1 The Hair Cell Analysis Toolbox: A machine learning-based whole cochlea analysis pipeline.

- 2 Christopher J. Buswinka, David B. Rosenberg, Artur A. Indzhykulian*
- 3 Mass Eye and Ear, Harvard Medical School.
- 4 * Corresponding author: Artur Indzhykulian, inartur@hms.harvard.edu

6 Abstract. Auditory hair cells, the whole length of the cochlea, are routinely visualized using light microscopy 7 techniques. It is common, therefore, for one to collect more data than is practical to analyze manually. There 8 are currently no widely accepted tools for unsupervised, unbiased, and comprehensive analysis of cells in an 9 entire cochlea. This represents a stark gap between image-based data and other tests of cochlear function. To close this gap, we present a machine learning-based hair cell analysis toolbox, for the analysis of whole 10 11 cochleae, imaged with confocal microscopy. The software presented here allows the automation of common 12 image analysis tasks such as counting hair cells, determining their best frequency, as well as quantifying single 13 cell immunofluorescence intensities along the entire cochlear coil. We hope these automated tools will remove 14 a considerable barrier in cochlear image analysis, allowing for more informative and less selective data analysis 15 practices. Keywords: hair cell, stereocilia, machine-learning, cochleogram, segmentation.

17 Introduction

5

16

18 Microscopy is an essential and very common tool in investigating the histology and pathology of the inner ear. Hearing loss phenotypes are often reflected in the histology of the sensory epithelium, compromising the organ 19 of Corti, and may be visualized with high magnification light microscopy. Collected micrographs can cover a wide 20 spatial area, capturing data on thousands of cells, often in three spatial dimensions. Auditory hair cells of the 21 22 cochlea are classified into two subtypes, inner hair cells (IHC) and outer hair cells (OHC), with varying geometric 23 locations, sizes, shapes, and protein expression 'fingerprints', all of which may vary with age and along the tonotopic axis¹. Each hair cell carries a bundle of actin-rich microvilli-like protrusions called stereocilia. The OHC 24 25 stereocilia bundles have a characteristic V-shape and are composed of thinner stereocilia as compared to those of IHCs. Each cell is a potential datum point which must be parsed from the image into a usable format for 26 analysis. Typically, this may have been done by hand, aided by a computer. While achievable for a small number 27 of cells, the slow speed and tedium of this analysis poses a significant barrier when faced with large datasets, 28 29 especially those generated through studies involving high-throughput screening². Alternatively, (and popular in 30 cochlear analysis) large datasets can be broken down into representative locations which can be analyzed in depth by hand. Often three small representative regions at the base, middle, and apex of the cochlear spiral³ are 31 32 chosen to reflect the tonotopic changes of biological variables. This approach may be enough to paint a general picture of the outcome of an experiment but scales poorly in comparison to other common techniques of cochlear 33 34 function testing, such as recording auditory brainstem responses or distortion product otoacoustic emissions, which could be easily collected, and analyzed, at an arbitrarily large number of frequencies^{4,5}. The disconnect in 35 the frequency resolution of histopathological analysis and cochlear function testing makes image analysis on a 36 single-cell level across the entire frequency spectrum of hearing desirable. To this end, an analysis software 37 must detect each hair cell, determine its type (IHC vs OHC), segment these cells to extract geometric and 38 39 fluorescence data, and assign a cell its best frequency based on its location along the cochlear turn.

40

49

Considerable work has been done on deep learning approaches for object detection. The predominant approach 41 42 is Faster RCNN⁶, a deep learning algorithm which guickly recognizes the location and position of objects in an 43 image. While designed for use with images collected by conventional means (cell phone or camera), there has been success in applying the architecture to biomedical image analysis tasks⁷⁻⁹. We apply this algorithm to detect 44 45 and classify hair cells at speeds orders of magnitude faster than manual analysis, while maintaining high 46 accuracy. High detection accuracy is paramount to generating cochleograms – a graph of hair cell count along the length of the cochlea - a common type of manual analysis reporting hair cell survival rates along the cochlear 47 coil most often used in human temporal bone studies. 48

50 Hair cell detection and classification alone, however, is not sufficient for the quantification of hair cell 51 fluorescence, which has applications in gene therapy, RNA scope quantification, and various fluorescent dye

52 loading assays. Instead, the pixels in an image must be assigned to the high-level cell object they represent in 53 a process known as *instance segmentation*. Cell segmentation has classically been tackled with a combination of manual thresholding and the watershed algorithm, a segmentation tool relying on intensity gradients between 54 objects^{10,11}. Often cells can be hard to detect in densely packed tissues with non-obvious delineation between 55 56 cells. Poor separation negatively impacts watershed segmentation, with optimal performance heavily reliant on manual fine-tuning, which is slow and can introduce bias. More recently, machine learning approaches have 57 58 been successful in supplanting the watershed algorithm for instance segmentation with increased accuracy¹²⁻¹⁵. Predicting the spatial embeddings of cells, a deep learning approach for generating instance segmentation 59 60 masks, is highly accurate and generalizes to a wide array of cell types¹⁶. This approach was popularized by the Cellpose algorithm¹⁷, and offers exceptional results for segmentation in two dimensions. 3D segmentation is 61 possible by applying Cellpose along different coordinate planes and using the 2D masks to generate the 3D 62 mask¹⁷. This is effective with isotropic and morphologically homogenous cells, however the algorithm's 63 performance in detecting hair cells was less impressive, likely due to nonobvious spatial separation. 64 65

Leveraging these recent deep learning advances, we present here a suite of tools for cochlear hair cell image analysis, the Hair Cell Analysis Toolbox (HCAT), a consolidated software for fully unsupervised hair cell detection and segmentation.

70 Results

71

69

72 Analysis Pipeline Overview: The software utilizes deep learning algorithms which have been trained to accept 73 input data of volumetric confocal micrographs of full cochlear turns, while also permissive to datasets containing 74 smaller cochlear fragments. The datasets contain at least two channels of information, anti-Mvo-VIIA labeling to visualize the hair cell body, and the phalloidin staining to visualize the stereocilia bundle, at a X-Y resolution of 75 289 nm/px (Figure 1). The user may either choose to run a cell detection analysis to generate a cochleogram of 76 77 inner and outer hair cells (from maximal projections of confocal data) or run a segmentation analysis (from 3D volumetric confocal data), which will extract fluorescence information of all segmented hair cells along the 78 79 tonotopic axis. Both algorithms scale and may be run on arbitrarily large confocal micrographs by repeatedly 80 analyzing small local crops of the image and merging the result to form a contiguous dataset. These local crops 81 are chosen to overlap such that each cell is guaranteed to be completely represented in at least one. Crops 82 which do not contain Myo-VIIA fluorescence above a certain threshold are skipped, increasing speed of large 83 image analysis and limiting false positive errors. In the case of a whole cochlear analysis, each cell is additionally assigned a best frequency via nonlinear regression and the Greenwood function as outlined below. The result 84 85 of every analysis is output as an associated csy data table to enable further data analysis or downstream postprocessing. 86

87

88 <u>Cochlear position detection:</u> For both, the segmentation and detection analysis, the software must determine the 89 place of each hair cell along the cochlea to analyze any location-based trends along the tonotopic axis. To do 90 this, we fit a Gaussian process nonlinear regression through the Myo-VIIA fluorescence image, effectively 91 treating each hair cell as a point in cartesian space. A line of best fit can be predicted through each hair cell and 92 in doing so approximate the curvature of the cochlea. We can then use the length of this curve as an 93 approximation for the length of the cochlea. For example, a cell that is 20% along the length of this curve could 94 be interpreted as one positioned at 20% along the length of the cochlea.

95

96 To optimally perform this regression and determine cell's best frequency, multiple preprocessing steps are 97 necessary. First, a maximum projection along the z axis of a previously predicted, whole cochlea segmentation 98 mask, is taken, reducing a three-dimensional volume to a single plane. The image is then down sampled by a 99 factor of ten using local averaging and converted to a binary image. Two final preprocessing steps are performed: 100 (1) binary hole closing which closes any gaps, and (2) binary erosion which reduces the effect of external 101 nonspecific staining. Each positive binary pixel of the resulting two-dimensional image is then treated as an X/Y 102 pair which may be regressed against.



HCAT







103 104

114

105 Figure 1. Overview of Algorithmic Approach. Top, Exemplar data used in this study. Cochlear coils 106 immunolabeled against Myo-VIIA (blue), stained with Phalloidin (red) and DAPI (grayscale), and genetically 107 expressing eGFP (green) were volumetrically imaged and input to software for analysis. Bottom, the user may 108 choose to run either a 3D segmentation analysis (blue), in which cells are volumetrically delineated from each other, 109 allowing for hair cell fluorescence intensity measurements on a single-cell level along the entire cochlear coil, or a 110 2D detection analysis (red), in which hair cells are counted and classified by type. These two analyses are paired 111 with a cochlear location algorithm (green) which assigns a best frequency to each hair cell, allowing for 112 cochleograms to be readily generated, or the fluorescence intensity measurements to be investigated as a function 113 of frequency or cochlear location.

115 The resulting image of an intact cochlea, when processed as is, would likely not form a mathematical function in cartesian space, and therefore be difficult to approximate. For example, the cochlea may curve over itself such 116 117 that for a single location on the X axis, there may be multiple clusters of cells at different Y values. To rectify this overlap, the data points are converted from cartesian, to polar coordinates by first shifting the points and 118 119 centering the cochlear spiral around the origin. From this, each X/Y pair can be converted to a corresponding angle/radius pair. In doing so, a gap is created as a cochlea is not a closed loop. This gap is detected, and these 120 points are shifted by one period, creating a continuous function. A Gaussian process¹⁸, a generalized nonlinear 121 function, is then fit to the spherical coordinates and a line of best fit is predicted. This line is then converted back 122 to cartesian coordinates and scaled to correct for the earlier down sampling. A linear interpolation of points is 123 124 performed to ensure each point is exactly 0.1% of the length of the cochlea. The apex of the cochlea is then inferred by comparing the curvature at each end of the line of best fit based on the observation that the apex has 125 a tighter curl when mounted on a slide. Next, the location of each hair cell along the cochlea as a function of its 126 total length (%) is determined by projecting the cell's center to the nearest point of the line of best fit. Finally, the 127 frequency at that location is then calculated using the Greenwood function¹⁹. 128

The final result is a curve of known length which tracks the hair cells of the cochlea. Any cell location may be 129 130 mapped to this curve as a function of the total cochlear length (%) and a best frequency calculated using the Greenwood function. Upon our careful examination, this process proves to be highly accurate. We also manually 131 mapped the cochlear length to cochlear frequency using a widely used *imageJ* plugin, developed by the 132 133 Histology Core at the Eaton-Peabody Laboratories, Mass Eye and Ear. The plugin is available for download at https://www.masseyeandear.org/research/otolaryngology/eaton-peabody-laboratories/histology-core. Over 134 135 eight manually analyzed cochleae, the maximum cell frequency error relative to a manually mapped best frequency result was under 10% of an octave, with the discrepancy between the two methods less than 5% for 136 137 most cells (60% of a semitone, see Supplementary Figure S2). If the curvature estimation fails, a manually annotated text file of points following the cochlear spiral may be additionally passed to the algorithm. A smoothing 138 spline fit through these points will be used instead for cell frequency estimation. It is also worth note, that both 139 hair cell detection and segmentation can be carried out by the HCAT software in isolated sections of the cochlea 140 in the absence of whole cochlear imaging. While this approach does not make use of the full functionality of the 141 software, it does provide researchers with a valuable tool for automated and unbiased cell counting and 142 fluorescence quantification, suitable for large datasets. 143

144

Hair Cell Detection: A maximum projection image containing cochlear hair cells can be iteratively evaluated in a deep learning model trained to detect and classify cochlear hair cells using this software. We leverage the Faster R-CNN deep learning model with a ResNet-50 backbone²⁰, trained on early postnatal cochlea labelled with phalloidin and antibodies against Myo-VIIA (**Figure 2**). While trained solely with these labels, the model can perform on cells labeled by other markers, provided the specimens contain both cytosolic hair cell, and stereocilia, labels.

151



152 153

Figure 2. Exemplar training data for the hair cell detection algorithm. Cochlear hair cells at varied frequency
 location were imaged using confocal microscopy. Images were immunolabeled against Myo-VIIA (blue) and stained
 with phalloidin (red). Bounding boxes were manually placed around inner (*white*) and outer (*green*) hair cell
 stereocilia bundles. These boxes and classifications were used to train the Faster R-CNN detection algorithm. *Scale bar*, 25 μm.

The deep learning model predicts three features for each predicted cell, a box encompassing the cell, a 158 classification label (IHC vs OHC), and a confidence score (Figure 3). After evaluation, two post-processing steps 159 are taken to refine the output and improve overall accuracy. First, some cells may be detected twice due to 160 redundancies in evaluating cropped subsections of a whole image. Redundant cell detections are removed and 161 the cell with the largest confidence score is kept. The second step is optional and relies on the estimation of the 162 cochlear spatial path outlined earlier, with cells whose distance away from this path exceeds a threshold value 163 164 are removed. This step can be enabled in cases with sub-optimal anti-MyoVIIA labeling outcomes with instances of non-specific labels away from the ribbon of hair cells along the cochlea, thus reducing the false-positive 165 detection rate. 166



Figure 3. Simplified overview of Faster R-CNN image detection pipeline. (A) Two-dimensional maximum projection images of three-dimensional Z-stacks (Myo-VIIA in *blue* and phalloidin in *red*) are encoded into a high-level representation by a trained convolutional neural network, schematized in *(B)*. (C) The encodings are transferred to an additional object proposal network which generates bounding boxes of predicted objects. Encoded crops, based on the predicted object proposals, are classified into outer and inner hair cells, and assigned a score. (D) Based on these scores and the overlap between boxes, objects are removed, resulting in the algorithm's best guess at the location and classification of every object detected on the image.

With the deep learning model optimally trained, followed by the post-processing steps, we see highly accurate performance on whole cochlea cell detection (Figure 4). Compared to cochleae analyzed manually we found a mean 95.2 ± 3.6% true positive accuracy for cell identification and a 0.5± 1% classification error (8 cochlear coils, each validation shown in Supplementary Figure S2). This algorithm is robust against tissue idiosyncrasies such as four rows of outer hair cells (Figure 4C) or atypical inner hair cell locations. Furthermore, when paired with cell frequency estimation, highly accurate cochleograms may be readily generated by the software (Figures 4F and 4G for IHCs and OHCs, respectively), taking just under a minute to run a full detection analysis.

Hair Cell Segmentation: To use the toolbox for the analysis of fluorescent signal in hair cells along the length of 185 the cochlear, rather than hair cell quantification, an alternative pathway and deep learning algorithm was 186 developed. In a similar approach to Cellpose^{16,17,21}, we elect to train a deep learning model to optimized spatial 187 embeddings of cells for instance segmentation. The U-Net architecture has been previously utilized successfully 188 for the performance of various biomedical segmentation tasks 14,17,22-29, and as such we employ a similar 189 architecture. Recent work on resource constrained segmentation mask has shown that recurrently applying 190 subsections of the U-Net architecture can improve performance without increasing model size³⁰. Building on this 191 idea, we recurrently apply each subsection of U-Net to improve performance while limiting the memory 192 193 complexity of the model. To account for the anisotropy in our dataset, we employ strided convolutions³¹ with different strides at each down sampling step of the architecture. A leaky ReLU³² has been shown to be 194 advantageous in the performance of residual neural networks^{33,34} and as such is chosen as the activation function 195 for our architecture. To maximize usability of the software and maximize flexibility in experimental design, we 196 train the model solely on the Myo-VIIA fluorescence signal (Figure 5) which were manually annotated in the 197 Amira software package³⁵. 198

An image volume with the Myo-VIIA fluorescence is first median filtered to remove noise and input into the machine learning model (**Figure 6A**) which has been trained to generate a set of spatial embedding vectors (**Figure 5B**) and a semantic segmentation mask (**Figure 6C**). From these, pixels likely belonging to a cell are

168

167

169 170

171

172

173

174

175

176



205 Figure 4. Validation output of the hair cell detection analysis. A validation output image is generated for each 206 detection analysis performed by the algorithm. Each image contains visual information on cell detection locations, 207 cell classifications and cochlear path estimation (if available). Additionally, each detected hair cell's ID, its location 208 along the cochlear turn (distance in μm from the apex), and best frequency are embedded. Such an image has 209 been automatically generated by the software here for an entire cochlea (A). While the vast majority of cells are 210 accurately detected (B, E), regions of poor performance are also highlighted in C and D. While the algorithm is 211 robust to four rows of outer hair cells, visual artifacts limit detection accuracy in C while low signal strength may 212 explain poor detection performance in D. When paired with single cell frequency estimation, accurate cochleograms 213 of inner hair cells (F) and outer hair cells (G) can be readily generated. This frequency estimation is highly accurate 214 with a maximum error on 8 different cochleae of 10% of an octave.

215

204

216 projected from their location in space via the predicted embedding vectors, forming clusters around the centroids object instances. When predicting the segmentation masks of new cells, their centroids are not known and 217 218 therefore must be inferred from the clusters of pixels (Figure 6D). We employ the DBSCAN clustering algorithm 219 to predict these clusters as it is robust against noise and scales well with large datasets. To improve performance, we found that down-sampling by a factor of two increases centroid detection speed with negligible loss in 220 221 detection accuracy. Pixels at the borders of objects tend to be more poorly projected to object centers, leading to sparse cluster formation. To reduce this effect, we disregard pixels near the edge of the semantic probability 222 map via binary erosion. From the resulting centroid prediction, probability maps for each instance a regenerated 223 based on equation (1). 224

$$\phi_k(e_i) = exp\left(-\frac{(e_{ix} - C_{kx})^2}{2\sigma_{kx}^2} - \frac{(e_{iy} - C_{ky})^2}{2\sigma_{ky}^2} - \frac{(e_{iz} - C_{kz})^2}{2\sigma_{kz}^2}\right)$$
(1)

226 To avoid double counting cells, we perform non maximum suppression on each instance probability map³⁶. To 227 ensure each segmentation mask predicted by the algorithm is of a complete cell, we remove any segmentation masks touching the edge of an evaluated image, in addition to removing any mask whose volume is below a 228 229 reasonable value of a cell chosen at 80% of the smallest volume in our training set. While the DBSCAN algorithm 230 is fast, it relies on tuning parameters to optimally perform. We have chosen these parameters to aggressively 231 reject loose clusters ensuring each detected cell is properly segmented, with the notable drawback of limiting 232 the ability of the algorithm to segment every cell. Finally, we disregard cells with mean Myo-VIIA fluorescence 233 signal intensity levels below 5% of maximum measured value, greatly reducing the number of false positive cell 234 mask detections.



241

242

Figure 5. Exemplar training data used to train the hair cell segmentation algorithm. To train a deep learning network to perform an instance segmentation task, 17 confocal Z-stacks containing cochlear hair cells immunolabeled against Myo7a (*top* panels) were manually segmented in three dimensions using the Amira Software package. The resulting 3D cell segmentation masks (*bottom* panels) were subsequently used for training of the deep learning model. Presented are single frames from the z-stacks representing individual tiles of the tile scans, with the tiles taken from a number of different datasets at various cochlear locations.



243

Figure 6: Overview of spatial embedding image segmentation. Volumetric confocal micrographs (A) are input to a machine learning model trained to predict spatial embedding vectors (B) in addition to a semantic segmentation mask (C). Pixels are projected using these vectors into clusters which form the "embedding space" (D) which are detected using a clustering algorithm. Each cluster gives rise to an object probability map which is used to form an instance segmentation mask of objects in the original image (E). The performance was fine-tuned by setting the thresholds in C and D (*red* "cutoff" lines). Every unique color in E represents a different hair cell.

250 In order to evaluate our method of instance segmentation we compared it against two other methods: 1) the 251 generalist Cellpose algorithm, and 2) a semantic segmentation approach coupled with the watershed segmentation algorithm. While Cellpose¹⁷ was not designed to natively segment volumes, the software provides 252 a pseudo 3D approach in which the algorithm is applied in 2D over different axis pairs of a volume. While the 253 254 segmentation performance is good, many cells were not segmented, in one case as low as only 17% of cells were segmented (Figure 7). As a generalist algorithm, Cellpose was not trained on examples of hair cells. With 255 256 the same training data and deep learning architecture, we trained a deep learning model to perform instance segmentation which was paired with a volumetric watershed algorithm. The results of this approach were poor, 257 258 with numerous detection and segmentation errors (Figure 7). We therefore find our method of instance segmentation outperform other methods of unsupervised volumetric biomedical instance segmentation. 259 260

While this algorithm relies on a strong Myo-VIIA fluorescence signal to perform, the resulting cell segmentation 261 mask can be further used to measure cell-specific fluorescence intensity data on any addition channels of 262 imaging information accompanying the dataset. As such, the main purpose of this segmentation process is to 263 264 delineate individual hair cells within their borders by assigning all pixels within a volume (i.e., voxels) representing 265 each hair cell to a single segmentation mask with a unique ID and illustrated with a unique color on Figure 8. These masks are then used to measure fluorescence intensity levels across different channels of imaging data 266 on a single-cell level along the cochlear coil. As each cell has a best frequency automatically inferred as 267 described above, not only do measurements of whole cell fluorescence intensity levels in 3-dimensional data 268 become automated, simple, and unbiased, they also allow to present these fluorescence measurements as a 269 function of location along the tonotopic axis of the cochlea (Figure 8). 270 271



Figure 7. Example segmentation outcomes of three different segmentation methods evaluated in this study.
 Comparison of multiple approaches at volumetric instance segmentation of confocal z-stack of cochlear hair cells.
 The spatial embedding approach more accurately detected hair cells compared to *U-Net + watershed* approach or
 the generalist *Cellpose* tool.



Figure 8: Single-cell fluorescence intensity analysis of hair cells along the cochlear coil. *Top*, exemplar maximum projection image of the Z-stack of a mouse cochlear coil analyzed with the hair cell segmentation algorithm. *Bottom*, Single-cell volumetric fluorescence intensity measurements of four fluorescence signals obtained from predicted instance segmentation masks. Paired with cell frequency estimation, the fluorescence intensity measurements are presented as a function of cochlear location from apex (0%) to base (100%) of the cochlea. Black dots represent single cell measurements, with an average fluorescence intensity value along the cochlear turn shown as a color line, with the population histogram presented on the right of each panel.

287 Discussion

288

278 279

280 281

282

283

284

285

286

Here we present the first fully automated cochlear hair cell analysis pipeline where users can enter multiple confocal micrographs of cochleae and quickly detect hair cells to generate cochleograms or segment hair cells enabling fluorescence intensity measurements on a single-cell level. We consider this to represent a considerable advancement in the unbiased analysis and quantification of hair cell imaging datasets.

293

294 Previous methods of extracting quantitative data from images of hair cells has been varied, and often adapted for a particular experimental outcome. For example, in hair cell survival studies it is often fastest to manually 295 count hair cells in order to generate cochleograms³⁷. While the accuracy that can be achieved in this task by a 296 trained individual is the target of any algorithmic approach, the significant time cost of performing such analysis 297 makes an automated solution desirable even with rare errors. Automated hair cell counting algorithms already 298 exist but are less feature rich, limiting adoption. One relies on the homogeneity of structure in the organ of Corti 299 and fails when irregularities, such as four rows of outer hair cells, are present³⁸. Another may count hair cells but 300 cannot differentiate between inner and outer hair cells³⁹. 301

302

303 The presence of a hair cell represents the most basic of quantification and therefore the quickest to do manually, although still represents a significant burden for large datasets. Image analysis of increasing complexity, such 304 as counting cells above a fluorescent intensity threshold, as is common in the study of protein expression levels 305 for example, take considerably longer to the point where they are no longer routinely performed over the entire 306 cochlea⁴⁰⁻⁴⁵. Analysis of fluorescence levels within hair cells represents even greater complexity. While some of 307 these analyses are being enabled more efficiently via general use software, such as ImageJ⁴⁶, the highly 308 309 specialized geometry of hair cells limits the effectiveness of these nonspecific tools for use in an automated, unsupervised manner. Furthermore, to carry out the highly favorable, more sophisticated quantification of 310 fluorescence in 3D cell volumes, discussed further below, requires volumetric segmentation of cells. Cellpose 311

represents the cutting edge of generalist cell segmentation, however, fails to achieve acceptable performance levels on our datasets.

314

The HCAT segmentation pipeline, presented here, offers fully automated volumetric segmentation of hair cells 315 along the entire length of the cochlea. This task would be prohibitively labor intensive to carry out manually. 316 however is carried out here at only marginally poorer segmentation performance when compared to manual 317 318 segmentation, removing the most significant barrier to analyzing cochlear-wide fluorescence within large datasets. The most accurate measurement of the fluorescence intensity signal within the cell are, three-319 320 dimensional, volume-normalized measurements. It is common for manual analysis methodology to be performed 321 on a two-dimensional projection of a three-dimensional micrograph, yet these projections can bias results. For example, a maximum projection, where any pixel is the maximum value of the entire z column of pixels, is 322 increasingly biased with poorer signal-to-noise ratios. Alternatively, analysis on a summed projection, where any 323 324 pixel is the summed value of all underlying z planes, may be particularly sensitive to hair cell orientation or morphology. The HCAT segmentation pipeline is a tool for the automated quantification of hair cell fluorescence 325 326 from three-dimensional imaging data, which can be carried out on full or partial cochlear samples. Furthermore, 327 in combination with the automated calculation of each hair cell's predicted best frequency in full cochleae. determined to be in good agreement with manual estimates of hair cell best frequency using the widely accepted 328 329 EPL method, this enables the creation of hair cell fluorescence cochleograms.

330

344

354

While our models were trained solely using the Myo-VIIa and phalloidin labels, the model can perform on 331 specimens labeled with other markers, provided they contain both cytosolic hair cell, and stereocilia, labels, An 332 example of cytosolic hair cell labels might include (i) hair-cell-specific expression of a genetically encoded 333 fluorescence marker, such as Atoh1-eGFP; (ii) cochlear samples collected from animals following AAV injection 334 335 resulting in strong expression of fluorescence markers in hair cells; or (iii) cochlear samples treated with FM1-43 styryl dye known to selectively permeate into hair cells with functional mechanotransduction channel when 336 337 briefly applied to live cochlea. An example of stereocilia-specific label is an immunolabeling against a highly enriched stereociliary protein espin, widely used to label hair cell stereocilia in adult cochlear preparations. 338 339

While our segmentation accuracy offers superior performance compared to other methods, this increase in segmentation accuracy comes at a cost of cell detection accuracy, as seen by the variability in the cochleogram seen in **Figure 5F**, likely due to the limitation in deep learning architecture complexity for volumetric analysis. Thus, in the HCAT software we offer two distinct pipelines for cell detection and segmentation analysis.

While our study shows that these algorithms perform well on data collected in-house, their generalizability across 345 datasets collected by other research groups is pending further validation. Deep learning models tend to 346 generalize poorly; when evaluated on community-supplied data, HCAT was no exception. While there were 347 several instances of highly accurate performance on community-provided datasets, the deep learning models 348 performed sub-optimally in cases where hair cells on community-supplied datasets were visually different than 349 350 those in our training data. To our surprise, while the hair cell detection pipeline was trained exclusively with 351 confocal microscopy data, it performed well on a set of community-provided images of hair cells collected using 352 widefield fluorescence. All community-provided datasets that resulted in poor HCAT performance will be used to retrain the model with a prior permission from the owner(s) of the data set. 353

Overall, we are planning to expand the training data employing a wider variety of hair cells, such as those from 355 adult mice and examples of cochlear coils with hair cells following a treatment with ototoxic drugs. Furthermore, 356 we will endeavor to continually update and maintain the software as well as periodically update the machine 357 358 learning model as the state of 3D cell segmentation and cell detection advances. Further development of this 359 software in the future will provide support for automatic segmentation of adult cochlear tissue, often imaged in pieces rather than as a contiguous piece of tissue, due to dissection limitations. Additionally, we will continue to 360 re-train the algorithm as more training data become available, improving its performance over time to a wider 361 362 array of tissue preparations and ages.

363

To our knowledge, this is the first whole cochlear analysis pipeline capable of accurately and quickly detecting or segmenting hair cells. Offering state of the art performance, this hair cell analysis toolbox (HCAT) enables expedited cochlear image data analysis while maintaining high accuracy. This accurate and unsupervised data analysis approach will both facilitate ease of research and improve experimental rigor.

368

369 Materials and Methods

370

Some of the procedures described below, including the staining authors used to prepare cochlear samples labeled with four different fluorescent markers as shown in **Figure 1** are not required for the use of the presented hair cell analysis toolbox. They do, however provide an example of a dataset for which the single-cell fluorescence intensity quantification feature following the application of the hair cell segmentation pipeline can be utilized.

376

377 Sample preparation, confocal microscopy. Postnatal day (P) 0 C57Bl6 mice of either sex were cryoanesthetized 378 and injected with AAVs encoding eGFP through the round window membrane of the cochlea as described previously. The animals were then returned to the dam for recovery. Organs of Corti were dissected at P5 in 379 Leibovitz's L-15 culture medium (21083-027, Thermo Fisher Scientific) and fixed in 4% formaldehyde for 1 hour. 380 The samples were permeabilized with 0.2% Triton-X for 30 minutes and blocked with 10% goat serum in calcium-381 382 free HBSS for two hours. To visualize the hair cells, samples were labeled with an anti-Myosin VIIA antibody (#25-6790 Proteus Biosciences, 1:400) and goat anti-rabbit CF568 (Biotium). Additionally, samples were labeled 383 384 with Phalloidin to visualize actin filaments (Biotium CF640R Phalloidin) and with DAPI to visualize cell nuclei (Molecular Probes DAPI, #D1306). Samples were then mounted on slides using ProLong® Gold Antifade 385 Mounting kit (P36931, Thermo Fisher Scientific,) and imaged with a Leica SP8 confocal microscope (Leica 386 Microsystems) using a 63x/1.3 NA objective. Confocal Z-stacks of 512x512 pixel images with an effective pixel 387 388 size of 288 nm were collected using the tiling functionality of the Leica LASX acquisition software. All experiments were carried out in compliance with ethical regulations and approved by the Animal Care Committees of 389 Massachusetts Eye and Ear. 390

391

392 <u>Computational Environment</u>: All scripts were run on a custom-built analysis computer running Ubuntu 20.04.1 393 LTS, an open-source Linux distribution from Canonical based on Debian. The workstation was equipped with an 394 AMD Ryzen 7 3700X 8-Core processor with 64 GB of RAM. Deep learning computations were accelerated by a 395 Nvidia 2080Ti graphics card with 11GB of DDR6 video memory. Many scripts were custom written in python 3.8 396 using open source scientific computation libraries including numpy⁴⁷, matplotlib⁴⁸, scikit-learn⁴⁹. All deep learning 397 architectures, training logic, and much of the data transformation pipeline was written in pytorch⁵⁰. All code has 398 been hosted on github and is available at: <u>https://github.com/buswinka/hcat</u>.

- 400 <u>*Training Data Annotation*</u>: Two deep learning models were trained on distinct datasets.
- 401

399

For the *hair cell segmentation task*, 17 training and 2 validation Z-stacks of cochlear hair cells were selected from different datasets and ensuring equal representation along the entire length of the cochlea, each containing ~50-100 hair cells. Each image was manually segmented using the Amira Software package (Thermo Fisher Scientific) running on a Lenovo ThinkPad X1 yoga touchscreen laptop. Cell outlines were delineated by hand using a stylus and saved with a unique cell ID to create cell masks.

407

For the *hair cell detection task* training data, bounding boxes for hair cells seen in maximum projected z-stacks were manually annotated in the labelImg software⁵¹ and saved as an xml file. For whole cochlear cell annotation, a "human in the loop" approach was taken, first evaluating the deep learning model on the entire cochlea, then manually correcting errors.

Training: Each deep learning model was trained to accurately perform its designed task. The procedure to train 413 the Faster R-CNN detection algorithm has been standardized and extensively documented in previous reports⁶. 414 while the training procedure for instance segmentation is more involved. A spatial embedding approach has 415 been shown to be a versatile method to perform biomedical instance segmentation^{16,17,21,52}, however the optimal 416 way to train a network for this task is less well defined. We optimize spatial embeddings based on the distance 417 of each pixel (i) in embedding space from the known centroid (C) of the underlying object. From this distance, 418 each pixel (i) can be assigned a probability (ϕ) of belonging to an object (k) by equation (1), regularized by the 419 parameter sigma (σ). Due to the anisotropy of the dataset, we elect to choose a proportionally smaller sigma 420 421 when computing the loss in Z compared to X and Y, thereby improving embedding accuracy along that dimension. For the spatial embedding instance segmentation, the architecture outputs spatial embedding vectors 422 in addition to a probability map. In the post processing of the spatial embedding, object detection heavily relies 423 on tight clusters of pixels, therefore it is advantageous to penalize false negative results more heavily than false 424 positive ones. This contrasts with the probability map, where false negative results should be more heavily 425 penalized to ensure the detection of cell borders. To rectify these opposing needs, we chose to optimize the 426 Tversky loss⁵³ of the spatial embedding and semantic probability map with differing penalties for each. 427

428

When using a probability map to aid in watershed segmentation, it is critical the borders of cells are well defined. To this end, we optimize binary cross entropy loss with a pixel-by-pixel penalty based on the distance from the border of an object, with pixels closer to the object penalized more heavily. Due to the class imbalance between background and foreground, optimizing over binary cross entropy alone leads the model to a local minimum where the entire output is predicted to be background. Therefore, we additionally include the dice coefficient in loss calculation.

435

For both tasks, the deep learning architectures were trained with the Adam optimizer with a learning rate starting 436 at 1e-4 and decaying based on cosine annealing with warm restarts with a period of 100 epochs. When training 437 the model for the spatial embedding task, we initialize sigma to be large values, and progressively decay by half 438 at set epochs. Due to the labor-intensive nature of the manual segmentation workflow, we are limited in the total 439 number of training volumes we can generate. In cases with a small number of training images, deep learning 440 models tend to fail to generalize and instead "memorize" the training data. To avoid this, we make heavy use of 441 image transformations which randomly add variability to our images and synthetically increase the variety of our 442 training data sets⁵⁴ (Supplementary Figure S1). 443

444

Acknowledgements. We would like to thank Dr. Marcelo Cicconet (Image and Data Analysis Core at Harvard Medical School) and Haobing Wang, MS (Mass Eye and Ear Light Microscopy Imaging Core Facility) for their assistance in this project. We thank Dr. Guy Richardson, Dr. Corne Kros, Dr. Bradley Walters, Dr. Albert Edge, Dr. Yushi Hayashi, Dr. Lisa Cunningham, Dr. Mark Rutherford, Dr. Tejbeer Kaur, Dr. Vijayprakash Manickam, Dr. Ksenia Gnedeva, the members of their teams and all other research groups, for providing their datasets to evaluate the HCAT. We also thank Dr. Richard Osgood and Evan Hale, MS for critical reading of the manuscript.

452 **Code availability.**

453 All code has been hosted on github and is available at: <u>https://github.com/buswinka/hcat</u>.

- 454
- 455

456 457	Refere	ences
458 459	1.	Ashmore J. Tonotopy of cochlear hair cell biophysics (excl. mechanotransduction). <i>Current opinion in physiology</i> , 2020:18:1-6.
460 461	2.	Potter PK, Bowl MR, Jeyarajan P, et al. Novel gene function revealed by mouse mutagenesis screens for models of age-related disease. <i>Nature communications</i> . 2016;7(1):1-13.
462 463 464	3.	Gray SJ, Foti SB, Schwartz JW, et al. Optimizing promoters for recombinant adeno-associated virus- mediated gene expression in the peripheral and central nervous system using self-complementary vectors. <i>Human gene therapy</i> 2011:22(9):1143-1153
465 466	4.	Kalluri R, Shera CA. Distortion-product source unmixing: A test of the two-mechanism model for DPOAE generation. <i>The Journal of the Acoustical Society of America</i> , 2001:109(2):622-637.
467 468	5.	Glasscock ME. <i>The ABR handbook : auditory brainstem response</i> . 2nd ed. ed. New York: Thieme Medical Publishers; 1987.
469 470	6.	Ren S, He K, Girshick R, Sun J. Faster R-CNN: towards real-time object detection with region proposal networks. <i>IEEE transactions on pattern analysis and machine intelligence</i> . 2016;39(6):1137-1149.
471 472 473	7.	Ezhilarasi R, Varalakshmi P. Tumor detection in the brain using faster R-CNN. Paper presented at: 2018 2nd International Conference on I-SMAC (IoT in Social, Mobile, Analytics and Cloud)(I-SMAC) I-SMAC (IoT in Social, Mobile, Analytics and Cloud)(I-SMAC), 2018 2nd International Conference on 2018.
474 475	8.	Yang S, Fang B, Tang W, Wu X, Qian J, Yang W. Faster R-CNN based microscopic cell detection. Paper presented at: 2017 International Conference on Security, Pattern Analysis, and Cybernetics (SPAC)2017.
476 477 478	9.	Zhang J, Hu H, Chen S, Huang Y, Guan Q. Cancer cells detection in phase-contrast microscopy images based on Faster R-CNN. Paper presented at: 2016 9th international symposium on computational intelligence and design (ISCID) 2016
479 480	10.	Strahler AN. Quantitative analysis of watershed geomorphology. <i>Eos, Transactions American Geophysical Union</i> , 1957:38(6):913-920.
481 482	11.	Atta-Fosu T, Guo W, Jeter D, Mizutani CM, Stopczynski N, Sousa-Neves R. 3D Clumped Cell Segmentation Using Curvature Based Seeded Watershed. <i>Journal of imaging</i> . 2016;2(4):31.
483 484	12.	Jones TR, Kang IH, Wheeler DB, et al. CellProfiler Analyst: data exploration and analysis software for complex image-based screens. <i>BMC bioinformatics</i> . 2008;9(1):482-482.
485 486 487	13.	Beliën JAM, van Ginkel HAHM, Tekola P, et al. Confocal DNA cytometry: a contour-based segmentation algorithm for automated three-dimensional image segmentation. <i>Cytometry (New York, NY)</i> . 2002;49(1):12-21.
488 489	14.	Xiaowei C, Xiaobo Z, Wong STC. Automated segmentation, classification, and tracking of cancer cell nuclei in time-lapse microscopy. <i>IEEE transactions on biomedical engineering</i> . 2006;53(4):762-766.
490 491	15.	He K, Gkioxari G, Dollár P, Girshick R. Mask r-cnn. Paper presented at: Proceedings of the IEEE international conference on computer vision2017.
492 493 494	16.	Neven D, Brabandere BD, Proesmans M, Gool LV. Instance segmentation by jointly optimizing spatial embeddings and clustering bandwidth. Paper presented at: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition2019.
495 496	17.	Stringer C, Wang T, Michaelos M, Pachitariu M. Cellpose: a generalist algorithm for cellular segmentation. <i>Nature methods.</i> 2021;18(1):100-106.
497 498	18.	Rasmussen CE. Gaussian processes in machine learning. Paper presented at: Summer school on machine learning2003.
499 500	19.	Greenwood DD. A cochlear frequency-position function for several species—29 years later. <i>The Journal of the Acoustical Society of America</i> . 1990;87(6):2592-2605.
501 502	20.	He K, Zhang X, Ren S, Sun J. Deep residual learning for image recognition. Paper presented at: Proceedings of the IEEE conference on computer vision and pattern recognition2016.
503 504	21.	Zhang D, Chun J, Cha SK, Kim YM. Spatial semantic embedding network: fast 3D instance segmentation with deep metric learning. <i>arXiv preprint arXiv:200703169.</i> 2020.

- 505 22. Yao C, Tang J, Hu M, et al. Claw U-Net: A Unet-based Network with Deep Feature Concatenation for 506 Scleral Blood Vessel Segmentation. *arXiv preprint arXiv:201010163.* 2020.
- 50723.Qamar S, Jin H, Zheng R, Ahmad P, Usama M. A variant form of 3D-UNet for infant brain segmentation.508Future Generation Computer Systems. 2020;108:613-623.
- 50924.Wang C, MacGillivray T, Macnaught G, Yang G, Newby D. A two-stage 3D Unet framework for multi-class510segmentation on full resolution image. *arXiv preprint arXiv:180404341.* 2018.
- Weigert M, Schmidt U, Haase R, Sugawara K, Myers G. Star-convex polyhedra for 3d object detection
 and segmentation in microscopy. Paper presented at: Proceedings of the IEEE/CVF Winter Conference
 on Applications of Computer Vision2020.
- 51426.Schmidt U, Weigert M, Broaddus C, Myers G. Cell detection with star-convex polygons. Paper presented515at: International Conference on Medical Image Computing and Computer-Assisted Intervention2018.
- 516 27. Sheridan A, Nguyen T, Deb D, et al. Local shape descriptors for neuron segmentation. *bioRxiv.* 2021.
- Ronneberger O, Fischer P, Brox T. U-net: Convolutional networks for biomedical image segmentation.
 Paper presented at: International Conference on Medical image computing and computer-assisted intervention2015.
- 52029.Lucchi A, Smith K, Achanta R, Knott G, Fua P. Supervoxel-based segmentation of mitochondria in em521image stacks with learned shape features. *IEEE transactions on medical imaging.* 2011;31(2):474-486.
- Pohlen T, Hermans A, Mathias M, Leibe B. Full-resolution residual networks for semantic segmentation
 in street scenes. Paper presented at: Proceedings of the IEEE Conference on Computer Vision and Pattern
 Recognition2017.
- 52531.Dumoulin V, Visin F. A guide to convolution arithmetic for deep learning. arXiv preprint arXiv:160307285.5262016.
- 52732.Xu B, Wang N, Chen T, Li M. Empirical evaluation of rectified activations in convolutional network. *arXiv*528*preprint arXiv:150500853.* 2015.
- 52933.Sun T, Tsang WM, Park WT, Cheng KJ, Merugu S. Modeling in vitro neural electrode interface in neural530cell culture medium. *Microsyst Technol.* 2015;21(8):1739-1747.
- 53134.He K, Zhang X, Ren S, Sun J. Identity mappings in deep residual networks. Paper presented at: European532conference on computer vision2016.
- 533 35. Stalling D, Westerhoff M, Hege H-C. Amira: A highly interactive system for visual data analysis. *The visualization handbook.* 2005;38:749-767.
- 53536.Neubeck A, Van Gool L. Efficient non-maximum suppression. Paper presented at: 18th International536Conference on Pattern Recognition (ICPR'06)2006.
- 537 37. Viberg A, Canlon B. The guide to plotting a cochleogram. *Hearing Res.* 2004;197(1-2):1-10.
- 53838.Urata S, Iida T, Yamamoto M, et al. Cellular cartography of the organ of Corti based on optical tissue539clearing and machine learning. *eLife*. 2019;8.
- 54039.Cortada M, Sauteur L, Lanz M, Levano S, Bodmer D. A deep learning approach to quantify auditory hair541cells. *Hearing Res.* 2021;409:108317.
- 54240.Li H, Liu H, Corrales CE, et al. Differentiation of neurons from neural precursors generated in floating543spheres from embryonic stem cells. BMC neuroscience. 2009;10(1):122-122.
- 41. Garza LA, Yang C-C, Zhao T, et al. Bald scalp in men with androgenetic alopecia retains hair follicle stem cells but lacks CD200-rich and CD34-positive hair follicle progenitor cells. *The Journal of clinical investigation.* 2011;121(2):613-622.
- 54742.György B, Sage C, Indzhykulian AA, et al. Rescue of Hearing by Gene Delivery to Inner-Ear Hair Cells Using548Exosome-Associated AAV. *Mol Ther.* 2017;25(2):379-391.
- 54943.Rhee C-K, He P, Jung JY, et al. Effect of low-level laser treatment on cochlea hair-cell recovery after550ototoxic hearing loss. Journal of biomedical optics. 2013;18(12):128003-128003.
- 551 44. Gersten BK, Fitzgerald TS, Fernandez KA, Cunningham LL. Ototoxicity and Platinum Uptake Following 552 Cyclic Administration of Platinum-Based Chemotherapeutic Agents. *Journal of the Association for* 553 *Research in Otolaryngology*. 2020;21(4):303-321.

- Liu Z, Owen T, Fang J, Srinivasan RS, Zuo J. In vivo notch reactivation in differentiating cochlear hair cells
 induces sox2 and prox1 expression but does not disrupt hair cell maturation. *Developmental Dynamics*.
 2012;192(4):684-696.
- 557 46. Rasband WS. ImageJ. Bethesda, MD; 1997.
- 558 47. Van Der Walt S, Colbert SC, Varoquaux G. The NumPy array: a structure for efficient numerical computation. *Computing in science & engineering.* 2011;13(2):22-30.
- 560 48. Hunter JD. Matplotlib: A 2D graphics environment. *IEEE Annals of the History of Computing.* 561 2007;9(03):90-95.
- 562 49. Pedregosa F, Varoquaux G, Gramfort A, et al. Scikit-learn: Machine learning in Python. *the Journal of machine Learning research*. 2011;12:2825-2830.
- 564 50. Paszke A, Gross S, Massa F, et al. Pytorch: An imperative style, high-performance deep learning library. 565 *arXiv preprint arXiv:191201703.* 2019.
- 566 51. Lin T. LabelImg. Github; 2015.
- 567 52. Sheridan A, Nguyen T, Deb D, et al. Local Shape Descriptors for Neuron Segmentation. *bioRxiv.* 2021:2021.2001.2018.427039.
- 53. Salehi SSM, Erdogmus D, Gholipour A. Tversky loss function for image segmentation using 3D fully convolutional deep networks. Paper presented at: International workshop on machine learning in medical imaging2017.
- 572 54. Shorten C, Khoshgoftaar TM. A survey on image data augmentation for deep learning. *Journal of Big* 573 *Data*. 2019;6(1):1-48.
- 574 575



Supplementary Figure S1: Training data augmentation pipeline. Training images for each deep learning 579 approach underwent identical data augmentation steps, increasing the variability of our dataset and improving 580 performance. Each of these augmentation steps were probabilistically applied sequentially (left to right).



Supplementary Figure S2: Validation of hair cell detection analysis and location estimation pipeline. Whole cochlear turns (*A*) were manually annotated and evaluated with the HCAT detection analysis pipeline. Each analysis generated highly accurate cochleograms (*B*), reporting the 'ground truth' result obtained from manual segmentation (*dark lines*) superimposed onto the cochleogram generated from hair cells detected by the HCAT detection analysis (*light lines*), reporting highly accurate results. The frequency estimation error was then calculates as an octave difference to quantify the accuracy of predicted frequency estimation for every hair cell vs their ground truth frequency (*C*). Optimal cell detection and non-maximum suppression thresholds were discerned via a grid search by maximizing the true positive rate penalized by the false positive and false negative rates (*D*). Black lines on the ROC curves (*E*) denote the optimal hyperparameter value.