

Transcriptomic immaturity of the hippocampus and prefrontal cortex in patients with alcoholism

Tomoyuki Murano<sup>1,2,3</sup>, Hisatsugu Koshimizu<sup>3</sup>, Hideo Hagihara<sup>3</sup>, Tsuyoshi Miyakawa<sup>2,3</sup>

<sup>1</sup>Department of Physiological Science, School of Life Science, The Graduate University for Advanced Studies, Japan [SOKENDAI]

<sup>2</sup>Section of Behavior Patterns, Center for Genetic Analysis of Behavior, National Institute for Physiological Sciences, Japan

<sup>3</sup>Division of Systems Medical Science, Institute for Comprehensive Medical Science, Fujita Health University, Toyoake, Japan

\*Correspondence: [miyakawa@fujita-hu.ac.jp](mailto:miyakawa@fujita-hu.ac.jp)

## Supplementary Information

### Computing overlap p-value of gene expression patterns from different datasets

NextBio compares the signatures in publicly available microarray data sets with a signature provided by the user using a “Running Fisher” algorithm, as previously described<sup>1</sup>. The overlap  $P$  value, i.e., the direction of the correlation between two given gene signature sets ( $b1$ ,  $b2$ ), and the  $P$  values between subsets of gene signatures, is calculated as follows:

First, each gene signature set was rank-ordered according to the absolute fold-change value. Upregulated and downregulated genes were denoted by positive and negative signs, respectively, to imply directionality. A directional subset was generated for each direction, such as  $b1+$ ,  $b1-$ ,  $b2+$ , and  $b2-$ . Second, all of the subset pairs were identified as  $b1Di$ ,  $b2Dj$ , where  $Di$  and  $Dj$  were the available directions (+ or -) in  $b1$  and  $b2$ , respectively. The Running Fisher algorithm was applied to each subset pair. The top ranking genes in the first subset  $b1Di$  were collected as a group,  $G$ , and the second subset  $b2Dj$  was scanned from top to bottom in rank order to identify each rank with a gene matching a member in group  $G$ . At each matching rank,  $K$ , the scanned portion of the second subset  $b2Dj$  consisted of  $N$  genes, and the overlap between group  $G$  and  $N$  genes was defined as  $M$ . A Fisher’s exact test was performed at rank  $K$  to evaluate the statistical significance of observing  $M$  overlaps between a set of size  $G$  and a set of size  $N$ , where the set of size  $G$  comes from platform  $P1$  and the set of size  $N$  comes from platform  $P2$ , given the sizes of  $P1$  and  $P2$  as well as the overlap between  $P1$  and  $P2$ . At the end of the scan, the best  $P$  value was retained, and a multiple hypothesis testing correction factor was applied. The negative log of the multiple testing corrected best  $P$  value ( $P_{b1Di \rightarrow b2Dj}$ ) was a score ( $S_{b1Di \rightarrow b2Dj}$ ) for the subset pair. Here, the subscript of  $b1Di \rightarrow b2Dj$  indicates that  $b1Di$  was the first subset used to define the top genes  $G$  and  $b2Dj$  was the second subset that is used for the scan.

$$S_{b1Di \rightarrow b2Dj} = -\ln P_{b1Di \rightarrow b2Dj} \quad (1)$$

Next, the Running Fisher algorithm was performed in the reverse direction. The same procedure in this reverse direction produced another score ( $S_{b2Dj \rightarrow b1Di}$ ) for the same subset pair. The two scores were averaged to represent the magnitude of the similarity between the two subsets.

$$S_{b1Dib2Dj} = \frac{S_{b1Di \rightarrow b2Dj} + S_{b2Dj \rightarrow b1Di}}{2} \quad (2)$$

The  $P$  value ( $P_{b1Dib2Dj}$ ) between  $b1Di$  and  $b2Dj$  was calculated using the following equation:

$$P_{b1Dib2Dj} = \exp(-S_{b1Dib2Dj}) \quad (3)$$

A positive sign was assigned to pairwise correlation scores ( $S_{b1+b2+}$  and  $S_{b1-b2-}$ ) for a subset pair of the same direction ( $b1+b2+$ ,  $b1-b2-$ ), and a negative sign was assigned to pairwise correlation scores ( $S_{b1+b2-}$  and  $S_{b1-b2+}$ ) for a subset pair of opposite directions ( $b1+b2-$ ,  $b1-b2+$ ). Then, the overall score ( $S_{b1b2}$ ) between  $b1$  and  $b2$  was calculated from the correlation scores ( $S_{b1+b2+}$ ,  $S_{b1-b2-}$ ,  $S_{b1+b2-}$ , and  $S_{b1-b2+}$ ) of subset pairs using the following equation:

$$S_{b1b2} = \frac{S_{b1+b1+} + S_{b1-b2-}}{2} - \frac{S_{b1+b1-} + S_{b1-b2+}}{2} \quad (4)$$

The sign of  $S_{b1b2}$  determined whether the two signatures were positively or negatively correlated. The overall  $P$  value ( $P_{b1b2}$ ) between  $b1$  and  $b2$  was calculated using the following equation:

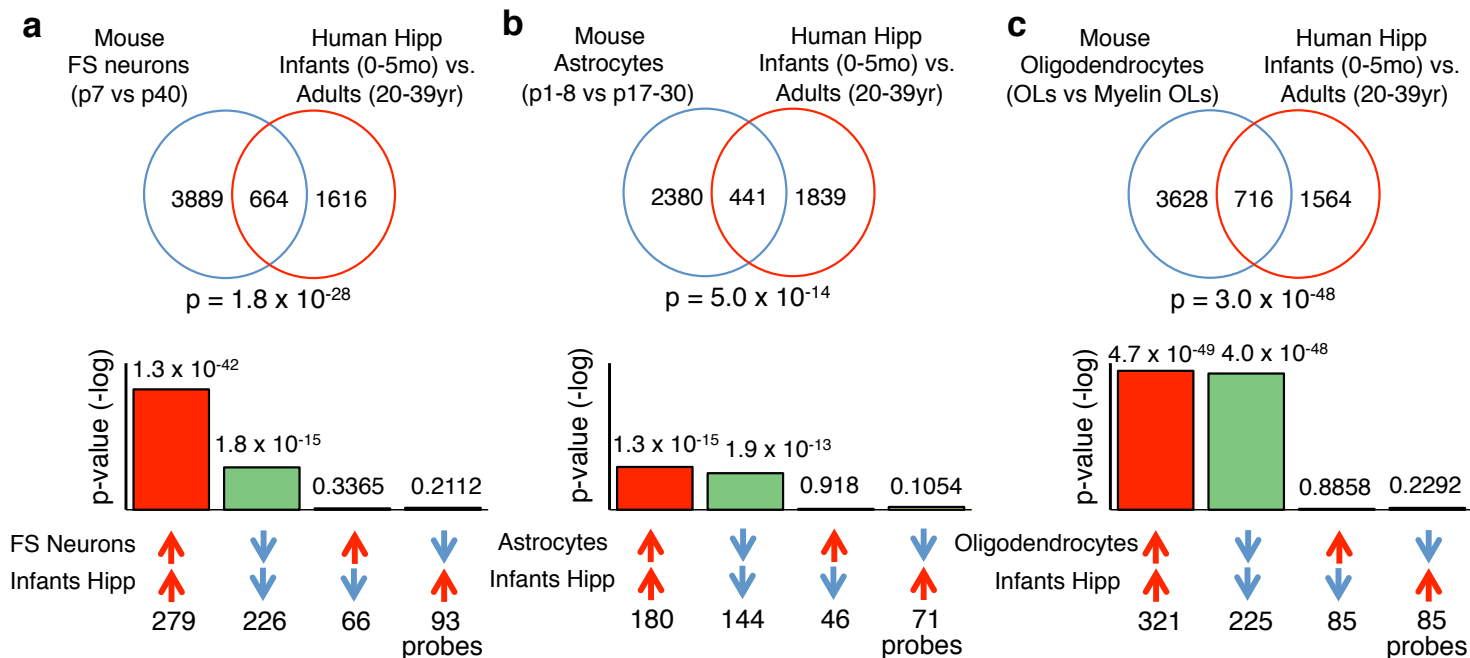
$$P_{b1b2} = \exp(-|S_{b1b2}|) \quad (5)$$

This overall  $P$  value was referred as an overlap  $P$  value between two gene expression patterns in this paper.

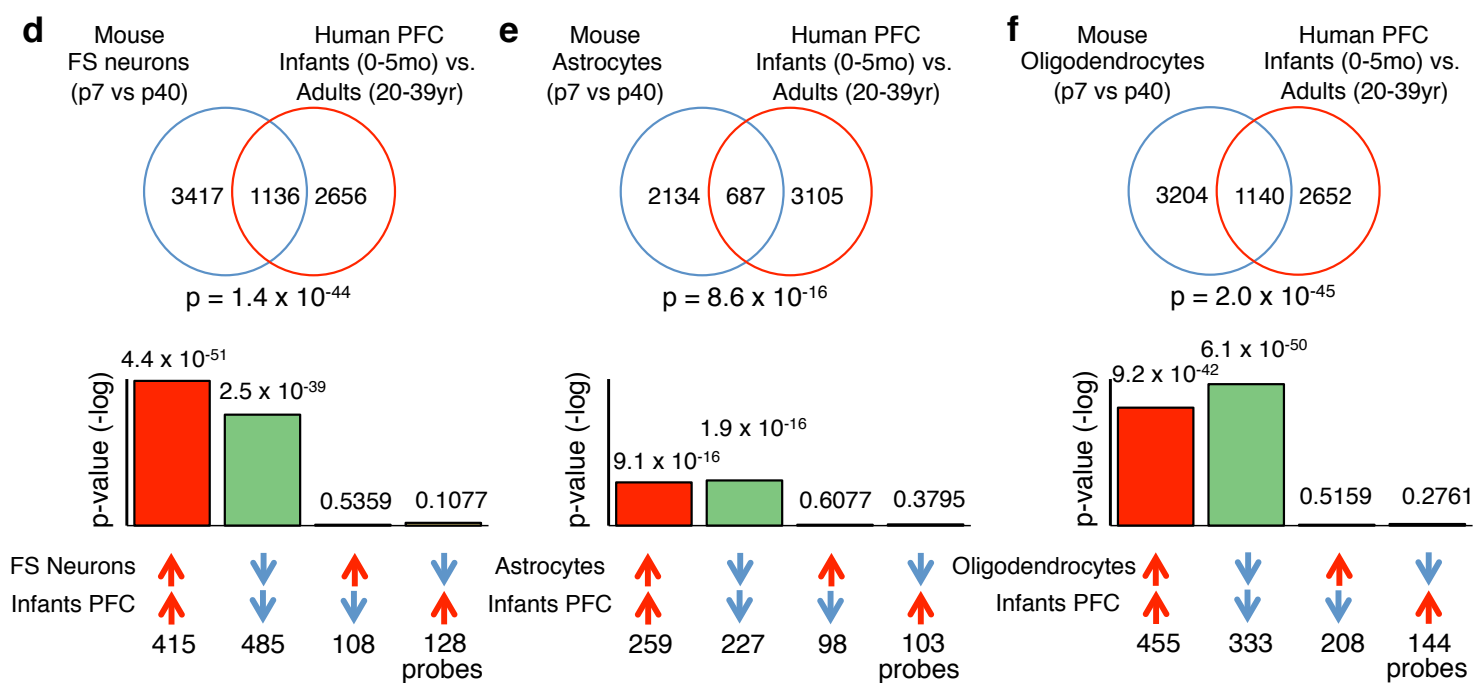
1. Kupersmidt, I. *et al.* Ontology-Based Meta-Analysis of Global Collections of

High-Throughput Public Data. *PLoS ONE* **5**, e13066 (2010).

## Hippocampus



## Prefrontal Cortex



**Supplementary Figure 1. Strong Correlation between datasets for mouse cell type-specific development and datasets for human hippocampi/PFCs development.**

Comparison between datasets from each cell type development (**a, d**, FS neurons [GSE17806]; **b, e**, astrocytes [GSE9566]; **c, f**, oligodendrocytes [GSE9566]) with datasets from human hippocampi development (**a, b, c**, [GSE44456]) and PFCs development (**d, e, f**, [GSE49376]). Venn diagrams illustrating the overlap in transcriptome-wide alterations in gene expression in the hippocampi/PFCs of male patients with cell type-specific gene expression. Bar graphs illustrate the overlapping *P*-value of genes up-regulated (red arrow) or down-regulated (blue arrows) by each condition, between the two conditions.