1 **Single-cell transcriptome analysis of fish immune cells**

2 **provides insight into the evolution of vertebrate immunity**

3 Santiago J. Carmona[1,2], Sarah A. Teichmann[3,4*], Lauren Ferreira[4,5,6], Iain C.

4 Macaulay[7], Michael J.T. Stubbington[3], Ana Cvejic[4,5,6*], David Gfeller[1,2*]

5

6

7 [1]Ludwig Center for Cancer Research, University of Lausanne, Lausanne, Switzerland

8 [2]Swiss Institute of Bioinformatics (SIB), 1015 Lausanne, Switzerland

9 [3]European Molecular Biology Laboratory European Bioinformatics Institute,

10 Wellcome Trust Genome Campus, Hinxton, Cambridge, UK

11 [4]Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Cambridge, UK

12 [5]Department of Haematology, University of Cambridge, Cambridge, UK

13 [6]Wellcome Trust – Medical Research Council Cambridge Stem Cell Institute,

14 Cambridge, CB2 1QR, UK

15 [7]Sanger Institute–EBI Single-Cell Genomics Centre, Wellcome Trust Genome

16 Campus, Hinxton, Cambridge, UK

17

18

19 *Joint senior authors and to whom correspondence should be addressed:

20 saraht@ebi.ac.uk, as889@cam.ac.uk, david.gfeller@unil.ch.

21

22

23

24

25

26

27

28

29

30

32   **Abstract**

33   The immune system of vertebrate species consists of many different cell types that

34   have distinct functional roles and are subject to different evolutionary pressures.

35   Here, we first analysed gene conservation of all major immune cell types in human

36   and mouse. Our results revealed higher gene turnover and faster evolution of trans-

37   membrane proteins in NK cells compared to other immune cell populations, and

38   especially T cells, but similar conservation of nuclear and cytoplasmic protein coding

39   genes. To validate these findings in a distant vertebrate species, we used single-cell

40   RNA-Sequencing of *lck:GFP* cells in zebrafish to obtain the first transcriptome of

41   specific immune cell types in a non-mammalian species. Unsupervised clustering

42   and single-cell TCR locus reconstruction identified three cell populations, T-cells, a

43   novel type of NK-like cells and a smaller population of myeloid-like cells. Differential

44   expression analysis uncovered new immune cell specific genes, including novel

45   immunoglobulin-like receptors, and neofunctionalization of recently duplicated

46   paralogs. Evolutionary analyses confirmed a higher gene turnover and lower

47   conservation of trans-membrane proteins in NK cells compared to T cells in fish

48   species, suggesting that this is a general property of immune cell types across all

49   vertebrates.

50

51

52    **Introduction**

53

54    The immune system of vertebrate species has evolved into a highly complex

55    structure, comprising many different subtypes of both innate and adaptive immune

56    cells. Adaptive immune cells are broadly classified into B and T lymphocytes that can

57    directly recognize and bind antigens with great specificity. Innate immune cells

58    include a variety of myeloid cells such as monocytes, neutrophils, basophils,

59    eosinophils and mast cells. A third major type of lymphocytes, the Natural Killer (NK)

60    cells, has also been historically classified among innate immune cells (Sun and

61    Lanier 2009; Sun et al. 2009). Traditionally, different immune cell types are

62    distinguished based on unique combinations of cell surface markers. In mouse and

63    human, many antibodies for these markers are available and can be used to isolate

64    homogeneous immune cell populations using flow cytometry. Gene expression

65    profiling studies of isolated immune cell populations have further allowed genome-

66    wide identification of cell-type specific genes (Shay and Kang 2013; Watkins et al.

67    2009; Chambers et al. 2007; Vu Manh et al. 2014). These studies revealed an overall

68    conservation of immune cells' gene expression between mouse and human (Shay et

69    al. 2013). However, beyond mouse and human, less is known about the

70    characteristics and evolution of immune cell types mainly due to the challenges of

71    isolating different immune cell populations.

72

73    Evolutionary studies based on mouse and human genes have shown that immune-

74    related genes tend to evolve faster than other genes (Bailey et al. 2013; Boehm

75    2012; Flajnik and Kasahara 2010; Kosiol et al. 2008). This faster evolution may

76    reflect a need of immune cells to adapt to a rapidly changing environment and

77    specific pathogens. In addition, different immune cell types are subject to different

78    evolutionary constrains. T and B lymphocytes can generate an extraordinary diverse

79    repertoire of antigen-specific receptors as a consequence of *rag*-mediated somatic

80    V(D)J rearrangement, and this process is conserved across all jawed vertebrates

81    (Boehm 2012). Many orthologs of T cell specific genes, like CD4, CD8 and TCR

82    genes, have been identified in all jawed vertebrates. In species like zebrafish, the

83    V(D)J variable regions has been recently annotated (Meeker et al. 2010; Schorpp et

84    al. 2006; Iwanami 2014). NK receptors instead are germline-encoded. Therefore,

85    selection pressure to generate different receptor specificities and transduce signals is

86    expected to operate at the population rather than at the individual cell level. Indeed,

3

87  mammalian NK cell receptors have expanded and diversified in a species-specific
88  fashion, like in the case of KIR receptors in primates and Ly49/Killer lectin-like
89  receptors in rodents (Carrillo-Bustamante et al. 2016). NK-like cells have been
90  identified in non-mammalian species such as chicken (Jansen et al. 2010), xenopus
91  (Horton et al. 1996) and catfish where spontaneous killing of allogeneic cells by non-
92  TCR expressing cytotoxic cells was demonstrated (Shen et al. 2004; Yoder 2004). A
93  recent study using single-cell qPCR based on known markers of blood cell lineages
94  revealed the presence of a small sub-population of immune cells in zebrafish, which
95  were proposed to represent putative NK-like cells based on expression of NK-lysin
96  genes (Moore et al. 2016). The identification of membrane receptors with similar
97  genomic organization as the KIR genes in human provided additional evidence for
98  the existence of NK cells in fish species. In zebrafish, these receptors include *nitr*
99  (Yoder et al. 2004) and *dicp* genes (Haire et al. 2012). However, pure T and NK cell
100 populations have so far not been isolated in zebrafish and no reliable antibody has
101 been developed against orthologs of mammalian T and NK cell receptors. Therefore
102 many properties of mammalian T and NK cell orthologs and their evolution in non-
103 mammalian species remain uncharacterized.

104

105 High-throughput single cell RNA-Seq (scRNA-Seq) has emerged as a promising
106 technology to unravel the landscape of cell types in heterogeneous cell populations
107 without relying on specific antibodies (Saliba et al. 2014). Simultaneous expression
108 of thousands of genes can be measured in each cell, thereby providing an unbiased
109 view of transcriptional activity at the cellular level and avoiding the averaging effect of
110 bulk gene expression studies (Shapiro et al. 2013). Cells can then be grouped into
111 biologically relevant clusters based on similarity of their gene expression profiles
112 rather than a handful of cell surface markers (Grün et al. 2015; Trapnell 2015;
113 Macaulay et al. 2016). Therefore, despite technical and biological noise and the
114 computational challenges associated with this variability (Brennecke et al. 2013;
115 Buettner et al. 2015), scRNA-Seq has the potential to uncover new immune cell
116 types that cannot be studied using traditional approaches.

117

118 To gain insight into the evolution of vertebrate immunity and non-mammalian
119 immune cell types, we first analysed the conservation of mouse and human immune
120 cell (i.e., T, B, NK and myeloid cells) specific genes. Next, we analysed immune cells
121 in zebrafish, a powerful model in biomedical research (Langenau and Zon 2005;

4

122 Renshaw and Trede 2012; Kaufman et al. 2016). To this end, we took advantage of
123 a transgenic line of zebrafish expressing GFP under the control of the *lck* promoter
124 (Langenau et al. 2004) and performed scRNA-Seq on *lck:GFP+* FACS sorted cells.
125 Our analysis revealed three distinct *lck+* cell populations: T cells, a novel type of NK-
126 like cells and a smaller population of myeloid-like cells. Our expression profiles
127 uncovered many immune signature genes both bony-fish specific and shared with
128 mammals, including innate immune receptors, cytokines, transcription factors,
129 proteases and antimicrobial peptides. In addition, we detected gene expression
130 divergence among bony-fish duplicated paralogs. Finally, evolutionary analysis of
131 differentially expressed genes showed higher gene turnover and lower conservation
132 for NK/NK-like cells specific genes compared to T cells specific genes in all three
133 species studied in this work.

134

135 **Results**

136

137 **Conservation analysis of mammalian T, B, NK and myeloid cell specific genes**
138 **across vertebrates**

139 Immune related genes tend to evolve more rapidly than other genes and between
140 functionally distinct immune cells the selective pressures might vary significantly.
141 Here we performed a conservation analysis of the most differentially expressed
142 genes in resting T, B, NK and myeloid cells in mouse and human at the genome-
143 wide level (Chambers et al. 2007; Watkins et al. 2009) (see Methods).

144

145 Our analysis revealed that among trans-membrane (TM) or secreted protein coding
146 genes, those specifically expressed in NK cells have proportionally less orthologs
147 across all vertebrates compared to other immune cells. The difference is most
148 evident between NK and T cells, although these are closer from a functional and
149 ontogenical point of view (Fig. 1A and C). No difference, however, was observed for
150 cytoplasmic or nuclear protein coding genes (Fig. 1B and D). As expected, the NK
151 receptor families Ly49 in mouse and KIR in human strongly contributed to this
152 difference. Interestingly, however, the differences between T and NK cell TM gene
153 conservation were still observed after removing these receptors from the analysis
154 (Supplemental Fig. S1). Examples of other mouse or human NK TM genes poorly
155 conserved across vertebrates include Fc receptors, granulysin, CD160, CD244 and
156 IFITM3. In addition, among conserved protein coding genes, NK cell specific genes

157  consistently had lower sequence identity across all vertebrates for TM genes, but not
158  for cytoplasmic ones (Supplemental Fig. S2).

159

160  The ratio between non-synonymous and synonymous substitutions (dN/dS ratio) of
161  one-to-one orthologs between human and mouse, can provide a good estimation of
162  the evolutionary pressure acting on a gene. Our results indicate that NKs' TM genes
163  evolve faster (i.e., present higher dN/dS values) compared to T cells' TM genes
164  (Supplemental Fig. S3). As expected, the lowest conservation for all immune cell
165  type specific genes was observed in lamprey (Fig. 1) since these organisms possess
166  a distinct adaptive immune system (Guo et al. 2009).

167

168  To further explore the conservation of immune cell types' specific genes and expand
169  our understanding of immune cell populations in an evolutionary distant non-
170  mammalian species, we set out to profile immune cell populations in zebrafish.

171

172  **Single-cell transcriptomics of zebrafish lck+ lymphocytes reveals three distinct**
173  **cell populations corresponding to T cells, NK-like cells and Myeloid-like cells.**
174  As reliable antibodies to isolate pure immune cell populations in fish species are not
175  available, we used single cell transcriptome analysis of zebrafish *Tg(lck:GFP)* cells.
176  This transgenic line expresses GFP under the control of the lymphocyte-specific
177  protein tyrosine kinase (*lck*) promoter, and it was shown to be mainly restricted to
178  zebrafish T cells (Langenau et al. 2004) However, as *Lck* in mouse and humans is
179  expressed in both T and NK cells, we speculated that its expression pattern could be
180  conserved in ray-finned fish. *Tg(lck:GFP)* zebrafish may therefore provide an ideal
181  model to investigate the large difference in conservation between T and NK cell
182  specific genes observed in mammalian species. To simultaneously obtain
183  information about cell morphology and high quality gene expression profiles, we used
184  high-throughput single-cell RNA sequencing combined with FACS (fluorescent
185  activated cell sorting) index sorting analysis of two adult zebrafish (three and ten
186  months old) spleen-derived *lck:GFP* cells.

187

188  We first generated and sequenced libraries from 278 single GFP+ cells isolated from
189  the spleen of two different fish from a different clutch and different age (see
190  Methods). Following quality controls (Supplemental Fig. S4, see Methods) 15 cells
191  were removed and gene expression profiles for the remaining 263 cells were

192  generated. Average single-cell profiles showed good correlation with independent
193  bulk samples (PCC=0.82, Supplemental Fig. S5). Correlations between single-cell
194  gene expression profiles were used to calculate cell-to-cell dissimilarities (see
195  Methods) and these were represented into low dimensional space using classical
196  Multidimensional Scaling (see Methods). Interestingly, a clear cell subpopulation
197  structure emerged (Fig. 2A) showing three distinct cell groups. The three groups
198  were confirmed by unsupervised whole-transcriptome clustering (see Methods, Fig. 3
199  and Supplemental Fig. S6E-F).
200
201  The first cluster (C1) seemed to correspond to T cells based on the expression of
202  *cd8a* and *cd4* genes (Fig. 2). To further support this hypothesis, we adapted a recent
203  method for detection of V(D)J recombination events of the TCR locus (Stubbington et
204  al. 2016) (see Methods). With a median of only 0.64 million gene-mapped reads per
205  cell, we were able to unambiguously detect V(D)J recombination events in 27 cells
206  (Fig. 3 and Supplemental Fig. S10). Occurrence of V(D)J recombination was
207  associated with Cluster 1 (p<0.01, see Methods), which provides additional genomic
208  evidences of the T cell identity. As expected, V(D)J recombined segments were also
209  strongly associated with expression of the TCR beta chain constant region (*trbc1*) (p
210  $< 10^{-5}$, Fig. 3). Interestingly, *cd8* and *cd4* displayed mutually exclusive expression (as
211  expected for mature T cells) (Fig. 2A and 3) and *cd4+* and *cd8+* cells clearly
212  separated when low dimensional projection is restricted to cells from C1
213  (Supplemental Fig. S8).
214
215  The second cluster (C2) displayed expression of NK lysins, which have been recently
216  proposed to mark a distinct sub-population of NK-like cells and are upregulated in
217  Recombination activation gene 1 deficient (*rag1−/−*) zebrafish (Moore et al. 2016;
218  Pereiro et al. 2015). In addition, multiple members of the teleost fish specific innate
219  immune receptor families *nitr* (novel immune type receptors) and *dicp* (diverse
220  immunoglobulin domain containing proteins) are highly expressed and specific to this
221  cluster. It has been suggested that these receptors play a similar role as mammalian
222  NK receptors (Yoder et al. 2004). Therefore, we hypothesized that this subpopulation
223  corresponds to a zebrafish equivalent of mammalian NK cells.
224

225 Finally, a third cluster (C3) showed high expression of the myeloid lineage specific
226 transcription factor (TF) *spi1b* (Ward et al. 2003). These data suggested that cells in
227 Cluster 3 had a myeloid cell-like identity.

229 The clustering structure of our fish immune cells was further validated in a set of
230 more than 300 single cells from a third fish and additional cells from the first fish,
231 where despite much lower coverage due to external RNA contamination of the
232 samples, the separation between cells expressing the different markers (*cd4*, *cd8*,
233 *nitr*, *dicp* and *spi1b*) is clearly visible (Supplemental Fig. S7 and Supplemental
234 Methods).

236 In addition to distinct transcriptional states, FACS analysis revealed that cells in
237 different clusters differ in their light scattering properties (Fig. 2B). In particular, side
238 scattered light (SSC), which is positively correlated with subcellular granularity or
239 internal complexity, was 25% higher in C2 than in C1 (Wilcoxon test p = 1.4e-05).
240 This is consistent with NK-like cells possessing dense cytoplasmic granules (Yoder
241 and Litman 2011). In addition, SSC of cells in C3 was 203% higher than in the other
242 two clusters together (p = 1.6e-05). The high granularity of cells in C3 further
243 supports the hypothesis that these cells originate from a subpopulation of *lck+*
244 myeloid cells, such as granulocytes (see (Gibbings and Befus 2009) for similar
245 findings in mammals).

247 Since *lck:GFP+* cells were sorted randomly from spleen, the number of cells within
248 each of the clusters could be used as an estimate of the frequency of each sub-
249 population in the spleen in zebrafish. Similar to what is known from mouse (including
250 *Lck:gfp* transgenic mice (Shimizu et al. 2001)) and human, T cells were more
251 frequently found (65.4% of cells fall in C1) than NK-like cells (30.8% of cells fall in
252 C2).

254 **Differential expression analysis identifies both known and novel genes specific**
255 **for each cell type**
256 To identify genes specific for each cell population we performed differential
257 expression analysis of each cluster versus the other two (see Methods and
258 Supplemental Table S1).

8

259    The T cell signature genes *cd4*, *cd8*a, *ctla4*, and cd28, the transcription factors

260    *bcl11b* and *tcf7* and the cytokine/chemokine receptors *il10rb*, *ccr7*, *ccr9* and cxcr4

261    were within the most differentially expressed genes in Cluster 1 (Fig. 3). We also

262    identified many T cell specific genes that were uncharacterized or did not have an

263    informative name or description in the zebrafish genome for which we assigned a

264    putative name, based on sequence similarity searches. These included the *cd8* beta

265    chain (ENSDARG00000058682) whose expression is highly correlated with the

266    alpha chain *cd8a* within Cluster 1, cd28 (ENSDARG00000069978) and an

267    uncharacterized Ig-like protein (*ENSDARG00000098787*) related to CD7 antigen

268    (Fig. 3).

269

270    Mammalian NK cells kill target cells by either of two alternative pathways: the

271    perforin/granzyme secretory pathway or the death receptor pathway. Our analysis

272    revealed differential expression of several members of both pathways in C2. For

273    instance, differential expression of known innate immune receptors *nitr* and *dicp*, *syk*

274    kinase, multiple granzymes, perforins and NK lysins is linked with activation of the

275    secretory pathway whereas differential expression of FAS ligand (*faslg*) indicates

276    activation of the death-receptor-ligand pathway (Dybkaer et al. 2007) in NK-like cells

277    (Supplemental Table S1 and Fig. 3). Expression of these genes further shows that

278    zebrafish presumably resting NK-like cells transcriptionally resemble effector CD8 T

279    cells, as observed in mammals (Bezman et al. 2012).

280

281    We also observed a high expression level of cytokines and cytokine receptors. For

282    example, differentially expressed genes in Cluster 2 included the sphingosine 1-

283    phosphate receptor *s1p5* (*s1pr5a*, whose homolog in mammalian NK cells is required

284    for homing), the interleukin-2 receptor beta (IL2 induces rapid activation of

285    mammalian NK cells), *tnfsf14* (*Tnfsf14* tumor necrosis factor (ligand) superfamily,

286    member 14) and chemokines of the families *ccr38* and *ccr34*. In addition, we

287    detected differential expression of putative activating NK receptors' adaptors (ITAMs)

288    Fc receptor gamma subunit FcRγ (*fcer1g*) DAP10 (*hcst*), CD3ζ/cd247-like (*cd247l*)

289    and multiple putative transcription factors (Supplemental Table S1). Finally, within

290    the top differentially expressed genes of these NK-like cells we found putative

291    homologs of mammalian granzyme B that is expanded in ray-finned fish genomes

292    (ENSDARG00000078451, ENSDARG00000093990, ENSDARG00000055986, see

293    also Fig. 4), and many uncharacterized putative Immunoglobulin-like receptors and

9

294   cytokines, such as immunoglobulin V-set domain containing proteins or interleukin-8-

295   like domain containing chemokines (Table 1). Altogether, these results add

296   confidence in our proposed classification of these cells as putative fish NK-like cells.

297

298   Regarding cells in Cluster 3, the small number of cells within this cluster limits the

299   power of differential expression analysis. Nevertheless, within the most differentially

300   expressed genes in Cluster 3, we found two myeloid specific genes: the TF *spi1b*

301   and the granulocyte/macrophage colony-stimulating factor receptor beta (*csf2rb*).

302   Other differentially expressed genes included an Fc-receptor gamma like protein

303   (*fcer1gl*), *hck*, a member of the Src family of tyrosine kinases mostly expressed by

304   phagocytes in mammals and potentially implicated in signal transduction of Fc

305   receptors and degranulation (Guiet et al. 2008), complement factor B (*zgc.158446*),

306   and *id2* (a transcription factor interacting with *spi1b*), Fig. 3.

307

308   We next compared differentially expressed genes in each cell population to human

309   transcriptomic data of homogeneous FACS sorted immune cells (Watkins et al. 2009;

310   Chambers et al. 2007). For genes differentially expressed in Cluster 1, our results

311   show a significant enrichment in differentially expressed in human T cells (P=0.008,

312   see Methods). Similarly, the comparison of differentially expressed genes in Cluster

313   2 with human gene expression data confirmed a significant enrichment in NK specific

314   genes (P=0.009) thus supporting the conservation of a core transcriptional program

315   between mammalian and zebrafish NK-like cells (see Methods). Finally, differentially

316   expressed genes in Cluster 3 were weakly enriched in human Myeloid-specific genes

317   (odds ratio=5.2, P=0.06, see Methods).

318

319   **Functional divergence of duplicated immune genes in zebrafish**

320   Gene duplication is a common event in eukaryotic genomes and plays a major role in

321   functional divergence. To systematically explore this functional divergence in fish

322   immune genes, we collected all duplicated genes pre- and post-ray finned fish

323   speciation (see Methods). Interestingly, genes more recently duplicated (ray finned

324   fish-specific) show lower expression in our dataset. For example, 53% of pre-

325   speciation duplicated genes showed expression in *lck+* cells, compared to 41% of

326   post-speciation duplicated paralogs. As expected, pre-speciation duplicated immune

327   genes were more likely (94%) to functionally diverge (i.e. show differential expression

328   in the immune subpopulations, see Methods) compared to the more recent post-

10

329    speciation paralogs (62%). Ray finned fish-specific duplicated genes with conserved

330    expression patterns included, for instance, the NK receptors *nitr* that, although

331    expanded in zebrafish, have kept their cell type specificity. In contrast, other fish-

332    specific paralogs show distinct expression, suggesting possible neo-functionalization

333    events (see Fig. 4). NK-lysins (*nkl.2*, *nkl.3*, *nkl.4*) provide an interesting example of

334    recent functional divergence. In our data *nkl.4* was expressed in both Myeloid- and

335    NK-like cells. However, *nkl.3* was only expressed in Myeloid-like cells, while *nkl.2*

336    expression was restricted to NK-like cells (Fig. 3 and 4). A second example of

337    neofunctionalization is the Fc receptor gamma subunit (FcRγ), which in mouse and

338    human, is highly expressed in myeloid and NK cells (Tassi et al. 2006). In zebrafish

339    *lck+* cells, fcr_gamma (*fcer1g*) was expressed in Myeloid- and NK-like cells, while its

340    paralog fcr_gamma-like (*fcer1gl*) expression was restricted to the Myeloid -like cells

341    (Fig. 4). Other examples of such neo- or sub-functionalization of recently duplicated

342    paralogs are shown in Fig. 4 and Supplemental Table S2.

343

344    **NK specific genes show lower conservation than T cell genes from mammals**

345    **to bony fish**

346    The immune system is constantly adapting to new pathogens and changes in

347    virulence mechanisms, and hence is one of the most rapidly evolving biological

348    systems in vertebrates (Fumagalli et al. 2011; Kosiol et al. 2008). To explore the

349    evolution of the newly identified zebrafish genes specific for T, NK-like and myeloid-

350    like cells, we performed the same conservation analysis as in Fig. 1 (see Methods).

351    Consistently, among TM or secreted proteins, 76% of differentially expressed genes

352    in zebrafish T cells had orthologs in mouse or human compared to ~36% of

353    differentially expressed genes in zebrafish NK-like cells ($p<10^{-4}$, Fig. 5), suggesting

354    higher rates of gene turn-over in NKs across vertebrate evolution. Among non-TM or

355    secreted protein coding genes, proportion of orthologs was similar between T, NK-

356    like and myeloid-like cell specific genes (Fig. 5), as observed in mammalian species.

357

358    Examples of TM genes with no assigned orthologs beyond bony fish include putative

359    chemokine receptors (e.g. ENSDARG00000105363), nk lysins and the NK receptors

360    *nitr* and *dicp* among NK-like specific genes (see also Table 1), as well as the Ig-like

361    protein coding genes (e.g. ENSDARG00000098787, ENSDARG00000092106 and

362    ENSDARG00000092106, Supp. Table 1) among T cell specific genes. Although *lck+*

363    myeloid cells represent only a sub-population of fish myeloid cells, their genes

11

364  consistently show intermediate level of conservation between T and NK cell specific

365  genes, as observed for myeloid cells in mammals (Fig. 1 and 5).

366

367  When compared at the sequence identity level, the conserved TM genes specifically

368  expressed in either T or NK-like cells had lower sequence identity than other genes

369  across vertebrates ($p<10^{-6}$, Supplemental Fig. S2C). Moreover, as in mammals,

370  zebrafish TM genes were on average less conserved in NK-like than T cells across

371  vertebrates (p=0.03, Supplemental Fig. S2C). In contrast, cytoplasmic and nuclear T

372  cell specific displayed similar sequence identity compared to other genes

373  (Supplemental Fig. S2C)

374

375  **Discussion**

376  The availability of fully sequenced genomes in several vertebrate species has

377  enabled analysis of the evolution of the immune system based on orthology of known

378  mammalian immune cell markers. However, a comparison of immune subtypes at

379  the cellular and molecular level has progressed slowly, mainly owing to the lack of

380  suitable antibodies that mark distinct immune cell subpopulations in lower

381  vertebrates. Here we used scRNA-Seq of *lck:GFP* cells to characterise immune cell

382  subpopulations in zebrafish and examine their conservation in other vertebrate

383  species. Our work establishes scRNA-Seq as a powerful technique to study immune

384  cell types across vertebrate species.

385

386  Using single cell from two fish, we find three consistent clusters of cells, each

387  comprising cells from both fish. The most abundant population of cells in our data set

388  had a clear molecular T cell signature. The cells in this cluster showed differential

389  expression of hallmark genes important in regulation of T cell development and

390  signalling, suggesting a conserved transcriptional program from mammals to fish.

391  Within this population we were able to detect TCR V(D)J recombination in 22 cells

392  (Fig. 3 and Supplemental Figure 10). Interestingly, a single TCR recombinant was

393  found in each cell (Supplemental Table 5), which is consistent with allelic exclusion.

394  Although V(D)J recombination was clearly correlated with T cell identity, five cells

395  with evidence of V(D)J recombination fall in Cluster 2 and three of them show clear

396  expression NK genes. It is tempting to speculate that these cells could be NKT cells.

397  However, in mammals, the process of TCR rearrangement first initiates in

398  uncommitted haematopoietic progenitors before NK/DC/B/T divergence. Therefore,

12

399　incomplete rearrangements are also observed in subpopulations of non T cells, such

400　as NKs (Pilbeam et al. 2008). This could explain the presence of V(D)J

401　rearrangements in NK-like cells at the transcriptional level, as well as expression of

402　single V or J segments in cells in Cluster 2 and Cluster 3.　Moreover, TCRs

403　expressed by NK T cells present a limited diversity while here we found no evidence

404　for preferential use of specific segments among these cells.

405

406　In mammals LCK is expressed in both T and NK cells and in our dataset one

407　population of *lck+* cells resembled NK cells. Although NK-like cells were first

408　identified in catfish (Shen et al. 2004) over a decade ago, very little is known about

409　the NK cell transcriptome beyond mammals. Our data revealed that the proposed

410　bony fish NK receptors of the *nitr* family showed restricted expression in a distinct

411　cell subpopulation of NK-like cells which also expressed granzymes, perforins, NK

412　lysins (Pereiro et al. 2015; Moore et al. 2016), FAS ligand, TNFSF14, IL2 receptor

413　beta, the homolog of chemokine receptor CCR2, the sphingosine 1-phosphate

414　receptor S1P5 (required for homing of mammalian NK cells), specific transcription

415　factors and multiple novel putative NK-specific receptors and chemokines (Fig. 3 and

416　Table 1).

417

418　Throughout evolution, animals and plants have developed complex immune defence

419　mechanisms to combat microbial infections.　However, pathogens experience strong

420　selective pressure to evade host recognition and thus impose selective pressure on

421　the host to re-establish immunity. As a consequence, immune-related genes have

422　been preferential targets of positive selection in vertebrates (Kosiol et al. 2008; Yoder

423　and Litman 2011). Using a genome-wide unbiased approach based on transcriptomic

424　data from two mammalian and one bony fish species, we showed that a lower

425　fraction of orthologs and lower protein sequence identity are observed for NK TM

426　genes compared to other immune cell type specific TM genes, and especially T cell

427　TM genes, even though T and NK cells are functionally more related (e.g., TCD8 and

428　NK cytotoxicity upon MHCI recognition). Importantly, the trend is not only due to

429　known NK receptors (i.e. Ly49 in mouse, KIR in humans and *nitr/dicp's* in zebrafish,

430　Supplemental Fig. S1). This suggests that rapid evolution of NK TM genes is key for

431　their function in all vertebrates. As NK genes cannot undergo somatic

432　rearrangement, we propose that this fast evolution reflects, at least partly, a need for

433　NK cells to possess a diverse repertoire of species-specific germline encoded

13

434   receptors and associated proteins to perform their functions. In particular, both T and

435   NK cells recognize the fast evolving and highly polymorphic MHC molecules. While T

436   cells do so by rearranging their TCR sequence, NK cells possess an expanded

437   family of receptors. The fast evolution of these receptors may be the result of a need

438   to adapt to MHC rapid evolutionary changes. Our observations also support a model

439   of high gene turnover and faster evolution of immune TM/secreted genes, but at the

440   same time conservation of core cytoplasmic immune genes from zebrafish to

441   mammalian species (Fig. 1 and 5). As such, it supports zebrafish as an appropriate

442   model organism for immune cell intracellular signalling studies.

443

444   Overall, our work expands the analysis of immune subpopulations and their evolution

445   to lower vertebrates. To our knowledge, this is the first study to characterize T and

446   NK cells at the transcriptome level in a non-mammalian species and one of the first

447   studies to analyse NK cells' gene expression at the single cell level (see (Moore et al.

448   2016) for qPCR analysis of *TG(lck:GFP)* zebrafish single cells). We confirmed cell-

449   type specific expression of expected zebrafish T and NK cell genes and predicted

450   new markers of these two cell types. We further found significant variability in highly

451   expressed NK receptors and identified multiple cases of immune genes duplications

452   followed by neofunctionalization. Global conservation analysis revealed more rapid

453   turnover of NK specific TM genes compared to other immune cell, and especially T

454   cell specific genes in mammals and fish suggesting that this is an essential property

455   of immune cells.

456

457

458   **Methods**

459

460   **Conservation analysis of mouse and human immune cell differentially**

461   **expressed genes**

462   Orthologs of mouse and human protein-coding genes and their sequence identities,

463   as well as transmembrane domains and signal peptide predictions were downloaded

464   from BioMart / Ensembl Genes 82. For genes having multiple orthologs, we

465   considered their average sequence identity. Mouse and human NK, T cell, B-cell,

466   granulocyte and monocyte microarray gene expression datasets were obtained from

467   (Chambers et al. 2007) and (Watkins et al. 2009). First we pre-filtered genes with low

468   expression levels among these cell types using a threshold on normalized

14

469  expression level of 5 for the mouse data (16060 genes), and 8 for the human data

470  (8242 genes). CD8 and CD4 T-cells samples were merged into a T-cell group and

471  Monocyte and Granulocyte samples were merged into a myeloid cells group. We

472  then obtained differentially expressed genes in each group compared to the others,

473  using limma (version 3.28.14). Significantly differentially expressed genes

474  (Benjamini-Hochberg adjusted p-value < 0.01) were ordered based on expression

475  fold-change and the top 100 genes unique for each cell type were selected as

476  'signature genes' for downstream analysis (Supplemental Table S4, sheets 2 and 3).

477  Results were robust to different cut-offs for the top N differentially expressed

478  (Supplemental Methods and Supplemental Fig. S9). Human and mouse dN/dS ratios

479  (Supplemental Fig. S3) of one-to-one orthologs between these two species were

480  obtained from Ensembl version 82. The two protein groups enriched in 1)

481  transmembrane and secreted proteins and 2) cytoplasmic and nuclear proteins were

482  defined based on the presence of predicted trans-membrane domains and/or signal

483  peptide. Statistical significances of differences in sequence identity and dN/dS

484  differences were assessed using Wilcoxon/Mann-Whitney tests. Statistical

485  significances of differences in proportion of orthologs in 1) a specific species (e.g.

486  'human' point in Figure S2A) were assessed by comparison against a null-model

487  distribution generated from 10,000 random permutations of gene – cell type

488  specificity class pairs, and 2) globally across all species (as in Fig. S2 C), using

489  paired Wilcoxon signed rank test (to evaluate 'consistency' of the difference in

490  conservation patterns between two cell types).

491

492  **Zebrafish strains and maintenance**

493  Wild type (Tubingen Long Fin) and transgenic zebrafish *Tg(lck:EGFP)* lines were

494  maintained as previously described (Bielczyk-Maczyńska et al. 2014), in accordance

495  with EU regulations on laboratory animals.

496

497  **Single cell sorting and whole transcriptome amplification**

498  The spleens from two heterozygote *Tg(lck:EGFP)* adult fish from a different clutch

499  and different age (3 and 10 months old) and one adult wild-type fish were dissected

500  and carefully passed through a 40μm cell strainer using the plunger of a 1-mL

501  syringe and cells were collected in cold 1xPBS/5% FBS. The non-transgenic line was

502  used to set up the gating and exclude autofluorescent cells. Propidium iodide (PI)

503  staining was used to exclude dead cells. Individual cells were sorted, using a Becton

504  Dickinson Influx sorter with 488- and 561 nm lasers(Schulte et al. 2015) and

505  collected in single wells of 96 well plates containing 2.3 uL of 0.2 % Triton X-100

506  supplemented with 1 U/uL SUPERase In RNAse inhibitor (Ambion). The size,

507  granularity and level of fluorescence for each cell were simultaneously recorded.

508  Seven wells were filled with 50 cells each, from the second fish to compare single-

509  cell with bulk RNA-Seq (Supplemental Figure S5). The Smart-seq2 protocol (Picelli

510  et al. 2014) was used to amplify the whole transcriptome and prepare libraries.

511  Twenty-five cycles of PCR amplification were performed. Similar analysis was

512  performed on two additional plates of the first fish and four plates from a third fish,

513  including 5 wells with 50 cells (see Supplemental Methods and Supplemental Fig.

514  S7).

515

516  **Single cell RNA-Seq data processing**

517  Following Illumina HiSeq2000 sequencing (125bp paired-end reads), single-cell

518  RNA-Seq reads were quality trimmed and cleaned from Nextera adaptor contaminant

519  sequences using BBduck (*http://sourceforge. net/projects/bbmap*

520  ) with parameters *minlen=25 qtrim=rl trimq=10 ktrim=r k=25 mink=11 hdist=1 tbo*.

521  An average of 2.1 million paired-end reads were obtained per single-cell

522  (Supplemental Fig. 4 B). Next, gene expression levels were quantified as

523  $E_{i,j}=\log_2(TPM_{i,j}+1)$, where $TPM_{i,j}$ refers to transcript-per-million (TPM) for gene *i* in

524  sample *j*, as calculated by RSEM 1.2.19 (Li and Dewey 2011). RSEM (which uses

525  Bowtie 2.2.4 for alignment) was run in paired-end non strand-specific mode with

526  other parameters by default using the latest zebrafish genome assembly and

527  transcript annotations (GRCz10 / GCA_000002035.3) combined with eGFP

528  sequence appended as an artificial chromosome. For each single-cell, ~0.8 million

529  reads on average (with a median of 0.65 million) were mapped to the transcriptome

530  (Supplemental Fig. S3A). On average, 1240 expressed genes per cell were detected

531  (Supplemental Fig. S3C). Cells having less than 500 detected genes or less than

532  10,000 reads mapped to transcripts were excluded from further analyses.

533

534  **Transcriptome dimensionality reduction, batch effect removal and cell**

535  **clustering**

536  In order to visualize cell heterogeneity at the transcriptomic level, we used classical

537  multidimensional scaling (MDS, *aka* Principal Coordinates Analysis; as implemented

538  in R's *cmdscale* function) for dimensionality reduction (Fig. 2A, Supplemental Fig.

16

539 6A). MDS attempts to preserve distances between points generated from any
540 dissimilarity measure. Pearson's correlation coefficients (PCC) between full
541 transcriptional profiles were used to define cell-to-cell similarities, and 1-PCC was
542 then used as MDS's input dissimilarity measure. Similar low-dimensionality projection
543 was obtained with Principal Components Analysis (PCA) on the expression levels
544 ($E_{i,j}$) of the 1500 most variable genes (Fig. 6B).

545

546 To correct for batch effects and remove unwanted variation between the first and
547 second fish, we used ComBat function from R Bioconductor's *sva* package (Parker et
548 al. 2014). After this procedure, variation between individuals was minimal (Supp. Fig.
549 6G).

550 To identify different cell populations, we performed hierarchical clustering using
551 Ward's criteria (as implemented in R's *hclust* using *Ward.D2* method) applied on the
552 first four Principal Coordinates generated by the MDS. The choice of the components
553 was based on the eigenvalue decomposition of the MDS (Supp. Fig. 6C).
554 Eigenvalues decrease smoothly after the fourth Component, *i.e.* contributing less
555 significantly to the overall variability. The three cell clusters C1, C2 and C3 (Fig. 2A,
556 Fig6 E and F) were obtained by maximising the mean silhouette coefficient for
557 different number K of clusters (Supp. Fig. 6D).

558

559 **TCR reconstruction**
560 All four TCR loci (α, β, δ, and γ) and Rag-dependent variable diversity joining V(D)J
561 recombination are found in zebrafish (Langenau and Zon 2005). However, only the
562 beta chain locus was fully annotated (Meeker et al. 2010). To explore TCR
563 recombination in our immune cell subpopulations, we adapted the recent method of
564 (Stubbington et al. 2016). Synthetic beta chain sequences containing all possible
565 combinations the 52 V and 33 J germline segments were generated, with the
566 addition of 20 'N' ambiguity bases in the 5' end, 7 'N's between V and J segments
567 and 50 'N's at the 3' end to account for unknown leader, possible D and constant
568 sequences, respectively. RNA-seq reads from each cell were aligned against the
569 collection of synthetic TCR beta chain sequences independently using the Bowtie 2
570 aligner, with low penalties for introducing gaps into either the read or the reference
571 sequence or for aligning against N nucleotides (parameters '--no-unal -k 1 --np 0
572 --rdg 1,1 --rfg 1,1'). Next, reads aligning to synthetic sequences were used as input
573 to the Trinity RNA-seq assembly software (Grabherr et al. 2011) using its default

17

574  parameters for *de novo* assembly. Contigs assembled by Trinity were used as input

575  to NCBI IgBlast 1.4 (Ye et al. 2013) using parameters '-qcov_hsp_perc 90 -evalue

576  0.001 -ig_seqtype TCR -perc_identity 99' and providing zebrafish V, D and J

577  segments, and the resulting output were processed with a custom parsing script.

578  Contigs with no stop codons and matching both a V and a J segment with at least

579  90% sequence identity against corresponding germline segments, and where at least

580  90% of the germline segment was recovered, were considered evidence for TCR

581  beta chain V(D)J recombination.

582

583  **Differential expression analysis and marker gene discovery**

584  Estimated gene counts obtained from RSEM were used as input for *SCDE* R

585  package v1.99 (Kharchenko et al. 2014) that explicitly accounts for high rate of

586  dropout events in scRNA-Seq. Differential expression between each cluster versus

587  the other two was assessed using 100 randomizations (Supplemental Table S1).

588  To assess transcriptional conservation between mammalian and zebrafish immune

589  cell types, we used the previously defined sets of human top 100 differentially

590  expressed genes in T, NK and myeloid cells. We then compared the proportion of

591  zebrafish genes with orthologs in T cell, NK cell and myeloid signature genes within

592  the differentially expressed genes in each cluster versus non-differentially expressed

593  genes. Statistical significance was assessed using Fisher's exact test.

594

595  **Expression analysis of duplicated immune genes in zebrafish**

596  A list of paralogs in zebrafish was obtained from Ensembl Compara GeneTrees

597  (Vilella et al. 2009) (version 82). We defined two groups of protein coding genes: 1)

598  14,342 genes that underwent 'recent' duplication, whose most recent common

599  ancestor was mapped to ray-finned fish (Actinopterygii) or any of its child nodes

600  (Neopterygii, Otophysa, Clupeocephala, *Danio rerio*), and 2) 19,499 genes that

601  underwent 'early' duplication, where their most recent common ancestor was

602  mapped to bony vertebrates (Euteleostomi) or any of its parent taxa (Bilateria,

603  Chordata, Vertebrata). Many of these genes suffered multiple duplication events both

604  before and after the fish common ancestor. Therefore, to compare differences in

605  expression between these two groups, we did not include the set of overlapped

606  genes and obtained 3,235 unique recently duplicated genes and 8,609 unique early

607  duplicated genes. From these, 1315 (41%) and 4,569 (53%) were detected in our

608  data (genes with > 0 TPM in at least 1% of the cells).

18

609  For the analysis of expression pattern divergence, we searched pairs of paralogs
610  where both genes show some specific expression pattern (therefore, likely to have
611  an immune-related function) according to one of the following criteria: 1) within the
612  top 100 differentially expressed genes in Cluster 1, Cluster 2, or Cluster 3; 2) within
613  the top 100 differentially expressed genes in Cluster 2 and Cluster 3 versus Cluster 1
614  (*i.e.*, depleted in Cluster 1), Cluster 1 and Cluster 3 versus Cluster 2 (*i.e.*, depleted in
615  Cluster 2), or Cluster 1 and Cluster 2 versus Cluster 3 (*i.e.*, depleted in Cluster 3);  3)
616  expressed in all the three clusters (in at least 10% of the cells of each cluster). In the
617  latter case, we only considered pairs of paralogs where only one gene is expressed
618  in the 3 clusters, and the second is either specifically expressed or depleted from the
619  major clusters 1 or 2 (pairs of paralogs where both genes are expressed in all 3
620  clusters were not considered since most of them are not immune-related genes, and
621  cluster 3 is too small to accurately assess enrichment/depletion).
622  Next, we identified cases where both paralogs belong to the same expression pattern
623  group (duplicate genes with conserved expression pattern) and cases where they
624  differ (cases of neofunctionalization due to different expression patterns). For
625  recently duplicated genes we found 23 pairs with distinct expression patterns and 14
626  pairs with the same expression patterns (i.e. 62% of paralogs' neofunctionalization),
627  while for early duplicated genes we found 121 pairs with distinct patterns and 8 pairs
628  with the same expression patterns (i.e. 94% of paralogs' neofunctionalization), as
629  shown in Supplemental Table S2.
630
631  **Gene sequence conservation analysis of zebrafish differentially expressed**
632  **genes**
633  Orthologous genes of zebrafish in vertebrate species and their sequence identities
634  were downloaded from BioMart / Ensembl Genes 82. For comparisons between
635  differentially expressed genes between Cluster 1 (T cells), Cluster 2 (NK-like cells)
636  and Cluster 3 (myeloid-like cells) we chose the top 100 differentially expressed
637  genes after filtering by Z-score > 1 and sorting by fold-change (Supplemental Table
638  S1). Results were robust to different cut-offs (Supplemental Methods and
639  Supplemental Fig. S10). To assess ortholog conservation among non differentially
640  expressed genes, we first excluded lowly expressed genes from the analysis (those
641  where its expression level E was below the global mean of 0.46). The reason for this
642  is that we observed a bias of higher gene conservation among highly expressed
643  genes compared to lowly expressed genes. After this filter, conservation of

19

644    differentially expressed genes could be compared to that of non-differentially (but

645    having equivalent expression levels) genes as in Figure 5.

646

647

648

649    **Data access**

650    The data reported in this paper was deposited in ArrayExpress under the accession

651    number E-MTAB-4617

652

653    **Acknowledgements**

661
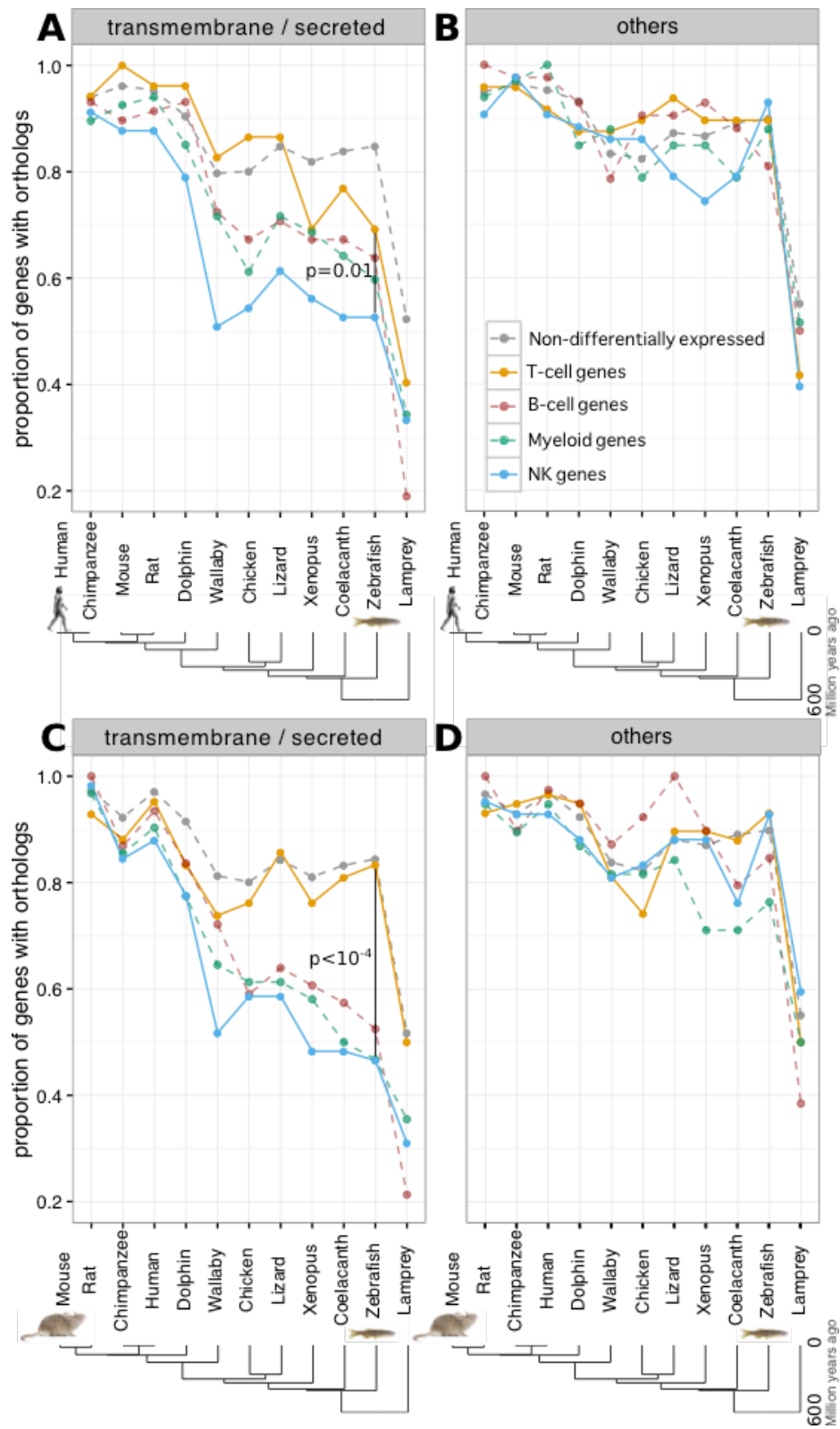
662    **Disclosure Declaration**

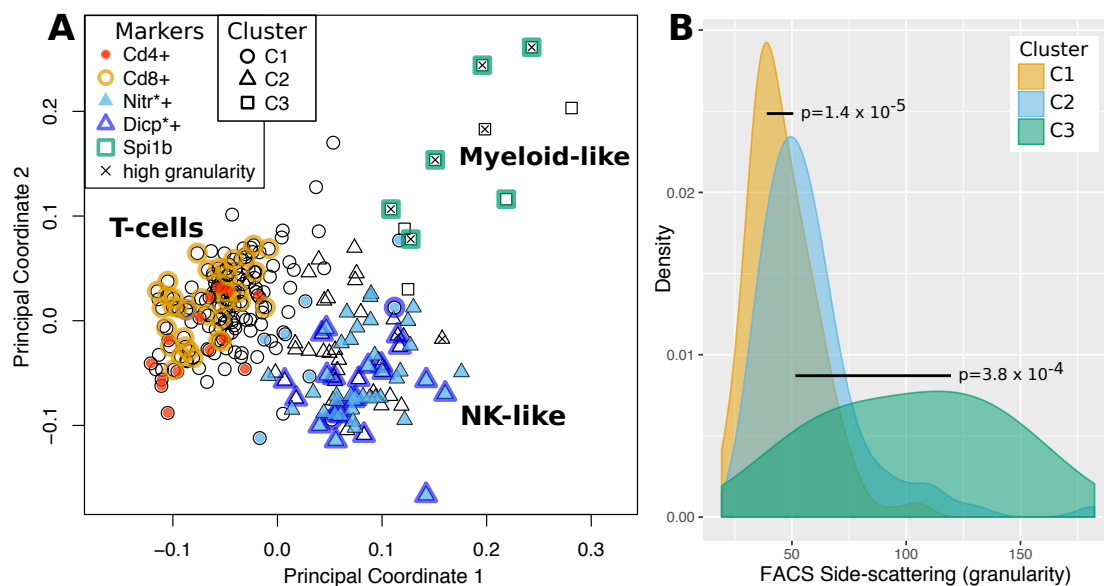663    The authors declare no competing financial interests.
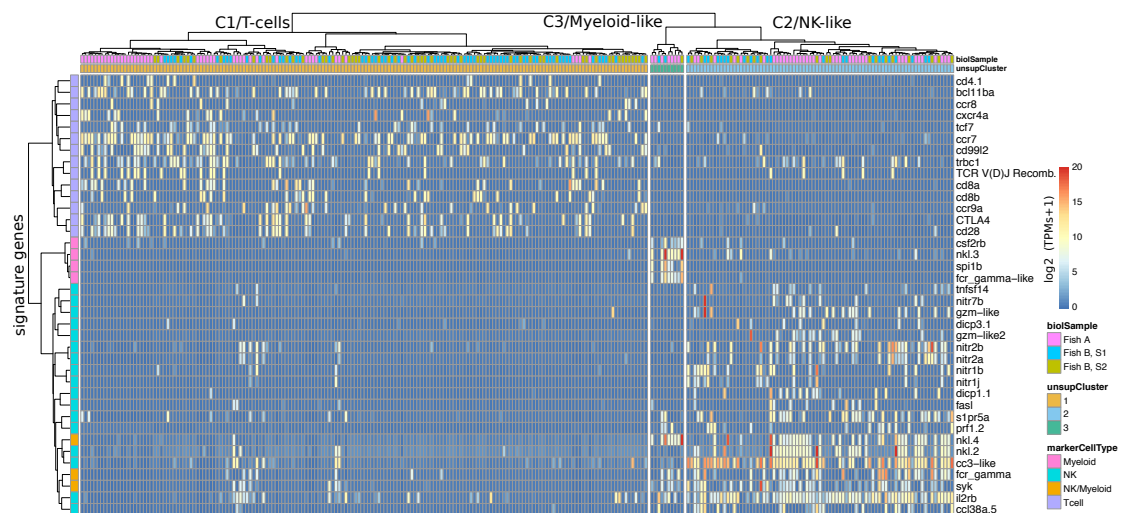
664

665

666 **Figures**

667

668

669 **Fig. 1: Conservation analysis of human and mouse genes differentially**
670 **expressed in major immune cell types. A,B:** Proportion of human genes specific
671 for distinct immune cell types (T, B, NK and myeloid cells) with orthologs in other
672 species. (A) shows the results for genes coding for transmembrane and secreted
673 proteins and (B) for cytoplasmic and nuclear proteins. **C,D**: Same analysis as in (A)
674 and (B) using mouse immune cell types' specific genes.

675
676
677



678
679 **Fig. 2: A:** Multidimensional scaling of zebrafish *lck*+ single-cell transcriptomes.
680 Unsupervised clustering revealed three subpopulations of cells: cluster 1 (C1),
681 cluster 2 (C2) and cluster 3 (C3), containing 65%, 31% and 4% of the cells,
682 respectively, and depicted with different symbols. A few examples of immune
683 signature genes are depicted (using an expression threshold of 5 TPMs): *cd4* and
684 *cd8* for T-cell specific genes, the innate immune receptors *nitr* and *dicp* for putative
685 NK-like specific genes in zebrafish and the myeloid associated transcription factor
686 *spi1b/pu.1* for myeloid-like cell specific genes. High granularity depicts cells with high
687 side scattered light. **B:** Distribution of side scattered light (proxy for cellular
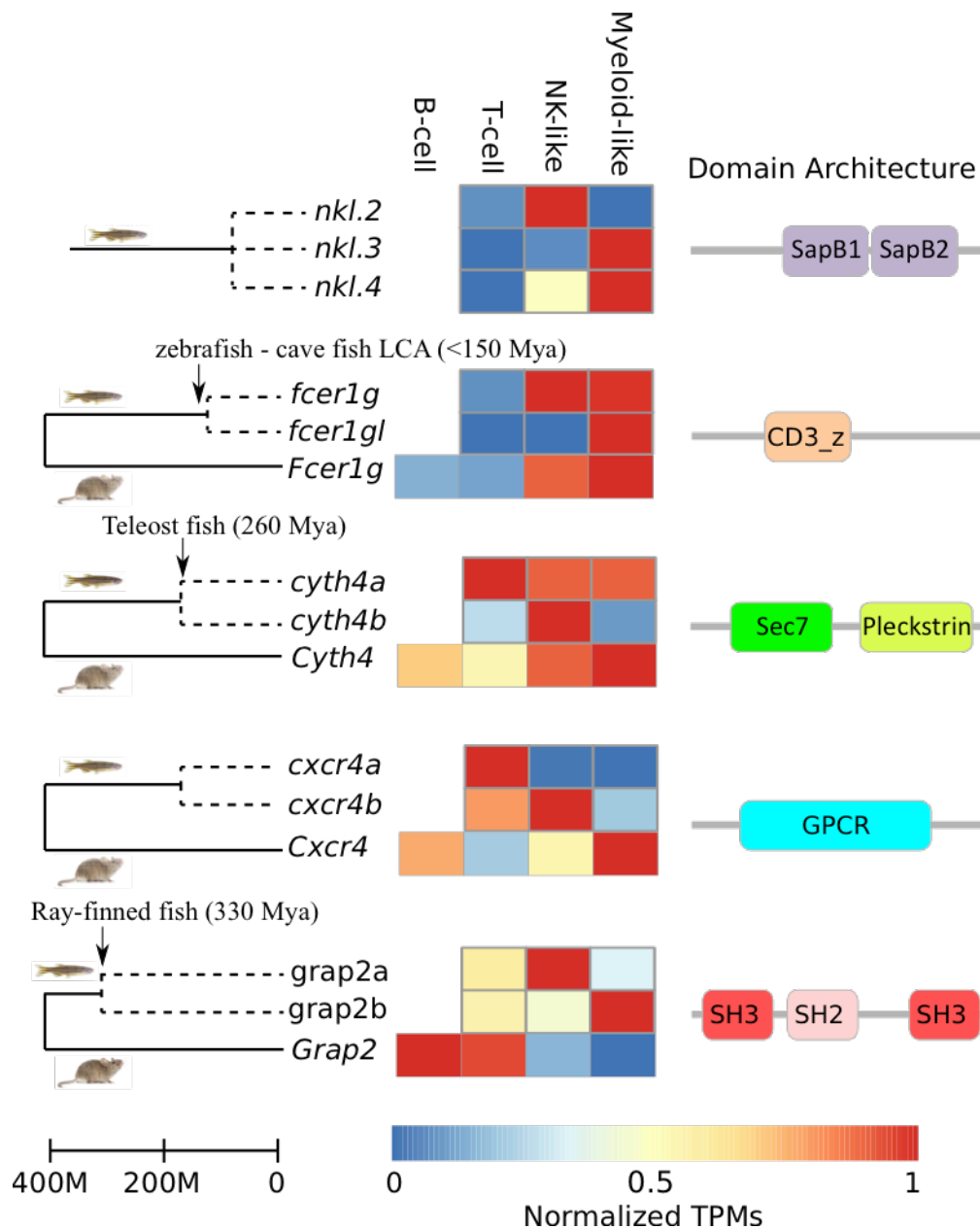688 granularity) for cells in each cluster.

689
690
691
692
693

22

694



695

**Fig. 3:** Heatmap showing the expression levels of some differentially expressed marker genes. Columns and rows represent cells and marker genes, respectively. Colours of the columns show the plates (top row) and the assigned clusters for each cell based on unsupervised whole-transcriptome clustering (second row, dendrogram shown on top). Colours of the rows (left-most column) indicate the known function of marker genes based on literature (T cell, NK or myeloid marker). The heatmap colour scale indicates the log2 transcripts per million (TPM, see Methods). Apart from a few cells in the T cell cluster that show expression of NK markers, the unsupervised whole-transcriptome clustering is very well recapitulated by expression of known and putative cell type markers.

706

707

708

**Fig. 4**: **Examples of ray-finned fish-specific duplicated genes with diverged expression patterns**. For genes with known mammalian orthologs, the expression in mouse is shown below. Times of gene duplication are indicated with arrows. Domain architectures were retrieved from PFAM (CD3_z: T cell surface glycoprotein CD3 zeta chain; GPCR: G-protein coupled receptor; SH2/3: Src-homology 2/3).

720

721

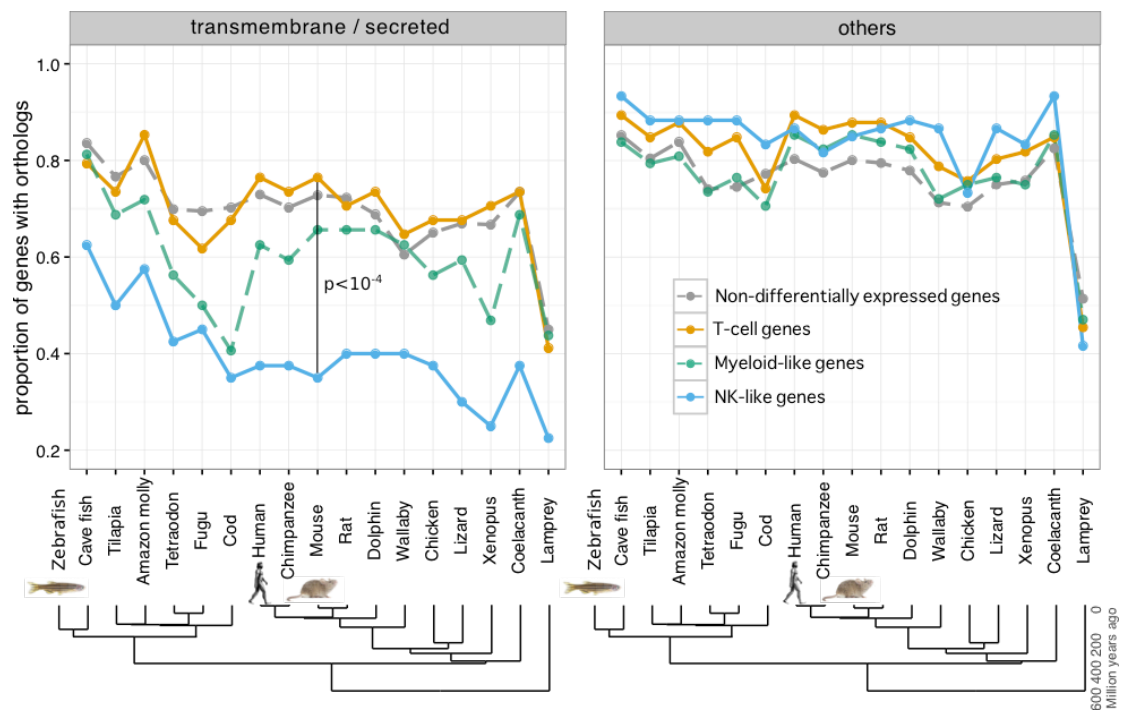

722

723

**Fig. 5: Conservation analysis of zebrafish immune genes across vertebrates**. The proportion of orthologs of protein-coding genes for non-differentially expressed genes (grey), differentially expressed genes in T cells (orange), NK-like cells (blue) and myeloid-like cells (green) are shown for both transmembrane or secreted proteins (left) and other proteins (right).

729

730

25

# Tables

| Ensembl Gene ID | Gene Name | Location | Domain Name | PFAM ID | Function |
|---|---|---|---|---|---|
| ENSDARG00000079353 | si:ch211-165d12.4 | Chr. 3 | | | |
| ENSDARG00000101860 | CABZ01034528.1 | Chr. 1 | | | |
| ENSDARG00000101860 | CABZ01034528.1 | Chr. 1 | Immunoglobulin V-set domain | PF07686 | |
| ENSDARG00000097065 | si:ch73-223p23.2 | Chr. 15 | | | |
| ENSDARG00000076358 | BX005329.1 | Chr. 22 | | | |
| ENSDARG00000103049 | CR392341.1 | Chr. 10 | Immunoglobulin V-set domain and Immunoglobulin C1-set | PF07654; PF07686 | Receptors |
| ENSDARG00000079387 | si:ch211-102c2.4 | Chr. 5 | Immunoglobulin-like domain | cd05716* | |
| ENSDARG00000071261 | BX248496.1 | Chr. 23 | | PF13895 | |
| ENSDARG00000090473 | si:ch211-269k10.5 | Chr. 16 | CD20-like family | PF04103 | |
| ENSDARG00000097847 | si:ch211-269k10.4 | Chr. 16 | | | |
| ENSDARG00000094002 | ccl34b.4 | Chr. 24 | | | |
| ENSDARG00000105263 | BX908792.2 | Chr. 7 | | | |
| ENSDARG00000041923 | ccl38.6 | Chr. 20 | Chemokine interleukin-8-like domain | PF00048 | Cytokines |
| ENSDARG00000071499 | cxcl32b.1 | Chr. 24 | | | |
| ENSDARG00000041835 | ccl38a.5 | Chr. 20 | | | |
| ENSDARG00000098656 | CT574575.1 | Chr. 24 | | | |
| ENSDARG00000095939 | si:ch73-226l13.2 | Chr. 11 | Interleukin-1 family | PF00340 | |
| ENSDARG00000101767 | si:dkey-183i3.6 | Chr. 21 | LAT2-like** | | Adaptors |
| ENSDARG00000093990 | si:ch211-165b19.8 | Chr. 9 | Peptidase S1 | PF00089 | Secreted Peptidases |

**Table 1**: List of novel zebrafish NK-specific membrane-bound or potentially secreted proteins, including putative receptors, cytokines and related proteins. Domain annotations were retrieved from PFAM except for * (NCBI Conserved domains database) and ** (PSI-BLAST search).

# References

Bailey M, Christoforidou Z, Lewis M. 2013. Evolution of immune systems: specificity and autoreactivity. *Autoimmun Rev* **12**: 643–647.

Bezman NA, Kim CC, Sun JC, Min-Oo G, Hendricks DW, Kamimura Y, Best JA, Goldrath AW, Lanier LL, Immunological Genome Project Consortium. 2012. Molecular definition of the identity and activation of natural killer cells. *Nat Immunol* **13**: 1000–1009.

Bielczyk-Maczyńska E, Serbanovic-Canic J, Ferreira L, Soranzo N, Stemple DL, Ouwehand WH, Cvejic A. 2014. A loss of function screen of identified genome-wide association study Loci reveals new genes controlling hematopoiesis. ed. L.I. Zon. *PLoS Genet* **10**: e1004450.

Boehm T. 2012. Evolution of vertebrate immunity. *Curr Biol* **22**: R722–32.

Brennecke P, Anders S, Kim JK, Kołodziejczyk AA, Zhang X, Proserpio V, Baying B, Benes V, Teichmann SA, Marioni JC, et al. 2013. Accounting for technical noise in single-cell RNA-seq experiments. *Nat Methods* **10**: 1093–1095.

Buettner F, Natarajan KN, Casale FP, Proserpio V, Scialdone A, Theis FJ, Teichmann SA, Marioni JC, Stegle O. 2015. Computational analysis of cell-to-cell heterogeneity in single-cell RNA-sequencing data reveals hidden subpopulations of cells. *Nat Biotechnol* **33**: 155–160.

Carrillo-Bustamante P, Keşmir C, de Boer RJ. 2016. The evolution of natural killer cell receptors. *Immunogenetics* **68**: 3–18.

Chambers SM, Boles NC, Lin K-YK, Tierney MP, Bowman TV, Bradfute SB, Chen AJ, Merchant AA, Sirin O, Weksberg DC, et al. 2007. Hematopoietic fingerprints: an expression database of stem cells and their progeny. *Cell Stem Cell* **1**: 578–591.

Dybkaer K, Iqbal J, Zhou G, Geng H, Xiao L, Schmitz A, d'Amore F, Chan WC. 2007. Genome wide transcriptional analysis of resting and IL2 activated human natural killer cells: gene expression signatures indicative of novel molecular signaling pathways. *BMC Genomics* **8**: 230.

Flajnik MF, Kasahara M. 2010. Origin and evolution of the adaptive immune system: genetic events and selective pressures. *Nat Rev Genet* **11**: 47–59.

Fumagalli M, Sironi M, Pozzoli U, Ferrer-Admetlla A, Ferrer-Admetlla A, Pattini L, Nielsen R. 2011. Signatures of environmental genetic adaptation pinpoint pathogens as the main selective pressure through human evolution. ed. J.M. Akey. *PLoS Genet* **7**: e1002355.

Gibbings D, Befus AD. 2009. CD4 and CD8: an inside-out coreceptor model for innate immune cells. *J Leukoc Biol* **86**: 251–259.

Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, Adiconis X, Fan L, Raychowdhury R, Zeng Q, et al. 2011. Full-length transcriptome assembly from RNA-Seq data without a reference genome. *Nat Biotechnol* **29**: 644–652.

Grün D, Lyubimova A, Kester L, Wiebrands K, Basak O, Sasaki N, Clevers H, van Oudenaarden A. 2015. Single-cell messenger RNA sequencing reveals rare intestinal cell types. *Nature* **525**: 251–255.

Guiet R, Poincloux R, Castandet J, Marois L, Labrousse A, Le Cabec V, Maridonneau-Parini I. 2008. Hematopoietic cell kinase (Hck) isoforms and phagocyte duties - from signaling and actin reorganization to migration and phagocytosis. *Eur J Cell Biol* **87**: 527–542.

Guo P, Hirano M, Herrin BR, Li J, Yu C, Sadlonova A, Cooper MD. 2009. Dual nature of the adaptive immune system in lampreys. *Nature* **459**: 796–801.

Haire RN, Cannon JP, O'Driscoll ML, Ostrov DA, Mueller MG, Turner PM, Litman RT, Litman GW, Yoder JA. 2012. Genomic and functional characterization of the diverse immunoglobulin domain-containing protein (DICP) family. *Genomics* **99**: 282–291.

Horton TL, Ritchie P, Watson MD, Horton JD. 1996. NK-like activity against allogeneic tumour cells demonstrated in the spleen of control and thymectomized Xenopus. *Immunol Cell Biol* **74**: 365–373.

Iwanami N. 2014. Zebrafish as a model for understanding the evolution of the vertebrate immune system and human primary immunodeficiency. *Exp Hematol* **42**: 697–706.

Jansen CA, van de Haar PM, van Haarlem D, van Kooten P, de Wit S, van Eden W, Viertlböck BC, Göbel TW, Vervelde L. 2010. Identification of new populations of

802    chicken natural killer (NK) cells. *Dev Comp Immunol* **34**: 759–767.

803    Kaufman CK, Mosimann C, Fan ZP, Yang S, Thomas AJ, Ablain J, Tan JL, Fogley
804        RD, van Rooijen E, Hagedorn EJ, et al. 2016. A zebrafish melanoma model
805        reveals emergence of neural crest identity during melanoma initiation. *Science*
806        **351**: aad2197–aad2197.

807    Kharchenko PV, Silberstein L, Scadden DT. 2014. Bayesian approach to single-cell
808        differential expression analysis. *Nat Methods* **11**: 740–742.

809    Kosiol C, Vinar T, da Fonseca RR, Hubisz MJ, Bustamante CD, Nielsen R, Siepel A.
810        2008. Patterns of positive selection in six Mammalian genomes. ed. M.H.
811        Schierup. *PLoS Genet* **4**: e1000144.

812    Langenau DM, Ferrando AA, Traver D, Kutok JL, Hezel J-PD, Kanki JP, Zon LI, Look
813        AT, Trede NS. 2004. In vivo tracking of T cell development, ablation, and
814        engraftment in transgenic zebrafish. *Proc Natl Acad Sci USA* **101**: 7369–7374.

815    Langenau DM, Zon LI. 2005. The zebrafish: a new model of T-cell and thymic
816        development. *Nat Rev Immunol* **5**: 307–317.

817    Li B, Dewey CN. 2011. RSEM: accurate transcript quantification from RNA-Seq data
818        with or without a reference genome. *BMC Bioinformatics* **12**: 323.

819    Macaulay IC, Svensson V, Labalette C, Ferreira L, Hamey F, Voet T, Teichmann SA,
820        Cvejic A. 2016. Single-Cell RNA-Sequencing Reveals a Continuous Spectrum of
821        Differentiation in Hematopoietic Cells. *Cell Rep* **14**: 966–977.

822    Meeker ND, Smith ACH, Frazer JK, Bradley DF, Rudner LA, Love C, Trede NS.
823        2010. Characterization of the zebrafish T cell receptor beta locus.
824        *Immunogenetics* **62**: 23–29.

825    Moore FE, Garcia EG, Lobbardi R, Jain E, Tang Q, Moore JC, Cortes M, Molodtsov
826        A, Kasheta M, Luo CC, et al. 2016. Single-cell transcriptional analysis of normal,
827        aberrant, and malignant hematopoiesis in zebrafish. *J Exp Med* **213**: 979–992.

828    Parker HS, Leek JT, Favorov AV, Considine M, Xia X, Chavan S, Chung CH, Fertig
829        EJ. 2014. Preserving biological heterogeneity with a permuted surrogate variable
830        analysis for genomics batch correction. *Bioinformatics* **30**: 2757–2763.

831    Pereiro P, Varela M, Diaz-Rosales P, Romero A, Dios S, Figueras A, Novoa B. 2015.
832        Zebrafish Nk-lysins: First insights about their cellular and functional
833        diversification. *Dev Comp Immunol* **51**: 148–159.

834    Picelli S, Faridani OR, Björklund AK, Winberg G, Sagasser S, Sandberg R. 2014.
835        Full-length RNA-seq from single cells using Smart-seq2. *Nat Protoc* **9**: 171–181.

836    Pilbeam K, Basse P, Brossay L, Vujanovic N, Gerstein R, Vallejo AN, Borghesi L.
837        2008. The ontogeny and fate of NK cells marked by permanent DNA
838        rearrangements. *J Immunol* **180**: 1432–1441.

839    Renshaw SA, Trede NS. 2012. A model 450 million years in the making: zebrafish
840        and vertebrate immunity. *Dis Model Mech* **5**: 38–47.

841    Saliba A-E, Westermann AJ, Gorski SA, Vogel J. 2014. Single-cell RNA-seq:
842        advances and future challenges. *Nucleic Acids Res* **42**: 8845–8860.

843    Schorpp M, Bialecki M, Diekhoff D, Walderich B, Odenthal J, Maischein H-M, Zapata
844        AG, Boehm T. 2006. Conserved functions of Ikaros in vertebrate lymphocyte
845        development: genetic evidence for distinct larval and adult phases of T cell
846        development and two lineages of B cells in zebrafish. *J Immunol* **177**: 2463–

847      2476.

848    Schulte R, Wilson NK, Prick JCM, Cossetti C, Maj MK, Göttgens B, Kent DG. 2015.
849        Index sorting resolves heterogeneous murine hematopoietic stem cell
850        populations. *Exp Hematol* **43**: 803–811.

851    Shapiro E, Biezuner T, Linnarsson S. 2013. Single-cell sequencing-based
852        technologies will revolutionize whole-organism science. *Nat Rev Genet* **14**: 618–
853        630.

854    Shay T, Jojic V, Zuk O, Rothamel K, Puyraimond-Zemmour D, Feng T, Wakamatsu
855        E, Benoist C, Koller D, Regev A, et al. 2013. Conservation and divergence in the
856        transcriptional programs of the human and mouse immune systems. *Proc Natl
857        Acad Sci USA* **110**: 2946–2951.

858    Shay T, Kang J. 2013. Immunological Genome Project and systems immunology.
859        *Trends Immunol* **34**: 602–609.

860    Shen L, Stuge TB, Bengtén E, Wilson M, Chinchar VG, Naftel JP, Bernanke JM,
861        Clem LW, Miller NW. 2004. Identification and characterization of clonal NK-like
862        cells from channel catfish (Ictalurus punctatus). *Dev Comp Immunol* **28**: 139–
863        152.

864    Shimizu C, Kawamoto H, Yamashita M, Kimura M, Kondou E, Kaneko Y, Okada S,
865        Tokuhisa T, Yokoyama M, Taniguchi M, et al. 2001. Progression of T cell lineage
866        restriction in the earliest subpopulation of murine adult thymus visualized by the
867        expression of lck proximal promoter activity. *Int Immunol* **13**: 105–117.

868    Stubbington MJT, Lönnberg T, Proserpio V, Clare S, Speak AO, Dougan G,
869        Teichmann SA. 2016. T cell fate and clonality inference from single-cell
870        transcriptomes. *Nat Methods*.

871    Sun JC, Beilke JN, Lanier LL. 2009. Adaptive immune features of natural killer cells.
872        *Nature* **457**: 557–561.

873    Sun JC, Lanier LL. 2009. Natural killer cells remember: an evolutionary bridge
874        between innate and adaptive immunity? *Eur J Immunol* **39**: 2059–2064.

875    Tassi I, Klesney-Tait J, Colonna M. 2006. Dissecting natural killer cell activation
876        pathways through analysis of genetic mutations in human and mouse. *Immunol
877        Rev* **214**: 92–105.

878    Trapnell C. 2015. Defining cell types and states with single-cell genomics. *Genome
879        Res* **25**: 1491–1498.

880    Vilella AJ, Severin J, Ureta-Vidal A, Heng L, Durbin R, Birney E. 2009.
881        EnsemblCompara GeneTrees: Complete, duplication-aware phylogenetic trees
882        in vertebrates. *Genome Res* **19**: 327–335.

883    Vu Manh T-P, Marty H, Sibille P, Le Vern Y, Kaspers B, Dalod M, Schwartz-Cornil I,
884        Quéré P. 2014. Existence of conventional dendritic cells in Gallus gallus
885        revealed by comparative gene expression profiling. *J Immunol* **192**: 4510–4517.

886    Ward AC, McPhee DO, Condron MM, Varma S, Cody SH, Onnebo SMN, Paw BH,
887        Zon LI, Lieschke GJ. 2003. The zebrafish spi1 promoter drives myeloid-specific
888        expression in stable transgenic fish. *Blood* **102**: 3238–3240.

889    Watkins NA, Gusnanto A, de Bono B, De S, Miranda-Saavedra D, Hardie DL,
890        Angenent WGJ, Attwood AP, Ellis PD, Erber W, et al. 2009. A HaemAtlas:
891        characterizing gene expression in differentiated human blood cells. *Blood* **113**:

892     e1–9.

893    Ye J, Ma N, Madden TL, Ostell JM. 2013. IgBLAST: an immunoglobulin variable
894         domain sequence analysis tool. *Nucleic Acids Res* **41**: W34–40.

895    Yoder JA. 2004. Investigating the morphology, function and genetics of cytotoxic
896         cells in bony fish. *Comp Biochem Physiol C Toxicol Pharmacol* **138**: 271–280.

897    Yoder JA, Litman GW. 2011. The phylogenetic origins of natural killer receptors and
898         recognition: relationships, possibilities, and realities. *Immunogenetics* **63**: 123–
899         141.

900    Yoder JA, Litman RT, Mueller MG, Desai S, Dobrinski KP, Montgomery JS, Buzzeo
901         MP, Ota T, Amemiya CT, Trede NS, et al. 2004. Resolution of the novel immune-
902         type receptor gene cluster in zebrafish. *Proc Natl Acad Sci USA* **101**: 15706–
903         15711.

904

905