

1 **Raman2RNA: Live-cell label-free prediction of single-cell RNA** 2 **expression profiles by Raman microscopy**

3 Koseki J. Kobayashi-Kirschvink^{1,2}, Shreya Gaddam^{1,9}, Taylor James-Sorenson¹, Emanuelle
4 Grody^{1,3}, Johain R. Ounadjela^{1,3}, Baoliang Ge⁴, Ke Zhang⁵, Jeon Woong Kang², Ramnik Xavier^{1,6},
5 Peter T. C. So^{2,4}, Tommaso Biancalani^{1,9,†,‡}, Jian Shu^{1,3,5,†,‡}, Aviv Regev^{1,7,8,9,†,‡}

6 ¹Broad Institute of MIT and Harvard, Cambridge, MA 02142, USA

7 ²Laser Biomedical Research Center, G. R. Harrison Spectroscopy Laboratory, Massachusetts Institute of
8 Technology, Cambridge, MA 02139, USA

9 ³Whitehead Institute for Biomedical Research, Cambridge, MA 02142, USA

10 ⁴Department of Mechanical and Biological Engineering, Massachusetts Institute of Technology, Cambridge, MA
11 02139, USA

12 ⁵Cutaneous Biology Research Center, Massachusetts General Hospital, Harvard Medical School, Boston, MA
13 02129, USA

14 ⁶Center for Computational and Integrative Biology and Department of Molecular Biology, Massachusetts General
15 Hospital, Boston, Massachusetts 02114, USA

16 ⁷Department of Biology, Massachusetts Institute of Technology, Cambridge, MA 02139, USA

17 ⁸Howard Hughes Medical Institute, Cambridge, MA 02142, USA

18 ⁹Present address: Genentech, 1 DNA Way, South San Francisco, CA 94080, USA

19 [†]These authors contributed equally

20 [‡]Correspondence: aviv.regev.sc@gmail.com; jian.shu@mgh.harvard.edu; tbiancal@broadinstitute.org

21

22

23 **Single cell RNA-Seq (scRNA-seq) and other profiling assays have opened new windows into**
24 **understanding the properties, regulation, dynamics, and function of cells at unprecedented**
25 **resolution and scale. However, these assays are inherently destructive, precluding us from**
26 **tracking the temporal dynamics of live cells, in cell culture or whole organisms. Raman**
27 **microscopy offers a unique opportunity to comprehensively report on the vibrational energy**
28 **levels of molecules in a label-free and non-destructive manner at a subcellular spatial**
29 **resolution, but it lacks in genetic and molecular interpretability. Here, we developed**
30 **Raman2RNA (R2R), an experimental and computational framework to infer single-cell**
31 **expression profiles in live cells through label-free hyperspectral Raman microscopy images**
32 **and multi-modal data integration and domain translation. We used spatially resolved single-**
33 **molecule RNA-FISH (smFISH) data as anchors to link scRNA-seq profiles to the paired**
34 **spatial hyperspectral Raman images, and trained machine learning models to infer**
35 **expression profiles from Raman spectra at the single-cell level. In reprogramming of mouse**
36 **fibroblasts into induced pluripotent stem cells (iPSCs), R2R accurately ($r>0.96$) inferred**
37 **from Raman images the expression profiles of various cell states and fates, including iPSCs,**
38 **mesenchymal-epithelial transition (MET) cells, stromal cells, epithelial cells, and fibroblasts.**
39 **R2R outperformed inference from brightfield images, showing the importance of**
40 **spectroscopic content afforded by Raman microscopy. Raman2RNA lays a foundation for**
41 **future investigations into exploring single-cell genome-wide molecular dynamics through**
42 **imaging data, *in vitro* and *in vivo*.**

43 **Keywords: Raman microscopy, single-cell transcriptomics, multi-domain translation**

44

45 **Main**

46 Cellular states and functions are determined by a dynamic balance between intrinsic and extrinsic
47 programs. Dynamic processes such as cell growth, stress responses, differentiation, and
48 reprogramming are not determined by a single gene, but by the orchestrated temporal expression
49 and function of multiple genes organized in programs and their interactions with other cells and
50 the surrounding environment¹. To understand how cells change their states in physiological and
51 pathological conditions it is essential to decipher the dynamics of the underlying gene programs.

52 Despite major advances in single cell genomics and microscopy, we still cannot track live cells
53 and tissues at the genomic level. On the one hand, single cell and spatial genomics have provided
54 a view of gene programs and cell states at unprecedented scale and resolution¹, but these
55 measurement methods are destructive, and involve tissue fixation and freezing and/or cell lysis,
56 precluding us from directly tracking the dynamics of full molecular profiles in live cells or
57 organisms. While advanced computational methods, such as pseudo-time algorithms (*e.g.*,
58 Monocle², Waddington-OT³) and velocity-based methods (*e.g.*, velocity⁴, scVelo⁵), can infer
59 dynamics from snapshots of molecular profiles, they rely on assumptions that remain challenging
60 to verify experimentally⁶. On the other hand, fluorescent reporters can be used to monitor the
61 dynamics of individual genes and programs within live cells, but are limited in the number of
62 targets they can report⁷, must be chosen ahead of the experiment and often involve genetically
63 engineered cells. Moreover, the vast majority of dyes and reporters require fixation or can interfere
64 with nascent biochemical processes and alter the natural state of the gene of interest⁷. Therefore,
65 it remains technically challenging to dynamically monitor the activity of a large number of genes
66 simultaneously.

67 Raman microscopy opens a unique opportunity for monitoring live cells and tissues, as it
68 collectively reports on the vibrational energy levels of molecules in a label-free and non-
69 destructive manner at a subcellular spatial resolution, thus providing molecular fingerprints of
70 cells⁸. Pioneering research has demonstrated that Raman microscopy can be used for
71 characterizing cell types and cell states⁸, non-destructively diagnosing pathological specimens
72 such as tumors⁹, characterizing the developmental states of embryos¹⁰, and identifying bacteria
73 with antibiotic resistance¹¹. However, the complex and high-dimensional nature of the spectra, the
74 spectral overlaps of biomolecules such as proteins and nucleic acids, and the lack of unified
75 computational frameworks have hindered the decomposition of the underlying molecular
76 profiles^{7,8}.

77 To address this challenge and leverage the complementary strengths of Raman microscopy and
78 scRNA-Seq, we developed Raman2RNA (R2R), an experimental and computational framework
79 for inferring single-cell RNA expression profiles from label-free non-destructive Raman
80 hyperspectral images (**Fig. 1**). R2R takes as input spatially resolved hyperspectral Raman images
81 from live cells, smFISH data of selected markers from the same cells, and scRNA-seq from the
82 same biological system. R2R then uses the smFISH data as an anchor to learn a model that links
83 spatially resolved hyperspectral Raman images to scRNA-seq. Finally, from this model, R2R then
84 computationally infers the anchor smFISH measurements from hyperspectral Raman images and
85 then the single-cell expression profiles. The result is a label-free live-cell inference of single-cell
86 expression profiles.

87 To facilitate data acquisition, we developed a high-throughput multi-modal spontaneous Raman
88 microscope that enables automated acquisition of Raman spectra, brightfield, and fluorescent
89 images. In particular, we integrated Raman microscopy optics to a fluorescence microscope, where

90 high-speed galvo mirrors and motorized stages were combined to achieve a large field of view
91 (FOV) scanning, and where dedicated electronics automate measurements across multiple
92 modalities (**Extended Data Fig. 1-2, Methods**).

93 We first demonstrated that R2R can infer profiles of two distinct cell types: mouse induced
94 pluripotent stem cells (iPSCs) expressing an endogenous *Oct4*-GFP reporter and mouse
95 fibroblasts¹². To this end, we mixed the cells in equal proportions, plated them in a gelatin-coated
96 quartz glass-bottom Petri dish, and performed live-cell Raman imaging, along with fluorescent
97 imaging of live-cell nucleus staining dye (Hoechst 33342) for cell segmentation and image
98 registration, and an iPSC marker gene, *Oct4*-GFP (**Fig. 2a**). The excitation wavelength for our
99 Raman microscope (785 nm) was distant enough from the GFP Stokes shift emission, such that
100 there was no interference with the cellular Raman spectra (**Extended Data Fig. 3**). Furthermore,
101 there was no notable photo-toxicity induced in the cells. After Raman and fluorescence imaging,
102 we fixed and permeabilized the cells and performed smFISH (with hybridization chain reaction
103 (HCR¹³), **Methods**) of marker genes for mouse iPSCs (*Nanog*) and fibroblasts (*Colla1*). We
104 registered the nuclei stains, GFP images, HCR images, and Raman images through either
105 polystyrene control bead images or reference points marked under the glass bottom dishes
106 (**Extended Data Fig. 4, Methods**).

107 The Raman spectra distinguished the two cell populations in a manner congruent with the
108 expression of their respective reporter (measured live or by smFISH in the same cells), as reflected
109 by a low-dimensional embedding of hyperspectral Raman data (**Fig. 2b**). Specifically, we focused
110 on the fingerprint region of Raman spectra (600-1800 cm^{-1} , 930 of the 1,340 features in a Raman
111 spectrum), where most of the signatures from various key biomolecules, such as proteins, nucleic
112 acids, and metabolites, lie⁸. After basic preprocessing, including cosmic-ray and background

113 removal and normalization, we aggregated Raman spectra that are confined to the nuclei, obtaining
114 a 930-dimensional Raman spectroscopic representation for each cell's nucleus. We then visualized
115 these Raman profiles in an embedding in two dimensions using Uniform Manifold Approximation
116 and Projection (UMAP)¹⁴ and labeled cells with the gene expression levels that were concurrently
117 measured by either an *Oct4*-GFP reporter or smFISH (**Fig. 2b**). The cells separated clearly in their
118 Raman profiles in a manner consistent with their gene expression characteristics, forming two
119 main subsets in the embedding, one with cells with high *Oct4* and *Nanog* expression (iPSCs
120 markers) and another with cells with relatively high *Coll1a1* expression (fibroblasts marker),
121 indicating that Raman spectra reflect cell-intrinsic expression differences (**Fig. 2b**).

122 We further successfully trained a classifier to classify the 'on' or 'off' expression states of *Oct4*,
123 *Nanog* and *Coll1a1* in each cell based on its Raman profile (**Methods**). We trained a logistic
124 regression classifier with 50% of the data and held out 50% for testing. We predicted *Oct4* and
125 *Nanog* expression states with high accuracy on the held-out test data (area under the receiver
126 operating characteristic curve (AUROC) = 0.98 and 0.95, respectively; **Fig. 2c**), indicating that
127 expression of iPSC markers can be predicted confidently from Raman spectra of live, label-free
128 cells. We also successfully classified the expression state of the fibroblast marker *Coll1a1*
129 (AUROC = 0.87; **Fig. 2c**), albeit with lower confidence, which is consistent with the lower contrast
130 in *Coll1a1* expression (**Fig. 2b**) between iPSC (*Oct4*⁺ or *Nanog*⁺ cells) vs. non-iPSCs, compared
131 to *Oct4* or *Nanog*. Most misclassifications occurred when the ground truth expression levels were
132 near the threshold of the classifier, showing that misclassifications were likely due to the
133 uncertainty in the ground truth expression level (**Extended Data Fig. 5**).

134 Next, we asked if the Raman images could predict entire expression profiles non-destructively at
135 single-cell resolution. To this end, we aimed to reconstruct scRNA-seq profiles from Raman

136 images by multi-modal data integration and translation, using multiplex smFISH data to anchor
137 between the Raman images and scRNA-seq profiles (**Fig. 3a**). As a test case, we focused on the
138 mouse iPSC reprogramming model system, where we have previously generated ~250,000
139 scRNA-seq profiles at ½ day intervals throughout an 18 day, 36 time point time course of
140 reprogramming³ (**Methods**). We used Waddington-OT³ (WOT) to select from the scRNA-seq
141 profiles nine anchor genes that represent diverse cell types that emerge during reprogramming
142 (iPSCs: *Nanog*, *Utf1* and *Epcam*; MET and neural: *Nnat* and *Fabp7*; epithelial: *Krt7* and *Peg10*;
143 stromal: *Bgn* and *Colla1*; **Methods**). We performed live-cell Raman imaging from day 8 of
144 reprogramming, in which distinct cell types begin to emerge³, up to day 14.5, at half-day intervals,
145 totaling 14 time points (**Methods**). We imaged ~500 cells per plate at 1µm spatial resolution.
146 Finally, we fixed cells immediately after each Raman imaging time point followed by smFISH on
147 the 9 anchor genes (**Methods**).

148 Strikingly, a low dimensional representation of the Raman profiles showed that they encoded
149 similar temporal dynamics to those observed with scRNA-seq during reprogramming (**Fig. 3b,c**,
150 **Extended Data Fig. 6**), indicating that they may qualitatively mirror scRNA-seq.

151 Integrating Raman and scRNA-seq profiles (**Methods**), R2R then learned a model that can infer
152 an scRNA-seq profile for each Raman imaged cell, by first predicting smFISH anchors from the
153 Raman profiles using Catboost¹⁵ (**Methods**) and then using our Tangram¹⁶ method to map from
154 the anchors to full scRNA-seq profiles (**Fig. 1, Fig. 3d-f**). In the first step, we averaged the smFISH
155 signal within a nucleus to represent a single nucleus's expression level. As we conducted smFISH
156 of 9 genes, the result was a 9-dimensional smFISH profile for each single nucleus. Then, Raman
157 profiles were translated to these 9-dimensional profiles with Catboost¹⁵, a non-linear regression
158 model, using 50% of the Raman and smFISH profiles as training data.

159 In the second step, we mapped these anchor smFISH profiles to full scRNA-seq profiles using
160 Tangram, yielding well-predicted single cell RNA profiles, as supported by several lines of
161 evidence. First, we performed leave-one-out cross-validation (LOOCV) analysis, in which we used
162 eight out of the nine anchor genes to integrate Raman with scRNA-seq, and compared the predicted
163 expression of the remaining genes to its smFISH measurements. The predicted left-out genes based
164 on scRNA-seq showed a significant correlation with the measured smFISH expression for any left-
165 out gene (Pearson $r \sim 0.7$, $p\text{-value} < 10^{-100}$, **Fig. 3d**). Notably, when we analogously applied a
166 modified U-net¹⁸ to infer smFISH profiles from brightfield (**Extended Data Fig. 15, Methods**),
167 we observed a poor, near-random prediction of expression profiles for all 9 genes in leave-one-out
168 cross-validation ($r < 0.15$), indicating that, unlike Raman spectra, brightfield z-stack images either
169 do not have the necessary information to infer expression profiles, or require more data. Second,
170 we compared the real (scRNA-seq measured) and R2R predicted expression profiles averaged
171 across cells of the same cell type (“pseudobulk” for each of iPSCs, epithelial cells, stromal cells,
172 and MET). Here, we obtained the “ground truth” cell types of the R2R profiles by transferring
173 scRNA-seq annotations to the matching smFISH profiles using Tangram’s label transfer function.
174 Then, based on the labels, we averaged R2R’s predicted profiles across the cells of a single cell
175 type. The two profiles (R2R-inferred and scRNA-seq pseudo-bulk per cell type) showed high
176 correlations (Pearson’s $r > 0.96$) (**Fig. 3e,f, Extended Data Fig. 7**), demonstrating the accuracy of
177 R2R at the cell type level. Furthermore, projecting the R2R predicted profiles of each cell onto an
178 embedding learned from the real scRNA-seq shows that the predicted profiles span the key cell
179 types as captured in real profiles (**Fig. 3g-j, Extended Data Fig. 8-12**). We note that the predicted
180 profiles had lower variance compared to real scRNA-seq. As this is observed even when co-
181 embedding only smFISH and scRNA-seq measurements (with no Raman data or projection,

182 **Extended Data Fig. 13**), we believe it mostly reflects the limited number and domain
183 maladaptation of the smFISH anchor genes used for integration. Given the similarity of the
184 separate embeddings of Raman and scRNA-seq profiles, future studies without anchors could
185 address this.

186 Lastly, we calculated feature importance scores in R2R predictions (**Methods**) and identified
187 Raman spectral features correlated with expression levels (**Fig. 3k, Extended Data Fig. 14**). For
188 example, Raman bands at approximately 752cm^{-1} (C-C, Try, cytochrome), 1004 cm^{-1} (CC, Phe,
189 Tyr), and 1445 cm^{-1} (CH_2 , lipids) contributed to predicting iPSCs-related expression profiles,
190 which is consistent with previous research that employed single cell Raman spectra to identify
191 mouse embryonic stem cells (ESCs)¹⁷ (**Fig. 3k**). The contributions of these bands were either
192 suppressed or increased for other cell types, such as stromal or epithelial cells (**Extended Data**
193 **Fig. 14**).

194 In conclusion, we reported R2R, a label-free non-destructive framework for inferring expression
195 profiles at single-cell resolution from Raman spectra of live cells, by integrating Raman
196 hyperspectral images with scRNA-seq data through paired smFISH measurements and multi-
197 modal data integration and translation. We inferred single-cell expression profiles with high
198 accuracy, based on both averages within cell types and co-embeddings of individual profiles. We
199 further showed that predictions using brightfield z-stacks had poor performance, indicating the
200 importance of Raman microscopy for predicting expression profiles.

201 R2R can be further developed in several ways. First, the throughput of single-cell Raman
202 microscopy is still limited. In this pilot study, we profiled $\sim 6,000$ cells in total. By using emerging
203 vibrational spectroscopy techniques, such as Stimulated Raman Scattering microscopy¹⁹ or photo-

204 thermal microscopy^{20,21}, we envision increasing throughput by several orders of magnitude, to
205 match the throughput of massively parallel single cell genomics. Second, because molecular
206 circuits and gene regulation are structured, with strong co-variation in gene expression profiles
207 across cells, we can leverage the advances in computational microscopy to infer high-resolution
208 data from low-resolution data, such as by using compressed sensing, to further increase
209 throughput²². Third, increasing the number of anchor genes (*e.g.*, by seqFISH²³, merFISH²⁴,
210 STARmap²⁵, or ExSeq²⁶) can increase our prediction accuracy and capture more single-cell
211 variance. Additionally, with single-cell multi-omics, we can project other modalities, such as
212 scATAC-seq from Raman spectra. Finally, given the similarity in the overall independent
213 embedding of Raman and scRNA-seq profiles, we expect computational methods such as multi-
214 domain translation²⁷ to allow mapping between Raman spectra and molecular profiles without
215 measuring any anchors *in situ*. Overall, with further advances in single-cell genomics, imaging,
216 and machine learning, Raman2RNA could allow us to non-destructively infer omics profiles at
217 scale *in vitro*, and possibly *in vivo* in living organisms.

218

219 **Materials and Methods**

220 **Mouse fibroblast reprogramming**

221 OKSM secondary mouse embryonic fibroblasts (MEFs) were derived from E13.5 female embryos
222 with a mixed B6;129 background. The cell line used in this study was homozygous for ROSA26-
223 M2rtTA, homozygous for a polycistronic cassette carrying *Oct4*, *Klf4*, *Sox2*, and *Myc* at the
224 *Coll1a1* 3' end, and homozygous for an EGFP reporter under the control of the *Oct4* promoter.
225 Briefly, MEFs were isolated from E13.5 embryos from timed-matings by removing the head,
226 limbs, and internal organs under a dissecting microscope. The remaining tissue was finely minced
227 using scalpels and dissociated by incubation at 37°C for 10 minutes in trypsin-EDTA
228 (ThermoFisher Scientific). Dissociated cells were then plated in MEF medium containing DMEM
229 (ThermoFisher Scientific), supplemented with 10% fetal bovine serum (GE Healthcare Life
230 Sciences), non-essential amino acids (ThermoFisher Scientific), and GlutaMAX (ThermoFisher
231 Scientific). MEFs were cultured at 37°C and 4% CO₂ and passaged until confluent. All procedures,
232 including maintenance of animals, were performed according to a mouse protocol (2006N000104)
233 approved by the MGH Subcommittee on Research Animal Care³.

234 For the reprogramming assay, 50,000 low passage MEFs (no greater than 3-4 passages from
235 isolation) were seeded in 14 3.5cm quartz glass-bottom Petri dishes (Waken B Tech) coated with
236 gelatin. These cells were cultured at 37°C and 5% CO₂ in reprogramming medium containing
237 KnockOut DMEM (GIBCO), 10% knockout serum replacement (KSR, GIBCO), 10% fetal bovine
238 serum (FBS, GIBCO), 1% GlutaMAX (Invitrogen), 1% nonessential amino acids (NEAA,
239 Invitrogen), 0.055 mM 2-mercaptoethanol (Sigma), 1% penicillin-streptomycin (Invitrogen) and
240 1,000 U/ml leukemia inhibitory factor (LIF, Millipore). Day 0 medium was supplemented with 2

241 mg/mL doxycycline Phase-1 (Dox) to induce the polycistronic OKSM expression cassette. The
242 medium was refreshed every other day. On day 8, doxycycline was withdrawn. Fresh medium was
243 added every other day until the final time point on day 14. One plate was taken every 0.5 days
244 after day 8 (D8-D14.5) for Raman imaging and fixed with 4% formaldehyde immediately after for
245 HCR.

246 **High-throughput multi-modal Raman microscope**

247 Due to the lack of commercial systems, we developed an automated high-throughput multi-modal
248 microscope capable of multi-position and multi-timepoint fluorescence imaging and point
249 scanning Raman microscopy (**Extended Data Fig. 1**). A 749 nm short-pass filter was placed to
250 separate brightfield and fluorescence from Raman scattering signal, and the fluorescence and
251 Raman imaging modes were switched by swapping dichroic filters with auto-turrets. To realize a
252 high-throughput Raman measurement, galvo mirror-based point scanning and stage scanning was
253 combined to acquire each FOV and multiple different FOVs, respectively.

254 To realize this in an automated fashion, a MATLAB (2020b) script that communicates with Micro-
255 manager²⁸, a digital acquisition (DAQ) board, and Raman scattering detector (Princeton
256 Instruments, PIXIS 100BR eXcelon) was written (**Extended Data Fig. 2**). A 2D point scan Raman
257 imaging sequence was regarded as a dummy image acquisition in Micro-manager, during which
258 the script communicated via the DAQ board with 1. the detector to read out a spectrum, 2. the
259 mirror to update the mirror angles, and 3. shutters to control laser exposure. All communications
260 were realized using transistor-transistor logic (TTL) signaling. Updating of the galvo mirror angles
261 was conducted during the readout of the detector. While the script ran in the background, Micro-

262 manager initiated a multi-dimensional acquisition consisting of brightfield, DAPI, GFP, and
263 dummy Raman channel at multiple positions and z-stacks.

264 An Olympus IX83 fluorescence microscope body was integrated with a 785 nm Raman excitation
265 laser coupled to the backport, where the short-pass filter deflected the excitation to the sample
266 through an Olympus UPLSAPO 60X NA 1.2 water immersion objective. The backscattered light
267 was collimated through the same objective and collected with a 50 μm core multi-mode fiber,
268 which was then sent to the spectrograph (Holospec f/1.8i 785 nm model) and detector. The
269 fluorescence and brightfield channels were imaged by the Orca Flash 4.0 v2 sCMOS camera from
270 Hamamatsu Photonics. The exposure time for each point in the Raman measurement was 20 msec,
271 and laser power at the sample plane was 212 mW. Each FOV was 100x100 pixels, with each pixel
272 corresponding to about 1 μm . The laser source was a 785 nm Ti-Sapphire laser cavity coupled to
273 a 532 nm pump laser operating at 4.7W.

274 The time to acquire Raman hyperspectral images was roughly 8 minutes per FOV. With 8 minutes,
275 it is unrealistic to image an entire glass-bottom plate. Therefore, we visually chose representative
276 FOVs that cover all representative cell types including iPSC-like, epithelial-like, stromal-like and
277 MET cells. 20 FOVs were chosen for each plate, where roughly 15 FOVs were from the boundaries
278 of colonies, five from non-colonies, and one from non-cells to use for background correction.

279 Due to the extended Raman imaging time, evaporation of the immersion water was no longer
280 negligible. Therefore, we developed an automated water immersion feeder using syringe pumps
281 and syringe needles glued to the tip of the objective lens. Here, water was supplied at a flow rate
282 of 1 $\mu\text{L}/\text{min}$.

283 **iPSC and MEF mixture experiment**

284 Low passage iPSCs were first cultured in N2B27 2i media containing 3 mM CHIR99021, 1 mM
285 PD0325901, and LIF. On the day of the experiment, 750,000 iPSCs and 750,000 MEFs were plated
286 on the same gelatin-coated 3.5cm quartz glass-bottom Petri dish. Cells were plated in the same
287 reprogramming medium as previously described (with Dox) with the exception of utilizing DMEM
288 without phenol red (Gibco) instead of KnockOut DMEM. 6 hours after plating, the quartz dishes
289 were taken for Raman imaging and fixed with 4% formaldehyde immediately after for HCR.

290 **Anchor gene selection by Waddington-OT**

291 To select anchor genes for connecting spatial information to the full transcriptome data,
292 Waddington-OT (WOT)³, a probabilistic time-lapse algorithm that can reconstruct developmental
293 trajectories, was used. We applied WOT to mouse fibroblast reprogramming scRNA-seq data
294 collected at matching time-points and culture condition (day 8-14.5 at ½ day intervals)³. For each
295 cell fate, we calculated the transition probabilities of each cell and selected the top 10 percentile
296 cells per time point (**Extended Data Fig. 6**). Based on this, we ran the *FindMarker* function in
297 Seurat²⁹ to find genes differentially expressed in these cell subsets per time point. Through this
298 approach, we chose two genes per cell type that are both found by Seurat and commonly used for
299 these cell types (iPSCs: *Nanog*, *Utf1*; epithelial: *Krt7*, *Peg10*; stromal: *Bgn*, *Coll1a1*; MET and
300 neural: *Fabp7*, *Nnat*), along with one gene that is an early marker of iPSCs, *Epcam*.

301 **smRNA-FISH by hybridization chain reaction (HCR)**

302 Fixed samples were prepared for imaging using the HCR v3.0 protocol for mammalian cells on a
303 chambered slide, incubating at the amplification step for 45 minutes in the dark at room

304 temperature. Three probes with amplifiers conjugated to fluorophores Alexa Fluor 488, Alexa
305 Fluor 546, and Alexa Fluor 647 were used. Samples were stained with DAPI prior to imaging.
306 After imaging, probes were stripped from samples by washing samples once for 5 minutes in 80%
307 formamide at room temperature and then incubating three times for 30 minutes in 80% formamide
308 at 37°C. Samples were washed once more with 80% formamide, then once with PBS, and reprobed
309 with another panel of probes for subsequent imaging.

310 **Image registration of Raman hyperspectral images and fluorescence/smFISH images**

311 Brightfield and fluorescence channels including DAPI and GFP, along with corresponding Raman
312 images, were registered by using 5 μm polystyrene beads deposited on quartz glass-bottom Petri
313 dishes (SF-S-D12, Waken B Tech) for calibration. The brightfield and fluorescence images of the
314 beads were then registered by the scale-invariant template matching algorithm of the OpenCV
315 (<https://github.com/opencv/opencv>) *matchTemplate* function followed by manual correction.

316 For the registration of smFISH and Raman images, four marks inscribed under the glass-bottom
317 Petri dishes were used as reference points (**Extended Data Fig. 4**). As the Petri dishes are
318 temporarily removed from the Raman microscope after imaging to do smFISH measurements, the
319 dishes cannot be placed back at the same exact location on the microscope. Therefore, the
320 coordinates of these reference points were measured along with the different FOVs. When the
321 dishes were placed again after smFISH measurements, the reference mark coordinates were
322 measured, and an affine mapping was constructed to calculate the new FOV coordinates. Lastly,
323 as smFISH consisted of 3 rounds of hybridization and imaging, the following steps were performed
324 to register images across different rounds with a custom MATLAB script:

325 1. Maximum intensity projection of nuclei stain and RNA images

- 326 2. Automatic registration of round 1 images to rounds 2 and 3 based on nuclei stain images
327 and MATLAB function *imregtform*. First, initial registration transformation functions were
328 obtained with a similarity transformation model passing the ‘multimodal’ configuration.
329 Then, those transformations were used as the initial conditions for an affine model-based
330 registration with the *imregtform* function. Finally, this affine mapping transformation was
331 applied to all the smFISH (RNA) images.
- 332 3. Use the protocol in (2) to register nuclei stain images obtained from the multimodal Raman
333 microscope and the 1st round of images used for smFISH. Then, apply the transformation
334 to the remaining 2nd and 3rd rounds.
- 335 4. Manually remove registration outliers in (3).

336 Fibroblast cells were mobile during the 2-class mixture experiment so that by the time Raman
337 imaging finished, cells had moved far enough from their original position that the above semi-
338 automated strategy could not be applied. Thus, we manually identified cells present in both nuclei
339 stain images before and after the Raman imaging.

340 **Hyperspectral Raman image processing**

341 Each raw Raman spectrum has 1,340 channels. Of those channels, we extracted the fingerprint
342 region (600-1800 cm^{-1}), which resulted in a total of 930 channels per spectrum. Thus, each FOV
343 is a 100x100x930 hyperspectral image. The hyperspectral images were then preprocessed by a
344 python script as follows:

- 345 1. Cosmic ray removal. Cosmic rays were detected by subtracting the median filtered spectra
346 from the raw spectra, and any feature above 5 was classified as an outlier and replaced with
347 the median value. The kernel window size for the median filter was 7.
- 348 2. Autofluorescence removal. The *baseline* function in *rampy*
349 (<https://github.com/charlesll/rampy>), a python package for Raman spectral preprocessing,
350 was used with the alternating least squares algorithm '*als*'.
- 351 3. Savitzky-Golay smoothing. The *scipy.signal.savgol_filter* function was used with window
352 size 5 and polynomial order 3.
- 353 4. Averaging spectra at the single-cell level. Nuclei stain images were segmented using
354 *NucleAIzer* (<https://github.com/spreka/biomagdsb>) and averaged pixel-level spectra that
355 fall within each nucleus.
- 356 5. Spectra standardization. Spectra were standardized to a mean of 0 and a standard deviation
357 of 1.

358 **Inferring anchor smFISH from Raman spectra or brightfield z-stacks**

359 For the two-class mixture and reprogramming experiment, we trained a decision tree-based non-
360 linear regression, *Catboost*¹⁵, to predict the 'on' or 'off' expression states for each anchor gene
361 from Raman spectra. We used 80% of the data as training and the remaining 20% as test data. The
362 early stopping parameter was set to 5.

363 For the brightfield z-stack to smFISH inference, we applied deep learning to the whole image level.
364 We trained a modified U-net with skip connections and residual blocks to estimate the
365 corresponding smFISH image¹⁸. Due to the small size of the available training dataset, we
366 augmented the data by rotation and flipping. Furthermore, a subsample of each brightfield image

367 was taken due to memory constraints (50x50 pixel region). Training was carried out on an NVIDIA
368 Tesla P100 GPU, the number of epochs was 100, the learning rate was 0.01, and the batch size
369 was 400. For each smFISH prediction, we chose the epoch that gave the best validation score.

370 **Inferring expression profiles from Raman images**

371 To infer expression profiles from Raman images, we used Tangram¹⁶. Tangram enables the
372 alignment of spatial measurements of a small number of genes to scRNA-seq measurements. After
373 using Catboost to infer anchor expression levels from Raman profiles, we aligned the inferred
374 expression levels to scRNA-seq profiles using the *map_cells_to_space* function
375 (learning_rate=0.1, num_epochs=1000) on an Nvidia Tesla P100 GPU, followed by the
376 *project_genes* function in Tangram.

377 When comparing different pseudo-bulk transcriptome predictions with the real scRNA-seq data,
378 we first transferred labels of annotated scRNA-seq profiles to the ground truth smFISH profiles
379 using Tangram's label transfer function *project_cell_annotations*. Then, the average expression
380 profiles across cells of a cell type were calculated by referring to the transferred labels and
381 compared with those from the real scRNA-seq data³.

382 **Dimensionality reduction, embedding and projection**

383 For dimension reduction and visualization of Raman and scRNA-seq profiles, we performed
384 forced layout embedding (FLE) using the *Pegasus* pipeline ([https://github.com/klarman-cell-](https://github.com/klarman-cell-observatory/pegasus)
385 [observatory/pegasus](https://github.com/klarman-cell-observatory/pegasus)). First, we performed principal component analysis on both Raman and
386 scRNA-seq profiles independently, calculated diffusion maps on the top 100 principal
387 components, and performed an approximated FLE graph using Deep Learning by *pegasus.net_fle*
388 with default parameters.

389 To project Raman profiles to a scRNA-seq embedding, we calculated a k-nearest neighbor graph
390 (k -NN, $k=15$) on the scRNA-seq top 50 principal components with the cosine metric, and UMAP
391 with the *scanpy.tl.umap* function in Scanpy³⁰ version 1.7.2 with default parameters. Then, the
392 Raman predicted expression profiles were projected on to the scRNA-seq UMAP embedding by
393 *scanpy.tl.ingest* using k-NN as the labeling method and default parameters.

394 **Feature importance analysis**

395 To evaluate the contributions of Raman spectral features to expression profile prediction, we used
396 the *get_feature_importance* function in Catboost with default parameters. As the dimensions of
397 Raman spectra were reduced by PCA prior to Catboost, feature importance scores were calculated
398 for each principal component, and the weighted linear combination of the Raman PCA eigen
399 vectors with feature scores as the weight were calculated to obtain the full spectrum.

400 **Author contributions**

401 KJKK, JS, TB and AR conceived the research and developed the methodology. JS, TB and AR
402 funded and supervised research. KJKK, JS, JO performed reprogramming experiments. KJKK
403 developed the multi-modal Raman microscope and control software with supervision from JWK
404 and PS. KJKK, EG, and KZ performed smFISH. KJKK, SG, TJS, and TB developed the Raman
405 spectral preprocessing and classification pipeline. KJKK developed the image registration
406 pipeline, and performed Waddington-OT, Tangram and feature importance analysis. KJKK and
407 BG performed U-net. KJKK, JS, and AR wrote the manuscript with input from all the authors.

408 **Competing interests statement**

409 AR is a co-founder and equity holder of Celsius Therapeutics, an equity holder in Immunitas, and
410 was a scientific advisory board member of ThermoFisher Scientific, Syros Pharmaceuticals,
411 Neogene Therapeutics and Asimov until 31 July 2020. AR, TB, and SG are employees of
412 Genentech from August 1, 2020, respectively. A patent application has been filed related to this
413 work.

414 **Acknowledgements**

415 KJKK was supported by the Japan Society for the Promotion of Science Postdoctoral Fellowship
416 for Overseas Researchers, and the Naito Foundation Overseas Postdoctoral Fellowship. BG was
417 supported by the MathWorks Fellowship. JS was supported by the Helen Hay Whitney Foundation
418 and NIH Pathway to Independence Award (1K99HD096049-01, 5K99HD096049-02,
419 4R00HD096049-03), and funds from the Broad Institute of MIT and Harvard and Massachusetts
420 General Hospital. This research was funded by NIH National Institute of Biomedical Imaging and
421 Bioengineering, grant P41EB015871 (JWK, PS), NIH grant U19 MH114821 (TB), HubMap
422 UH3CA246632 (TB), and HHMI and the Klarman Cell Observatory (AR). AR was a Howard
423 Hughes Medical Institute Investigator when this work was initiated. We thank Eric Lander, Rudolf
424 Jaenisch, Doeke Hekstra, Joseph Kirschvink for their helpful discussion and insights. We thank
425 Leslie Gaffney for creating and editing figures.

426

427 **References**

- 428 1. Tanay, A. & Regev, A. Scaling single-cell genomics from phenomenology to mechanism.
429 *Nature* **541**, 331–338 (2017).
- 430 2. Trapnell, C. *et al.* The dynamics and regulators of cell fate decisions are revealed by
431 pseudotemporal ordering of single cells. *Nat. Biotechnol.* **32**, 381–386 (2014).
- 432 3. Schiebinger, G. *et al.* Optimal-Transport Analysis of Single-Cell Gene Expression Identifies
433 Developmental Trajectories in Reprogramming. *Cell* **176**, 928-943.e22 (2019).
- 434 4. La Manno, G. *et al.* RNA velocity of single cells. *Nature* **560**, 494–498 (2018).
- 435 5. Bergen, V., Lange, M., Peidli, S., Wolf, F. A. & Theis, F. J. Generalizing RNA velocity to
436 transient cell states through dynamical modeling. *Nat. Biotechnol.* **38**, 1408–1414 (2020).
- 437 6. Wagner, D. E. & Klein, A. M. Lineage tracing meets single-cell omics: opportunities and
438 challenges. *Nat. Rev. Genet.* **21**, 410–427 (2020).
- 439 7. Wei, L. *et al.* Super-multiplex vibrational imaging. *Nature* **544**, 465–470 (2017).
- 440 8. Kobayashi-Kirschvink, K. J. *et al.* Linear Regression Links Transcriptomic Data and
441 Cellular Raman Spectra. *Cell Systems* vol. 7 104-117.e4 (2018).
- 442 9. Singh, S. P. *et al.* Label-free characterization of ultra violet-radiation-induced changes in
443 skin fibroblasts with Raman spectroscopy and quantitative phase microscopy. *Sci. Rep.* **7**,
444 10829 (2017).
- 445 10. Ichimura, T. *et al.* Visualizing cell state transition using Raman spectroscopy. *PLoS One* **9**,
446 e84478 (2014).
- 447 11. Ho, C.-S. *et al.* Rapid identification of pathogenic bacteria using Raman spectroscopy and
448 deep learning. *Nat. Commun.* **10**, 4927 (2019).
- 449 12. Stadtfeld, M., Maherali, N., Borkent, M. & Hochedlinger, K. A reprogrammable mouse
450 strain from gene-targeted embryonic stem cells. *Nat. Methods* **7**, 53–55 (2010).

- 451 13. Choi, H. M. T. *et al.* Third-generation in situ hybridization chain reaction: multiplexed,
452 quantitative, sensitive, versatile, robust. *Development* **145**, (2018).
- 453 14. McInnes, L., Healy, J., Saul, N. & Großberger, L. UMAP: Uniform Manifold
454 Approximation and Projection. *J. Open Source Softw.* **3**, 861 (2018).
- 455 15. Prokhorenkova, L., Gusev, G., Vorobev, A., Dorogush, A. V. & Gulin, A. CatBoost:
456 unbiased boosting with categorical features.
- 457 16. Biancalani, T. *et al.* Deep learning and alignment of spatially-resolved whole transcriptomes
458 of single cells in the mouse brain with Tangram. *bioRxiv* 2020.08.29.272831 (2020)
459 doi:10.1101/2020.08.29.272831.
- 460 17. Germond, A., Panina, Y., Shiga, M., Niioka, H. & Watanabe, T. M. Following Embryonic
461 Stem Cells, Their Differentiated Progeny, and Cell-State Changes During iPS
462 Reprogramming by Raman Spectroscopy. *Anal. Chem.* **92**, 14915–14923 (2020).
- 463 18. He, K., Zhang, X., Ren, S. & Sun, J. Deep residual learning for image recognition. in *2016*
464 *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (IEEE, 2016).
465 doi:10.1109/cvpr.2016.90.
- 466 19. Freudiger, C. W. *et al.* Label-free biomedical imaging with high sensitivity by stimulated
467 Raman scattering microscopy. *Science* **322**, 1857–1861 (2008).
- 468 20. Bai, Y. *et al.* Ultrafast chemical imaging by widefield photothermal sensing of infrared
469 absorption. *Sci Adv* **5**, eaav7127 (2019).
- 470 21. Tamamitsu, M., Toda, K., Horisaki, R. & Ideguchi, T. Quantitative phase imaging with
471 molecular vibrational sensitivity. *Opt. Lett.* **44**, 3729–3732 (2019).

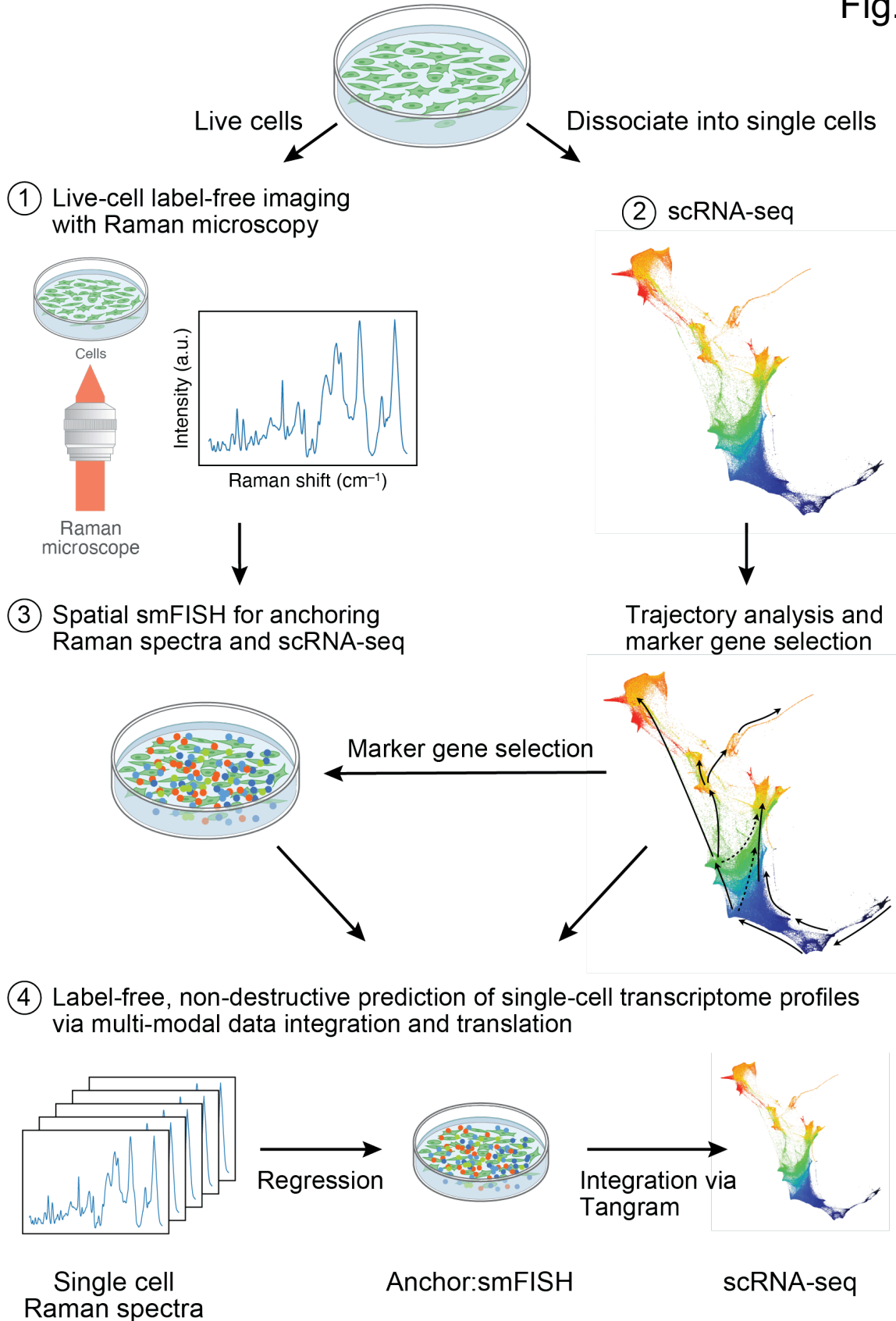
- 472 22. Cleary, B., Cong, L., Cheung, A., Lander, E. S. & Regev, A. Efficient Generation of
473 Transcriptomic Profiles by Random Composite Measurements. *Cell* **171**, 1424-1436.e18
474 (2017).
- 475 23. Eng, C.-H. L. *et al.* Transcriptome-scale super-resolved imaging in tissues by RNA
476 seqFISH. *Nature* **568**, 235–239 (2019).
- 477 24. Chen, K. H., Boettiger, A. N., Moffitt, J. R., Wang, S. & Zhuang, X. RNA imaging.
478 Spatially resolved, highly multiplexed RNA profiling in single cells. *Science* **348**, aaa6090
479 (2015).
- 480 25. Wang, X. *et al.* Three-dimensional intact-tissue sequencing of single-cell transcriptional
481 states. *Science* (2018) doi:10.1126/science.aat5691.
- 482 26. Alon, S. *et al.* Expansion sequencing: Spatially precise in situ transcriptomics in intact
483 biological systems. *Science* **371**, (2021).
- 484 27. Yang, K. D. *et al.* Multi-domain translation between single-cell imaging and sequencing
485 data using autoencoders. *Nat. Commun.* **12**, 31 (2021).
- 486 28. Edelstein, A., Amodaj, N., Hoover, K., Vale, R. & Stuurman, N. Computer control of
487 microscopes using μ Manager. *Curr. Protoc. Mol. Biol.* **Chapter 14**, Unit14.20 (2010).
- 488 29. Stuart, T. *et al.* Comprehensive Integration of Single-Cell Data. *Cell* **177**, 1888-1902.e21
489 (2019).
- 490 30. Wolf, F. A., Angerer, P. & Theis, F. J. SCANPY: large-scale single-cell gene expression
491 data analysis. *Genome Biol.* **19**, 15 (2018).

492

493

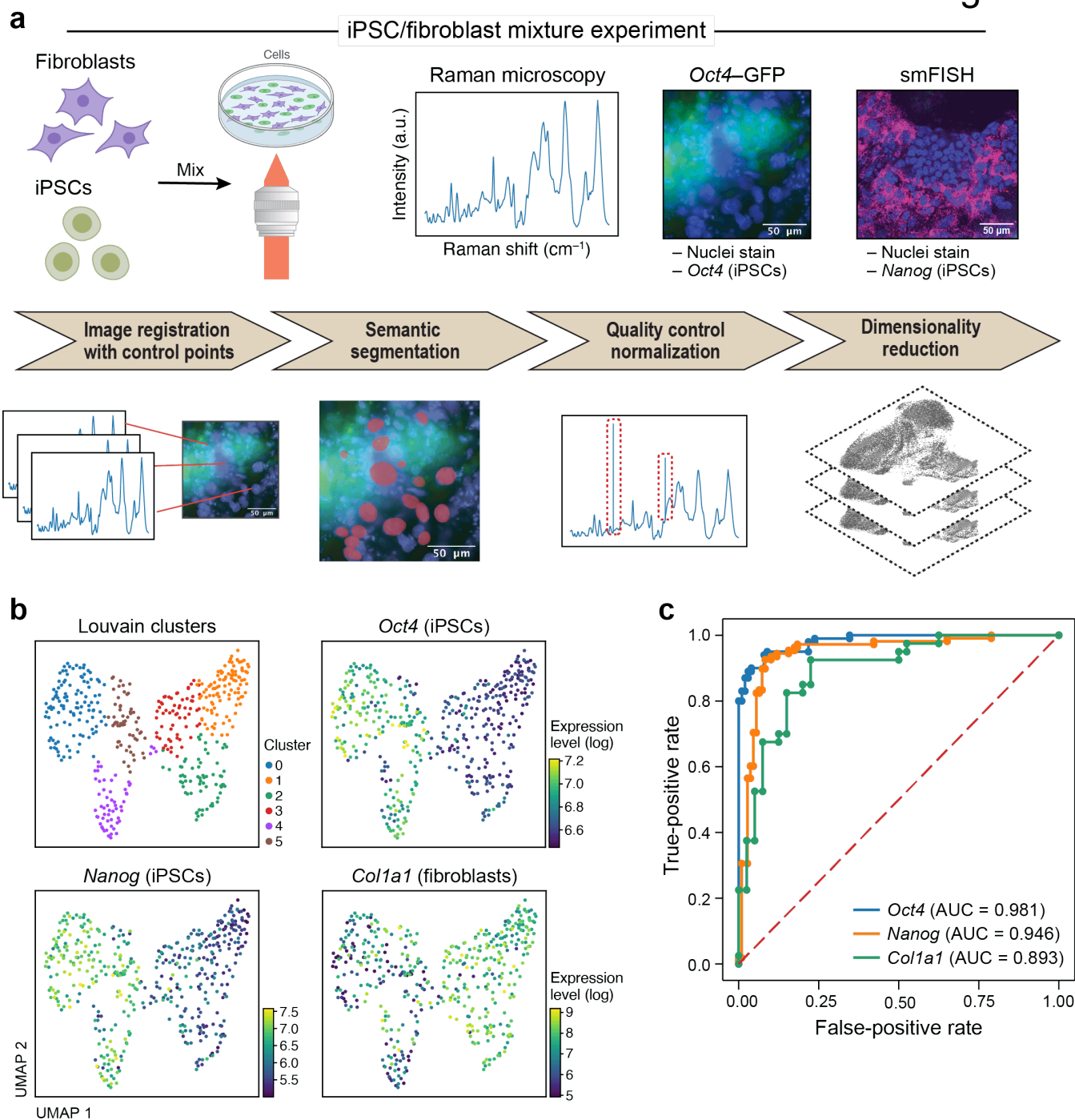
Live cell culture

Fig. 1



494 **Fig. 1 | Raman2RNA.** Live cells are cultured on gelatin-coated quartz glass-bottom plates (top) and
495 Raman spectra are then measured at each pixel (at spatial sub-cellular resolution) within an image frame
496 (1), followed by smFISH imaging in the same area (3). From parallel plates, cells are dissociated into a
497 single cell suspension and profiled by scRNA-seq (2). scRNA-seq profiles are used to select 9 marker
498 genes for 5 major cell clusters, and those are measured with spatial smFISH (3). Lastly, a regression
499 model is trained (4) to predict anchor smFISH profiles from Raman spectra, followed by integration via
500 Tangram¹ to predict whole single-cell transcriptome profiles from smFISH profiles.

Fig. 2



501

502 **Fig. 2 | Raman2RNA accurately distinguishes cell types and predicts binary expression of marker**

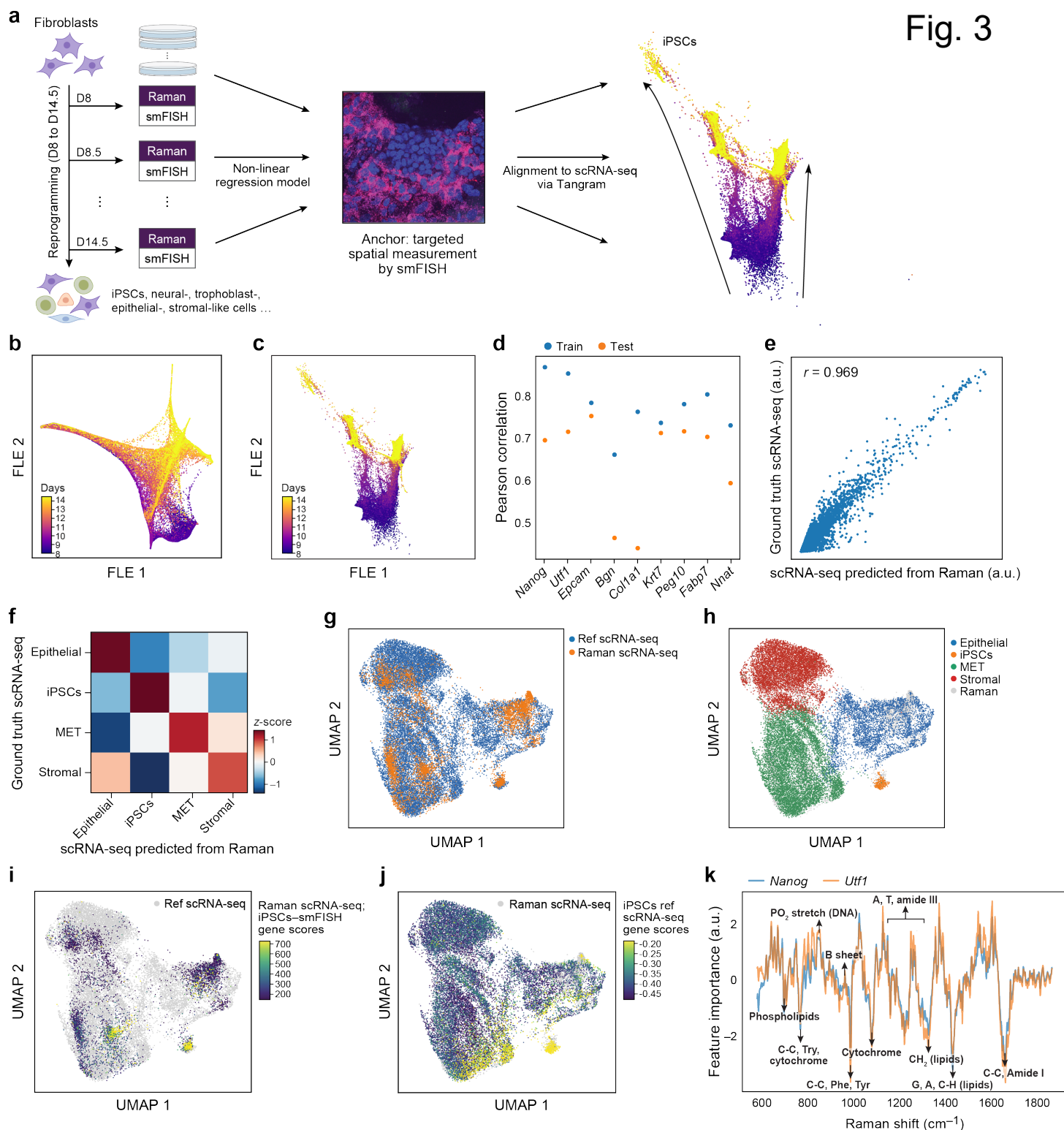
503 **genes in a mixture of mouse fibroblasts and iPSCs. a.** Overview. Top: Experimental procedures.

504 Mouse fibroblasts and iPSCs were mixed 1:1 and plated on glass-bottom plates, followed by Raman

505 imaging of live cells, nuclei staining and measurement of endogenous *Oct4*-GFP (iPSC marker) reporter)

506 by fluorescence imaging, and cell fixation and processing for smFISH with DAPI and probes for *Nanog*

507 (iPSCs, magenta) and *Coll1a1* (fibroblasts). Bottom: Preprocessing and analysis. From left: Image
508 registration with control points (**Methods**), was followed by semantic cell segmentation, outlier
509 removal/normalization and dimensionality reduction. **b.** Raman2RNA distinguishes cell states from
510 Raman spectra. 2D UMAP embedding of single-cell Raman spectra (dots) colored by Louvain clustering
511 labels (top left) or smFISH measured expression of *Oct4* (top right), *Nanog* (bottom left) and *Coll1a1*
512 (bottom right). **c.** Raman2RNA accurately predicts binary (on/off) expression of marker genes. Receiver
513 operating characteristic (ROC) plots and area under the curve (AUC) obtained by classifying the ‘on’ and
514 ‘off’ states of *Oct4* (blue), *Nanog* (orange) and *Coll1a1* (green).



515

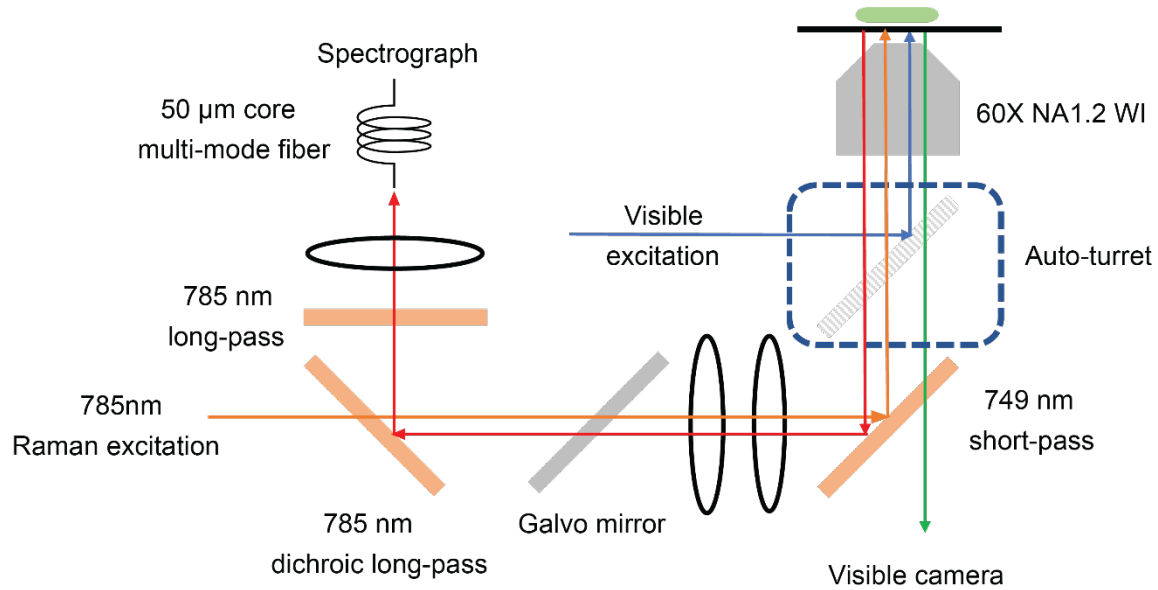
516 **Fig. 3 | Raman2RNA predicts single-cell RNA profiles across cell types during reprogramming of**

517 **mouse fibroblasts to iPSCs. a.** Approach overview. From left: Mouse fibroblasts were reprogrammed

518 into induced pluripotent stem cells (iPSCs) over the course of 14.5 days ('D'), and, at half-day intervals
519 from days 8 to 14.5, spatial Raman spectra, smFISH for nine anchor genes, and nuclei stain by
520 fluorescence imaging were measured for each plate. Machine learning and multi-modal data integration
521 methods (Catboost and Tangram) were used to predict single-cell RNA-seq profiles from Raman spectra
522 using smFISH as anchor. **b,c.** Low dimensionality embedding of single-cell Raman spectra captures
523 progress in reprogramming. Force-directed layout embedding (FLE) of Raman spectra (b, dots) or
524 scRNA-seq (c, dots) colored by days of measurement (colorbar). **d.** Correct prediction of smFISH anchors
525 from Raman spectra. Pearson correlation coefficient (y axis) between measured (smFISH) and Raman-
526 predicted levels for each smFISH anchor (x axis) in leave-one-out cross-validation where 8 out of 9
527 smFISH anchor genes were used for training, and the left-out gene was predicted. **e.f.** Raman2RNA
528 accurately predicts pseudo-bulk expression profiles of major cell types. **e.** scRNA-seq measured (y axis)
529 and R2R-predicted (x axis) for each gene (dot) in pseudo-bulk RNA profiles averaged across iPSCs. **f.**
530 Pair-wise correlation (color bar) between Raman-predicted and scRNA-seq measured pseudo-bulk
531 profiles in each cell types (rows, columns). **g-j.** Co-embedding highlights agreement between real and
532 R2R inferred single cell profiles. UMAP co-embedding of Raman predicted RNA profiles and measured
533 scRNA-seq profiles (dots) colored by data source (**g**, Raman predicted in orange; measured scRNA-seq in
534 blue), cell type annotations (**h**) or by iPSC gene signature scores (calculated by averaging expression of
535 genes *Nanog* and *Utf1*, and subtracting the average of a randomly selected set of reference genes;
536 **Methods**) of Raman-predicted profiles (**i**) or of real scRNA-seq (**j**). **k.** Feature importance scores of
537 Raman spectra in predicting expression profiles. Feature scores for iPSC related marker genes (y axis)
538 along the Raman spectrum (x axis). Known Raman peaks² were annotated.
539

540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562

Supp. Fig. 1

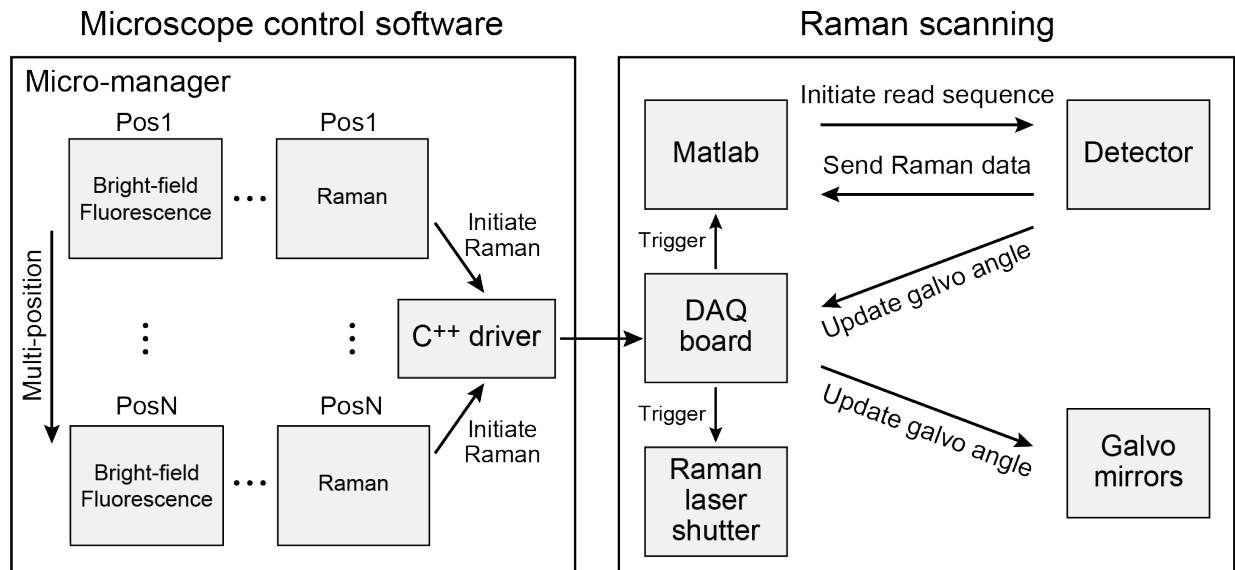


Extended Data Fig. 1 | A multi-modal Raman microscope capable of fluorescence imaging and

Raman microscopy. Schematic of a Raman microscope integrated with a wide-field fluorescence microscope for simultaneous detection of nuclei staining, bright field, fluorescence channels, and Raman images.

563
564
565

Supp. Fig. 2



566

567 **Extended Data Fig. 2 | Overview of high-throughput Raman imaging software used in the study.** A
568 general-purpose microscope control software Micro-manager and a custom MATLAB script were
569 combined to enable automated multi-modal measurements. Under Micro-manager, a Raman channel was
570 registered as a ‘dummy’ channel along with brightfield and fluorescence channels. Micro-manager was
571 responsible for changing the field of view (FOV) and imaging modality. During the Raman sequence,
572 Micro-manager communicated with a digital acquisition (DAQ) board, through which a transistor-to-
573 transistor logic (TTL) signal was generated to initiate the scanning sequence. Upon detection of the TTL
574 signal, the MATLAB script controlled the Raman detector, laser shutter, and updated the galvo mirror
575 angles through the DAQ board.

576

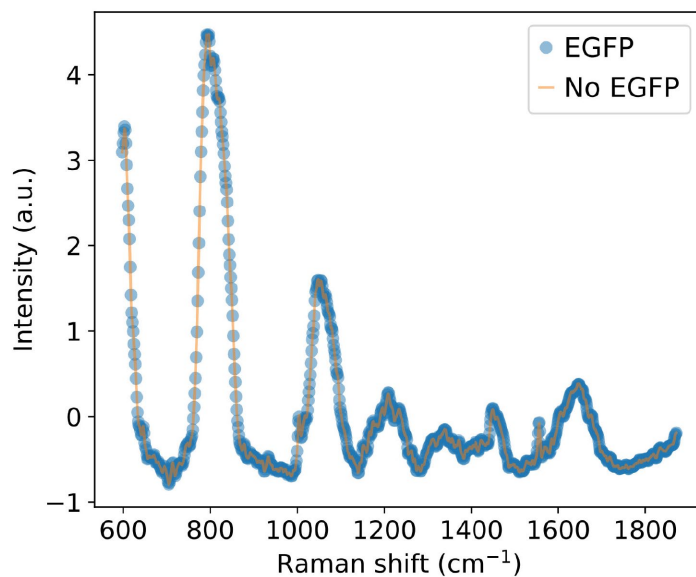
577

578

579

580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599

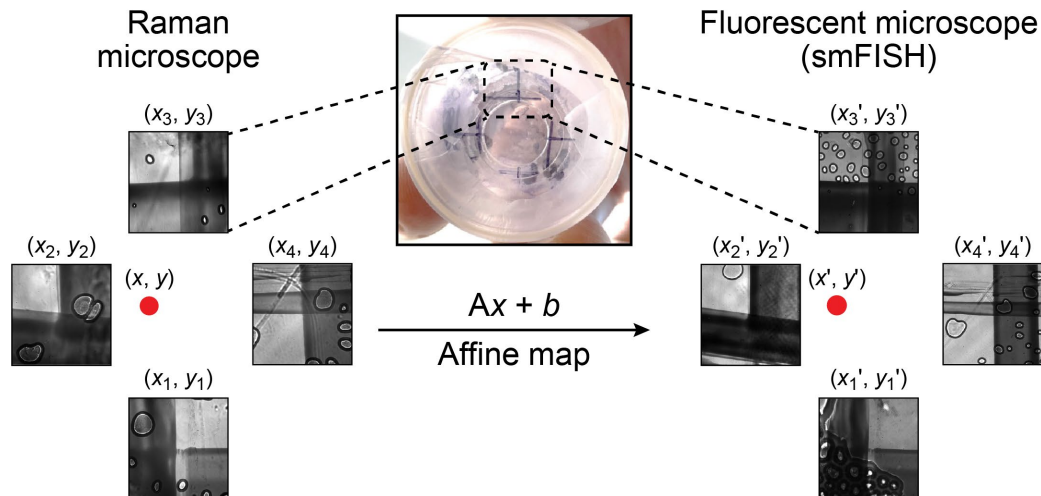
Supp. Fig. 3



Extended Data Fig. 3 | GFP does not interfere in Raman spectra measurement. Raman spectra of culture media with (blue) and without (orange) GFP at physiological concentration.

600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622

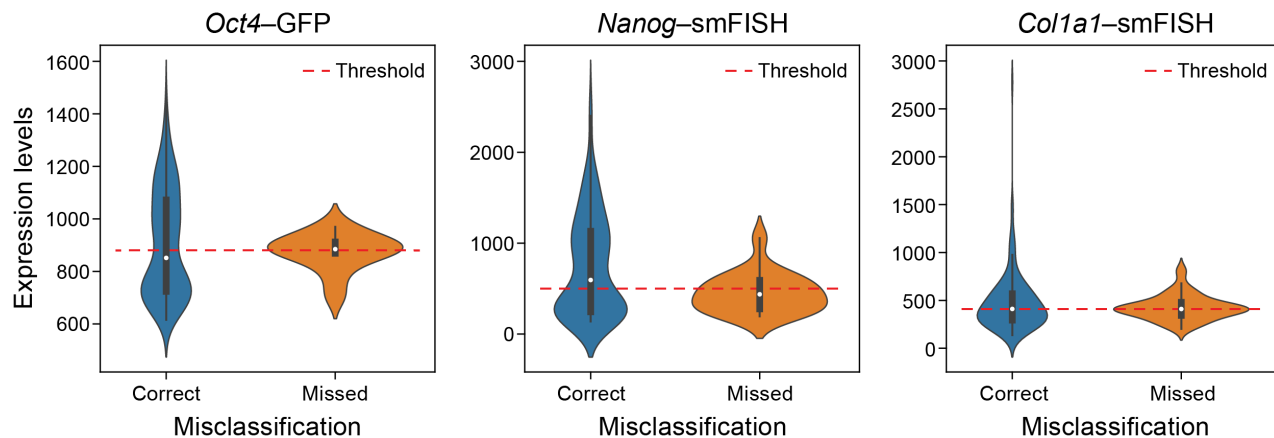
Supp. Fig. 4



Extended Data Fig. 4 | Image registration between the Raman and smFISH microscope using control points. Control points were inscribed under petri dishes with permanent markers and the coordinates were measured prior to any data acquisition. After Raman measurement and smFISH processing, samples were placed back to the microscope and control point coordinates were remeasured. Then, affine mapping was used to update the FOV coordinates to locate the exact same cells.

623
624
625
626

Supp. Fig. 5

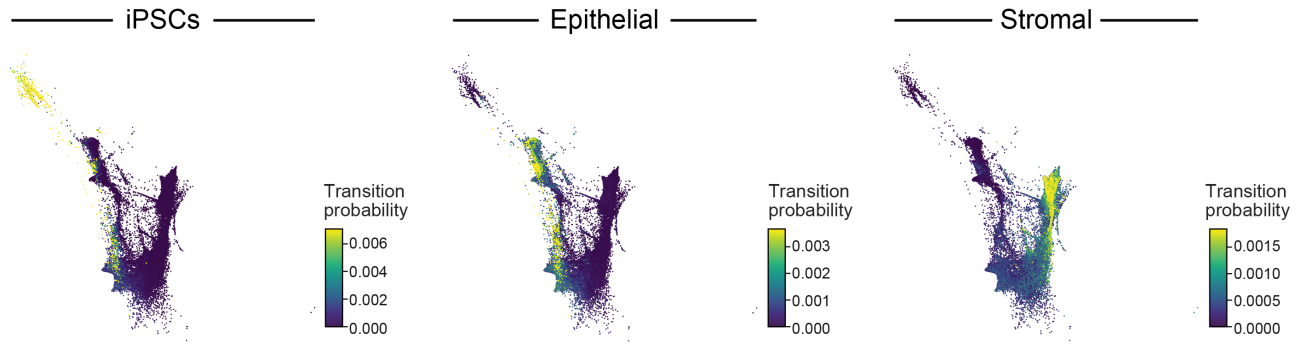


627
628
629
630
631
632
633
634
635

Extended Data Fig. 5 | Misclassification of genes in the cell mixture classification experiment occurs when the ground truth smFISH is near the expression threshold. Distribution of measured smFISH expression level (y axis) for cells correctly (blue) or incorrectly (orange) classified by their Raman spectra for the expression of that gene. Horizontal line: an example threshold used for the logistic regression classifier.

636
637
638
639
640

Supp. Fig. 6

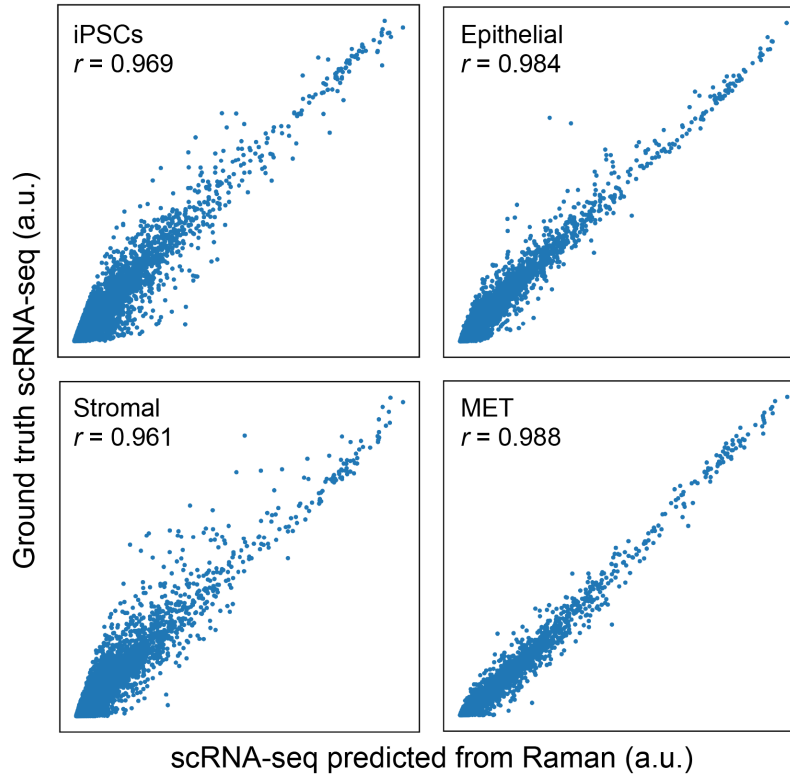


641
642
643
644
645
646
647
648
649

Extended Data Fig. 6 | Cell transition probabilities inferred by Waddington-OT from scRNA-seq during reprogramming. Force-directed layout embedding (FLE) of scRNA-seq profiles (dots) from days 8 to 14.5 of reprogramming (dots) colored by the transition probability of each cell as inferred by Waddington-OT to be an ancestor of iPSCs (left), epithelial cells (middle) or stromal cells (right) at day 14.5.

650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675

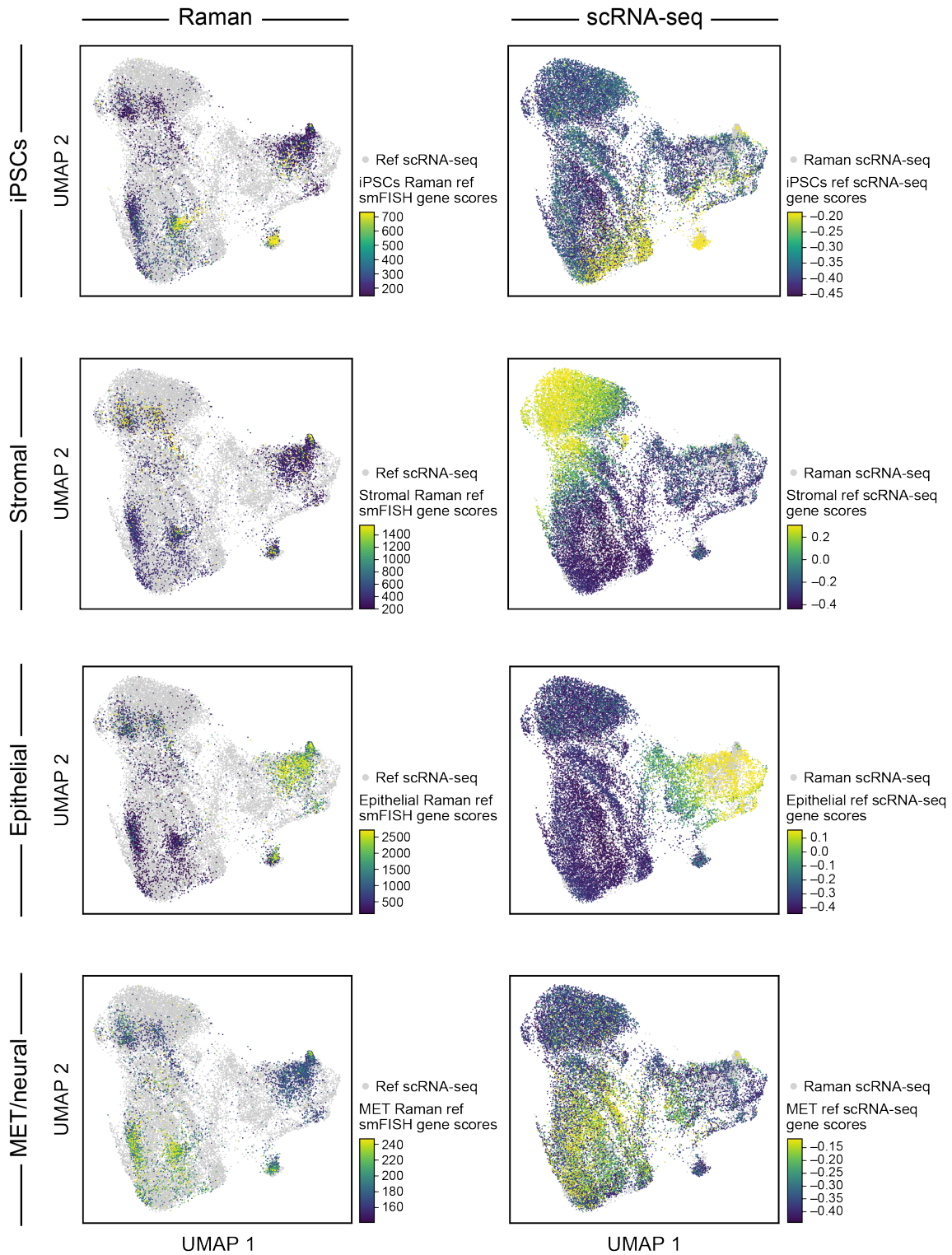
Supp. Fig. 7



Extended Data Fig. 7 | Raman-predicted and scRNA-seq measured pseudo-bulk profiles are well correlated across cell types. ScRNA-seq measured (y axis) and R2R-predicted (x axis) expression for each gene (dot) in pseudo-bulk RNA profiles averaged across cells labeled as iPSC (top left), epithelial (top right), stromal (bottom left) and MET (bottom right). Pearson's r is denoted at the top left corner.

Supp. Fig. 8

676
677
678

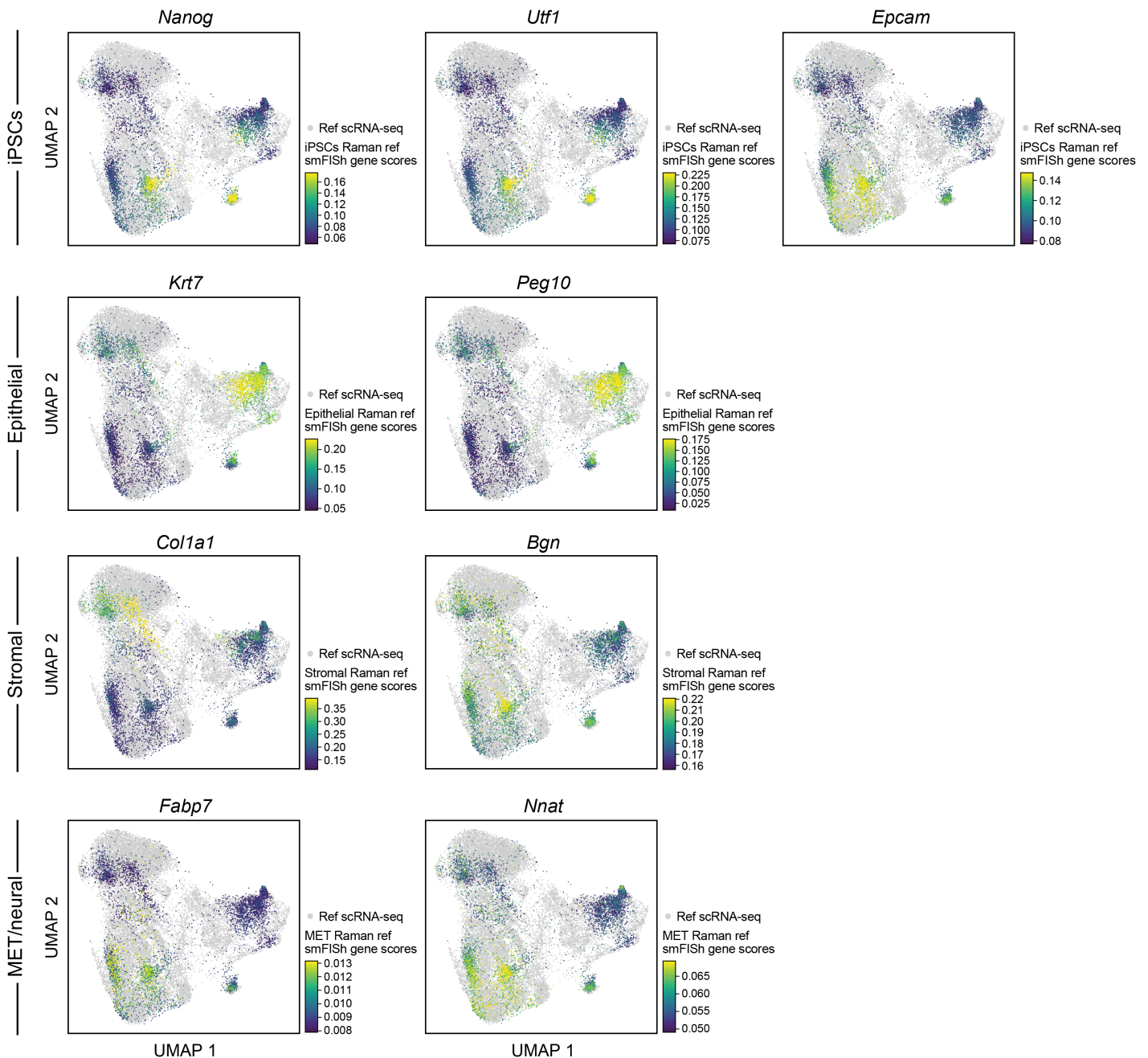


679 **Extended Data Fig. 8 | Measured and Raman-predicted single cell profiles co-embed well as**
680 **reflected by gene scores for each cell type. UMAP co-embedding of Raman predicted RNA profiles and**
681 **measured scRNA-seq profiles (dots) colored by scores of marker gene for different cell types (rows)**

682 determined by smFISH measurements (left, for cells with Raman-predicted profiles) or real scRNA-seq
683 measurements (right, for cells with scRNA-seq profiles).

684
685
686

Supp. Fig. 9

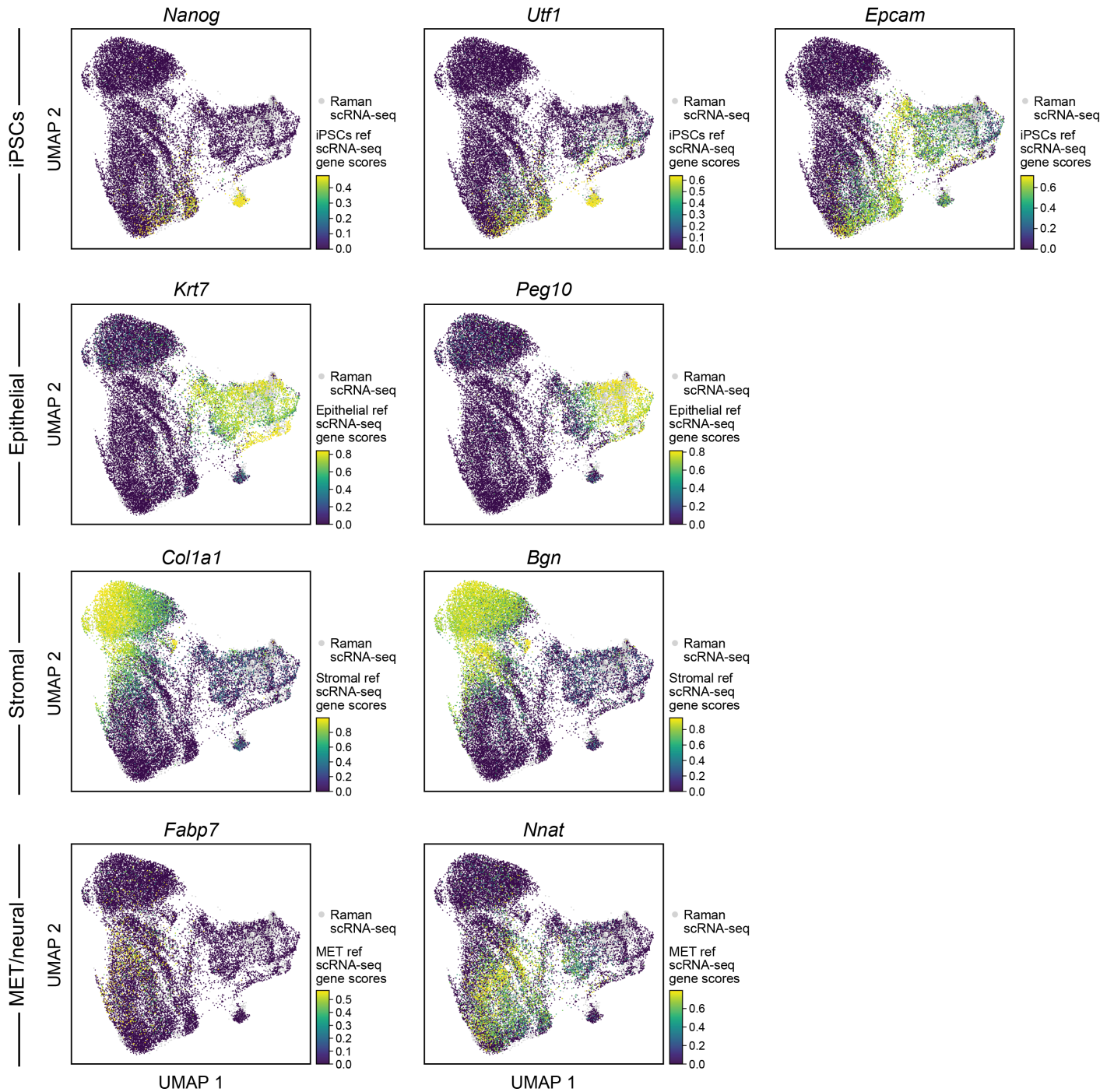


687 **Extended Data Fig. 9 | Measured and Raman-predicted single cell profiles co-embed well as**
688 **reflected by smFISH measurement of Raman cells.** UMAP co-embedding of Raman predicted RNA
689 profiles and measured scRNA-seq profiles (dots) where the Raman cells are colored by smFISH
690 measurement of each of nine anchor genes.

691

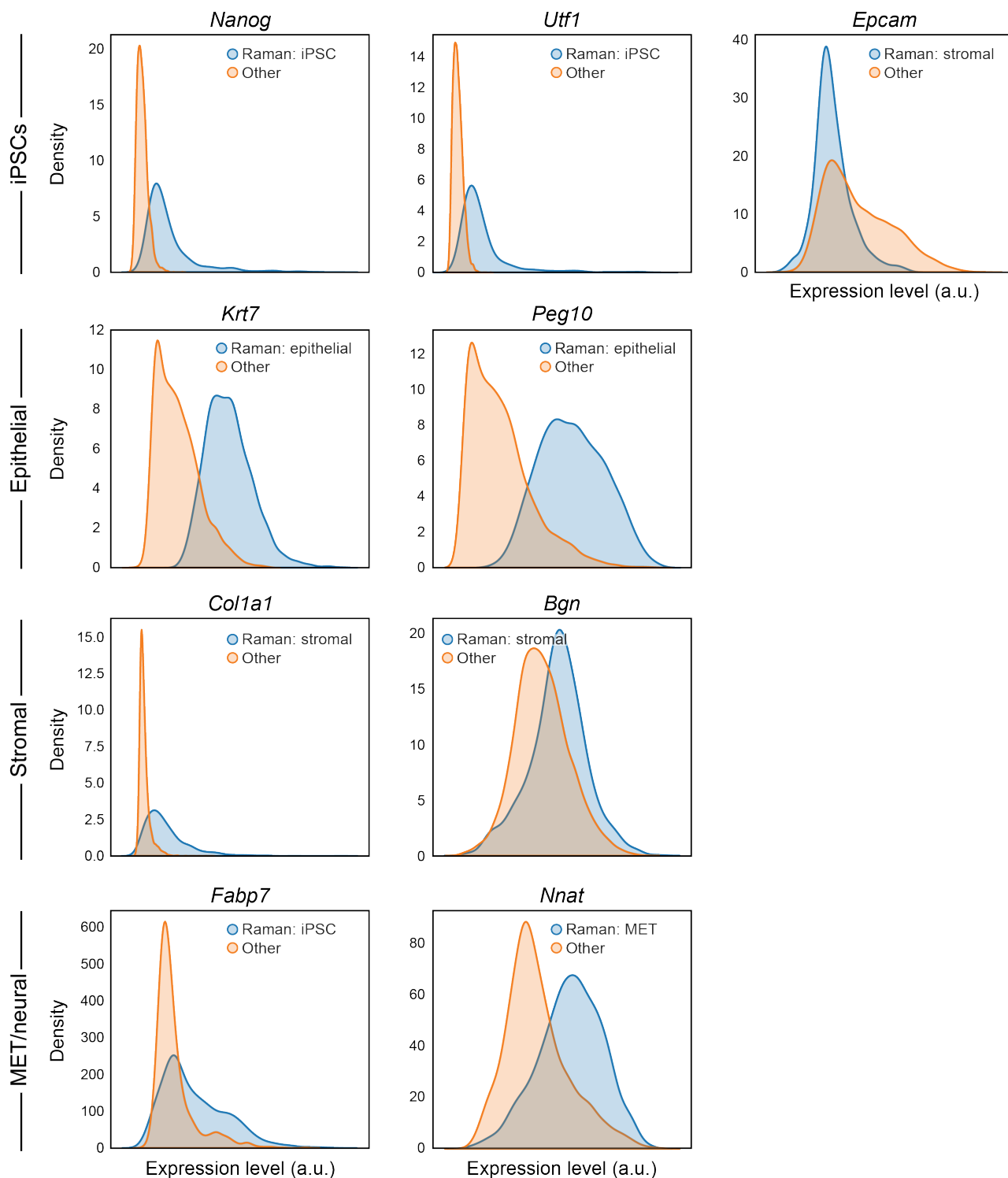
692

Supp. Fig. 10



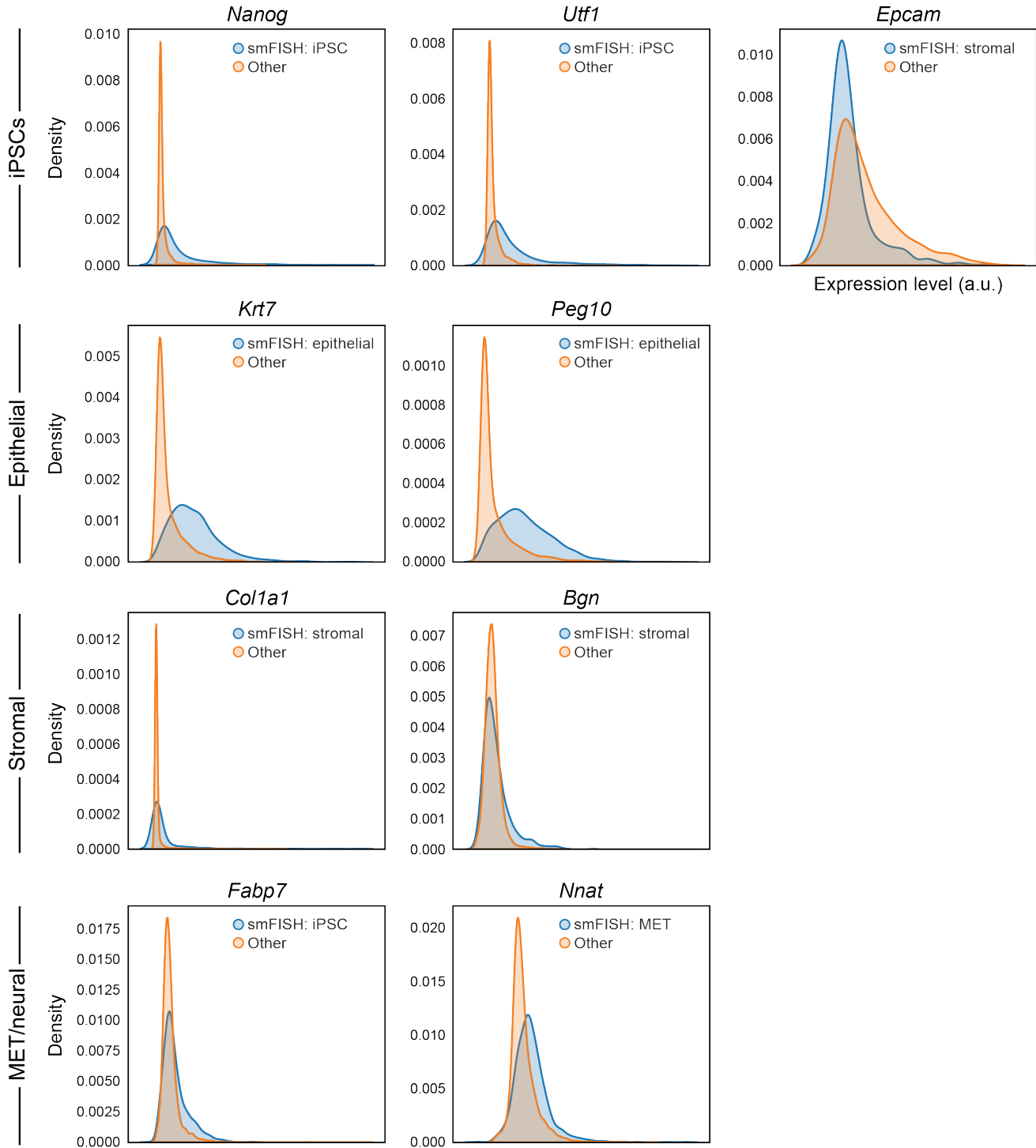
693 **Extended Data Fig. 10 | Measured and Raman-predicted single cell profiles co-embed well as**
694 **reflected by scRNA-seq based expression of nine anchor genes.** UMAP co-embedding of Raman
695 predicted RNA profiles and measured scRNA-Seq profiles (dots) where the scRNA-seq profiled cells are
696 colored by scRNA-seq measured expression of each of nine anchor genes.

Supp. Fig. 11



697 **Extended Data Fig. 11 | Distributions of expression of marker genes based on R2R-predicted**
 698 **profiles.** Distributions (density plots) of the predicted expression in Raman2RNA inferred profiles for
 699 each marker gene (panel) in its expected corresponding cell type (blue, based on the predicted expression
 700 profiles) and all other cells (orange).

Supp. Fig. 12



701

702 **Extended Data Fig. 12 | Distributions of expression of marker genes based on real smFISH profiles.**

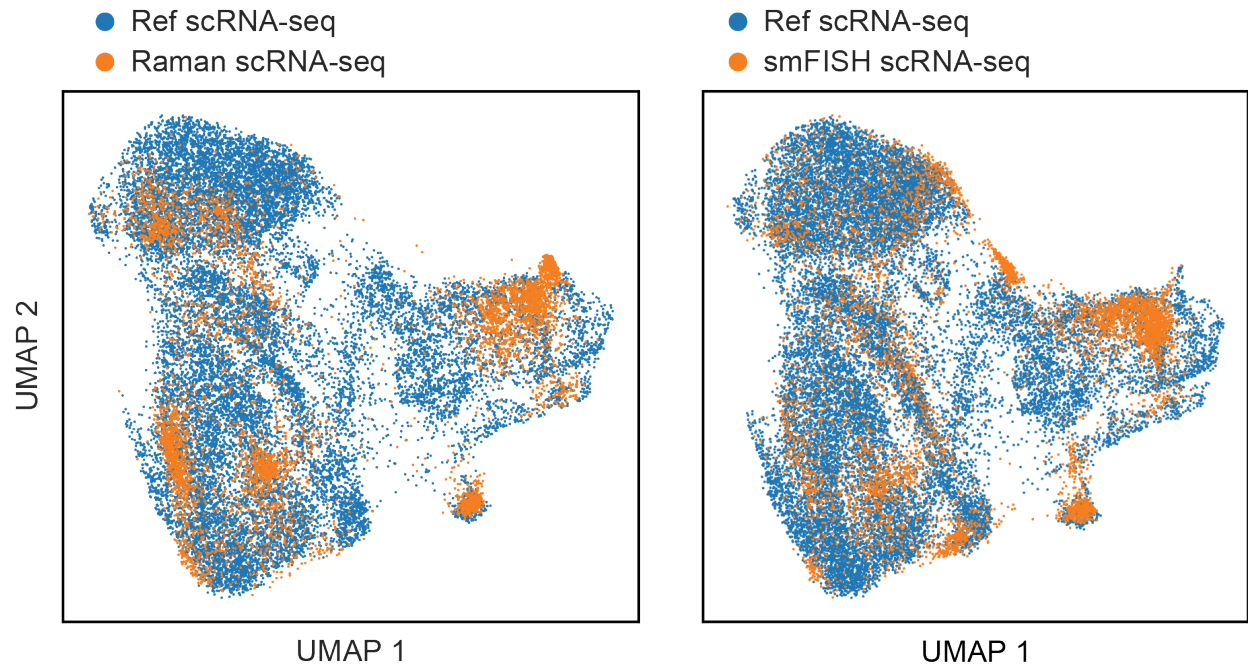
703 Distributions (density plots) of the real smFISH profiles for each marker gene (panel) in its expected

704 corresponding cell type (blue, based on the R2R *predicted* expression profiles) and all other cells

705 (orange).

706

Supp. Fig. 13



707

708

709 **Extended Data Fig. 13 | RNA profiles predicted directly from 9 anchor smFISH measurements lead**

710 **to reduced variance compared to scRNA-seq.** UMAP co-embedding of cells from scRNA-seq (blue)

711 and Raman (orange) experiments, with the latter based on either the Raman-predicted RNA profiles (left)

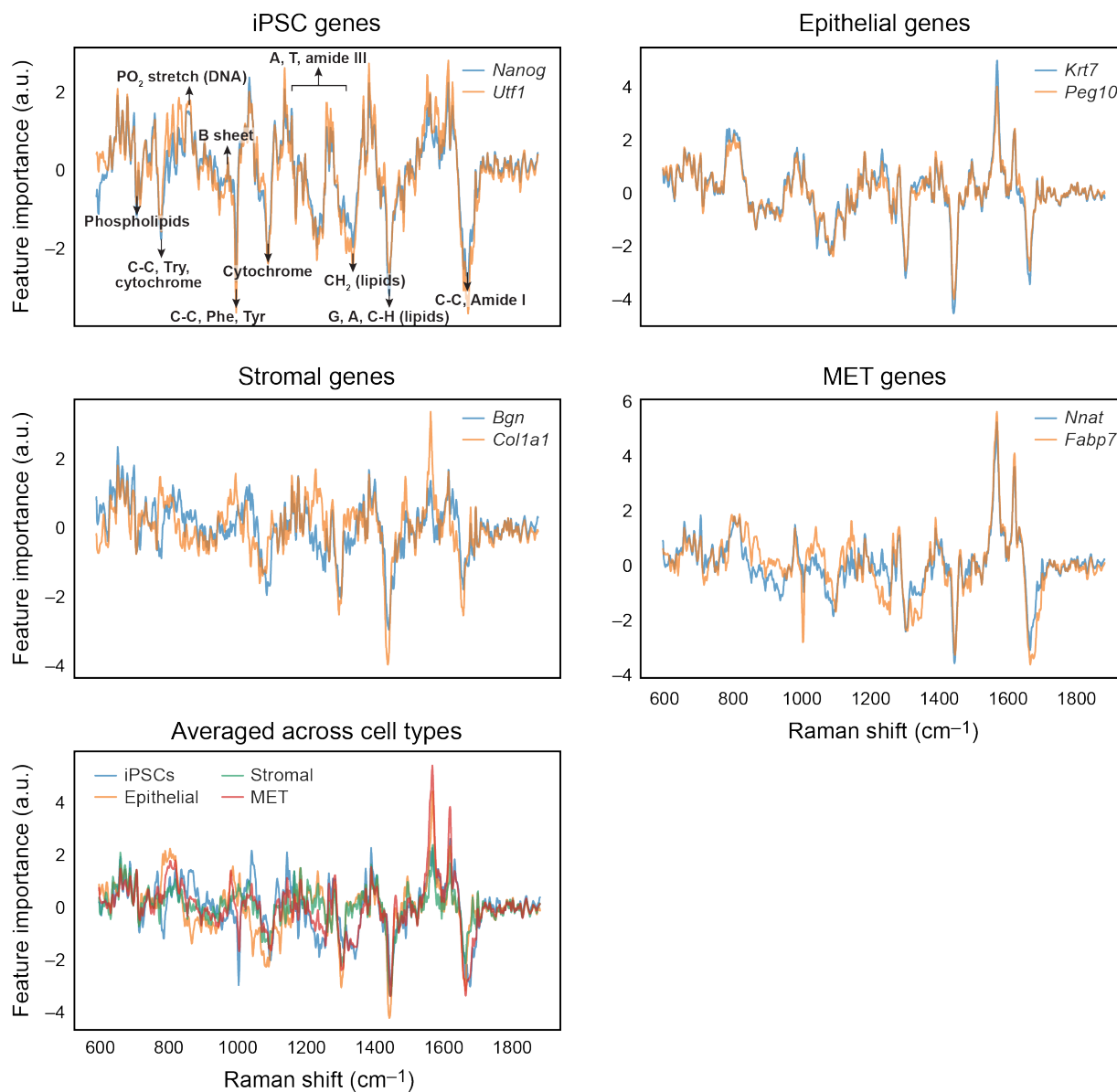
712 or only smFISH-predicted RNA profiles (right).

713

714

715

Supp. Fig. 14



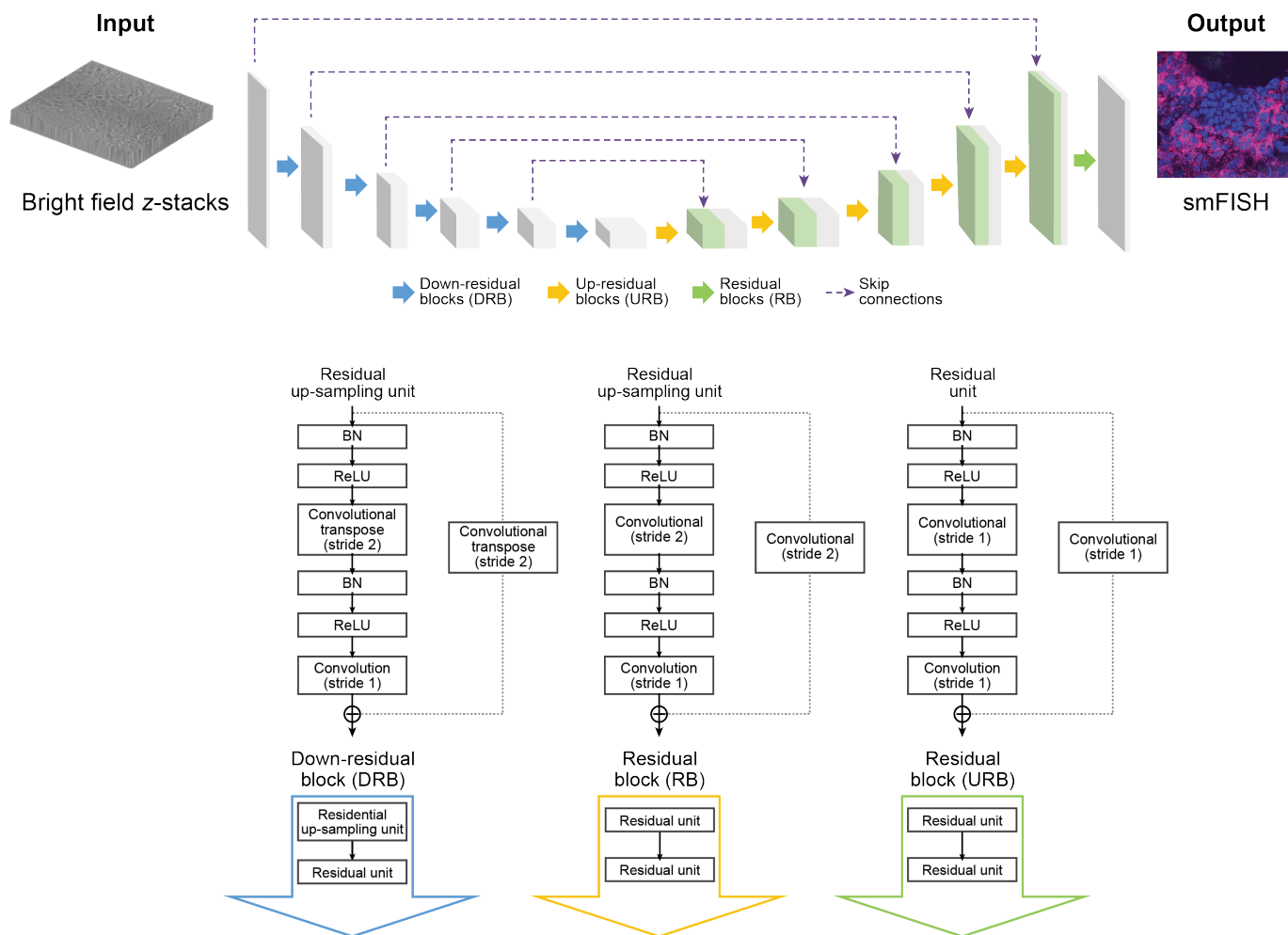
716

717 **Extended Data Fig. 14 | Raman spectral feature importance scores for each smFISH anchor gene**
718 **and its average across all genes for a cell type.** Feature importance scores (y axis) for marker genes of
719 each cell type (top two rows), and for all cell types (bottom row), along the Raman spectrum (x axis).
720 Known signals² are annotated in the top left panel (identical to **Fig. 3k**).

721

722

Supp. Fig. 15



723

724 **Extended Data Fig. 15 | Neural network-based prediction of smFISH using brightfield z-stacks.**

725

726 **References**

727

728 1. Biancalani, T. *et al.* Deep learning and alignment of spatially-resolved whole transcriptomes of

729 single cells in the mouse brain with Tangram. *bioRxiv* 2020.08.29.272831 (2020)

730 doi:10.1101/2020.08.29.272831.

731 2. Germond, A., Panina, Y., Shiga, M., Niioka, H. & Watanabe, T. M. Following Embryonic Stem

732 Cells, Their Differentiated Progeny, and Cell-State Changes During iPS Reprogramming by Raman

733 Spectroscopy. *Anal. Chem.* **92**, 14915–14923 (2020).

734