# Protocol 2: Assessing cell type replicability against a pre-trained reference taxonomy

Protocol 2 demonstrates how to assess cell types of a newly annotated dataset against a reference cell type taxonomy. Here we consider the cell type taxonomy established by the Brain Initiative Cell Census Network (BICCN) in the mouse primary motor cortex. The BICCN taxonomy was defined across a compendium of datasets sampling across multiple modalities (transcriptomics and epigenomics), it constitutes one of the richest neuronal resources currently available. When matching against a reference taxonomy, we assume that the reference is of higher resolution than the query dataset, i.e. the query dataset samples the same set or a subset of cells compared to the reference.

## Step 1 - Pre-train a reference MetaNeighbor model.

1. We start by importing utility packages and setting up the default behavior for plots.

```
] a d c ˙f b h  a d U m g b d
] a d c ˙f d h U b X ˙W g j d X
] a d c ˙f g h W U b ˙d U m g g W
] a d c ˙f a h U h d ` c h ` ] V ˮ U d g m d d ` ` h c h
] a d c ˙f g h Y U V c U f d g b ˙g
] a d c ˙f d h m a ˙b
] a d c ˙f f h Y
```

```
ı a U h d ` c h ] b V ] b Y
```

```
˳ H \ Y g Y ˙ g U j Y ˙ W \ U f U W h Y f g ˙ U g ˙ h Y I h ˙ ] b ˙ D 8 : g
] a d c ˙f a h U h d ` c h ` ] V
a U h d ` c h ˙h f ` W D V U f U f ] a d g X Z ˮ Z c b h O h 1 m d & f ]
a U h d ` c h ˙h f ` W D V U f U f ] a d g g ˮ Z c b h O ˙ h c d & f ]
˙
˳ H \ Y g Y ˙ W \ U b [ Y ˙ d ` c h ˙ U Y g h \ Y h ] W g
˙
g b ˙g g Y h h g h m 1 f Y k \ ] h z Y Z f c b h S g W U ˮˮ & Y )
d ` ˙h f ˙ W U l Y g ˮ g d z ˙ b c h g ˮ U ˙ g Y ] [ \ h U ` ˙ b Y
d ` ˙h f ˙ W f ] l h ] W U _ V f c h h 1 c l U f ] Ł Y
d ` ˙h f ˙ W f ] m h ] W U _ ` f Y Z 1 h l f ] Ł Y
```

1. We load an already merged Anndata object containing the BICCN dataset. The full code for generating the dataset is available here, the dataset itself can be downloaded directly from Figshare using the link below.

```
˙ W i f ` ˙ ! @ ˙ ! c ˙ V ] W W b S \ j [ ˮ \ ) U X ˙ \ h h d g . # # b X c k b ` c U X Y f ˮ Z
```

```
= b s O *   U X U h 1 U g W f Y U X S \f l j i U X WW b S \ j [ Ł " \ ) U X fi
```

```
# I g Y f g # ` Y c b # a ] b ] Wc b X U ' # Y b j g # 6 = 7 5 B S a c i g Y # ` ] V # d mh
c f Y # U b b X U h U " d m . %, &, . ` I g Y f K U f b ] b [ . ` C V g Y f j U h ] c b ` b l
Y a ` i b ] e i Y ž ` WU ` ` ` T " c V g S b U a Y g S a U _ Y S i b ] e i Y T "
` ` i h ] ` g " k U f b S b U a Y g S X i d `` ] WU h Y g fl ˝ c V g ˝ Ł
```

```
= b s O +   U X U h d U V ˝ g Wc ` i á b ˝ g X U h d U V ˝ g Wc ` i á U b g h m f U Y h f
```

The BICCN data contains 7 datasets totaling 482,712 cells. There are multiple sets of cell type labels depending on resolution (class, subclass, cluster) or type of labels (independent labels or labels defined from joint clustering). Note that, to reduce memory usage, we have already computed and restricted the dataset to a set of 319 highly variable genes.

1. We create pre-trained models with the *trainModel* function, which has identical parameters as the *MetaNeighborUS* function. Here, we choose to focus on two sets of cell types: subclasses from the joint clustering (medium resolution, e.g., Vip interneurons, L2/3 IT excitatory neurons), and clusters from the joint clustering (high resolution, e.g., Chandelier cells).

Since the dataset has already been subsetted to the highly variable genes we can make a column of all Trues under ` " j U f O fi \ ] [ \ ` m S j U f ] U V ` Y fi Q `

```
= b s O ,   U X U h j U U O fi \ ] [ \ ` m S j U f Q 1 U H V f ` i Y Y fi
```

```
= b s O -   d h f U ] b Y X S g i 1 W m b ` g b U ] b A d U XX Y U h U i g h i X m š ] fi X f i ] b h S g i V W` U g i
             d h f U ] b Y X S g i " W W S W f g j d f Y h f U ] b Y X S V ] WWb S g i V W` U g g Y g " Wg j f
             d h f U ] b Y X S W ` i d g m h U Y ff U ] b A d U XX Y U h U i g h i X m š ] fi X f i ] b h S W ` i g h Y Ł f S
             d h f U ] b Y X S W h i c g h W W i f d f Y h f U ] b Y X S V ] WWb S W Ł ` i g h Y f g " Wg j fi
```

For simplicity of use, we store the pretrained models to file using the "write_csv" function in pandas.

# Step 2 - Compare annotations to pre-trained taxonomy

1. We start by loading our query dataset (Tasic 2016, neurons from mouse primary visual cortex, available for download using curl) and our pre-trained subclass and cluster taxonomies.

   > Tasic data was aquired using the R scRNAseq package. You can see the code for aquiring and processing the data using a combination of these two R and python scripts

Note that we add a "study_id" column to the Tasic metadata, as this information will be needed later by MetaNeighbor.

1. To run MetaNeighbor, we use the "MetaNeighborUS" function but, compared to Protocol 1, we provide a pre-trained model instead of a set of highly variable genes (which are already contained in the pre-trained model). We start by checking if Tasic cell types are consistent with the BICCN subclass resolution.

1. We visualize AUROCs as a rectangular heatmap, with the reference taxonomy as columns and query cell types as rows.

As in Protocol 1, we start by looking for evidence of global structure in the dataset. Here we recognize 3 red blocks, which correspond to non-neurons (top left), inhibitory neurons (middle) and excitatory neurons (bottom right). The presence of sub-blocks inside the 3 global blocks suggest that cell types can be matched more finely. For example, inside the inhibitory block, we can recognize sub-blocks corresponding to CGE- derived interneurons (Vip, Sncg and Lamp5 in the BICCN taxonomy) and MGE-derived interneurons (Pvalb and Sst in the BICCN taxonomy).

1. We refine AUROCs by focusing on inhibitory neurons. We use two utility functions ("splitTrainClusters" and "splitTestClusters") to select the relevant cell types.

```
= b s O %    [ U V U Y f [ ] W S 1 V d   m A  W p  d  `  ]  h H Y g h 7  f l h i  U g n 1 Y  W  g ž  g U j  Y S 1: b U g  b g  8 Q
             [ U V U Y f [ ] W S 1 h d U m g a  g W d  `  ]  h H f U ]  b 7 f l h i U g h ž Y W  g ž  g U j  Y S 1: b U g  b g P%Q

             _ Y Y d S W Y 1 ` b g ] b %X
             ` ` d m a" b c ] b S ` U f h Y U g g c W g f i g h i  X m S ž  h X U f i g ] c W g f i d f ] a U f m S O h L m ž d Y fi
             ` ` ` [ U V U Y f [ ] W S h U g ] W
```

The heatmap suggests that there is a broad agreement at the subclass level between the BICCN MOp taxonomy and the Tasic 2016 dataset. For example, the Ndnf subtypes, Igtp and Smad3 cell types from the Tasic dataset match with the BICCN Lamp5 subclass.

1. The previous heatmaps suggest that all Tasic cell types can be matched with one BICCN subclass. We now go one step further and ask whether inhibitory cell types correspond to one of the BICCN clusters.

# I g Y f g # ` Y c b # a ] b ] Wc b X U ' # Y b j g # 6 = 7 5 B S a c i g Y # ` ] V # d mh
g " d m . - , . ` I g Y f K U f b ] b [ . ` F Y d ` U W] b [ ` U b m ` p ` k ] h \ ` U ` " ` ]
` ` k U f b ] b [ g " k U f b fl ˝ F Y d ` U W] b [ ` U b m ` p ` k ] h \ ` U ` '" ` ] b ` g h i
# I g Y f g # ` Y c b # a ] b ] Wc b X U ' # Y b j g # 6 = 7 5 B S a c i g Y # ` ] V # d mh
h f ] I " d m . %% & ( . ` I g Y f K U f b ] b [ . ` T T g e i U f Y 1 H f i ` Y T T ` ] [ b c f
` ` k U f b ] b [ g " k U f b fl a g [ Ł

Here the heatmap is difficult to interpret due to the large number of BICCN cell types
(output omitted here). Because there is a limited number of cell types in the query dataset,
we directly investigate the top hits for each query cell type.

= b s O &   f Y g i ` 1 h h U g ] WS g i " i V b g f i A Y h U B Y ] [ \ O c f I G fi
f Y g i " ` h c O V h U g ] Wp G g h ž 7 g c f h S j U f U i g W b X ] b [ g " X Y U f f % $

We note two properties of matching against a pre-trained reference. First, replicable cell types have a clear top match in each of the reference dataset. Sst Chodl (long-projecting interneurons) match to similarly named clusters in the BICCN with an AUROC > 0.9999, Pvalb Cpne5 (Chandelier cells) match with the Pvalb Vipr2_2 cluster with AUROC > 0.93. Second, we have to be beware of false positives. For example, Sst Chodl secondarily matches with the L6b Ror1 cell types with AUROC > 0.98, an excitatory cell type only distantly related with long-projecting interneurons. When we use a pre-trained model, we only compute AUROCs with the reference data as the train data, so we cannot identify reciprocal hits. If we had been able to use "Tasic|Sst Chodl" as the training cluster, its votes would have gone heavily in favor of the BICCN's Sst Chodl, making L6b Ror1 a low AUROC match on average. Because of the low dimensionality of gene expression space, we expect false positive hits to occur just by chance (e.g., cell types reusing similar pathways) when a cell type is missing in the query dataset. Here L6b Ror1 (an excitatory type) had no natural match with the Tasic inhibitory cell types and voted for its closest match, long-projecting interneurons.

There are three alternatives to separate true hits from false positive hits. First, if a cell type is highly replicable, it will have a clear top matching cluster in the reference dataset. Second, if the query dataset is known to be a particular subset of the reference dataset (e.g., inhibitory neurons, as is the case here), we recommend restricting the reference taxonomy to that subset. Third, if the first two solutions don't yield clear results or cannot be performed, it is possible to go back to reciprocal testing by using the full BICCN dataset instead of the pre-trained reference.

We illustrate the first solution in the case of Chandelier cells.

To illustrate AUROC differences, we chose a logarithmic scaling to reflect that AUROC

values do not scale linearly: when AUROCs are close to 1, a difference of 0.05 is substantial. Here, the best matching BICCN cluster ("Pvalb Vipr2_2") is order of magnitudes better than other clusters, suggesting very strong replicability.

1. The second solution to avoid false positive hits is to subset the reference to cell types that reflect the composition of the query datasets. Since we are looking at inhibitory neurons, we can restrict the BICCN taxonomy to inhibitory cell types, which names all start with "Pvalb", "Sst", "Lamp5", "Vip" or "Sncg".

Again, we note that there is a significant gap between the best hit and the secondary hit, but now secondary hits are closely related cell types (Sst subtype for Sst Chodl, secondary Chandelier cell type Pvalb Vipr2_1 for Pvalb Cpne5).

1. To obtain a more stringent mapping between the query cell types and reference cell types, we use one-vs-best AUROC, which will automatically match the best hit against the best secondary hit.

Now the hit structure is much sparser, which helps identify 1:1 and 1:n hits. The heatmap suggests that most Tasic cell types match with one or several BICCN clusters, which we can further inspect by looking at top hits.

Using this more stringent assessment, we confirm that Sst Chodl strongly replicates inside the BICCN (one-vs-best AUROC ~ 1, best secondary hit = 0.43), same for Pvalb Cpne5 (one-vs-best AUROC > 0.83, best secondary hit = 0.58), while for example Sst Tacstd2 corresponds to multiple BICCN subtypes (including Sst C1ql3_1, Sst C1ql3_2, AUROC > 0.95).

Pre-training a MetaNeighbor model thus provides a rigorous, fast and simple way to query a large reference dataset and obtain quantitative estimations of the replicability of newly annotated clusters.