

502 **Estimation of cell type reference signatures from scRNA-seq.** Given cell type
 503 annotation for each cell, the corresponding reference cell type signatures $g_{f,g}$, which represent
 504 the average mRNA count of each gene g in each cell type $f = \{1, \dots, F\}$, can be estimated
 505 using a negative binomial regression model, which allows for combining data across batches
 506 and technologies (see below and Suppl. methods).

507 **Cell2location model.** An untransformed spatial expression count matrix $d_{s,g}$ is used
 508 for input, as obtained from the 10X SpaceRanger software (10X Visium data). Cell2location
 509 models the elements of $d_{s,g}$ as Negative Binomial (NB) distributed, given an unobserved gene
 510 expression level (rate) $\mu_{s,g}$ and a gene-specific over-dispersion α_g :

$$511 \quad 512 \quad d_{s,g} \sim NB(\mu_{s,g}, \alpha_g).$$

514 The expression level of genes $\mu_{s,g}$ in the mRNA count space is modelled as a linear
 515 function of expression signatures of reference cell types $g_{f,g}$:

$$517 \quad \mu_{s,g} = \underbrace{m_g}_{\text{technology sensitivity}} \cdot \underbrace{\left(\sum_f w_{s,f} g_{f,g} \right)}_{\text{cell type contributions}} + \underbrace{l_s + s_g}_{\text{additive shift}},$$

518 where, $w_{s,f}$ denotes regression weight of each reference signature f at location s , which can
 519 be interpreted as the number of cells at location s that express reference signature f ; m_g is a
 520 gene-specific scaling parameter, which adjusts for global differences in sensitivity between
 521 technologies; l_s and s_g are additive variables that account for gene- and location-specific shift,
 522 such as due to contaminating or free-floating RNA.

524 To account for the similarity of location patterns across cell types, $w_{s,f}$ is modelled
 525 using another layer of decomposition (factorization) using $r = \{1, \dots, R\}$ groups of cell types,
 526 that can be interpreted as cellular compartments or tissue zones (Suppl. Methods). Unless
 527 stated otherwise, R is set to 50.

528 Approximate Variational Inference is used to estimate all model parameters,
 529 implemented in the pymc3 framework⁵², which supports GPU acceleration. For full details see
 530 Suppl. Methods.

531 **Note on selecting scRNA-seq profiles for constructing reference cell type data.**
 532 It is important to aim for a comprehensive and detailed cell-type reference, which includes as
 533 many of the cell types and subpopulations that are present *in-situ* as possible, for example,
 534 by generating a paired snRNA-seq reference from the same tissue sample. However,
 535 imperfect matching of cell populations is often acceptable (see Fig 4, Fig S4D). In such
 536 instances, the stability of the model fit, which can be assessed using multiple random restarts,
 537 can serve as diagnostic criteria (see Suppl. Methods).

538 **Note on selecting the method for estimating reference signatures of cell types.**
 539 The first step of our model is to estimate reference cell type signatures from sc/snRNA-seq
 540 profiles, by providing the model with annotated cell type and subpopulation labels for each
 541 cell. The cell2location software comes with two implementations for this estimation step: 1) a
 542 statistical method based on Negative Binomial regression and 2) hard-coded computation of
 543 per-cluster average mRNA counts for individual genes. We generally recommend using NB
 544 regression, which allows to robustly combine data across technologies and batches (Fig S23),

545 which results in improved spatial mapping accuracy (Fig S22B). However, when the batch
546 effects are small a faster hard-coded method of computing per cluster averages provides
547 similarly high accuracy (Fig S22A). We also recommend the hard-coded method for non-UMI
548 technologies such as Smart-Seq 2.

549 **Hyperparameter selection.** The cell2location model has 4 hyper-priors, which can be
550 set by the user taking known experimental and biological characteristics of a given dataset
551 into consideration:

- 552 1) Expected number of cells per location \hat{N}
- 553 2) Expected number of cell types per location \hat{Y}
- 554 3) Expected number co-abundance cell type groups per location \hat{A}
- 555 4) Expected mean of gene-specific technology sensitivity parameter μ_m

556 The Fig S24 provides a flowchart of how the values of these hyper-priors can be determined.

557 *Expected cell abundance \hat{N} per location* is a tissue-level global estimate, which can be
558 derived from histology images (H&E or DAPI), ideally paired to the spatial expression data or
559 at least representing the same tissue type. This parameter can be estimated by manually
560 counting nuclei in a 10-20 locations in the histology image (e.g. using 10X Loupe browser, Fig
561 S8), and computing the average cell abundance. An appropriate setting of this prior is
562 essential to inform the estimation of absolute cell type abundance values, however, the model
563 is robust to a range of similar values (Fig S5). In settings where suitable histology images are
564 not available, the size of capture regions relative to the expected size of cells can be used to
565 estimate \hat{N} (Slide-Seq V2, Fig S24). For all analysis in this manuscript, a single tissue-level
566 estimate was used, however, as an advanced feature, cell2location can utilise the per-location
567 number of cells.

568 *Expected number of cell types per location \hat{Y} and expected number co-abundance cell
569 type groups per location \hat{A} .* The value of these hyper-priors has minimal effect on model
570 accuracy (Fig S5). Consequently, we recommend setting their values to 7, a single global
571 estimate.

572 *The difference in technology sensitivity mean μ_m and variance σ_m^2 parameters* can be
573 chosen by comparing the average total number of mRNA per cell in the reference cell type
574 data to the average total number of mRNA per location in the spatial data divided by \hat{N} (Fig
575 S24).

576 While good choices of these hyper-parameters can have a positive impact on
577 accuracy, the estimate of relative cell abundance is robust to a range of suboptimal choices
578 (Fig S5). The estimation of absolute cell abundance requires appropriate settings of \hat{N} and
579 μ_m in particular.

580 **Constructing a synthetic spatial transcriptomics data set**

581 Simulated spatial transcriptomics data were generated by combining expression
582 profiles of cells drawn from each one of 49 cell types in the mouse brain snRNA-seq reference
583 data (see below), to generate abundance profiles at 2,500 locations. snRNA data from the two
584 most homogenous mouse brain snRNA-seq samples were split into one dataset used to
585 generate the synthetic data (50% of cells) and a second dataset used to evaluate cell2location
586 and alternative approaches (50% of cells), similarly to the strategy proposed by Andersson et
587 al³. Hyperparameters for data simulation were chosen to mimic (the typically low) cell counts
588 observed for cell types in real tissues, additionally matching sparsity profiles as observed in
589 real data. Cell type abundances were simulated according to either a spatially ubiquitous

590 pattern (8 cell types), or a regional pattern (41 cell types). Regional patterns are represented
 591 by 12 tissue zones defined by co-located cell types that mimic the organisation of real tissues.
 592 The assignment of 41 regional cell types to the 12 tissue zones is shown in Fig S2, with each
 593 cell type belonging to 1-3 tissue zones and each tissue zone containing 2-8 cell types. The
 594 number of cell types present at each location, as well as the absolute abundance (the number
 595 of cells per location), were simulated according to either low or high average cell type
 596 abundance (Fig 1B), stratified by ubiquitous and regional location pattern (see below). The
 597 mathematical description and the step-by-step procedure to simulate abundance of cell types
 598 across locations and to generate multi-cell mRNA counts is described in detail below in three
 599 sections: 1) generating abundance of cell types across locations, 2) generating expected multi-
 600 cell mRNA expression of genes across locations, 3) generating multi-cell mRNA integer counts
 601 weighted by technology difference effect.

602
 603 **First**, follow this step-by-step procedure to generate ground truth spatial abundance
 604 $w_{s,f}$, integer cell count $count_{s,f}$ and fraction of mRNA captured $frac_{s,f}$ for cell types f across
 605 locations s :

- 606 1. Assign cell types to ubiquitous (n=8) and regional (n=41) abundance patterns (denoted
 607 as r).
- 608 2. Perform binary assignment of 41 sparse cell types to 12 tissue zones and 8 ubiquitous
 609 cell types to 8 ubiquitous patterns (total n=20), shown in Fig S2A and denoted as $x_{r,f}$.
- 610 3. Stratified by location pattern, randomly assign up to 20% of cell types to high
 611 abundance groups and all other cell types to low abundance groups. Generate per cell
 612 type average abundance d_f , which is different for 4 groups shown in Fig 1B:

613 **a) Ubiquitous and low density:** 5 cell types present in most locations at
 614 density:

$$615 d_f \sim Gamma(\mu = 1.0, \sigma^2 = \mu / 5).$$

616 **b) Ubiquitous and high density:** 3 cell types present in most locations at
 617 density:

$$618 d_f \sim Gamma(\mu = 2.8, \sigma^2 = \mu / 5).$$

619 **c) Sparse and low density:** 32 cell types present in sparse tissue zones at
 620 density:

$$621 d_f \sim Gamma(\mu = 1.0, \sigma^2 = \mu / 5).$$

622 **d) Sparse and high density:** 9 cell types present in sparse tissue zones at
 623 density:

$$624 d_f \sim Gamma(\mu = 2.8, \sigma^2 = \mu / 5).$$

625 By following this procedure, sparsity and density parameters for each cell type were
 626 generated that produced an average total number of cells per location close to 10,
 627 mimicking cell count observed by nuclear segmentation of the mouse brain histology
 628 images (Fig S8, Suppl Methods).
 629 Per cell type maximum abundance d_f was used to scale $x_{r,f} = x_{r,f} * d_f$, thus defining
 630 the average abundance of each cell type across patterns r .

- 631 4. Generate spatial abundance $z_{s,r}$ for locations s for 20 location patterns (denoted as r)
 632 representing 12 tissue zones and 8 ubiquitous cell types. Gaussian Process in 50x50
 633 grid of locations was used with randomly generated bandwidth parameters:

$$634 bw \sim Gamma(\mu = 8.0, \sigma^2 = \mu / 1.2)$$

- a) 8 ubiquitous patterns with for non-zero density in most locations:

642
636
643
637
638
639
640
641
644
645
646
647
648
649
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
676
677

$$z_{s,r} \sim GP(\mu = 0, \text{eta} = 0.5, \text{bw} = \text{bw})$$

b) 12 tissue zones with $z_{s,r} \sim GP(\mu = 0, \text{eta} = 1.5)$ for sparse locations

$$z_{s,r} \sim GP(\mu = 0, \text{eta} = 1.5, \text{bw} = \text{bw})$$

To ensure positive scale, cell abundance for each pattern $z_{s,r}$ were softmax-transformed $z_{s,r} = \exp(z_{s,r}) / \sum_r \exp(z_{s,r})$. Next, to ensure that maximum for each location pattern r is equal to 1 further normalisation was applied: $z_{s,r} = z_{s,r} / \sum_s z_{s,r}$, which is needed to use abundances established in step 3 as average value for each cell type.

5. Per cell type abundance for each location s was generated as $w_{s,f} = (\sum_r z_{s,r} x_{r,f}) q_{s,f}$ (shown in Fig S2B), where $\log(q_{s,f}) \sim Normal(\mu = 0, \sigma = 0.35)$ introduces randomness to abundances of individual cell types. This additional variability mimics the observation that co-located cell types in real tissues do not have perfectly correlated abundance within tissue zones.
6. Cell abundance $w_{s,f}$ was used to generate integer cell count $count_{s,f}$ and fraction of mRNA captured $frac_{s,f}$ for each location and cell type as follows:
 - a) generate $count_{s,f}$ by rounding $w_{s,f}$ to the smallest integer, such that $count_{s,f} \geq w_{s,f}(\lceil w_{s,f} \rceil)$.
 - b) Compute the fraction of mRNA captured as $frac_{s,f} = w_{s,f} / count_{s,f}$.

Second, follow this step-by-step procedure to use 1) the integer cell count $count_{s,f}$ and 2) the fraction of mRNA captured $frac_{s,f}$ for each cell type f in a given location s to generate expected multi-cell mRNA count profiles $ed_{s,g}$ for every gene g in a given location s by combining cells c drawn from reference cell types f in the snRNA-seq data $j_{c,g}$ as follows:

1. Randomly select indices of cells $c \in f$ that form a subset $p \subset c$ containing $n = count_{s,f}$ cells.
2. Construct per cell type expected mRNA abundance profiles for a given location and cell type:

$$ed_{s,f,g} = (\sum_{c \in p} j_{c,g}) frac_{s,f}$$
3. Construct multi-cell expected mRNA abundance profiles by adding mRNA across all cell types:

$$ed_{s,g} = \sum_f ed_{s,f,g}$$

Third, follow this step-by-step procedure to generate multi-cell mRNA integer counts $d_{s,g}$. In this step, gene-specific scaling was applied, denoted as m_g , to mimic the difference in sensitivity between technologies and counts were samples from Poisson distribution:

$$d_{s,g} \sim Poisson(ed_{s,g} m_g),$$

where m_g characterises the difference between the mouse brain Visium data and single nucleus RNA-seq reference (Fig S2C, estimated by cell2location). Using these values makes the simulated data representative for mapping single nucleus RNA-seq derived reference cell types.

Under this simulation, the total number of mRNA per location mimics that observed in mouse brain data (Fig S2D).