

1 **Antiviral activity of a human placental protein of retroviral origin**

2

3 **John A. Frank<sup>1</sup>, Manvendra Singh<sup>1</sup>, Harrison B. Cullen<sup>1</sup>, Raphael A. Kirou<sup>1</sup>, Carolyn**

4 **B. Coyne<sup>2</sup>, Cédric Feschotte<sup>1</sup>**

5

6 **Affiliations**

7 **<sup>1</sup>Department of Molecular Biology and Genetics, Cornell University, Ithaca, NY,**

8 **USA**

9 **<sup>2</sup>Department of Pediatrics, University of Pittsburgh, Pittsburgh, PA, USA**

10

11 **Summary**

12 **Viruses circulating in wild and domestic animals pose a constant threat to human**

13 **health<sup>1</sup>. Identifying human genetic factors that protect against zoonotic infections**

14 **is a health priority. The RD-114 and Type-D retrovirus (RDR) interference group**

15 **includes infectious viruses that circulate in domestic cats and various Old World**

16 **monkeys (OWM), and utilize ASCT2 as a common target cell receptor<sup>2</sup>. While**

17 **human ASCT2 can mediate RDR infection in cell culture, it is unknown whether**

18 **humans and other hominoids encode factors that restrict RDR infection in**

19 **nature<sup>2,3</sup>. Here we test the hypothesis that Suppressyn, a truncated envelope**

20 **protein that binds ASCT2 and is derived from a human endogenous retrovirus<sup>4,5</sup>,**

21 **restricts RDR infection. Transcriptomics and regulatory genomics reveal that**

22 **Suppressyn expression initiates in the preimplantation embryo. Loss and gain of**

23 **function experiments in cell culture show Suppressyn expression is necessary and**

24 **sufficient to restrict RDR infection. Evolutionary analyses show Suppressyn was**  
25 **acquired in the genome of a common ancestor of hominoids and OWMs, but**  
26 **preserved by natural selection only in hominoids. Restriction assays using modern**  
27 **primate orthologs and reconstructed ancestral genes indicate that Suppressyn**  
28 **antiviral activity has been conserved in hominoids, but lost in most OWM. Thus in**  
29 **humans and other hominoids, Suppressyn acts as a restriction factor against**  
30 **retroviruses with zoonotic capacity. Transcriptomics data predict that other virus-**  
31 **derived proteins with potential antiviral activity lay hidden in the human genome.**

32

### 33 **Main**

34 Viral zoonosis poses a constant threat to human health and has led to devastating  
35 epidemics such as those caused by Influenza<sup>6</sup>, HIV<sup>7</sup>, Ebola<sup>8</sup>, and SARS  
36 coronaviruses<sup>9,10</sup>. Some zoonotic viruses have gained access to new host species by  
37 acquiring envelope (env) glycoproteins that mediate target-cell entry by binding to host  
38 cell surface receptors<sup>6,11</sup>. Notably, beta- and gamma-retroviruses have captured the so-  
39 called RDR env, which has enabled them to infect and transfer across diverse mammalian  
40 hosts<sup>11,12</sup>. For instance, the feline leukemia virus RD-114, an infectious endogenous  
41 retrovirus from the domestic cat, emerged from the *Felis catus* endogenous virus by  
42 acquiring the RDR *env* from the Baboon endogenous virus<sup>11</sup>. Because all RDR env bind  
43 to the highly conserved and broadly expressed amino acid transporter ASCT2 (also  
44 known as SLC1A5) to mediate cell entry, RDR env-mediated infection poses a zoonotic  
45 threat to humans<sup>2,13,14</sup>. Thus, it is critical to assess whether humans are equipped with  
46 mechanisms to protect against RDR infection.

47

48 Previous reports have shown that endogenous retroviral env are capable of restricting  
49 retroviral infection by a mechanism of receptor interference in multiple vertebrates,  
50 including chicken, mouse, sheep, and cat<sup>15</sup>. Some of these env-derived restriction factors  
51 acquired truncating mutations resulting in loss of their C-terminal membrane-anchoring  
52 transmembrane domain<sup>15</sup>, but retention of receptor-binding activity. Suppressyn  
53 (*SUPYN*) is a protein that is derived from a human endogenous retroviral env, lacks a  
54 transmembrane domain, and is expressed throughout human placenta development<sup>4,5</sup>.  
55 Previous *in vitro* studies have shown that SUPYN, like Syncytin-1 (SYN1), binds ASCT2  
56 and thereby modulates the fusogenic activity of SYN1 during placenta development<sup>4,5</sup>.  
57 Here we investigate whether *SUPYN* confers resistance to RDR infection.

58

59 **SUPYN embryonic expression is driven by pluripotency and placentation**  
60 **regulatory factors.**

61 To obtain a detailed view of *SUPYN* expression and regulation during human embryonic  
62 development, we analyzed publicly available scRNA-seq, ATAC-seq, DNase-seq and  
63 ChIP-seq datasets generated from preimplantation embryos and human embryonic stem  
64 cells (hESC) (**Supplementary Table 1**). *SUPYN* mRNA appears after the onset of  
65 embryonic genome activation at the eight-cell stage and peaks in morula (**Fig 1a**). By  
66 blastula formation, *SUPYN* expression persists in the inner cell mass, epiblast, ESCs,  
67 and in the trophectoderm, which gives rise to the placenta (**Fig 1a, Extended Data 1a**).  
68 Consistent with this expression pattern, the *SUPYN* locus is marked by open chromatin  
69 from 8-cell to blastocyst stages (**Extended Data 1b**). In hESCs, the *SUPYN* promoter

70 region is marked by H3K4me1 and H3K27ac histone modifications characteristic of  
71 'active' chromatin, and bound by core pluripotency (OCT4, NANOG, KLF4, SMAD1) and  
72 self-renewal (SRF, OTX2) transcription factors (**Fig 1b**). Together, these data indicate  
73 *SUPYN* is robustly expressed throughout early embryonic development and likely under  
74 the control of pluripotency factors. By contrast, we found little evidence for *SYN1*  
75 expression in preimplantation embryos and hESCs (**Fig 1a, Extended Data 1a**).

76

77 To examine *SUPYN* expression throughout placentation, we interrogated RNA-seq and  
78 ChIP-seq datasets generated from in vitro trophoblast (TB) differentiation models and  
79 placenta explants isolated at multiple stages of pregnancy (**Supplementary Table 1**).  
80 During hESC to TB differentiation, we observed that pluripotency factors NANOG and  
81 OCT4 occupying the *SUPYN* promoter region are replaced by trophoblast-specific  
82 transcription factors TFAP2A and GATA3 (**Fig 1b**). *SUPYN* expression likely persists  
83 through the TB differentiation process because *SUPYN* transcripts and active chromatin  
84 marks (H3K27ac, H3K4me3, H3K9ac) are maintained across all analyzed TB cell  
85 lineages (**Fig 1b**). By contrast, expression of other envelope-derived genes *SYN1*, *SYN2*,  
86 and *ERVV1/V2* is only detectable in differentiated trophoblasts (**Extended Data 1c**). We  
87 next mined scRNA-seq data generated from placenta at multiple developmental stages  
88 to examine the cell-type specific expression of *SUPYN* (**Supplementary Table 1**). After  
89 classifying cell clusters based on known markers (**Fig 1c, d, Extended Data 2a, b, c**),  
90 we found *SUPYN* and *ASCT2* expression specifically in the TB lineage (**Fig. 1e, f**;  
91 **Extended Data 2c, d**). Consistent with previous reports<sup>5</sup>, *SUPYN* expression was  
92 relatively high in cytotrophoblasts (CTB) and extra-villous trophoblasts (EVTB), but also

93 detectable in syncytiotrophoblasts (STB) (**Fig 1e, Extended Data 2c, d**). *SUPYN*  
94 expression in EVTB was maintained throughout placental development (**Fig 1f**). *SYN1*  
95 expression appears restricted to CTB and STB lineages (**Fig 1e, Extended Data 2c, d**),  
96 as previously reported<sup>5,16-18</sup>. To confirm these transcriptomic observations, we performed  
97 immunostaining of second (21w gestation) and third (31w gestation) trimester placenta  
98 with SUPYN antibody. These stains show SUPYN is widely expressed in STB, and  
99 perhaps CTB within the lumen of 2<sup>nd</sup> trimester placental villi (**Fig 1g, Extended Data 3**).  
100 Together these analyses indicate *SUPYN* is expressed throughout human fetal  
101 development and has a broader expression pattern than *SYN1*.

102

### 103 **SUPYN confers resistance to RD114env-mediated infection**

104 SUPYN expression during early embryonic and placental development, which coincides  
105 with that of ASCT2 (**Fig 1a**), suggests SUPYN may interact with ASCT2 throughout fetal  
106 development and confer RDR resistance to the developing embryo. To begin testing this  
107 hypothesis, we first examined whether human placenta-derived cell lines Jar and JEG3,  
108 and hESC H1 cells are resistant to RDR env-mediated infection. We generated HIV-GFP  
109 viral particles pseudotyped with either the feline RD114env (HIV-RD114env) or the  
110 glycoprotein G of vesicular stomatitis virus (HIV-VSVg, which uses the LDL receptor<sup>19</sup>) to  
111 monitor the level of infection in cell culture based on GFP expression (**Fig 2a, Extended**  
112 **Data 4**)<sup>20</sup>. These experiments revealed that Jar, JEG3, and H1 cells were susceptible to  
113 HIV-VSVg, as previously reported<sup>21-25</sup>, but resistant to HIV-RD114env infection (**Fig 2b,**  
114 **c**). Concurrently infected 293T cells, which do not express *SUPYN* (**Extended Data Fig**  
115 **1a**), were similarly susceptible to infection by HIV-RD114env and HIV-VSVg (**Fig 2b, c**).

116

117 To test whether SUPYN contributes to the HIV-RD114env resistance phenotype, we  
118 repeated these infection experiments in Jar cells engineered to stably express short  
119 hairpin RNAs depleting ~80% of *SUPYN* or *SYN1* mRNAs (**Extended Data Fig 5a**).  
120 Depletion of *SUPYN* in Jar cells (shSUPYN) resulted in a significant increase in  
121 susceptibility to HIV-RD114env infection (**Fig 2d**), but did not affect infection by HIV-VSVg  
122 (**Fig 2d**). Importantly, *SYN1* depletion from Jar cells did not increase susceptibility to HIV-  
123 RD114env infection (**Fig 2d**). These results suggest that *SUPYN* expression contributes  
124 to the RD114 resistance phenotype of Jar placental cells.

125

126 To account for possible off-target effects of *SUPYN*-targeting siRNAs, we transfected Jar-  
127 shSUPYN cells with two siRNA-resistant, HA-tagged SUPYN rescue constructs and  
128 examined their susceptibility to HIV-RD114env infection. Briefly, the siRNA-targeted  
129 SUPYN signal peptide sequence was replaced with a luciferase (SUPYN-lucSP) or  
130 modified signal peptide sequence (SUPYN-rescSP) that disrupts siRNA-binding but  
131 retains codon-identity. Transfection with either SUPYN-rescSP or SUPYN-lucSP restored  
132 a significant level of resistance to HIV-RD114env infection (**Fig 2e**). Western Blot analysis  
133 of transfected cell lysates showed SUPYN-rescSP was more abundantly expressed than  
134 SUPYN-lucSP (**Fig 2f**), which may account for the stronger HIV-RD114env resistance  
135 conferred by SUPYN-rescSP (**Fig 2e**). These results corroborate the notion that SUPYN  
136 restricts RD114env-mediated infection in Jar cells.

137

138 To test if SUPYN expression alone is sufficient to confer protection against RD114env-  
139 mediated infection, we transfected 293T cells with a SUPYN overexpression construct  
140 and subsequently infected with HIV-RD114env and HIV-VSVg respectively. As a positive  
141 control, we also transfected 293T cells with a RD114env overexpression construct, which  
142 is predicted to confer resistance to HIV-RD114env, but not to HIV-VSVg. Expression of  
143 either SUPYN or RD114env resulted in ~80% reduction in the level of HIV-RD114env  
144 infection (**Fig 2g, Extended Data Fig 6a**), but had no significant effect on HIV-VSVg  
145 infectivity (**Extended Data Fig 6b**). Taken together, our knockdown and overexpression  
146 experiments indicate SUPYN expression is both necessary and sufficient to confer  
147 resistance to RD114env-mediated infection.

148

#### 149 **SUPYN restricts RDR infection likely through receptor interference**

150 Our RD114env-specific resistance phenotype (**Extended Data Fig 6a, b**) strongly  
151 suggests SUPYN functions by receptor interference. If so, this protective effect should  
152 extend to infection mediated by other RDR env<sup>2,26,27</sup> since they all use ASCT2 as receptor.  
153 To test this prediction, we generated HIV-GFP reporter virions pseudotyped with Squirrel  
154 Monkey Retrovirus (SMRV) env (HIV-SMRVenv)<sup>2</sup> (**Fig 2a**) and infected 293T cells  
155 previously transfected with SUPYN, SMRVenv (as a positive control) or an empty vector.  
156 Cells expressing SUPYN or SMRVenv showed an ~80% reduction of HIV-SMRVenv  
157 infected cells (**Fig 2h**). Thus, SUPYN can restrict infection mediated by multiple RDRenv.

158

159 Another prediction of RDR restriction via receptor interference is that it should be a  
160 property of env proteins recognizing ASCT2, such as SUPYN, but not those binding other

161 cellular receptors. Consistent with this prediction, expressing HA-tagged env from  
162 amphotrophic murine leukemia virus (aMLV) or human endogenous retrovirus H, neither  
163 of which are expected to interact with ASCT2<sup>28-30</sup>, had no effect on HIV-RD114env nor  
164 HIV-VSVg infection in 293T cells, while HA-tagged SUPYN strongly restricted HIV-  
165 RD114env (**Extended Data Fig 6a, b**). Importantly, all tested env proteins were  
166 expressed at comparable levels (**Extended Data Fig 6c**). Furthermore, we observed that  
167 SUPYN overexpression did not significantly impact ASCT2 expression levels in 293T  
168 cells (**Extended Data Fig 6c**). This result suggests that if SUPYN acts by receptor  
169 interference, its interaction with ASCT2 does not result in receptor degradation, which is  
170 consistent with some instances of receptor interference<sup>31-33</sup>. In agreement with previous  
171 observations<sup>5</sup>, we noted that *SUPYN* knockdown in Jar cells resulted in the specific loss  
172 of a putative non-glycosylated isoform of ASCT2 (**Extended Data 5b**). We speculate that  
173 the presence of SUPYN-dependent non-glycosylated ASCT2 may be the result of SUPYN  
174 sterically interfering with the glycosylation machinery within the secretory pathway. It has  
175 not been reported whether ASCT2 glycosylation impacts RDR env-mediated infection in  
176 human cells, but it is known that receptor glycosylation may interfere with RDR infection  
177 in mouse and hamster cells<sup>34,35</sup>. Collectively, these observations converge on the model  
178 that SUPYN restricts against RDR infection through receptor interference (**Fig 2i**).

179

## 180 ***SUPYN* emerged in a catarrhine ancestor and evolved under functional constraint**

### 181 **in hominoids**

182 Little is known about the evolutionary origin of *SUPYN*. The gene was originally identified  
183 as derived from a copy of the HERV-Fb family of endogenous retroviruses (also known



184 as HERVH48<sup>36</sup>) located on human chromosome 21q22.3 with an ortholog in  
185 chimpanzee<sup>4</sup>. Using comparative genomics (see Methods), we found that this HERVH48  
186 element is shared at orthologous position across the genomes of all available hominoids  
187 (i.e. apes) and most Old World monkeys (OWM), but precisely lacking in New World  
188 monkeys and prosimians (**Fig 3a, Extended Data 7, 8**). Thus, the provirus copy that gave  
189 rise to *SUPYN* inserted in the common ancestor of catarrhine primates ~20-38 million  
190 years ago<sup>37</sup> (**Fig 3a**).

191

192 All primates with HERVH48 orthologs also share a nonsense mutation which would have  
193 truncated the ancestral env protein at site 185 in the common ancestor of catarrhine  
194 primates (**Fig 3a, Extended Data 8**). Hominoids share an additional nonsense mutation  
195 further truncating the protein to the 160-aa *SUPYN*-encoding ORF annotated in the  
196 human reference genome (**Fig 3a, Extended Data 8**). The *SUPYN* ORF is almost  
197 perfectly conserved in length across hominoids, but not in OWM where some species  
198 display further frameshifting and truncating mutations (**Extended Data 9**), suggesting  
199 *SUPYN* may have evolved under different evolutionary regimes in hominoids and OWMs.  
200 To test this idea, we analyzed the ratio ( $\omega$ ) of nonsynonymous (dN) to synonymous (dS)  
201 substitution rates using codeml<sup>38</sup>, which provides a measure of selective constraint acting  
202 on codons. Log-likelihood ratio tests comparing models of neutral evolution with selection  
203 indicate *SUPYN* evolved under purifying selection in hominoids ( $\omega = 0.38$ ;  $p = 1.47E-02$ ),  
204 but did not depart from neutral evolution in OWMs ( $\omega = 1.44$ ;  $p = 0.29$ ) (**Fig 3a**). For  
205 comparison, we performed the same type of analysis for *SYN1* and *SYN2*, two primate-  
206 specific *env*-derived genes thought to be involved in placentation<sup>17,39,40</sup>. Consistent with

207 previous reports<sup>41,42</sup>, we found that both *SYN1* ( $\omega = 0.64$ ;  $p = 1.80E-02$ ) and *SYN2* ( $\omega =$   
208  $0.29$ ;  $p = 3.22E-08$ ) evolved under purifying selection during hominoid evolution (**Fig. 3a**).  
209 In OWMs, *SYN2* also evolved under purifying selection ( $\omega = 0.22$ ,  $p = 2.78E-08$ ), while  
210 *SYN1* was lost through an ancestral deletion<sup>14</sup> (**Fig. 3a**). These results suggest that the  
211 level of functional constraint acting on *SUPYN* during hominoid evolution is comparable  
212 to that seen on other *env*-derived genes with placental function.

213

### 214 **SUPYN antiviral activity is conserved across hominoid primates**

215 To assess whether primate *SUPYN* orthologs have antiviral activity, we generated and  
216 transfected 293T cells with HA-tagged overexpression constructs for the orthologous  
217 *SUPYN* sequences of chimpanzee, siamang, African green monkey, pigtailed macaque,  
218 crab-eating macaque, rhesus macaque, and olive baboon and challenged these cells with  
219 HIV-RD114env virions. Both chimpanzee and siamang *SUPYN* proteins displayed  
220 antiviral activity with potency comparable to or greater than human *SUPYN*, respectively  
221 (**Fig 3b**). By contrast, only one (African green monkey) of the five OWM orthologous  
222 *SUPYN* proteins exhibited a modest but significant level of antiviral activity (**Fig 3b, c**).  
223 The lack of restriction activity for some of the OWM proteins may be attributed to their  
224 relatively low expression level in these human cells (**Fig 3b**) and/or their inability to bind  
225 the human ASCT2 receptor due to *SUPYN* sequence divergence (**Extended Data 8, 9**).  
226 To gain further insight into the evolutionary origins of *SUPYN* antiviral activity, we  
227 expressed *SUPYN* sequences predicted for the common ancestor of hominoid and OWM  
228 (**Extended Data 9**) (see Methods) and assayed their antiviral activity in 293T cells. Both  
229 ancestral proteins were expressed at levels comparable to human and green monkey

230 SUPYN and exhibited strong restriction activity (**Fig 3c**). These data indicate that SUPYN  
231 antiviral activity against RDR env-mediated infection is an ancestral trait preserved over  
232 ~20 million years of hominoid evolution, but apparently lost in some OWM lineages.

233

#### 234 **Truncated HERV env-coding sequences are expressed in diverse cell types**

235 The antiviral activity of SUPYN against RDR env-mediated infection raises the possibility  
236 that the human genome may harbor other *env*-derived protein-coding genes with antiviral  
237 function. To assess the potential pool of such sequences, we scanned the human  
238 genome for *env*-derived open reading frames (*envORF*) that minimally encode at least  
239 70 aa predicted to include the receptor-binding surface domain (see Methods). This  
240 search identified a total of 1,507 unique *envORFs*, including ~20 *env*-derived sequences  
241 currently annotated as human genes such as *SUPYN* and *SYN1*<sup>4,43,44</sup>. We then mined  
242 transcriptome datasets generated from human preimplantation embryos and various  
243 tissues (**Supplementary Table 1**), and observed that ~44% (668/1507) of *envORFs*  
244 showed evidence of RNA expression in at least one of the cell types surveyed (**Fig 4**).  
245 These analyses revealed three general insights about expressed non-annotated  
246 *envORFs*: (1) like known *env*-derived genes, *envORFs* exhibit tissue-specific expression  
247 patterns (**Fig 4; Extended Data 11, 12**); (2) the majority of *envORFs* are expressed  
248 during human fetal development or in stem and progenitor cells (**Fig 4; Extended Data**  
249 **11**); (3) *envORFs* are rarely expressed in differentiated tissues under normal conditions,  
250 with the exception of brain (**Fig 4; Extended Data 11, 12**). These analyses suggest that  
251 the human genome harbors a vast reservoir of *env*-derived sequences with coding and

252 receptor-binding potential, many of which are transcribed in a tissue-specific fashion and  
253 may encode receptor-binding antiviral peptides.

254

255

## 256 **Discussion**

257 Our expression and selection analyses firmly establish that *SUPYN* is a bona fide gene  
258 encoding a truncated envelope of retroviral origin that is expressed in the human  
259 preimplantation embryo and throughout placental development. Virologic assays in  
260 human cell culture show *SUPYN* is necessary and sufficient to confer resistance to RDR  
261 env-mediated infection, likely by competing for the receptor (ASCT2) utilized by this  
262 diverse set of retroviruses (see model, **Fig 2i**). The expression profile of *SUPYN* and the  
263 RD114 resistance phenotype of human ESCs and placental cells suggest *SUPYN*  
264 provides a layer of protection against RDR infection to the developing embryo, including  
265 the nascent germline. The observation that infectious and replication-competent RDRs  
266 are known to currently circulate in several mammals, both as exogenous or endogenous  
267 viruses, but not in hominoids<sup>2</sup> lends further support to a model in which *SUPYN* has  
268 contributed to hominoid resistance to RDR infection and possibly endogenization.

269

270 Like *SYN1*, *SUPYN* emerged in the common ancestor of catarrhine primates and was  
271 preserved by natural selection in hominoids. The parallel evolutionary path of *SYN1* and  
272 *SUPYN* and their pattern of expression in the placenta are compatible with a model in  
273 which *SUPYN* acts as a negative modulator of *SYN1* fusogenic activity during STB  
274 development<sup>4,5</sup>. The developmental and antiviral functions of *SUPYN* are not mutually

275 exclusive, but may even be interlocked. Syncytins, including SYN1, are fully functional  
276 envelopes that can be incorporated into heterologous retroviral particles and exosomes  
277 originating from the placenta<sup>50-55,45</sup>. Because ASCT2 is broadly expressed, SYN1-  
278 pseudotyped particles produced in the developing placenta have the potential to infiltrate  
279 a wide range of surrounding cell types. Thus, the physiological benefits afforded by  
280 Syncytins in promoting cell-cell fusion during STB development may have come with the  
281 cost of exposing the developing embryo to a wide variety of endogenous and exogenous  
282 invasive genetic elements, including but not limited to RDRs, that could be serendipitously  
283 enveloped by SYN1 throughout pregnancy. It is tempting to speculate that *SUPYN* has  
284 been maintained by natural selection to shield the developing embryo from the constant  
285 threat and adverse effects of SYN1-mediated infections. The conserved antiviral activity  
286 of ancestral hominoid and OWM *SUPYN* suggest resistance against RDR infection may  
287 have precipitated the initial retention of *SUPYN* in a catarrhine ancestor, and  
288 subsequently facilitated the domestication of *SYN1* in hominoids.

289

290 More broadly, this study serves as a proof of principle that truncated envelope peptides  
291 expressed from relics of ancient retroviruses integrated in the human genome can exert  
292 and retain antiviral activities for millions of years. We identified hundreds of candidate  
293 env-derived genes in the human genome that may encode peptides with receptor-binding  
294 activity and antiviral function. Furthermore, Gag (capsid)-derived proteins encoded by  
295 endogenous retroviruses are also capable of retroviral restriction<sup>46</sup>. Thus, it is possible  
296 that our genome holds a vast reservoir of retrovirus-derived proteins with protective  
297 activity against various zoonotic agents, including non-retroviral pathogens.

298

299 **Main References**

- 300 1. Warren, C. J. & Sawyer, S. L. How host genetics dictates successful viral  
301 zoonosis. *PLoS Biol.* **17**, e3000217 (2019).
- 302 2. Sinha, A. & Johnson, W. E. ScienceDirect Retroviruses of the RDR superinfection  
303 interference group: ancient origins and broad host distribution of a promiscuous  
304 Env gene. *Current Opinion in Virology* **25**, 105–112 (2017).
- 305 3. Montiel, N. A. An updated review of simian betaretrovirus (SRV) in macaque  
306 hosts. *J. Med. Primatol.* **39**, 303–314 (2010).
- 307 4. Sugimoto, J., Sugimoto, M., Bernstein, H., Jinno, Y. & Schust, D. A novel human  
308 endogenous retroviral protein inhibits cell-cell fusion. *Sci Rep* **3**, 1462–8 (2013).
- 309 5. Sugimoto, J. *et al.* Suppressyn localization and dynamic expression patterns in  
310 primary human tissues support a physiologic role in human placentation. *Sci Rep*  
311 **9**, 19502–12 (2019).
- 312 6. Mostafa, A., Abdelwhab, E. M., Mettenleiter, T. C. & Pleschka, S. Zoonotic  
313 Potential of Influenza A Viruses: A Comprehensive Overview. *Viruses* **10**, 497  
314 (2018).
- 315 7. Sharp, P. M. & Hahn, B. H. Origins of HIV and the AIDS pandemic. *Cold Spring*  
316 *Harb Perspect Med* **1**, a006841–a006841 (2011).
- 317 8. Baseler, L., Chertow, D. S., Johnson, K. M., Feldmann, H. & Morens, D. M. The  
318 Pathogenesis of Ebola Virus Disease. *Annu Rev Pathol* **12**, 387–418 (2017).
- 319 9. Song, Z. *et al.* From SARS to MERS, Thrusting Coronaviruses into the Spotlight.  
320 *Viruses* **11**, 59 (2019).

- 321 10. Andersen, K. G., Rambaut, A., Lipkin, W. I., Holmes, E. C. & Garry, R. F. The  
322 proximal origin of SARS-CoV-2. *Nat. Med.* **100**, 1–3 (2020).
- 323 11. Henzy, J. E. & Johnson, W. E. Pushing the endogenous envelope. *Philos. Trans.*  
324 *R. Soc. Lond., B, Biol. Sci.* **368**, 20120506–20120506 (2013).
- 325 12. Diehl, W. E., Patel, N., Halm, K. & Johnson, W. E. Tracking interspecies  
326 transmission and long-term evolution of an ancient retrovirus using the genomes  
327 of modern mammals. *Elife* **5**, e12704 (2016).
- 328 13. Green, B. J., Lee, C. S. & Rasko, J. E. J. Biodistribution of the RD114/mammalian  
329 type D retrovirus receptor, RDR. *J Gene Med* **6**, 249–259 (2004).
- 330 14. Haeussler, M. *et al.* The UCSC Genome Browser database: 2019 update. *Nucleic*  
331 *Acids Res.* **47**, D853–D858 (2019).
- 332 15. Johnson, W. E. Origins and evolutionary consequences of ancient endogenous  
333 retroviruses. *Nat. Rev. Microbiol.* **72**, 5955 (2019).
- 334 16. Holder, B. S., Tower, C. L., Abrahams, V. M. & Aplin, J. D. Syncytin 1 in the human  
335 placenta. *Placenta* **33**, 460–466 (2012).
- 336 17. Mi, S. *et al.* Syncytin is a captive retroviral envelope protein involved in human  
337 placental morphogenesis. *Nature* **403**, 785–789 (2000).
- 338 18. Frendo, J. L. *et al.* Direct Involvement of HERV-W Env Glycoprotein in Human  
339 Trophoblast Cell Fusion and Differentiation. *Molecular and Cellular Biology* **23**,  
340 3566–3574 (2003).
- 341 19. Finkelshtein, D., Werman, A., Novick, D., Barak, S. & Rubinstein, M. LDL receptor  
342 and its family members serve as the cellular receptors for vesicular stomatitis  
343 virus. *Proc. Natl. Acad. Sci. U.S.A.* **110**, 7306–7311 (2013).

- 344 20. Sandrin, V., Muriaux, D., Darlix, J. L. & Cosset, F. L. Intracellular Trafficking of  
345 Gag and Env Proteins and Their Interactions Modulate Pseudotyping of  
346 Retroviruses. *J. Virol.* **78**, 7153–7164 (2004).
- 347 21. Dottori, M., Tay, C. & Hughes, S. M. Neural development in human embryonic  
348 stem cells-applications of lentiviral vectors. *J. Cell. Biochem.* **112**, 1955–1962  
349 (2011).
- 350 22. Gropp, M. *et al.* Stable genetic modification of human embryonic stem cells by  
351 lentiviral vectors. *Mol. Ther.* **7**, 281–287 (2003).
- 352 23. Sakata, M. *et al.* Analysis of VSV pseudotype virus infection mediated by rubella  
353 virus envelope proteins. *Sci Rep* **7**, 11607 (2017).
- 354 24. Delorme-Axford, E. *et al.* Human placental trophoblasts confer viral resistance to  
355 recipient cells. *Proc. Natl. Acad. Sci. U.S.A.* **110**, 12048–12053 (2013).
- 356 25. Vidricaire, G., Tardif, M. R. & Tremblay, M. J. The low viral production in  
357 trophoblastic cells is due to a high endocytic internalization of the human  
358 immunodeficiency virus type 1 and can be overcome by the pro-inflammatory  
359 cytokines tumor necrosis factor-alpha and interleukin-1. *J. Biol. Chem.* **278**,  
360 15832–15841 (2003).
- 361 26. Johnson, W. E. Endogenous Retroviruses in the Genomics Era. *Annu Rev Virol*  
362 **2**, 135–159 (2015).
- 363 27. Sommerfelt, M. A. & Weiss, R. A. Receptor interference groups of 20 retroviruses  
364 plating on human cells. *Virology* **176**, 58–69 (1990).
- 365 28. de Parseval, N., Casella, J.-F., Gressin, L. & Heidmann, T. Characterization of the  
366 Three HERV-H Proviruses with an Open Envelope Reading Frame



- 367 Encompassing the Immunosuppressive Domain and Evolutionary History in  
368 Primates. *Virology* **279**, 558–569 (2001).
- 369 29. van Zeijl, M. *et al.* A human amphotropic retrovirus receptor is a second member  
370 of the gibbon ape leukemia virus receptor family. *Proc Natl Acad Sci USA* **91**,  
371 1168–1172 (1994).
- 372 30. Miller, D. G., Edwards, R. H. & Miller, A. D. Cloning of the cellular receptor for  
373 amphotropic murine retroviruses reveals homology to that for gibbon ape  
374 leukemia virus. *Proc Natl Acad Sci USA* **91**, 78–82 (1994).
- 375 31. Liu, M. & Eiden, M. V. The receptors for gibbon ape leukemia virus and  
376 amphotropic murine leukemia virus are not downregulated in productively infected  
377 cells. *Retrovirology* **8**, 53–14 (2011).
- 378 32. Jobbagy, Z., Garfield, S., Baptiste, L., Eiden, M. V. & Anderson, W. B. Subcellular  
379 redistribution of Pit-2 P(i) transporter/amphotropic leukemia virus (A-MuLV)  
380 receptor in A-MuLV-infected NIH 3T3 fibroblasts: involvement in superinfection  
381 interference. *J. Virol.* **74**, 2847–2854 (2000).
- 382 33. Kim, J. W. & Cunningham, J. M. N-linked glycosylation of the receptor for murine  
383 ecotropic retroviruses is altered in virus-infected cells. *J. Biol. Chem.* **268**, 16316–  
384 16320 (1993).
- 385 34. Marin, M., Taylor, C. S., Nouri, A. & Kabat, D. Sodium-dependent neutral amino  
386 acid transporter type 1 is an auxiliary receptor for baboon endogenous retrovirus.  
387 *J. Virol.* **74**, 8085–8093 (2000).
- 388 35. Marin, M., Lavillette, D., Kelly, S. M. & Kabat, D. N-linked glycosylation and  
389 sequence changes in a critical negative control region of the ASCT1 and ASCT2

- 390 neutral amino acid transporters determine their retroviral receptor functions. *J.*  
391 *Virology* **77**, 2936–2945 (2003).
- 392 36. Hubley, R. *et al.* The Dfam database of repetitive DNA families. *Nucleic Acids*  
393 *Res.* **44**, D81–9 (2016).
- 394 37. Perelman, P. *et al.* A molecular phylogeny of living primates. *PLoS Genet* **7**,  
395 e1001342 (2011).
- 396 38. Yang, Z. PAML 4: phylogenetic analysis by maximum likelihood. *Molecular*  
397 *Biology and Evolution* **24**, 1586–1591 (2007).
- 398 39. Malassiné, A. *et al.* Expression of the fusogenic HERV-FRD Env glycoprotein  
399 (syncytin 2) in human placenta is restricted to villous cytotrophoblastic cells.  
400 *Placenta* **28**, 185–191 (2007).
- 401 40. Mallet, F. *et al.* The endogenous retroviral locus ERVWE1 is a bona fide gene  
402 involved in hominoid placental physiology. *Proc Natl Acad Sci USA* **101**, 1731–  
403 1736 (2004).
- 404 41. Bonnaud, B. *et al.* Evidence of selection on the domesticated ERVWE1 env  
405 retroviral element involved in placentation. *Molecular Biology and Evolution* **21**,  
406 1895–1901 (2004).
- 407 42. de Parseval, N. *et al.* Comprehensive search for intra- and inter-specific sequence  
408 polymorphisms among coding envelope genes of retroviral origin found in the  
409 human genome: genes and pseudogenes. *BMC Genomics* **6**, 117–11 (2005).
- 410 43. Heidmann, O. *et al.* HEMO, an ancestral endogenous retroviral envelope protein  
411 shed in the blood of pregnant women and expressed in pluripotent stem cells and  
412 tumors. *Proc. Natl. Acad. Sci. U.S.A.* **114**, E6642–E6651 (2017).

413 44. Heidmann, A. D. C. L. T., Lavalie, C. & Heidmann, T. From ancestral infectious  
414 retroviruses to bona fide cellular genes: Role of the captured syncytins in  
415 placentation. *Placenta* **33**, 663–671 (2012).

416 45. Tang, Y. *et al.* Endogenous Retroviral Envelope Syncytin Induces HIV-1  
417 Spreading and Establishes HIV Reservoirs in Placenta. *Cell Rep* **30**, 4528–  
418 4539.e4 (2020).

419 46. Frank, J. A. & Feschotte, C. Co-option of endogenous viral sequences for host  
420 cell function. *Current Opinion in Virology* **25**, 81–89 (2017).

421

422

423

424

425

426

427

428

429

430

431

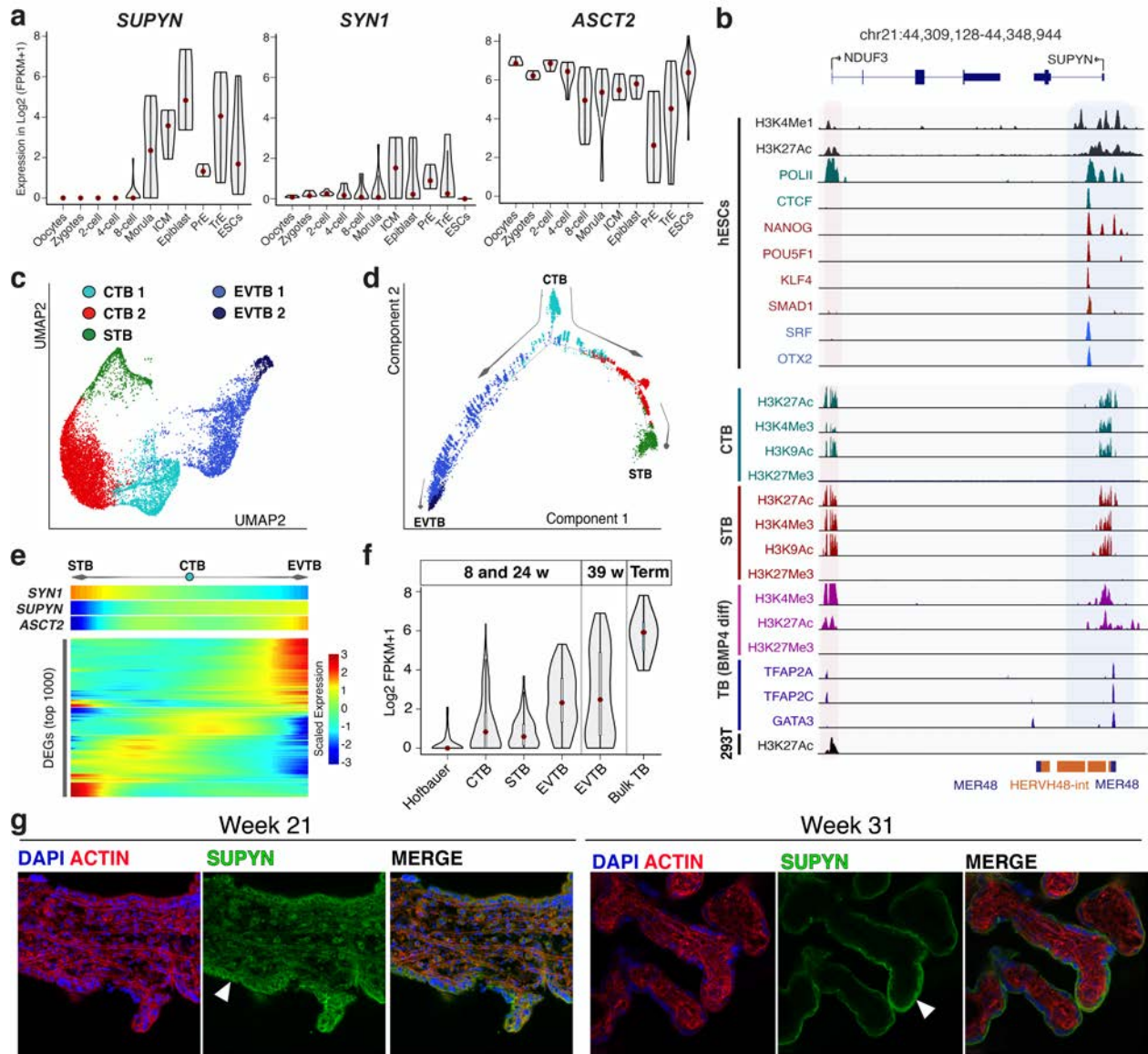
432

433

434

435

436 **Figures**



437

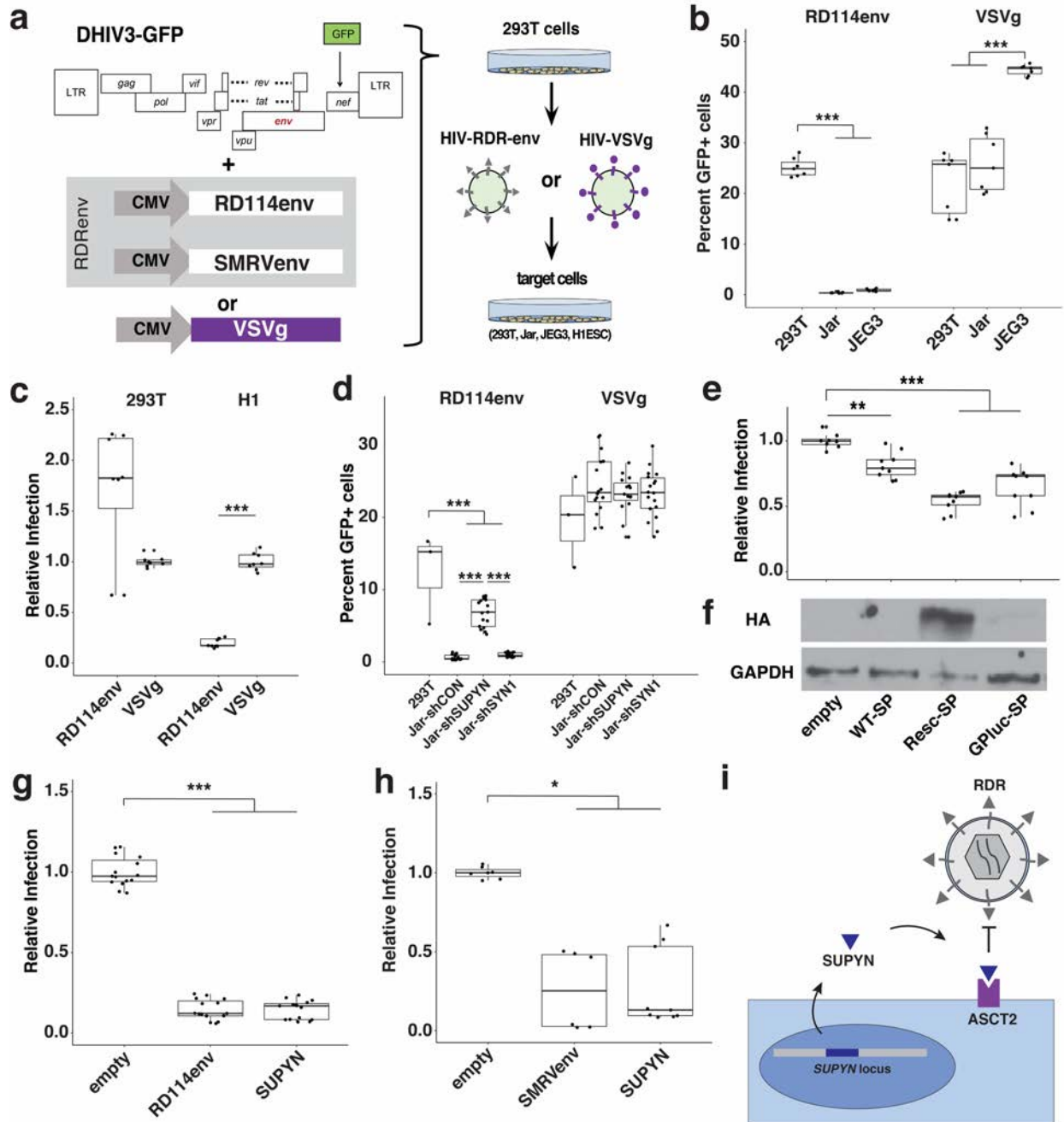
438 **Figure 1: Pluripotency and placentation regulatory factors drive SUPYN expression during fetal**  
 439 **development.**

440 **(a)** Violin plots summarizing *SUPYN*, *SYN1* and *ASCT2* expression in human preimplantation  
 441 embryo and ESC single-cell RNA-seq data. **(b)** Genome browser view of the regulatory elements  
 442 around the *SUPYN* locus in hESCs, cyto- (CTB), syncytio- (STB), BMP4 differentiated trophoblasts  
 443 (TB), and 293T cells. ChIP-seq profiles are shown for indicated transcription factors and histone

444 modifications with shaded area highlighting regions of active chromatin. **(c)** UMAP visualization  
445 of TB cell clusters, shown in **Extended Data 2**. **(d)** Monocle2 pseudotime analysis of cell clusters  
446 in **c** illustrates the developmental trajectory of CTBs that give rise to STB and EVTB respectively.  
447 Color codes in **c** and **d** denote cell identity. **(e)** Heatmap represents the top 1000 differentially  
448 expressed genes (row) of single cells (column) sorted according to pseudotime analyzed in **c** and  
449 **d**. Cells are ordered according to the pseudotime progression of CTB (middle) to STB (left) and  
450 EVTB (right). SYN1, SUPYN, and ASCT2 were fetched from the heatmap below. **(g)** Confocal  
451 microscopy of placental villi explants stained for SUPYN (green), Actin (red) and DAPI. STBs are  
452 marked by arrowheads.

453

454



455

456 **Figure 2: SUPYN confers resistance to RDR env-mediated infection**

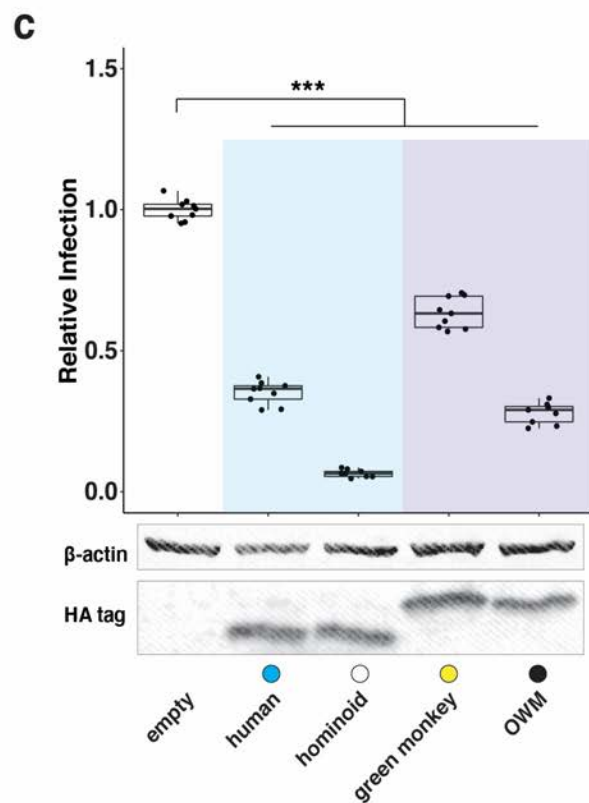
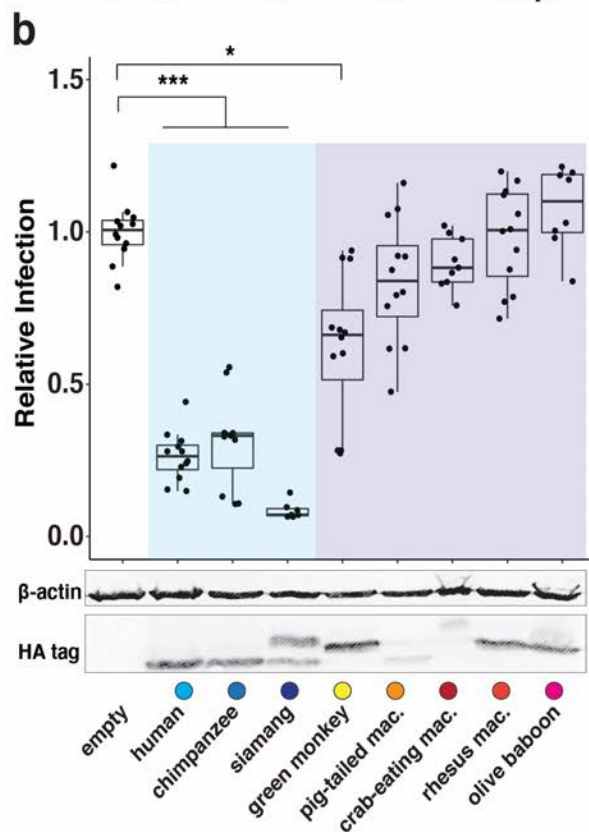
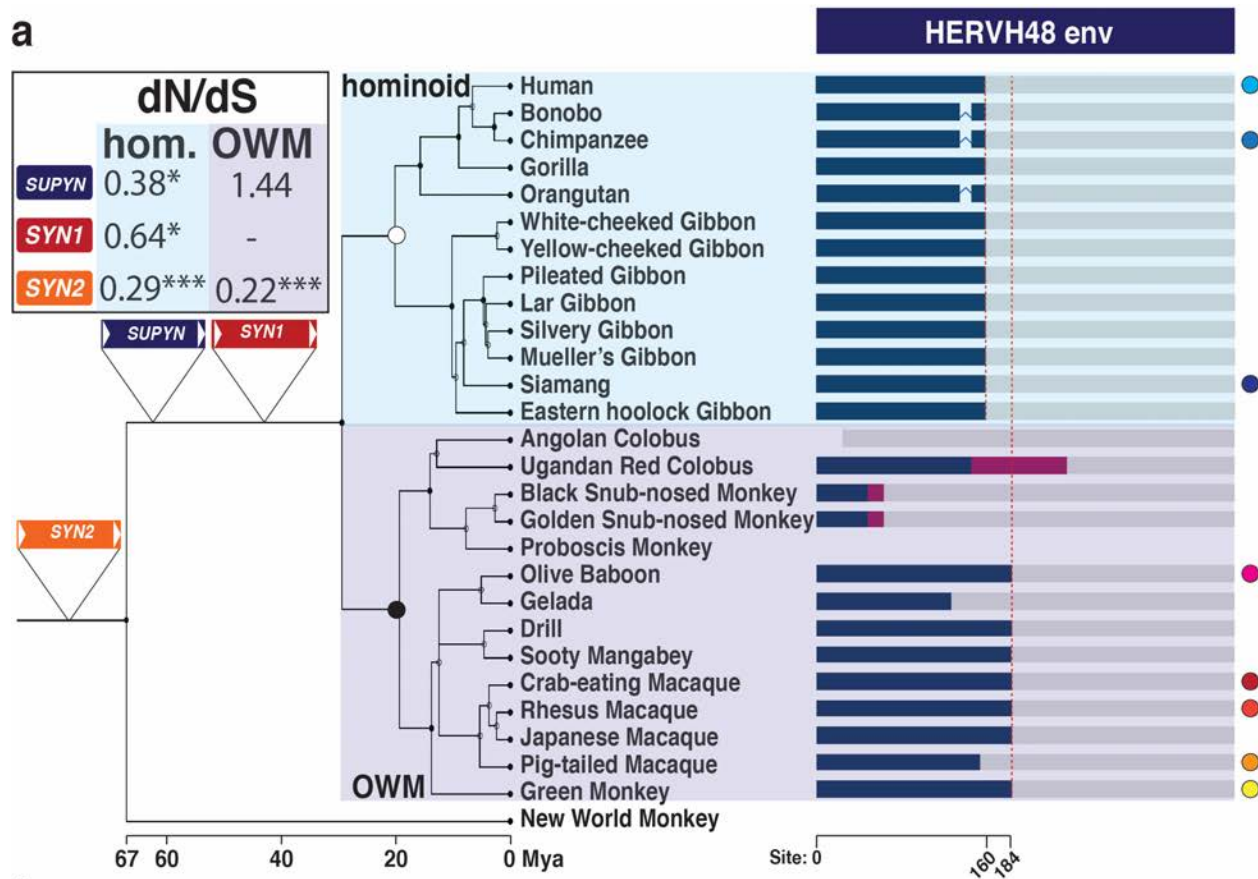
457 **(a)** Virus production and infection assay approach (see Methods). **(b, d)** Proportion of GFP+ 293T,

458 JEG3, Jar, and shRNA-transduced Jar cells infected with HIV-RD114env or HIV-VSVg. **(c)** Relative

459 infection rate of 293T and H1-ESCs normalized to mean proportion of HIV-VSVg-infected cells. **(e,**

460 **g, h)** Relative infection rates of GFP+ 293T cells transfected with **(e)** wild-type (WT-SP), rescue

461 (Resc-SP), luciferase signal peptide (GPluc-SP), unmodified **(e, g)** SUPYN, **(g)** RD114env, or **(h)**  
462 SMRVenv overexpression constructs. Relative infection was determined by normalizing indicated  
463 constructs to empty vector ( $n \geq 3$  with  $\geq 2$  technical replicates; \*\*\*adj.  $p < 0.001$ ; \*\*adj.  $p < 0.01$ ;  
464 \*adj.  $p < 0.05$ ; ANOVA with Tukey HSD). **(f)** Western Blot analysis ( $\alpha$ HA,  $\alpha$ GAPDH) of 293T cell  
465 lysates transfected with indicated constructs. **(i)** Model of SUPYN-dependent RDR infection  
466 restriction.  
467





469 **Figure 3: SUPYN emerged in a Catarrhine ancestor and has conserved antiviral activity in**  
470 **Hominoids.**

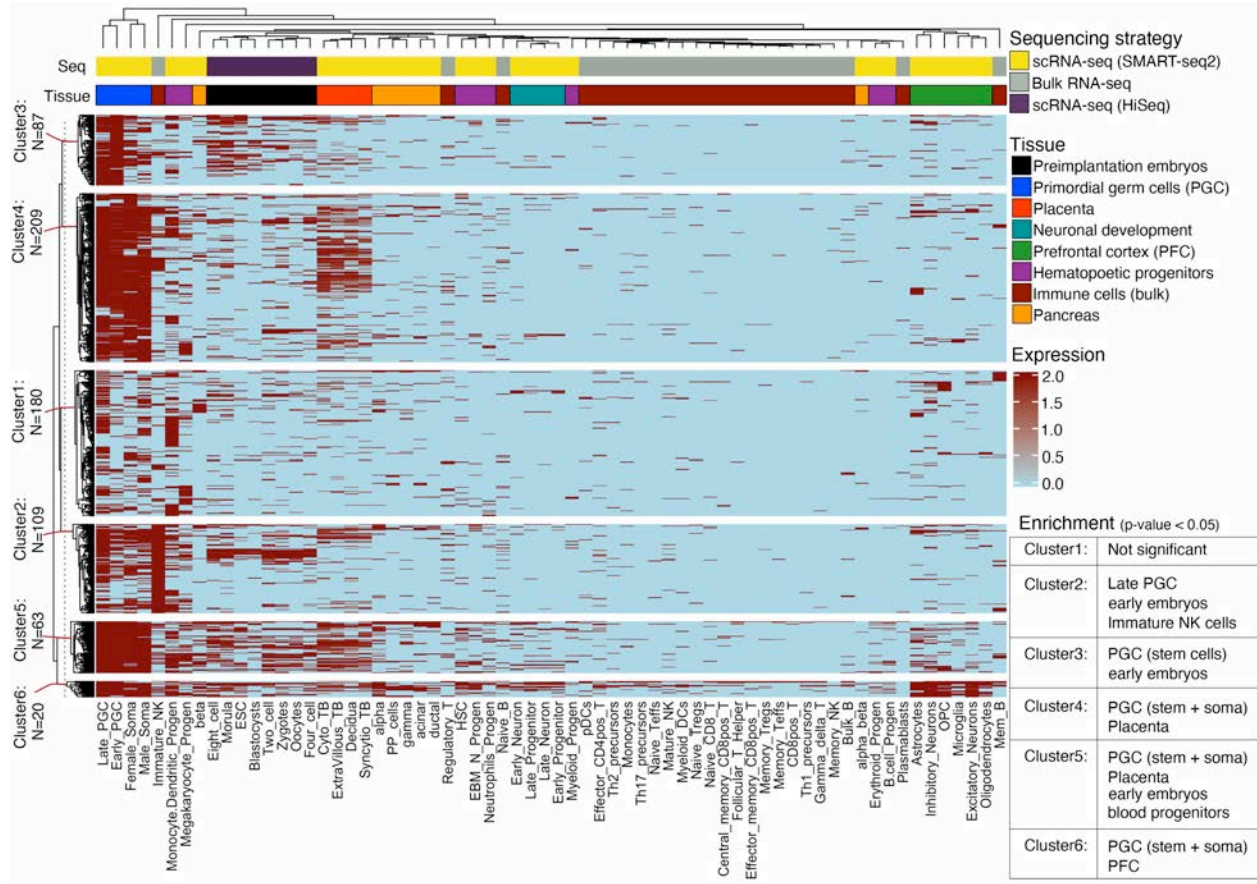
471 **(a)** Consensus primate phylogeny with cartoon representation of *SUPYN* ORFs (blue box).  
472 Magenta boxes represent frame-shifts in *SUPYN* ORFs. Red dashed lines denote conserved  
473 premature stop codon positions. Grey bars represent degraded HERVH48env sequence. Labeled  
474 triangles denote ancestral lineage where HERVH48env was acquired. Colored circles indicate  
475 species used in **b** and **c**. *SUPYN*, *SYN1* and *SYN2* dN/dS values are shown in box (\*\* $p < 0.001$ ; \* $p$   
476  $< 0.05$ ; LRT). **(b, c)** Relative infection rates and Western Blot of 293T cells transfected with primate  
477 **(b)** or ancestral **(c)** *SUPYN*-HA constructs and infected with HIV-RD114env are shown. Relative  
478 infection rates were determined by normalizing GFP+ counts to empty vector. ( $n \geq 3$  with  $\geq 2$   
479 technical replicates \*\*\*adj.  $p < 0.001$ ; \*adj.  $p < 0.05$ ; ANOVA with Tukey HSD)

480

481

482

483



484

485 **Figure 4: Expression profile of env-derived transcripts over a subset of human cell types.**

486 Heatmap shows expression of 668 *envORF* loci in 66 distinct cell-types from 8 independent

487 datasets shown in Extended Data 11. Rows represent individual *envORF* loci expressed ( $\log_2$  CPM

488 > 1) in at least one cell type (columns). Upper bars, located above the heatmap, denote the

489 sequencing strategy (top) and tissue source (bottom) with the same color scheme shown in

490 Extended Data 11. Rows are clustered into six distinct groups based on their expression.

491 Significant envORF enrichment in cell-types was calculated using a hypergeometric test.

492

493

494

495 **Methods**

496

497 **RNA-seq analyses**

498 We mined published single cell transcriptome datasets (scRNA-seq) of human pre-  
499 implantation embryos isolated at developmental stages ranging from oocyte to  
500 blastocyst<sup>47</sup> (GSE36552) and from human placenta<sup>48,49</sup> (GSE89497, GSE87726). Reads  
501 were mapped to the human genome (hg19) with STAR<sup>50</sup> using the following settings --  
502 *alignIntronMin 20 --alignIntronMax 1000000 --chimSegmentMin 15 --*  
503 *chimJunctionOverhangMin 15 --outFilterMultimapNmax 20*. Only uniquely mapped reads  
504 were considered for expression calculations. Gene level counts were obtained using  
505 *featureCounts*<sup>51</sup> run with RefSeq annotations. Gene expression levels were calculated at  
506 Transcript Per Million (TPM) from counts mapped over the entire gene (defined as any  
507 transcript located between Transcription Start Site (TSS) and Transcription End Site  
508 (TES)). Only genes and cells that met the following criteria were included in this analysis:  
509 (1) each cell must express at least 5,000 genes; (2) each gene must be expressed in at  
510 least 1% of cells; (3) each gene must be expressed with log<sub>2</sub> TPM >1. We clustered cells  
511 meeting these criteria using the default parameters of the Seurat<sup>52</sup> package (v3.1.1)  
512 implemented in R (v3.6.0). Seurat applies the most variable genes to get top principal  
513 components that are used to discriminate cell clusters in tSNE or UMAP plots. We set  
514 Seurat to use 10 principal components in this cluster analysis. For the placental  
515 scRNAseq data (**Fig1, Extended Data 2**), the 2000 most differentially expressed genes  
516 were used to define cell clusters. Major clusters corresponding to CTB, STB, EVTB,  
517 macrophages, and stromal cells were identified based on the expression of known marker

518 genes. Monocle2<sup>53</sup> was used to perform single-cell trajectory analysis and cell ordering  
519 along an artificial temporal continuum. The top 500 differentially expressed genes were  
520 used to distinguish between CTB, STB and EVTB cell populations. The transcriptome  
521 from each single cell represents a pseudo-time point along an artificial time vector that  
522 denotes the progression of CTB to STB or EVTB respectively.

523 Data generated on the 10X Genomics scRNA-seq platforms were processed in the  
524 following way. The processed data matrix from Vento-Tormo et al.<sup>54</sup> was retrieved from  
525 the E-MTAB-6701 entry. The normalized counts and cell-type annotations were used as  
526 provided by the original publications. Seurat was used for filtering, normalization and cell-  
527 type identification. The following data processing steps were performed: (1) Cells were  
528 filtered based on the criteria that individual cells must have between 1,000 and 5,000  
529 expressed genes with a count  $\geq 1$ ; (2) cells with more than 5% of counts mapping to  
530 mitochondrial genes were filtered out; (3) data was normalized by dividing uniquely  
531 mapping read counts (defined by Seurat as unique molecular identified (UMI)) for each  
532 gene by the total number of counts in each cell and multiplying by 10,000. These  
533 normalized values were then natural-log transformed. Cell types were defined by using  
534 the top 2000 variable features expressed across all samples. Clustering was performed  
535 using the “FindClusters” function with largely default parameters; except resolution was  
536 set to 0.1 and the first 20 PCA dimensions were used in the construction of the shared-  
537 nearest neighbor (SNN) graph and the generation of 2-dimensional embeddings for data  
538 visualization using UMAP. Cell types were assigned based on the annotations provided  
539 by the original publication.

540 Bulk RNAseq datasets generated from placenta<sup>55</sup>, 293T<sup>56,57</sup> and human immune cells<sup>58</sup>,  
541 were processed as described above. Briefly, reads were mapped with STAR and uniquely  
542 mapped reads were counted with featureCounts.

543

#### 544 **ChIP-seq, DNase-seq and ATAC-seq data analysis**

545 Various ChIP-seq datasets representing histone modifications and transcription factors in  
546 human embryonic stem cells and their differentiation were retrieved from<sup>59,60</sup> (GSE61475,  
547 GSE99631). We obtained the H3K27Ac<sup>61</sup> (GSE127288) for CTB to STB primary cultures,  
548 H3K4Me1 for trophoblasts<sup>62</sup> (GSE118289), H3K4Me3, H3K27Me3 for differentiated  
549 trophoblasts<sup>63</sup> (GSE105258), and GATA2/3, TFAP2A/C<sup>63</sup> (GSE105081) ChIP-seq  
550 datasets in raw fastq format. DNase-seq and ATAC-seq datasets were retrieved from  
551 Gao et al.<sup>64</sup> (GSA:CRA000297) and Wu et al.<sup>65</sup> (GSE101571) respectively.

552 Reads from the above described datasets were aligned to the hg19 human reference  
553 genome using Bowtie2<sup>66</sup> run in *--very-sensitive-local* mode. All reads with MAPQ < 10  
554 and PCR duplicates were removed using Picard and *samtools*<sup>67</sup>. All of the peaks were  
555 called by MACS2<sup>68</sup> (<https://github.com/macs3-project/MACS>) with the parameters in  
556 narrow mode for TFs and broad mode for histone modifications keeping FDR < 1%.  
557 ENCODE-defined blacklisted regions<sup>69</sup> were excluded from called peaks. We then  
558 intersected these peak sets with repeat elements from hg19 repeat-masked coordinates  
559 using bedtools *intersectBed*<sup>70</sup> with a 50% overlap. To visualize over Refseq genes (hg19)  
560 using IGV<sup>71</sup>, the raw signals of ChIP-seq were obtained from MACS2, using the  
561 parameters: *-g hs -q 0.01 -B*. The conservation track was visualized through UCSC

562 genome browser<sup>14</sup> under net/chain alignment of given non-human primates (NHPs) and  
563 merged beneath the IGV tracks.

564

#### 565 **Cell culture**

566 293T cells were cultured in DMEM (GIBCO, 11995065) containing 10% Fetal Bovine  
567 Serum (FBS) (GIBCO, 10438026). Jar cells (provided by Carolyn Coyne) were cultured  
568 in RPMI (GIBCO, 11875093) containing 10% FBS. JEG3 cells were cultured in MEM  
569 (GIBCO, 11095080) containing 10% FBS. Culture medium for these cell lines was  
570 supplemented with sodium pyruvate (GIBCO, 11360070), glutamax (GIBCO, 35050061),  
571 and Penicillin Streptomycin (GIBCO, 15140122) according to manufacturer  
572 specifications. H1-ESCs (obtained from WiCell) were grown on Matrigel (Corning,  
573 356277) coated plates in MTESR+ (Stemcell, 05825) growth-media and sub-cultured  
574 using Accutase (Innovative Cell Technologies, AT-104) and MTESR+ supplemented with  
575 CloneR (Stemcell, 05888). All cell lines were cultured at 37°C and 5% CO<sub>2</sub>.

576

#### 577 **Vector cloning**

578 DHIV3-GFP, phCMV-RD114env, psi(-)-amphoMLV plasmids were provided by Vicente  
579 Planelles (University of Utah). pCGCG-SMRVenv plasmid was provided by Welkin  
580 Johnson (Boston University). psPAX2 and pVSVg plasmids were provided by John Lis  
581 (Cornell University). The following cloning approaches were performed using primers and  
582 constructs described in Supplementary Table 3. HERVH1env ORF was PCR-amplified  
583 using Q5 polymerase (NEB, M0491L) from HeLa and 293T genomic DNA respectively  
584 and cloned into a TOPO vector (ThermoFisher, 450245).

585 To generate stable *SYN1* and *SUPYN* knock-down cell lines, pHIV lentiviral constructs  
586 containing shRNAs targeting *SYN1* and *SUPYN* respectively were cloned using the  
587 following strategy. The shRNA encoded in pHIV7-U6-shW3, generously provided by Lars  
588 Aagaard (Aarhus University), targets *SYN1*<sup>72</sup>. *SUPYN*-targeting shRNAs were designed  
589 using siRNA sequences employed by Jun Sugimoto<sup>4</sup> as a template. pHIV7 lentiviral  
590 constructs were cloned using the pHIV7-U6-shW3 plasmid<sup>72</sup> as a template. pHIV7-U6-  
591 shSup-cer, pHIV7-U6-shSup-puro, pHIV7-U6-shC-cer, pHIV7-U6-shC-puro, pHIV7-U6-  
592 shSyn1-cer, pHIV7-U6-shSyn1-puro were generated using a Gibson assembly approach.  
593 To replace the native GFP marker of pHIV7-U6-shW3 with a Cerulean reporter or  
594 puromycin resistance marker, we digested pHIV7-U6-shW3 with NheI (NEB, R3131S)  
595 and KpnI (NEB, R3142S). This digest resulted in the production of three DNA fragments:  
596 pHIV7 backbone, GFP-, and WPRE-containing fragments. We separately PCR amplified  
597 each selection marker and WPRE containing pHIV7 fragment. InFusion cloning was then  
598 used to ligate the digested pHIV7 backbone to the Cerulean or puromycin cassette and  
599 WPRE containing PCR product. shRNAs were cloned into the pHIV7-  
600 Cerulean/puromycin transfer construct previously digested with NotI (NEB, R0189S) and  
601 NheI. U6-promoter containing shRNA cassettes and the CMV promoter driving marker  
602 cassette expression were PCR amplified and subsequently InFusion cloned into the  
603 NotI/NheI digested pHIV7-cerulean/puromycin backbone.

604 All pHCMVenv and *SUPYN* expression constructs, described in this study, were  
605 generated as follows: HA-tagged and untagged ORFs with pHCMV homologous  
606 overhanging sequence were either PCR amplified using Q5 polymerase (NEB, M0491S)  
607 or synthesized (IDT) (see **Supplementary Table 2**), and cloned into EcoRI (NEB,

608 R3101T) digested pHCMV backbones using the InFusion cloning kit (Takara Bio,  
609 638920). To generate siRNA-resistant SUPYN rescue constructs, we replaced the native  
610 signal peptide sequence<sup>4</sup> (which is targeted by siRNAs used in this study) with (1) a  
611 *Gaussia princeps* luciferase SP (SUPYN-lucSP)<sup>73,74</sup> and (2) a shSUPYN resistant  
612 SUPYN rescue construct (SUPYN-rescSP) in which the codons were modified to retain  
613 the codon identity but disrupt siRNA binding.

614

### 615 **Antibodies**

616 All antibodies used in this study are commercially available.  $\alpha$ -GAPDH (D4C6R,  
617 D16H11),  $\alpha$ - $\beta$ actin (D6A8),  $\alpha$ -HA (C29F4),  $\alpha$ -ASCT2 (V501) primary antibodies were  
618 purchased from Cell Signaling Technology.  $\alpha$ -Mouse (#7076) and  $\alpha$ -Rabbit (#7074) HRP  
619 conjugated secondary antibodies were purchased from Cell Signaling Technology. IRDye  
620 secondary antibodies were purchased from Licor (925-32211, 925-68072, 925-32210,  
621 925-68073).  $\alpha$ -SUPYN primary antibody was purchased from Phoenix Pharma (H-059-  
622 052). Alexa-fluor conjugated secondary antibody was purchased from Invitrogen.

623

### 624 **Western Blot**

625 Whole cell extracts from cultured cell lines were prepared using 1x GLO lysis buffer  
626 (Promega, E266A). One third volume of 4x Laemli buffer was added to one volume whole  
627 cell extract samples, then incubated at 95°C for 5 minutes, and sonicated for 15 minutes  
628 at 4°C (amplitude 100; pulse interval 15 seconds on, 15 seconds off). Approximately 30  
629  $\mu$ g of protein were separated by SDS-PAGE (BioRad, 1610175), transferred to PVDF  
630 membrane (BioRad, 1620177), blocked according to antibody manufacturers



631 specification, and incubated overnight in appropriate primary antibody then incubated in  
632 IRDye or peroxidase conjugated goat anti-mouse or anti-rabbit secondary antibodies for  
633 1 hour at room temperature. Protein was then detected using ECL reagent (BioRad,  
634 1705061) or the Licor Odyssey imaging system.

635

### 636 **IF microscopy**

637 Human second trimester placental tissue that resulted from elective terminations was obtained  
638 from the University of Pittsburgh Health Sciences Tissue Bank through an honest broker system  
639 after approval from the University of Pittsburgh Institutional Review Board (IRB) and in  
640 accordance with the University of Pittsburgh's tissue procurement guidelines. Tissue was  
641 excluded in cases of fetal anomalies or aneuploidy. Third trimester placental tissue was obtained  
642 through the Magee Obstetric Maternal & Infant (MOMI) Database and Biobank after approval from  
643 the University of Pittsburgh IRB. Women who had previously consented for tissue donation and  
644 underwent cesarean delivery were included. Placental tissues were fixed in 4% PFA (in 1x  
645 PBS) for 30 minutes, permeabilized with 0.25% Triton X-100 for 30 minutes (on a rocker),  
646 washed with 1x PBS and then incubated with primary anti-Suppressyn antibody at 1:200  
647 in 1x PBS for 2-4 hours at room temperature. These samples were incubated with Alexa-  
648 fluor conjugated secondary antibody (Invitrogen) diluted 1:1000 and counterstained with  
649 actin. DAPI was included in our PBS and then mounted in Vectashield mounting medium  
650 with DAPI (Vector Laboratories, H-1200).

651

### 652 **Virus production**

653 Low passage 293T cells were used to produce lentiviral particles. DHIV3-GFP and env-  
654 expression plasmids were co-transfected at a mass ratio of 2:1 using lipofectamine 2000

655 (ThermoFisher, 11668030). shRNA encoding lentiviral particles were produced by co-  
656 transfecting pHIV7, psPAX2, pVSVg according to BROAD institute lentiviral production  
657 protocol (<https://portals.broadinstitute.org/gpp/public/resources/protocols>) using  
658 Lipofectamine 2000. Growth media was replaced on transfected cells after overnight  
659 incubation. At 72 hours post-transfection, virus containing supernatant was harvested,  
660 centrifuged to remove cell debris, filtered through a 0.45 um pore filter, and stored at -  
661 80°C.

662

### 663 **Infection Assays**

664 293T cells were transfected with env-overexpression constructs using Lipofectamine  
665 2000 and incubated 24 hours. Transfected cells were infected with reporter virus by  
666 applying virus (HIV-RD114env, HIV-VSVg, HIV-SMRVenv) stocks in the presence of  
667 polybrene (Santa Cruz Bio, sc-134220) at a final concentration of 4 ug/mL. After 6-8  
668 hours, virus stock was replaced with fresh growth media. Infected cells were maintained  
669 for 72 hours, replacing media when necessary, and harvested with trypsin. Detached cells  
670 were suspended in fresh growth media, strained and analyzed by flow cytometry. For the  
671 H1-ESC infection experiment, relative infection rates were calculated by normalizing the  
672 percent GFP+, HIV-RD114env infected cells to the percent GFP+ HIV-VSVg infected  
673 cells. For env/SUPYN overexpression experiments, relative infection rates were  
674 calculated by normalizing the percent GFP+ env/SUPYN-transfected cells to the percent  
675 GFP+ empty vector transfected cells. ANOVA with Tukey HSD tests were implemented  
676 in R (v3.6.3).

677

## 678 **Placental cell shRNA transduction**

679 Placenta-derived cell lines were treated with pHIV-shRNA-virus-containing supernatant  
680 and incubated for 72 hours as described in Infection Assays. Cerulean positive cells were  
681 sorted using the BD FACS Aria cytometer. Cells transduced with puroR cassette were  
682 treated with Puromycin (GIBCO, A1113802) at a final concentration of 3.5 ug/mL for 7  
683 days, then cultured in regular growth media.

684

## 685 **RT-qPCR**

686 RNA was isolated from cultured cells using the RNeasy Mini Kit (Qiagen, 74104) and an  
687 on column dsDNAse digestion (Qiagen, 79254) was performed. 1-3 ug of total RNA were  
688 used to generate cDNA with the maxima cDNA synthesis with dsDNAse kit  
689 (ThermoFisher, K1681). qPCR reactions were performed using the LC480 Instrument  
690 with Sybr Green PCR master mix (Roche, 04707516001) according to manufacturer's  
691 protocol and using primers indicated in Supplementary Table 2. Gene expression was  
692 then quantified using the  $\Delta\Delta$ CT method<sup>75</sup>. 18S expression was used as a reference  
693 housekeeping gene. Wilcox rank sum tests were performed using R (v3.6.3).

694

## 695 **Envelope evolutionary sequence analyses**

696 Orthologous *SUPYN*, *SYN1*, and *SYN2* sequences were extracted from the 30-species  
697 MULTIZ alignment<sup>14</sup> and formatted for sequence alignment using the phast package<sup>76</sup>.  
698 These and additional syntenic *SUPYN* and *SYN2* open reading frame sequences were  
699 validated/identified by BLASTn<sup>77</sup> search with default settings of publicly available  
700 Catarrhine primate genomes (ncbi.nih.gov). Mariam Okhvat of the Carbone Lab (Oregon

701 Health and Science University) generously provided BAM files containing read alignment  
702 information for *SUPYN*, *SYN1*, and *SYN2* generated from whole genome sequencing of  
703 *Hoolock leuconedys* (Hoolock Gibbon), *Symphalangus syndactylus* (Siamang),  
704 *Hylobates muelleri* (Müller's Gibbon), *Hylobates lar* (Lar Gibbon), *Hylobates moloch*  
705 (Silvery Gibbon), *Hylobates pileatus* (Pileated Gibbon), and *Nomascus gabriellae*  
706 (Yellow-cheeked Gibbon). Where multiple individuals were sequenced, a consensus  
707 sequence was generated using samtools<sup>67</sup> and JalView<sup>78</sup>.

708 To perform *dN/dS* analyses, orthologous env sequences (>90bp length) encoding the  
709 mature sequence downstream of the signal peptide cleavage site, were aligned using  
710 MEGA7<sup>79</sup> and manually converted to PHYLIP format. A Newick tree was generated based  
711 on this alignment using the maximum likelihood algorithm implemented in MEGA7.  
712 *Codeml*, implemented in the PAML package, was run to calculate *dN/dS* values and log  
713 likelihood (LnL) scores generated under models M0, M1, M2, M7 and M8<sup>38</sup>. Chi-square  
714 tests comparing LnL scores generated under models of neutral evolution and selection  
715 were performed.

716 We used two approaches to reconstruct ancestral hominoid and OWM *SUPYN*  
717 sequences. First, we reconstructed ancestral *SUPYN* sequences using the majority rule  
718 consensus sequence (calculated in JalView) of the hominoid and OWM clade  
719 respectively. At positions where nucleotide identity was ambiguous, the dominant  
720 nucleotide identity in the neighboring clade was used as a tiebreaker. These sequences  
721 were used for our infection assays shown in **Fig 3c**. We also employed a maximum  
722 likelihood approach using the *baseml* program, implemented in PAML<sup>38</sup>. We  
723 reconstructed ancestral *SUPYN* sequences using the hominoid species, shown in **Fig 3a**,

724 and the 6 OWM monkeys with the most complete *SUPYN*-coding open reading frame  
725 (olive baboon, drill, crab-eating macaque, rhesus macaque, japanese macaque, green  
726 monkey) as our input sequences. Because PAML requires a Newick tree as an input, the  
727 MEGAX<sup>80-82</sup> maximum likelihood algorithm was used to generate a Newick tree with the  
728 above described *SUPYN* sequences. *Baseml* was run using models 3-7 (F84, HKY85,  
729 T92, TN93, REV). As shown in **Extended Data 10**, both the consensus- and maximum  
730 likelihood reconstructions were identical for the OWM *SUPYN* sequences. The  
731 consensus-based hominoid sequence reconstruction differed from our maximum  
732 likelihood-based reconstruction by two amino acids. These two positions are unlikely to  
733 affect the function of the resulting protein because these sites are identical to siamang  
734 *SUPYN*, which restricts RD114env-mediated infection (**Fig 3b, c**).

735

### 736 **Genome-wide search for endogenous retrovirus derived envelope open reading** 737 **frames**

738 Candidate envelope open reading frames (envORF) were identified by performing  
739 tBLASTn<sup>83</sup> searches of the hg19 human genome assembly using envelope amino acid  
740 sequences, obtained from Rebase<sup>84</sup> and published retroviral envelope sequences, as a  
741 query. Collected hits were used as a query to repeat a tBLASTn search, initially yielding  
742 82715 candidate envORF sequences. This list of candidates was filtered using the  
743 following criteria. (1) EnvORFs must have a length of  $\geq 70$  aa. (2) Hits starting at position  
744  $\geq 300$  aa were removed because such open reading frames are predicted to encode a  
745 portion of the envelope transmembrane domain, which is not expected to play a role in  
746 receptor binding. (3) After these processing steps, our list was further concatenated to

747 only include unique open reading frame sequences (n=1507) in the hg19 genome  
748 assembly.

749 To estimate envORF expression, we mined publicly available scRNA-seq datasets  
750 generated on the SMART-seq2 platform because this technology yields higher genomic  
751 coverage compared with methods employing poly(A) enrichment. We gathered scRNA-  
752 seq datasets in fastq format from preimplantation embryos<sup>47</sup> (GSE36552), placenta<sup>48</sup>  
753 (GSE89497), primordial germ cells<sup>85</sup> (GSE63818), hematopoiesis<sup>86</sup> (GSE75478),  
754 neuronal differentiation<sup>87</sup> (GSE93593), prefrontal cortex<sup>88</sup> (GSE67835), and pancreas<sup>89</sup>  
755 (GSE81547) (Supplemental Table 1). We also included bulk RNA-seq of immune cells  
756 from six individuals comprising 25 subtypes<sup>58</sup> (GSE118165) (Supplemental Table 1).

757 We curated our envORF sequences (in fasta format) and hg19 refseq gene models<sup>14</sup> (gtf  
758 format) to calculate envORF expression with high confidence. First, we masked envORF  
759 DNA sequences from the hg19 genome. These sequences were added, as individual  
760 fasta contigs, to the hg19 reference assembly. We then modified the transcriptome model  
761 by appending the envORFs models (gtf format) to the hg19 refseq gene models. Next,  
762 we indexed the modified reference sequences using STAR default parameters, by  
763 providing the new transcriptome models that include the envORFs. This approach  
764 enabled us to simultaneously calculate envORF and protein-coding gene expression.  
765 After aligning and retaining uniquely mappable sequencing reads, we calculated envORF  
766 and gene expression (see Single cell RNA-seq analysis procedures above). We  
767 processed the individual datasets independently, then constructed the expression matrix  
768 using *Seurat* and checked if the obtained clusters were consistent with the original  
769 studies. Despite including the envORFs, we identified each cell type in similar proportions,

770 as shown in their corresponding studies (Extended Data 11). Then we merged all of the  
771 datasets into a single pool for our downstream analysis while maintaining cell type and  
772 tissue annotations. We scaled and normalized these datasets using Seurat within the R  
773 environment and checked the normalization status of merged datasets with a UMAP  
774 biplot. This was done by using the most variable genes to ensure that cells did not cluster  
775 based on platform or batch differences (Extended Data 11 and Supplemental Table 1).  
776 We calculated gene and envORF expression (TPM) using uniquely mapped reads. Any  
777 envORF with TPM < 1 in > 99% of cells across the data frame was discarded from further  
778 analysis. This strategy identified 668 unique envORFs with evidence for expression.  
779 Finally, expression of each envORF across all cell types was used to cluster envORF  
780 according to their expression profile, which were visualized in Fig. 4.

781 To survey the expression of known envelope genes, we utilized Gtex<sup>90</sup> bulk RNA-seq  
782 datasets (phs000424.v6.p1), which were generated from distinct post-mortem tissues.

783

#### 784 **Methods References**

- 785 47. Yan, L. *et al.* Single-cell RNA-Seq profiling of human preimplantation embryos  
786 and embryonic stem cells. *Nat. Struct. Mol. Biol.* **20**, 1131–1139 (2013).
- 787 48. Liu, Y. *et al.* Single-cell RNA-seq reveals the diversity of trophoblast subtypes and  
788 patterns of differentiation in the human placenta. *Cell Res.* **28**, 819–832 (2018).
- 789 49. Pavličev, M. *et al.* Single-cell transcriptomics of the human placenta: inferring the  
790 cell communication network of the maternal-fetal interface. *Genome Res.* **27**,  
791 349–361 (2017).

- 792 50. Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15–  
793 21 (2013).
- 794 51. Liao, Y., Smyth, G. K. & Shi, W. featureCounts: an efficient general purpose  
795 program for assigning sequence reads to genomic features. *Bioinformatics* **30**,  
796 923–930 (2014).
- 797 52. Stuart, T. *et al.* Comprehensive Integration of Single-Cell Data. *Cell* **177**, 1888–  
798 1902.e21 (2019).
- 799 53. Qiu, X. *et al.* Single-cell mRNA quantification and differential analysis with  
800 Census. *Nat. Methods* **14**, 309–315 (2017).
- 801 54. Vento-Tormo, R. *et al.* Single-cell reconstruction of the early maternal-fetal  
802 interface in humans. *Nature* **563**, 347–353 (2018).
- 803 55. Zadora, J. *et al.* Disturbed Placental Imprinting in Preeclampsia Leads to Altered  
804 Expression of DLX5, a Human-Specific Early Trophoblast Marker. *Circulation*  
805 **136**, 1824–1839 (2017).
- 806 56. D Antonio, M. *et al.* Identifying DNase I hypersensitive sites as driver distal  
807 regulatory elements in breast cancer. *Nat Commun* **8**, 436 (2017).
- 808 57. Chung, H. *et al.* Human ADAR1 Prevents Endogenous RNA from Triggering  
809 Translational Shutdown. *Cell* **172**, 811–824.e14 (2018).
- 810 58. Calderon, D. *et al.* Landscape of stimulation-responsive chromatin across diverse  
811 human immune cells. *Nature Publishing Group* **51**, 1494–1505 (2019).
- 812 59. Tsankov, A. M. *et al.* Transcription factor binding dynamics during human ES cell  
813 differentiation. *Nature* **518**, 344–349 (2015).



- 814 60. Barakat, T. S. *et al.* Functional Dissection of the Enhancer Repertoire in Human  
815 Embryonic Stem Cells. *Cell Stem Cell* **23**, 276–288.e8 (2018).
- 816 61. Kwak, Y.-T., Muralimanoharan, S., Gogate, A. A. & Mendelson, C. R. Human  
817 Trophoblast Differentiation Is Associated With Profound Gene Regulatory and  
818 Epigenetic Changes. *Endocrinology* **160**, 2189–2203 (2019).
- 819 62. Dunn-Fletcher, C. E. *et al.* Anthropoid primate-specific retroviral element THE1B  
820 controls expression of CRH in placenta and alters gestation length. *PLoS Biol.* **16**,  
821 e2006337 (2018).
- 822 63. Krendl, C. *et al.* GATA2/3-TFAP2A/C transcription factor network couples human  
823 pluripotent stem cell differentiation to trophoctoderm with repression of  
824 pluripotency. *Proc. Natl. Acad. Sci. U.S.A.* **114**, E9579–E9588 (2017).
- 825 64. Gao, L. *et al.* Chromatin Accessibility Landscape in Human Early Embryos and Its  
826 Association with Evolution. *Cell* **173**, 248–259.e15 (2018).
- 827 65. Wu, J. *et al.* Chromatin analysis in human early development reveals epigenetic  
828 transition during ZGA. *Nature* **557**, 256–260 (2018).
- 829 66. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat.*  
830 *Methods* **9**, 357–359 (2012).
- 831 67. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics*  
832 **25**, 2078–2079 (2009).
- 833 68. Gasper, J.M. Improved peak-calling with MACS2. *bioRxiv.org* doi:  
834 <https://doi.org/10.1101/496521> (2018)
- 835 69. Amemiya, H. M., Kundaje, A. & Boyle, A. P. The ENCODE Blacklist: Identification  
836 of Problematic Regions of the Genome. *Sci Rep* **9**, 9354 (2019).

- 837 70. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing  
838 genomic features. *Bioinformatics* **26**, 841–842 (2010).
- 839 71. Robinson, J. T. *et al.* Integrative genomics viewer. *Nat. Biotechnology* **29**, 24-26  
840 (2011).
- 841 72. Aagaard, L. *et al.* Silencing of endogenous envelope genes in human  
842 choriocarcinoma cells shows that envPb1 is involved in heterotypic cell fusions.  
843 *J. Gen. Virol.* **93**, 1696–1699 (2012).
- 844 73. Luft, C. *et al.* Application of Gaussia luciferase in bicistronic and non-conventional  
845 secretion reporter constructs. *BMC Biochem.* **15**, 14 (2014).
- 846 74. Knappskog, S. *et al.* The level of synthesis and secretion of Gaussia princeps  
847 luciferase in transfected CHO cells is heavily dependent on the choice of signal  
848 peptide. *J. Biotechnol.* **128**, 705–715 (2007).
- 849 75. Schmittgen, T. D. Livak, J. L. Analyzing real-time PCR data by the comparative  
850 CT method. *Nat. Protocols* **3**, 1101-1108 (2008).
- 851 76. Hubisz, M. J., Pollard, K. S. & Siepel, A. PHAST and RPHAST: phylogenetic  
852 analysis with space/time models. *Brief. Bioinformatics* **12**, 41–51 (2011).
- 853 77. Ye, J., McGinnis, S. & Madden, T. L. BLAST: improvements for better sequence  
854 analysis. *Nucleic Acids Res.* **34**, W6–9 (2006).
- 855 78. Waterhouse, A. M., Procter, J. B., Martin, D. M. A., Clamp, M. & Barton, G. J.  
856 Jalview Version 2--a multiple sequence alignment editor and analysis workbench.  
857 *Bioinformatics* **25**, 1189–1191 (2009).

- 858 79. Kumar, S., Stecher, G. & Tamura, K. MEGA7: Molecular Evolutionary Genetics  
859 Analysis Version 7.0 for Bigger Datasets. *Molecular Biology and Evolution* **33**,  
860 1870–1874 (2016).
- 861 80. Tamura, K., Nei, M. Estimation of the number of nucleotide substitutions in the  
862 control region of mitochondrial DNA in humans and chimpanzees. *Mol Biol Evol.*  
863 **10(3)**, 512-526 (1993).
- 864 81. Kumar, S., Stecher, G., Li, M., Knyaz, C., Tamura, K. MEGA X: molecular  
865 evolutionary genetics analysis across computing platforms. *Mol Biol Evol.* **35(6)**,  
866 1547-1549 (2018).
- 867 82. Stecher, G. *et al.* Molecular evolutionary genetics analysis (MEGA) for macOS.  
868 *Mol Biol Evol.* **37(4)**, 1237-1239 (2020).
- 869 83. Gertz, E. M., Yu, Y.-K., Agarwala, R., Schäffer, A. A. & Altschul, S. F.  
870 Composition-based statistics and translated nucleotide searches: improving the  
871 TBLASTN module of BLAST. *BMC Biol.* **4**, 41 (2006).
- 872 84. Bao, W., Kojima, K. K. & Kohany, O. Repbase Update, a database of repetitive  
873 elements in eukaryotic genomes. *Mob DNA* **6**, 11–6 (2015).
- 874 85. Guo, F. *et al.* The Transcriptome and DNA Methylome Landscapes of Human  
875 Primordial Germ Cells. *Cell* **161**, 1437–1452 (2015).
- 876 86. Velten, L. *et al.* Human haematopoietic stem cell lineage commitment is a  
877 continuous process. *Nat. Cell Biol.* **19**, 271–281 (2017).
- 878 87. Close, J. L. *et al.* Single-Cell Profiling of an In Vitro Model of Human Interneuron  
879 Development Reveals Temporal Dynamics of Cell Type Production and  
880 Maturation. *Neuron* **93**, 1035–1048.e5 (2017).

- 881 88. Darmanis, S. *et al.* A survey of human brain transcriptome diversity at the single  
882 cell level. *Proc. Natl. Acad. Sci. U.S.A.* **112**, 7285–7290 (2015).
- 883 89. Enge, M. *et al.* Single-Cell Analysis of Human Pancreas Reveals Transcriptional  
884 Signatures of Aging and Somatic Mutation Patterns. *Cell* **171**, 321–330.e14  
885 (2017).
- 886 90. GTEx Consortium. *et al.* Genetic effects on gene expression across human  
887 tissues. *Nature* **550**, 204–213 (2017).

888

### 889 **Acknowledgements**

890 This work was funded by Cornell University and by the National Institutes of Health to CF  
891 (R01 GM112972, R35 GM122550). We thank Dr. Vicente Planelles for sharing reporter  
892 virus plasmids and training; Dr. Welkin Johnson for sharing SMRVenv expression  
893 plasmid; Dr. Lars Aagaard for sharing shRNA transduction constructs; Dr. John Lis for  
894 providing lentiviral packaging construct; Dr. Amnon Koren and Rita Rebello for providing  
895 embryonic stem cell cultures and technical support respectively; Dr Lucia Carbone and  
896 Dr. Mariam Okhvat for sharing unpublished gibbon genome sequences. We thank Ray  
897 Malfavon-Borja for producing an initial list of envelope open reading frames in the human  
898 genome. We thank Maia Clare for her contribution to evolutionary sequence analyses.  
899 We thank members of the Feschotte lab for helpful advice and discussion throughout the  
900 project.

901

902

903

904 **Author contributions**

905 CF conceived of this project. JAF and CF designed and developed this project. JAF  
906 designed and conducted all experiments, evolutionary sequence analyses, and analyzed  
907 all experimental data. MS performed all gene expression and regulation analyses. HBC  
908 and RAK helped perform infection assays and evolutionary sequence analyses. CBC  
909 performed placental stains. JAF, CF and MS wrote this manuscript.

910

911 **Competing interest declaration**

912 The authors declare no competing interests.

913

914 **Additional Information**

915

916

917

918

919

920

921

922

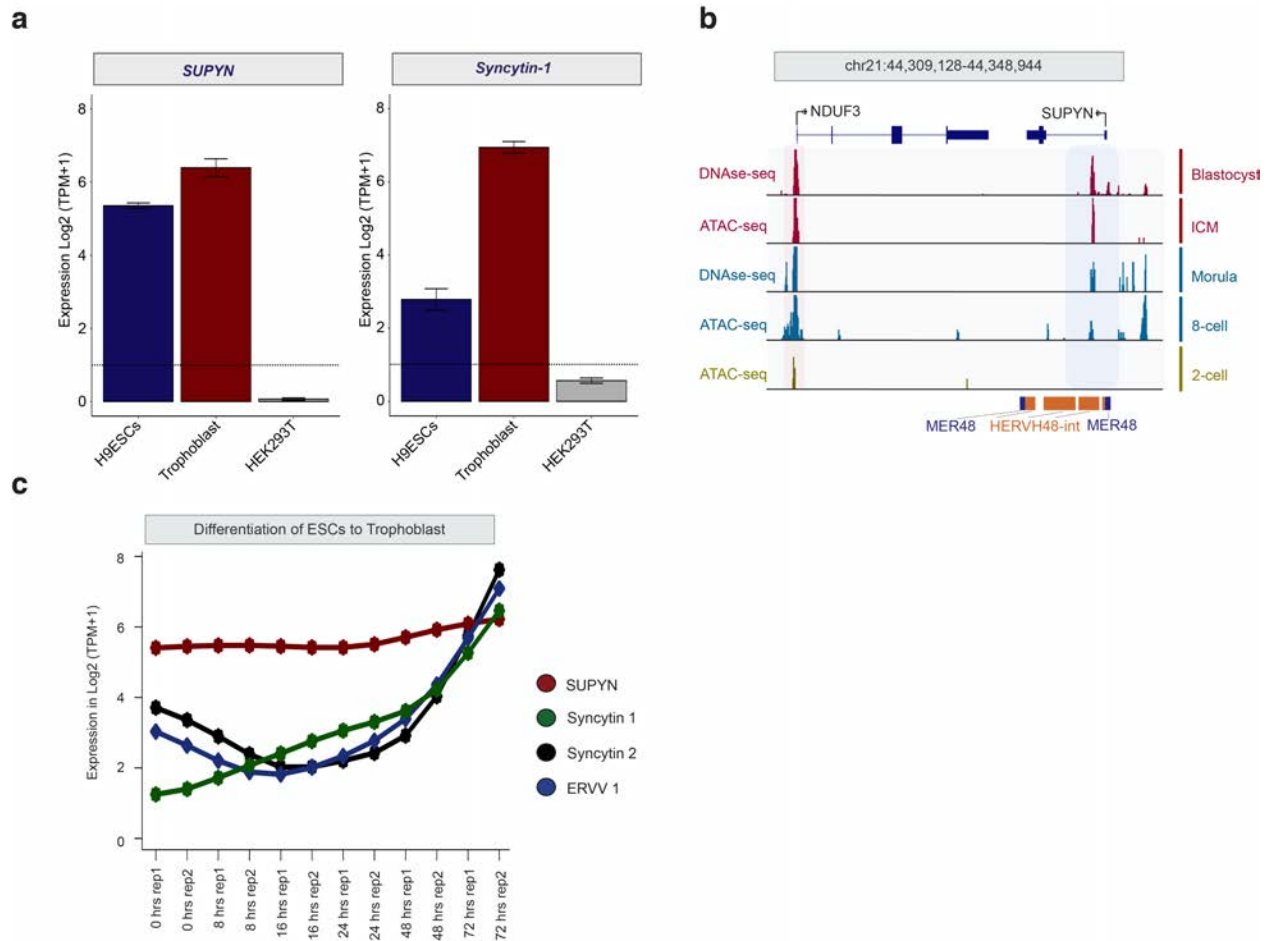
923

924

925

926

927 **Extended Data**

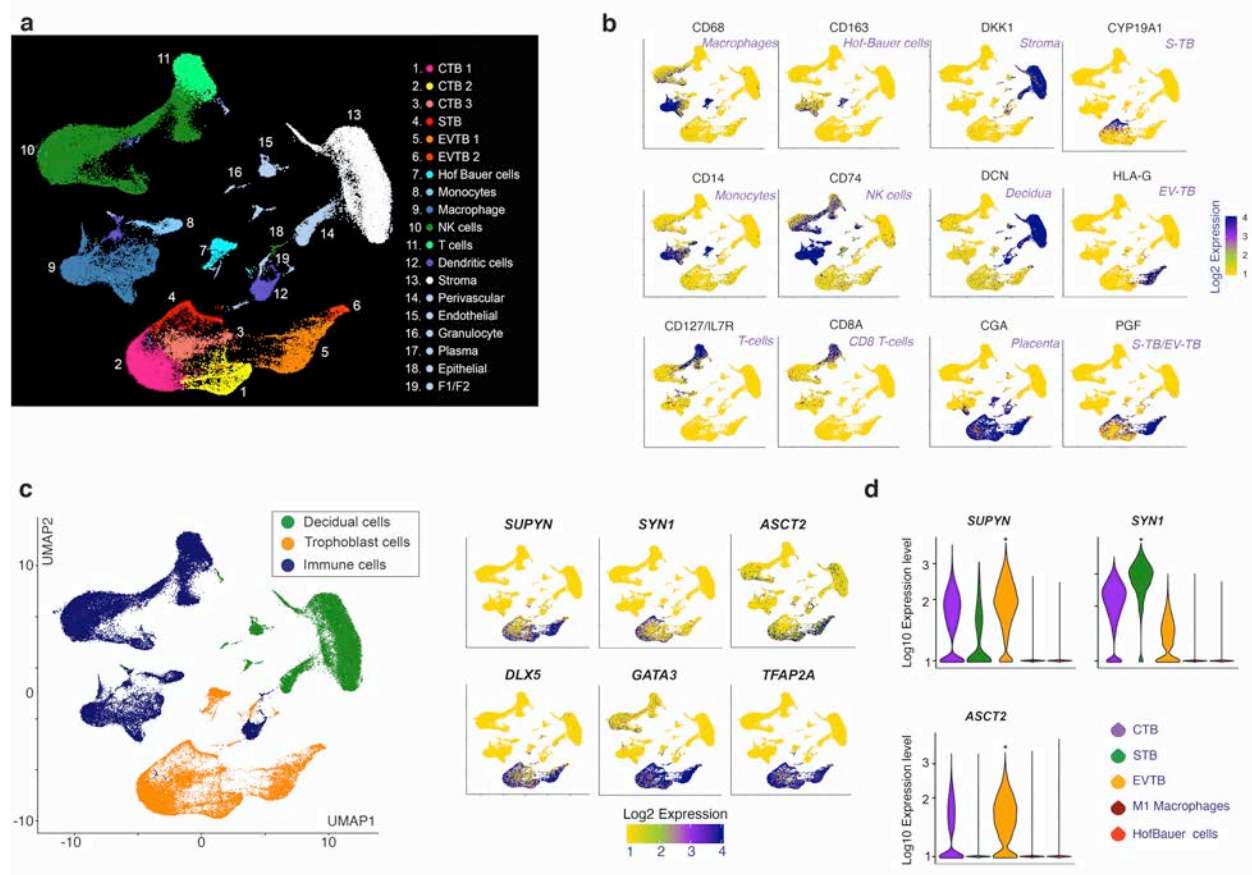


928

929 **Extended Data Figure 1: SUPYN is expressed in human early embryo and placenta but in**  
930 **293T cells.**

931 **(a)** Bar graphs show SUPYN and SYN1 expression in indicated cell types. **(b)** Genome browser  
932 view showing ATAC-seq and DNase-seq signals at the SUPYN locus, including upstream and  
933 downstream sequences. Framed region highlights overlapping peaks at the SUPYN locus. **(c)** Line  
934 plot depicts HERVenv gene expression level during BMP4-mediated *in vitro* hESCs to trophoblast  
935 differentiation. Time points correspond to cells harvested 8hr, 24hr, 48hr, and 72hr post BMP4  
936 treatment.

937

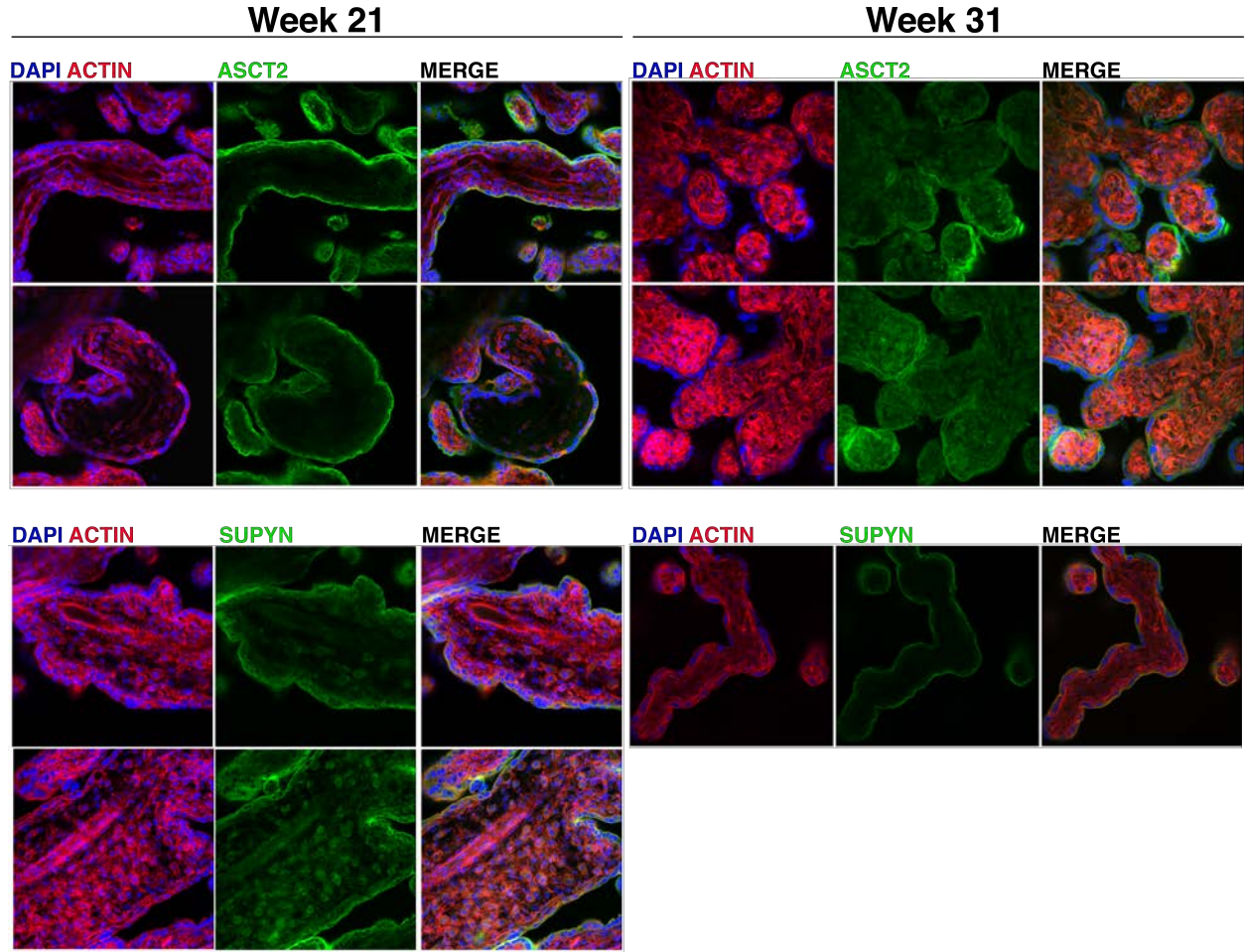


938

939 **Extended Data Figure 2: Defining lineage-specific SYN1, SUPYN, and ASCT2 expression from**  
940 **placental single-cell transcriptomics**

941 **(a)** UMAP plot generated from published scRNA-seq data generated from 1<sup>st</sup> trimester placental  
942 explants. Colors denote placental (pink, red, orange, and yellow) immune (blue and green) and  
943 maternal cell lineages (white and grey). **(b)** Feature plots visualize single-cell expression level of  
944 lineage-defining marker genes. **(c)** Simplified UMAP plot, shown in **a**, of scRNAseq data displaying  
945 trophoblast (yellow), decidual (green) and immune (purple) cell identity. Sub-panels display  
946 single-cell-level expression of indicated genes. **(d)** Violin plots denote single-cell *SUPYN* and  
947 *ASCT2* expression in multiple placental-cell lineages.

948



949

950 **Extended Data Figure 3: ASCT2 and SUPYN expression in 2<sup>nd</sup> and 3<sup>rd</sup> trimester human placenta.**

951 Confocal microscopy of 2<sup>nd</sup> (week 21) and 3<sup>rd</sup> (week 31) trimester placental villi explants. Villi

952 were stained for ASCT2 (green upper panels) or SUPYN (green lower panels) and Actin (red). Cell

953 nuclei are marked with DAPI (blue).

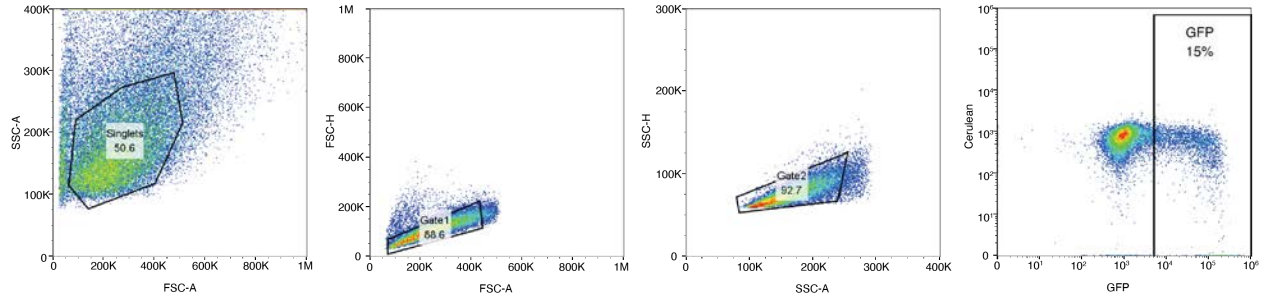
954

955

956

957





958

959 **Extended Data Figure 4: Flow Cytometry analysis scheme.** Representative sequential gating

960 scheme to assess reporter virus infection rate.

961

962

963

964

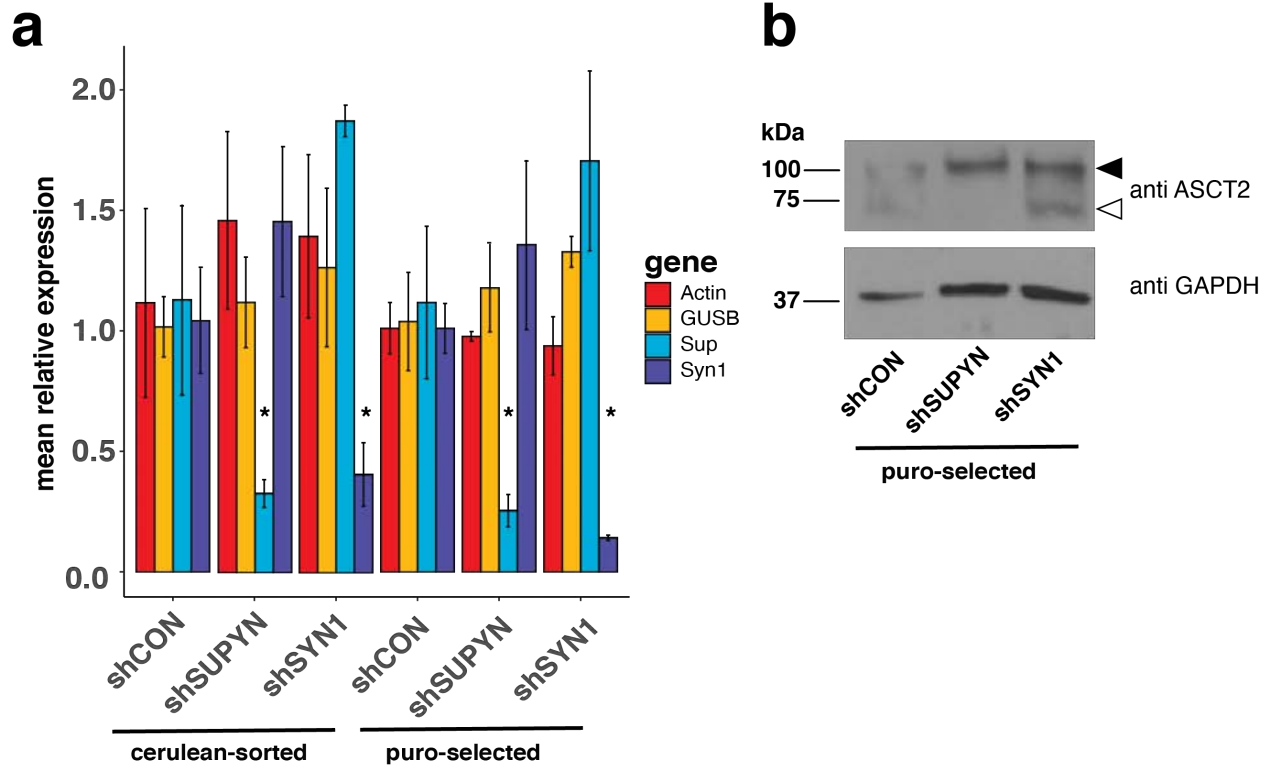
965

966

967

968

969



970

971 **Extended Data Figure 5: Characterization of shRNA transduced Jar cells and validation of env**  
972 **overexpression constructs.**

973 **(a, b)** *SUPYN* and *SYN1* knock down was validated by qPCR. Bar plots represent mean relative  
974 gene expression normalized to shCON in cerulean-sorted and puromycin-selected cell lines  
975 respectively (n = 3). Error bars represent  $\pm$  standard error mean (\* $p < 0.1$ ; Wilcoxon rank sum test).

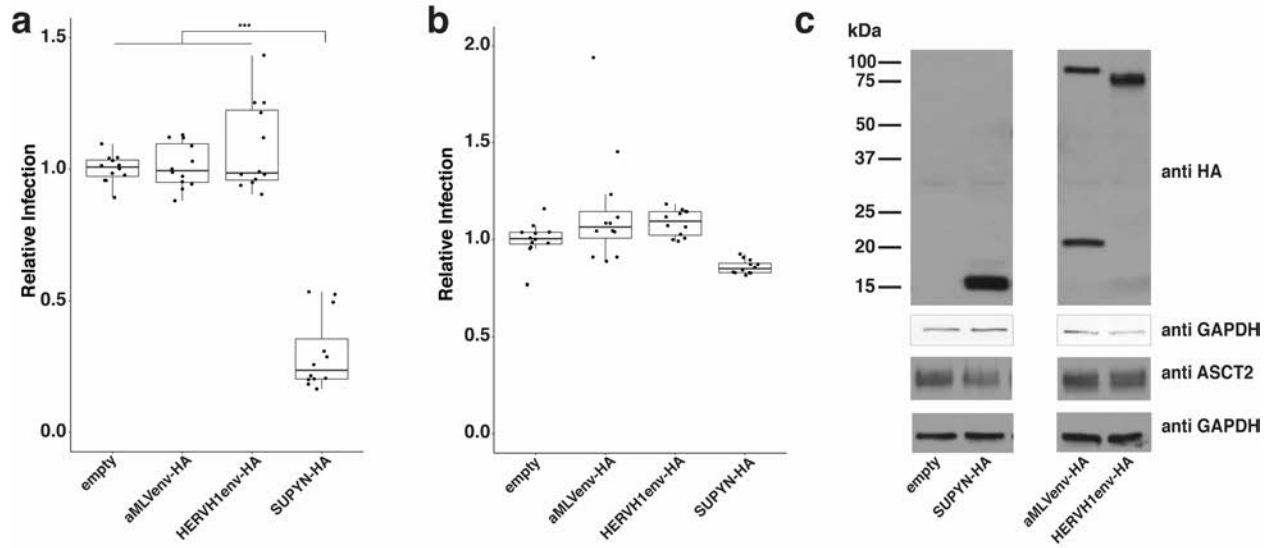
976 **(b)** Western Blot analysis ( $\alpha$ GAPDH,  $\alpha$ ASCT2) of shRNA-transduced Jar cell lysates. Putatively  
977 glycosylated and unglycosylated ASCT2 are marked by filled and empty arrowheads respectively.

978

979

980

981



982

983 **Extended Data Figure 6: SUPYN expression is sufficient to specifically restrict RDR env-**  
984 **mediated infection**

985 **(a, b)** 293T cells, transfected with HA-tagged SUPYN and env constructs, were infected with HIV-  
986 RD114env **(a)** and -VSVg **(b)** respectively. Relative infection rates were determined by normalizing  
987 GFP+ counts to empty vector. ( $n \geq 3$  with  $\geq 1$  technical replicate; \*\*\*adj.  $p < 0.001$ ; \*adj.  $p < 0.05$ ;  
988 Tukey HSD). **(c)** Western Blot analysis ( $\alpha$ HA,  $\alpha$ GAPDH,  $\alpha$ ASCT2) of 293T cell lysates following  
989 transfection with indicated constructs.

990

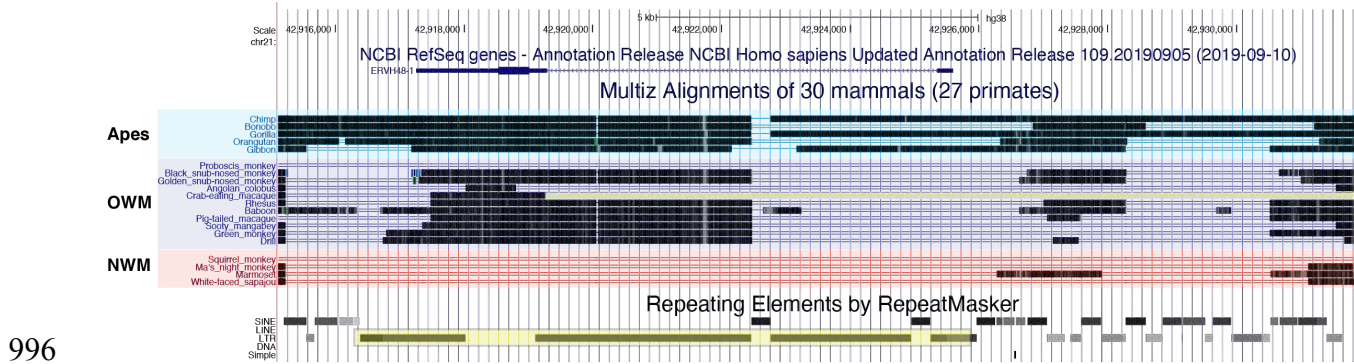
991

992

993

994

995



996

997 **Extended Data Figure 7: SUPYN locus conservation in primates.**

998 UCSC genome Browser snapshot of *SUPYN*-coding locus with surrounding sequence.

999 Modified NCBI RefSeq gene, simian whole genome alignment (from Multiz 30-species track), and

1000 RepeatMasker repetitive element tracks are shown. The *SUPYN*-coding ERVH48 provirus is

1001 highlighted by the yellow box.

1002

Sup\_Human1-555  
 Sup\_Bonobo1-549  
 Sup\_Chrimp1-549  
 Sup\_Gorilla1-555  
 Sup\_Orangutan1-549  
 Sup\_Northern\_white-cheeked\_Gibbon1-555  
 Sup\_Yellow-cheeked\_Gibbon1-555  
 Sup\_Platyotus\_Gibbon1-554  
 Sup\_Lar\_Gibbon1-555  
 Sup\_Silvery\_Gibbon1-555  
 Sup\_Milleri1\_Gibbon1-555  
 Sup\_Hooked\_Gibbon1-555  
 Sup\_Angolan\_Colibou1-201  
 Sup\_Ugandan\_Colibou1-547  
 Sup\_Black\_Stub-nosed\_Monkey1-525  
 Sup\_Guinea\_Stub-nosed\_Monkey1-525  
 Sup\_Prothomaei\_Monkey1-2  
 Sup\_Olive\_Baboon1-555  
 Sup\_Gelata1-552  
 Sup\_Dhr1-555  
 Sup\_Soxy\_Mangabey1-552  
 Sup\_Che-wating\_Macaque1-555  
 Sup\_Rhesus\_Macaque1-555  
 Sup\_Japanese\_Macaque1-555  
 Sup\_Southern\_Pipistrelle\_Macaque1-555  
 Sup\_Green\_Monkey1-555

150  
 151  
 152  
 153  
 154  
 155  
 156  
 157  
 158  
 159  
 160  
 161  
 162  
 163  
 164  
 165  
 166  
 167  
 168  
 169  
 170  
 171  
 172  
 173  
 174  
 175  
 176  
 177  
 178  
 179  
 180  
 181  
 182  
 183  
 184  
 185  
 186  
 187  
 188  
 189  
 190  
 191  
 192  
 193  
 194  
 195  
 196  
 197  
 198  
 199  
 200  
 201  
 202  
 203  
 204  
 205  
 206  
 207  
 208  
 209  
 210  
 211  
 212  
 213  
 214  
 215  
 216  
 217  
 218  
 219  
 220  
 221  
 222  
 223  
 224  
 225  
 226  
 227  
 228  
 229  
 230  
 231  
 232  
 233  
 234  
 235  
 236  
 237  
 238  
 239  
 240  
 241  
 242  
 243  
 244  
 245  
 246  
 247  
 248  
 249  
 250  
 251  
 252  
 253  
 254  
 255  
 256  
 257  
 258  
 259  
 260  
 261  
 262  
 263  
 264  
 265  
 266  
 267  
 268  
 269  
 270  
 271  
 272  
 273  
 274  
 275  
 276  
 277  
 278  
 279  
 280  
 281  
 282  
 283  
 284  
 285  
 286  
 287  
 288  
 289  
 290  
 291  
 292  
 293  
 294  
 295  
 296  
 297  
 298  
 299  
 300  
 301  
 302  
 303  
 304  
 305  
 306  
 307  
 308  
 309  
 310  
 311  
 312  
 313  
 314  
 315  
 316  
 317  
 318  
 319  
 320  
 321  
 322  
 323  
 324  
 325  
 326  
 327  
 328  
 329  
 330  
 331  
 332  
 333  
 334  
 335  
 336  
 337  
 338  
 339  
 340  
 341  
 342  
 343  
 344  
 345  
 346  
 347  
 348  
 349  
 350  
 351  
 352  
 353  
 354  
 355  
 356  
 357  
 358  
 359  
 360  
 361  
 362  
 363  
 364  
 365  
 366  
 367  
 368  
 369  
 370  
 371  
 372  
 373  
 374  
 375  
 376  
 377  
 378  
 379  
 380  
 381  
 382  
 383  
 384  
 385  
 386  
 387  
 388  
 389  
 390  
 391  
 392  
 393  
 394  
 395  
 396  
 397  
 398  
 399  
 400  
 401  
 402  
 403  
 404  
 405  
 406  
 407  
 408  
 409  
 410  
 411  
 412  
 413  
 414  
 415  
 416  
 417  
 418  
 419  
 420  
 421  
 422  
 423  
 424  
 425  
 426  
 427  
 428  
 429  
 430  
 431  
 432  
 433  
 434  
 435  
 436  
 437  
 438  
 439  
 440  
 441  
 442  
 443  
 444  
 445  
 446  
 447  
 448  
 449  
 450  
 451  
 452  
 453  
 454  
 455  
 456  
 457  
 458  
 459  
 460  
 461  
 462  
 463  
 464  
 465  
 466  
 467  
 468  
 469  
 470  
 471  
 472  
 473  
 474  
 475  
 476  
 477  
 478  
 479  
 480  
 481  
 482  
 483  
 484  
 485  
 486  
 487  
 488  
 489  
 490  
 491  
 492  
 493  
 494  
 495  
 496  
 497  
 498  
 499  
 500

Extended Data Figure 8: Nucleic acid sequence alignment of primate Suppressyn orthologs.

Suppressyn encoding nucleotide sequences are shaded blue based on a minimum sequence identity threshold of 45% (light), 75% (medium) and 80% (dark). Conserved ape-specific and

ancestral stop codons are highlighted in red.

1003  
1004  
1005  
1006  
1007  
1008

SUPYN_Hominoid_Reconstruction_consensus	1 MACIYPTTCYTSLPTKSLNTGISLTTLLLSAVALLSTAAPLSGHECYQSLYRGKMDQYFTYHTHIERSCYGTLIEECVEGKSGKSYKVKNLGVSGSRNGAIC	103
SUPYN_Human	1 MACIYPTTCYTSLPTKSLNTGISLTTLLLSAVALLSTAAPLSGHECYQSLYRGKMDQYFTYHTHIERSCYGTLIEECVEGKSGKSYKVKNLGVSGSRNGAIC	103
SUPYN_Bonobo	1 MACIYPTTCYTSLPTKSLNTGISLTTLLLSAVALLSTAAPLSGHECYQSLYRGKMDQYFTYHTHIERSCYGTLIEECVEGKSGKSYKVKNLGVSGSRNGAIC	103
SUPYN_Chimp	1 MACIYPTTCYTSLPTKSLNTGISLTTLLLSAVALLSTAAPLSGHECYQSLYRGKMDQYFTYHTHIERSCYGTLIEECVEGKSGKSYKVKNLGVSGSRNGAIC	103
SUPYN_Gorilla	1 MACIYPTTCYTSLPTKSLNTGISLTTLLLSAVALLSTAAPLSGHECYQSLYRGKMDQYFTYHTHIERSCYGTLIEECVEGKSGKSYKVKNLGVSGSRNGAIC	103
SUPYN_Orangutan	1 MACIYPTTCYTSLPTKSLNTGISLTTLLLSAVALLSTAAPLSGHECYQSLYRGKMDQYFTYHTHIERSCYGTLIEECVEGKSGKSYKVKNLGVSGSRNGAIC	103
SUPYN_Northern_whitecheeked_Gibbon	1 MACIYPTTCYTSLPTKSLNTGISLTTLLLSAVALLSTAAPLSGHECYQSLYRGKMDQYFTYHTHIERSCYGTLIEECVEGKSGKSYKVKNLGVSGSRNGAIC	103
SUPYN_Yellowcheeked_Gibbon	1 MACIYPTTCYTSLPTKSLNTGISLTTLLLSAVALLSTAAPLSGHECYQSLYRGKMDQYFTYHTHIERSCYGTLIEECVEGKSGKSYKVKNLGVSGSRNGAIC	103
SUPYN_Pileated_Gibbon	1 MACIYPTTCYTSLPTKSLNTGISLTTLLLSAVALLSTAAPLSGHECYQSLYRGKMDQYFTYHTHIERSCYGTLIEECVEGKSGKSYKVKNLGVSGSRNGAIC	103
SUPYN_Lar_Gibbon	1 MACIYPTTCYTSLPTKSLNTGISLTTLLLSAVALLSTAAPLSGHECYQSLYRGKMDQYFTYHTHIERSCYGTLIEECVEGKSGKSYKVKNLGVSGSRNGAIC	103
SUPYN_Silvery_Gibbon	1 MACIYPTTCYTSLPTKSLNTGISLTTLLLSAVALLSTAAPLSGHECYQSLYRGKMDQYFTYHTHIERSCYGTLIEECVEGKSGKSYKVKNLGVSGSRNGAIC	103
SUPYN_Müller's_Gibbon	1 MACIYPTTCYTSLPTKSLNTGISLTTLLLSAVALLSTAAPLSGHECYQSLYRGKMDQYFTYHTHIERSCYGTLIEECVEGKSGKSYKVKNLGVSGSRNGAIC	103
SUPYN_Siamang	1 MACIYPTTCYTSLPTKSLNTGISLTTLLLSAVALLSTAAPLSGHECYQSLYRGKMDQYFTYHTHIERSCYGTLIEECVEGKSGKSYKVKNLGVSGSRNGAIC	103
SUPYN_Hoolock_Gibbon	1 MACIYPTTCYTSLPTKSLNTGISLTTLLLSAVALLSTAAPLSGHECYQSLYRGKMDQYFTYHTHIERSCYGTLIEECVEGKSGKSYKVKNLGVSGSRNGAIC	103
SUPYN_OWM_Reconstruction_consensus	1 MACIYPTTCYTSLPTKSLNTGISLTTLLLSAVALLSTAAPLSGHECYQSLYRGKMDQYFTYHTHIERSCYGTLIEECVEGKSGKSYKVKNLGVSGSRNGAIC	103
SUPYN_Ugandan_Colobus	1 MACIYPTTCYTSLPTKSLNTGISLTTLLLSAVALLSTAAPLSGHECYQSLYRGKMDQYFTYHTHIERSCYGTLIEECVEGKSGKSYKVKNLGVSGSRNGAIC	103
SUPYN_Black_Snubnosed_Monkey	1 MAGTYPTACYTSLPPKSLNTGISLTPPLLLLSAVALLSTAAPLSGHECYQSLYRGKMDQYFTYHTHIERSCYGTLIEECVEGKSGKSYKVKNLGVSGSRNGAIC	103
SUPYN_Golden_Snubnosed_Monkey	1 MAGTYPTACYTSLPPKSLNTGISLTPPLLLLSAVALLSTAAPLSGHECYQSLYRGKMDQYFTYHTHIERSCYGTLIEECVEGKSGKSYKVKNLGVSGSRNGAIC	103
SUPYN_Olive_Baboon	1 MAGTYPTACYTSLPPKSLNTGISLTPPLLLLSAVALLSTAAPLSGHECYQSLYRGKMDQYFTYHTHIERSCYGTLIEECVEGKSGKSYKVKNLGVSGSRNGAIC	103
SUPYN_Gelada	1 MAGTYPTACYTSLPPKSLNTGISLTPPLLLLSAVALLSTAAPLSGHECYQSLYRGKMDQYFTYHTHIERSCYGTLIEECVEGKSGKSYKVKNLGVSGSRNGAIC	102
SUPYN_Drill	1 MAGTYPTACYTSLPPKSLNTGISLTPPLLLLSAVALLSTAAPLSGHECYQSLYRGKMDQYFTYHTHIERSCYGTLIEECVEGKSGKSYKVKNLGVSGSRNGAIC	103
SUPYN_Sooty_Mangabey	1 MAGTYPTACYTSLPPKSLNTGISLTPPLLLLSAVALLSTAAPLSGHECYQSLYRGKMDQYFTYHTHIERSCYGTLIEECVEGKSGKSYKVKNLGVSGSRNGAIC	102
SUPYN_Crab-eating_Macaque	1 MAGTYPTACYTSLPPKSLNTGISLTPPLLLLSAVALLSTAAPLSGHECYQSLYRGKMDQYFTYHTHIERSCYGTLIEECVEGKSGKSYKVKNLGVSGSRNGAIC	103
SUPYN_Rhesus_Macaque	1 MAGTYPTACYTSLPPKSLNTGISLTPPLLLLSAVALLSTAAPLSGHECYQSLYRGKMDQYFTYHTHIERSCYGTLIEECVEGKSGKSYKVKNLGVSGSRNGAIC	103
SUPYN_Japanese_Macaque	1 MAGTYPTACYTSLPPKSLNTGISLTPPLLLLSAVALLSTAAPLSGHECYQSLYRGKMDQYFTYHTHIERSCYGTLIEECVEGKSGKSYKVKNLGVSGSRNGAIC	103
SUPYN_Southern_Pigtailed_Macaque	1 MAGTYPTACYTSLPPKSLNTGISLTPPLLLLSAVALLSTAAPLSGHECYQSLYRGKMDQYFTYHTHIERSCYGTLIEECVEGKSGKSYKVKNLGVSGSRNGAIC	103
SUPYN_Green_Monkey	1 MAGTYPTACYTSLPPKSLNTGISLTPPLLLLSAVALLSTAAPLSGHECYQSLYRGKMDQYFTYHTHIERSCYGTLIEECVEGKSGKSYKVKNLGVSGSRNGAIC	103
SUPYN_Hominoid_Reconstruction_consensus	104 PRGKQWLCFTKIGOWGNTQVLEDIKREQIIAKAKASKPTTTPPENHPRHFFHSFIQKL*	161
SUPYN_Human	104 PRGKQWLCFTKIGOWGNTQVLEDIKREQIIAKAKASKPTTTPPENHPRHFFHSFIQKL*	161
SUPYN_Bonobo	104 PRGKQWLCFTKIGOWGNTQVLEDIKREQIIAKAKASKPTTTPPENHPRHFFHSFIQKL*	159
SUPYN_Chimp	104 PRGKQWLCFTKIGOWGNTQVLEDIKREQIIAKAKASKPTTTPPENHPRHFFHSFIQKL*	159
SUPYN_Gorilla	104 PRGKQWLCFTKIGOWGNTQVLEDIKREQIIAKAKASKPTTTPPENHPRHFFHSFIQKL*	161
SUPYN_Orangutan	104 PRGKQWLCFTKIGOWGNTQVLEDIKREQIIAKAKASKPTTTPPENHPRHFFHSFIQKL*	159
SUPYN_Northern_whitecheeked_Gibbon	104 PRGKQWLCFTKIGOWGNTQVLEDIKREQIIAKAKASKPTTTPPENHPRHFFHSFIQKL*	161
SUPYN_Yellowcheeked_Gibbon	104 PRGKQWLCFTKIGOWGNTQVLEDIKREQIIAKAKASKPTTTPPENHPRHFFHSFIQKL*	161
SUPYN_Pileated_Gibbon	104 PRGKQWLCFTKIGOWGNTQVLEDIKREQIIAKAKASKPTTTPPENHPRHFFHSFIQKL*	161
SUPYN_Lar_Gibbon	104 PRGKQWLCFTKIGOWGNTQVLEDIKREQIIAKAKASKPTTTPPENHPRHFFHSFIQKL*	161
SUPYN_Silvery_Gibbon	104 PRGKQWLCFTKIGOWGNTQVLEDIKREQIIAKAKASKPTTTPPENHPRHFFHSFIQKL*	161
SUPYN_Müller's_Gibbon	104 PRGKQWLCFTKIGOWGNTQVLEDIKREQIIAKAKASKPTTTPPENHPRHFFHSFIQKL*	161
SUPYN_Siamang	104 PRGKQWLCFTKIGOWGNTQVLEDIKREQIIAKAKASKPTTTPPENHPRHFFHSFIQKL*	161
SUPYN_Hoolock_Gibbon	104 PRGKQWLCFTKIGOWGNTQVLEDIKREQIIAKAKASKPTTTPPENHPRHFFHSFIQKL*	161
SUPYN_OWM_Reconstruction_consensus	104 PRGKQWLCFTKIGOWGNTQVLEDIKREQIIAKAKASKPTTTPPENHPRHFFHSFIQKL*	185
SUPYN_Ugandan_Colobus	104 PRGKQWLCFTKIGOWGNTQVLEDIKREQIIAKAKASKPTTTPPENHPRHFFHSFIQKL*	206
SUPYN_Black_Snubnosed_Monkey	104 PRGKQWLCFTKIGOWGNTQVLEDIKREQIIAKAKASKPTTTPPENHPRHFFHSFIQKL*	185
SUPYN_Golden_Snubnosed_Monkey	104 PRGKQWLCFTKIGOWGNTQVLEDIKREQIIAKAKASKPTTTPPENHPRHFFHSFIQKL*	184
SUPYN_Olive_Baboon	104 PRGKQWLCFTKIGOWGNTQVLEDIKREQIIAKAKASKPTTTPPENHPRHFFHSFIQKL*	185
SUPYN_Gelada	104 PRGKQWLCFTKIGOWGNTQVLEDIKREQIIAKAKASKPTTTPPENHPRHFFHSFIQKL*	184
SUPYN_Drill	104 PRGKQWLCFTKIGOWGNTQVLEDIKREQIIAKAKASKPTTTPPENHPRHFFHSFIQKL*	185
SUPYN_Sooty_Mangabey	104 PRGKQWLCFTKIGOWGNTQVLEDIKREQIIAKAKASKPTTTPPENHPRHFFHSFIQKL*	184
SUPYN_Crab-eating_Macaque	104 PRGKQWLCFTKIGOWGNTQVLEDIKREQIIAKAKASKPTTTPPENHPRHFFHSFIQKL*	185
SUPYN_Rhesus_Macaque	104 PRGKQWLCFTKIGOWGNTQVLEDIKREQIIAKAKASKPTTTPPENHPRHFFHSFIQKL*	185
SUPYN_Japanese_Macaque	104 PRGKQWLCFTKIGOWGNTQVLEDIKREQIIAKAKASKPTTTPPENHPRHFFHSFIQKL*	185
SUPYN_Southern_Pigtailed_Macaque	104 PRGKQWLCFTKIGOWGNTQVLEDIKREQIIAKAKASKPTTTPPENHPRHFFHSFIQKL*	185
SUPYN_Green_Monkey	104 PRGKQWLCFTKIGOWGNTQVLEDIKREQIIAKAKASKPTTTPPENHPRHFFHSFIQKL*	185
SUPYN_Hominoid_Reconstruction_consensus	207 HPYLLASQNP SLTFTPNARSPGHLPTQ*	235
SUPYN_Human		
SUPYN_Bonobo		
SUPYN_Chimp		
SUPYN_Gorilla		
SUPYN_Orangutan		
SUPYN_Northern_whitecheeked_Gibbon		
SUPYN_Yellowcheeked_Gibbon		
SUPYN_Pileated_Gibbon		
SUPYN_Lar_Gibbon		
SUPYN_Silvery_Gibbon		
SUPYN_Müller's_Gibbon		
SUPYN_Siamang		
SUPYN_Hoolock_Gibbon		
SUPYN_OWM_Reconstruction_consensus		
SUPYN_Ugandan_Colobus		
SUPYN_Black_Snubnosed_Monkey		
SUPYN_Golden_Snubnosed_Monkey		
SUPYN_Olive_Baboon		
SUPYN_Gelada		
SUPYN_Drill		
SUPYN_Sooty_Mangabey		
SUPYN_Crab-eating_Macaque		
SUPYN_Rhesus_Macaque		
SUPYN_Japanese_Macaque		
SUPYN_Southern_Pigtailed_Macaque		
SUPYN_Green_Monkey		

1009

1010 **Extended Data Figure 9: Amino Acid sequence alignment of primate SUPYN orthologs.**

1011 Primate and ancestral SUPYN peptide sequences are shown. Sequences are shaded blue based

1012 on a minimum sequence identity threshold of 45% (light), 75% (medium) and 80% (dark)

1013 Ancestral SUPYN sequences are based on our consensus-based sequence reconstruction (see

1014 Methods).

1015

```
SUPYN_Hominoid_Reconstruction_consensus 1 MACIYPTTCYTSLPTKSLNTGI SLTTI L I LSVAVLLSTAAPPSCHECYOSLHYRGKMOQYFTYHTH I ERSCYGT L I EECVE SGKSYKVKNLGVSGSR 98
SUPYN_Hominoid_Reconstruction_PAML_Model3 1 MACIYPTTCYTSLPTKSLNTGI SLTTI L I LSVAVLLSTAAPPSCHECYOSLHYRGKMOQYFTYHTH I ERSCYGT L I EECVE SGKSYKVKNLGVSGSR 98
SUPYN_Hominoid_Reconstruction_PAML_Model4 1 MACIYPTTCYTSLPTKSLNTGI SLTTI L I LSVAVLLSTAAPPSCHECYOSLHYRGKMOQYFTYHTH I ERSCYGT L I EECVE SGKSYKVKNLGVSGSR 98
SUPYN_Hominoid_Reconstruction_PAML_Model5 1 MACIYPTTCYTSLPTKSLNTGI SLTTI L I LSVAVLLSTAAPPSCHECYOSLHYRGKMOQYFTYHTH I ERSCYGT L I EECVE SGKSYKVKNLGVSGSR 98
SUPYN_Hominoid_Reconstruction_PAML_Model6 1 MACIYPTTCYTSLPTKSLNTGI SLTTI L I LSVAVLLSTAAPPSCHECYOSLHYRGKMOQYFTYHTH I ERSCYGT L I EECVE SGKSYKVKNLGVSGSR 98
SUPYN_Hominoid_Reconstruction_PAML_Model7 1 MACIYPTTCYTSLPTKSLNTGI SLTTI L I LSVAVLLSTAAPPSCHECYOSLHYRGKMOQYFTYHTH I ERSCYGT L I EECVE SGKSYKVKNLGVSGSR 98
SUPYN_OWM_Reconstruction_consensus 1 MACTYPTACYTSLPPKSLNTGI SLTPI L I LSVAVLLSAAAPPSCRECYOSFHYRGK IQQSFTYHTH I ERSCYGT L I EECVE SGKSYKVKNLGVSGSR 98
SUPYN_OWM_Reconstruction_PAML_Model3 1 MACTYPTACYTSLPPKSLNTGI SLTPI L I LSVAVLLSAAAPPSCRECYOSFHYRGK IQQSFTYHTH I ERSCYGT L I EECVE SGKSYKVKNLGVSGSR 98
SUPYN_OWM_Reconstruction_PAML_Model4 1 MACTYPTACYTSLPPKSLNTGI SLTPI L I LSVAVLLSAAAPPSCRECYOSFHYRGK IQQSFTYHTH I ERSCYGT L I EECVE SGKSYKVKNLGVSGSR 98
SUPYN_OWM_Reconstruction_PAML_Model5 1 MACTYPTACYTSLPPKSLNTGI SLTPI L I LSVAVLLSAAAPPSCRECYOSFHYRGK IQQSFTYHTH I ERSCYGT L I EECVE SGKSYKVKNLGVSGSR 98
SUPYN_OWM_Reconstruction_PAML_Model6 1 MACTYPTACYTSLPPKSLNTGI SLTPI L I LSVAVLLSAAAPPSCRECYOSFHYRGK IQQSFTYHTH I ERSCYGT L I EECVE SGKSYKVKNLGVSGSR 98
SUPYN_OWM_Reconstruction_PAML_Model7 1 MACTYPTACYTSLPPKSLNTGI SLTPI L I LSVAVLLSAAAPPSCRECYOSFHYRGK IQQSFTYHTH I ERSCYGT L I EECVE SGKSYKVKNLGVSGSR 98

SUPYN_Hominoid_Reconstruction_consensus 99 NGAICPRGKQWLCFTK I GQGWVNTQVLEDI KREQI I AKAKA SKPTT PENHPRH FHSF I OKL ADA SLPKPKGN L FVD LGEP SRLP 185
SUPYN_Hominoid_Reconstruction_PAML_Model3 99 NGAICPRGKQWLCFTK I GQGWVNTQVLEDI KREQI I AKAKA SKPTT PENHPRH FHSF I OKLQADA SLPKPKGYL FVD LGEP SRLP 184
SUPYN_Hominoid_Reconstruction_PAML_Model4 99 NGAICPRGKQWLCFTK I GQGWVNTQVLEDI KREQI I AKAKA SKPTT PENHPRH FHSF I OKLQADA SLPKPKGYL FVD LGEP SRLP 184
SUPYN_Hominoid_Reconstruction_PAML_Model5 99 NGAICPRGKQWLCFTK I GQGWVNTQVLEDI KREQI I AKAKA SKPTT PENHPRH FHSF I OKLQADA SLPKPKGYL FVD LGEP SRLP 184
SUPYN_Hominoid_Reconstruction_PAML_Model6 99 NGAICPRGKQWLCFTK I GQGWVNTQVLEDI KREQI I AKAKA SKPTT PENHPRH FHSF I OKLQADA SLPKPKGYL FVD LGEP SRLP 184
SUPYN_Hominoid_Reconstruction_PAML_Model7 99 NGAICPRGKQWLCFTK I GQGWVNTQVLEDI KREQI I AKAKA SKPTT PENHPRH FHSF I OKLQADA SLPKPKGYL FVD LGEP SRLP 184
SUPYN_OWM_Reconstruction_consensus 99 NGAICPQGKQWLCFTK I GQGWVNTQVLEDI KREQI I AKAKA SKPTT PENHPRY FHSF I RKLOADA SLPKPKGYL FVD LGEP SRLP 185
SUPYN_OWM_Reconstruction_PAML_Model3 99 NGAICPQGKQWLCFTK I GQGWVNTQVLEDI KREQI I AKAKA SKPTT PENHPRY FHSF I RKLOADA SLPKPKGYL FVD LGEP SRLP 184
SUPYN_OWM_Reconstruction_PAML_Model4 99 NGAICPQGKQWLCFTK I GQGWVNTQVLEDI KREQI I AKAKA SKPTT PENHPRY FHSF I RKLOADA SLPKPKGYL FVD LGEP SRLP 184
SUPYN_OWM_Reconstruction_PAML_Model5 99 NGAICPQGKQWLCFTK I GQGWVNTQVLEDI KREQI I AKAKA SKPTT PENHPRY FHSF I RKLOADA SLPKPKGYL FVD LGEP SRLP 184
SUPYN_OWM_Reconstruction_PAML_Model6 99 NGAICPQGKQWLCFTK I GQGWVNTQVLEDI KREQI I AKAKA SKPTT PENHPRY FHSF I RKLOADA SLPKPKGYL FVD LGEP SRLP 184
SUPYN_OWM_Reconstruction_PAML_Model7 99 NGAICPQGKQWLCFTK I GQGWVNTQVLEDI KREQI I AKAKA SKPTT PENHPRY FHSF I RKLOADA SLPKPKGYL FVD LGEP SRLP 184
```

1016

1017 **Extended Data Figure 10: Amino Acid sequence alignment of ancestral SUPYN sequences.**

1018 Peptide sequences of consensus and maximum likelihood-based SUPYN sequence

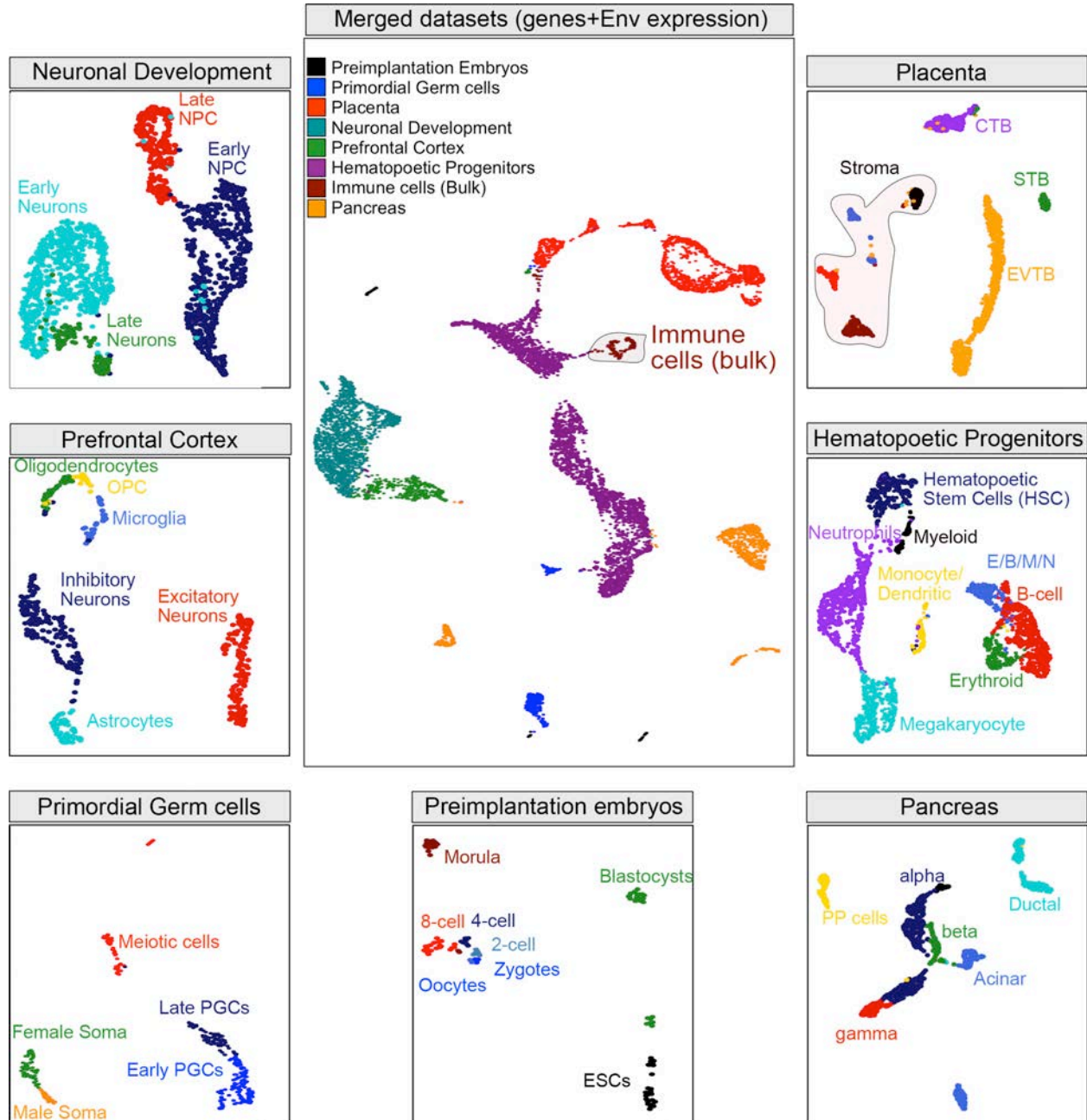
1019 reconstructions (see methods) are aligned.

1020

1021

1022

1023



1024

1025 **Extended Data Figure 11: Analysis of envORF expression in single cell RNA-seq taken from**

1026 **various human tissue sources**

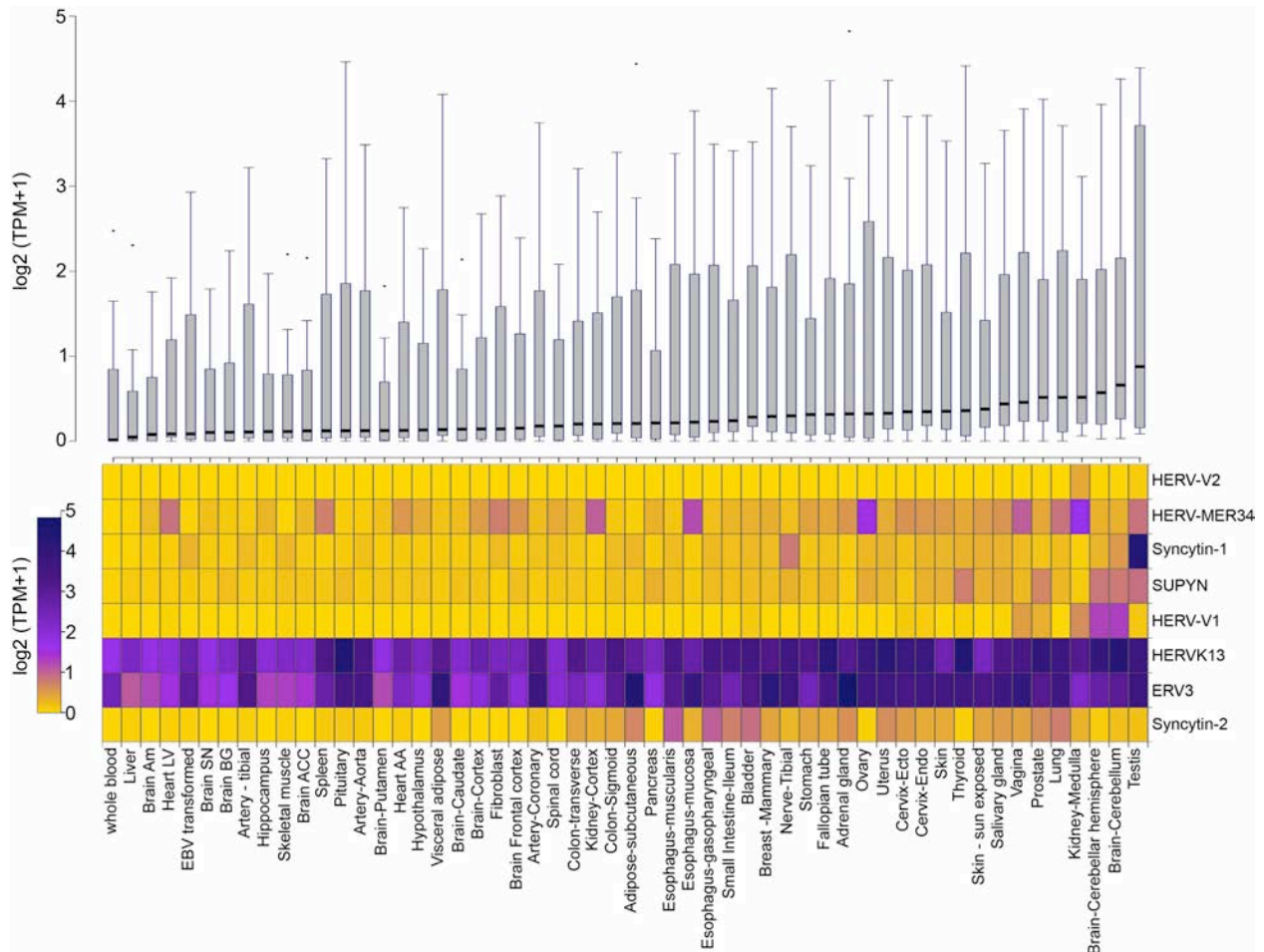
1027 **(a)** UMAP plots illustrate the Louvain clustering of seven independent sc-RNA-seq datasets

1028 corresponding to the development of human embryos and somatic tissues (see Methods). The

1029 large central UMAP plot represents our integrative analysis of seven scRNA-seq and one bulk



1030 RNA-seq datasets obtained from independent studies (Supplemental Table 1). Surrounding  
1031 UMAP plots represent the combined expression of human RefSeq and envORF genes  
1032 (supplemental Table 3) with cell identities labeled. **(b)** Multiple violin plots demonstrate the  
1033 expression of annotated ERV envelopes in each dataset shown in **a**.  
1034



1035  
1036 **Extended Data 12: Annotated human endogenous retrovirus envelope expression in various**  
1037 **human tissues.**  
1038 **(a)** The boxplot shows the transcript expression distribution of annotated ERV envelopes in  
1039 tissues assayed by the GTEx project (Supplementary Table 1). **(b)** The heatmap displays the

1040 expression (log2 TPM) of individual ERV envelope genes (rows) in each tissue (columns). The color  
 1041 scheme ranges gradually from no expression (gold) to higher expression (midnight blue).

1042

1043

1044 **Tables**

1045 **Supplementary Table 1: External data sources**

Description	Author	Year	Publication	Methods Reference	Dataset	Tissue source	Figure
scRNAseq	Yan et al.	2013	PMID: 23934149	46	GSE36552	embryo	1a, 4d, Ext 11
	Vento-Tormo et al.	2018	PMID: 30429548	53	E-MTAB-6701	placenta	1d-f, Ext 2
	Liu et al.	2018	PMID: 30042384	47	GSE89497	placenta	1f, 4d, Ext 11
	Pavlicev	2017	PMID: 28174237	48	GSE87726	placenta	1f
	Close	2017	PMID: 28279351	85	GSE93593	Neuro_diff	4d, Ext 11
	Guo	2015	PMID: 26046443	83	GSE63818	PGC	4d, Ext 11
	Velten	2017	PMID: 28319093	84	GSE75478	HSC	4d, Ext 11
	Enge	2017	PMID: 28965763	87	GSE81547	Pancreas	4d, Ext 11
	Darmanis	2015	PMID: 26060301	86	GSE67835	PFC	4d, Ext 11
Bulk RNA-seq	Zadora	2017	PMID: 28904069	54	available via approval	placenta	1f
	GTEEx Consortium	2017	PMID: 29022597	88	phs000424.v6.p1	Human tissues	Ext 12
	D Antonio	2017	PMID: 28874753	55	E-MTAB-5714	293T	Ext 1a
	Chung	2018	PMID: 29395325	56	GSE99249	293T	Ext 1a
	Calderon	2019	PMID: 31570894	57	GSE118165	Immune cells	4d, Ext 11

ChIPseq	Barakat	2018	PMID: 30033119	59	GSE99631	hESC	1b
	Krendl	2017	PMID: 29078328	62	GSE105258	hESC	1c, Ext 1a-b
	Krendl	2017	PMID: 29078328	62	GSE105081	hESC	1c, Ext 1a- b
	Tsankov	2015	PMID: 25693565	58	GSE61475	hESC	1b, Ext 1d
	Dunn- Fletcher	2018	PMID: 30231016	61	GSE118289	placenta	1c
	Kwak	2019	PMID: 31294776	60	GSE127288	placenta	1c
DNase-seq	Gao	2018	PMID: 29526463	63	GSA:CRA000297	embryo	Ext 1c
ATAC-seq	Wu	2018	PMID: 29720659	64	GSE101571	embryo	Ext 1c

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064 **Suppelemental table 2: Primer list**

<b>InFusion cloning primers</b>		
<b>ID</b>	<b>sequence</b>	<b>description</b>
JF023	TTTTGGCAAAGAATTCATGGCCTGTATCTACCCA	pHCMV human suppressyn fwd
JF024	CCTGAGGAGTGAATTCTTATAGTTTTGTATAAA GGAATGG	pHCMV human suppressyn rev
JF025	TTTTGGCAAAGAATTCATGGCCTGTATCTACCCA ACC	pHCMV human suppressyn-HA fwd
JF026	CCTGAGGAGTGAATTCTTAAGCGTAATCTGGAAC ATCG	pHCMV human suppressyn-HA rev
JF046	TTTTGGCAAAGAATTCATGGCGCGTTCAACGCTC T	pHCMV amphiMLVenv fwd
JF047	CCTGAGGAGTGAATTCTCATGGCTCGTACTCTAT GGGT	pHCMV amphiMLVenv rev
JF081	GCATAGTAGTCTCATTTGCTACCACA	HERVH-1 (H62) provirus fwd
JF082	CATGACTCGGATCAGGGGAC	HERVH-1 (H62) provirus rev
JF091	ATCATTTTTGGCAAAGaattCATGATCTTTGCTGG CAAGGCACC	pHCMV HERVH-1env fwd
JF093	CAGCCTGCACCTGAGGAGTgaattCctaagcgta atctggaacatcgtatgggtaAGCTGAAGGGAGG TCTTGTGGTAAG	pHCMV HERVH-1env-HA rev
JF130	GGCGCCTAAGCTGGTATTCTTAACTATGTTGCTC C	pHIV7 eGFP downstream sequence fwd
JF131	CGCAGAGCCGGCAGCAGGCCGCGGAAGGAAGGT CCGCTGGATTG	pHIV7 eGFP downstream sequence rev
JF135b	tttggggatcctgagcggccgcatccaaggtcgg gcagga	U6 shRNA cloning fwd (pHIV-shSup cloning)
JF136	ccaaccactttctatactacacaaatatagaag tggttgggtagatacactttcgtcctttccacaa gatataaaagccaagaa	U6-shSup rev (pHIV-shSup cloning)

JF138	CTACTTCCTTTTCGATACTACACAAATATCGAAAG GAAGTAGAAGTCCGGGCTTTCGTCCTTCCACAA GATATATAAAGCCAAGAA	U6-shC rev (pHIV-shC cloning)
JF145	gtagtatagaaagtgggttagatacact ttttgtaccgagctcggatccactagagatgg	shSup-CMV promoter fwd (pHIV7-shSup cloning)
JF147	GTAGTATCGAAAGGAAGTAGAAGTCCGGGCT ttttgtaccgagctcggatccactagagatgg	shC-CMV promoter fwd (pHIV7-shC cloning)
JF152b	TGAACCGTCAGATCCGCTAGCATGGTGAGCAAGG GCGAGG	cerulean fwd (pHIV7-shRNA-cerulean cloning)
JF153b	AGAATACCAGttACTTGTACAGCTCGTCCATGCC	cerulean rev (pHIV7-shRNA-cerulean cloning)
JF154b	CTGTACAAGTAACTGGTATTCTTAACTATGTTGC TCC	cerulean-WPRE fwd (pHIV7-shRNA-cerulean cloning)
JF155	TGAACCGTCAGATCCGCTAGCATGGCCACCGAGT ACAAG	puroR fwd (pHIV7-shRNA-puroR cloning)
JF156	AGAATACCAGTTATTTGTAACCATTATAAGCTGC	puroR rev (pHIV7-shRNA-puroR cloning)
JF157	TTACAAATAACTGGTATTCTTAACTATGTTGCTC C	puroR-WPRE fwd (pHIV7-shRNA-puroR cloning)
JF165	cggtggccatGctagcggatctgacggttc	pHIV7-CMV promoter rev (pHIV-shRNA puroR cloning)
JF166	tgctcaccatgctagcggatctgacggttc	pHIV7-CMV promoter rev (pHIV-shRNA Cerulean cloning)
RAK010	TTTTGGCAAAGAATTCATGCTCTGCATCCTCATC CTCCT	pHCMV SMRVenv cloning primer fwd
RAK011	CCTGAGGAGTGAATTCCTACAGTCTGCCATATTC TAGGTCACGA	pHCMV SMRVenv cloning primer rev
<b>qPCR primers</b>		
JF042	GCAATTATTCCTCCATGAACG	18S fwd
JF043	GGCTCACTAAACCATCCAA	18S rev
JF108	CTGTGCATGCACATCGGTCACTG	Suppressyn fwd
JF109	GAGAAATTGGCCCAGACAAACT	Suppressyn rev
JF112	ACCACGAACGGACATCCAAAG	Syncytin-1 fwd
JF113	GCCACTTTAACCGCAGTTGG	Syncytin-1 rev
Act-F	CGACAGGATGCAGAAGGAG	Actin fwd
Act-R	GTAAGTGGCTCAGGAGGAG	Actin rev
GUSB-F	AGAGTGGTGCTGAGGATTGG	GUSB fwd

1065

GUSB-R	CCCTCATGCTCTAGCGTGTC	GUSB rev
--------	----------------------	----------

