# SUPPLEMENTARY INFORMATION

## Supplementary Note 1: Image Data Overview

A typical CODEX dataset is composed of images broken down by the following dimensions (**Supplementary Fig. 1a**):

- **# of regions:** e.g., four different regions of tissue or four different samples
- Region size: e.g., 7x9 tiles (with 30% overlap)
- **# of Z stacks:** e.g., 15 (at 1.5 µm each), which ensures in-focus image capture across a large tissue region
- **# of channels**: e.g., CH1: Hoechst, CH2: Alexa488, CH3: ATTO550, CH4: Alexa647
- **# of cycles:** e.g., Cycle 3: Hoechst, CD3, CD4, CD8; Cycle 4: Hoechst, CD31, CD19, podoplanin.



**Supplementary Figure 1. CODEX data analysis pipeline. a**, A visual representation of the data structure for CODEX imaging files. **b**, The user interfaces for the CODEX Uploader and Segmentation software. **c**, A visual representation of the single-cell segmentation masks for a CODEX image. **d**, Key steps for manual cleaning the segmented data using a flow cytometry software platform.

Data will be generated from the scope as individual images across all dimensions. These images need to be concatenated to form a 7x9 imaging region, contain all z stacks, and all channels for each cycle. There is minor drift of the microscope in between cycles, so drift compensation must then be done to align markers across all channels. This is done with the Hoechst or DRAQ5 (nuclear staining) used for all cycles as a reference.

Once all dimensions have been merged, it is necessary to implement cell segmentation to quantitatively describe the image in single-cell resolution. This results in a file with each row representing a segmented cell and each column indicating a feature of that cell (e.g., X/Y positions, Channel intensities, original tile number).

**Supplementary Note 2: Image Processing**

To automate and standardize these image processing steps, our laboratory has published an open source software, CODEXSetup (available at https://github.com/nolanlab/CODEX). There is a single software interface for combining image dimensions and drift compensation (CODEXUploader.exe) and another for segmentation (CODEXSegm.exe). Within CODEXUploader the necessary parameters that must be set are input and output directories and X and Y region size (**Supplementary Fig. 1b**). Other parameters should be auto-populated based on the Experiment.json, channelnames.txt, and expsouretimes.txt files found in the directory with your images. There are additional options for background subtraction (blank cycle), including processing of H&E staining (if completed), processing for multipoint imaging (especially useful for tissue microarrays), deconvolution (if there is a microevolution license), and exporting as MAV sequence (downstream analysis with Akoya Biosciences's ImageJ MAV plugin). Finally, there is an option for previewing the processing to make sure that parameters used are appropriate; this is important as the image processing can be time consuming.

Once run, the program will produce the following outputs in the specified output directory:

- Concatenated individual tiles across all dimensions
- bestFocus folder with:
  - Concatenated individual tiles with best focal plane selected
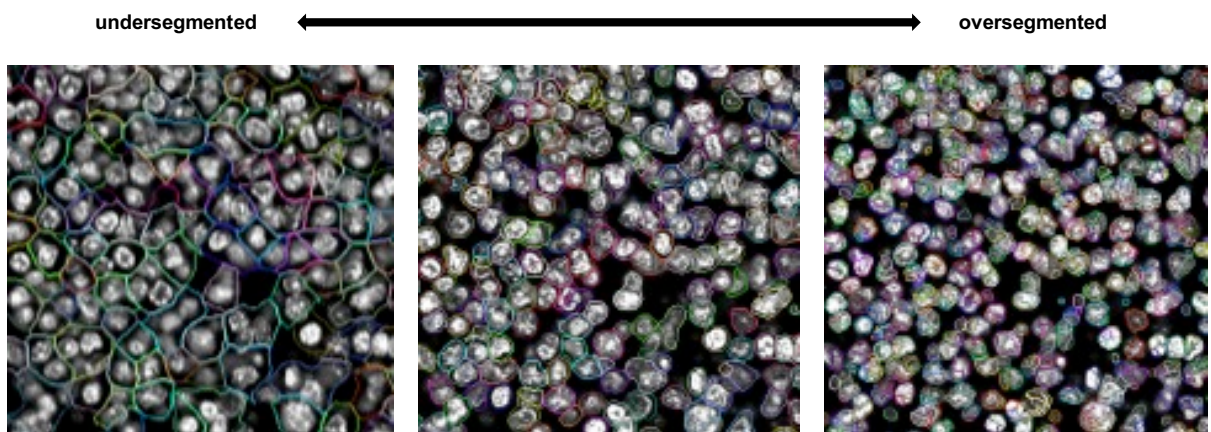  - Montages of stitched concatenated individual tiles with best focal planes selected (**Supplementary Fig. 1a**)

**Supplementary Note 3: Single-cell Segmentation**

After initial data processing, segmentation can be run on the output folder generated. Within CODEXSegm the following parameters must be set based on the dataset for volumetric watershed segmentation (good starting point values are given in parentheses) (**Supplementary Fig. 1c**):

- Radius: The estimated radius of nuclei size (6)
- maxCutoff: The highest signal to be considered a nuclei to filter potential debris (0.99)
- minCutoff: The lowest signal to be considered a part of a nuclei (0.07)
- relativeCutoff – Used for region growth relative to peak intensity of each nuclei; usually not changed (0.2)
- sizeCutOffFactor – Filters cells smaller than expected cell size based on radius (0.5)
- Nuclear stain channel: The channel number of nuclear stain

● Nuclear stain cycle: The number of the nuclear stain cycle
● Membrane stain channel: The channel number of membrane stain
●Membrane stain cycle: The number of the membrane stain cycle; if there is no membrane stain, then segmentation can be done with nuclear channel only and the number is -1
●Anisotropic region growth: If checked, this restricts pos-neg separation between crowded cell types (typically not checked)
●Single Plane Segmentation: If checked, this integrates the signal only based on one plane instead of volumetrically (typically checked)

CODEXSegm also provides an option to preview the segmentation parameters on one tile. It is usually an iterative process to optimize the segmentation mask for representative regions of the tissue. It is generally better to oversegment than to undersegment because single cell data will get extracted from within the mask. Undersegmentation will include multiple cells in the same mask and thus cells will appear "positive" for both cell-type markers confounding the cell-type identification results (**Supplementary** Fig. 2). Oversegmentation will occur when multiple masks are drawn through one cell, while not as problematic for cell-type identification, this will pose problems for cell frequency assessments. Optimal settings would occur where each single cell is individually segmented with its own mask.



**Supplementary Figure 2. Demonstration of poor and good segmentation results.** Undersegmented, optimally segmented, and oversegmented images are represented with masks overlaid on one tile from (Hoechst nuclear staining) from the multicycle.
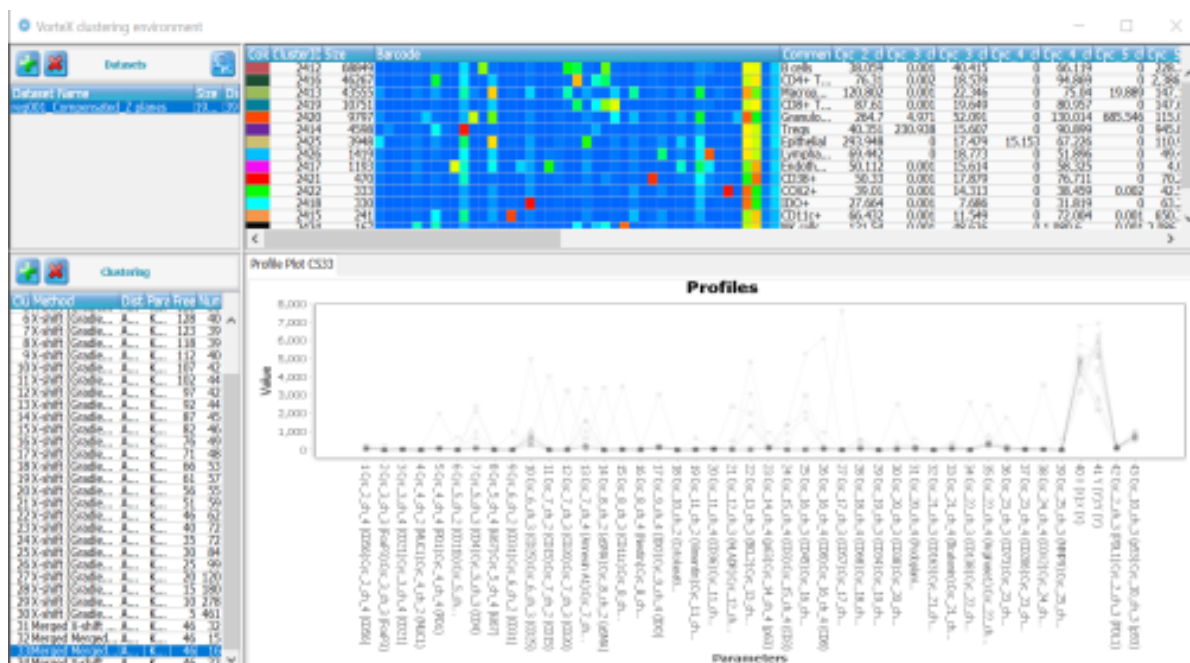
Once run, the program will produce the following outputs in the directory where segmentation was performed:

●Dataframe for each region using either compensated or uncompensated settings that include channel intensities and location for each segmented cell
■FCS and CSV formats are generated
●Individual tile segmentation masks
■TXT and PNG formats are generated

**Supplementary Note 4: Cell type Identification**

The individual data frames can be processed and analyzed with methods used for other single-cell techniques like flow cytometry and mass cytometry. Cell-type identification can be performed by manual gating or unsupervised cell type clustering.

First, it is necessary to import compensated files into a manual gating program like CellEngine or FlowJo. This allows selection of nucleated cells by gating for dual-positive nucleated cells (**Supplementary Fig. 3**). Selection of cells from a best focal plane not on the extremes of the Z stacks is necessary to eliminate noise in the data prior to unsupervised clustering.



**Supplementary Figure 3. Unsupervised single-cell clustering software.** The user interface for VorteX clustering.

The cleaned datasets can be imported into Vortex, which is another unsupervised clustering (available at https://github.com/nolanlab/vortex/wiki/Getting-Started) (28). The data are imported into the clustering database. Markers that should be included for clustering (blue), those not used for clustering but kept for analysis (yellow), and markers to be removed (gray) are selected by clicking and changing the color. Parameters are the following (good starting points for CODEX data are given):

- Numerical transformation: none
- Noise threshold: uncheck
- Feature rescaling: none

- Normalization: none
- Minimal Euclidean length: uncheck
- Distance Measure: Angular Distance
- Clustering Algorithm: X-shift
- Density Estimate: N nearest neighbors (fast)
- Num. Neighbors for mode finding N: Determine Automatically

The program calculates clusters for each *K* between the default parameters. This can be used to find the elbow point for optimal clustering results. Alternatively, one can select to over-cluster by choosing a low *K* with a high corresponding cluster number and then manually merge clusters. Using the cluster investigation interface with expression profiles, minimal spanning trees, and divisive marker trees, cell type annotation can be done for each cluster based on marker expression. Finally, spatial locations of these clusters should be verified to ensure accuracy of the clustering.