1  Automated Interpretation of Blood Culture Gram Stains using a Deep Convolutional Neural

2  Network

3

4  Running Title: Gram stain interpretation with deep learning

5

6  Kenneth P. Smith†[a,b], Anthony D. Kang†[a,b,c], and James E. Kirby[a,b#]

7

8  [a]Department of Pathology, Beth Israel Deaconess Medical Center, Boston, MA

9  [b]Harvard Medical School, Boston, MA, USA

10  [c]United States Army Medical Department Center and School, Fort Sam Houston, TX.

11

12  †These authors contributed equally to this work

13

14  [#]Corresponding Author

15  James E. Kirby

16  Beth Israel Deaconess Medical Center

17  330 Brookline Avenue - YA309

18  Boston, MA 02215

19  jekirby@bidmc.harvard.edu

20  Phone: 617-667-3648

21  Fax: 617-667-4533

22

24    **Abstract**

25    Microscopic interpretation of stained smears is one of the most operator-dependent and time

26    intensive activities in the clinical microbiology laboratory. Here, we investigated application of

27    an automated image acquisition and convolutional neural network (CNN)-based approach for

28    automated Gram stain classification. Using an automated microscopy platform, uncoverslipped

29    slides were scanned with a 40x dry objective, generating images of sufficient resolution for

30    interpretation. We collected 25,488 images from positive blood culture Gram stains prepared

31    during routine clinical workup. These images were used to generate 100,213 crops containing

32    Gram-positive cocci in clusters, Gram-positive cocci in chains/pairs, Gram-negative rods, or

33    background (no cells). These categories were targeted for proof-of-concept development as they

34    are associated with the majority of bloodstream infections. Our CNN model achieved

35    classification accuracy of 94.9% on a test set of image crops. Receiver operating characteristic

36    curve (ROC) analysis indicated a robust ability to differentiate between categories with area

37    under the curve >0.98 for each. After training and validation, we applied the classification

38    algorithm to new images collected from 189 whole slides without human intervention.

39    Sensitivity/specificity was 98.4/75.0% for Gram-positive cocci in chains/pairs; 93.2/97.2% for

40    Gram-positive cocci in clusters; and 96.3/98.1% for Gram-negative rods. Taken together, our

41    data support proof-of-concept for a fully automated classification methodology for blood-culture

42    Gram-stains. Importantly, the algorithm was highly adept at identifying image crops with

43    organisms and could be used to present prescreened, classified crops to technologists to

44    accelerate smear review. This concept could potentially be extended to all Gram stain

45    interpretive activities in the clinical laboratory.

46

47      **Introduction**

48              Bloodstream infections (BSI) are rapidly progressive infections with mortality rates up to

49      nearly 40% (1, 2). Each day delay in institution of active antimicrobial therapy is associated with

50      up to a ~10% increase in mortality (3, 4). Due to relatively low bacterial burden (<10 CFU mL$^-$

51      $^1$)(5), patient blood is pre-incubated in broth culture to detect presence of bacteria, typically by

52      semi-continuous measurement of $CO_2$ production or pH with an automated blood culture

53      instrument. If organism growth is detected, an aliquot of broth (now containing >$10^6$ CFU mL$^{-1}$)

54      is removed for Gram stain smear and subculture. The Gram stain provides the first critical piece

55      of information that allows a clinician to tailor appropriate therapy and optimize outcome (6).

56              Despite recent advances in automation in other stages of the BSI diagnosis process

57      (automated blood culture incubators and Gram staining systems) (7), Gram stain interpretation

58      remains labor and time intensive, and highly operator-dependent. With consolidation of hospital

59      systems, increasing workloads, and potential unavailability of highly trained microbiologists on

60      site (8), automated image collection paired with computational interpretation of Gram stains to

61      augment and complement manual testing would provide benefit. However, there has been a

62      dearth of scientific exploration in this area, and several technical difficulties need to be

63      overcome.

64              Practically, automated Gram stain interpretation requires both automated slide imaging

65      and automated image analysis. Although automated slide scanners and microscopes are being

66      used in anatomic pathology, for example, telepathology (9), their application in clinical

67      microbiology has been limited based on several technical challenges. First, Gram stained slides

68      are typically read using 100X objectives, greatly complicating image acquisition due to the need

69      for addition of oil during scanning. Second, microbiology smear material can adequately be

70    imaged only in a very narrow field of focus, a challenge for existing slide scanners. Third, Gram

71    stained slides exhibit ubiquitous and highly variable background staining. This background may

72    cause autofocus algorithms to target areas that are either devoid of bacteria or miss the

73    appropriate focal plane entirely. Image analysis to identify Gram stain characteristics presents

74    separate hurdles.  Importantly, background and staining artifacts, both fairly ubiquitous, often

75    mimics the shape and color of bacterial cells. Therefore, algorithms relying on color intensity

76    thresholding and shape detection will provide suboptimal accuracy.

77        Here, we provide proof-of-concept for automated, deep learning-based Gram stain

78    analysis. The major conceptual and technical innovations were twofold. First, we developed an

79    imaging protocol using an automated slide imaging platform equipped with a 40X air objective

80    to collect highly resolved data from Gram-stained blood culture slides. Second, image data were

81    used to train a convolutional neural network (CNN)-based model to recognize morphologies

82    representing the most common causative agents of BSI: Gram-negative rods, Gram-positive

83    cocci in clusters, and Gram-positive cocci in pairs or chains (1). CNNs are modeled based on the

84    organization of neurons within the mammalian visual cortex, and were applied here based on

85    their ability to excel in image recognition tasks without requiring time-intensive selective feature

86    extraction by humans (10). Our trained model was subsequently evaluated for accuracy in

87    comparison to manual classification.

88

89    **Results**

90        **Slide collection and manual classification.** Blood culture Gram stain slides prepared

91    manually during the course of normal laboratory operation were used for analysis. Slides were

92    selected based on the presence of any of the three most common morphotypes observed in

93    bloodstream infection: Gram-positive cocci in clusters, Gram-positive cocci in pairs and chains,

94    and Gram-negative rods. Less common morphotypes (e.g. Gram-positive rods or yeast) and

95    polymicrobial infections were excluded. To capture real-world variability, slides were not pre-

96    screened for suitability for automated microscopy or deep learning, and had characteristic slide-

97    to-slide variability in staining intensity, staining artifacts, and sample distribution. We

98    anticipated that inherent variability would pose a real-world challenge to slide classification

99    models.

100    **Automated image collection.** CNN-based deep learning models require large datasets

101    for training, typically at least on the order of thousands of images (and ideally at least an order of

102    magnitude more). Therefore, an automated microscopy image acquisition strategy was used. We

103    performed image acquisition on the MetaFer Slide Scanning and Imaging Platform

104    (MetaSystems Group, Inc., Newton, MA) based on a robust Gram stain-compatible autofocus

105    system, ability to sample multiple distributed positions on a slide to account for variations in

106    specimen distribution, and automated slide loading capability to enable high throughput slide

107    scanning.

108    Clinically, Gram stains are read under oil immersion. However, semi-continuous addition

109    of oil during automated microscopy was undesirable. In preliminary experiments with

110    uncoverslipped slides (data not shown), we determined that the 40x dry objective provided

111    sufficient resolution for machine-learning applications based on our prior experience (11).

112    Therefore, we selected use of the 40x air objective for image acquisition, thus avoiding the

113    requirement for oil immersion and allowing us to capture a larger field of view in each image.

114    **Deep convolutional neural network training.** For CNN training, a total of 25,488

115    images were automatically collected from distributed locations on 180 slides. A representative

116    image is shown in Fig. 1. This image demonstrates features typical of blood culture Gram stain

117    smears including: (A) intense background staining; (B) stain crystallization artifact; (C) diffuse

118    background staining; (D) individually resolvable, high-contrast Gram-negative cells; and (E)

119    individually resolvable, low-contrast Gram-negative cells. Of note, ubiquitous background

120    material was often similar in color, intensity, and/or shape to bacterial cells.

121          Highly experienced medical technologists can readily differentiate bacteria from this

122    background. However, it is prohibitively difficult to manually define computational rules for

123    Gram-stain classification that would adequately distinguish signal from noise in highly variable

124    Gram-stain preparations. Therefore, we chose instead to use a deep learning approach, more

125    specifically, a CNN, for image analysis. CNNs do not interpret raw images directly. Rather, they

126    consist of a number of layers, each of which convolutes regions of the image to detect specific

127    features. During each step of the learning process, a subset of images is presented to the network,

128    allowing function parameters to be changed such that the CNN identifies features important for

129    classification based on optimization of output accuracy. The final model is defined by a set of

130    weights and biases that control the flow of information through the network such that the most

131    discriminatory features in the images are used for classification.

132          Each CNN model has a unique architecture that differs in organization, function and

133    number of convolutional layers (10). The model used in our analysis, Inception v3, has

134    previously been shown to perform robustly on complex image classification tasks including

135    accurate classification of 1,000 different objects (12). The Inception v3 model is composed of a

136    series of small convolutional networks termed "inception modules" and was designed to be less

137    computationally intensive than comparable networks (13). Nevertheless, it is still a highly

138    complex model requiring weeks to train even with state-of-the-art computational infrastructure

139    (12). However, training the entire network is not always necessary. Many image classification

140    tasks can be addressed using pre-computed parameters from a network trained to classify an

141    unrelated image set, a method called transfer learning (14). To this end, we used an Inception v3

142    model previously trained to recognize 1,000 different image classes from the 2012 ImageNet

143    Large Scale Visual Recognition Competition  dataset (15), and re-trained the final layer to

144    identify our Gram stain categories of interest.

145        From an image analysis perspective, blood culture Gram stains are mostly background.

146    This excessive background increases the chance that a CNN will learn features during training

147    that are unrelated to bacterial Gram-stain classification. This is termed overfitting and results in a

148    model with high accuracy in classifying images on which it was trained (the training set), but

149    poor accuracy when presented with an independent validation set. Therefore, we enriched the

150    training data through use of selected image crops rather than whole slide images. A training crop

151    selection tool was created using the Python programming language which allowed the trainer to

152    select areas of an image containing bacteria with a single mouse click. This allowed us to train

153    our model on regions of images containing bacteria without inclusion of excessive background.

154        For model training (Fig. 2), we used our training crop selection tool to generate a total of

155    100,213 manually classified image crops from 180 slides. Training and validation accuracy were

156    indistinguishable (Fig. 2A), implying robust ability of the model to evaluate data on which it had

157    not previously been trained. It further confirmed success in minimizing overfitting. During

158    training, predictions made by our model were compared to the observed data, and differences

159    between these values were quantified using a metric called cross-entropy (16). In practice, low

160    cross-entropy indicates that the model fits the observed data well. Cross entropy decreased

161    during training and plateaued after 12,000 iterations (Fig. 2B). Additional training iterations

162   beyond what is shown in Fig. 2 did not reduce cross-entropy or therefore improve model

163   accuracy.

164        **Evaluation of model performance on a per-crop basis.** Our CNN outputs relative

165   probabilities that an image crop belongs to each of four categories of training data: specifically,

166   Gram-positive cocci in chains/pairs, Gram-positive cocci in clusters, Gram-negative rods, and

167   background (i.e., no bacteria) (17). Per convention (10), the class with the highest probability is

168   assigned as the predicted class. Using this method, we tested our model using a test set of image

169   crops not used during model training, and achieved a classification accuracy of 94.9%, providing

170   an initial estimate of model performance. However, this metric may be impacted by the fact that

171   the test set was not wholly independent of the training set, as it may still contain crops from the

172   same slide or images used in developing the training and validation sets.

173        Therefore, to rigorously evaluate ability of our model to generalize to an entirely

174   independent dataset, we evaluated performance on an evaluation set of 4,000 manually classified

175   image crops (n = 1,000 crops per class) from 59 slides that were not a component of the training,

176   validation, or test sets.  Here, we achieved a similar overall 93.1% image crop classification

177   accuracy. Importantly, the evaluation set also allowed us to calculate sensitivity and specificity

178   on a per-category basis. Sensitivity/specificity was 96.6/99.4% for Gram-positive clusters,

179   97.7/99.0% for Gram-positive chains, 80.1/99.4% for Gram-negative rods, and 97.4/93.0% for

180   background. Calculation of the area under the receiver operating characteristic (ROC) curve

181   (AUC) for each category (Fig. 3) further indicated robust ability to differentiate between

182   categories (AUC > 0.98 for all).

183        **Development of whole-slide classification algorithm.** To this point, we performed

184   classifications on manually selected cropped images based on category assignment using the

185    highest probability output from the classification. However, we hypothesized that it was not the

186    optimal way to interpret our results for whole-slide classification. Specifically, a whole slide

187    classification task differs from our evaluation experiments in that it would necessarily examine a

188    much larger number of crops that were not preselected and only consist of background. Given

189    that background may simulate bacterial cells (Fig. 1), we expected a greater likelihood of false

190    positive calls.

191        To test this possibility during whole-slide classification, we decided to set a very

192    stringent probability cutoff (0.99) for category calls to minimize false positives at the image crop

193    level and maximize specificity at the whole slide level. Using this stringent cutoff, 65.6% of

194    evaluated crops had a prediction with confidence of ≥0.99, and 99.6% of these were correctly

195    classified. Classification accuracy was 99.9% for Gram-positive clusters, 100% for Gram-

196    positive chains, and 97.4% for Gram-negative rods.

197        To investigate how this stringent cutoff would impact false-positive rate on a per slide

198    basis when applied to images cropped automatically, we collected 350 whole images containing

199    no visible cells and which were not part of the training, validation, and evaluation datasets.

200    Images were cropped into 192 non-overlapping crops (n = 67,200) using a custom Python script

201    and evaluated using our trained model with the classification threshold described above. For each

202    category, false positive rates were ≤0.006% on a per image crop basis. Based on an assumed

203    normal distribution of false positives calls, we set a minimal threshold for slide classification of 6

204    positive crops per category in order to achieve a desired ≤ 0.1% false positive whole-slide

205    classification rate.

206        Our whole-slide classification algorithm was then tested on 189 slides previously

207    classified manually by a microbiologist and not a component of the training, validation, test, or

208    evaluation sets. Each of 54 images scanned per slide was divided into 192 non-overlapping 146 x

209    146 pixel crops and evaluated using the parameters described above for a total of 10,368 crops

210    per slide. We first qualitatively evaluated performance on automated image crops. This was

211    achieved by writing a Python program (called TA for technologist assist) that would output

212    images corresponding to crop calls by the CNN allowing for specific review. Fig. 4 shows

213    examples of correctly classified image crops corresponding to each of the four classification

214    labels.

215        We then quantitatively evaluated our whole slide classification accuracy in comparison to

216    manual classification by constructing a table that shows each slide's manual classification and

217    corresponding automated prediction (Table 1). We found that bacteria were detected in 84.7% (n

218    = 160) of slides by our automated algorithm. For those slides where bacteria were detected, we

219    calculated classification accuracy, sensitivity, and specificity. Classification accuracy was 92.5%

220    across all categories. Sensitivity was >97% for Gram-negative rods and Gram-positive clusters.

221    Sensitivity was lower for Gram-positive chains, largely owing to misclassifications as Gram-

222    positive clusters across a relatively lower overall number of slides (n = 40). Further, manual

223    inspection of Gram-positive chains misclassified as clusters revealed that these slides were

224    somewhat ambiguous owing to substantial clumping of cells. Specificity for Gram-positive

225    chains and Gram-negative rods was >96%. Specificity was slightly lower (93.2%) for Gram-

226    positive clusters, again owing to misclassification of Gram-positive chains as clusters. Despite

227    qualitative difference in background staining, accuracy of slides from aerobic bottles (88.8%) or

228    anaerobic bottles (92.9%) was not significantly different (Fisher's exact test, $P > 0.05$).

229        Overall, the most common error was misclassification of slides as background,

230    representing 70.7% (n= 29) of all misclassifications. On manual review of images from these

231    slides, we found that 44.8% (n = 13) had insufficient crops with bacteria to make a positive call

232    based on our pre-established thresholds. We found an additional 48.3% (n = 14) had organisms

233    that were either out of focus or very low contrast, and of these, the majority (78.6%, n = 11)

234    contained Gram-negative organisms, as expected based on superficial similarity to background

235    material. The remaining 6.9% (n = 2) of slides contained highly elongated Gram-negative rods or

236    minute Gram-negative coccobacilli. Neither morphology was a component of our training set.

237    Gram stain category miscalls (n = 5) other than conflation of Gram-positive cocci in chains and

238    Gram-positive cocci in clusters, were related to a combination of poor representation of the

239    causal organism in crops and excessive background artifact.

240

241    **Discussion**

242         The Gram stain smear provides the first microbiological data to guide treatment for BSI.

243    Notably, earlier results are correlated with positive patient outcome (6). However, interpretation

244    of Gram stains is time intensive and strongly operator dependent, requiring a skilled technologist

245    for interpretation. Concerningly, the most recent survey from the American Society for Clinical

246    Pathology indicates that, as of 2014, trained microbiology technologist jobs in the United States

247    have a vacancy rate of ~9%, and nearly 20% of technologists plan to retire in the next 5 years

248    (8). This finding highlights the need for development of solutions to make the current work force

249    more efficient. However, there has been relatively little progress in automation of tests requiring

250    subjective interpretation such as the Gram stain.

251         Lack of progress in this area is related to technical issues with automated microscopy and

252    need for imaging interpretation algorithms that are robust to identifying rare organisms in the

253    presence of variable background. Here, we demonstrated that the MetaFer Slide Scanning and

254 Imaging Platform provides a robust automated image acquisition system, capable of providing

255 sufficient resolution for Gram stain analysis using a 40X dry objective. For such analysis, we

256 chose to use a CNN based on its ability to excel in image analysis tasks with minimal human

257 intervention. A summary of workflow for implementation, testing, and validation of our platform

258 is provided in Fig. 5.

259 This work adds to the examples of successful CNN use in several areas of image-based

260 diagnostics. These include detection of skin cancer (18); interpretation of echocardiograms (19);

261 and detection of metastatic cancer in lymph nodes (20) in which combined contributions of

262 pathologists and CNN increased sensitivity for diagnosis (21). A CNN has also previously been

263 used by our group for early prediction of antibiotic minimal inhibitory concentrations in

264 microscopy-based microdilution assays (11).

265 Importantly, CNNs improve in performance as more image data is added to the training

266 set. Unlike other machine learning models, however, training on more data neither increases the

267 size of a CNN model, nor the complexity of model implementation. Nevertheless, training of an

268 entire CNN model requires substantial computational infrastructure. Here, we took advantage of

269 an existing trained CNN and re-trained its final layers, a method called transfer learning (14, 18).

270 In this way, we were able to train and implement our model using a standard office computer

271 containing an Intel Core i7 CPU, 32GB RAM with no GPU (graphics processing unit, the

272 computational workhorse for image analysis).

273 Not surprisingly, implementation of the trained CNN for whole slide analysis using this

274 computer infrastructure was relatively slow. We therefore piloted whole slide classification using

275 a system containing an Nvidia GTX 1070 GPU. Though still underpowered compared to other

276 currently available GPUs, it improved whole-slide classification time by a factor of 6, resulting

277    in a classification time of ~9 minutes. The best available GPUs are markedly more powerful than

278    the GTX 1070 and are expected to provide even better performance (<5 minutes per slide), not

279    even considering the ability of CNN algorithms to distribute computations across multiple GPUs.

280           Overall, we found that our trained model performed well on whole-slide image

281    classification. Where cells were detected, we achieved overall classification accuracy of 92.5%

282    and specificity of >93% for all classification labels with no human intervention. The most

283    common classification error from our model was misclassification of slides containing rare

284    bacteria as background, representing the majority (70.7%) of all classification errors. In practice,

285    these misclassifications would be flagged for direct technologist review, making these low-

286    consequence errors. We also note that our sensitivity/specificity in whole slide image

287    classification accuracy was modestly lower than on a per-image-crop basis. This is likely due in

288    part to inclusion of slides with very few bacteria and therefore higher propensity for false-

289    positives. Optimization of data collection or slide preparation would likely bring our whole-slide

290    accuracy close to that of per-image-crop accuracy.

291           Our study had several limitations. As a proof-of-principle examination, we included only

292    the most common BSI pathogens and omitted several important, but less common bacterial

293    morphologies, largely due to limitation in availability of training data. However, given an

294    appropriate amount of training data, these could easily be incorporated into the Inception v3

295    model, which can distinguish 1000 different categories and will be a future goal. Similarly,

296    discrimination of polymicrobial infections could be incorporated by inclusion of "mixed"

297    categories into our algorithm.

298           We also recognize that there are several steps that could be taken to improve

299    classification. Foremost, the number of slides (and therefore image crops) used for training is

300    relatively modest and could be increased to improve CNN accuracy. In addition, our whole slide

301    scanning protocol was based on selecting pre-defined positions for imaging that were invariant

302    between slides. This contributed to inadequate sampling in a significant subset of slides, which

303    we believe was the greatest contributor to reduction in model accuracy. This hypothesis is

304    supported by the observation that misclassified whole slide calls were typically from slides with

305    very few bacteria or poor sample spread. Notably, to address this issue, it is possible with the

306    existing microscope platform to perform an automated rapid scan for areas of appropriate

307    staining intensity and thereby pre-select regions of the slide that are more likely to have

308    sufficient Gram stained sample for image acquisition.

309        Gram stain smear preparation is also expected to have a significant impact on automated

310    slide imaging. Here, we used slides prepared by technologists during the course of normal

311    laboratory operation. Slides exhibited a high degree of variability in smear area, thickness,

312    location, and staining intensity. We anticipate that standardization of these variables will

313    improve ability of an automated microscope to consistently sample microscopic fields with

314    evaluable organisms. Further, use of an automated Gram stain device for staining would also

315    increase reproducibility of staining characteristics and further enhance accuracy. We plan to

316    investigate all of these areas in the future.

317        We envision a potential role of our technology in augmenting technologist classification.

318    Given that manual interpretation of blood culture Gram stains by trained technologists are very

319    accurate (22-24), our model could be used to enhance productivity by selectively presenting

320    crops containing bacteria to local or remote technologists. This would increase efficiency of

321    classification by sparing the operator the need to manually locate fields of interest among a

322    preponderance of background. This would also conceivably reduce technologist read time from

323  minutes to seconds. Upon further development and intensive algorithm training, the platform

324  could potentially also be used as a fully automated classification platform with no human

325  intervention.

326      In the era of laboratory consolidation and limitations in the number of skilled

327  technologists (8), we believe our system could provide enhanced opportunities for rapid Gram

328  stain classification at the site of care or during understaffed shifts in conjunction with later

329  analysis at a central laboratory or day shifts. We further envision extension of CNN analysis to

330  other smear-based microbiological diagnostics in the parasitology, mycobacteriology, and

331  mycology laboratories. We believe that this technology could form the basis of a future

332  diagnostic platform that provides automated smear classification results and augments

333  capabilities of clinical laboratories.

334

335  **Materials and Methods**

336      **Slide collection and manual slide classification.** A total of 468 de-identified Gram-

337  stained slides from positive blood cultures were collected from the clinical microbiology

338  laboratory at Beth Israel Deaconess Medical center between April and July, 2017 under an IRB-

339  approved protocol. Slides were prepared during the course of normal clinical workup. No pre-

340  selection of organism identity, organism abundance, or staining quality was performed prior to

341  collection. Positive blood culture broth Gram stains included those prepared from both non-lytic,

342  BD BACTEC Standard Aerobic (n = 232) and lytic, BD BACTEC Lytic Anaerobic Medium (n =

343  196) (BD, Sparks, MD).

344      All slides were imaged without coverslips using a MetaFer Slide Scanning and Imaging

345  platform (MetaSystems Group, Inc., Newton, MA) with a 140-slide capacity automated slide

346    loader equipped with a 40x magnification Plan-Neofluar objective (0.75 Numerical Aperture,

347    Zeiss, Oberkochen, Germany). For each slide, 54 images were collected from defined positions

348    spanning the entirety of the slide. The first 279 slides collected were used in training, validation,

349    and evaluation of our deep-learning model. The remaining 189 slides were classified manually as

350    Gram-negative rods, Gram-positive chains/pairs or Gram-positive clusters using a Nikon

351    Labophot 2 (Nikon Inc., Tokyo, Japan) microscope equipped with a 100x oil objective. Results

352    were recorded for later use in evaluation of our whole-slide classification algorithm.

353        **Training a Deep Convolutional Neural Network.** A training dataset consisting of 146 x

354    146 pixel image crops was generated manually with the assistance of a custom Python script.

355    The script allowed crop selection, classification, and file archiving with a single mouse click

356    allowing large numbers of annotated crops to be saved in a short period of time in a manner

357    directly accessible to the deep learning training program. Each crop was assigned to one of four

358    classifications: Gram-positive cocci in pairs or chains, Gram-positive cocci in clusters, Gram-

359    negative rods, or background (no cells). Prior to training, the dataset was randomly divided into

360    three subsets: 70% of image crops were used to train the model, 10% were reserved for hold-out

361    validation during model training, and 20% were reserved for testing to evaluate model

362    performance after completion of training. We used a transfer learning technique based on the

363    Inception v3 convolutional neural network (CNN) architecture pre-trained on the ImageNet

364    Large Scale Visual Recognition Competition (ILSVRC) 2012 image database (12). We used the

365    Python language (version 3.5) and the TensorFlow library (25)(version 1.0.1) to retrain the final

366    layer of the model using a custom graphical user interface (GUI) controlling a modified script

367    ("retrain.py") found in the TensorFlow GitHub repository (25, 26). Training was performed

368    using mini-batch gradient descent (batch size 200) with Nesterov momentum (momentum = 0.9)

369    (27) and cross-entropy as the loss function (16). The initial learning rate was 0.001 and decayed

370    exponentially at a rate of 0.99 per epoch. The output layer was a 4-way softmax classification

371    which assigned probabilities to each of the four categories described above.

372          **Analysis of model performance on a per-crop basis.** Using our trained CNN, we

373    evaluated model performance on a per-image-crop basis using an evaluation set of 1,000

374    manually selected crops from each class (total crops = 4,000), all of which were independent of

375    the training, validation, and testing datasets. For each category, true positives were defined as a

376    crop correctly classified as the category of interest; false positives were defined as crops that

377    were incorrectly classified as the category of interest; true negatives were defined as crops

378    correctly classified as a category other than the category of interest; and false negatives were

379    defined as crops incorrectly classified as a category other than the category of interest.

380    Sensitivity and specificity were modeled as receiver operating characteristic (ROC) curves for

381    each classification label by varying the softmax classification thresholds required for positivity.

382    Sensitivity    was    defined    as    $\frac{\text{True Positive}}{\text{True Positive+False Negative}}$.    Specificity    was    defined    as

383    $\frac{\text{True Negative}}{\text{True Negative+False Positive}}$. Area under the ROC curve (AUC) was calculated for each label using

384    the trapezium rule as implemented in the scipy library (28). ROC curves were visualized using

385    the matplotlib library (29).

386          **Development of whole-slide classification algorithm.** False positive rates for

387    automatically cropped images containing only background were determined by analysis of 350

388    whole images from 40 different slides. Images contained no visible cells and were independent

389    of the training, validation, testing, and evaluation datasets. Each image was automatically

390    segmented into 192 non-overlapping crops of 146 x 146 pixels using a custom Python script

391    (total crops = 67,200) and classified with our trained CNN using a stringent cutoff for positivity

392 (cutoff = 0.99). If no label achieved a probability greater than or equal to the cutoff, the

393 associated crop was called background. False positive rates were recorded for each classification

394 label.

395       **Whole-slide classification.** Using the automated imaging protocol outlined in the

396 "Automated Image Collection" section, we evaluated whole slide classification accuracy using

397 images collected from 189 slides which were previously manually classified (outlined in the

398 "slide collection and manual slide classification" section). For each slide, a custom Python script

399 was employed to automatically divide each image of the 54 images collected from predefined

400 locations into 192 crops of 146 x 146 pixels. Each crop was evaluated by our trained deep-

401 learning model and probabilities assigned to each category (Gram-negative rods, Gram-positive

402 chains/pairs, Gram-positive clusters, or background) with a stringent cutoff for classification

403 (cutoff = 0.99). If no label met the classification cutoff, the crop was classified as background.

404       After classification of all crops from a slide, the category corresponding to the greatest

405 number of predicted crops was selected; however, only if the number of crops in the selected

406 category exceeded the number of expected false positives (calculated in the "Determination of

407 False Positive Rate" section). If none of the three label categories representing organisms were

408 selected based on these criteria, the slide was classified as background. All results were recorded

409 and used to construct a confusion matrix tabulation per convention in the deep learning field

410 (30). Whole-slide sensitivity and specificity were defined and calculated as in the "Analysis of

411 model performance on a per-crop basis" section. Classification accuracy for slides from aerobic

412 or anaerobic bottles was compared using Fisher's exact test with significance defined as $P < 0.05$

413 (JMP Pro version 13.0).

414

**References**

435  **References**

436  1.    Laupland KB. 2013. Incidence of bloodstream infection: a review of population-based

437        studies. Clin Microbiol Infect 19:492-500.

438  2.    Wisplinghoff H, Bischoff T, Tallent SM, Seifert H, Wenzel RP, Edmond MB. 2004.

439        Nosocomial bloodstream infections in US hospitals: analysis of 24,179 cases from a

440        prospective nationwide surveillance study. Clin Infect Dis 39:309-17.

441  3.    Schwaber MJ, Carmeli Y. 2007. Mortality and delay in effective therapy associated with

442        extended-spectrum beta-lactamase production in Enterobacteriaceae bacteraemia: a

443        systematic review and meta-analysis. J Antimicrob Chemother 60:913-20.

444  4.    Kang CI, Kim SH, Kim HB, Park SW, Choe YJ, Oh MD, Kim EC, Choe KW. 2003.

445        Pseudomonas aeruginosa bacteremia: risk factors for mortality and influence of delayed

446        receipt of effective antimicrobial therapy on clinical outcome. Clin Infect Dis 37:745-51.

447  5.    Wain J, Diep TS, Ho VA, Walsh AM, Nguyen TT, Parry CM, White NJ. 1998.

448        Quantitation of bacteria in blood of typhoid fever patients and relationship between

449        counts and clinical features, transmissibility, and antibiotic resistance. J Clin Microbiol

450        36:1683-7.

451  6.    Barenfanger J, Graham DR, Kolluri L, Sangwan G, Lawhorn J, Drake CA, Verhulst SJ,

452        Peterson R, Moja LB, Ertmoed MM, Moja AB, Shevlin DW, Vautrain R, Callahan CD.

453        2008. Decreased mortality associated with prompt Gram staining of blood cultures. Am J

454        Clin Pathol 130:870-6.

455  7.    Bourbeau PP, Ledeboer NA. 2013. Automation in clinical microbiology. J Clin Microbiol

456        51:1658-65.

457    8.    Garcia E, Ali AM, Soles RM, Lewis DG. 2015. The American Society for Clinical

458          Pathology's 2014 vacancy survey of medical laboratories in the United States. Am J Clin

459          Pathol 144:432-43.

460    9.    Meyer J, Pare G. 2015. Telepathology Impacts and Implementation Challenges: A

461          Scoping Review. Arch Pathol Lab Med 139:1550-7.

462    10.   LeCun Y, Bengio Y, Hinton G. 2015. Deep learning. Nature 521:436-44.

463    11.   Smith KP, Richmond DL, Brennan-Krohn T, Elliott HL, Kirby JE. 2017. Development of

464          MAST: A Microscopy-Based Antimicrobial Susceptibility Testing Platform. SLAS

465          Technol doi:10.1177/2472630317727721:2472630317727721.

466    12.   Szegedy C, Vanhoucke V, Ioffe S, Shlens J, Wojna Z.  Rethinking the Inception

467          Architecture for Computer Vision. https://arxiv.org/abs/1512.00567. Accessed Sept. 12,

468          2017.

469    13.   Szegedy C, Wei L, Yangqing J, Sermanet P, Reed S, Anguelov D, Erhan D, Vanhoucke

470          V, Rabinovich A. Going deeper with convolutions, p 1-9. *In* (ed),

471    14.   Shin H-C, Roth HR, Gao M, Lu L, Xu Z, Nogues I, Yao J, Mollura D, Summers RM.

472          2016. Deep Convolutional Neural Networks for Computer-Aided Detection: CNN

473          Architectures, Dataset Characteristics and Transfer Learning. IEEE Transactions on

474          Medical Imaging 35:1285-1298.

475    15.   Russakovsky O, Deng J, Su H, Krause J, Satheesh S, Ma S, Huang Z, Karpathy A,

476          Khosla A, Bernstein M, Berg AC, Fei-Fei L. 2015. ImageNet Large Scale Visual

477          Recognition Challenge. International Journal of Computer Vision 115:211-252.

478    16.   de Boer P-T, Kroese D, Reuven S, Rubinstein RY. 2005. A Tutorial on the Cross-

479          Entropy Method. Annals of Operations Research 134:19-67.

480    17.    Bridle JS. 1989. Probabilistic Interpretation of Feedforward Classification  Network
481           Outputs,   with   Relationships   to Statistical Pattern Recognition. Neurocomputing
482           F68:227-236.

483    18.    Esteva A, Kuprel B, Novoa RA, Ko J, Swetter SM, Blau HM, Thrun S. 2017.
484           Dermatologist-level  classification  of  skin  cancer  with  deep  neural  networks. Nature
485           542:115-118.

486    19.    Madani A, Arnaout R, Mofrad M, Arnaout R. 2017.  Fast and accurate classification of
487           echocardiograms               using               deep               learning.
488           https://arxiv.org/ftp/arxiv/papers/1706/1706.08658.pdf. Accessed Sept. 22, 2017.

489    20.    Litjens G, Sanchez CI, Timofeeva N, Hermsen M, Nagtegaal I, Kovacs I, Hulsbergen-van
490           de Kaa C, Bult P, van Ginneken B, van der Laak J. 2016. Deep learning as a tool for
491           increased accuracy and efficiency of histopathological diagnosis. Sci Rep 6:26286.

492    21.    Wang D, Khosla A, Gargeya R, Irshad H, Beck AH. 2016.  Deep Learning for
493           Identifying Metastatic Breast Cancer. https://arxiv.org/abs/1606.05718. Accessed Sept.
494           12, 2017.

495    22.    Samuel LP, Balada-Llasat JM, Harrington A, Cavagnolo R. 2016. Multicenter
496           Assessment of Gram Stain Error Rates. J Clin Microbiol 54:1442-7.

497    23.    Sogaard M, Norgaard M, Schonheyder HC. 2007. First notification of positive blood
498           cultures and the high accuracy of the gram stain report. J Clin Microbiol 45:1113-7.

499    24.    Rand KH, Tillan M. 2006. Errors in interpretation of Gram stains from positive blood
500           cultures. Am J Clin Pathol 126:686-90.

501    25.    Anonymous. 2017.    TensorFlow: An  open-source  software  library  for  Machine
502           Intelligence https://www.tensorflow.org/. Accessed Sept. 12, 2017.

503    26.    Anonymous.         2017.         TensorFlow         GitHub         Repository.

504         https://github.com/tensorflow/tensorflow. Accessed Sept. 12, 2017.

505    27.    Nesterov Y. 1983. A method of solving a convex programming problem with

506         convergence rate O(1/sqr(k)). Soviet Mathematics Doklady 27:372-376.

507    28.    Jones E, Oliphant E, Peterson P. 2001. SciPy: Open Source Scientific Tools for Python.

508         http://www.scipy.org/. Accessed Sept. 12, 2017.

509    29.    Hunter JD. 2007. Matplotlib: A 2D graphics environment. Computing in Science and

510         Engineering 9:90-95.

511    30.    Stehman SV. 1997. Selecting and interpreting measures of thematic classification

512         accuracy. Remote Sensing of Environment 62:77-89.

513

514

515     **Figures Legends**

516

517     **Figure 1. Representative image collected using our automated imaging protocol.** This image

518     shows several features characteristic of blood culture Gram stains including (A) area of intense

519     background staining, (B) artifact from stain crystallization, (C) diffuse background staining, and

520     individually resolved Gram-negative rods with (D) high and (E) low contrast compared to

521     background.

522

523     **Figure 2. CNN Model Training Results.** (A) Training and validation accuracy increased

524     exponentially, plateauing at ~95%. There was no observable difference in training and validation

525     accuracy, implying negligible overfitting during training. (B) Cross entropy is a metric used for

526     comparing model predictions to observed data. Lower cross entropy values indicate a better fit of

527     the model to the data. During training, we observed that cross entropy decreased to a final value

528     of ~0.1. Cross entropy plateaued at approximately 12,000 training iterations indicating that

529     additional learning was not possible without increasing the number of input images, a goal of

530     future work.

531

532     **Figure 3. Receiver operating characteristic (ROC) curve.** Curves were generated for each

533     category by varying threshold for positivity. Area under the curve is indicated in parentheses.

534

535     **Figure 4. Automatically classified crops.** Each image represents a correctly classified crop that

536     was automatically extracted from an image during whole slide classification. Rows of images

537     represent (A) background, (B) Gram-positive chains/pairs, (C) Gram-positive clusters, or (D)

538 Gram-negative rods. One practical application of the platform would be to present such organism

539 enriched images to a technologist to expedite smear review.

540

541 **Figure 5. Summary of CNN training and evaluation.** Prior to CNN training, we collected

542 images using an automated microscopy protocol (image example shown in Fig. 1). For CNN

543 training and preliminary testing, 100,213 image crops were manually selected, classified, and

544 randomly partitioned into training, validation, and test sets. Sizes of boxes correlate to relative

545 size of each data set. During iterative model training, accuracy was monitored using the training

546 and validation sets (Fig 2.). After completion of training, model accuracy was initially assessed

547 by quantification of accuracy on the test set (as discussed in text). However, the test set image

548 crops came from the same slides as the training set. We therefore further assessed performance

549 using a completely independent evaluation set to obtain a more reliable, real-world readout of

550 image crop classification accuracy and to generate receiver operating characteristics (ROC)

551 shown in Fig 3. Finally, we used a second independent dataset of automatically generated image

552 crops from 189 slides to evaluate whole slide classification accuracy. Each whole slide

553 classification was based on aggregate CNN categorizations of all image crops from a given slide

554 (examples of such crops are shown Fig. 4). Accuracy was determined in comparison to manual

555 slide interpretation (Table 1).

556

557 **Table 1. Confusion matrix of whole-slide classification results.**

| Human Classification | Predicted Classification (n) | | | | Sensitivity % % (CI)[a] | Specificity % (CI)[a] |
|---|---|---|---|---|---|---|
| | Gram-negative | Gram-positive pairs or chains | Gram-positive clusters | Background | | |
| **Gram-negative** | 51 | 1 | 0 | 17 | 98.1 (94.3-100) | 96.3 (93.7-98.9) |
| **Gram-positive pairs or chains** | 3 | 27 | 6 | 4 | 75.0 (60.9-89.0) | 98.4 (90.8-100) |
| **Gram-positive clusters** | 1 | 1 | 70 | 8 | 97.2 (93.4-100) | 93.2 (89.7-96.6) |

558

559 CI = 95% confidence interval

560 [a]Based on slides where bacteria were detected

561

A


B

```
┌─────────────────────────────────────────┐  ┌──────────┐
│           Training set                    │  │Validation│
│  (image example shown in Fig. 1)          │  │   set    │
└─────────────────────────────────────────┘  └──────────┘
┌──────────┐
│   Test   │              Train CNN and monitor
│   set    │              training/validation accuracy
└──────────┘                      (Fig. 2)

                    ╭───────────────────────────╮
                    │        Trained CNN         │
                    ╰───────────────────────────╯

                    ┌──────────┐  ┌──────────────┐
                    │Evaluation│  │ Whole slide  │
                    │   set    │  │classification set│
                    └──────────┘  └──────────────┘

     Test              Evaluation          Whole slide
accuracy (see text)  accuracy (see text)   accuracy (Table 1) in
 in comparison       and ROC (Fig. 3)       comparison to
 to manually          in comparison        manual slide classification
classified crops    to manually classified (example of automatically
                         crops              classified crops, Fig. 4)
```