



OPEN

## Pseudotime analysis reveals novel regulatory factors for multigenic onset and monogenic transition of odorant receptor expression

Mohammad Hussainy<sup>1,2</sup>, Sigrun I. Korsching<sup>2,5</sup> & Achim Tresch<sup>1,3,4,5</sup>✉

During their maturation from horizontal basal stem cells, olfactory sensory neurons (OSNs) are known to select exactly one out of hundreds of olfactory receptors (ORs) and express it on their surface, a process called monogenic selection. Monogenic expression is preceded by a multigenic phase during which several OR genes are expressed in a single OSN. Here, we perform pseudotime analysis of a single cell RNA-Seq dataset of murine olfactory epithelium to precisely align the multigenic and monogenic expression phases with the cell types occurring during OSN differentiation. In combination with motif analysis of OR gene cluster-associated enhancer regions, we identify known and novel transcription (co-)factors (Ebf1, Lhx2, Ldb1, Fos and Ssbp2) and chromatin remodelers (Kdm1a, Eed and Zmynd8) associated with OR expression. The inferred temporal order of their activity suggests novel mechanisms contributing to multigenic OR expression and monogenic selection.

The sense of smell is tasked with the daunting challenge of making sense of potentially billions of different chemical stimuli to enable a multitude of different behaviors such as food search, prey hunting, predator evasion, mating and other social interactions<sup>1</sup>. This task is solved by several different receptor families, of which the best studied is both the evolutionary oldest and the largest (odorant receptor genes, OR genes)<sup>2</sup>. The analytical power of this system is maximal when information gathered from activation of individual receptors is kept separate at the peripheral level. Indeed for both vertebrates and insects, it has been shown that individual olfactory sensory neurons (OSNs) express only a single OR gene out of the entire olfactory receptor repertoire, which has been christened as monogenic expression<sup>3,4</sup>. Sensory neurons expressing the same receptor are distributed across the olfactory sensory surface, but their axons converge into a single target region in the olfactory bulb, the first relay station of olfactory information processing<sup>5–8</sup>. These target regions (so-called glomeruli) show a stereotyped arrangement, resulting in a receptotopic map on the olfactory bulb (or antennal lobe in the case of insects). Thus monogenic expression has a central importance for the olfactory coding logic. In fact, expression of ORs is even monoallelic, i.e. restricted to one allele of the OR selected in monogenic expression<sup>3,4,9,10</sup>. The molecular path towards monogenic and monoallelic expression is still not well understood, and the relative timing of these processes is not clear.

To reach monogenic and monoallelic expression presents a massive challenge in the case of very large gene families such as those of mouse and rat ORs, which both number well over one thousand intact genes<sup>2</sup>. A striking feature of genomic arrangement of OR genes is the occurrence in several clusters, which contain from a single to over one hundred different OR genes<sup>11</sup>. For mouse 68 such clusters have been identified, with the largest cluster containing 269 OR genes<sup>11</sup>. Another large cluster contains all 145 class I OR genes, which show a spatially restricted expression pattern in the olfactory epithelium<sup>12–14</sup>. These observations have prompted the search for cluster-specific regulatory elements. In a seminal publication, 63 genomic regions containing such elements were identified and named after Greek islands<sup>11,15</sup>. Forty-two class II OR clusters and the single class I OR cluster are associated with these Greek islands, which lie proximal to and sometimes even inside the clusters<sup>11</sup>. A common

<sup>1</sup>Institute of Medical Statistics and Computational Biology, Faculty of Medicine, University of Cologne, Cologne, Germany. <sup>2</sup>Institute of Genetics, Faculty of Mathematics and Natural Sciences, University of Cologne, Cologne, Germany. <sup>3</sup>Cologne Excellence Cluster On Cellular Stress Responses in Aging-Associated Diseases (CECAD), University of Cologne, Cologne, Germany. <sup>4</sup>Center for Data and Simulation Science, University of Cologne, Cologne, Germany. <sup>5</sup>These authors contributed equally: Sigrun I. Korsching and Achim Tresch. ✉email: achim.tresch@uni-koeln.de

feature of Greek islands is the presence of closely adjacent Lhx2 and Ebf1-binding motifs, which are also found individually in promoter regions of individual OR genes<sup>16–19</sup>.

Beyond individual cluster-associated regulatory elements the chromatin structure itself appears to play an essential role in regulating OR expression. OSNs possess a unique nuclear architecture compared to other cell types including the basal cells giving rise to the OSN lineage. In the silent phase before onset of expression OR genes are aggregated in constitutive heterochromatin and are associated with its molecular hallmarks, H3K9me3 and H4K20me3<sup>20,21</sup>. Onset of expression is concomitant with selective de-methylation (of H3K9me3), tri-methylation of H3K27 and re-location into expression-competent territory<sup>21–24</sup>. Moreover, of the two alleles of an active OR only one is found in the more plastic facultative heterochromatin<sup>22–24</sup>, i.e. amenable to expression, whereas the other remains blocked inside the constitutive heterochromatin, resulting in monoallelic expression. This suggests an involvement of chromatin remodelers in regulation of expression of OR genes. Furthermore, the stabilization of monogenic expression appears to require negative feedback from an active OR gene<sup>25,26</sup>, which may be mediated by silencing of the activating demethylase LSD1, synonym Kdm1a<sup>27,28</sup>. Recent progress in deep sequencing techniques has allowed to obtain high quality single cell transcriptomes (scRNA-Seq), resulting in the surprising observation that monogenic expression of ORs found in mature OSNs is preceded by a multigenic phase in immature OSN<sup>29,30</sup>.

It is so far mostly unclear how these OR expression phases align to the developmental stages of OSN differentiation. Furthermore, although the basic stages (stem cell, dividing precursor cell, immature neuron, mature neuron) were known previously<sup>29,30</sup>, deep sequencing techniques allow an unbiased ordering of individual cells along pseudotime trajectories according to their entire transcriptome. This enables a more precise and more stringent categorization of developmental stages compared to previous attempts.

Here, we re-analyzed a scRNA-Seq dataset obtained by Fletcher et al.<sup>31</sup> with the goal to precisely determine the timing of multigenic and monogenic expression during OSN differentiation. A combination of sequence binding motif and time series analysis then identifies novel regulatory components involved in establishing OR gene expression patterns. We ascertain the transcription (co-)factors and chromatin remodelers that are specifically correlated with the onset of multigenic and of monogenic expression (e.g., Fos, Ssbp2, Eed and Zmynd8). Finally, we suggest potential mechanisms for multigenic and monogenic selection.

## Results

### Re-analysis of a single cell RNA-Seq data set reveals four lineages originating from quiescent globose basal cells.

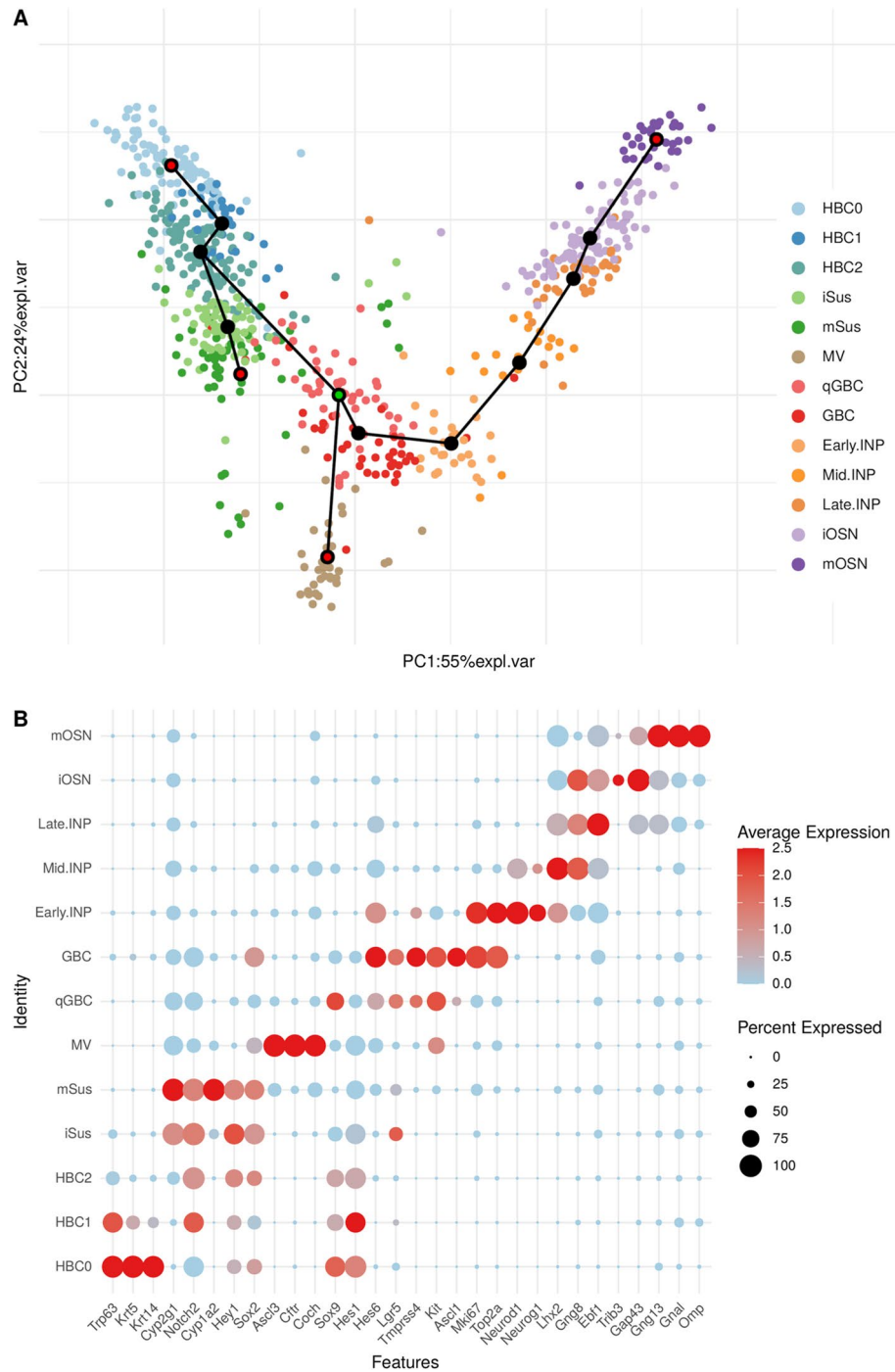
We have searched the GEO and ENA repositories for murine and human single-cell RNA-seq data sets that contain cells from the olfactory neuronal lineage (see STable 1 for a comprehensive list). These datasets were either generated by the SMART-Seq or the 10× Genomics technology. While the latter typically allows the analysis of many cells, the former provides a higher number of unique mapped sequence reads per cell<sup>32</sup>. For our purpose, it is most important to detect the sporadic (and often low level) of OR genes with sufficient sensitivity. Therefore, we selected the SMART-seq dataset by Fletcher et al.<sup>31</sup> for our re-analysis. This dataset contains the most significant number of neuronal lineage cells as our primary source. For validation, we compare our results to the 10× Genomics dataset richest in OSNs<sup>33</sup>.

We re-analyzed a scRNA-Seq dataset obtained from Fletcher et al.<sup>31</sup>. After quality filtering and pre-processing data from Fletcher et al.<sup>31</sup> (Methods, SFigure 1), 687 cells were included into further analysis and grouped into 13 clusters using Seurat KNN clustering on the top 15 principal components (Methods). Dimension reduction and visualization was performed using principal components analysis (Fig. 1A) and tSNE / UMAP. Using an extensive set of known marker genes, we assigned clusters to cell types of the main olfactory epithelium (MOE) (Fig. 1B, Methods SFigure 2a and STable 2). We detected all cell types described by Fletcher et al.<sup>31</sup>, and additionally we could subdivide the globose basal stem cell cluster into quiescent cells (qGBC) and active cells (GBC). Active GBC were identified by their expression of *Ascl1/Mash1*<sup>34–37</sup> and by the presence of cell cycle genes such as *Mki67* and *Top2a*. GBC are known as the adult OSN stem cells responsible for a sustained self-renewal of the OSNs throughout life<sup>37</sup>.

Further, trajectory inference by Slingshot<sup>38</sup> (Methods) revealed a tree with four leaves (Fig. 1A). Slingshot does not provide information on the direction of development but merely a tree topology. We chose qGBC as a root node for pseudotime analysis because qGBC has been reported as a general stem cell population in the MOE following injury<sup>37</sup>. This results in four lineages starting from qGBC (Fig. 1A):

- (1) The basal stem cells lineage, which connects the quiescent horizontal basal stem cells (HBC0, represented by 90 cells in the data) via two transient populations of horizontal basal stem cells (HBC1, 32 cells, and HBC2, 115 cells) with the qGBC (50 cells).
- (2) The supporting cells lineage ranges from mature sustentacular cells (mSus, 44 cells) to qGBCs, and includes immature sustentacular cells (iSus, 75 cells) and HBC2.
- (3) The microvillous cells lineage contains merely microvillous cells (MV, 36 cells). No transient cell types have been detected in this trajectory.
- (4) The neuronal lineage ends with mature olfactory sensory neurons (mOSN, 32 cells). It spans a range of several stages, namely GBC (38 cells), three intermediate neuronal precursors Early.INP (26 cells), Mid.INP (20 cells), Late.INP (34 cells) and immature olfactory sensory neurons (iOSN, 95 cells).

A previous analysis by Fletcher et al.<sup>31</sup> did not recognize qGBC as a separate population and set their root node as HBC0, which is a leaf node of the tree, and therefore leads to merely three lineages. Apart from this, their reconstruction is essentially identical to ours. Notably, both reconstructions agree on the neuronal lineage, i.e.,



**Figure 1.** Cell type identification, trajectory inference and pseudotime assignment. **(A)** 2D PCA projection of MOE cells shows the four predicted lineages starting from qGBC (center, pink, the starting point of the four trajectories is marked in green). The differentiation trajectory end points (red dots at the end of paths) are HBC0 (top left, pale blue), mSus (center left, green), MV (bottom, brown) and mOSN (top right, purple). The transition stages in the four lineages (black dots) are HBC1 (light blue), HBC2 (blue), iSus (light green), GBC (red), early.INP (pale orange), mid.INP (light orange), late.INP (dark orange) and iOSN (light purple) from the left to the right. **(B)** DotPlot representation of the average expression (red for high and pale blue for low expression) of known marker genes (x.axis) corresponding to cell types of MOE (y.axis), and the size of each dot represents the percentage of cells that expresses a corresponding marker gene in a given cell type.

the sequence of cell types qGBCs - GBC - Early.INP - Mid.INP - Late.INP - iOSN - mOSN. In the following, we restrict our analysis to the neuronal lineage.

For comparison purposes, we selected the main cell types of the olfactory epithelium from a dataset by Wang et al.<sup>33</sup>. We kept more than 35,000 cells annotated as HBC, SUS, MV, GBC, INP, iOSN and OSN (see SFigure 3a). The median number of genes per cell is 2448 (compared to 4164 in Fletcher) and the median number of counts per cell is 5694 (460,561 in Fletcher). Cells from the neuronal lineage (branch 4) were re-clustered (Methods, see SFigure 3b for a UMAP plot of the cells and their annotation), which lead to the cell type assignment used in the following.

**OR gene expression is limited to the last three stages of OSN differentiation: Sudden onset of multigenic OR expression in Late.INP is followed by transition to monogenic expression in immature OSN stage.**

Expectedly, OR expression is essentially unique to the neuronal lineage (Fig. 2 and data not shown). Next we used Slingshot to assign a pseudotime to each cell, thereby providing a linear order of all 295 cells in the neuronal lineage from qGBC to the terminal cell cluster (Methods). Our analysis could detect 157 cells of neuronal lineage that express at least one OR gene at relevant levels, i.e.  $\geq 50$  normalized counts. 132 of them belonged to the last three stages of OSN differentiation. Most of Late.INP cells (28 of 34), iOSN (76 of 95) and mOSN (28 of 32) express at least one OR. Many cells express more than one OR (multigenic expression), in particular in the Late.INP stage (26 of the 28 cells expressing OR genes, 92.8%). The frequency of multigenic expression drops sharply in later stages, 42% and 32% for iOSN and mOSN, respectively. We found 212 different OR genes that were expressed at least once in a single cell of the neuronal lineage (STable 3). Figure 2A shows the total number of reads that were assigned to OR genes, separately for each cell. While aggregate OR expression levels are almost zero for qGBC/GBC, early and mid INPs, there is a steep onset of OR expression in Late.INP. Then, overall OR expression stays at similarly high levels in iOSN and mOSN (Fig. 2A). A few OSN do not appear to express any OR at a relevant level. More precisely, 19 out 95 iOSN cells (20%) and 4 out 32 mOSN (12.5%) have less than 50 normalized OR counts. While it cannot be excluded that this is caused by incomplete annotation of the OR repertoire, or the exclusion of reads due to multiple mapping to closely similar OR genes, it is also possible that some OSN can not solve the challenging task of selecting an OR gene.

It is known that mature OSNs express only a single OR gene<sup>3,4,9,39</sup>, after a transient period of multigene expression<sup>29,30</sup>. We therefore decided to rank OR genes by expression level, separately in each cell. We then investigated the temporal behavior of the top four ranked OR genes in each cell. These genes account for 99% of all reads (413,388 out of 416,033) mapped to OR genes in cells from the neuronal lineage. The top-ranked gene of each stage will be referred to as ‘winner’ and the others as the ‘runners-up’. While the abundance of all runners-up drops sharply after Late.INP, the winner does not drop and in fact absolute levels keep increasing several fold until the mOSN stage (Fig. 2C,D). As a consequence the distance between winner and runners-up increases considerably in the iOSN stage and even more so in the mature neurons (mOSN). Since we observe each cell only once, we cannot be sure that the winner gene observed in one cell at a certain stage will be the highest expressed OR gene when that cell matures. However, this is the most plausible explanation, since a rank switch between the winner and a runner-up before/during the iOSN stage would require a coordinated switch of expression between these two specific OR genes, from high to almost zero and vice versa, an unlikely scenario.

Taken together, the pseudotime analysis of neuronal lineage cells suggests three main phases for OR expression (Fig. 2C,D):

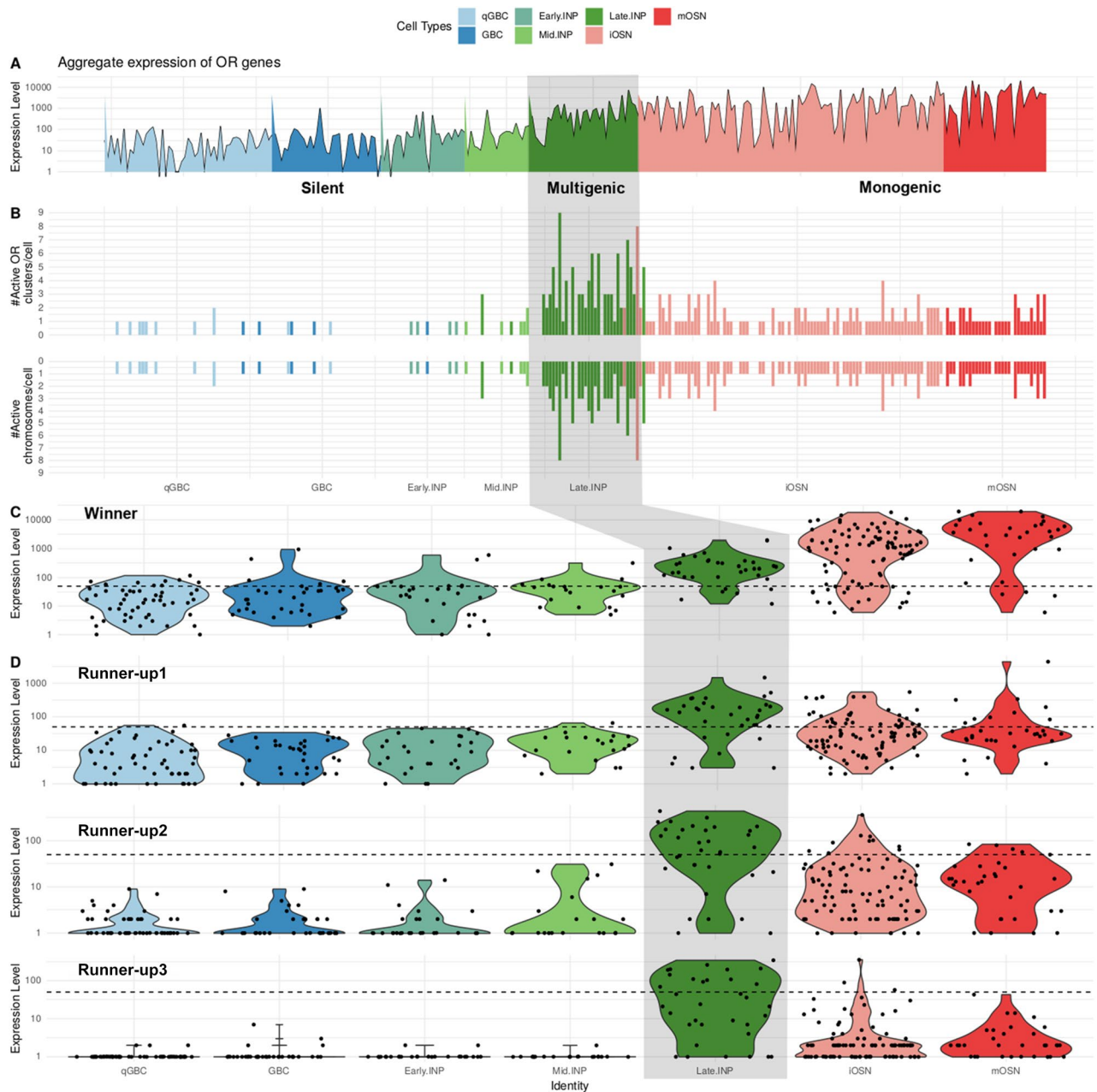
- (1) The silent phase exhibits virtually no OR expression, which is the case for the four early stages of neuronal lineage (qGBC, GBC, Early.INP and Mid.INP). In the present data, the silent phase is represented by 134 cells.
- (2) The onset of OR expression (multigenic phase) is characterized by the simultaneous expression of several OR genes per cell at relatively similar levels; this phase is represented by 34 cells and contains specifically the Late.INP stage and the beginning of the iOSN stage.
- (3) Finally, the monogenic phase includes the end of the iOSN stage and the mOSN stage, where each cell expresses essentially one functional OR allele (henceforth called the “winner”, Fig. 2C) at a very high level while the remaining ORs (the “runners-up”, Fig. 2D) show no or very low expression. Our data contains 127 cells in this phase.

Our analysis reveals that Late.INP is a crucial stage in stochastic selection.

We performed the same analysis for the Wang data. It turns out that the coverage per cell is not enough to reveal the characteristic expression time course of the runners-up during the multigenic phase (SFigure 4a), which is the reason why we chose Smart-seq data as a primary source.

**OR gene expression during multigenic phase shows no sign of OR cluster-specific activation.**

Next, in an attempt to infer the mechanism of activation in the multigenic stage and transition to monogenic stage, we analyzed the joint OR gene expression per cell (Fig. 2B, STable 4). We found that each of the cells expressing more than one OR gene in the Late.INP stage had at least two active OR clusters (i.e., clusters with an OR gene expressed at a level of at least 50 counts). The number of active clusters reached up to 9 for some cells. We performed a permutation test to assess whether OR genes that are jointly active in one cell have the tendency to be located in the same cluster (see SCode 1 and the description therein). We calculated, across all cells, the average frequency of clusters with more than one active OR gene. The null hypothesis we challenge is that this number is due to chance. We calculate the probability of observing this or a higher number if one randomly shuffles the genes  $\times$  cells matrix such that the marginal frequency of active genes per cell as well as the marginal frequency of cells expressing a given gene stays constant. The results however show that the average number of clusters with more than 1 active OR gene is consistent with the null model of random cluster allocation of



**Figure 2.** OR expression dynamics along neuronal lineage. (A) Aggregate expression (normalized pseudocounts) of all OR genes for each cell. Cells are sorted according to pseudotime and colored according to cell type. (B) Number of active OR clusters and chromosomes per cell, sorted by pseudotime. Since active clusters can be located on the same chromosome, the number of active chromosomes is less or equal to the number of active clusters per cell (e.g., in Late.INP all cells with 2 active clusters have also 2 active chromosomes, but in iOSN and mOSN cells with a single active chromosome may have 2 active clusters). Note that three cells classified as iOSN according to our marker genes are placed slightly before some cells classified as Late.INP in the pseudotime ordering. (C) Expression of the OR gene with the highest expression level in each cell (“winner”). For each stage of the neuronal lineage, we show the distribution of the corresponding expression values (pseudocounts) as a violin plot. (D) Using the same representation as in C, the expression of the OR gene with second, third and fourth highest expression (“runner-up” 1–3) in each cell is shown in the top, middle, and bottom row, respectively. Note that the y-axis is in log scale, and that the scale of the winner expression is on average one to two orders of magnitude higher than that of the runners-up. The dotted horizontal line in (C) and (D) marks an expression level of 50 normalized counts.

OR genes (SFigure 5,  $p=0.112$ ). Moreover, we found that the top two active OR clusters (ranked by expression level) always belonged to different chromosomes in each cell of the Late.INP stage (Fig. 2B). Thus the onset of

OR gene expression in the multigenic stage cannot be caused by activation of an individual chromosome or a particular OR cluster. Conversely, the transition to the monogenic stage could be partially caused by the restriction of expression to a single chromosome and cluster. However, this may not be the only selection mechanism involved, as some cells express up to 3–4 different OR genes simultaneously within a single cluster. Specifically, for the 34 Late.INP cells we found 20 cells which expressed at least 2 OR genes from the same cluster (expression defined as  $\geq 50$  normalized counts). Thus, additional steps are required to restrict expression to a single gene within a cluster.

**Motif search in Greek island enhancers identifies novel transcription factors.** Our next goal was to identify DNA-binding proteins potentially involved in the onset of multiple OR gene expression or the transition to monogenic expression. Such candidates should have a characteristic temporal gene expression pattern along the neuronal lineage. We thus require such candidates in our scRNA-seq dataset to be expressed at detectable levels. A factor is considered detectable if it has at least 35 counts in at least 15 out of 295 cells in the neuronal lineage, leaving us with 1358 (co)TFs. A characteristic expression pattern, however, might merely be the consequence and not the cause of the OR selection and differentiation process. Therefore, to make our search more specific for factors causally involved in OR selection, we confine our search to factors whose binding motif is enriched in enhancer regions of OR gene clusters.

Previously, 63 intergenic enhancer regions, termed Greek islands, have been identified inside or near OR clusters using DNase I hypersensitivity-sequencing and chromatin immunoprecipitation sequencing<sup>15</sup> and ATAC-seq<sup>11</sup>. SFigure 6 shows the co-localization of the Greek islands and the OR clusters on a map of the murine genome. Chromatin conformation capture experiments have revealed that Greek islands extensively contact OR clusters, remarkably both in cis and trans<sup>40,41</sup>.

We performed a de novo motif search on all Greek island enhancer regions as annotated by<sup>11</sup> using MEME<sup>42,43</sup> (Methods). Ungapped motif analysis of Greek islands identified one known motif, TYCCYWKGGVCTHATTARM (reported in Monahan et al.<sup>11</sup>), and two novel motifs GVDHCYCACGRGAV and TBYTCHTCTCYMC-AGDGWBNY, with E-value  $1.7e^{-057}$ ,  $3.7e^{-027}$  and  $4.1e^{-008}$ , respectively. Almost all Greek islands contain each of these motifs exactly once, except 8 Greek islands which are missing the third motif. TOMTOM was employed to align these motifs with known transcription factor motifs from the JASPAR database<sup>44</sup> (Methods). TOMTOM did not predict any significant TF binding for the two novel motifs, therefore we do not discuss them further. We found 65 significant target binding sites for transcription factors inside the first motif (see STable 5). From those, 9 TFs were expressed at a detectable level.

The most significant motif, TYCCYWKGGVCTHATTARM is composed of two adjacent submotifs, which are overlapped by a third submotif (Fig. 3A). The first submotif is targeted by the COE1 DNA-binding domain which is found exclusively in the Ebf transcription factor family (Ebf1–4). The second submotif is bound by homeodomain TFs such as Lhx2, Emx2 and Uncx (Fig. 3A). These results are consistent with previously reported Ebf1 and Lhx2 motifs to be positioned next to each other in most Greek islands<sup>11</sup>. Furthermore, the second submotif is expected to interact with transcription factors from three other families, the homeobox domain TF family, the Pou TF family (Pou6f1, this family has a strong enrichment in OR genes) and the ARID (AT-Rich Interaction Domain) domain TF family (Arid3a).

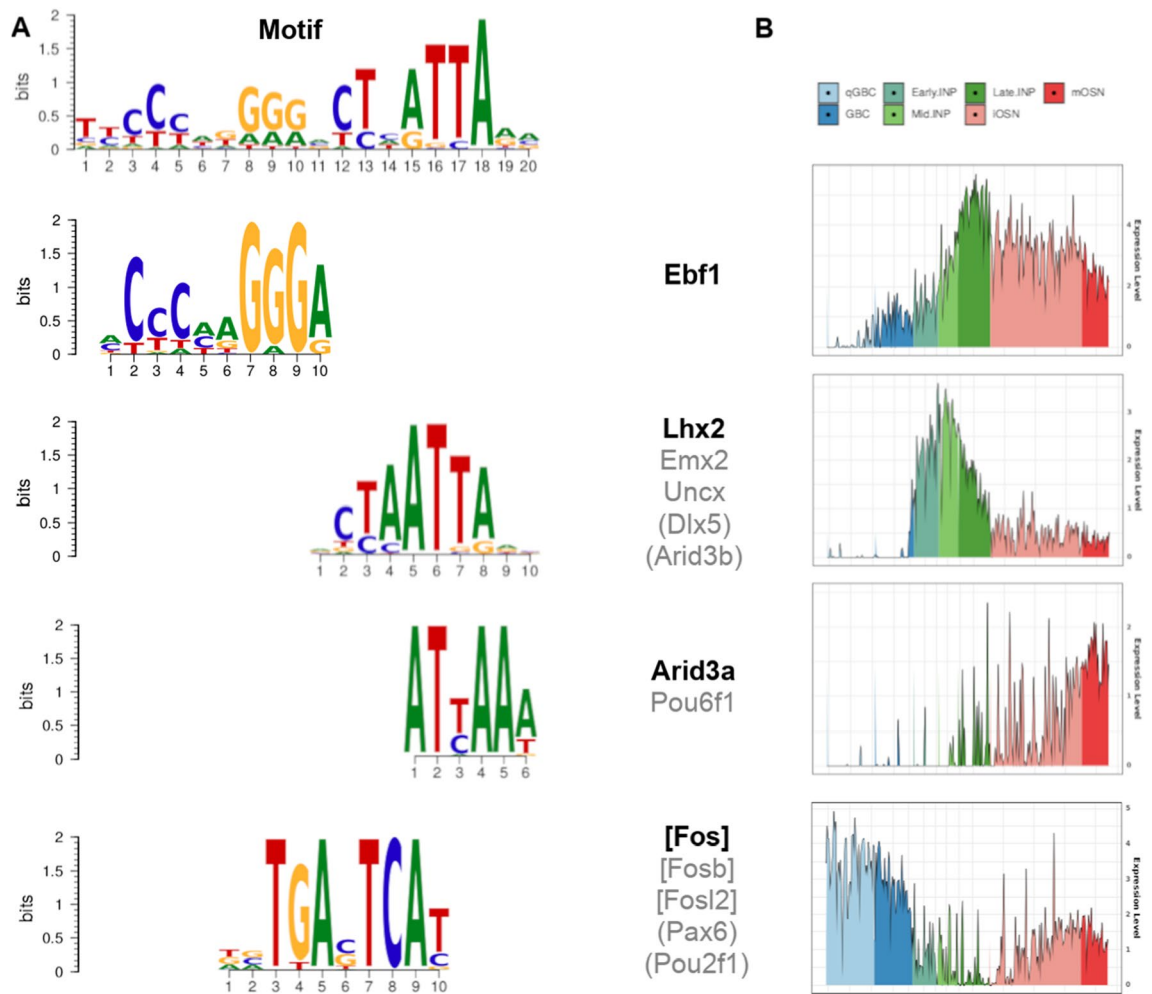
Noteworthy, our analysis predicted a third submotif, which is a possible binding site for Fos, Fosb and Fosl2. Fos and Fosb are well-known early response transcription factors, which in turn regulate a broad variety of other transcription factors thus regulating many physiological processes<sup>45–47</sup>. This Fos-binding motif overlaps with the end of the first and the beginning of the second submotif, suggesting a cooperative/inhibitory interaction of the respective binding factors. This possibility will be investigated further in the context of the pseudotime analysis.

Complementary to the de novo motif analysis, we did a forward motif search with AME<sup>48</sup> looking for known TF binding sites enriched in the 63 Greek Islands (Methods). This returned a total of 120 TFs (STable 6), of which 10 had a detectable expression (at least 35 count per cell in 15/295 cells of neuronal lineage) in our scRNA-seq dataset (SFigure 7). Six of those are also detected in the de novo motif search (see above), the additional TFs are Pax6, Dlx5, Pou2f1 and Arid3b (Fig. 3). Dlx5 is part of the same TF family as Lhx2, which was found in the de novo motif search. Pax6 and Pou2f1 have a homeobox domain, whereas Arid3b is part of the same TF family as Arid3a.

Taken together we describe 13 TFs that are found at detectable levels and predicted to bind to Greek islands by de novo and/or forward motif search (SFigure 7). Next we determined the pseudotime profiles of these TFs and found clear and distinct temporal expression patterns, providing additional evidence for their active involvement in the regulation of OR expression.

**Pseudotime analysis suggests transcription factors involved in OR expression.** The sorting of cells according to pseudotime (Methods) generates, for each gene, a time course of its expression (see above). Notably, all TFs found by motif search in the previous paragraph show a pronounced temporal expression pattern, which belongs to one of three groups (Fig. 3B and SFigure 7): The first group is active early in the silent phase, but strongly downregulated in late silent phase to reach a minimum in the multigenic phase (Fos, Fosb, Fosl2, Pax6 and Pou2f1). Some, but not all, are upregulated again in the monogenic phase (Fos, Fosb). The second group peaks within the multigenic phase (Ebf1, Lhx2, Emx2, Uncx and Dlx5). The third group is specifically upregulated during the monogenic phase (Arid3a and Pou6f1). Hereafter we will refer to these group definitions.

The fact that all TFs with a known Greek island binding site show a clear temporal pattern prompted us to perform a systematic search for TFs that change their expression upon transition between the three phases of OR expression. We also include co-factors in this analysis, because co-factors such as LDB1 have been found to be selectively associated with Greek islands and were suggested to initiate OR expression<sup>40</sup>. We searched for



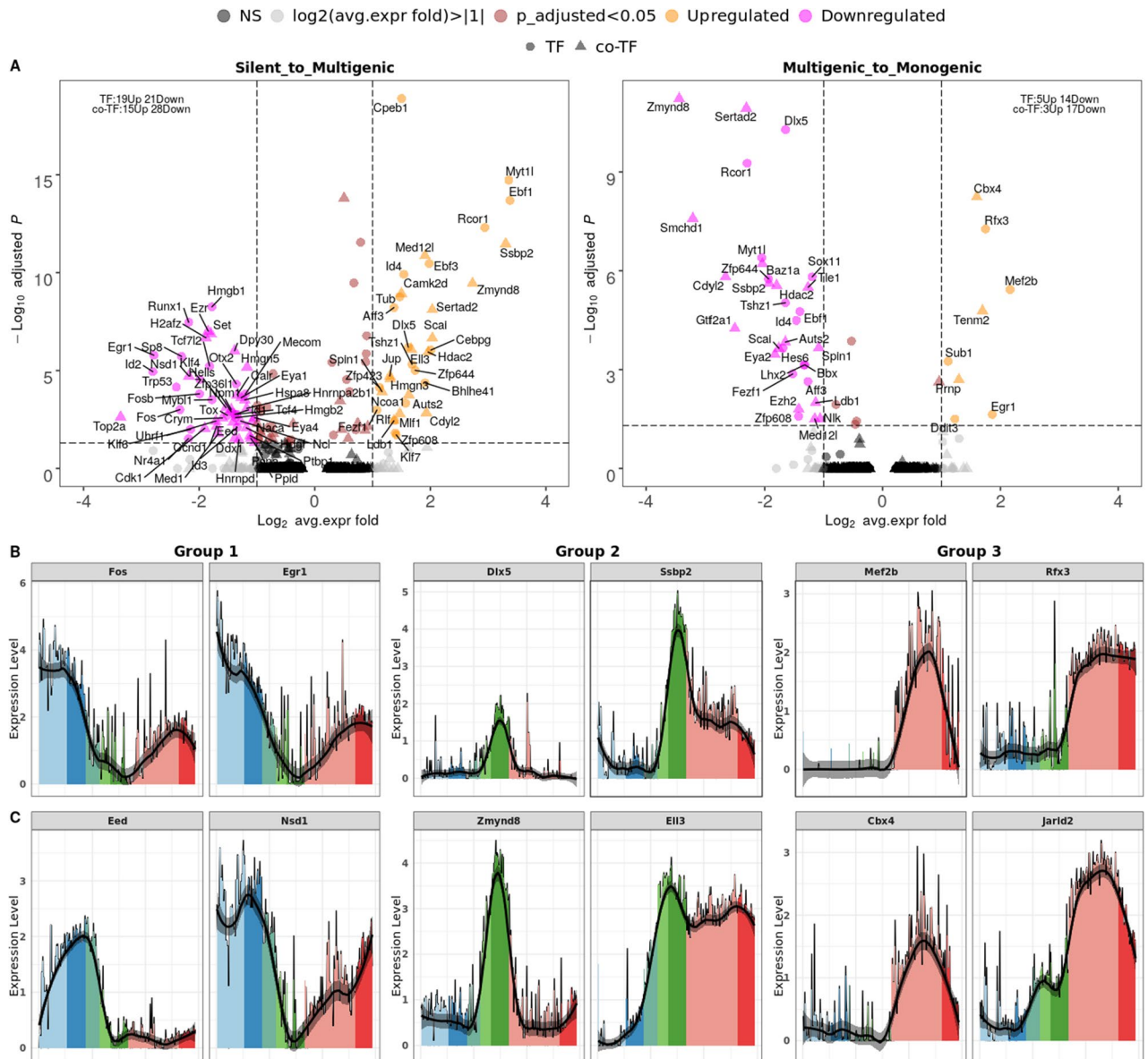
**Figure 3.** Greek islands binding motifs and representative pseudotime courses of respective TFs. **(A)** Binding motifs found in Greek islands. The top row shows the motif found de novo. In the rows below, known binding sites of transcription factors that partly align with the de novo motif in different sites (TOMTOM) and/or are enriched in Greek islands (AME). The transcription factors given in square brackets refer to the TFs found only by Tomtom alignment whereas the round brackets refer to TFs found only by forward motif search (see SFigure 6). **(B)** Single cell pseudotime courses of the TFs in bold print show characteristic trends along OSN differentiation (y-axis is a natural log of normalized counts). The grouping of the TFs in gray follows their pseudotime profile. It agrees with their grouping according to binding sites, except for Pax6 and Pou2f1 that have a homeobox domain like Lhx2.

(co-)TFs with a significant differential gene expression of at least twofold between silent vs. multigenic phase or between multigenic vs. monogenic phase (Methods), resulting in 83 (34 going up, 49 down) respectively 39 (8 up, 31 down) relevant hits (Fig. 4A, see STable 7 and STable 8).

Several of these differentially expressed factors were also identified by motif analysis (e.g., Ebf1, Lhx2, Dlx5, Fosb and Fos), but many are not. We manually inspected the pseudotime patterns of all differentially expressed (co-)TFs, and for detailed discussion, we selected those factors whose pseudotime expression pattern falls clearly into one of the three groups described above. We limit ourselves in the following to discussion of novel (co-) factors with additional evidence from motif analysis or a previous link to OR expression<sup>11,18,19,40</sup>. All other factors with a characteristic expression time course are shown in SFigure 8A.

For group 1 (active early in silent phase, downregulated in late silent phase and minimum in the multigenic phase) the differential expression search newly identified Egr1, whose expression resembles that of Fos and Fosb, which were already identified by motif analysis (see above). Therefore we searched explicitly for the known Egr1 binding motif (Methods) in Greek islands, which could be identified in 16 of 63 Greek islands. Fos, Fosb and Egr1 are immediate early genes, which are rapidly upregulated in response to external stimuli, immune response, and cellular stress<sup>47</sup>. Egr1, Fos and Fosb are specifically downregulated during the multigenic phase (Fig. 4B and SFigure 8A). This suggests combinatorial interactions with the other components that regulate OR expression and will be discussed later.

The second group peaks specifically within the multigenic phase and 6 factors have been identified by motif analysis (see above). Differential expression analysis further obtains the cofactor Ssbp2 (Fig. 4B) and three factors, Cebpγ, Rcor1 and Ldb1, which have been reported previously to be involved in OR expression<sup>11,40,49,50</sup>. Ssbp2



**Figure 4.** Differential expression analysis and pseudotime expression profiles of candidate factors. **(A)** Volcano plots revealing the up- (peach) and downregulated (magenta) factors in the two transitions from silent to multigenic (left) and from multigenic to monogenic phase (right). TFs are shown as filled circles, cofactors are triangles. We selected (co-)TFs\* with an adjusted p-value less than 0.05 (y-axis,  $\log_{10}$  Bonferroni adjusted p-value) and an average expression change of at least twofold (x-axis,  $\log_2$  fold). (Co-)TFs with only significant changes are shown in brown, those with only relevant expression fold > 2 are colored in grey, all others are in black (not significant). **(B)** Single cell pseudotime courses of selected transcription factors and cofactors that show characteristic trends along OSN differentiation (cell stages are indicated by colors). Left to right: each group of 2 columns shows examples of transcription (co-)factors in group 1, 2, and 3 as defined in the text. **(C)** Same as **(B)** for chromatin remodelers. \*Note that co-TFs in **(A)** include chromatin remodelers.

binds to Ldb1 and thereby prevents Ldb1 from degradation<sup>51,52</sup>. While Ldb1 and Lhx2 were shown to bind to Greek island enhancers to regulate OR expression in trans<sup>40</sup>, Ssbp2 is a novel candidate with such a function.

The third group is specifically upregulated during the monogenic phase and two factors from this group have been identified by motif analysis. Differential analysis identifies additionally the TFs Mef2b, Rfx3 and Sub1, also known as PC4 (Fig. 4B and SFigure 8A). Among these factors, only Rfx3 has a known motif, for which we performed a strict motif search in Greek islands (Methods). We report that 56 out of 63 Greek islands contain the binding motif for Rfx3. Moreover, note that Mef2a, which shares a similar SRF binding domain with Mef2b<sup>53</sup>, is found to be strongly bound to OR promoters<sup>19</sup>.

The expression time courses of the transcription factors shown in Fig. 4B were found to be in good agreement with the Wang data (SFigure 4b).



Taken together, our pseudotime analysis recovers a large proportion of candidate (co-)TFs identified by motif analysis—both for initiating onset of OR expression, and for the transition to monogenic stage. Moreover it extends the range of candidates whose time course correlates with these two transitions, and consequently the regulatory repertoire for these transitions.

**Changes in chromatin remodeler expression accompany both transitions in the OR selection process.** It is known that chromatin changes accompany the selection of OR genes<sup>20,21,23</sup>. We therefore searched our data for chromatin remodelers that show expression changes during OR selection (Methods, Excel files 5,6). We confirmed previous observations that the chromatin remodelers Lbr and Cbx5 (SFigure 8B) are expressed at earlier stages and are downregulated in the course of OSN differentiation<sup>15,21</sup>. Furthermore, we discovered novel candidates for silencing OR genes, for onset of (multigenic) expression, and for transition to monogenic expression (Fig. 4C):

Among the genes whose expression profiles fall into group 1 (minimum in multigenic phase), we found Eed, one of the constitutive subunits of the polycomb repressive complex 2, PRC2 (Fig. 4C). Eed is required to maintain repressive H3K27me3 marks<sup>54,55</sup> and its downregulation may lead to de-repression of OR expression in Late.INP stage. Note that another PRC2 subunit, Ezh2 is expressed during the silent phase as well, but decreases later, at the transition to monogenic phase (SFigure 9). Nsd1 is a histone methyltransferase that demethylates H3K36me2<sup>56</sup> (Fig. 4C). All remodelers found with a group 1 pseudotime profile (Hells, H2afz and Set) are predicted to play a repressive role in the silent phase of OR selection (we only show H2afz as example SFigure 8B).

For group 2 (peak in multigenic phase), we found prominent chromatin remodelers such as Zmynd8, Ell3, Sertad2, Med12l and Scai (Fig. 4C, SFigure 8B). We also investigated the expression profile of Kdm1a which was known before as a regulator of OR expression. Kdm1a alias LSD1 is a Lysine demethylase and functions both as a coactivator by demethylation of mono- or di-methylated H3K9 and as a corepressor through demethylation of mono- or di-methylated H3K4<sup>57–60</sup>. There have been contradictory reports on the function of Kdm1a in OR expression as an activator<sup>27</sup> or repressor of transcription<sup>50</sup>. The present data sheds light on this debate: While Kdm1a expression sharply peaks directly before the multigenic phase (arguing for its role as activator), it can be part of the Co-REST repressor complex<sup>61</sup>. Two components of the Co-REST repressor complex, Rcor1 and Hdac2, sharply peak during multigenic phase (Fig. 5C, SFigure 10A), arguing for a change of function of Kdm1a by recruitment to the Co-REST complex at the transition to monogenic phase<sup>57,62,63</sup>.

Of the four novel remodelers with a group 2 pseudotime profile, Zmynd8 and Ell3 are highly differentially expressed (Fig. 4C). Zmynd8, a chromatin reader, is a particularly appealing candidate, since it is also known to play a role in the selective expression of the immunoglobulin heavy chain (*Igh*) regions in immune cells (B cells). Its product ZMYND8 binds *Igh* super-enhancers known as 3' regulatory region (3'RR). ZMYND8 thereby controls the 3'RR activity by modulating the enhancer transcriptional status<sup>64</sup>. Consistent with an activating role during the multigenic phase, Ell3 does not only bind to active enhancers, but also marks the enhancers that are in a poised or inactive state in ES cells<sup>65</sup>.

Remodelers whose pseudotime profiles fall into group 3 (specifically upregulated during the monogenic phase) are *Cbx4* (Chromobox 4) and *Jarid2* (Fig. 4C). *Cbx4* is a component of a Polycomb group (PcG) multiprotein PRC1-like complex, which is required to maintain the transcriptionally repressive state of many genes<sup>66,67</sup>. *Jarid2* (Jumonji and AT-Rich Interaction Domain Containing 2) is required to repress expression of cyclin-D1 (CCND1) in cardiac cells by setting H3K9 methylation marks<sup>68</sup>, and it is upregulated upon transition from Late.INP to iOSN (SFigure 2b), i.e. upon exit from the cell cycle.

We investigated the expression of the chromatin remodelers shown in Fig. 4C in the Wang data (SFigure 4b) and we found a good agreement with our time courses.

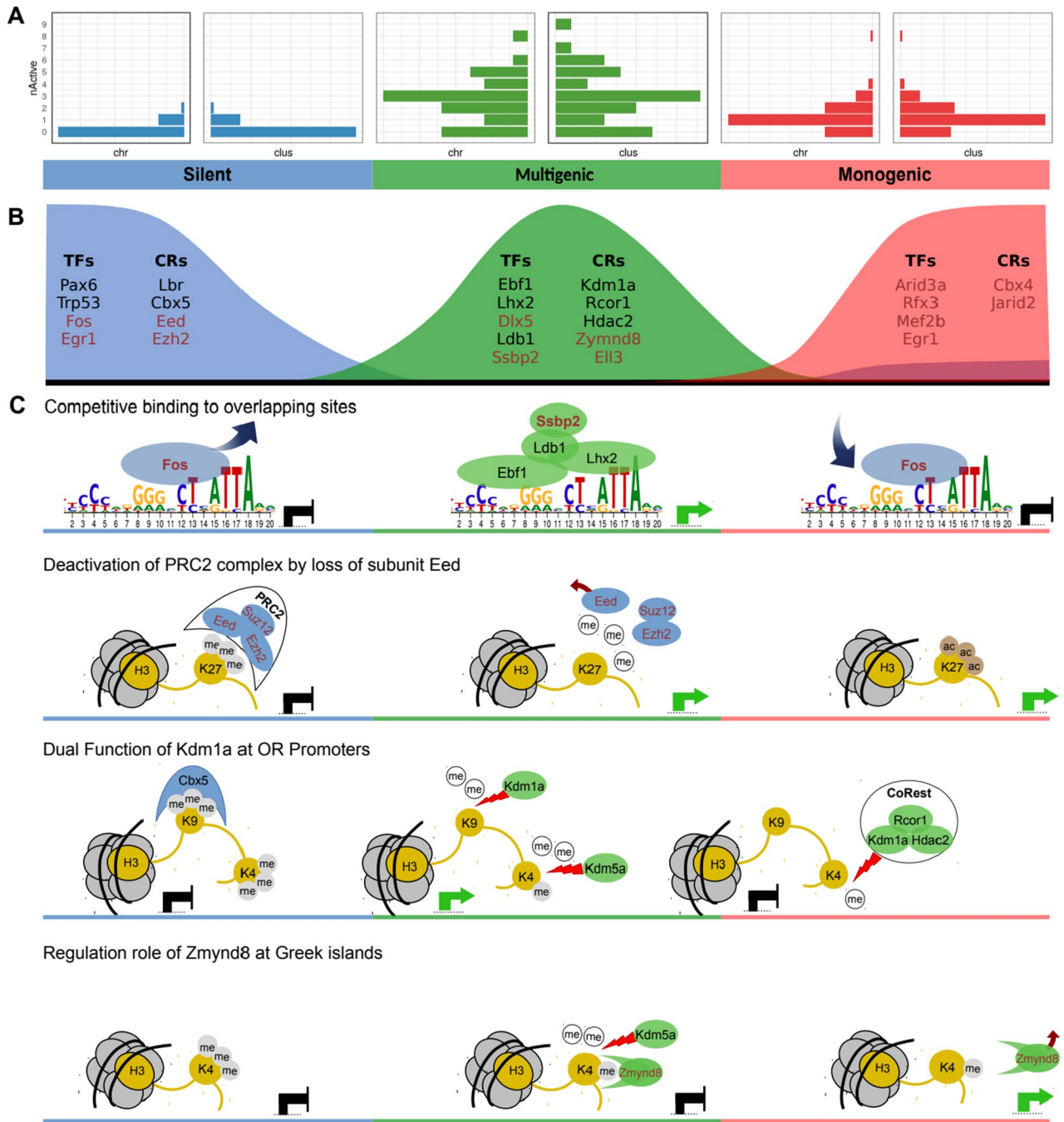
So far, we identified several chromatin remodelers that add to the regulatory repertoire for the onset of OR expression, and for the transition to monogenic stage. Another important feature which requires chromatin remodeling is the monoallelic expression of ORs in mature OSNs. This feature appears to be established from the very beginning of OR expression<sup>4,23</sup>. Factors involved in generating allelic exclusion therefore would be expected to peak at least as early as factors regulating the onset of (multigenic) expression (group 2 factors). We found two remodeling factors with a very early onset within group 2 which could potentially play such a role: *Smchd1* (structural maintenance of chromosomes flexible hinge domain-containing protein 1) and *Cdyl2* (chromodomain Y-like protein 2) (see SFigure 11).

Based on the interactions and temporal coordination of remodelers and (co-)TFs, we generate and discuss some hypotheses about OR selection below.

## Discussion

The monogenic expression of ORs presents a big challenge for the olfactory system, since it requires a random selection of exactly one OR per cell from the large family of OR genes. OR genes are known to be aggregated in silenced chromatin clusters during the silent phase, before the selection is made<sup>21</sup>. The identification of a multigenic state in early immature OSN<sup>29,30</sup> by scRNA-seq made clear that the initial escape from silencing is not limited to a single OR gene per cell. In a single cell, we observed up to nine different OR clusters and more than a dozen different OR genes concomitantly active. Both transcription factors and chromatin remodelers have been identified as regulators of OR expression.

Here we have employed pseudotime analysis of a single cell transcriptome data set<sup>31</sup> focusing on OR expression. We aligned the three main phases of OR expression, silent, multigenic and monogenic<sup>29,30</sup>, with the stages of OSN differentiation as defined in<sup>31</sup>. By analysis of winner vs. runners-up OR expression, we could precisely assign the onset of multigenic phase to Late.INP, and the onset of monogenic selection to iOSN. Most cells in the multigenic phase express ORs from more than one cluster and more than one chromosome (Figs. 2B, 5A).



**Figure 5.** Graphical summary of OR gene expression and hypothetical selection mechanisms. (A) Frequency distribution of active OR clusters and chromosomes per cell in the three phases of OR expression (Silent: most cells have 0 active cluster/chromosome, multigenic: 3 active cluster/chromosome is the most frequent case, and in the monogenic phase: most cells have only 1 active cluster/chromosome). (B) Pseudotime ordering of (co-)TFs and remodelers according to peak of expression in one of the three phases. Factors that so far have not been reported are highlighted in red. (C) Four mechanisms potentially contributing to OR selection. Factors are colored according to their phase of peak expression. The placement of the nucleosomes is according to the pseudotime at which the respective processes take place. Green arrows indicate an activating effect on target genes, black blunt end arrows indicate an inhibitory effect. Top row: Competitive binding of Fos vs. Ebf1 and Lhx2 to different parts of the Greek island consensus motif, resulting in different enhancer activities during the three OR selection phases. Fos blocks the binding of Ebf1 and Lhx2 to Greek islands in silent phase; in multigenic phase, Fos is significantly downregulated and allows the binding of the super-enhancer-forming complex (Ebf1, Lhx2, Ldb1, Ssbp2). Finally, when the OR selection is made during the monogenic phase, Fos is expressed again. Second row: Downregulation of the Eed subunit dissolves the PRC2 complex at the beginning of the multigenic phase and thus enables access of demethylases to H3K27me3. Third row: Methylated H3K9 is recognized, bound and protected by Cbx5 during the silent phase. At the end of this phase Cbx5 is downregulated and Kdm1a is upregulated. It participates in the process of H3K9 demethylation by removing the H3K9me1/2 mark. Later, after H3K4me3 has been demethylated to H3K4me2, presumably by Kdm5a during the multigenic phase, Kdm1a has a second role as part of the CoRest complex and demethylates H3K4me2. Bottom row: After demethylation of H3K4me3/2 by Kdm5a at Greek island enhancers, the chromatin reader Zmynd8, which peaks during the multigenic phase, can recognize H3K4me and repress enhancer activity. Activity is restored upon downregulation of Zmynd8 in monogenic phase.

We tested whether ORs co-expressed in one cell tend to lie in identical OR clusters and did not find evidence for that (SCode 1, SFigure 5). We conclude that the selective expression of OR genes during multigenic phase is not caused by escape from compaction of merely one OR cluster. In the monogenic phase, the vast majority of cells express only a single OR gene (Fig. 5A).

It is worth noting that using data with a high coverage per cell (Smart-seq) was instrumental to detecting the multigenic phase, as it was not detectable in data with lower coverage ( $10\times$  Genomics). Moreover, the  $10\times$  Genomics technology primarily queries the 3' end of a gene, while Smart-seq captures reads from the whole transcript region. OR genes have a high sequence similarity, which leads to exclusion of many 3' end reads that map to multiple loci. This may be another reason for the disproportionately low count numbers for OR genes in  $10\times$  Genomics data.

We then identified candidate factors (TFs and cofactors including chromatin remodelers) differentially expressed between these stages and thus potentially involved in the transitions silent-to-multigenic (40 TFs, 43 cofactors) and multigenic-to-monogenic (19 TFs, 20 cofactors) (Methods, Excel files 5,6). Many of these differentially expressed factors are likely not directly involved in OR selection, since cells undergo many substantial changes along OSN differentiation. Thus, we performed an independent de novo motif analysis of Greek islands, which should enrich TFs involved in OR selection. This search revealed one motif that could be decomposed into three consecutive submotifs, only two of which were described previously<sup>11</sup>. We additionally recognized a central Fos binding motif overlapping the previously described two elements (Fig. 3A). All but one of the factors that bind to these motifs show characteristic pseudotime expression time courses, which could be classified in three groups (Fig. 5B). The temporal coordination of these factors together with the precise location of their respective binding sites enables us to generate hypotheses about their molecular interactions and possible functional consequences. We provide extensive evidence from the literature to provide additional support for our hypotheses. Of course, this cannot replace pending experimental verification.

In the silent phase the transcriptional regulator Fos is binding to the central motif and may competitively prevent binding of the known activators of OR expression, Lhx2 and Ebf1<sup>11</sup>, which bind to the left and right submotifs (Fig. 3A, top row of Fig. 5C). The strong downregulation of Fos expression in the multigenic phase then would allow Lhx2 and Ebf1 to access their binding motifs, and recruit their known binding partner Ldb1 to Greek islands<sup>40</sup>. In situ Hi-C experimental data indicates that Ldb1 mediates trans interactions between different Greek islands, creating super-enhancer hubs that include neighboring OR clusters<sup>40</sup>. Our pseudotime analysis additionally predicted the co-factor Ssbp2 to play a role in this process. Ssbp2 is in other contexts known to bind to Ldb1 and thereby prevents its degradation by the proteasome<sup>51,52</sup>. Thus elevated expression of Ssbp2 in the multigenic phase would amplify the effect of Ldb1.

All four components of the super-enhancer-forming complex (Lhx2, Ebf1, Ldb1 and Ssbp2) peak during the multigenic phase (Fig. 4A,B, Fig. 5B and SFigure 8A), and thus are anti-correlated to Fos (rs: -0.58, -0.591, -0.264 and -0.258, respectively). In the monogenic phase the persistent expression of Ssbp2 would maintain Ldb1 protein levels<sup>51,52</sup>. Thus, the lower levels of Ldb1 transcript in monogenic phase presumably are still sufficient to keep the expression of the selected OR at high levels, the relative high activity of Fos in the monogenic phase notwithstanding. However, the relatively high Fos levels in the monogenic phase, together with lower levels of Lhx2 and Ldb1, may counteract formation of additional super-enhancer complexes thus ensuring stably monogenic expression—in line with the competitive interaction hypothesis outlined above (first row of Fig. 5C).

Our results found no OR expression during the Mid.INP stage despite significant expression of Lhx2, Ebf1 and Kdm1a (known activators of OR genes). We note that all components of the polycomb repressive complex 2 (PRC2), Eed, Ezh2 and Suz12, are active in Mid.INP, which could explain the absence of OR gene expression in Mid.INP despite the presence of the activators (second row of Fig. 5C). Moreover, an essential subunit of PRC2, Eed, is significantly reduced during onset of OR expression in Late.INP (Fig. 4 and SFigure 9). This elimination of Eed is sufficient to disassemble the PRC2 complex, which then can no longer maintain the repressive H3K27me3 mark<sup>54,55</sup>. Furthermore, we showed a dramatic reduction in expression of Ezh2 and Suz12 subunits of PRC2 along OSN differentiation (second row of Fig. 5C and SFigure 9). We conclude that PRC2 activity may be involved in repression of OR expression in Mid.INP. The disassembly of PRC2 in Late.INP then enables the Greek Island hubs to transiently activate the cis-corresponding OR gene/s, which enables the expression of multiple OR genes in most of Late.INP stage cells at the same time with relatively low levels compared to later stages of OSN differentiation (monogenic phase).

Heterochromatic silencing of ORs throughout OSN differentiation is enforced by the (interchromosomal) convergence of OR loci to OSN-specific, highly compacted nuclear bodies<sup>21</sup>. It has been shown that in the monogenic phase individual active OR genes require de-silencing by lysine demethylase Kdm1a<sup>27</sup> and spatial segregation of the single chosen OR allele towards euchromatic nuclear territories<sup>21</sup>. Another study however showed that the deletion of Kdm1a leads to persistent multigenic expression, suggesting a silencing role of Kdm1a rather than activating one<sup>50</sup>. Here, pseudotime analysis sheds light on the seemingly contradictory role of Kdm1a (third row of Fig. 5C):

The CBX5 protein is responsible for silencing of OR genes during the silent phase by binding to and thereby protecting the repressive H3K9me3 mark in gene bodies<sup>41</sup>. It vanishes upon transition to multigenic phase (SFigure 8B). After H3K9me3 has lost one methyl group (e.g., through the action of Kdm4a), Kdm1a can demethylate H3K9me2. Kdm1a peaks at Mid.INP stage and acts on di-methylated lysines only<sup>57</sup>. Thus the action of Kdm1a on H3K9me3 leads to activation in Mid.INP (Fig. 5C).

In contrast, another methylation site, on H3K4 is a mark of an active promoter/enhancer in the *methylated* stage (trimethylated for promoter, monomethylated for enhancer). Kdm5a peaks during multigenic phase (Fig. 5C, SFigure 10A) and can demethylate tri- or di-methylated H3K4 to its monomethylated form. Kdm1a by itself cannot act on H3K4me, but in a complex with Rcor1 and Hdac2 (CoRest complex) it is able to demethylate H3K4me2/1<sup>57,62,63</sup> (Fig. 5C). Rcor1 and Hdac2 have their peak expression during Late.INP stage, i.e. shortly after

onset of the Kdm1a peak (Fig. 5B, SFigure 10A). Thus, in Late.INP but not in Mid.INP, Kdm1a can demethylate H3K4me2/1, resulting in repression. This amounts to the multigenic phase (Late.INP) beginning to build up the molecular machinery to downregulate all but one of expressed ORs.

Taken together, the same enzymatic activity of the same factor (demethylation by Kdm1a) results in opposing effects on transcription due to co-factors Rcor1 and Hdac2 modulating substrate specificity of Kdm1a. Moreover our pseudotime analysis supports the hypothesis that OR expression issues a negative feedback signal on Kdm1a which is mediated by Atf5 and Adcy3 during transition from multigenic to monogenic phase<sup>25,27,28</sup> (see SFigure 10B).

During the monogenic phase, a super-enhancer is formed by the trans interaction between multiple Greek islands mediated by Lhx2-Ebf1-Ldb1-Ssbp2. Among the OR genes associated with this super-enhancer, merely one OR is expressed at very high levels<sup>40</sup>. During the multigenic phase, H3K4me3 marks of Greek islands<sup>20</sup> are converted to H3K4me, e.g. by Kdm5a (Fig. 5C). The latter histone mark can be recognized by the chromatin reader Zmynd8<sup>69</sup>. Zmynd8 has been shown to participate in another, highly specific selection process, namely the expression of *Igh* genes<sup>64</sup>. There, it recognizes H3K4me and represses super-enhancer activity in B cells. We therefore speculate that it might play a similar role in OR expression.

Expression of Zmynd8 peaks simultaneously with expression of the super-enhancer forming complex Ebf1-Lhx2-Ldb1-Ssbp2 during multigenic phase (Fig. 4C for Zmynd8 and Ssbp2, Fig. 3 for Ebf1 and Lhx2, SFigure 8A for Ldb1). Upon transition to the monogenic phase, ZMYND8 expression ceases while Ebf1 and Ssbp2 maintain intermediate expression. The disappearance of the ZMYND8 protein might abolish the suppression of the super-enhancer and could allow very high expression of the selected OR via the Lhx2-Ebf1-Ldb1-Ssbp2 complex (Fig. 5C bottom row).

So far it is unknown whether the monoallelic expression characteristic for mature OSNs is preceded by a biallelic phase. The dataset evaluated here does not allow us to analyze this question directly, because it does not exhibit enough sequence diversity between alleles to distinguish them. We searched for epigenetic factors known to be involved in allelic selection in other contexts, which show significant and relevant expression changes during OR differentiation. We found two factors, Smchd1 and Cdy12 (SFigure 11), which were discussed as stabilizing monoallelic expression. They play a role in epigenetic silencing, spermatogenesis, random inactivation of X chromosome, and stabilize monoallelic expression<sup>70–72</sup>. Both factors peak sharply in mid to late INP. Assuming these two factors initiate monoallelic expression, this would place monoallelic selection before or concomitant with the onset of (multigenic) expression in Late.INP, in other words, there might be no biallelic stage at all.

The dissection of the OSN maturation process into different stages allowed us to reveal three phases of OR gene selection. This in turn enabled an in-depth analysis of pseudotime expression profiles<sup>73</sup>, leading to several promising candidates and to testable hypotheses on the mechanisms involved in OR gene selection. However, the provision of additional, independent experimental evidence is beyond the scope of the present work, and will form the basis of future studies. For example, it will be interesting to complement the current data with single-cell chromatin accessibility data (single-cell ATAC-seq) or single-cell chromatin conformation data (Hi-C). Those experiments could also be carried out in the context of conditional knockouts of the factors identified above. Experiments should be carried out in hybrid crosses to additionally monitor allelic expression. Our study has narrowed down the cell stages and the time window that need to be analyzed for these purposes, thereby enhancing future research on this topic.

## Methods

Most statistical analyses and visualization were done in RStudio using R version 3.6.3.

**Data processing and quality control.** Our analysis of OR expression patterns during OSN differentiation is based on a scRNA-seq dataset generated by (Fletcher et al.<sup>31</sup>, GEO: GSE95601 “GSE95601\_oeHBCdiff\_Cufflinks\_eSet.Rda.gz” file). Ngai group investigated the homeostatic differentiation in the postnatal olfactory epithelium. Horizontal basal stem cells (HBC) were released from quiescence by a conditional knockout of the Trp63 transcription factor. Briefly, cells were FACS (fluorescence-activated cell sorting) selected for Sox2-EGFP-positive, ICAM1-negative, SCARB1/F3-negative expression to enrich for the cell population of interest (GBCs, later neuronal intermediates, and microvillous cells over sustentacular cells). Then scRNA-seq was done using the Fluidigm C1 microfluidics cell capture platform followed by Illumina multiplex sequencing. Processing of the raw data was done in Fletcher et al.<sup>31</sup> by RefSeq transcript annotations, which were used to align reads to the GRCh38 mouse genome with Tophat2, followed by Trimmomatic, featureCounts, and then Cufflinks. This resulted in 849 cells with 47,083 transcripts (before read and cell quality control). Cells were filtered according to Fletcher et al.<sup>31</sup> to remove contaminants, doublets and other technical artifacts, which resulted in 687 cells. Our transcript filtering slightly differs from Fletcher et al.<sup>31</sup>: We included all genes that have more than 40 counts in at least one cell to ensure retrieval of OR genes with sporadic expression (keeping 18,558 transcripts, among them 222 OR transcripts from an initial set of 1654 quantified OR transcripts). The filtering can be reproduced using an R script (Supplementary file ‘modified\_filtering.R’) which is modified from (<https://github.com/rufletch/p63-HBC-diff>). Count normalization and further transcript filtering was performed by SCTransform<sup>72</sup> with default parameters as implemented in Seurat (version 3.1.4)<sup>74–76</sup>. This resulted in 687 cells with a mean library size of 460 k unique reads and a median number of 4164 genes per cell and 17,179 quantified transcripts including 222 OR genes, then followed the LogNormalization seurat workflow. The library size distribution of these cells is shown in (SFigure 1).

**Dimension reduction, clustering and cell type assignment.** We followed the Seurat clustering workflow. First, dimension reduction was done using Principal Component Analysis (PCA). The number of

principal components kept was set to 15, after assessing the goodness of approximation by JackStraw and Elbow-Plot functions. A shared k-nearest neighbor graph was built by the FindNeighbors function. Afterwards, the Louvain algorithm was applied to define 13 distinct clusters from the shared nearest neighbor graph using the FindClusters function and the Jaccard index as a similarity measure. This number matches the number of clusters identified in (Fletcher et al.<sup>31</sup>) for this data. The expression of cell type marker genes that were collected from the literature (STable 2) served to assign cell clusters manually to known cell types according to the Seurat guidelines<sup>74</sup>. At least two expressed markers were required to confidently annotate a specific cell type. Visualization of the data was performed by PCA, UMAP<sup>77</sup> and tSNE<sup>78</sup>.

**Trajectory inference and pseudotime assignment.** Slingshot<sup>38</sup> was used to construct a minimum spanning tree (MST) based on the top 10-dimensional representation of the cells obtained above. The topology of the MST is independent of the root choice. For biological reasons, we selected qGBC as the root for the assignment of pseudo time, where it can differentiate to any cell type of MOE<sup>37</sup>. For each cell, a pseudotime between 0 (cells at the root node) and 1 (leaf node cells) was assigned by the slingPseudotime function.

**Differential expression analysis.** For each cell stage, we identified marker genes showing differential expression compared to all other cell stages using FindAllMarkers in Seurat, using the Wilcoxon rank sum test. Supplemental Fig. 3 shows a heatmap of the top 10 differentially expressed genes (i.e., putative markers) for each cell stage of MOE.

Among 2712 (co-)TFs (include chromatin remodelers) obtained from the GO.db package and AnimalTFDB3.0 (<http://bioinfo.life.hust.edu.cn/AnimalTFDB/#/>), we found 2004 (co-)TFs expressed (at least one count in one cell) in the neuronal lineage. In later differential expression analysis we compared the expression profiles of cell stages that were placed consecutively along the neuronal trajectory (i.e., the maturation of OSN) to identify genes that change their expression upon transition between cell stages (SFigure 2b). Volcano plots for each cell type transition were generated by a slightly modified EnhancedVolcano function (<https://github.com/kevinblighe/EnhancedVolcano>).

Differential expression analysis for silent to multigenic and multigenic to monogenic phase transitions were performed by comparison of pre-Late.INP cells vs. Late.INP cells and Late.INP vs. post-Late.INP cells, respectively. Genes with a Bonferroni adjusted  $p$ -value  $< 0.05$  (Wilcoxon rank sum test, FindMarkers function in Seurat) and an average absolute FC  $\geq 2$  were considered differentially expressed. This yielded 83 respectively 39 differentially expressed (co-)TFs for the two transitions.

**Motif analysis.** The genomic ranges of 68 OR clusters and 63 Greek islands were compiled from Monahan et al.<sup>11</sup>, which allows the matching of OR clusters and Greek islands. UCSC genome browser tools were used to find all genes inside OR clusters. We performed a motif search on the 63 Greek island sequences (STable 9 and FASTA file) and the approximate promoter regions (500 bp upstream the transcription starting sites) of all OR genes were obtained by using the "ucsc-twobittofa" bioconda package and the "biomart" R package respectively. The MEME suite web server for motif search and analysis<sup>42,43</sup> was used to predict the transcription factors (TF) that bind to Greek islands. We applied MEME using default values for all parameters to find the novel, ungapped motifs inside Greek islands with the following command: `meme greek_islands.fa -dna -oc . -nostatus -time 14,400 -mod zoops -nmotifs 3 -minw 6 -maxw 50 -objfun classic -revcomp -markov_order 0`.

Then we performed motif comparison between each motif found in the above-mentioned analysis against a database of known TFs motifs (JASPAR2018\_CORE\_non-redundant and uniprobe\_mouse databases) using Tomtom tool<sup>44</sup>. The Pearson correlation coefficient was used to measure the similarity between columns of position weight matrices (PWMs) and we restricted the results by setting  $q$ -value  $\leq 0.1$  (rather than 0.5 by default) as a threshold (10% FDR) using the following command: `tomtom -no-ssc -oc . -verbosity 1 -min-overlap 5 -mi 1 -dist pearson -thresh 0.1 -time 300 query_motifs db/MOUSE/uniprobe_mouse.meme db/JASPAR/JASPAR2018_CORE_non-redundant.meme`.

We also investigated the enrichment motifs in 63 Greek islands sequences using AME tool<sup>48</sup> by using an average odds score method and Fisher's exact test as a motif enrichment test through the following command: `ame -verbose 1 -oc . -scoring avg -method fisher -hit-lo-fraction 0.25 -evalue-report-threshold 1.0 -control -shuffle -kmer 2 greek_islands.fa db/MOUSE/uniprobe_mouse.meme db/JASPAR/JASPAR2018_CORE_non-redundant.meme`.

Finally, a strict motif search in Greek islands for selected TFs was done by "ucsc-findmotif" bioconda package, allowing for 3 mismatches.

**Visualization of time series.** Grouped time series<sup>79</sup> was used to visualize pseudotime series of individual genes and to calculate and visualize aggregated groups of genes, e.g. all OR genes. Since the original expression count matrix is sparse (75.45% zero count entries), we first applied ALRA80, which has specifically been designed for the imputation of missing values in scRNA-Seq data. The imputed expression matrix retrieved ~2403 missing values, reducing the fraction of zero count entries to 61.50%. The median number of expressed genes per cell was 6715 (see SFigure 1b). Note that the imputed expression matrix was used only for visualization, for all analysis steps we used normalized counts without data imputation.

### Data availability

The single cell data has been generated by<sup>31</sup> and is available at GEO under the accession number GSE95601 [<https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE95601>]. The processed data and all additional data used in our analysis are available as Supplemental Materials.

## References

- Korsching, S. I. Olfaction. in *The Physiology of Fishes* 256, Chapter 14 (CRC Press, 2020).
- Niimura, Y. Olfactory receptor multigene family in vertebrates: from the viewpoint of evolutionary genomics. *Curr. Genomics* **13**, 103–114 (2012).
- Buck, L. & Axel, R. A novel multigene family may encode odorant receptors: a molecular basis for odor recognition. *Cell* **65**, 175–187 (1991).
- Chess, A., Simon, I., Cedar, H. & Axel, R. Allelic inactivation regulates olfactory receptor gene expression. *Cell* **78**, 823–834 (1994).
- Mombaerts, P. *et al.* Visualizing an Olfactory Sensory Map. *Cell* **87**, 675–686 (1996).
- Ressler, K. J., Sullivan, S. L. & Buck, L. B. Information coding in the olfactory system: Evidence for a stereotyped and highly organized epitope map in the olfactory bulb. *Cell* **79**, 1245–1255 (1994).
- Buck, L. B. Information coding in the vertebrate olfactory system. *Annu. Rev. Neurosci.* **19**, 517–544 (1996).
- Feinstein, P. & Mombaerts, P. A contextual model for axonal sorting into glomeruli in the mouse olfactory system. *Cell* **117**, 817–831 (2004).
- Serizawa, S., Miyamichi, K. & Sakano, H. One neuron–one receptor rule in the mouse olfactory system. *Trends Genet.* **20**, 648–653 (2004).
- Mombaerts, P. Odorant receptor gene choice in olfactory sensory neurons: the one receptor–one neuron hypothesis revisited. *Curr. Opin. Neurobiol.* **14**, 31–36 (2004).
- Monahan, K. *et al.* Cooperative interactions enable singular olfactory receptor expression in mouse olfactory neurons. *eLife* **6**, e28620 (2017).
- Zhang, X. *et al.* High-throughput microarray detection of olfactory receptor gene expression in the mouse. *Proc. Natl. Acad. Sci.* **101**, 14168–14173 (2004).
- Miyamichi, K., Serizawa, S., Kimura, H. M. & Sakano, H. Continuous and overlapping expression domains of odorant receptor genes in the olfactory epithelium determine the dorsal/ventral positioning of glomeruli in the olfactory bulb. *J. Neurosci.* **25**, 3586–3592 (2005).
- Tsuboi, A., Miyazaki, T., Imai, T. & Sakano, H. Olfactory sensory neurons expressing class I odorant receptors converge their axons on an antero-dorsal domain of the olfactory bulb in the mouse. *Eur. J. Neurosci.* **23**, 1436–1444 (2006).
- Markenscoff-Papadimitriou, E. *et al.* Enhancer interaction networks as a means for singular olfactory receptor expression. *Cell* **159**, 543–557 (2014).
- Michaloski, J. S., Galante, P. A. F. & Malnic, B. Identification of potential regulatory motifs in odorant receptor genes by analysis of promoter sequences. *Genome Res.* **16**, 1091–1098 (2006).
- Hirota, J. & Mombaerts, P. The LIM-homeodomain protein Lhx2 is required for complete development of mouse olfactory sensory neurons. *Proc. Natl. Acad. Sci.* **101**, 8751–8755 (2004).
- Clowney, E. J. *et al.* High-throughput mapping of the promoters of the mouse olfactory receptor genes reveals a new type of mammalian promoter and provides insight into olfactory receptor gene regulation. *Genome Res.* **21**, 1249–1259 (2011).
- Plessy, C. *et al.* Promoter architecture of mouse olfactory receptor genes. *Genome Res.* **22**, 486–497 (2012).
- Magklara, A. *et al.* An epigenetic signature for monoallelic olfactory receptor expression. *Cell* **145**, 555–570 (2011).
- Clowney, E. J. *et al.* Nuclear aggregation of olfactory receptor genes governs their monogenic expression. *Cell* **151**, 724–737 (2012).
- Lyons, D. B. *et al.* Heterochromatin-mediated gene silencing facilitates the diversification of olfactory neurons. *Cell Rep.* **9**, 884–892 (2014).
- Armelin-Correa, L. M., Gutiyama, L. M., Brandt, D. Y. C. & Malnic, B. Nuclear compartmentalization of odorant receptor genes. *Proc. Natl. Acad. Sci.* **111**, 2782–2787 (2014).
- Armelin-Correa, L. M., Nagai, M. H., Silva, A. G. L. & Malnic, B. Nuclear architecture and gene silencing in olfactory sensory neurons. *BioArchitecture* **4**, 160–163 (2014).
- Serizawa, S. *et al.* Negative feedback regulation ensures the one receptor-one olfactory neuron rule in mouse. *Science* **302**, 2088–2094 (2003).
- Lewcock, J. W. & Reed, R. R. A feedback mechanism regulates monoallelic odorant receptor expression. *Proc. Natl. Acad. Sci.* **101**, 1069–1074 (2004).
- Lyons, D. B. *et al.* An epigenetic trap stabilizes singular olfactory receptor expression. *Cell* **154**, 325–336 (2013).
- Dalton, R. P., Lyons, D. B. & Lomvardas, S. Co-opting the unfolded protein response to elicit olfactory receptor feedback. *Cell* **155**, 321–332 (2013).
- Hanchate, N. K. *et al.* Single-cell transcriptomics reveals receptor transformations during olfactory neurogenesis. *Science* <https://doi.org/10.1126/science.aad2456> (2015).
- Tan, L., Li, Q. & Xie, X. S. Olfactory sensory neurons transiently express multiple olfactory receptors during development. *Mol. Syst. Biol.* **11**, 844 (2015).
- Fletcher, R. B. *et al.* Deconstructing olfactory stem cell trajectories at single-cell resolution. *Cell Stem Cell* **20**, 817–830.e8 (2017).
- Wang, X., He, Y., Zhang, Q., Ren, X. & Zhang, Z. Direct comparative analyses of 10X genomics chromium and smart-seq2. *Genomics Proteomics Bioinformatics* **19**, 253–266 (2021).
- Wang, I.-H. *et al.* Spatial transcriptomic reconstruction of the mouse olfactory glomerular map suggests principles of odor processing. *Nat. Neurosci.* **25**, 484–492 (2022).
- Caggiano, M., Kauer, J. S. & Hunter, D. D. Globose basal cells are neuronal progenitors in the olfactory epithelium: A lineage analysis using a replication-incompetent retrovirus. *Neuron* **13**, 339–352 (1994).
- Chen, X., Fang, H. & Schwob, J. E. Multipotency of purified, transplanted globose basal cells in olfactory epithelium. *J. Comp. Neurol.* **469**, 457–474 (2004).
- Chen, M. *et al.* Wnt-responsive Lgr5+ globose basal cells function as multipotent olfactory epithelium progenitor cells. *J. Neurosci.* **34**, 8268–8276 (2014).
- Jang, W., Chen, X., Flis, D., Harris, M. & Schwob, J. E. Label-retaining, quiescent globose basal cells are found in the olfactory epithelium. *J. Comp. Neurol.* **522**, 731–749 (2014).
- Street, K. *et al.* Slingshot: cell lineage and pseudotime inference for single-cell transcriptomics. *BMC Genomics* **19**, 477 (2018).
- Monahan, K. & Lomvardas, S. Monoallelic expression of olfactory receptors. *Annu. Rev. Cell Dev. Biol.* **31**, 721–740 (2015).
- Monahan, K., Horta, A. & Lomvardas, S. LHX2- and LDB1-mediated trans interactions regulate olfactory receptor choice. *Nature* **565**, 448–453 (2019).
- Lomvardas, S. *et al.* Interchromosomal interactions and olfactory receptor choice. *Cell* **126**, 403–413 (2006).
- Bailey, T. L. *et al.* MEME suite: Tools for motif discovery and searching. *Nucleic Acids Res.* **37**, W202–W208 (2009).
- Bailey, T. L., Johnson, J., Grant, C. E. & Noble, W. S. The MEME suite. *Nucleic Acids Res.* **43**, W39–W49 (2015).
- Gupta, S., Stamatoyannopoulos, J. A., Bailey, T. L. & Noble, W. S. Quantifying similarity between motifs. *Genome Biol.* **8**, R24 (2007).

45. Wang, H.-N. *et al.* Inhibition of c-Fos expression attenuates IgE-mediated mast cell activation and allergic inflammation by counteracting an inhibitory AP1/Egr1/IL-4 axis. *J. Transl. Med.* **19**, 261 (2021).
46. Ray, N. *et al.* c-Fos suppresses systemic inflammatory response to endotoxin. *Int. Immunol.* **18**, 671–677 (2006).
47. Bahrami, S. & Drablos, F. Gene regulation in the immediate-early response process. *Adv. Biol. Regul.* **62**, 37–49 (2016).
48. McLeay, R. C. & Bailey, T. L. Motif enrichment analysis: A unified framework and an evaluation on ChIP data. *BMC Bioinformatics* **11**, 165 (2010).
49. Zhang, G., Titlow, W. B., Biecker, S. M., Stromberg, A. J. & McClintock, T. S. Lhx2 Determines odorant receptor expression frequency in mature olfactory sensory neurons. *eNeuro* **3**, (2016).
50. Vyas, R. N., Meredith, D. & Lane, R. P. Lysine-specific demethylase-1 (LSD1) depletion disrupts monogenic and monoallelic odorant receptor (OR) expression in an olfactory neuronal cell line. *Mol. Cell. Neurosci.* **82**, 1–11 (2017).
51. Wang, Y. *et al.* SSBP2 is an in vivo tumor suppressor and regulator of LDB1 stability. *Oncogene* **29**, 3044–3053 (2010).
52. Wang, H. *et al.* Crystal structure of human LDB1 in complex with SSBP2. *Proc. Natl. Acad. Sci.* **117**, 1042–1048 (2020).
53. Wu, Y. *et al.* Structure of the MADS-box/MEF2 domain of MEF2A bound to DNA and its implication for myocardin recruitment. *J. Mol. Biol.* **397**, 520–533 (2010).
54. Cao, Q. *et al.* The central role of EED in the orchestration of polycomb group complexes. *Nat. Commun.* **5**, 3127 (2014).
55. Potjewyd, F. *et al.* Degradation of polycomb repressive complex 2 with an EED-targeted bivalent chemical degrader. *Cell Chem. Biol.* **27**, 47–56.e15 (2020).
56. Qiao, Q. *et al.* The structure of NSD1 reveals an autoregulatory mechanism underlying histone H3K36 methylation. *J. Biol. Chem.* **286**, 8361–8368 (2011).
57. Bannister, A. J. & Kouzarides, T. Regulation of chromatin by histone modifications. *Cell Res.* **21**, 381–395 (2011).
58. Wang, J. *et al.* Opposing LSD1 complexes function in developmental gene activation and repression programmes. *Nature* **446**, 882–887 (2007).
59. Kilinc, S., Savarino, A., Coleman, J. H., Schwob, J. E. & Lane, R. P. Lysine-specific demethylase-1 (LSD1) is compartmentalized at nuclear chromocenters in early post-mitotic cells of the olfactory sensory neuronal lineage. *Mol. Cell. Neurosci.* **74**, 58–70 (2016).
60. Rusconi, F., Grillo, B., Toffolo, E., Mattevi, A. & Battaglioli, E. NeuroLSD1: Splicing-generated epigenetic enhancer of neuroplasticity. *Trends Neurosci.* **40**, 28–38 (2017).
61. Coleman, J. H., Lin, B. & Schwob, J. E. Dissecting LSD1-dependent neuronal maturation in the olfactory epithelium. *J. Comp. Neurol.* **525**, 3391–3413 (2017).
62. Song, Y. *et al.* Mechanism of crosstalk between the LSD1 demethylase and HDAC1 deacetylase in the CoREST complex. *Cell Rep.* **30**, 2699–2711.e8 (2020).
63. Shi, Y.-J. *et al.* Regulation of LSD1 histone demethylase activity by its associated factors. *Mol. Cell* **19**, 857–864 (2005).
64. Delgado-Benito, V. *et al.* The chromatin reader ZMYND8 regulates Igh enhancers to promote immunoglobulin class switch recombination. *Mol. Cell* **72**, 636–649.e8 (2018).
65. Lin, C., Garruss, A. S., Luo, Z., Guo, F. & Shilatifard, A. The RNA Pol II elongation factor Ell3 marks enhancers in ES cells and primes future gene activation. *Cell* **152**, 144–156 (2013).
66. Levine, S. S. *et al.* The core of the polycomb repressive complex is compositionally and functionally conserved in flies and humans. *Mol. Cell. Biol.* <https://doi.org/10.1128/MCB.22.17.6070-6078.2002> (2002).
67. Vandamme, J., Völkel, P., Rosnoble, C., Faou, P. L. & Angrand, P.-O. Interaction Proteomics Analysis of Polycomb Proteins Defines Distinct PRC1 Complexes in Mammalian Cells. *Mol. Cell. Proteomics* **10**, (2011).
68. Shirato, H. *et al.* A Jumonji (Jarid2) protein complex represses cyclin D1 expression by methylation of histone H3–K9\*. *J. Biol. Chem.* **284**, 733–739 (2009).
69. Li, N. *et al.* ZMYND8 reads the dual histone mark H3K4me1–H3K14ac to antagonize the expression of metastasis-linked genes. *Mol. Cell* **63**, 470–484 (2016).
70. Mould, A. W. *et al.* Smchd1 regulates a subset of autosomal genes subject to monoallelic expression in addition to being critical for X inactivation. *Epigenetics Chromatin* **6**, 19 (2013).
71. Qin, R. *et al.* CDYL deficiency disrupts neuronal migration and increases susceptibility to epilepsy. *Cell Rep.* **18**, 380–390 (2017).
72. Liu, S. *et al.* Chromodomain protein CDYL acts as a crotonyl-CoA hydratase to regulate histone crotonylation and spermatogenesis. *Mol. Cell* **67**, 853–866.e5 (2017).
73. Tanay, A. & Regev, A. Scaling single-cell genomics from phenomenology to mechanism. *Nature* **541**, 331–338 (2017).
74. Satija, R., Farrell, J. A., Gennert, D., Schier, A. F. & Regev, A. Spatial reconstruction of single-cell gene expression data. *Nat. Biotechnol.* **33**, 495–502 (2015).
75. Hafemeister, C. & Satija, R. Normalization and variance stabilization of single-cell RNA-seq data using regularized negative binomial regression. *Genome Biol.* **20**, 296 (2019).
76. Stuart, T. & Satija, R. Integrative single-cell analysis. *Nat. Rev. Genet.* **20**, 257–272 (2019).
77. McInnes, L., Healy, J., Saul, N. & Großberger, L. UMAP: Uniform manifold approximation and projection. *J. Open Source Softw.* **3**, 861 (2018).
78. van der Maaten, L. & Hinton, G. Visualizing data using t-SNE. *J. Mach. Learn. Res.* **9**, 2579–2605 (2008).
79. Hyndman, R. & Athanasopoulos, G. Forecasting: Principles and Practice. (2021).
80. Linderman, G. C., Zhao, J. & Kluger, Y. Zero-preserving imputation of scRNA-seq data using low-rank approximation. *bioRxiv* 397588 (2018) <https://doi.org/10.1101/397588>.

## Author contributions

A.T. and S.I.K. conceived and designed the analysis; M.H. collected the data, conducted the analysis and prepared the figures. All authors wrote, read and approved the manuscript.

## Funding

Open Access funding enabled and organized by Projekt DEAL.

## Competing interests

The authors declare no competing interests.

## Additional information

**Supplementary Information** The online version contains supplementary material available at <https://doi.org/10.1038/s41598-022-20106-w>.

**Correspondence** and requests for materials should be addressed to A.T.

**Reprints and permissions information** is available at [www.nature.com/reprints](http://www.nature.com/reprints).

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2022