

Scaling up reproducible research for single cell transcriptomics using MetaNeighbor (Protocol 3)

1 Protocol 3: Functional characterization of replicating clusters

Protocol 3 demonstrates how to characterize functional gene sets contributing to cell type identity. In this section, we will focus on the characterization of inhibitory neuron subclasses as provided by the BICCN. The BICCN has shown that subclasses are strongly replicable across datasets and provided marker genes that are specific to each subclass. MetaNeighbor can be used to further quantify which pathways contribute to the subclasses' unique biological properties.

1.1 Step 1 - Creation of biologically relevant gene sets

1. To compute the functional characterization of clusters, we first need an ensemble of gene sets sampling relevant biological pathways. In this protocol we will consider the Gene Ontology (GO) annotation for mouse. The scripts used to build up-to-date gene sets are here XXX, gene sets can be downloaded directly here XXX.

```
go_sets = readRDS("go_mouse.rds")
```

Gene sets are stored as a named list, each element of the list corresponds to a gene set and contains a vector of gene symbols.

2. Then we load our dataset containing inhibitory neurons from the BICCN. The scripts used to build the dataset can be found here XXX, the dataset can be downloaded here XXX.

```
library(SingleCellExperiment)
biccn_gaba = readRDS("biccn_gaba.rds")
dim(biccn_gaba)
```

```
## [1] 24140 71368
```

3. Next we restrict our gene sets to genes that are present in the dataset. We then filter gene sets to obtain gene sets of meaningful size: large enough to learn expression profiles (> 10), small enough to represent specific biological function or processes (< 100).

```
known_genes = rownames(biccn_gaba)
go_sets = lapply(go_sets, function(gene_set) { gene_set[gene_set %in% known_genes] })
min_size = 10
max_size = 100
go_set_size = sapply(go_sets, length)
go_sets = go_sets[go_set_size >= min_size & go_set_size <= max_size]
length(go_sets)
```

```
## [1] 6488
```

1.2 Step 2: Functional characterization with supervised MetaNeighbor

4. Once the gene set list is ready, we run the supervised “MetaNeighbor” function. Its inputs are similar to “MetaNeighborUS”, but it assumes that cell types have already been matched across datasets (i.e., they have identical names). Here we will use joint BICCN subclasses, for which names have been normalized across datasets (“Pvalb”, “Sst”, “Sst Chodl”, “Vip”, “Lamp5”, “Sncg”). Note that, because we are testing close to 6,500 gene sets, this step is expected to take a long time for large datasets. We recommend using this function inside a script and always save results to a file as soon as it’s done.

```
library(MetaNeighbor)
#devtools::load_all("~/projects/metaneighbor/MetaNeighbor")

system.time({
aurocs = MetaNeighbor(dat = biccn_gaba,
                      experiment_labels = biccn_gaba$study_id,
                      celltype_labels = biccn_gaba$joint_subclass_label,
                      genesets = go_sets[1:10],
                      fast_version = TRUE, bplot = FALSE, batch_size = 50)
})

## [1] "GO:0000002|mitochondrial genome maintenance|BP"
## [1] "GO:0000012|single strand break repair|BP"
## [1] "GO:0000018|regulation of DNA recombination|BP"
## [1] "GO:0000027|ribosomal large subunit assembly|BP"
## [1] "GO:0000028|ribosomal small subunit assembly|BP"
## [1] "GO:0000038|very long-chain fatty acid metabolic process|BP"
## [1] "GO:0000041|transition metal ion transport|BP"
## [1] "GO:0000045|autophagosome assembly|BP"
## [1] "GO:0000050|urea cycle|BP"
## [1] "GO:0000054|ribosomal subunit export from nucleus|BP"

##      user  system elapsed
## 14.044    8.778    5.760

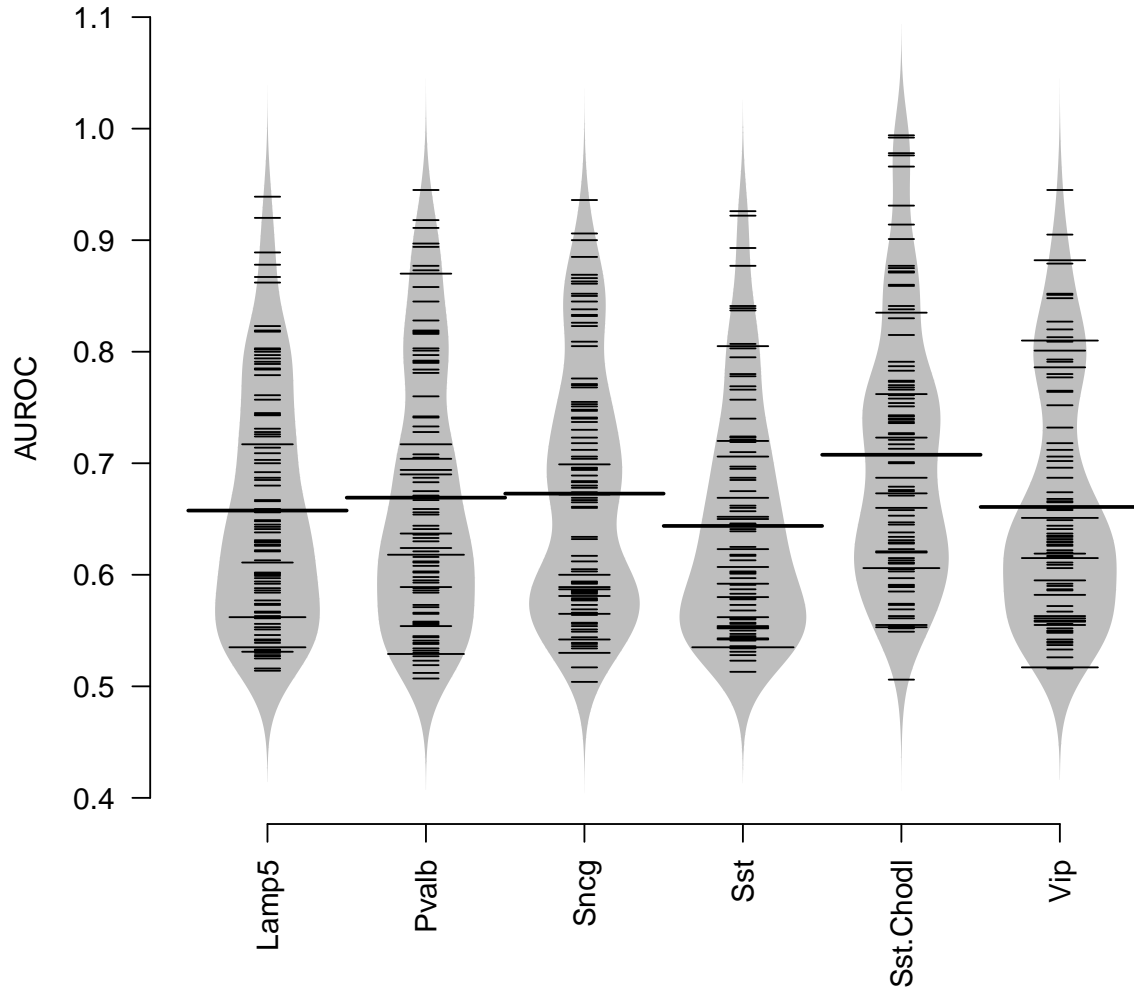
if (FALSE) {
write.table(aurocs, "functional_aurocs.txt")
}
```

Later, results can be retrieved with the “read.table” function:

```
aurocs = read.table("functional_aurocs.txt")
```

5. We use the “plotBPlot” function on the first 100 gene sets to rapidly visualize how replicability depends on gene sets.

```
plotBPlot(head(aurocs, 100))
```



We note that, on average, gene sets contribute moderately to replicability (AUROC ~ 0.7), numerous gene sets have a performance close to random (AUROC $\sim 0.5 - 0.6$) and some gene sets have exceedingly high performance (AUROC > 0.8).

6. To focus on gene sets that contribute highly to specificity, we create a summary table containing, for each gene set, cell type specific AUROCs, average AUROC across all cell types and gene set size.

```
gs_size = sapply(go_sets, length)
aurocs_df = as.data.frame(aurocs)
aurocs_df$average_auroc = rowMeans(aurocs)
aurocs_df$gs_size = gs_size[rownames(aurocs)]
```

We then order gene set by AUROC and look at top scoring gene sets:

```
head(aurocs_df[order(aurocs_df$average_auroc, decreasing = TRUE),], 10)
```

```
##                                                                 Lamp5
## G0:0007215|glutamate receptor signaling pathway|BP             0.967
## G0:0051966|regulation of synaptic transmission, glutamatergic|BP 0.960
```

## G0:0060076 excitatory synapse CC	0.960
## G0:0033555 multicellular organismal response to stress BP	0.952
## G0:0098839 postsynaptic density membrane CC	0.924
## G0:0099565 chemical synaptic transmission, postsynaptic BP	0.968
## G0:0008306 associative learning BP	0.967
## G0:0099601 regulation of neurotransmitter receptor activity BP	0.955
## G0:0060079 excitatory postsynaptic potential BP	0.967
## G0:0010771 negative regulation of cell morphogenesis involved in differentiation BP	0.979
##	Pvalb
## G0:0007215 glutamate receptor signaling pathway BP	0.985
## G0:0051966 regulation of synaptic transmission, glutamatergic BP	0.971
## G0:0060076 excitatory synapse CC	0.967
## G0:0033555 multicellular organismal response to stress BP	0.975
## G0:0098839 postsynaptic density membrane CC	0.970
## G0:0099565 chemical synaptic transmission, postsynaptic BP	0.977
## G0:0008306 associative learning BP	0.976
## G0:0099601 regulation of neurotransmitter receptor activity BP	0.975
## G0:0060079 excitatory postsynaptic potential BP	0.976
## G0:0010771 negative regulation of cell morphogenesis involved in differentiation BP	0.975
##	Sncg
## G0:0007215 glutamate receptor signaling pathway BP	0.971
## G0:0051966 regulation of synaptic transmission, glutamatergic BP	0.975
## G0:0060076 excitatory synapse CC	0.989
## G0:0033555 multicellular organismal response to stress BP	0.975
## G0:0098839 postsynaptic density membrane CC	0.983
## G0:0099565 chemical synaptic transmission, postsynaptic BP	0.973
## G0:0008306 associative learning BP	0.962
## G0:0099601 regulation of neurotransmitter receptor activity BP	0.955
## G0:0060079 excitatory postsynaptic potential BP	0.969
## G0:0010771 negative regulation of cell morphogenesis involved in differentiation BP	0.971
##	Sst
## G0:0007215 glutamate receptor signaling pathway BP	0.979
## G0:0051966 regulation of synaptic transmission, glutamatergic BP	0.965
## G0:0060076 excitatory synapse CC	0.964
## G0:0033555 multicellular organismal response to stress BP	0.947
## G0:0098839 postsynaptic density membrane CC	0.984
## G0:0099565 chemical synaptic transmission, postsynaptic BP	0.952
## G0:0008306 associative learning BP	0.959
## G0:0099601 regulation of neurotransmitter receptor activity BP	0.951
## G0:0060079 excitatory postsynaptic potential BP	0.951
## G0:0010771 negative regulation of cell morphogenesis involved in differentiation BP	0.959
##	Sst.Chodl
## G0:0007215 glutamate receptor signaling pathway BP	0.996
## G0:0051966 regulation of synaptic transmission, glutamatergic BP	0.992
## G0:0060076 excitatory synapse CC	0.987
## G0:0033555 multicellular organismal response to stress BP	0.999
## G0:0098839 postsynaptic density membrane CC	0.983
## G0:0099565 chemical synaptic transmission, postsynaptic BP	0.987
## G0:0008306 associative learning BP	0.990
## G0:0099601 regulation of neurotransmitter receptor activity BP	0.994
## G0:0060079 excitatory postsynaptic potential BP	0.986
## G0:0010771 negative regulation of cell morphogenesis involved in differentiation BP	0.986
##	Vip
## G0:0007215 glutamate receptor signaling pathway BP	0.986

## G0:0051966 regulation of synaptic transmission, glutamatergic BP	0.973
## G0:0060076 excitatory synapse CC	0.957
## G0:0033555 multicellular organismal response to stress BP	0.975
## G0:0098839 postsynaptic density membrane CC	0.970
## G0:0099565 chemical synaptic transmission, postsynaptic BP	0.955
## G0:0008306 associative learning BP	0.954
## G0:0099601 regulation of neurotransmitter receptor activity BP	0.976
## G0:0060079 excitatory postsynaptic potential BP	0.954
## G0:0010771 negative regulation of cell morphogenesis involved in differentiation BP	0.925
##	average_auroc
## G0:0007215 glutamate receptor signaling pathway BP	0.9806667
## G0:0051966 regulation of synaptic transmission, glutamatergic BP	0.9726667
## G0:0060076 excitatory synapse CC	0.9706667
## G0:0033555 multicellular organismal response to stress BP	0.9705000
## G0:0098839 postsynaptic density membrane CC	0.9690000
## G0:0099565 chemical synaptic transmission, postsynaptic BP	0.9686667
## G0:0008306 associative learning BP	0.9680000
## G0:0099601 regulation of neurotransmitter receptor activity BP	0.9676667
## G0:0060079 excitatory postsynaptic potential BP	0.9671667
## G0:0010771 negative regulation of cell morphogenesis involved in differentiation BP	0.9658333
##	gs_size
## G0:0007215 glutamate receptor signaling pathway BP	92
## G0:0051966 regulation of synaptic transmission, glutamatergic BP	75
## G0:0060076 excitatory synapse CC	75
## G0:0033555 multicellular organismal response to stress BP	98
## G0:0098839 postsynaptic density membrane CC	93
## G0:0099565 chemical synaptic transmission, postsynaptic BP	91
## G0:0008306 associative learning BP	100
## G0:0099601 regulation of neurotransmitter receptor activity BP	61
## G0:0060079 excitatory postsynaptic potential BP	83
## G0:0010771 negative regulation of cell morphogenesis involved in differentiation BP	98

Without surprise, replicability is mainly driven by gene sets related to neuronal functions that are immediately relevant to the physiology of inhibitory neurons, such as “glutamate receptor signaling pathway”, “regulation of synaptic transmission, glutamatergic”, or “chemical synaptic transmission, postsynaptic”. Note that most of the top scoring gene sets have a large number of genes. To obtain even more specific biological function, we can further filter gene sets that have fewer than 20 genes.

```
small_sets = aurocs_df[aurocs_df$gs_size < 20,]
head(small_sets[order(small_sets$average_auroc, decreasing = TRUE),],10)
```

##	Lamp5
## G0:0004970 ionotropic glutamate receptor activity MF	0.901
## G0:0035235 ionotropic glutamate receptor signaling pathway BP	0.823
## G0:0032230 positive regulation of synaptic transmission, GABAergic BP	0.844
## G0:0007216 G protein-coupled glutamate receptor signaling pathway BP	0.890
## G0:1905874 regulation of postsynaptic density organization BP	0.830
## G0:0099150 regulation of postsynaptic specialization assembly BP	0.833
## G0:0150052 regulation of postsynapse assembly BP	0.833
## G0:0021889 olfactory bulb interneuron differentiation BP	0.808
## G0:0070679 inositol 1,4,5 trisphosphate binding MF	0.921
## G0:1902711 GABA-A receptor complex CC	0.821
##	Pvalb
## G0:0004970 ionotropic glutamate receptor activity MF	0.917
## G0:0035235 ionotropic glutamate receptor signaling pathway BP	0.822

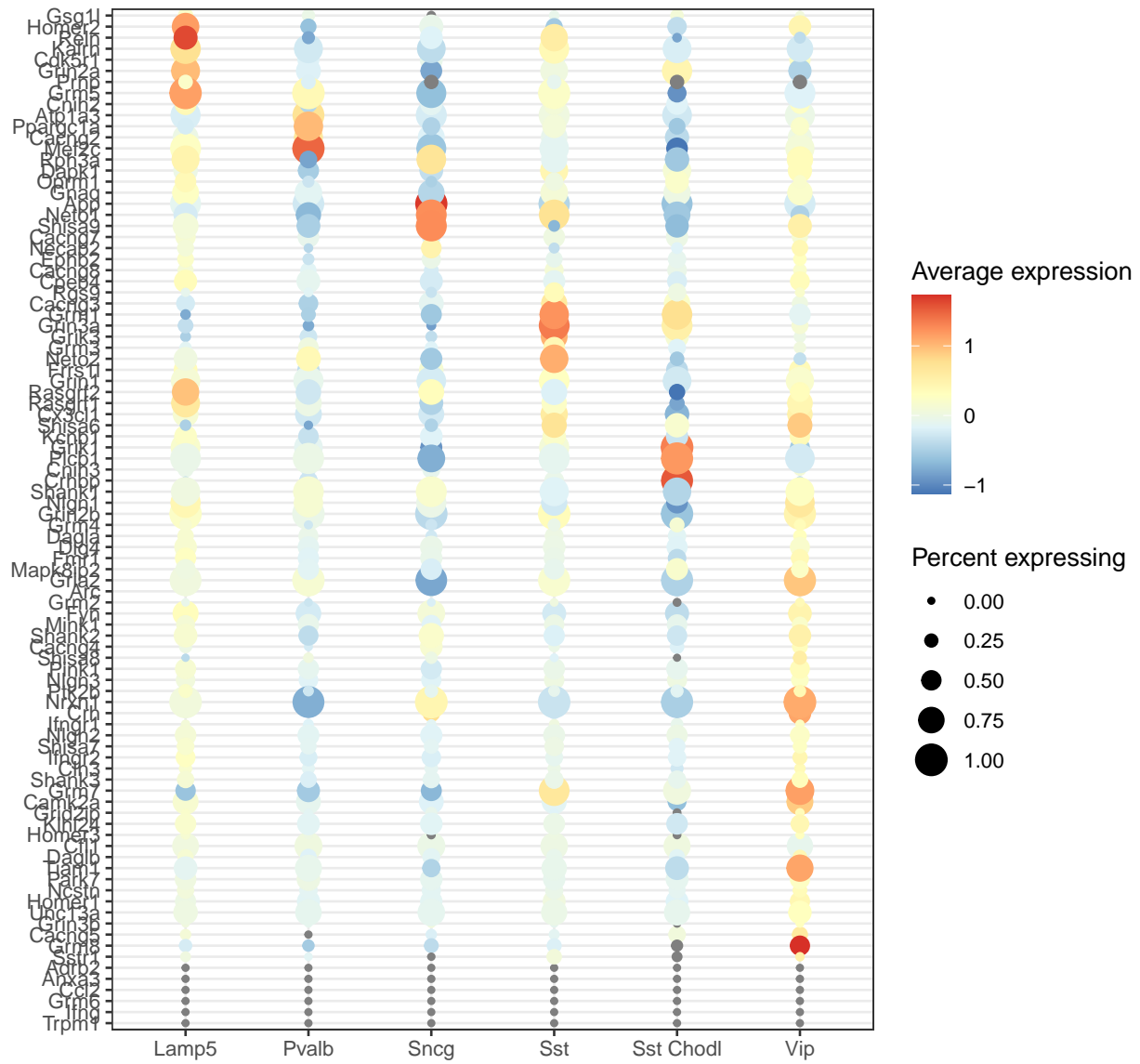
## G0:0032230 positive regulation of synaptic transmission, GABAergic BP	0.856
## G0:0007216 G protein-coupled glutamate receptor signaling pathway BP	0.849
## G0:1905874 regulation of postsynaptic density organization BP	0.862
## G0:0099150 regulation of postsynaptic specialization assembly BP	0.893
## G0:0150052 regulation of postsynapse assembly BP	0.893
## G0:0021889 olfactory bulb interneuron differentiation BP	0.914
## G0:0070679 inositol 1,4,5 trisphosphate binding MF	0.940
## G0:1902711 GABA-A receptor complex CC	0.869
##	Sncg
## G0:0004970 ionotropic glutamate receptor activity MF	0.914
## G0:0035235 ionotropic glutamate receptor signaling pathway BP	0.906
## G0:0032230 positive regulation of synaptic transmission, GABAergic BP	0.816
## G0:0007216 G protein-coupled glutamate receptor signaling pathway BP	0.764
## G0:1905874 regulation of postsynaptic density organization BP	0.866
## G0:0099150 regulation of postsynaptic specialization assembly BP	0.855
## G0:0150052 regulation of postsynapse assembly BP	0.855
## G0:0021889 olfactory bulb interneuron differentiation BP	0.824
## G0:0070679 inositol 1,4,5 trisphosphate binding MF	0.791
## G0:1902711 GABA-A receptor complex CC	0.870
##	Sst
## G0:0004970 ionotropic glutamate receptor activity MF	0.957
## G0:0035235 ionotropic glutamate receptor signaling pathway BP	0.929
## G0:0032230 positive regulation of synaptic transmission, GABAergic BP	0.920
## G0:0007216 G protein-coupled glutamate receptor signaling pathway BP	0.916
## G0:1905874 regulation of postsynaptic density organization BP	0.903
## G0:0099150 regulation of postsynaptic specialization assembly BP	0.910
## G0:0150052 regulation of postsynapse assembly BP	0.910
## G0:0021889 olfactory bulb interneuron differentiation BP	0.884
## G0:0070679 inositol 1,4,5 trisphosphate binding MF	0.806
## G0:1902711 GABA-A receptor complex CC	0.802
##	Sst.Chodl
## G0:0004970 ionotropic glutamate receptor activity MF	0.966
## G0:0035235 ionotropic glutamate receptor signaling pathway BP	0.943
## G0:0032230 positive regulation of synaptic transmission, GABAergic BP	0.985
## G0:0007216 G protein-coupled glutamate receptor signaling pathway BP	0.951
## G0:1905874 regulation of postsynaptic density organization BP	0.917
## G0:0099150 regulation of postsynaptic specialization assembly BP	0.910
## G0:0150052 regulation of postsynapse assembly BP	0.910
## G0:0021889 olfactory bulb interneuron differentiation BP	0.891
## G0:0070679 inositol 1,4,5 trisphosphate binding MF	0.863
## G0:1902711 GABA-A receptor complex CC	0.991
##	Vip
## G0:0004970 ionotropic glutamate receptor activity MF	0.923
## G0:0035235 ionotropic glutamate receptor signaling pathway BP	0.866
## G0:0032230 positive regulation of synaptic transmission, GABAergic BP	0.829
## G0:0007216 G protein-coupled glutamate receptor signaling pathway BP	0.842
## G0:1905874 regulation of postsynaptic density organization BP	0.826
## G0:0099150 regulation of postsynaptic specialization assembly BP	0.802
## G0:0150052 regulation of postsynapse assembly BP	0.802
## G0:0021889 olfactory bulb interneuron differentiation BP	0.864
## G0:0070679 inositol 1,4,5 trisphosphate binding MF	0.853
## G0:1902711 GABA-A receptor complex CC	0.796
##	average_auroc
## G0:0004970 ionotropic glutamate receptor activity MF	0.9296667

## G0:0035235 ionotropic glutamate receptor signaling pathway BP	0.8815000
## G0:0032230 positive regulation of synaptic transmission, GABAergic BP	0.8750000
## G0:0007216 G protein-coupled glutamate receptor signaling pathway BP	0.8686667
## G0:1905874 regulation of postsynaptic density organization BP	0.8673333
## G0:0099150 regulation of postsynaptic specialization assembly BP	0.8671667
## G0:0150052 regulation of postsynapse assembly BP	0.8671667
## G0:0021889 olfactory bulb interneuron differentiation BP	0.8641667
## G0:0070679 inositol 1,4,5 trisphosphate binding MF	0.8623333
## G0:1902711 GABA-A receptor complex CC	0.8581667
##	gs_size
## G0:0004970 ionotropic glutamate receptor activity MF	19
## G0:0035235 ionotropic glutamate receptor signaling pathway BP	16
## G0:0032230 positive regulation of synaptic transmission, GABAergic BP	16
## G0:0007216 G protein-coupled glutamate receptor signaling pathway BP	16
## G0:1905874 regulation of postsynaptic density organization BP	19
## G0:0099150 regulation of postsynaptic specialization assembly BP	18
## G0:0150052 regulation of postsynapse assembly BP	18
## G0:0021889 olfactory bulb interneuron differentiation BP	15
## G0:0070679 inositol 1,4,5 trisphosphate binding MF	15
## G0:1902711 GABA-A receptor complex CC	19

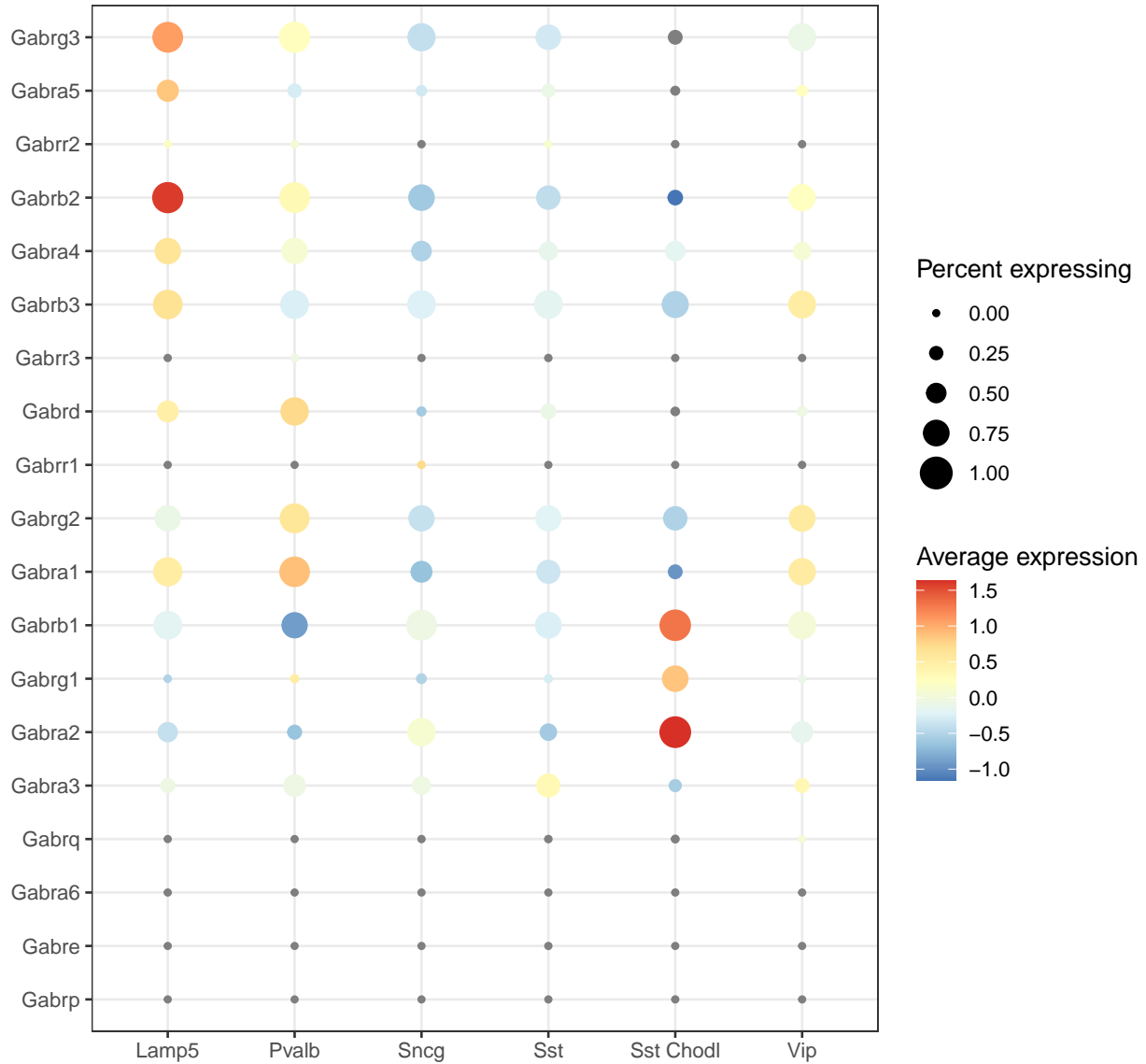
Again, the top scoring gene sets are dominated by biological functions immediately relevant to inhibitory neuron physiology, such as “ionotropic glutamate receptor signaling pathway”, “positive regulation of synaptic transmission, GABAergic”, or “GABA-A receptor complex”.

7. To understand how individual genes contribute to gene set performance, we use the “plotDotPlot” function, which will show the expression of all genes in a gene set of interest, averaged over all datasets.

```
plotDotPlot(dat = biccn_gaba,
            experiment_labels = biccn_gaba$study_id,
            celltype_labels = biccn_gaba$joint_subclass_label,
            gene_set = go_sets[["G0:0007215|glutamate receptor signaling pathway|BP"]])
```



```
plotDotPlot(dat = biccn_gaba,
  experiment_labels = biccn_gaba$study_id,
  celltype_labels = biccn_gaba$joint_subclass_label,
  gene_set = go_sets[["GO:1902711|GABA-A receptor complex|CC"]])
```

High scoring gene sets are characterized by the differential usage of genes from a given gene set. For example, when looking at the GABA-A receptor complex composition, Lamp5 will preferentially use the Gabrb2 and Gabrg3 receptors, Pvalb the Gabra1 receptor, and Sst Chodl the Gabra2, Gabrb1 and Gabrg1 receptors.

Once the overall replicability of clusters has been established with unsupervised MetaNeighbor (as in Protocols 1 and 2), supervised MetaNeighbor enables the functional interpretation of the biology contributing to each cell type's unique identity.