1 **TITLE**

2 Jointly leveraging spatial transcriptomics and deep learning models for pathology image

3 annotation improves cell type identification over either approach alone.

4

5 **AUTHORS AND AFFILIATIONS**

6 Asif Zubair[1], Richard H. Chapple[1], Sivaraman Natarajan[1], William C. Wright[1], Min Pan[1],

7 Hyeong-Min Lee[1], Heather Tillman[2], John Easton[1], Paul Geeleher[1, #].

8

9 [1] Department of Computational Biology, St. Jude Children's Research Hospital, Memphis, TN

10 38105, USA

11 [2] Department of Pathology, St. Jude Children's Research Hospital, Memphis, TN 38105, USA

12 [#] Correspondence: paul.geeleher@stjude.org

13

14 **ABSTRACT** (Approx. 150 words)

15 The disorganization of cell types within tissues underlies many human diseases and has been

16 studied for over a century using the conventional tools of pathology, including tissue-marking

17 dyes such as the H&E stain. Recently, spatial transcriptomics technologies were developed that

18 can measure spatially resolved gene expression directly in pathology-stained tissues sections,

19 revealing cell types and their dysfunction in unprecedented detail. In parallel, artificial

20 intelligence (AI) has approached pathologist-level performance in computationally annotating

21 H&E images of tissue sections. However, spatial transcriptomics technologies are limited in their

22 ability to separate transcriptionally similar cell types and AI-based pathology has performed less

23 impressively outside their training datasets. Here, we describe a methodology that can

24 computationally integrate AI-annotated pathology images with spatial transcriptomics data to

25 markedly improve inferences of tissue cell type composition made over either class of data

26 alone. We show that this methodology can identify regions of clinically relevant tumor immune

1

27    cell infiltration, which is predictive of response to immunotherapy and was missed by an initial

28    pathologist's manual annotation. Thus, combining spatial transcriptomics and AI-based image

29    annotation has the potential to exceed pathologist-level performance in clinical diagnostic

30    applications and to improve the many applications of spatial transcriptomics that rely on

31    accurate cell type annotations.

32

33    **BACKGROUND**

34    The traditional tools of histopathology, such as tissue-marking dyes like the hematoxylin and

35    eosin (H&E) stain, remain the primary tool used to study the disorganization and dysfunction of

36    cells within diseased tissue, representing a principal diagnostic tool in medicine. Although these

37    tools are very widely applied, they are limited by their reliance on cell morphology[1]. In the last

38    five years, sequencing-based spatial transcriptomics technologies[2-6] have emerged as a

39    powerful tool to measure spatially resolved genome-wide gene expression directly within

40    pathology-stained tissue sections, offering the potential to interrogate diseased tissue biology in

41    unprecedented detail[7,8]. Novel computational methods have already begun to address several

42    analytical challenges posed by these new data, with specific tools developed to identify spatially

43    varying genes[9,10], spatial gene expression patterns[11,12], and cell-cell interactions[13,14]. However,

44    the most fundamental problem posed by spatial transcriptomics data—upon which almost all

45    other applications of the data depend—is that of identifying the location and abundance of

46    different cell types (herein referred to as "cell type decomposition"). Several methods have

47    already been developed for this task and all function by leveraging the expression of a set of cell

48    type-specific marker genes to infer the abundance of each cell type at each slide region[15-18].

49    Notably, while mRNA is typically captured from one side of a permeabilized tissue section,

50    sequencing-based spatial transcriptomics technologies also allow images of the opposite side of

51    the tissue section to be obtained (e.g. H&E or immunohistochemical stains). Recent advances

52    in artificial intelligence, specifically deep convolutional neural networks, have profoundly

53    improved our ability to computationally extract meaningful information from these types of

54    stained tissue images[19]. For example, it was recently shown that deep learning algorithms

55    applied to H&E-stained pathology slides from The Cancer Genome Atlas (TCGA) could identify

56    diagnostically informative features of tumors, including clinically relevant estimates of cell-type

57    composition, chromosomal ploidy and signaling pathway activity[20,21].

58    However, deep learning models and spatial transcriptomics platforms each have limitations and

59    neither technology alone has displaced conventional pathology techniques. For example,

60    methods for cell type decomposition in spatial transcriptomics data will always struggle to

61    differentiate between cell types that are transcriptionally similar due to statistical

62    multicollinearity[22] and deep learning-based models for pathology have often failed to

63    recapitulate their expected performance when deployed on out-of-test-set data in real-world

64    settings[23,24].

65    Here, we present a conceptually novel computational methodology termed Guiding-Image

66    Spatial Transcriptomics (GIST). This method improves cell type decomposition in spatial

67    transcriptomics data by jointly leveraging gene expression data obtained from the spatial

68    transcriptomics platform with image-derived information from the same tissue section, for

69    example, the output of deep learning models applied to images of histopathology stains. We

70    applied this computational tool to integrate spatial transcriptomics data with deep learning-

71    derived cell type annotations in breast cancer pathology slides where we identified

72    prognostically relevant immune cell infiltration that was missed by an initial pathologist's manual

73    annotation. The methodology presented is generalizable to any sequencing-based spatial

74    transcriptomics platform where informative image-derived cell-type compositional estimates can

75    be obtained. Thus, combining spatial transcriptomics and paired pathology images has potential
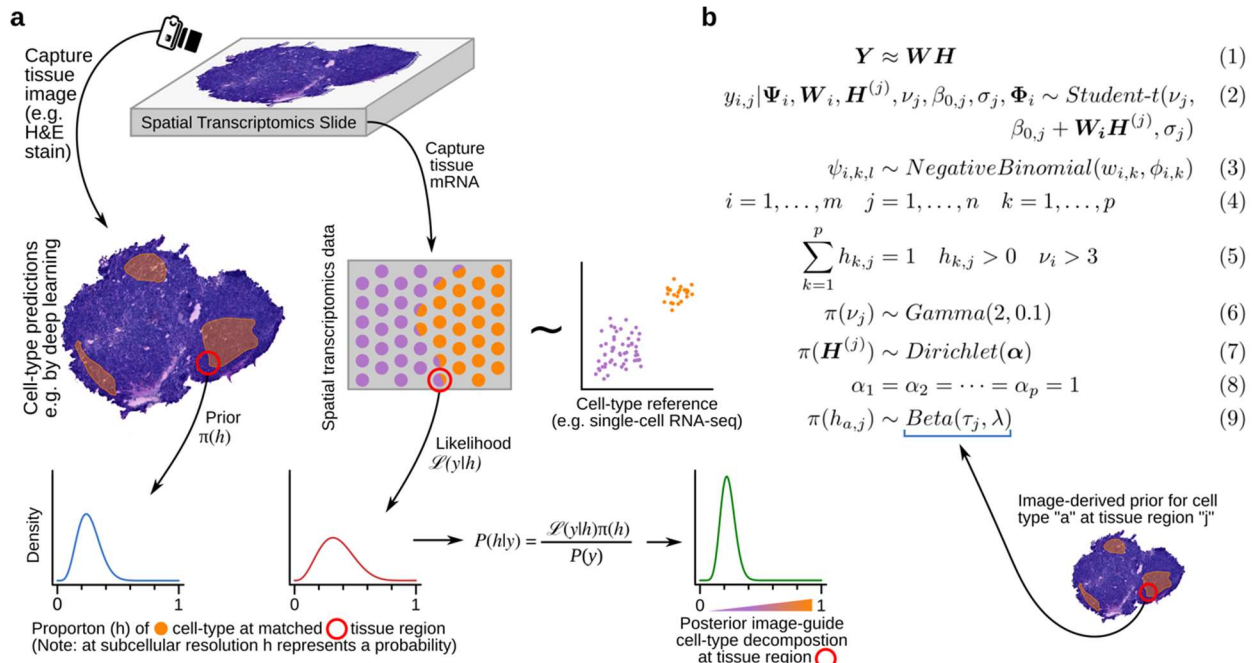
3

76    applications in clinical diagnostics and can also improve all analytical applications of spatial

77    transcriptomics data that rely on the correct annotation of cell types.

78  **RESULTS**

79  **Guiding-Image Spatial Transcriptomics (GIST) jointly leverages spatial transcriptomics**

80  **and paired tissue images to improve cell type decomposition.**

81  GIST attempts to improve cell type decomposition in spatial transcriptomics data by leveraging

82  prior estimates of cell type composition from paired pathology images. The method relies on

83  Bayesian probabilistic modeling, a statistical approach that naturally lends itself to integrating

84  multiple sources of information, jointly leveraging spatial transcriptomics and imaging

85  information to improve cell type decomposition estimates over either approach applied alone.

86  Intuitively, the approach uses the imaging data to provide an initial "suggestion" as to the cell

87  types in a particular region of the spatial transcriptomics slide, but this suggestion can be

88  overcome if outweighed by the evidence from the transcriptomic data (schematic representation

89  in Fig. 1a, model formation in Fig. 1b; see Methods for further technical details).

90



$$Y \approx WH \tag{1}$$

$$y_{i,j} | \Psi_i, W_i, H^{(j)}, \nu_j, \beta_{0,j}, \sigma_j, \Phi_i \sim Student\text{-}t(\nu_j, \tag{2}$$
$$\beta_{0,j} + W_i H^{(j)}, \sigma_j)$$

$$\psi_{i,k,l} \sim NegativeBinomial(w_{i,k}, \phi_{i,k}) \tag{3}$$

$$i = 1, \ldots, m \quad j = 1, \ldots, n \quad k = 1, \ldots, p \tag{4}$$

$$\sum_{k=1}^{p} h_{k,j} = 1 \quad h_{k,j} > 0 \quad \nu_i > 3 \tag{5}$$

$$\pi(\nu_j) \sim Gamma(2, 0.1) \tag{6}$$

$$\pi(H^{(j)}) \sim Dirichlet(\boldsymbol{\alpha}) \tag{7}$$

$$\alpha_1 = \alpha_2 = \cdots = \alpha_p = 1 \tag{8}$$

$$\pi(h_{a,j}) \sim Beta(\tau_j, \lambda) \tag{9}$$

92  **Figure 1: Overview of Guiding-Image Spatial Transcriptomics (GIST) methodology.**
93      a)  Schematic representation of GIST. The schematic shows a hypothetical tissue section, where we
94          wish to identify the location of a hypothetical cell type (colored orange); this could represent, for

95  example, immune cell infiltration in a tumor. Estimates of this cell type's proportions from a deep

96  learning model applied to an H&E stain image (left) are used to optimize the estimates derived

97  from the spatial transcriptomics data (right), yielding improved estimates over what could be

98  achieved from either approach alone (bottom right).

99  b)  Mathematical notation describing the GIST model: We assume that the spatial transcriptomics

100  data $Y_{m \times n}$ can be approximately factorized as a cell type signature matrix $W_{m \times p}$ and a matrix of

101  cell type compositional estimates $H_{p \times n}$ (eqn. (1)). We propose estimating cell type composition $H$

102  using the model in eqn. (2-9). A single-cell RNA-seq dataset from the same tissue type is

103  represented by $\Psi$. Each element of $W$ is estimated from $\Psi$ using a negative binomial distribution

104  (with overdispersion parameter $\phi_{i,k}$) estimated for each gene $i$, in each cell type $k$, from the

105  expression in each single-cell $l$. Eqn. (5) shows the model constraints. Eqn. (6-9) show the priors,

106  denoted by $\pi$. Other parameters are assigned weakly informative priors. The key informative prior

107  is shown in eqn. (9), where the image-derived prior estimate of cell type composition for a cell

108  type of interest, contained in row $a$ of $H$, is specified as a beta distribution. For each tissue region

109  (e.g. unique barcoded spot), this beta distribution is parameterized by its mean, $\tau_j$, specifying the

110  prior cell type composition estimate from the image, and the hyperparameter $\lambda$, a scalar that

111  determines how much weight to place on the image data and how much to place on the

112  transcriptomic data. Notes: Superscript notation (e.g. $H^{(j)}$) denotes the columns of a matrix.

113  Vectors are shown using boldface and matrices bold capital letters. All equations herein assume

114  $m$ genes (indexed by $i$), $n$ tissue regions (e.g. slide mRNA capture spots, indexed by $j$), $p$ cell
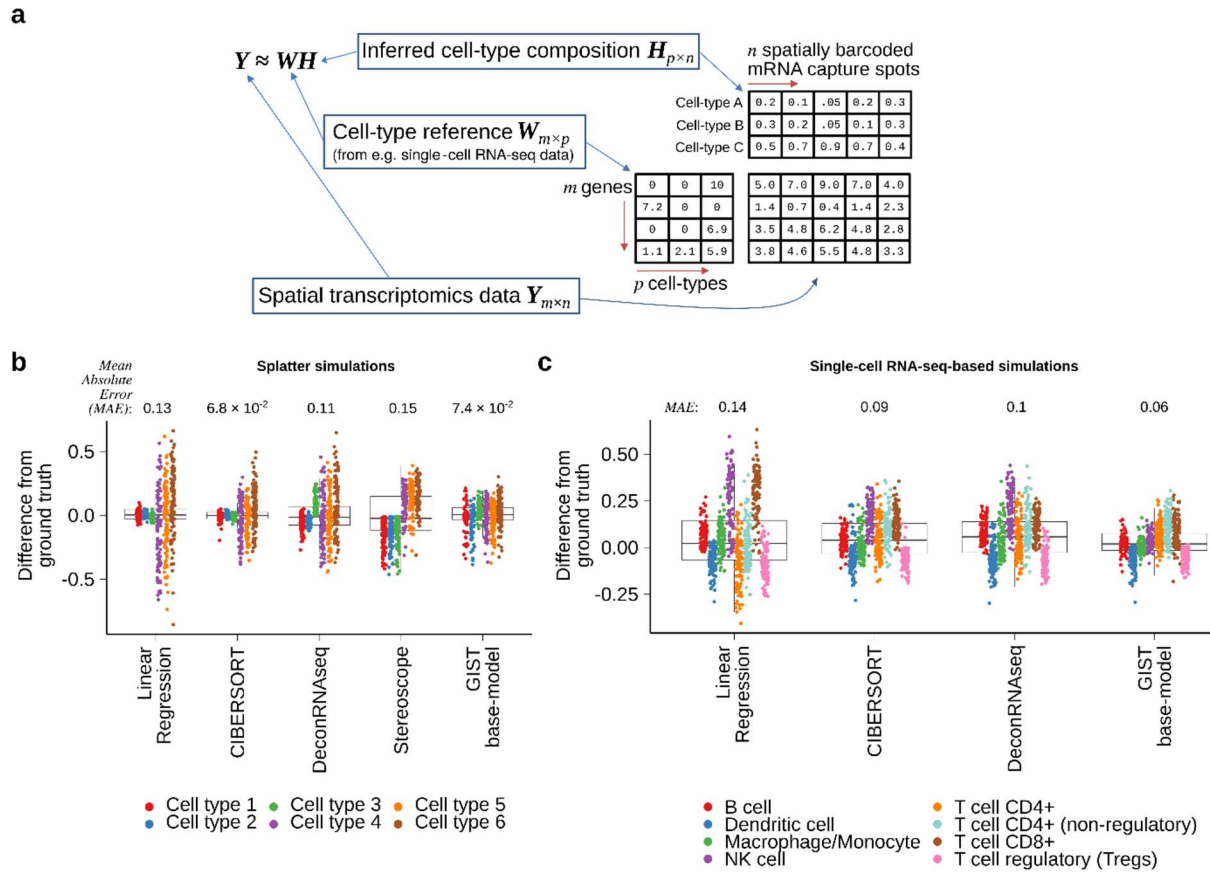
115  types (indexed by $k$).

116

117  **A Bayesian probabilistic model for cell type decomposition performs competitively when**

118  **compared to existing methods in simulations when no paired image information is**

119  **leveraged.**

120  Existing methods for cell type decomposition in spatial transcriptomics data are related to

121  previous models for bulk gene expression deconvolution and can be broadly conceptualized as

122  a matrix decomposition, where some reference basis matrix of expression data from purified

123  cells $W$ (e.g. derived from single-cell RNA-seq) is used to estimate the proportion of each cell

124  type $H$ in the bulk mixture $Y$ (Fig. 2a for schematic representation). At subcellular resolution, the

125  $H$ matrix can be thought of as probability estimates, rather than proportion estimates[16], although

126  for simplicity we use the term "proportion" throughout this manuscript.

127  The statistical model underlying GIST is related to these existing approaches but includes the

128  ability to leverage prior information derived from paired tissue images. Thus, we were first

6

129    interested in assessing whether our model performed competitively when compared to existing

130    approaches in the absence of prior information derived from images (*henceforth referred to as*

131    *the "GIST base-model"*). To test this, we first developed two complementary unbiased

132    benchmarking simulations, one based on the existing tool Splatter[25] and one based on a

133    published benchmarking dataset[26], which evaluates methods on a simulated mixture of immune

134    cell types from a real single-cell RNA-seq dataset. We compared the GIST base-model to two

135    methods originally designed for bulk gene expression data (CIBERSORT[27], DeconRNASeq[28]), a

136    method tailored specifically for spatial transcriptomics data (Stereoscope[18]), and linear

137    regression (the simplest conceivable model.) Based on the mean absolute error (*MAE*),

138    CIBERSORT performed slightly better on the Splatter simulations (Fig. 2b, Supplementary

139    Figure 1, Supplementary Table S1; *MAE* = 6.8 × 10$^{-2}$ for CIBERSORT and 7.4 × 10$^{-2}$ for the

140    GIST base-model), while the GIST base-model performed best on the other benchmarking

141    dataset (Fig. 2c, Supplementary Figure 2, Supplementary Table S2; *MAE* = 0.09 for

142    CIBERSORT and 0.06 for the GIST base-model). However, given the conceptual similarity of

143    the underlying models, it is not surprising that none of these existing methods produce markedly

144    dissimilar results in either simulation, suggesting that, rather than further model tweaking and

145    optimization, a new conceptual advance may be necessary to achieve meaningful progress on

146    the cell type decomposition problem.

147

**Figure 2: A Bayesian probabilistic model performs similarly to existing cell type decomposition methods when no prior information is available.**

a) Schematic representation of the cell type decomposition problem posed as a matrix decomposition. Spatial transcriptomics expression data is arranged in an $m$ genes by $n$ mRNA-capture-spots matrix $Y$. This matrix is decomposed into a basis matrix $W$ and a matrix $H$ that contains the proportion of each of $p$ cell types on each spot or (at subcellular resolution) the probability that a spot matches a cell type (shown for three hypothetical cell types A, B, and C). The basis matrix $W$ is typically known and can be derived for example from single-cell RNA-seq data from the same or similar tissue. Given this, all existing cell type decomposition algorithms, be they designed specifically for spatial transcriptomics data or not, aim to estimate $H$.

b) Boxplot showing the results of five cell type decomposition methods on simulated mixture gene expression data, for a mixture of 6 cell types, generated using the tool Splatter (see Methods). Points have been colored by the simulated cell type and the y-axis shows the deviation from ground truth, quantified by the difference between the estimated cell type proportions in a sample and the true proportion used as ground truth for the simulation. The Mean Absolute Error ($MAE$), summarizing the overall performance of each method is as follows (lower values imply better performance): Linear regression = 0.13, CIBERSORT = $6.8 \times 10^{-2}$, DeconRNAseq = 0.11, Stereoscope = 0.15, GIST base-model = $7.4 \times 10^{-2}$.

c) Similar to (b) but based on the simulated dataset obtained from the benchmarking procedure outlined in Strum et al.[26]. Points have been colored by the immune cell type and the y-axis shows the deviation from ground truth, quantified by the difference between the estimated cell type proportions in a sample and the true proportion used as ground truth for the simulation. The Mean Absolute Error ($MAE$), summarizing the overall performance of each method is as follows (lower

8

171        values imply better performance):  Linear regression = 0.14, CIBERSORT = 0.09, DeconRNAseq

172        = 0.1, GIST base-model = $6.4 \times 10^{-2}$. Note Stereoscope is not included in this second set of

173        simulations as it was not possible to pass the CIBERSORT LM22 signature matrix, which is used

174        as the cell-type reference in this simulation, to Stereoscope (see Methods).

175        In all boxplots, the center line represents the median, bound of box is upper and lower quartiles

176        and the whiskers are 1.5× the interquartile range.

177

178    **The GIST base-model performs competitively on spatial transcriptomics data obtained**

179    **from mouse brain sections when cell type specific immunofluorescence markers are**

180    **treated as a ground truth.**

181    We were next interested in comparing the performance of the GIST base-model to other

182    methods using real spatial transcriptomics data. To do this, we leveraged a publicly available

183    dataset (see Data Availability), which measured gene expression in the mouse brain using the

184    10x Genomics Visium spatial transcriptomics platform, and where immunofluorescence (IF)

185    staining was performed on the reverse side of the tissue section. These IF stains were

186    conducted for two proteins, RBFOX3 and GFAP, which are protein markers unique to neurons

187    and glia respectively (Fig. 3a). We calculated the average pixel intensity of each of these two

188    markers in all image pixels overlapping each spatially barcoded mRNA capture spot on the

189    Visium slide (Fig. 3b; see Methods), then we used these spot-level intensity estimates to

190    represent an independent ground-truth approximating the abundance of neurons and glia in

191    regions of the slide overlapping each of the Visium array's 4,992 spots.

192    Next, using the GIST base-model, we estimated the cell type composition on each spot from the

193    spatial transcriptomics data by leveraging a single-cell RNA-seq dataset that was available from

194    a similar region of a mouse brain, allowing us to estimate the abundance of glial and neuronal

195    cell types from the spatial transcriptomics expression data alone (Fig. 3c). We compared the

196    results obtained from the GIST base-model to popular spatial transcriptomics cell type

197    decomposition methods Spotlight[15], RCTD[16], Stereoscope[18], and Cell2location[29], treating the IF-

198    derived estimates of neurons and glia at each spot as ground truth. Consistent with our

199    simulations, the GIST base-model, RCTD, Cell2location, and Spotlight all performed quite

200    similarly in these benchmarks on real data; however, we note that the GIST base-model had

201    slightly better performance than the other methods, achieving Spearman's rank correlations of

202    0.49 and 0.77, compared to 0.33 and 0.77 for RCTD (the second best performing method), for

203    the glial and neuronal comparisons respectively (Fig. 3d; $P < 2.2 \times 10^{-16}$ from Spearman's

204    correlation against IF-derived ground truth for all five methods; Supplementary Figures 3-7).

205    Overall, these results suggest that the GIST base-model performs competitively when

206    compared to existing methods for cell type decomposition in real spatial transcriptomics data.
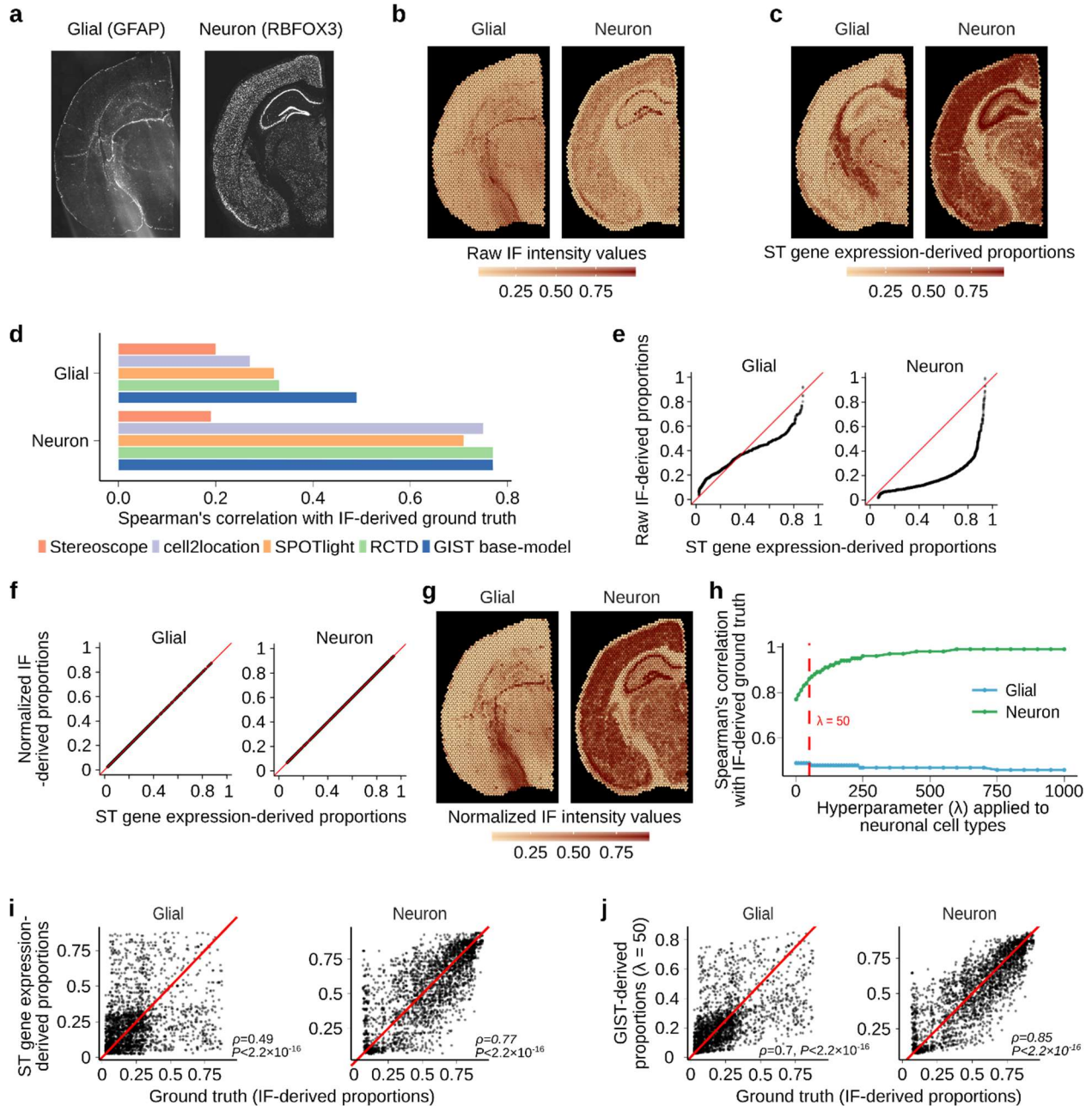
207

208    **Incorporating image-derived prior information from matched immunofluorescence stains**

209    **improves cell type decomposition in spatial transcriptomics data generated from a**

210    **mouse brain section.**

211    Even though our GIST base-model performed well compared to existing methods, the results

212    above also showed that the best-performing methods were not markedly different and fall well

213    short of an optimal performance when compared to the IF-derived ground truth. Thus, we next

214    hypothesized that it should be possible to markedly improve our performance by leveraging our

215    model's Bayesian implementation and supplying the model with informative image-derived prior

216    information (*henceforth referred to as the "GIST model"*). We reasoned that we could first

217    demonstrate this principle on this mouse brain dataset, leveraging the IF-derived estimates of

218    cell type abundance. However, IF-derived pixel intensity estimates do not represent proportions

219    on a 0-1 scale and thus it is not obvious how this information could be leveraged as prior

220    estimates of cell type composition in the GIST model. To solve this problem, we first normalized

221    the IF-derived estimates by mapping them onto the quantiles of the spatial transcriptomics-

222    derived cell type proportion estimates, generated by an initial round of model fitting using the

223    GIST base-model (Fig. 3e-g; see Methods). We then refit our GIST model, incorporating this

224    prior knowledge derived from the RBFOX3 IF data, providing "suggestions" of the abundance of

225    neuronal cell types over each spatial transcriptomics spot. We specified these priors using a

226    beta distribution applied to the appropriate group of model parameters corresponding to

227    neuronal cell type estimates. The beta distribution was parameterized by its mean ($\tau$; the point

228    estimate of the normalized cell type proportion estimate from the IF image) and the total-count

229    parameter ($\lambda$; the strength of the prior, corresponding to the weight placed on the IF image)—

230    any beta distribution is naturally constrained to a 0-1 scale, meaning it is appropriate for

231    specifying image-derived prior estimates of cell type composition. The key modeling question is

232    then determining how much weight to place on these image-derived priors and how much to

233    place on the spatial transcriptomics data itself. This must be determined by tuning the

234    hyperparameter $\lambda$, where selecting a value that is too small will mean there is little to no

235    influence of the image-derived cell type information on the model's output but selecting a value

236    that is too large will overfit the model to the image and degrade performance on unseen test

237    data.

238    We chose this hyperparameter $\lambda$ by observing how the estimates of glial cell type composition

239    compared to IF-derived glial-cell ground-truth (GFAP stain) when fitting the model with ever-

240    increasing values of $\lambda$ for the IF-derived neuronal cell type prior (RBFOX3 stain), only placing

241    priors on the neuronal cell types. As expected, when increasing the value of $\lambda$ and placing more

242    weight on the image-derived prior for neuronal cells, the model's output progressively more

243    closely matched these IF-derived estimates for the neuronal cell types (Fig. 3h). However, as

244    we continued to increase $\lambda$, placing more and more weight on the image-derived estimates of

245    neuronal cells, we eventually observed a precipitous drop-off in the model's performance, as

246    measured by the agreement between the glial cell type estimates from the GIST model and the

247    IF-derived ground truth from the GFAP glial marker protein (Fig. 3h). This drop-off begins at $\lambda =$

248     $50$, suggesting that beyond this point the model has been overfitted, providing us a reasonable

249     value of $\lambda$ for image-derived priors. This value of $\lambda$ concentrates most of the prior probability

250     mass within approximately ±10% of the mean. Notably, at this $\lambda$ value, the Spearman's rank

251     correlation between the model-derived neuronal cell type estimates and the IF-derived ground

252     truth increased from 0.7 to 0.85, substantially better than any method that does not leverage the

253     images and approaching an optimal performance (Figs. 3i and 3j). Overall, these results support

254     the notion that applying informative prior information, derived from matched images has the

255     potential to improve the performance of cell type decomposition in spatial transcriptomics data

256     and provides a reasonable initial estimate of the key hyperparameter $\lambda$ to now be applied to out-

257     of-batch test datasets.

**Figure 3: Incorporating image-derived prior information from matched immunofluorescence stains in mouse brain spatial transcriptomics data.**

a) Raw immunofluorescence image of the mouse brain tissue section showing the glial (GFAP) and neuronal (RBFOX3) cell markers.

b) Spatial distribution of raw IF intensity values for GFAP (glial) and RBFOX3 (neuronal) when fluorescence intensity has been averaged over pixels corresponding to each spatial transcriptomics spot's location. Intensity values were rescaled from 0 to 1.

c) Spatial distribution of glial and neuronal proportions estimated from the spatial transcriptomics gene expression data using the GIST base-model.

d) Bar plot showing Spearman's correlation between IF-derived ground truth cell type proportions and cell type proportions estimated from five different gene expression-based spatial

13

270    transcriptomics cell type decomposition methods (Stereoscope, cell2location, SPOTlight, RCTD,
271    and the GIST base-model).

272  e)  Quantile-quantile plot (QQ plot) of image-based IF-derived values for total glial and neuronal
273    content for each spot (y-axis) versus values obtained for total glial and neuronal content from the
274    spatial transcriptomics gene expression data only using the GIST base-model (x-axis).

275  f)  Same as in (d) except that this QQ plot is generated after post-mapping normalization where the
276    distribution of cell type compositional estimates from the IF images were mapped onto the
277    distribution of cell type compositional estimates from the spatial transcriptomics gene expression
278    data generated using the GIST base-model.

279  g)  Spatial distribution of IF intensity values for the glial and neuronal channel where the values have
280    now been mapped to a distribution estimated from the gene expression data using the GIST base-
281    model.

282  h)  Line plot showing the change in GIST model performance as we increase the key hyperparameter
283    $\lambda$ (x-axis). Performance is quantified by Spearman correlation with IF-derived ground truth (y-
284    axis) and is shown for both neuronal (green) and glial (blue) cell types. The RBFOX3 IF image-
285    derived prior is only applied to the neuronal cell type. A non-informative prior is applied to the
286    glial cell type. The vertical dashed red line indicates a stopping point ($\lambda = 50$) where performance
287    in the glial channel begins to deteriorate, indicating the model has been overfitted to the RBFOX3
288    IF data.

289  i)  Scatter plots showing the cell type compositional estimates against IF-derived ground truth (x-
290    axis) in the mouse brain for glia (left) and neurons (right) derived from the spatial transcriptomics
291    gene expression data using the GIST base-model (y-axis) when no prior information is leveraged.
292    P-values from Spearman's correlation test.

293  j)  Similar to (i) but showing the improved agreement with ground truth (x-axis) when the IF-derived
294    cell type compositional estimates are incorporated as prior information using the GIST model
295    with a $\lambda$ hyperparameter value of 50 (y-axis). P-values from Spearman's correlation test.

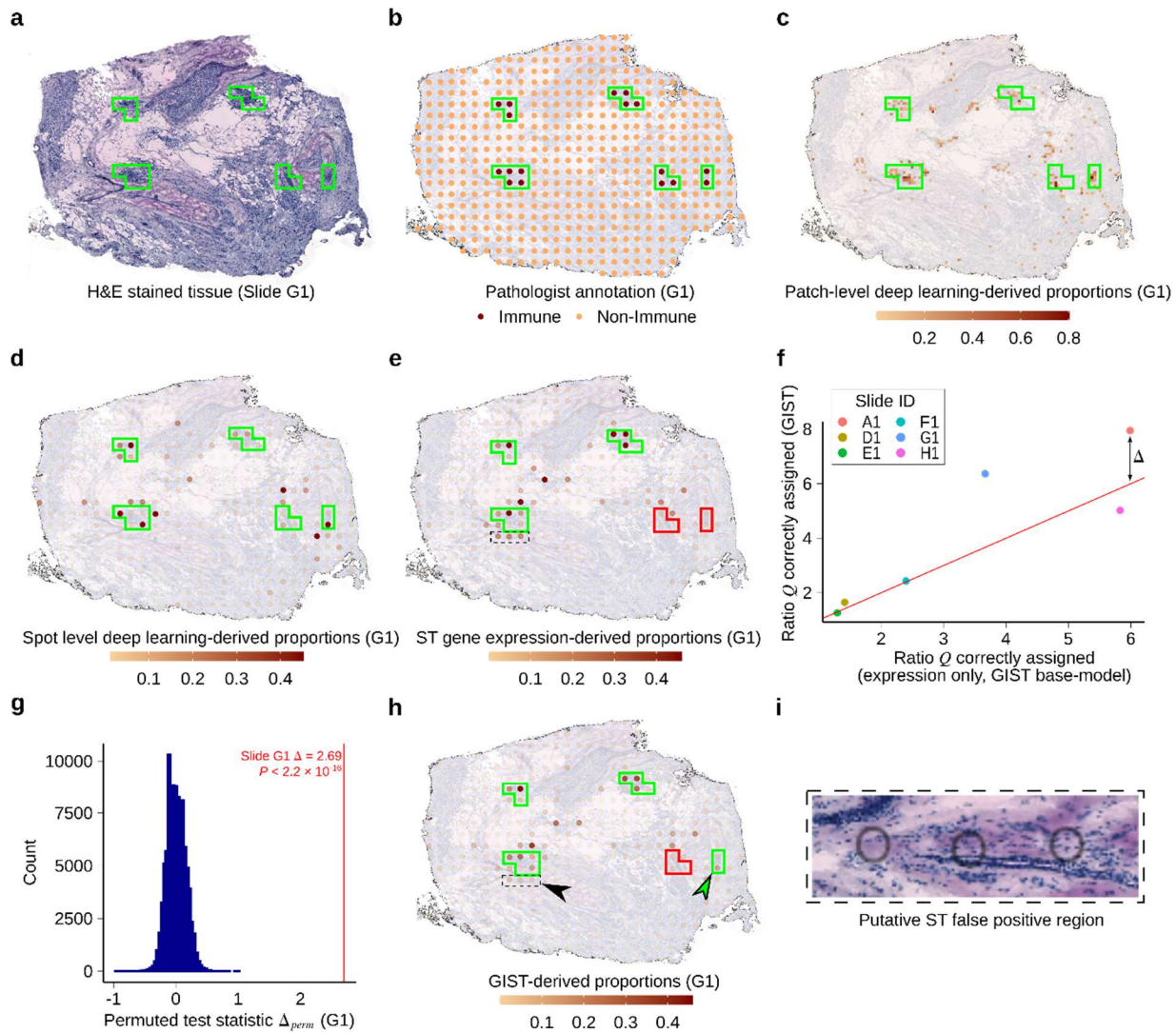296  Abbreviations: ST, Spatial Transcriptomics; IF: Immunofluorescence.

297

298  **Incorporating prior information derived from deep learning models applied to matched**

299  **H&E-stained images improves estimates of immune cell infiltration in breast cancer**

300  **spatial transcriptomics data.**

301  The results above provide a convincing proof-of-principle that it should be possible to improve

302  cell type decomposition in spatial transcriptomics data by leveraging matched images. However,

303  while IF stains can provide reliable markers of cell types, they are restricted to a small number

304  of proteins and are much less commonly collected than the H&E stain. Thus, we wondered

305  whether it would be possible to leverage image information derived from deep learning models

306  applied to H&E stains—the principal pathology stain that is collected as a part of almost all

14

307    sequencing-based spatial transcriptomics protocols. Deep learning models have already been

308    developed that can output numerous clinically relevant annotations from H&E-stained tissue

309    section images alone, which could theoretically be usefully propagated in the spatial

310    transcriptomics assay. These annotations include cell type composition, expression of signaling

311    pathways, chromosomal ploidy, and immune cell infiltration[20,21,30]. To test whether such

312    information could be usefully exploited in spatial transcriptomics assays, we obtained 8

313    previously published spatial transcriptomics tissue slides, which had measured gene expression

314    in biologically independent breast cancer tumors. Critically, each of these tissue sections had

315    also been H&E stained (Fig. 4a, panel (a) in Supplementary Figures S8-S12), and regions of

316    immune cell infiltration had been annotated by a previous pathologist (Fig. 4b, panel (b) in

317    Supplementary Figure S8-S12), providing an independent ground truth against which to assess

318    our model predictions. Identifying immune cell infiltration has prognostic value[31] and is predictive

319    of response to cancer immunotherapy[32], hence represents a particularly interesting use case of

320    the GIST model.

321    Thus, we applied a previously published deep convolutional neural network[30], which had been

322    trained using images collected as part of TCGA to identify regions of tumor-infiltrating

323    lymphocytes from H&E stained tumor tissue sections. This yielded patches of deep learning-

324    derived predictions of immune cell infiltration across each of our breast cancer tumor tissue

325    sections (Fig. 4c, panel (c) in Supplementary Figures S8-S12), where gene expression had also

326    been measured using spatial transcriptomics. We then averaged these deep learning derived

327    predictions over the pixels overlapping each of the spatial transcriptomics mRNA capture spots,

328    yielding a deep-learning-derived per-spot estimate of immune cell composition in each tumor

329    (Fig. 4d, panel (d) in Supplementary Figures S8-S12, similar to the approach applied above for

330    IF data; see Methods). Initial immune cell proportions at each spot were then estimated using

331    the GIST base-model (Fig. 4e, panel (e) in Supplementary Figures S8-S12). We applied a

15

332    similar normalization approach as we described for the IF data, mapping the deep learning

333    derived estimates to the quantiles of the initial gene expression derived estimates, then applied

334    these deep-learning-derived immune cell compositional estimates as informative priors, again

335    specified as a beta distribution on the appropriate GIST model parameters. We used a $\lambda$ value

336    of 50, which was derived from the previous independent mouse brain dataset, avoiding any

337    potential issues with overfitting to this new dataset (Fig. 3h). If the GIST model performs better

338    than the expression-only GIST base-model, the expectation is that we should identify more

339    immune cells in pathologist-annotated immune cell regions, but less in other regions of the

340    slides. Thus, we quantified model performance by the ratio of immune cells identified within the

341    pathologist's annotated regions of immune infiltration, compared to all other regions of the tissue

342    slide (this ratio is defined herein as $Q$ (see Methods); note that regions of immune cells had

343    been identified by the pathologist in six of eight slides). When compared to the pathologist-

344    derived ground truth, the GIST model, leveraging deep learning-derived prior information,

345    performed better than the expression-only GIST base-model in four out of the six slides (Fig. 4f,

346    panel (f) in Supplementary Figures S8-S12). The performance increase over the GIST base-

347    model was particularly large for two slides (Fig. 4g, panel (g) in Supplementary Figures S8-S12;

348    increase in $Q$ for GIST vs GIST base-model (defined herein as $\Delta$) of 1.95 and 2.69, $P = 7.2 \times 10^{-3}$

349    and $P < 2.2 \times 10^{-16}$ for slides A1 and G1 respectively; empirical $P$-values were calculated by

350    permutation, see Methods). Visual inspection of the results revealed examples of clear regions

351    where leveraging the deep learning-derived prior information correctly decreased the estimates

352    of immune cell composition in regions where the pathologist marked an absence of immune

353    cells (Fig. 4h, black arrowhead, and Fig. 4i) and regions where estimates of immune cell

354    composition increased to match the pathologist (Fig. 4h, green arrowhead). Thus, leveraging

355    deep learning derived prior information has the potential to markedly improve cell type

356    decomposition in data generated from spatial transcriptomics technologies.

**Figure 4: Tissue image-derived cell type compositional estimates can be leveraged to improve estimates of immune cell infiltration in breast cancer tissue sections profiled using spatial transcriptomics.**

a) H&E stained tissue image obtained from the reverse side of the breast cancer spatial transcriptomics slide G1. Green outline shows regions containing ST spots annotated as containing immune cells by the pathologist.

b) Pathologist annotation for slide G1 showing regions containing spatial transcriptomics spots that were labeled immune cell infiltrated (marked by dark-colored spots and green outlines).

c) Output from the deep learning model for slide G1 overlayed on top of the breast cancer tissue section H&E image. The color scale indicates deep learning-derived predictions for the proportions of immune cells made on 50×50 micron patches of the tissue. Green boxes outline regions of pathologist's annotated immune spots.

d) Slide G1 showing the patch level deep learning predictions converted to spot level predictions, so that they can be used as priors in the GIST model. Spot level predictions are a sum of patch level predictions weighted by their percent overlap with the spot. Boxes outline regions of pathologist's annotated immune spots.

17

374 e) Slide G1 showing the gene expression-derived immune cell proportions from the GIST base
375 model. Solid boxes indicate the regions of the pathologist's annotated immune spots. Green
376 indicates that the model reasonably identifies immune-infiltrated spots. Red indicates that the
377 immune spots were not captured by the model. The dashed black box indicates a region of
378 interest that likely is a false positive (see panels (h) and (i)).
379 f) Scatterplot showing the performance of the GIST model (y-axis) versus the performance of a
380 base-model based on only gene expression data (x-axis) for six pathologist-annotated spatial
381 transcriptomics slides. Performance is defined as the ratio of the median proportion of immune
382 cells in pathologist labeled immune cell slide spots, versus the median proportion of immune cells
383 in the other slide spots ($Q$, see Methods). Points are colored by slide ID. The red line is the
384 identity line (intercept of 0, slope of 1), and the distance between this line and each point (black
385 arrow) represents the observed test statistic $\Delta$ for that sample.
386 g) Histogram showing the empirical null distribution of ratio-based test statistic ($\Delta_{perm}$, see
387 Methods) generated using a permutation procedure (x-axis). The test statistic is a measure of
388 improvement in model performance, versus the pathologist-annotated ground truth, when deep-
389 learning derived prior cell type annotations are incorporated. The observed test statistic $\Delta$ is
390 shown using a vertical red line. $P$-value from permutation test.
391 h) Slide G1 showing the GIST model-derived immune cell proportions, when the deep learning
392 immune cell type annotation has been used as an informative prior. Solid boxes indicate regions
393 of pathologist's annotated immune spots. Green indicates that immune spots were successfully
394 identified, and red indicates that immune spots were not well captured. The dashed black box,
395 highlighted by the black arrowhead, indicates the same region of interest as in (e), where the false
396 positive immune cell predictions have been mitigated. The green arrowhead highlights a region
397 where the correct identification of a pathologist annotated immune-infiltrated region has
398 improved.
399 i) Tissue image showing the region of interest highlighted by a dashed black box in panels (e) and
400 (h). The H&E stain shows minimal evidence of immune infiltration in the areas overlapping the
401 three spatial transcriptomics spots, whose location is shown by black circles.
402 Abbreviations: ST, Spatial Transcriptomics.

403

**The GIST model identified large regions of immune cell infiltration that were missed by**
404
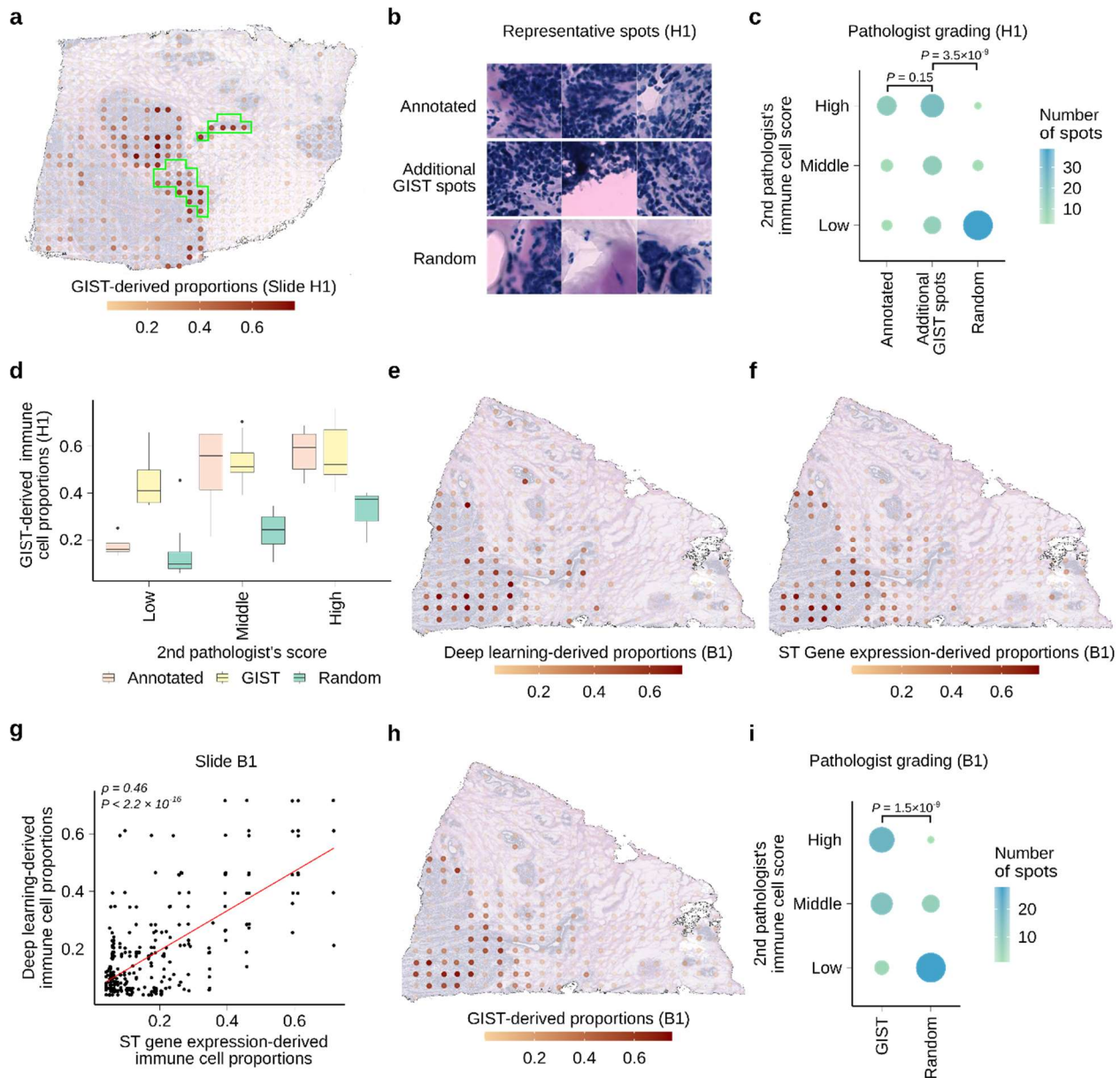**the initial pathologist.**
405

406 Surprisingly, one of the six breast cancer slides assessed demonstrated a statistically significant

407 *decrease* in performance when we leveraged the image-derived prior estimates of immune cell

408 infiltration (slide H1 in Fig. 4f, $P = 3.56 \times 10^{-11}$, Supplementary Figure S12g). However, closer

409 inspection of this slide's results revealed that there was a large region of this tumor that was

410 identified as immune cell infiltrated by both the spatial transcriptomics assay and the deep

411 learning model, but this region was not marked by the initial pathologist's annotation

412    (Supplementary Figure S12a-S12e and S12h)). Unsurprisingly, this region was predicted as

413    heavily immune cell infiltrated by the GIST model, which also correctly identified the original

414    pathologist's annotated regions of immune infiltration in this slide (Fig. 5a, Supplementary

415    Figure S12f).

416    Thus, we hypothesized that the apparent decrease in performance may have represented an

417    oversight in the initial pathologist's annotation, and thus a deficiency in the assumed ground

418    truth, rather than a deficiency in the GIST model's prediction. To test this, we devised a

419    procedure that would allow a second independent pathologist (see Author's Contributions) to re-

420    examine the relevant regions of this slide, while remaining blinded to the GIST model's output

421    and the original pathologist's annotation. The second pathologist was presented with ($n$ = 115)

422    100 × 100-micron subregions from this slide and asked to categorize them as either low, middle,

423    or high levels of immune cell infiltration. These subregions were chosen either from (i) the first

424    pathologist's annotated immune cell regions (ii) high-confidence immune cell regions identified

425    by the GIST model but not the first pathologist or (iii) other randomly chosen regions

426    (representative examples shown in Fig. 5b; see Methods). Remarkably, the second

427    pathologist's reannotation determined no statistical difference between the high-confidence

428    regions of immune cell infiltration annotated by the first pathologist and the additional high-

429    confidence regions identified by the GIST model, which were missed by the first pathologist

430    (Fig. 5c; $P$ = 0.15 from Wilcoxon rank-sum test). However, the high-confidence regions of

431    immune cell infiltration identified by GIST were much more likely to be marked as high

432    probability regions of immune cell infiltration when compared to randomly chosen slide regions

433    (Fig. 5c, $P$ = 3.5 × $10^{-9}$). Additionally, the second pathologist's high confidence immune

434    infiltrated regions were mirrored by higher estimated proportions by GIST (Fig. 5d). These

435    results support the notion that the additional regions identified by the GIST model were true

436    regions of immune cell infiltration and that the poor performance on this slide arose from an

19

437 omission in the original pathologist's annotation, not falsely identified regions by the GIST

438 model.

439 We also reexamined the two available spatial transcriptomics slides where the original

440 pathologist's annotation of the H&E images had not identified any regions of immune cell

441 infiltration (Supplementary Figure S13a-b). Surprisingly, for both slides the deep learning model

442 (Fig. 5e, Supplementary Figure 13c) and the expression-only cell type predictions from the

443 spatial transcriptomics assay (Fig. 5f, Supplementary Figure 13d) agreed that there were in fact

444 regions of immune cell infiltration (Fig. 5g, Spearman's correlation = 0.46, $P < 2.2 \times 10^{-16}$;

445 Supplementary Figure 13e, Spearman's correlation = 0.25, $P < 2.2 \times 10^{-4}$). Unsurprisingly, these

446 same regions were identified by the GIST model (Fig. 5h, Supplementary Figure 13f) and thus it

447 seemed plausible that the initial pathologist had also missed these immune infiltrated regions in

448 their initial examination of these two slides. We used the same scoring procedure outlined

449 above to reannotate these slides by the second pathologist, who convincingly annotated these

450 predicted regions as true regions of immune cell infiltration (Figure 5i, $P = 1.5 \times 10^{-9}$;

451 Supplementary Figure 13g, $P = 4.5 \times 10^{-2}$; see Methods), which were also mirrored by higher

452 proportions estimated by GIST (Supplementary Figure 13h-i). Taken together, these results

453 suggest that our GIST model, which can jointly leverage image-derived deep learning

454 predictions with spatial transcriptomics data, has the potential to outperform a human

455 pathologist in identifying predictive and prognostically important features in human tissue

456 sections.

20

**Figure 5: GIST model identifies regions of immune cell infiltration that were missed by an initial pathologist's annotation.**

a) GIST model-derived proportions plotted on top of tissue from slide H1. Green outline indicates the original annotation of immune infiltrated spot regions identified by the initial pathologist.

b) Three representative 100×100 micron images showing spots from the first pathologist's annotated regions of immune cell infiltration (top), additional high confidence immune infiltrated regions identified by the GIST model (middle), and additional randomly selected regions (bottom). Spots are taken from slide H1.

c) Dot plot showing the second pathologist's immune infiltration grading with a score of low, middle, and high (y-axis) for spots from different regions of the tissue (x-axis). Spots were taken from slide H1 from regions previously annotated by the first pathologist as immune-rich, additional high confidence regions from the GIST model, and additional random regions on the slide. *P*-values from Wilcoxon rank sum test.

d) Boxplot showing distribution of GIST model predicted immune cell proportions (y-axis) broken down by immune infiltration grade (x-axis) provided by the second pathologist. For each pathologist grade (low, middle & high), GIST scores are shown for spots from annotated, GIST high confidence, and random regions. Spots taken from slide H1.

e) Deep learning-derived proportions for spots on slide B1. The color scale shows the predicted proportion of immune cells at a spot.

f) Gene expression-derived proportions for slide B1 from GIST base-model. The color scale shows the predicted proportion of immune cells at a spot.

g) Scatter plot showing the per-spot correlation between deep learning-derived predictions (y-axis) and ST gene expression-derived proportions (x-axis) for slide B1. Each dot is a spot and the red line is the regression line. *P*-value from Spearman's correlation test.

h) GIST model-derived proportions for slide B1. The color scale shows the predicted proportion of immune cells at a spot.

i) Dot plot showing the second pathologist's immune infiltration grading with a score of low, middle, and high (y-axis) for spots from different regions of the tissue (x-axis). Spots were taken from slide B1 from high confidence regions from the GIST model and random regions on the slide. *P*-value from Wilcoxon rank sum test.

In all boxplots, the center line represents the median, bound of box is upper and lower quartiles and the whiskers are 1.5× the interquartile range.

Abbreviations: ST, Spatial Transcriptomics.

## DISCUSSION

We have presented a conceptually novel computational methodology that can leverage data derived from paired tissue images to improve inferences of cell type composition in spatial transcriptomics data. The most exciting application of such a methodology may be the ability to leverage inferences from deep-learning models applied to pathology images, which itself has recently reached close to pathologist level performance in annotating clinically relevant features of tissue sections[20,21,30]. However, the methodology is highly generalizable and could be easily extended to any image-derived prior information, which we have demonstrated for immunofluorescence. Our proposed integrated approach may have clinical applications in areas of prognostics and diagnostics that rely on cell type information but also has the potential to improve all downstream applications of spatial transcriptomics that rely on accurate cell type annotations, including identification of cell-cell or gene-gene interactions.

504    Our framework will also spur the development of future similar computational approaches.

505    Indeed, any cell type decomposition method that could be re-implemented in a Bayesian

506    framework could be adapted to leverage image-derived prior information and this is likely

507    possible for most of the existing models used in our comparisons-of-methods (Figs. 2 and 3).

508    Thus, there is scope for future model development and optimization within our novel framework.

509    We also anticipate that our framework will lead to new modes of spatial transcriptomics

510    experimental design. For example, we showed that IF data could also be informatively

511    leveraged. This opens the possibility of *a priori* staining for a few particularly informative protein

512    markers, knowing that such markers can be used in downstream analyses to directly influence

513    and improve the results of the spatial transcriptomics data analysis. This may be particularly

514    useful for separating cell types when multicollinearity affects the performance of conventional

515    models for cell type decomposition[22].

516    Additionally, while we have shown some illustrative examples, the Bayesian implementation

517    allows enormous flexibility in how prior information is specified. It is theoretically possible to, for

518    example, apply one prior to groups of cell types, or apply multiple partially overlapping priors

519    derived from various sources of information. For the breast cancer dataset shown, we also fixed

520    the $\lambda$ hyperparameter to 50, using information obtained in the previous dataset. While certainly

521    avoiding overfitting, this is likely a conservative means by which to choose this key value and

522    also assumes that the influence of the image and gene expression data should be treated as

523    equal at each spot—almost certainly an oversimplification. Methods could likely be devised to

524    adaptively adjust the value of the $\lambda$ hyperparameter, such that, for example, the differences in

525    uncertainty associated with the deep learning-based outputs could be accounted for at each

526    tissue region. Thus, it is likely that creative applications within the described framework will

527    eventually yield improvements over the results presented here.

528    In conclusion, we anticipate that jointly leveraging spatial transcriptomics and pathology images

529    collected from the same tissue section will represent an important conceptually novel

530    computational methodology, which has the potential to improve many applications of emerging

531    spatial transcriptomics technologies, including potential translational applications in clinical and

532    diagnostic pathology.

533

534    **METHODS**

535    **Technical details of the GIST statistical model.**

536    The expression of gene $i$ at each spatial transcriptomics mRNA capture spot $j$ is assumed to be

537    approximately a weighted sum of the average expression of that gene in each of the cell types

538    captured by that spot. If our spatial transcriptomics data are arranged in a matrix $\boldsymbol{Y}$, where the

539    rows represent $i = 1, ..., m$ genes and the columns represent $j = 1, ..., n$ spots, then this

540    relationship can be summarized by the following equation (see also Fig. 2a):

541    $$Y \approx WH$$

542    where $\boldsymbol{W}$ is an $m \times p$ matrix of cell type specific gene expression signatures, approximating the

543    average expression of each gene in each cell type in this tissue, with each column of $\boldsymbol{W}$

544    representing one of the $p$ cell types and each row representing one of the $m$ genes. $\boldsymbol{H}$ is a

545    $p \times n$ matrix of cell type proportions (or probabilities if the data are subcellular resolution)

546    where each column $\boldsymbol{H}^{(j)}$ represents the proportions of each of $p$ cell types at spot $j$.

547    Each element of $\boldsymbol{W}$ is best modeled from $\boldsymbol{\Psi}$ (a reference single-cell RNA-seq dataset) using a

548    negative binomial distribution estimated for each gene $i$, in each cell type $k$, from the expression

549    data of the available single-cells indexed by $l$. $\phi_{i,k}$ represents the overdispersion parameter of

550    such a distribution:

551    $$\Psi_{i,k,l} \sim NegativeBinomial(w_{i,k}, \phi_{i,k})$$

24

552
$$i = 1, \dots, m; \ k = 1, \dots, p$$

553 For practical reasons, we approximated the elements of $W$ by taking the mean normalized

554 (details below) expression of each gene in each cell type in the reference single-cell RNA-seq

555 dataset $\Psi$.

556 Given $Y$ and $W$, the following model is then used for estimating $H$:

557
$$y_{i,j} | \ \Psi_i, W_i, H^{(j)}, v_j, \beta_{0,j}, \sigma_j \ \sim \ t\left(v_j, \beta_{0,j} + W_i H^{(j)}, \sigma_j\right)$$

558
$$i = 1, \dots, m; \ j = 1, \dots, n; \ k = 1, \dots, p$$

559 We place a gamma prior (priors are denoted herein by $\pi$) on the degrees of freedom parameter

560 $v$ of the $t$-distribution, using shape and rate parameter values previously proposed by Juarez

561 and Steele[33]:

562
$$\pi(v_j) \sim Gamma(2, 0.1)$$

563
$$v_i > 3$$

564 We constrain the elements of $H$ to be positive and to sum to one within each spot:

565
$$\sum_{k=1}^{p} h_{k,j} = 1$$

566
$$h_{k,j} > 0$$

567 This is achieved by placing a non-informative Dirichlet prior on the columns of $H$:

568
$$\pi(H^{(j)}) \sim Dirichlet(\alpha)$$

569
$$\alpha_1 = \alpha_2 = \cdots = \alpha_p = 1$$

570 All other parameters are assigned non-informative priors.

25

571    We use the image data to generate a prior estimate of the abundance of some cell type $a$ (e.g.

572    immune cells) at each spot $j$ (details below), then we place a beta distribution prior on the

573    corresponding proportion of cell type $a$ at spot $j$:

574    $$\pi(h_{a,j}) \sim Beta(\tau_j, \lambda)$$

575    Here, $\tau_j$ is the mean of the beta distribution representing the image-derived prior estimate for

576    the proportion of this cell type $a$ at spot $j$. $\lambda$ is a hyperparameter, representing the total count

577    parameter of the beta distribution, determining how much weight is to be placed on the image

578    data and how much to place on the transcriptomic data.

579    In the notation above, vectors are shown using boldface and matrices bold capital letters. We

580    assume $m$ genes (indexed by $i$), $n$ spots, (indexed by $j$), and $p$ cell types (indexed by $k$).

581

582    **Fitting the GIST and GIST base-model.**

583    The statistical model described above was implemented in the Stan programming language

584    using the *rstan* package. The Hamiltonian Monte Carlo (HMC) algorithm was used to estimate

585    the model parameters. The HMC algorithm was run for 2000 iterations where the first 1000

586    iterations were discarded as burn-in. The posterior mean was used as final parameter

587    estimates.

588

589    **Prior construction.**

590    Mouse brain dataset: To avoid outlier bias in the IF image data the pixel-level image intensity

591    values were first capped at the 99th percentile and values below the 1st percentile were set to

592    zero. These pixel-level intensity values were then rescaled from 0 to 1, by dividing all values by

593  the maximum capped value. Pixels overlapping each spatial transcriptomics mRNA capture spot

594  were defined as those centered around the middle of the spatial transcriptomics spot in a 70-

595  pixel radius—the center of the spot was defined in an annotation file that was output by the 10x

596  Genomics SpaceRanger software. The rescaled pixel-level intensity values were then averaged

597  over the slide regions corresponding to each spatial transcriptomics spot to obtain a single

598  intensity value for each spot. This procedure was repeated for both IF channels—RBFOX3

599  (Neuron) and GFAP (Glia). Finally, the intensity values for each spot in each channel were

600  mapped onto the quantiles of the cell type proportion estimates obtained from a first round of

601  model fitting using the GIST base-model. These IF image-derived mapped spot level intensity

602  values, which act as a proxy for the abundance of neurons or glia, were used as priors on the

603  appropriate parameters in the GIST model.

604  Breast cancer dataset: The deep learning models used in the breast cancer analyses were

605  previously published by Saltz *et al*.[30] and were obtained from the Quantitative Imaging in

606  Pathology (QuIP) group's website (https://sbu-bmi.github.io/quip_distro). These are

607  convolutional neural network-based deep learning models, which had been pre-trained to

608  recognize tumor-infiltrating lymphocytes. The original authors had trained these models using

609  pathologist annotated H&E-stained tissues sections from TCGA.  We used the VGG16-based

610  model provided by the group. The breast cancer H&E images were converted from JPEG format

611  to tiled TIFF format and the software suite VIPS was used to encode the TIFF files with a micron

612  per pixel (MPP) value for each slide. The encoded TIFF files were processed using QuIP's deep

613  learning pipeline to generate a probability map over the entirety of each breast cancer H&E

614  stained slide image. The deep learning model assigned probability values to patches of 50×50

615  microns. For a given spot, the assigned patch-level probability values were converted to spot-

616  level probability values by taking a weighted sum of the patches, where the weight is the pixel

617  overlap between the patch and the spot. This generated values for each spatial transcriptomics

27

618    spot that approximately corresponded to the probability of immune cell infiltration. Similarly to

619    the mouse brain IF dataset, these probability values were then mapped onto the distribution of

620    total lymphocyte (T cell and B cell) content estimated from gene expression-derived proportions

621    alone, obtained by an initial round of model fitting using the GIST base-model. These mapped

622    values were used as informative priors on the appropriate model parameters in the GIST model.

623    The image processing code was implemented in Python using imaging libraries PIL.Image and

624    imageio. Visualization and analysis of imaging data were carried out using the NumPy, pandas,

625    and Matplotlib libraries.

626

627    **Quantifying the improvement achieved by the GIST model, compared to an expression-**

628    **only model, by benchmarking against a pathologist-defined ground truth.**

629    For each slide in the breast cancer dataset, we quantified a model's ability to accurately

630    estimate regions of immune cells by the median of immune cell proportions in spots labeled as

631    immune-infiltrated by the original pathologist, divided by the median of immune cell proportions

632    estimated in the other remaining spots:

633
$$Q = \frac{median(h_{ImmuneSpots})}{median(h_{OtherSpots})}$$

634    $h_{ImmuneSpots}$ is a vector of model-estimated immune cell proportions for spots annotated by the

635    pathologist as containing immune cells, and $h_{OtherSpots}$ are the immune cell proportions

636    estimated at the other spots on the same slide.

637    With better performance, the scalar value $Q$ will increase, as the model's output better matches

638    the pathologist-defined ground truth for this slide. Having defined this performance metric, we

639    defined the improvement of the GIST model over the expression-only GIST base-model below

28

640     as $\Delta$, a scalar representing the difference between this ratio statistic $Q$ when immune cell

641     proportions were estimated with the GIST model ($Q_{GIST}$) or the GIST base-model

642     ($Q_{GISTBaseModel}$):

643     $$\Delta = Q_{GIST} - Q_{GISTBaseModel}$$

644     To assess whether the improved performance $\Delta$ observed for the GIST model over the GIST

645     base-model was statistically significant, we used a permutation-based strategy, building a null

646     distribution by randomly shuffling the pathologist's spot level annotations. Specifically, for each

647     permutation, the spots were randomly assigned as either immune infiltrated or non-immune,

648     fixing the total number of immune infiltrated spots to the same number as the pathologist's

649     annotation of that slide; we then computed the improvement in the performance $\Delta_{perm}$ of the

650     GIST model over the GIST base-model using the same procedure that was applied to the real

651     arrangement of the pathologist's annotations. This was repeated for 100,000 permutations,

652     generating a null distribution against which to compare the observed test statistic $\Delta$. A *P*-value

653     was then calculated by the proportion of permuted values $\Delta_{perm}$ that achieved a value at least

654     as extreme as $\Delta$, the test statistic observed in the pathologist's real annotations. In the cases

655     where no permutated value more extreme than the original test statistic was observed (G1 and

656     H1), a *P*-value was calculated by approximating the null distribution using a normal distribution,

657     with a mean and standard deviation equal to that of the $\Delta_{perm}$ values from the 100,000

658     permutations.

659

660     **Second pathologist's re-annotation of the breast cancer spatial transcriptomics slides.**

661     A second pathologist was asked to assign new immune infiltration grades from H&E images of

662     spots for three spatial transcriptomics breast cancer slides – B1, C1, and H1. The pathologist

663     (co-author Dr. Heather Tillman) was asked to *blindly* score H&E images of slide regions

29

664 overlapping the spatial transcriptomics mRNA capture spots from three groups of spots: These

665 were (i) spots that were annotated as immune cell infiltrated by the original pathologist (slide H1

666 only), (ii) spots that were identified as high-confidence immune infiltrated by the GIST model, or

667 (iii) other randomly chosen spots. High-confidence immune-cell-infiltrated spots from the GIST

668 model were selected as the spots having a predicted proportion of immune cells that was

669 greater than the upper quartile plus 1.5 times the interquartile range of the data, a *de facto*

670 metric used to define outliers. For each slide, the number of random spots selected was equal

671 to the number of spots included from the GIST model. This second pathologist was then asked

672 to score/grade an H&E stain image of each spot, scoring immune cell infiltration levels as low,

673 middle, or high, while remaining blinded to the group from which the spot image was selected.

674 This provided a new score for each spot from each of the three groups (annotated, GIST,

675 random). We then applied a one-sided Wilcoxon rank-sum test to assess whether these scores

676 were significantly higher in the group of spots predicted as high confidence immune infiltrated by

677 the GIST model compared to the randomly selected spots or the immune infiltrated spots from

678 the initial pathologist's annotation, where low, middle and high scores were encoded on an

679 ordinal scale as 1, 2 and 3 respectively.

680

681 **Simulations to assess the ability of the GIST base-model to accurately identify cell type**

682 **composition in gene expression data from a mixture of cell types.**

683 Splatter

684 The accuracy of cell type proportions estimated from the various computational methods was

685 compared to the GIST base-model by first creating synthetic mixtures of gene expression data

686 using the popular Splatter model[25]. We used the Splatter model with a slight modification, which

687 was recently proposed by Zhang *et al*.[34], who reported that the native Splatter model did not

688    capture the empirical distribution of log fold changes observed in real data. The enhanced

689    Splatter model was obtained from the GitHub repository of Zhang *et al*.

690    (https://github.com/Irrationone/splatter), where the author's had learned the simulation

691    parameters from the counts matrix of a publicly available PBMC single-cell RNA-seq dataset

692    generated by 10X Genomics. The parameters for log fold changes were learned by fitting a

693    truncated student's t-distribution to the log fold changes between B cells and CD4 T cells in this

694    same PBMC dataset.

695    Using the enhanced Splatter framework, we generated a dataset with 100 gene expression

696    samples, each created from mixtures of cell types, along with a simulated paired reference

697    single-cell RNA-seq dataset. The paired single-cell RNA-seq data were collapsed by their mean

698    to create the required reference signature matrix $W$, which was passed to each of the

699    computational methods. Each expression mixture sample was generated by taking a weighted

700    average of gene expression across 100 cells (generated independently of the reference single-

701    cell RNA-seq data) from each of six synthetic cell types. Ground truth cell type proportions for

702    the 100 simulated mixture samples were randomly generated from a Dirichlet distribution, where

703    each cell type was assigned equal weight.

704

705    <u>Immune cell deconvolution.</u>

706    We performed a second set of benchmarking simulations using the framework developed by

707    Strum *et al.*[26], which rather than relying entirely on simulation, created a mixture gene

708    expression dataset by computationally mixing real single-cell RNA-seq data, previously

709    generated by Schelker *et al.*[35]. In this benchmark, ground truth was established by mixing gene

710    expression counts from 500 single-cells from each of eight immune cell types in known

711    proportions and the simulated mixture was created by taking an average across cells. For the

31

712    fairest comparison, we supplied each of the methods the LM22 cell type signature matrix[36]

713    (corresponding to $W$ in our notation herein), which is a signature matrix created by the

714    developers of CIBERSORT that represents average gene expression values in each of 22

715    immune cell types. Note this was not possible for Stereoscope, which only accepts single-cell

716    RNA-seq data as the reference input, from which it estimates the cell type signature matrix

717    internally. Because the LM22 cell types do not have a strict one-to-one correspondence with the

718    cell types annotated in Schelker *et al*., the results were mapped to the most relevant cell type

719    using the same mappings previously employed by Strum *et al*.

720    In all simulations, the performance of each method was summarized by the mean absolute error

721    (MAE), which is the average of the absolute value of the difference between each predicted cell

722    type proportion and the known simulated ground truth proportion:

723    $$\text{MAE} = \frac{\sum_{i=1}^{n}|y_i - x_i|}{n} = \frac{\sum_{i=1}^{n}|e_i|}{n}$$

724    Where $y_i$ is a predicted cell type proportion, $x_i$ is a predicted proportion, $e_i$ is the error

725    associated with the prediction, and $n$ is the total number of predicted data points generated by a

726    given method.

727

728    **Datasets used in the analyses.**

729    Mouse Brain

730    The mouse brain spatial transcriptomics Visium data with associated IF images were

731    downloaded from the 10X Genomics website:  https://support.10xgenomics.com/spatial-gene-

732    expression/datasets/1.1.0/V1_Adult_Mouse_Brain_Coronal_Section_2

32

733    As a cell type reference $W$ for these data, we used the curated mouse brain single-cell RNA-seq

734    data provided by Andersson *et al.*[18]. This data had been originally retrieved from

735    http://www.mousebrain.org and was processed by Andersson *et al.* for use in spatial

736    transcriptomics analysis: https://github.com/almaan/stereoscope/tree/master/data/mousebrain

737

738    <u>Breast Cancer</u>

739    The eight separate breast cancer spatial transcriptomics slides, previously generated by

740    Andersson *et al.*, were downloaded from https://github.com/almaan/her2st. This repository

741    contained count matrices generated from the spatial transcriptomics assays, H&E images of the

742    tissue sections (with and without pathologist annotation), and matrices detailing the location of

743    the spots.

744    The single-cell RNA-seq breast cancer dataset, used to generate the cell type reference matrix

745    $W$ for all breast cancer analyses, was previously generated by Karaayvaz *et al.*[37] and obtained

746    from: https://github.com/Michorlab/tnbc_scrnaseq.

747

748    **Data preprocessing, filtering, normalization and imputation.**

749    All public datasets were obtained as preprocessed counts matrices, which had been processed

750    according to the previous authors. Generally, spatial transcriptomics data displayed greater

751    sparsity than the single-cell RNA-seq data, which arises because of differences in platform-

752    specific mRNA capture efficiency. To alleviate this difference, we used a non-parametric

753    imputation approach. Specifically, we used the knnSmooth[38] algorithm (available at the GitHub

754    repository https://github.com/yanailab/knn-smoothing) to impute the spatial transcriptomics data.

755 For the IF mouse brain dataset, we set the "number of nearest neighbors to aggregate"

756 parameter $k$ to 5 and the "number of principal components" parameter $d$ to 10 (author's

757 suggested default). For the breast cancer dataset, we used the same approach with slight

758 modifications. The resolution of spots on the breast cancer slides was coarser than on the

759 Visium array and transcript capture was poorer. Thus, to overcome these limitations, we

760 combined the spots from all the breast cancer spatial transcriptomics slides and imputed them

761 together using the knnSmooth algorithm with a $k$ parameter of 10, mitigating the lower transcript

762 capture efficiency in the breast cancer dataset.

763 Thereafter, both the spatial transcriptomics and single-cell RNA-seq data were normalized

764 separately by using Seurat's SCTransform[39], which importantly removes technical effects such

765 as library size effects. We restricted the single-cell RNA-seq and spatial transcriptomics data to

766 the intersection of their 2,000 most highly variable genes, yielding totals of 1,024 and 837 genes

767 used for GIST model fitting in the mouse and breast cancer datasets respectively.

768

769 **Software and code availability.**

770 The GIST model has been made available as an R package, which can be obtained at:

771 https://github.com/asifzubair/GIST

772 All the code for the analyses presented in this manuscript are available on GitHub:

773 https://github.com/asifzubair/GIST-paper

774 Note: These are private repositories accessible by the links above for peer review, these

775 repositories will be made publicly accessible upon completion of manuscript review.

776

777

**Author Contributions.**

P.G. conceived and directed the project. A.Z. and P.G. wrote the code, with additional input

from R.C. P.G. and A.Z wrote the manuscript. H.T. blindly re-annotated the breast cancer

pathology slides to resolve the discrepancies between the GIST model and the original

pathologist's annotations. S.N., W.C.W, M.P., H.M.L, and J.E. provided additional support in

data analysis and interpretation. All author's edited and approved the final manuscript.

**References.**

1    Peck, M., Moffat, D., Latham, B. & Badrick, T. Review of diagnostic error in anatomical pathology and the role and value of second opinions in error prevention. *Journal of clinical pathology* **71**, 995-1000 (2018).

2    Vickovic, S. *et al.* High-definition spatial transcriptomics for in situ tissue profiling. *Nature methods* **16**, 987-990 (2019).

3    Stickels, R. R. *et al.* Highly sensitive spatial transcriptomics at near-cellular resolution with Slide-seqV2. *Nature biotechnology* **39**, 313-319 (2021).

4    Rodriques, S. G. *et al.* Slide-seq: A scalable technology for measuring genome-wide expression at high spatial resolution. *Science* **363**, 1463-1467 (2019).

5    Liu, Y. *et al.* High-spatial-resolution multi-omics sequencing via deterministic barcoding in tissue. *Cell* **183**, 1665-1681. e1618 (2020).

6    Chen, A. *et al.* Large field of view-spatially resolved transcriptomics at nanoscale resolution. *bioRxiv* (2021).

7    Van de Velde, L.-A. *et al.* Neuroblastoma formation requires unconventional CD4 T cells and myeloid amino acid metabolism. *bioRxiv* (2021).

8    Moncada, R. *et al.* Integrating microarray-based spatial transcriptomics and single-cell RNA-seq reveals tissue architecture in pancreatic ductal adenocarcinomas. *Nature biotechnology* **38**, 333-342 (2020).

9    Svensson, V., Teichmann, S. A. & Stegle, O. SpatialDE: identification of spatially variable genes. *Nature methods* **15**, 343-346 (2018).

10   Sun, S., Zhu, J. & Zhou, X. Statistical analysis of spatial expression patterns for spatially resolved transcriptomic studies. *Nature methods* **17**, 193-200 (2020).

11   Pham, D. T. *et al.* stLearn: integrating spatial location, tissue morphology and gene expression to find cell types, cell-cell interactions and spatial trajectories within undissociated tissues. *bioRxiv* (2020).

12   Maaskola, J. *et al.* Charting tissue expression anatomy by spatial transcriptome decomposition. *BioRxiv*, 362624 (2018).

13   Tanevski, J., Gabor, A., Flores, R. O. R., Schapiro, D. & Saez-Rodriguez, J. Explainable multi-view framework for dissecting inter-cellular signaling from highly multiplexed spatial data. *BioRxiv* (2020).

14   Arnol, D., Schapiro, D., Bodenmiller, B., Saez-Rodriguez, J. & Stegle, O. Modeling cell-cell interactions from spatial molecular data with spatial variance component analysis. *Cell reports* **29**, 202-211. e206 (2019).

15   Elosua-Bayes, M., Nieto, P., Mereu, E., Gut, I. & Heyn, H. SPOTlight: seeded NMF regression to deconvolute spatial transcriptomics spots with single-cell transcriptomes. *Nucleic acids research* **49**, e50-e50 (2021).

16   Cable, D. M. *et al.* Robust decomposition of cell type mixtures in spatial transcriptomics. *Nature Biotechnology*, 1-10 (2021).

17   Biancalani, T. *et al.* Deep learning and alignment of spatially-resolved whole transcriptomes of single cells in the mouse brain with Tangram. *bioRxiv* (2020).

18   Andersson, A. *et al.* Single-cell and spatial transcriptomics enables probabilistic inference of cell type topography. *Communications biology* **3**, 1-8 (2020).

19   van der Laak, J., Litjens, G. & Ciompi, F. Deep learning in histopathology: the path to the clinic. *Nature medicine* **27**, 775-784 (2021).

20   Kather, J. N. *et al.* Pan-cancer image-based detection of clinically actionable genetic alterations. *Nature Cancer* **1**, 789-799 (2020).

836  21  Fu, Y. *et al.* Pan-cancer computational histopathology reveals mutations, tumor composition and
837      prognosis. *Nature Cancer* **1**, 800-810 (2020).
838  22  Li, B., Liu, J. S. & Liu, X. S. Revisit linear regression-based deconvolution methods for tumor gene
839      expression data. *Genome biology* **18**, 1-5 (2017).
840  23  Stenzinger, A. *et al.* in *Seminars in cancer biology.*   (Elsevier).
841  24  Haibe-Kains, B. *et al.* Transparency and reproducibility in artificial intelligence. *Nature* **586**, E14-
842      E16 (2020).
843  25  Zappia, L., Phipson, B. & Oshlack, A. Splatter: simulation of single-cell RNA sequencing data.
844      *Genome biology* **18**, 1-15 (2017).
845  26  Sturm, G. *et al.* Comprehensive evaluation of transcriptome-based cell-type quantification
846      methods for immuno-oncology. *Bioinformatics* **35**, i436-i445 (2019).
847  27  Chen, B., Khodadoust, M. S., Liu, C. L., Newman, A. M. & Alizadeh, A. A. in *Cancer systems*
848      *biology*   243-259 (Springer, 2018).
849  28  Gong, T. & Szustakowski, J. D. DeconRNASeq: a statistical framework for deconvolution of
850      heterogeneous tissue samples based on mRNA-Seq data. *Bioinformatics* **29**, 1083-1085 (2013).
851  29  Kleshchevnikov, V. *et al.* Comprehensive mapping of tissue cell architecture via integrated single
852      cell and spatial transcriptomics. *bioRxiv* (2020).
853  30  Saltz, J. *et al.* Spatial organization and molecular correlation of tumor-infiltrating lymphocytes
854      using deep learning on pathology images. *Cell reports* **23**, 181-193. e187 (2018).
855  31  Jochems, C. & Schlom, J. Tumor-infiltrating immune cells and prognosis: the potential link
856      between conventional cancer therapy and immunity. *Experimental biology and medicine* **236**,
857      567-579 (2011).
858  32  Taube, J. M. *et al.* Association of PD-1, PD-1 ligands, and other features of the tumor immune
859      microenvironment with response to anti–PD-1 therapy. *Clinical cancer research* **20**, 5064-5074
860      (2014).
861  33  Juárez, M. A. & Steel, M. F. Model-based clustering of non-Gaussian panel data based on skew-t
862      distributions. *Journal of Business & Economic Statistics* **28**, 52-66 (2010).
863  34  Zhang, A. W. *et al.* Probabilistic cell-type assignment of single-cell RNA-seq for tumor
864      microenvironment profiling. *Nature methods* **16**, 1007-1015 (2019).
865  35  Schelker, M. *et al.* Estimation of immune cell content in tumour tissue using single-cell RNA-seq
866      data. *Nature communications* **8**, 1-12 (2017).
867  36  Newman, A. M. *et al.* Robust enumeration of cell subsets from tissue expression profiles. *Nature*
868      *methods* **12**, 453-457 (2015).
869  37  Karaayvaz, M. *et al.* Unravelling subclonal heterogeneity and aggressive disease states in TNBC
870      through single-cell RNA-seq. *Nature communications* **9**, 1-10 (2018).
871  38  Wagner, F., Yan, Y. & Yanai, I. K-nearest neighbor smoothing for high-throughput single-cell
872      RNA-Seq data. *BioRxiv*, 217737 (2018).
873  39  Hafemeister, C. & Satija, R. Normalization and variance stabilization of single-cell RNA-seq data
874      using regularized negative binomial regression. *Genome biology* **20**, 1-15 (2019).

875