Nat Methods* **9**, 671-675 (2012).

599    13    Amos-Landgraf, J. M. *et al.* Sex disparity in colonic adenomagenesis involves

600           promotion by male hormones, not protection by female hormones. *Proc Natl*

601           *Acad Sci U S A* **111**, 16514-16519, doi:10.1073/pnas.1323064111 (2014).

602    14    Chae, W. J. & Bothwell, A. L. IL-17F deficiency inhibits small intestinal

603           tumorigenesis in ApcMin/+ mice. *Biochem Biophys Res Commun* **414**, 31-36,

604           doi:10.1016/j.bbrc.2011.09.016 (2011).

605    15    Maywald, R. L. *et al.* IL-33 activates tumor stroma to promote intestinal

606           polyposis. *Proc Natl Acad Sci U S A* **112**, E2487-2496,

607           doi:10.1073/pnas.1422445112 (2015).

608    16    Kornprat, P. *et al.* Value of tumor size as a prognostic variable in colorectal

609           cancer: a critical reappraisal. *Am J Clin Oncol* **34**, 43-49,

610           doi:10.1097/COC.0b013e3181cae8dd (2011).

611    17    Suzuki, C. *et al.* The initial change in tumor size predicts response and

612           survival in patients with metastatic colorectal cancer treated with

613           combination chemotherapy. *Ann Oncol* **23**, 948-954,

614           doi:10.1093/annonc/mdr350 (2012).

615    18    Friedman, J., Hastie, T. & Tibshirani, R. Regularization Paths for Generalized

616          Linear Models via Coordinate Descent. *J Stat Softw* **33**, 1-22 (2010).

617    19    Moser, A. R., Dove, W. F., Roth, K. A. & Gordon, J. I. The Min (multiple

618          intestinal neoplasia) mutation: its effect on gut epithelial cell differentiation

619          and interaction with a modifier system. *J Cell Biol* **116**, 1517-1526 (1992).

620    20    Kettunen, H. L., Kettunen, A. S. & Rautonen, N. E. Intestinal immune responses

621          in wild-type and Apcmin/+ mouse, a model for colon cancer. *Cancer Res* **63**,

622          5136-5142 (2003).

623    21    McAlpine, C. A., Barak, Y., Matise, I. & Cormier, R. T. Intestinal-specific

624          PPARgamma deficiency enhances tumorigenesis in ApcMin/+ mice. *Int J*

625          *Cancer* **119**, 2339-2346, doi:10.1002/ijc.22115 (2006).

626    22    Desile, M. *Hugin 2013.0.0*, 2013).

627    23    Team, R. C. *R: A language and environment for statistical computing.*, 2013).

628    24    Husson, F, Josse, J., Le, S. & Mazet, J. *FactoMineR: Multivariate Exploratory*

629          *Data Analysis and Data Mining with R.  http://CRAN.R-*

630          *project.org/package=FactoMineR.*, 2013).

631    25    Venables, W. N. & Ripley, B. D. *Modern Applied Statistics with S. Fourth*

632          *Edition. .* (Springer, 2002).

633    26    Manuilova, E., Schuetzenmeister, A. & Model, F. *MCR: method comparison*

634          *regression. http://CRAN.R-project.org/package=mcr.*, 2014).

635

636

637

**Author Contributions**

648     ALS developed and performed gut preparations, counted adenomas manually and classified

649     the automatically identified features, and wrote the manuscript. AATS carried out statistical

650     analyses of image features. KAW carried out adenoma validation and generated adenoma

651     data. SK and JY carried out TRAD quantifications. DAE developed the *FeatureCounter*

652     macro. FR supervised the project, provided support and suggestions for investigations, and

653     edited the manuscript. All authors provided suggestions and approved the final manuscript.

654

**Competing interests**

656     The authors declare no competing interests.

657

658

**Figure 1.  Schematic of the tumour measurement methods described in this paper.**
Flow chart illustrating each step needed to perform the TRAD, DRAW, CALL, and LDA intestinal adenoma identification methods described in this paper. The icons represent the tools required to perform each step; estimated time costs per step are indicated. Please refer to the Materials & Methods and Results for a detailed description of the workflow for each method.
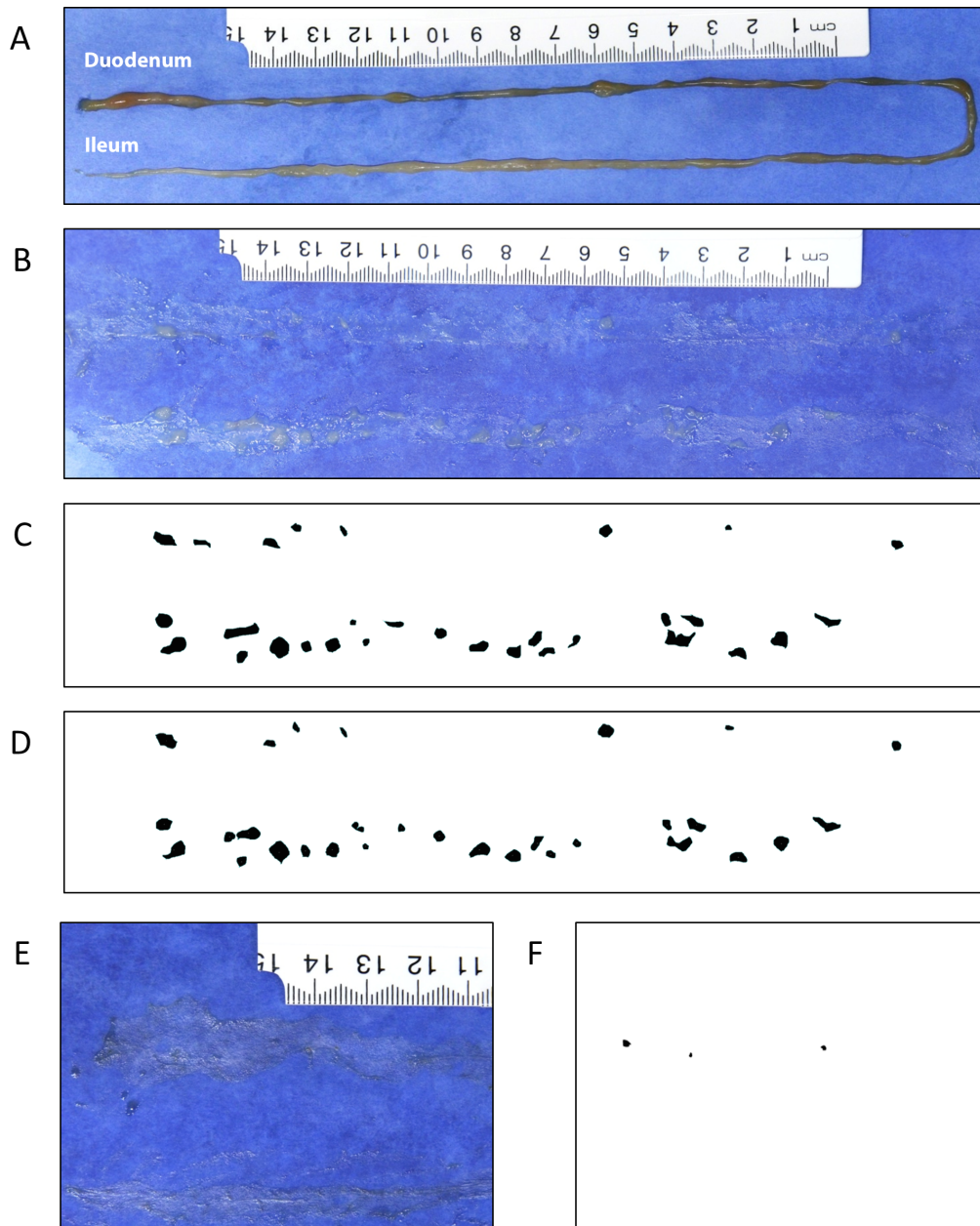
**Figure 2. The image features (adenomas) identified by the automated *FeatureCounter* macro mostly correspond to adenomas as identified using the manual DRAW method**

(**A**) Freshly collected SI from an *Apc^{Min}* mouse placed on blue cardboard. (**B**) The same SI after being cut longitudinally, spread and and cleaned with PBS to expose tumours. (**C**) *FeatureCounter*-generated tumour mask for the same sample. (**D**) Manually-generated tumour mask for the same sample. (**E**) A representative partial picture of a control SI. (**F**) *FeatureCounter*-generated mask, showing features picked up on the section shown in E. No additional features were picked up from the complete image.

**Figure 3. The image features identified using the DRAW method are adenomas.**

Fresh SI tissue was isolated from $Apc^{Min}$ mice, immediately set up on blue paper (**C**) and examined using the DRAW method in *FeatureCounter* to generate the mask in (**D**). Two relatively isolated features were chosen (marked by orange lines and magnified in **B, E**) excised from the paper support using a scalpel, and processed by formalin fixation, paraffin embedding and H&E staining to generate the images in (**A**) and (**F**). Data are from one of 3 mice and 7 SI tumours that were similarly treated and analysed.

**Figure 4. Spleen weight, body weight, number of SI tumors and their total area differ significantly between *Apc<sup>Min</sup>* mice and their WT littermates.**

*Apc<sup>Min</sup>* mice (n=27, 13 females and 14 males) were sacrificed when anemic and their body and spleen weights were determined. SI tumor numbers and total area were determined as shown in Figure 2 using the DRAW method. WT littermates (n=22, 9 females and 13 males) were sacrificed together with, or soon after, the last surviving *Apc<sup>Min</sup>* littermate. Average ages ± SD were 149 ± 36 days for *Apc<sup>Min</sup>* mice, and 177 ± 21 days for WT controls. Bar graphs show mean ± SD, each dot represents one mouse. P values were calculated using a Mann-Whitney test, ****: p<0.0001.

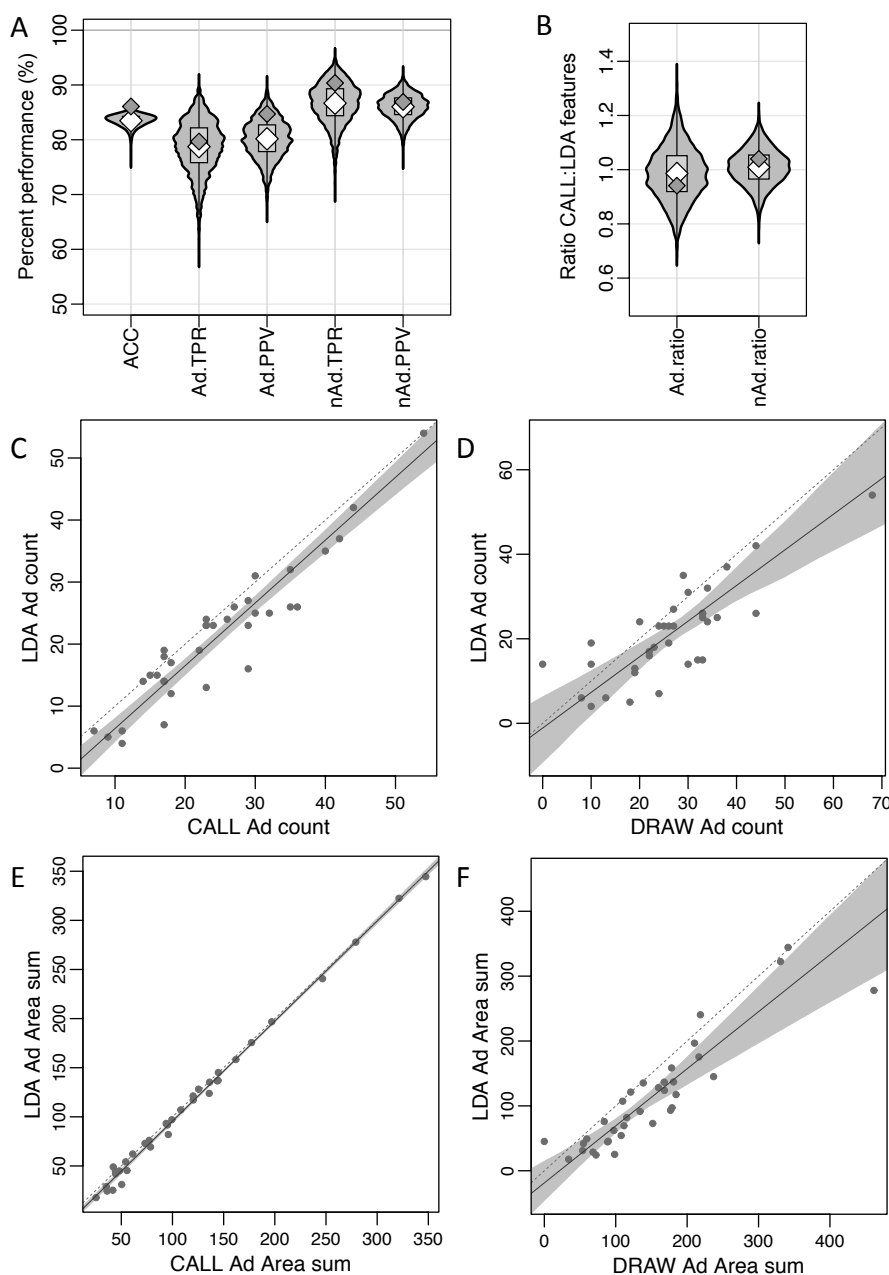**Figure 5. LDA predicts $Apc^{Min}$ tumour count and area with similar accuracy to the CALL and DRAW approaches.**

(**A, B**) Violin plots illustrating the distribution of the selected LDA performance indicators across the 4000 cross-validation datasets from 117 mice, each including 750-959 image features, when compared to the CALL-defined adenomas. The light grey violins are representative of the distribution of values obtained across the CV datasets; central grey boxes indicate the middle 50% of values; white diamonds represent median values for the CV datasets; dark grey diamonds represent the values observed in the full LDA. (**A**) shows Accuracy (ACC); Ad True Positive Rate (TPR, or sensitivity); Ad Positive Predictive Value (PPV); nAd TPR (or specificity); and nAd PPV distributions.

(**B**) The Ad.ratio is the ratio between the number of CALL Ad and LDA Ad, with a value of 1 indicating a perfect match. The nAd.ratio is determined similarly for nAd features.

(**C-F**) Deming regression plots comparing mouse-level adenoma number and total area values obtained through different approaches, for 35 mice. (**C**) compares adenoma counts generated by the LDA and CALL methods, (**D**) compares adenoma counts generated by the LDA and DRAW methods, (**E**) compares total adenoma area generated by the LDA and CALL methods, and (**F**) compares total adenoma area generated by the LDA and DRAW methods. Each dot corresponds to one mouse. Dotted grey line represents equality between measures. Solid grey line represents the regression line. Shaded grey area represents 95% confidence interval around the regression line.
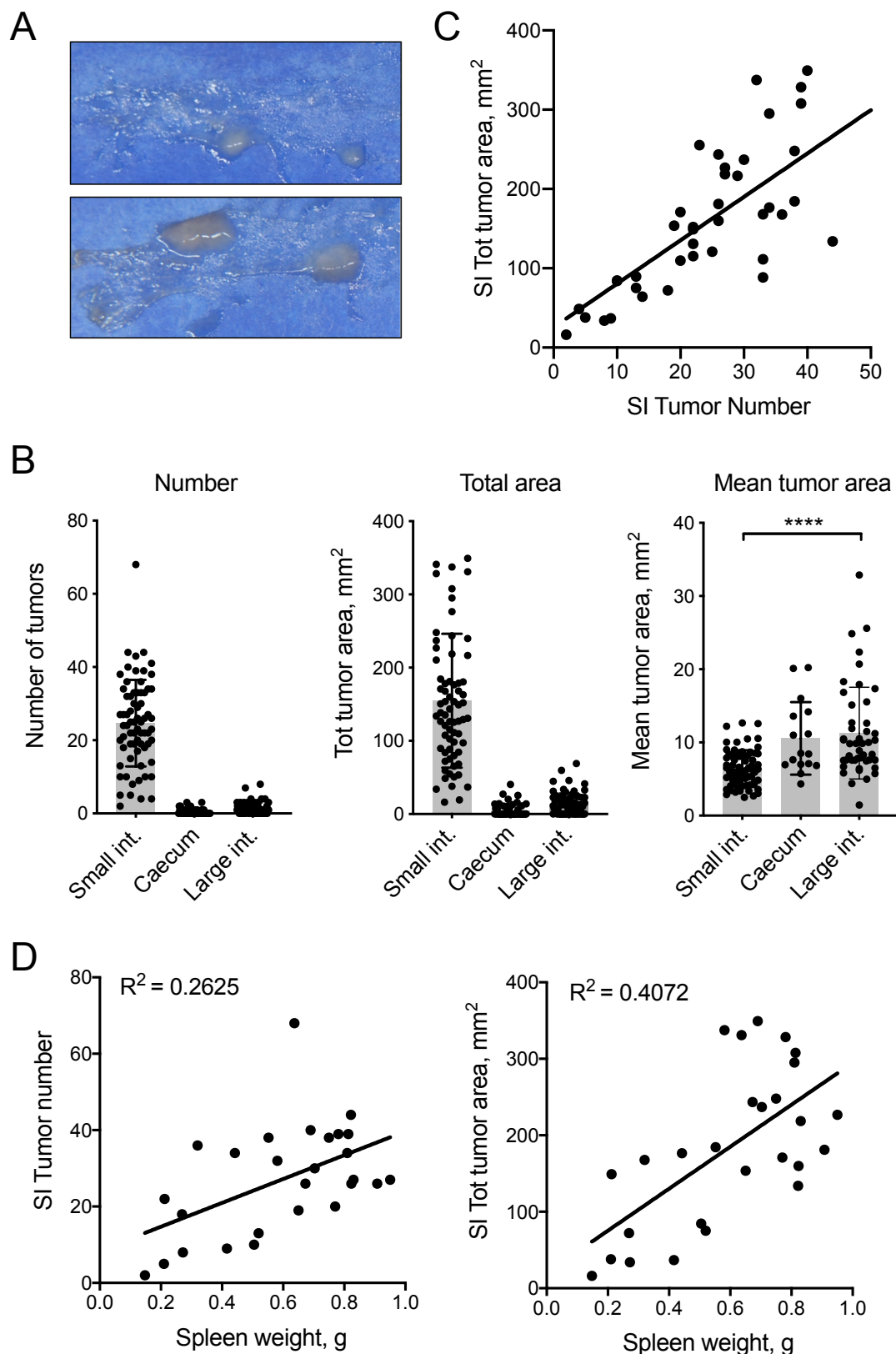
6

**Figure 6. Total tumor area is an informative measure of tumour burden in $Apc^{Min}$ mice.**

(**A**) Duodenal samples from two $Apc^{Min}$ mice, each with two tumours. Note the large difference in

tumour sizes between the two samples. (**B**) Bar graphs show the mean Number, Total area and Mean tumor area of tumours in different locations of the intestinal tract, +/- SD. Tumors were identified and measured using the DRAW method in a sample of 70 Apc$^{Min}$ mice. Each dot represents a single mouse. ****: $p<0.0001$ as determined using a Kruskal-Wallis test with Dunn's multiple comparison test. (**C**) Correlation between tumour count and size in LDA-called features in the SI of 35 mice. The dotted line represents the regression line. (**D**) Linear regression analysis of spleen weight vs. SI tumor number (left panel) or total area (right panel) in the SI of 27 mice for which spleen weight was available. Each dot represents one mouse. Data are from Figure 4.
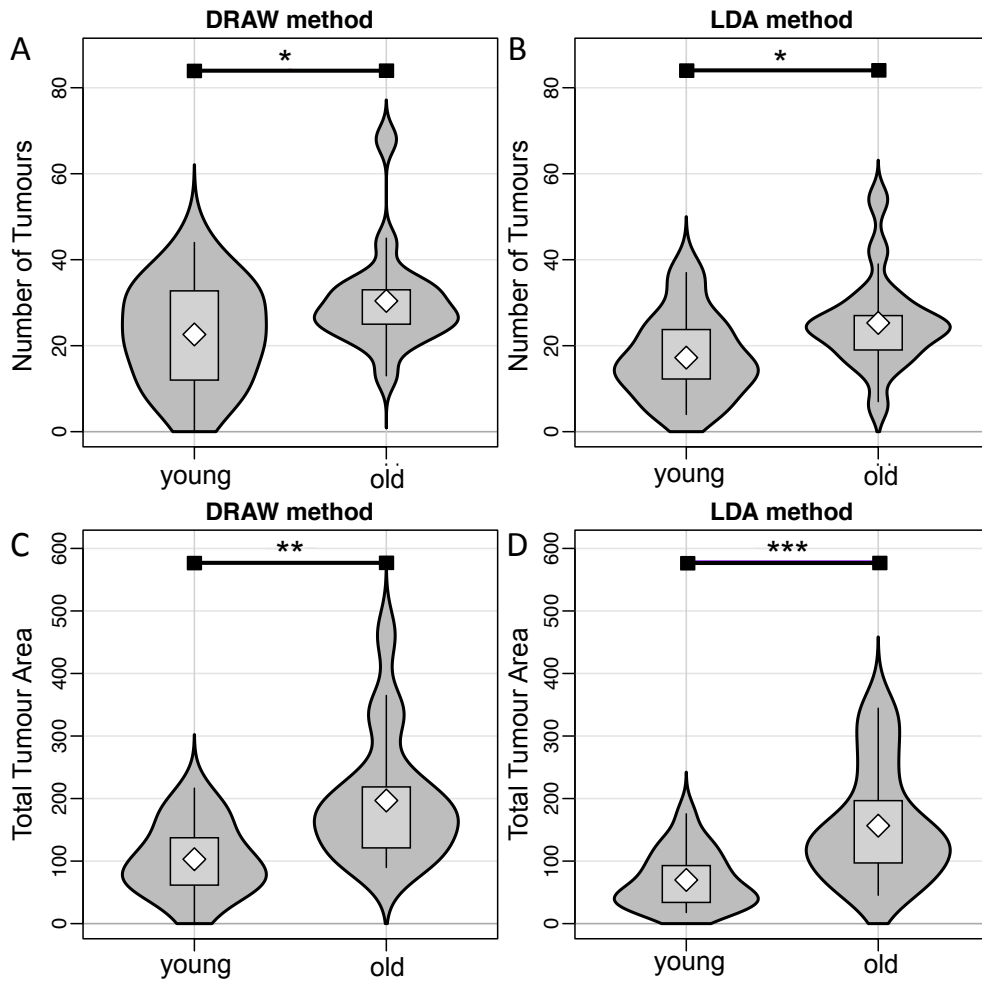
**Figure 7. The DRAW and LDA methods both differentiate tumour number and total tumour area in young vs. old mice.**

35 $Apc^{Min}$ mice were sacrificed when anaemic and then split by age: 'Young' (n=18) range from 1-147 days, while 'Old' (n=17) range from 147-214 days. (**A, B**): Violin plots of the number of tumours enumerated by the DRAW method (**A**) and the LDA method (**B**). (**C, D**): Violin plots of total area of tumours calculated by the DRAW method (**C**) and the LDA method (**D**). Stars indicate significance at the 5% level for approximate one-tailed Mann-Whitney-Wilcoxon tests (*: $p<0.05$, **: $p<0.01$, ***= $p<0.001$).