

1 *November 26, 2021*

2
3 **Evolution of Human-specific Alleles Protecting Cognitive Function of Grandmothers**

4
5 Sudeshna Saha¹, Naazneen Khan¹, Troy Comi², Andrea Verhagen¹, Aniruddha Sasmal¹,
6 Sandra Diaz¹, Hai Yu³, Xi Chen³, Joshua M. Akey², Martin Frank⁴,
7 Pascal Gagneux^{1*}, and Ajit Varki^{1*}

8
9
10 ¹. Departments of Medicine, Pathology, Anthropology and Cellular and Molecular Medicine,
11 Center for Academic Research and Training in Anthropogeny and Glycobiology Research and
12 Training Center, University of California San Diego, San Diego, California 92093, USA

13 ². Department of Genetics, Princeton University, Princeton, New Jersey 08544, USA

14 ³. Department of Chemistry, University of California Davis, Davis, California 95616, USA

15 ⁴. Biognos AB, Gothenburg, SE-402 74, Sweden

16
17
18 *Address correspondence to Ajit Varki: alvarki@ucsd.edu or Pascal Gagneux:
19 pgagneux@ucsd.edu.

20
21
22
23
24

25
26

Summary (250 Words)

27 Late-onset Alzheimer’s Disease (LOAD) pathology is rare in our closest living evolutionary
28 relatives (chimpanzees), which also express much lower microglial levels of CD33(Siglec-3)—a
29 myelomonocytic receptor inhibiting innate immune reactivity by extracellular V-set domain
30 recognition of sialic acid(Sia)-containing “self-associated molecular patterns” (SAMPs). We
31 earlier showed that V-set domain-deficient *CD33*-variant allele, protective against LOAD, is
32 derived and specific to hominin-lineage. We now report that *CD33* also harbors multiple hominin-
33 specific V-set domain mutations and explore selection forces that may have favored such genomic
34 changes. *N*-glycolylneuraminic acid (Neu5Gc), the preferred Sia-ligand of ancestral CD33 is
35 absent in humans, due to hominin-specific, fixed loss-of-function mutation in CMAH, which
36 generates CMP-Neu5Gc from its precursor, CMP-*N*-acetylneuraminic acid (Neu5Ac). Extensive
37 mutational analysis and MD-simulations indicate that fixed change in amino acid 21 of hominin
38 V-set domain and conformational changes related to His45 corrected for Neu5Gc-loss by
39 switching to Neu5Ac-recognition. Considering immune-evasive “molecular mimicry” of SAMPs
40 by pathogens, we found that human-specific pathogens *Neisseria gonorrhoeae* and Group B
41 *Streptococcus* (affecting fertility and fetuses/neonates respectively) selectively bind huCD33 and
42 this binding is significantly impacted by amino acid 21 modification. Alongside LOAD-protective
43 *CD33* alleles, humans harbor additional, derived, population-universal, cognition-protective
44 variants absent in “great ape” genomes. Interestingly, 11 of 13 SNPs in these human genes
45 (including *CD33*), that protect the cognitive health of elderly populations, are not shared by
46 genomes of archaic hominins: Neanderthals and Denisovans. Finally, we present a plausible
47 evolutionary scenario to compile, correlate and comprehend existing knowledge about huCD33
48 evolution and suggest that grandmothing emerged in humans.

49

50 **Keywords (up to 10 words):** Siglec3/CD33, pathogens, sialic acids, archaic genome, molecular
51 dynamics simulation, phylogenetic analysis, menopause, grandmother.

52

53

Introduction

54 In keeping with the fundamental importance of reproduction for the process of biological evolution
55 via natural selection, loss of fecundity generally coincides with the end of lifespan in almost all
56 species studied to date. Humans and certain toothed whales like orcas are so far the only mammals
57 known to manifest prolonged post-reproductive lifespans under naturalistic conditions [1–6]. One
58 current explanation for such prolonged post-reproductive survival is late-life kin selection of
59 grandmothers and other elderly caregivers of helpless young; apparently contrary to the concept
60 of “antagonistic pleiotropy”, which posits that natural selection does not operate in late-life [7, 8].
61 An interesting exception is a human-specific derived allele of CD33 associated with direct or
62 indirect protection against late-onset Alzheimer’s Disease (LOAD) [9]. Furthermore, we noted
63 that humans harbor additional examples of such derived, population-universal gene variants that
64 directly or indirectly impact late-life cognitive decline, which were not found in other “great ape”
65 genomes. This was considered as genomic evidence for the evolution of human postmenopausal
66 longevity [10]. Here we further explore the human-specific, derived alleles of genes that protect
67 against late-life cognitive decline, and ask when and how these emerged in hominins?

68 In vertebrates, glycan-binding proteins of the immunoglobulin (Ig) superfamily called
69 sialic acid (Sia)-binding Ig-like lectins (Siglecs) form a major component of the immune system
70 [11]. As the name indicates, Siglecs recognize Sias on cell surface or secreted glycoproteins and
71 glycolipids. Siglec-3, commonly known as CD33, is the eponymous member of the rapidly
72 evolving subgroup of Siglecs called CD33-related Siglecs (or CD33rSiglecs) [12, 13]. In contrast,
73 other Siglecs (Siglecs 1, 2, 4 and 15) show evolutionary conservation [14]. CD33 is a type-I
74 transmembrane protein with an amino terminal Ig-like V-set domain followed by one Ig-like C2-
75 set domain proximal to the transmembrane region [13]. Its cytoplasmic tail contains
76 immunoregulatory signaling motifs called immunoreceptor tyrosine-based inhibitory motif
77 (ITIM)s, which upon ligand binding to the extracellular V-set domain, undergo phosphorylation
78 and recruit effector molecules like tyrosine phosphatases, SHP-1/2, which inhibit the cellular
79 immune response. Human CD33 (huCD33) binds α 2-3- and α 2-6-linked *N*-acetylneuraminic acid
80 (Neu5Ac), the predominant Sia in humans, associated either with N- and O-glycosylated
81 molecules or sialylated glycolipids (gangliosides). HuCD33 undergoes alternative splicing,
82 resulting in two isoforms – full length CD33M containing the ligand binding V-set domain and
83 truncated D2-CD33 (or CD33m) lacking this domain [15]. The elimination of the terminal V-set

84 domain is mediated through differential splicing affected by two co-inherited single nucleotide
85 polymorphisms (SNPs) at positions rs3865444 in huCD33 promoter and rs12459419 located
86 within exon 2 [16]. The two isoforms, CD33M and D2-CD33, differ not only in their molecular
87 weights, but also in their cellular localization and functionality which are associated with Sia-
88 interacting V-set domain [9, 12, 17].

89 HuCD33 is extensively studied for its role in different immune responses, under both
90 normal and pathophysiological conditions including cancers [16, 18, 18–21]. Furthermore, the
91 microglial expression of CD33 is linked with neurological pathologies like LOAD. Incidence of
92 LOAD has been strongly associated with varied expression of CD33 isoforms in the brain of
93 affected individuals [20, 21], where the LOAD-protective CD33 allele increases the ratio of D2-
94 CD33 isoform relative to CD33M. CD33 is reported in almost all vertebrates, including nonhuman
95 primates [6, 14]. While there is often high similarity in the sequence and overall genomic location,
96 CD33 has undergone various species-specific changes. For example, murine CD33 which shows
97 about 54% identity with huCD33 V-set and 72% identity with C2 domain, has markedly different
98 Sia-binding and cellular expression patterns from human CD33 protein [22]. CD33 expression has
99 greatly diverged in humans even in comparison to our closest living evolutionary relatives, the
100 great apes. Examination of CD33 in peripheral blood showed significantly increased production
101 of CD33M in human monocytes relative to those of chimpanzees [9]. Furthermore, the abundance
102 of CD33 was also markedly higher in the human brain. Interestingly, although LOAD-associated
103 neurological pathologies, for example, buildup of A β proteins, hyperphosphorylated tau proteins
104 as neurofibrillary tangles, have been observed in aged nonhuman primate brains, AD has largely
105 been regarded as a uniquely human disease [23, 24]. Interspecies variations in CD33 have also
106 been studied in other apes like gorilla and bonobo, in comparison to huCD33 [25].

107 The presence of two physiologically significant isoforms, their distinct cellular localization and
108 association with uniquely human pathologies like LOAD have made huCD33 a target of much
109 evolutionary interest. The Sia-binding V-set domain of CD33rSigslecs including CD33 itself show
110 high sequence variability amongst different species, often making it difficult to identify their
111 orthologs. The selective pressure for this accelerated evolution of the V-set domains has been
112 attributed to evasion of infectious pathogens that exploit these human innate receptors. The
113 surfaces of each vertebrate cell are layered with tens to hundreds of million Sia-terminating
114 glycans, forming as “self-associated molecular patterns” (SAMPs), which prevent erroneous

115 activation of innate immune responses against the body's own cells [26]. However, several human
116 pathogens e.g., *Neisseria gonorrhoeae*, *Neisseria meningitidis*, *Haemophilus influenzae*, *E. coli*
117 K1, Group B *Streptococcus*, and *Trypanosoma cruzi* cloak themselves with sialoglycans,
118 effectively mimicking host SAMPs, and thereby avoiding the immune response [27]. Conversely,
119 other infectious agents like influenza virus recognize SAMPs and utilize them as receptors to
120 initiate binding and subsequent infections [28]. CD33 has also been shown to interact with
121 Hepatitis B viral surface sialoglycans, thereby impacting its pathogenesis [29]. SAMPs and their
122 interacting partners, Siglecs (primarily the V-set domains) are therefore continually evolving to
123 maintain their distinctive "self-recognition" properties, while also avoiding exploitation by Sia-
124 cloaked pathogens and parasites – a powerful example of the "Red Queen Effect" [30].
125 In this work, with a focus on CD33, a post-reproductive cognitive health associated human protein,
126 we attempt to explore the evolutionary pressures that selected for unique changes in huCD33.
127 Using human-specific pathogens like *Neisseria gonorrhoeae*, *Group B Streptococcus* and *E. coli*
128 K1, we demonstrate differential impact of these mutations on the bacterial interactions with
129 huCD33. We also determine the effect of these mutations on huCD33-sialoglycan binding and
130 identify that the amino acid at position 21 within the V-set domain plays a critical role in Sia-
131 specificity of human and chimpanzee CD33. Furthermore, we extend our study to archaic hominin
132 genomes and show that the human-specific CD33 mutations (except the presence of truncated
133 isoform) are shared evolutionary changes of human, Neanderthal and Denisovan common
134 ancestor. We also expanded our analysis to include other human-specific derived genomic changes
135 associated with cognitive health of post-reproductive human grandmothers and other elderly
136 caregivers. Finally, we draw an evolutionary scenario to connect the current knowledge of CD33
137 sialoglycan recognition and pathogen engagement to propose a role for the infectious pathogens
138 as key selective agents in human-specific CD33 evolution, generating new alleles protective
139 against infections, that could secondarily have come under selection for their protective effects
140 against cognitive pathologies like LOAD.

141

142

Results

143

144

145

Sequences of human CD33 extracellular domains show many changes distinct from closely related great apes. Previous investigations have identified unique properties of huCD33 that influence the functionality of this molecule in humans. The presence of a huCD33 V-set truncated

146 isoform as well as its overall expression difference in microglia has been associated with the
147 protection against the occurrence of neurological pathologies like LOAD in humans. Like other
148 CD33rSiglecs, CD33 immunomodulatory roles depend both on its ligand-interacting extracellular
149 domains and signaling motif-containing cytoplasmic tail. To gain a comprehensive understanding
150 of different CD33 domain variations, we compared the amino acid residues of full-length CD33
151 from human and related nonhuman primates including chimpanzee, gorilla and bonobo (Figure
152 1A). While the regions encoding the C2-set domain and cytoplasmic tail are highly conserved, the
153 amino acid residues within huCD33 V-set domain differ from their nonhuman counterparts in as
154 many as 10 positions. Since different amino acid residues in Sia-binding V-set domain could
155 potentially impact huCD33-sialoglycan interactions and subsequent downstream signaling
156 pathways, we further examined the overall frequency of these changes (Figure 1B). We analyzed
157 human sequences from the 1000 Genome database [31] and compared them with 44 gorilla, 59
158 chimpanzee and 10 bonobo sequences [32–34]. Most of these amino acid residues (except at
159 positions 66 and 148) are conserved in all the great apes and appeared to have changed only in the
160 human lineage. Interestingly, the amino acid residues at positions 66 and 148 in huCD33 are
161 isoleucine and leucine respectively, similar to CD33 of chimpanzee and bonobo. The
162 corresponding amino acids in its more distant evolutionary relative, gorilla, are phenylalanine
163 (Phe) (at position 66) and valine (at position 148). The presence of the same amino acid in human,
164 chimpanzee and bonobo at these positions suggests a more ancient occurrence of these two
165 changes, possibly prior to the divergence of chimpanzee about 6-8 million years ago (mya).
166 Previously it has been shown that the two linked SNPs, resulting in the splicing of the V-set
167 truncated isoform represent a derived evolutionary modification of the CD33 proteins in humans
168 and are absent in chimpanzees [9].

169 To further understand the selection pressure, we calculated the nonsynonymous to synonymous
170 substitution rate ratio (omega, $\omega = d(N)/d(S)$) for the CD33 V-set domains of human and other
171 great apes. The omega value of CD33 V-set domain is greater than >1 ($\omega = 1.49$) which reflects V-
172 set domain evolution under positive selection. Subsequently, we also analyzed the Ka/Ks ratios of
173 exon 2 sequences in every species. Except for gorilla, the other two great apes (chimpanzee and
174 bonobo) showed Ka/Ks ratios greater than one indicating that high Ka/Ks ratio of exon 2 is not an
175 accidental event but an evolutionary phenomenon. Taken together, these results demonstrate that

176 CD33 in humans has been rapidly evolving possibly under positive selection, distinct from its
177 orthologs in the great apes.

178 **Archaic Neanderthal and Denisovan genomes share most human CD33 protein changes,**
179 **except for the SNPs for the LOAD-protective allele.** Divergence of humans from other ancient
180 hominin lineage such as Neanderthals and Denisovans has been estimated to date back
181 approximately 0.5 mya [35]. Although full length CD33 itself is an ancient molecule, we noted
182 that the AD-protective CD33 truncated isoform is recently derived in humans, postdating our
183 divergence from Neanderthals and Denisovans [9]. Since huCD33 extracellular domains showed
184 high accumulation of changes compared to the great apes, we wanted to determine if these changes
185 were present in the common ancestor of the hominin lineage. We therefore compared CD33 protein
186 coding sequences from 6 Neanderthal and 2 Denisovan archaic genomes obtained from the Max
187 Planck Institute for Evolutionary Anthropology [36] (<http://cdna.eva.mpg.de>) with the
188 corresponding human sequences of the 1000 Genome database (Figure 1B). Interestingly, all the
189 amino acid residues in huCD33 that are different from the great apes are present in the ancient
190 genomes, suggesting their occurrence in a common ancestor. These observations thus suggest that
191 the complete loss of Sia-binding V-set domain is the latest evolutionary modification of huCD33,
192 likely succeeding the individual amino acid changes within its extracellular domain.

193 **A single amino acid change facilitated CD33 engagement to the uniquely human pathogen**
194 ***Neisseria gonorrhoeae*.** In addition to microglial expression in the brain, CD33 is also present on
195 tissue macrophages and peripheral blood monocytes [9]. These cells are important components of
196 innate immune responses throughout the body, including the reproductive tract. The human female
197 reproductive tract is also a unique niche for the microbiome, which can be invaded by important
198 pathogens like *Neisseria gonorrhoeae* (Ng). Ng is a uniquely human infectious agent, responsible
199 for the second most prevalent, sexually transmitted infection causative for the disease gonorrhea
200 in human populations. Gonorrhea affects both males and females and if untreated, can have
201 detrimental effects on reproductive health [37]. We have previously shown that Ng interacts with
202 human CD33 but not the chimpanzee ortholog [38]. The bacterium is incapable of endogenous
203 Neu5Ac synthesis, but instead scavenges the molecule from its host [39, 40]. Once inside the
204 female reproductive tract, Ng utilizes the host sugar nucleotide CMP-Neu5Ac from its
205 microenvironment to transfer Neu5Ac onto its own bacterial lipooligosaccharide. Sialylated Ng
206 then successfully interacts with several human Siglecs including 3 (CD33), 5, 9, 11, 14 and 16

207 [38]. However, unlike other Siglec interactions, Ng binding to CD33 appears to be entirely Sia-
208 dependent. Interestingly, of all the *Neisseria* species currently known, only Ng and *Neisseria*
209 *meningitidis* are pathogenic to humans and both are thought to be evolutionarily young compared
210 to others [41]. Since reproductive health/success of an organism is the key determinant of
211 Darwinian fitness, we hypothesized that highly infectious disease like gonorrhea could potentially
212 impact the evolution of humans, mediated through binding immune modulating proteins like
213 CD33.

214 To explore our hypothesis, we examined the binding of sialylated Ng to different recombinant
215 CD33 protein mutants, each containing the two extracellular domains with an amino acid residue
216 changed from human to chimpanzee at the corresponding positions identified in Figure 1B.
217 Fluorescently labelled Ng was allowed to interact with human recombinant Fc-chimeric constructs
218 of the CD33 proteins that were immobilized onto protein A-coated plates (Figure 2A). Sia-
219 dependence of the interaction was confirmed by comparing binding with bacteria grown in
220 presence and absence of CMP-Neu5Ac (Supplemental Figure S1A). We observed significant
221 reduction in bacterial binding to chimpanzee CD33 (chCD33) compared to human protein
222 containing both V- and C2- domains (Figure 2B). However, in the absence of the V-set domain in
223 the truncated form of huCD33 (CD33m), bacterial binding was lost. Except for the residue at
224 position 21, all the other amino acid alterations from human to chimpanzee CD33 maintained high
225 bacterial binding. In fact, changing the amino acid residues at positions 22, 65 (of the V-set
226 domain), 152, and 154 (of C2 domain) increased the binding significantly compared to wildtype
227 huCD33. In contrast, mutating the amino acid at position 21 from human to chimpanzee residue
228 completely abolished huCD33 binding of sialylated Ng. Interestingly, mutating the chimpanzee
229 CD33 amino acid at position 21 to its corresponding human residue enabled Ng to now engage
230 chimpanzee CD33 (Figure 2C). Considering that Ng and its closest relative meningococcus are
231 both uniquely human pathogens thought to have evolved from commensal *Neisseria* [42], our data
232 suggest important implications of CD33 amino acid change at position 21 on Ng-huCD33
233 interaction and their mutual evolution.

234 **Many amino acid changes in CD33 extracellular domains impact GBS engagement.** While
235 the association of *Neisseria* with CD33 is a case of Sia-mediated interaction, there are other
236 examples of human pathogens that engage Siglecs in Sia-independent manner. One such example
237 is Group B *Streptococcus* (GBS) which has been widely studied for its various ways of engaging

238 host Siglecs [43]. GBS is an encapsulated pathogen commonly associated with pneumonia, sepsis
239 and meningitis in infants and neonates. It comprises nine serologically distinct groups (Ia, Ib and
240 II -VIII), differing in their capsular sialoglycan structures, but all containing α 2-3-linked terminal
241 Neu5Ac. Certain GBS strains have been shown to bind human Siglecs 5 and 7 in a Sia-independent
242 manner through cell wall anchored β -protein [44], whereas some Sia-dependent binding was
243 observed for CD33 and Siglec-9. Human Siglec-9 binding is also thought to be partially Sia-
244 independent. Interestingly, some GBS strains are also known to interact with nonhuman primate
245 Siglecs, for example, Siglec-9 from chimpanzee [25]. Since infections by GBS mostly impact
246 newborns and infants, we hypothesized that it could also play a role in overall Siglec evolution in
247 humans. Similar to the Ng-CD33 binding assay (as in Figure 2A), we examined the interactions
248 between the recombinant CD33 proteins and GBS group III strain, COHI (Figure 2D). While the
249 bacteria bound strongly with full-length extracellular domains of huCD33, the binding was
250 significantly reduced in the truncated human isoform (CD33m) and the chimpanzee protein. Like
251 Ng, GBS COHI interaction was also markedly disrupted by amino acid changes at position 21.
252 Additionally, changing the residues at positions 20 and 65 from human to chimpanzee significantly
253 reduced the bacterial interaction with CD33. However, GBS COHI engagement with the CD33
254 mutants was not entirely Sia-dependent (Figure 2E). Using GBS COHI Δ neuA, a mutant strain
255 lacking its sialyltransferase enzyme (NeuA) and hence incapable of surface sialylation, we
256 observed that about 50% of the bacterial binding to CD33 could be attributed to Sia-independent
257 interactions. Interestingly, the CD33 binding profile of COHI was not uniform for the other
258 serogroups of GBS, for example GBS group Ia strain, A909 (Figure 2F). None of the amino acid
259 changes showed significant effects on CD33 interaction with GBS A909, relative to the wildtype
260 human protein. Even the truncated human CD33 isoform (CD33m) displayed similar binding
261 suggesting that the CD33 binding for A909 is primarily Sia-independent. Unlike Ng and GBS, we
262 did not observe any differential sialoglycan binding with *E. coli* K1, another uniquely human
263 pathogen of newborn infants, which contains Sia polymers on its surface (Supplemental Figure
264 S1B). Altogether, the data demonstrate the diverse nature of CD33-interactions in three major
265 pathogens and suggest an impact of uniquely human pathogens in the evolution of CD33 ligand-
266 binding domain.

267 **Ancestral sialoglycan preference of CD33 is disrupted by amino acid change at position 21.**

268 A key change in the evolution of humans was the loss of CMP-Neu5Ac hydroxylase (CMAH), the

269 enzyme that converts CMP-Neu5Ac to CMP-Neu5Gc resulting in a primarily Neu5Ac-rich
270 sialome in humans, unlike any other Old-World primates, which express both Neu5Ac and
271 Neu5Gc. This change is dated to ~2-3 mya when human ancestors were evolving from ancestral
272 hominins. Since we observed numerous changes mainly in huCD33 V-set domain which is critical
273 in sialoglycan interaction and therefore important for its downstream signaling pathways, we
274 wanted to specifically understand the effect on CD33 sialoglycan interactions. We used a
275 microarray of chemoenzymatically synthesized glycans with defined structures, terminally capped
276 with either Neu5Ac or Neu5Gc in different glycosidic linkages and examined their relative
277 interactions with recombinant, soluble CD33 proteins (Figure 3). Human CD33 with V- and C2-
278 domains bound to both Neu5Ac and Neu5Gc-terminating sialoglycans and showed maximum
279 binding when the Sia was α 2-6-linked to an underlying lactose or lactosamine glycan
280 (Supplemental Figure S2). Most of this binding was lost in the truncated huCD33 lacking the Sia-
281 binding V-set domain, indicating that the interactions are Sia-dependent. Conversely, the
282 chimpanzee protein (which is identical to the bonobo orthologs and differs by only two amino
283 acids from the gorilla) demonstrated strong preference towards Neu5Gc-terminating sialoglycans
284 and showed almost no binding for Neu5Ac-epitopes. Considering the varied sialoglycan profiles
285 of the two organisms, these distinct binding preferences of human and chimpanzee CD33 are
286 interesting and suggest functional implications of the evolutionary changes in their extracellular
287 domains.

288 Because the impact of amino acid residue at position 21 was most pronounced in both of our
289 bacterial-CD33 binding assays (Figure 2B and 2D), we next examined the influence of this change
290 on CD33-sialoglycan binding (Figure 3 and Supplemental Figure S2). Indeed, changing the amino
291 acid at position 21 completely altered Sia-epitope preference of CD33 for both human and
292 chimpanzee. The presence of human amino acid residue at position 21 enabled strong binding of
293 Neu5Ac-epitopes by chCD33, unlike its entirely Neu5Gc-preferring wildtype counterpart. On the
294 other hand, the chimpanzee amino acid at the same position in human CD33 abolished its Neu5Ac
295 binding. To determine if the Sia-binding changes are specific for position 21 and not an arbitrary
296 effect of any amino acid change in V-set domain, we also looked at the Sia-epitopes of position 20
297 amino acid substitutions. Unlike position 21, amino acid modifications at position 20 did not have
298 any major impact on the Sia-binding of CD33, which maintained the overall wildtype profile.
299 Interestingly, modifications at position 22 demonstrated Neu5Ac-preferred binding for huCD33,

300 while chimpanzee amino acid residues at 65 and 66 of huCD33 almost abolished any sialoglycan
301 binding. Altogether, the data emphasized the importance of different amino acid changes in
302 huCD33 V-set domain for its sialoglycan binding and identified the amino acid at position 21 to
303 be critical in the functionality of CD33 protein.

304 **MD simulations provide structural insights for the differences in Sia-binding preference**

305 **between human and chimpanzee CD33.** We performed an extensive theoretical investigation

306 based on molecular dynamics simulations. A detailed analysis of several available crystal

307 structures of huCD33 revealed that the V-set domain is dynamic. For example, the C-C' loop as

308 well as the side chains of phenylalanine at position 21 (Phe21) and histidine at position 45 (His45)

309 are resolved in two different conformations in PDB entry 5ihb (Supplemental Figure S3). Of all

310 the amino acids that differ between human and chimpanzee, only the side chain of Phe21 is in

311 direct contact with a bound Neu5Ac residue in the crystal structures of huCD33 (through the

312 methyl group at position 5). Based on the assumption that Neu5Gc binds to the same binding site

313 as Neu5Ac, the change in binding preference from Neu5Gc (in chimpanzee) to Neu5Ac (in human)

314 cannot be explained by a simple I21F mutation. Both amino acids have hydrophobic side chains

315 that cannot establish favorable interactions with the polar glycolyl group. Consequently, there is

316 probably a more complex reason for the shift of binding preference. Based on data derived from

317 47 molecular dynamics (MD) simulations covering an accumulated timescale of more than 100 μ s

318 we conclude that in chCD33 His45 adopts mainly the 'up' conformation (Figure 4A), which allows

319 favorable hydrogen bonding with the glycolyl group. MD simulations (as well as x-ray

320 crystallography) show that in huCD33 His45 can also exist in the 'up' conformation (Figures 4C

321 and Supplemental Figure S4), which would be compatible with favorable Neu5Gc binding.

322 However, when His45 is in the 'down' conformation Phe21 can stack partly with tyrosine (Tyr) at

323 position 127 (Figure 4B) forming a small hydrophobic pocket, which allows the methyl group of

324 Neu5Ac to bind favorably. To demonstrate if the binding affinity difference between Neu5Ac and

325 Neu5Gc may be indeed correlated to the up/down conformational equilibrium of His45, we

326 performed a series of MD simulations of chCD33 on the microsecond timescale where

327 Neu5AcOMe or Neu5GcOMe molecules are present in the solution. The lifetimes of the

328 complexes spontaneously formed during the MD with Neu5Gc are on average much longer when

329 His45 is 'up' (Figure 4D top, Supplemental Figure S4). In contrast the lifetimes of the complexes

330 spontaneously formed with Neu5Ac are much shorter independent of the conformational state of

331 His45 (Figure 4D bottom), which would explain the lack of measurable binding affinity of Neu5Ac
332 to chCD33. In summary, our extensive MD simulations - including unbiased simulation of
333 carbohydrate binding and unbinding events - could provide a reasonable explanation for a change
334 in binding specificity that is likely to be caused by an alteration of the protein-ligand interaction
335 pattern remote from the mutated amino acid.

336 **Human-specific polymorphisms in cognitive-health related genomic variants are present in**
337 **all human populations.** In an earlier study we observed several genes, directly associated with
338 neurodegenerative diseases or correlated with aggravation of the cognitive decline in aged-
339 population, are derived alleles in humans [9]. Increasing evidence of correlation between cognitive
340 health and non-neurological, metabolic conditions, e.g., diabetes [45, 46] suggest that such derived
341 alleles could be important in the maintenance of cognitive health in human grandparents. Here, we
342 expanded this list of cognition-protective gene variants through literature and database
343 (<https://alzoforum.org>) searches [20, 47–58] to include additional gene variants, namely *BINI*,
344 *ARID5B*, *PICALM*, *PILRA*. Supplemental Table S1 describes the characteristics of 13 human
345 genes that are implicated in diseases including dementia, cardiovascular diseases (CVD),
346 hypertension and AD. While some of these physiological abnormalities like salt retention,
347 hypertension, diabetes, appear non-neurological, they have been associated with the aggravation
348 of the pathologies resulting in late-life cognitive decline [59]. Notably, the derived alleles are
349 common and found in globally diverse human populations, indicating that they predate the
350 common ancestor of modern humans (Supplemental Table S2).

351 **SNPs associated with human-specific cognitive protective alleles are unusual in their absence**
352 **in the archaic hominin genomes.** With the availability of genomes from extinct archaic hominins
353 [36, 60, 61], a set of SNPs can be assessed as to whether their protective phenotypes arose recently
354 in the evolutionary history of anatomically modern humans. We previously showed many other
355 human-chimpanzee differences were shared with archaic hominins (Denisovan/Neanderthal)
356 genomes; for example, genomic changes in CD33rSiglecs [62]. To gain similar insights about the
357 evolutionary origin of these cognitive-protective loci, we analyzed the Neanderthal and Denisovan
358 reference genomes and compared them with modern human sequences. Analysis of the 1000
359 Genomes dataset shows the presence of protective alleles in human populations with variable
360 frequency (Table 1). Analysis of the available genomic data from Neanderthal and Denisovan
361 genomes showed that only two derived variants (rs2975760 and rs2588969; Table 1) are present

362 in these archaic genomes, suggesting the remaining eleven derived, protective variants arose after
363 the divergence of modern and archaic hominins approximately 0.5 mya [35, 63]. This is in striking
364 contrast to most human-chimpanzee genomic differences in which the archaic hominins are similar
365 to humans. In fact, majority of the Sia-related genes lack positive selection signatures and rather
366 show neutral evolution in the modern human lineage [64]. To more formally assess whether the
367 high frequency, global distribution, and recent origin observed for eleven of the thirteen SNPs is
368 unusual, we performed a resampling analysis of variants in the genome. Variants in the 1000
369 genomes dataset that met the following criteria were considered: 1) present in both Altai
370 Neanderthal and Denisovan minimal filters, 2) derived in at least one modern individual from non-
371 admixed African populations, 3) called in both archaic samples, and 4) have an ancestral allele
372 matching the reference or alternative allele. To eliminate any bias in the analysis and match the
373 allele frequency (AF) of these SNPs compared with that of any random SNPs, we first matched
374 our universe of SNPs to the 13 SNPs of interest by AF, ± 2 derived haplotypes (Figure 5).
375 Resampling was then performed by drawing a SNP from each of the 13 matched sets and assessing
376 how many derived alleles were observed, resulting in a p -value = 0.08333 ± 0.00003 . As a less
377 conservative estimate, directly sampling from the universe of SNPs and estimating the probability
378 of observing at most two derived SNPs and a mean allele frequency as large as the empirical
379 variants of interest produced a highly significant p -value = 0.00487 (Supplemental Figure S5).
380 Repeating either analysis on the set of other Siglec-related SNPs indicates they are consistent with
381 a random draw from the genome [62]. Regardless of the individual limitations, taken together our
382 phylogenetic analyses demonstrate the unique patterns of allele frequencies in worldwide
383 populations distribution of these thirteen late-life cognitive decline linked SNPs (Figure 5 and
384 Supplemental Figure S5). Interestingly, co-inherited CD33 SNPs associated with the cognitive
385 health in LOAD are present only in modern human genomes [9]. A noteworthy example in our list
386 is the human gene encoding the protein, apolipoprotein E (APOE), involved in fat metabolism in
387 mammals. *APOE* gene exists in three allelic variants (E2, E3 and E4) where APOE4 is associated
388 with high risk of LOAD and other allele like APOE2 is protective against the cognitive decline in
389 elderly caregivers [65]. Interestingly the presence of APOE4 is also correlated with the protection
390 from severe diarrhea in children [66]. While conclusive determination of the positive selection of
391 these alleles in modern human requires further analysis, our data suggest that the evolutionary
392 origin of most of these cognitive-health protective changes followed the divergence of modern

393 humans from archaic genomes. This is also supported by the presence of grandparents, uniquely
394 in humans. Regardless, the process of evolutionary emergence of each of these alleles is likely to
395 be distinct and deserves further investigation.

396

397

Discussion

398 Fossil evidence and genomic comparisons leave little doubt about the fact that our species evolved
399 from an African hominin. However, a detailed understanding of modern human origins is plagued
400 by numerous uncertainties, with regard to the identity of the ancestral lineage and precise
401 geographic locations. Evolution of modern humans was accompanied by many anatomical and
402 behavioral changes, but increasing evidence suggests it also included uniquely human-derived
403 modifications in the genome compared to the archaic genomes (Neanderthal/Denisovan) or the
404 genome of the great apes [9, 62]. Taken together with our previous study [9], we have identified
405 many such human-specific genes associated with cognitive health of grandmothers and other
406 human elders who are often involved in the caregiving of the young. These findings, which appear
407 paradoxical to the concept of senescence due to antagonistic pleiotropy, have lent much additional
408 support to the “Grandmother hypothesis” [1] bolstering the case for selection of human female
409 post-reproductive survival and the existence of grandmothers. Unlike in any other mammals
410 (except orcas and some other toothed whales), the occurrence of this prolonged post-reproductive
411 life span in humans has stirred scientific interest. While deciphering the precise evolutionary
412 course of any gene/protein is challenging and the proposed schemes/players are not entirely
413 verifiable, here we attempt to compile the current evolutionary and experimental findings of one
414 such protein associated with late-life cognitive-decline: CD33.

415 A ratio of high wildtype human CD33 and a low truncated isoform of CD33 have been implicated
416 in the progression of LOAD associated with the cognitive health of elderly population. In contrast,
417 LOAD is unknown in chimpanzees, although evidence of LOAD pathologies has been observed
418 in some chimpanzee brains. We found that human CD33, which is highly expressed in microglia
419 of the human but not chimpanzee brain, recognizes Neu5Ac – the predominant Sia synthesized in
420 humans – as self-associated molecular patterns (SAMPs). In contrast, our closest evolutionary
421 relative, the apes and other Old-World primates contain both Neu5Ac and Neu5Gc. We found that
422 the ancestral form of CD33 in chimpanzees and other great apes selectively recognizes only
423 Neu5Gc-glycans as SAMPs (Figure 3). Notably, Neu5Gc – the ligand recognized by chCD33 is

424 rare in chimpanzee brain, and there is also significantly less chCD33 protein compared to CD33
425 in humans [9]. On the other hand, SNPs resulting in the truncated CD33 have only been observed
426 in the human genome and not any of the archaic or great ape genomes. We also find that the
427 truncated human CD33 does not interact with Sia (Figure 3). Taken together, these observations
428 suggest that full-length CD33-Sia interactions are stronger in human brain compared to
429 chimpanzee and the human-specific SNPs in CD33 resulting in the truncated protein abolish this
430 interaction. The question remains what could have possibly led to the selection of the truncated
431 isoform of human CD33 that does not interact with Sia. In this regard, CD33 on macrophages
432 plays crucial roles in different immune responses as well as during infections. Human CD33 has
433 also recently been shown to be involved in immunomodulation during infection with hepatitis B
434 virus [29]. Our previous and current data show that uniquely human pathogens like *Neisseria* and
435 GBS display Neu5Ac that is recognized as ‘self’ by human but not chimp CD33 [38]. In the current
436 work, we further found that the Sia-binding-domain-depleted, truncated human CD33 isoform
437 doesn’t bind and thus escape exploitation by sialylated pathogens (Figure 2). This suggests that
438 this truncated CD33 may have been an adaptation to counter the CD33-exploiting, immune-
439 evasive behavior of pathogens like *Neisseria* and GBS.

440 Taking together all currently available experimental data (including this study) we attempt to draw
441 a plausible evolutionary scenario for CD33 protein evolution in humans and present in the context
442 of relevant evolutionary events (Figure 6). We hypothesize that the scarcity of the strongly
443 preferred Neu5Gc ligand of ancestral CD33 in the brains of chimpanzee (and other great apes) was
444 associated with low microglial expression. Subsequent hominin loss of CMAH (i.e., complete loss
445 of Neu5Gc ligand) could then have selected for the upregulation of CD33 levels perhaps to
446 compensate for the loss of ligands, a change to Neu5Ac-binding preference, and functional
447 recruitment of CD33 to human microglia. Alongside the microglial CD33, the corresponding
448 changes in the tissue macrophage proteins might have facilitated the emergence of Neu5Ac-coated
449 pathogens (for example, *N. gonorrhoeae* and Group B *Streptococcus*) that evolved “molecular
450 mimicry” of Neu5Ac-SAMP ligands to manipulate the immune response. Appearance of the
451 truncated isoform lacking the ligand-binding domain (CD33m), then probably allowed CD33 to
452 escape the immune evasion by these sialylated pathogens (Figure 2). This selection pressure to
453 stop manipulation by sialylated pathogens could have also altered splicing towards a higher level
454 of truncated CD33, which also gets diverted to peroxisomes [12]. While the significance of this

455 diversion is unclear, decrease of full-length CD33 would facilitate escape from Neu5Ac-coated,
456 CD33-engaging pathogens. Finally, sometime during the last 1 million years, increased brain size
457 presumably selected for early, short interbirth interval in human, which might have resulted in
458 more helpless young, requiring cooperative breeding and caregiving. However, the value of
459 postmenopausal grandmothers and other elderly caregivers would then have been blunted by the
460 appearance of LOAD. The synthesis of the truncated isoform of CD33 protects from *Neisseria*
461 during reproductive age and a higher ratio of truncated to full-length isoforms correlates to
462 decrease of LOAD in grandmothers. However, a small amount of the full-length isoform remains,
463 likely to downregulate hyper-inflammation that might arise during prolonged absence of SAMP-
464 recognition. Notably when an elderly caregiver gets LOAD, not only are the evolutionary benefits
465 of the individual lost, but this also presents an increased burden to care for that elder individual.
466 Altogether under this proposed scenario, the current state in the evolution of human CD33 protein
467 represents a trade-off between the evolutionary response to exploitation by pathogens in early life
468 and cognitive maintenance in post-reproductive late life.

469 A similar evolutionary scenario appears to underly the case of the human *APOE* gene
470 where variants include both risk alleles (*APOE4*) and protective alleles, (*APOE2*, and *APOE3*) for
471 CVD and LOAD [65]. In this instance, the ancestral *APOE4* allele is associated with increased
472 risks of loss of cognitive functions and the derived alleles may serve to protect the cognition of the
473 elderly caregivers. Interestingly the *APOE4* allele is also correlated with the protection from severe
474 diarrhea in early years of life [66]. Given these examples like APOE and CD33, it remains to be
475 seen how widespread this evolutionary pattern is wherein variants conveying survival advantages
476 in early life coexist with other variants that protect cognition late in life.

477

478 **Acknowledgments**

479 We are very thankful to Dr Sanjay Ram (University of Massachusetts, Worcester) and Dr. Victor
480 Nizet (University of California, San Diego) for the bacterial strains that have greatly benefited the
481 project. This work was supported by NIH grant R01GM32373 and Cure Alzheimer's Fund grant
482 (to A.V.).

483 **Author Contribution**

484 Experimentation (S.S., N.K., T.C., A.V., A.S., S.D., M.F.), Data analysis (S.S., N.K., T.C.,
485 A.V., A.S., S.D., J.M.A., M. F., P. G., A.V.), Critical reagents (H.Y., X.C.), Original draft (S.S.,

486 P.G., A.V.), Writing (S.S., N.K., T.C., A.V., A.S., X.C., J.M.A., M.F., P.G., A.V.), Overall
487 supervision (J.M.A., P.G., A.V.), Funding acquisition (A.V.).

488 **Data Availability**

489 The data for the resampling analysis is available at Code Ocean.

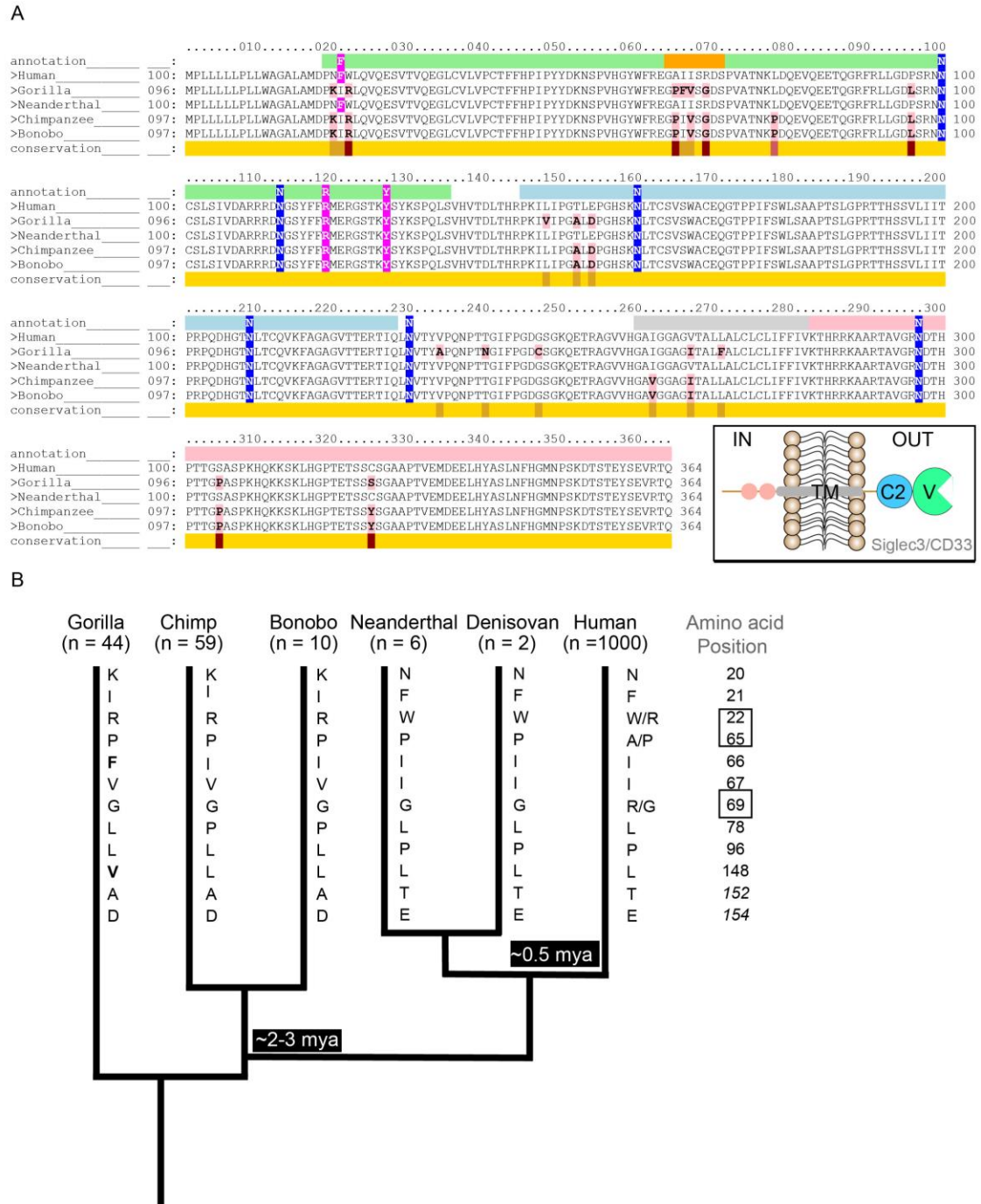
490

491 **Declaration of interest**

492 The authors have declared that no conflict of interest exists.

493

Figures



494

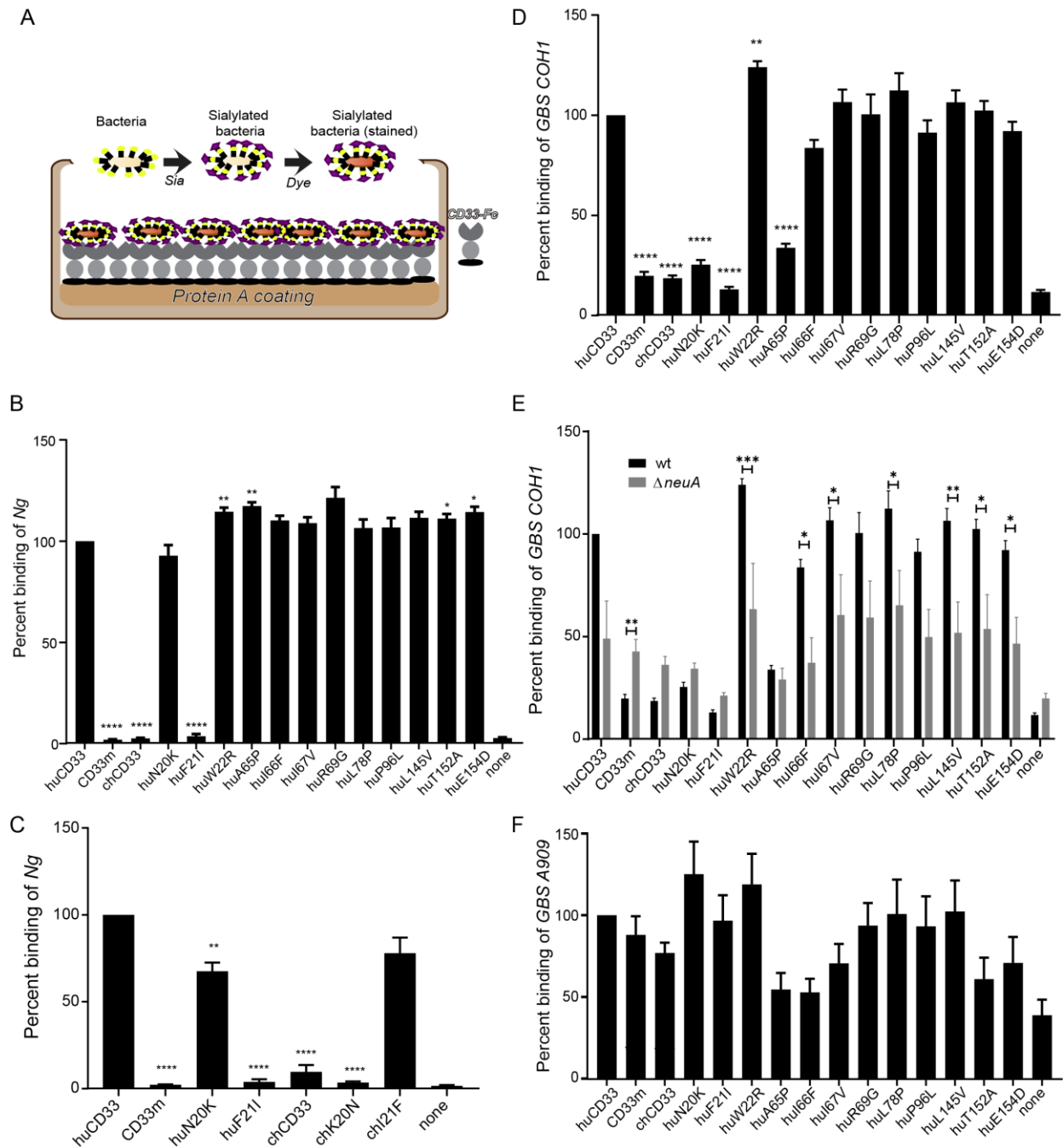
495 **Figure 1: Human specific changes in CD33 are primarily present in the Sia binding V-set**

496 **domain.** (A) Comparison of amino acid sequences of CD33 from humans and “great apes” was

497 performed using Conformational Analysis Tools software. The Great ape genomes included in the

498 analysis are gorilla, chimpanzee and bonobo and were compared against the human protein as the

499 template. The conservation of the sequence is indicated with yellow being the most and red being
500 the least conserved regions. Amino acids that are different from huCD33 are highlighted in pink.
501 Amino acids encoding the different CD33 domains are indicated above the sequence with different
502 colors corresponding to schematic in inset, namely, V-set domain in green, C2 domain in light
503 blue, transmembrane domain in grey and cytoplasmic end in light pink. The flexible C-C' loop is
504 indicated in orange. Amino acids that are in contact with Neu5Ac in huCD33 are highlighted in
505 magenta and N-glycosylation sites in blue. **(B)** Phylogenetic analysis of the evolution of the
506 extracellular domains of CD33 proteins from human, Great apes and archaic genomes. The number
507 of genomes (n) for each group included in the analysis is indicated. Human and great ape CD33
508 sequences were compared with six Neanderthal and two Denisovan genomes. Amino acid changes
509 present in human CD33 were also present in the ancient genomes. The positions of the amino acids
510 that are different between human and the apes are mentioned, and the identity of the amino acid
511 present in the corresponding positions for each group is indicated by the single letter abbreviations
512 along the branch. Amino acids at positions 152 and 154 are within the C2 domain of CD33 protein
513 and italicized. Polymorphisms within the human population at positions 22, 65 and 69 of CD33
514 protein are indicated. Amino acids in gorilla CD33 at positions 66 and 148 are different from other
515 apes and are bold. Possible timeline for the diversion of the hominin lineage is indicated in the
516 tree. Length of the branches in the tree is not to scale. Mya = million years ago.



517

518 **Figure 2: Human specific amino acid changes in CD33 affects bacterial binding. (A)**

519 Schematic of the ELISA-based assay using recombinant CD33-Fc chimeric proteins immobilized

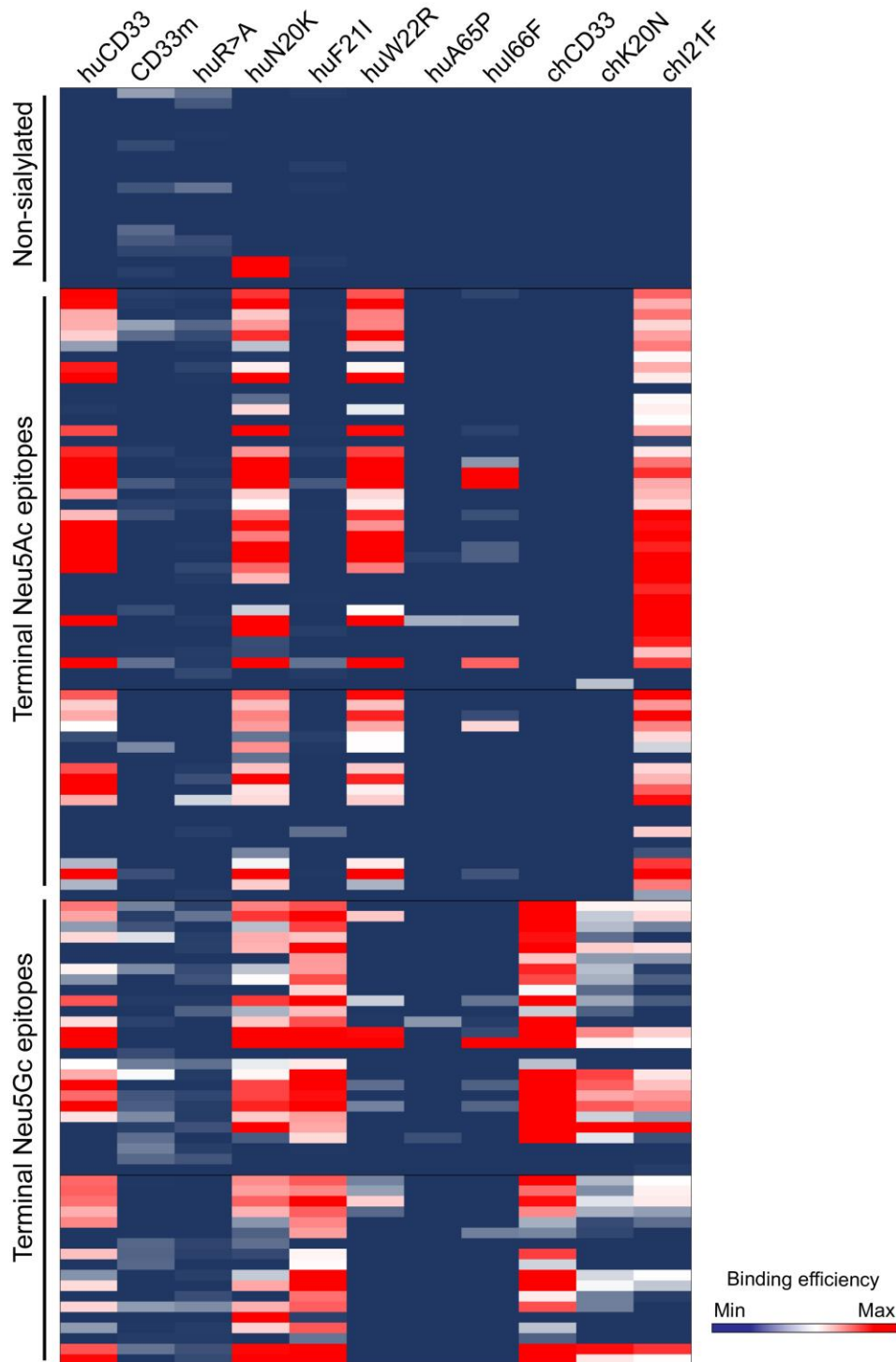
520 on protein A coated plates used to determine binding of the sialylated bacteria is shown. (B)

521 Binding of fluorescently labelled *Neisseria gonorrhoeae* (*Ng*) was determined. The position of the

522 amino acid different from the wildtype human CD33 protein is indicated below each bar in the x-

523 axis. The bacterial binding to each individual CD33 mutant was normalized to the binding of

524 wildtype human CD33 for that assay. “None” indicates no protein control for the background
525 bacterial binding to the plate. **(B)** Binding of *Neisseria* to immobilized recombinant CD33 proteins
526 containing the corresponding amino acid mutation (position 20 or 21) in either in human or chimp
527 CD33 protein backbone. **(D)** Binding of Group B *Streptococcus* (GBS) COH1 strain to different
528 CD33 mutant proteins in an ELISA based assay with immobilized recombinant CD33 proteins.
529 **(E)** Sialic acid dependence of the binding was determined using wildtype and $\Delta neuA$ mutant strains
530 of COH1. **(F)** Interaction of CD33 proteins among different GBS strains was compared using A909
531 and COH1 strains. ‘hu’ indicates the corresponding amino acid change in human CD33 backbone
532 and ‘ch’ using chimp CD33. The graphs show the cumulative result from 3 independent
533 experiments, each done in triplicate. Statistical analysis was performed in Prism software using
534 one-way ANOVA with Durrett post comparison test. * < 0.01, ** < 0.001, *** <0.0001.

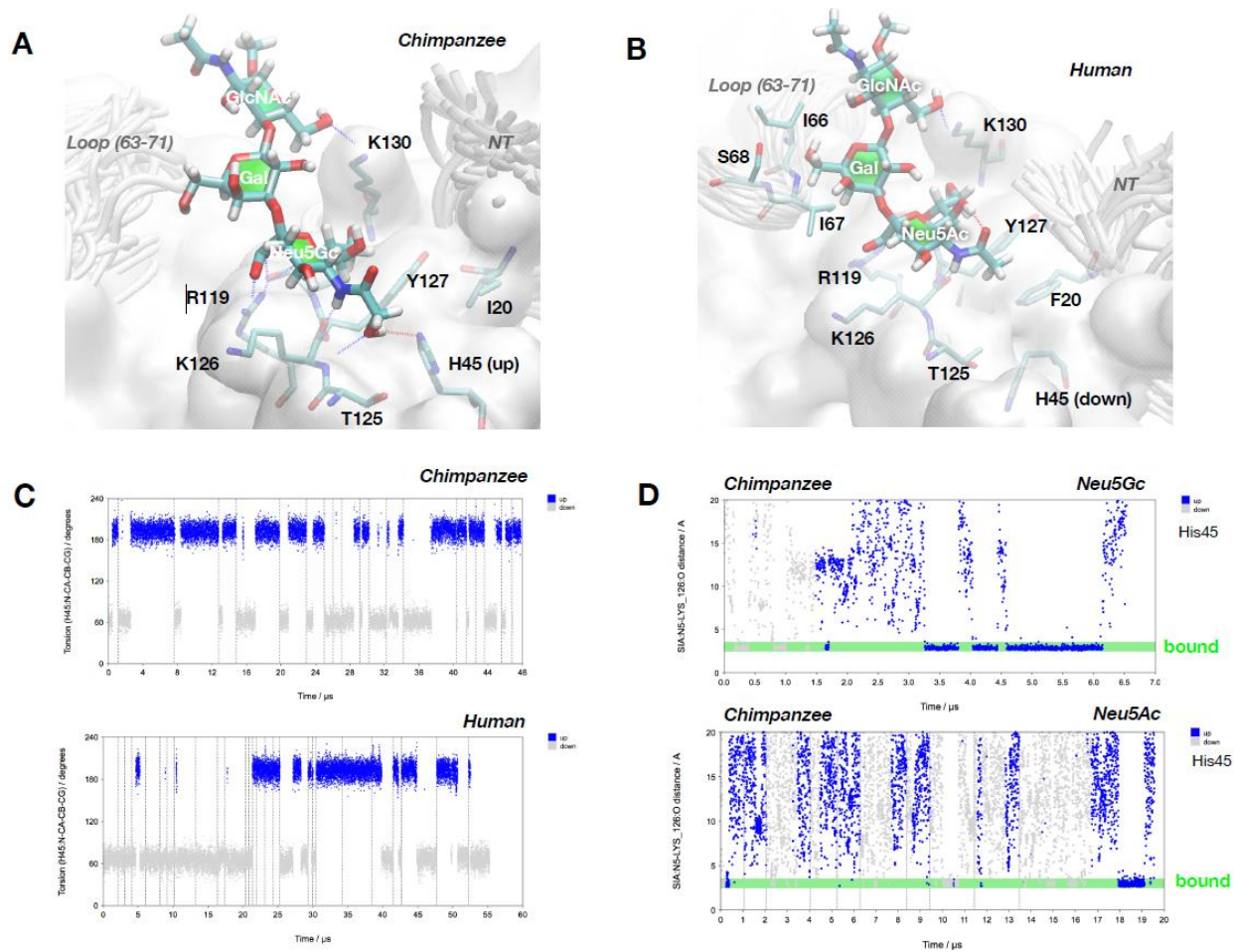


535

536

537 **Figure 3: Single amino acid changes affect CD33 sialoglycan binding.** Sialoglycan binding
538 profile of purified, soluble, recombinant CD33 proteins was determined using a sialoglycan
539 microarray containing defined, chemically synthesized glycans. Non-sialylated, Neu5Ac- and

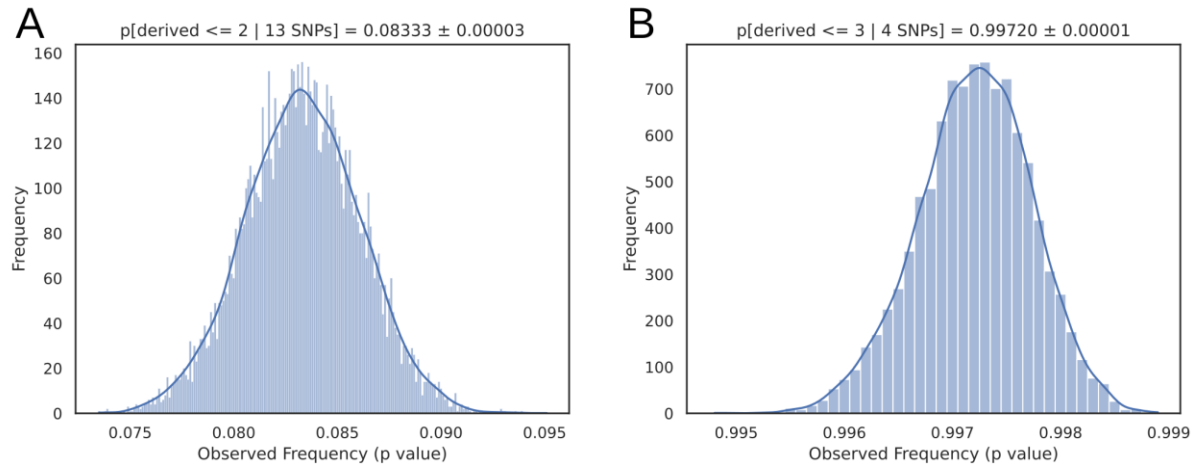
540 Neu5Gc-terminating glycans were grouped together in the heatmap as shown in the left. Each
541 column indicates the binding profile of the protein indicated on the top and each row represent a
542 distinct glycan. Blue indicates no binding and red indicates very strong binding preferences
543 characterized by an average relative fluorescence unit (RFU) of more than 90th percentile. The
544 result of the heatmap is summarized in Supplemental Figure S2 and the names of the individual
545 glycans are presented in Supplemental File S1.



546

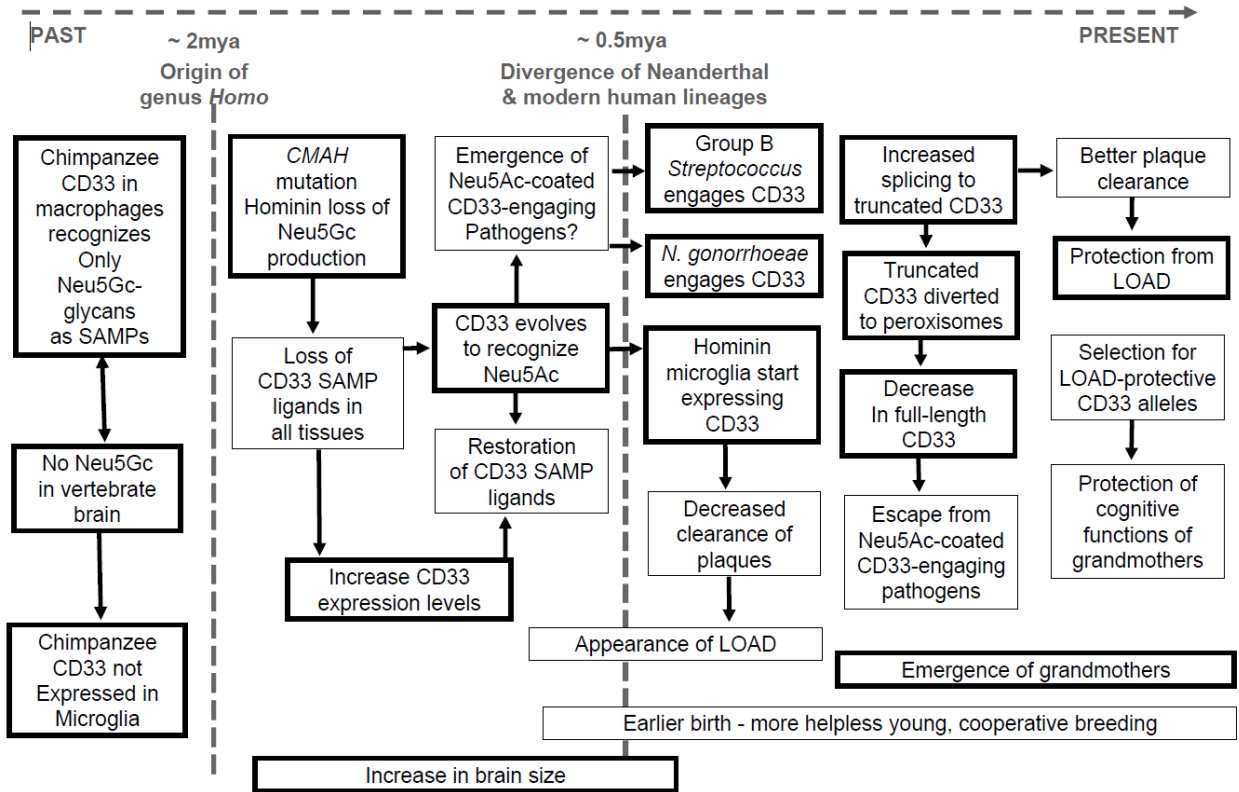
547 **Figure 4: Structural modeling to understand the differential binding preference of human**
 548 **and chimpanzee CD33 proteins.** (A) 3D model of the complex between Neu5Gc α 2-3Gal β 1-
 549 4GlcNAc β OME and chCD33. The increased affinity of Neu5Gc may be explained by
 550 intermolecular hydrogen bonds involving the OH-group of Gc. It should be noted that the number
 551 of favorable interactions is maximal when His45 is in ‘up’ conformation. (B) 3D model of the
 552 complex between Neu5Ac α 2-3Gal β 1-4GlcNAc β OME and huCD33. The methyl group of Ac is
 553 located in a small hydrophobic pocket formed by the side chains of Tyr127 and Phe20. It should
 554 be noted that His45 is in ‘down’ conformation because otherwise – in the conformation shown -
 555 the bulky side chain of Phe20 would overlap partly with His45 in ‘up’ conformation. (C) Molecular
 556 dynamics of His45 side chain orientation. Accumulated MD trajectories of torsion angle N-C α -
 557 C β -C γ are shown. The ‘up’ conformation is present when torsion values are fluctuating around
 558 200 degrees and the ‘down’ conformation is characterized by values around 70 degrees. For
 559 chimpanzee, it can be observed reproducibly that simulations started with His45 in ‘down’

560 conformation undergo a transition to the ‘up’ conformation on the microsecond timescale. In
561 contrast the ‘down’ conformation appears to be more stable in huCD33, which would make
562 Neu5Ac binding more likely. **(D)** MD simulation of unbiased binding and unbinding events of
563 Neu5Gc (top) and Neu5Ac (bottom) to chCD33. For Neu5Gc the lifetime of the complex is
564 significantly longer when His45 is in ‘up’ conformation, as can be seen from the 6.5 μ s MD
565 simulation shown on the top. Also, for Neu5Ac multiple binding and unbinding events occurred
566 on a timescale of about 20 μ s, however in general (with one exception) the lifetimes of the
567 complexes formed are significantly shorter than for Neu5Gc.



568

569 **Figure 5: Resampling analysis with matched allele frequency SNPs from 1000 genomes**
570 **variants.** Frequency distribution of SNPs with similar properties to the LOAD protective set (**A**)
571 and other Siglec SNPs (**B**).



572

573 **Figure 6: Scenario for evolution of human CD33 in relationship to cell surface sialic acids,**
 574 **infectious disease, brain microglia and cognitive maintenance of grandmothers and other**
 575 **elderly caregivers.** This schematic presentation combines the known/likely facts (thick-outlined
 576 boxes, including data from this manuscript) as well as suggested possibilities (thin-outlined boxes)
 577 into the most likely evolutionary scenario for human-specific evolution of CD33. Starting from
 578 the left, the likely chronological order of occurrence is indicated (by arrowheads) with the
 579 approximate timeline on the top, along the dotted lines. ‘?’ indicates our reasonable assumption
 580 leading to the event. See text for further discussion.

581
582

Table

583 **Table 1:** Gene variants directly or indirectly affecting cognitive function. Allele has the derived
584 allele as the lower, bolded entry. Archaic genotypes are reported for SNPs passing all quality
585 filters.

Gene	Associated Disease ^a	SNP ID	hg19 Position	Allele	Derived Global AF ^b	Archaic Genotype	
						Altai	Denisovan
<i>AGT</i>	Sodium retention	rs699	1:230845794	G	0.295	0/0	0/0
				A			
<i>BIN1</i>	AD	rs7561528	2:127889637	A	0.8	0/0	0/0
				G			
<i>SGC2</i>	Hypertension	rs1017448	2:224466344	T	0.879	0/0	0/0
				C			
<i>CAPN10</i>	Type II Diabetes	rs2975760	2:241531163	C	0.882	1/1	1/1
				T			
<i>PPARG</i>	Type II Diabetes	rs1801282	3:12393125	C	0.070	0/0	0/0
				G			
<i>CYP3A5</i>	Sodium retention	rs776746	7:99270539	T	0.621	0/0	0/0
				C			
<i>ARID5B</i>	AD	rs2588969	10:63611354	A	0.532	0/0	1/1
				C			
<i>SPON1</i>	Dementia	rs2618516	11:14021639	C	0.341	0/0	0/0
				T			
<i>PICALM</i>	AD	rs10792832	11:85867875	G	0.313	0/0	0/0
				A			
		rs3851179	11:85868640	C	0.315	0/0	0/0
				T			
<i>APOE</i>	LOAD	rs429358	19:45411941	C	0.849	0/0	0/0
				T			
		rs7412	19:45412079	C	0.075	0/0	0/0
				T			
<i>CD33</i>	LOAD	rs3865444	19:51727962	C	0.211	0/0	0/0
				A			
<i>PILRA</i>	AD	rs1859788	7:99971834	G	0.341	-	-
				A			
<i>TCFLC2</i>	Type II Diabetes	rs7903146	10:114758349	T	0.772	-	-
				C			
<i>CD33</i>	LOAD	rs12459419	19:51728477	C	0.211	-	-
				T			

586 ^a See supplemental Table S1 for details and the primary literature citations.

587 ^b See supplemental Table S2 for the global population distribution.

588

Methods

589 **Bacterial culture and cell lines.** The bacterial strains used were *Neisseria gonorrhoeae* F62Δ*lgtD*
590 (generous gift from Sanjay Ram, University of Massachusetts Worcester), Group B *Streptococcus*
591 (GBS) strains COH1wt, COH1Δ*neuA*, A909wt and A909Δ*neuA* (generous gifts from Victor Nizet,
592 University of California San Diego). *Neisseria* were grown overnight on chocolate II agar plate
593 and GBS on Todd Hewitt agar plate at 37 °C and 5% CO₂ from the respective frozen glycerol
594 stocks. Prior to the assay, GBS was grown in Todd Hewitt broth at 37 °C and 5 % CO₂ without
595 shaking. The *E. coli* K1 strain was grown in LB. For the CD33 protein purification, HEK293A
596 cells were grown in DMEM media (Invitrogen) containing 10% FCS at 37 °C and 5 % CO₂.

597 **Sialylation of *Neisseria*.** Following overnight growth on chocolate agar plate, the bacteria were
598 grown in GC broth supplemented with IsoVitaleX at 37 °C, 5% CO₂ and shaking at 200 rpm in
599 presence or absence of 30 μM CMP-Neu5Ac (Nacalai USA, Inc.) until OD600 equivalent to 0.4–
600 0.5.

601 **Bacterial staining.** Following appropriate growth, the bacteria were washed with pre-warmed
602 HBSS and stained with 2 μM SYTO13 (Thermo Scientific) for 30 min at 37 °C and shaking at 200
603 rpm in dark. After incubation, the stained bacteria were washed with HBSS and resuspended to a
604 final concentration of OD600 = 1/ml in HBSS for the binding assay.

605 **Generation of CD33 mutant proteins.** A genomic fragment (1228 bp) of human or Chimpanzee
606 CD33(M), including the first 4 exons (2 Ig domains) was fused with pcDNA3.1(-) containing a C-
607 terminal FLAG (EK) sequence followed by a hIgG1-Fc genomic fragment (hinge + 2 Ig-like
608 domains) and described elsewhere [67, 68]. Sixteen mutant variants were made from either
609 construct above using New England Biolabs Q5 site directed mutagenesis Kit according to the
610 manufacturer's instructions (Supplemental Table 3). Mutagenesis primers listed were designed
611 using NEBaseChanger software.

612 **Truncated CD33(CD33m) _EK_Fc Construction:** U937 cells were cultured in RPMI 1640
613 supplemented with 10% FCS. Total mRNA was isolated using Qiagens Oligotex Direct mRNA
614 Mini Kit according to the manufacturer's instructions. CD33m was amplified by PCR using
615 SuperScript III One-Step RT-PCR (Invitrogen) and Gene-specific primers 5'-
616 TTATATGCTAGCGCCACCATGCCGCTGCTGCTACTGCTGC-3', NheI site underlined and
617 5'-GCGCGCGATATCATGAACCACTCCTGCTCTGGTCTCTTG-3', EcoRV site underlined.
618 PCR products were run on 2% agarose gel and the 396 bp bands corresponding to CD33(m) were

619 excised and cut with NheI/EcoRV restriction enzymes. Digested bands were sub-cloned into
620 pcDNA3.1(-) containing a C-terminal FLAG (EK) sequence followed by a hIgG1-Fc genomic
621 fragment (hinge + 2 Ig-like domains).

622 **Purification of CD33 mutants.** Transfection supernatants were collected and spun down at 500
623 g for 5 mins to remove cellular debris. The pH of each supernatant was adjusted to pH 8.0 for
624 optimal binding of protein A-Sepharose beads to hIgG Fc fusion protein. Protein A-Sepharose 4
625 Fast Flow suspension (GE Healthcare) was washed with Tris-Buffered Saline (TBS) pH 8.0, and
626 a 1:500 ratio of beads:media added to each supernatant. Each tube was subsequently incubated for
627 24 hrs on a roller in the cold-room. After 24 hours supernatants plus beads were transferred to
628 disposable columns until all liquid has run thru. Beads were washed 3x with TBS pH 8.0 before
629 being eluted directly in 0.3 ml of 1 M Tris-HCl pH 8.0 using 0.1 M Glycine Buffer pH 2.8. Each
630 eluate was put into an Amicon Ultra-15 filter unit with MWCO 30 K for each full length CD33-
631 EK_Fc variant and MWCO 10 K for huCD33m-EK_Fc. Tubes were centrifuged at 4,000 g for 20
632 mins. Run-through was discarded and the columns washed 3x with TBS pH 8.0. After the last
633 wash, each retentate was recovered from the column and stored at -80°C.

634 **Binding assay with the bacteria.** Bacterial binding with the CD33 proteins were done with the
635 recombinant Fc-chimeric proteins of CD33. Briefly, protein A coated black 96-well plate (Pierce,
636 Thermo Scientific) was washed thrice with TBS containing 0.05% Tween 20 (TBS-T) and coated
637 with 200 ng/well of the respective CD33 protein diluted in 200 mM Tris-HCl pH 8.0, 150 mM
638 NaCl and 1% BSA at 4 °C overnight. Following incubation, the coated plate was washed with 200
639 mM Tris pH 8.0, 150 mM NaCl to eliminate the unbound proteins. Stained bacteria equivalent to
640 OD600 = 0.1 was added to each well of the plate and allowed to interact with the proteins for 30
641 min at 37 °C and 5% CO₂ without shaking. Following incubation, the plate was washed with TBS-
642 T to eliminate any unbound bacteria and the residual fluorescence was measured upon excitation
643 at 488 nm and emission at 530 nm. The data were analyzed using the excel and Prism software.

644 **Evolutionary analysis and Detection of positive selection.** The protein coding sequences of
645 CD33 were aligned using CLUSTAL W program implemented in MEGA7 and then back
646 translated to obtain a codon alignment. The phylogenetic tree of CD33 protein coding sequences
647 were reconstructed with neighbor-joining method which was implemented in MEGA7 (Figure 1),
648 1000 bootstrap replicates [69]. The unrooted neighbor joining tree was used for the subsequent
649 analysis.

650 VCF files were accessed from <ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/> for
651 1000 genomes project, <http://cdna.eva.mpg.de/neandertal/altai/AltaiNeandertal/VCF/> for Altai
652 Neanderthal and http://cdna.eva.mpg.de/denisova/VCF/hg19_1000g/ for Denisovan. Quality
653 filters were obtained from https://bioinf.eva.mpg.de/altai_minimal_filters/ for Altai and
654 Denisovan. Individuals in 1000 genomes datasets were assigned to populations using
655 http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/working/20130606_sample_info/20130606_s
656 [ample_info.txt](#). First, all vcf files were filtered by intersecting with the quality bed files using
657 bedtools intersect (v2.28.0). The filtered vcf files were then combined, per chromosome, to match
658 their position, reference and alternative allele using a custom python script. VCF information was
659 retained along with per-population allele frequencies and archaic genotypes. Next, ancestral alleles
660 were obtained from ensembl (https://rest.ensembl.org/variation/homo_sapiens) by querying each
661 SNP id and appending to the joint vcf entries. The joint vcf files were used as input for further
662 processing in a jupyter notebook to perform resampling analysis. Each SNP of interest is used to
663 select collections of SNPs with matching global allele frequencies, +/- 0.004, or +/- 2 observed
664 haplotypes. A single draw consists of selecting one SNP from each collection to produce a
665 simulated observation and the number of SNPS with derived archaic haplotypes are recorded.
666 After 10,000 such draws, the fraction of draws with fewer or equal numbers of derived SNPs is
667 used to produce a p-value estimate. The process is repeated 10,000 times to produce a histogram
668 and provide a confidence estimate on the reported p values (+/- SEM). Methods to replicate the
669 analysis can be found on Code Ocean.

670 Non-synonymous/ synonymous substitution ratios ($\omega = dN/dS$, or Ka/Ks) have become a useful
671 means for quantifying the impact of natural selection on molecular evolution. In general, the ratio
672 $\omega = dN/dS$ is less than one if the gene is undergoing purifying selection, equal to one if the gene
673 is evolving neutrally, and greater than one if positive selection has accelerated the fixation of non-
674 synonymous substitutions that resulted in amino acid changes. The pair-wise computation of
675 Ka/Ks between V-set exon of each species were performed using the program DnaSp v.0 6.0. The
676 initial unrooted tree fed to the program in the format of Newick was: ((Chimpanzee5:0.00000000,
677 Bonobo:0.00000000):0.00222522, Gorilla:0.00979959, Human:0.01351877).

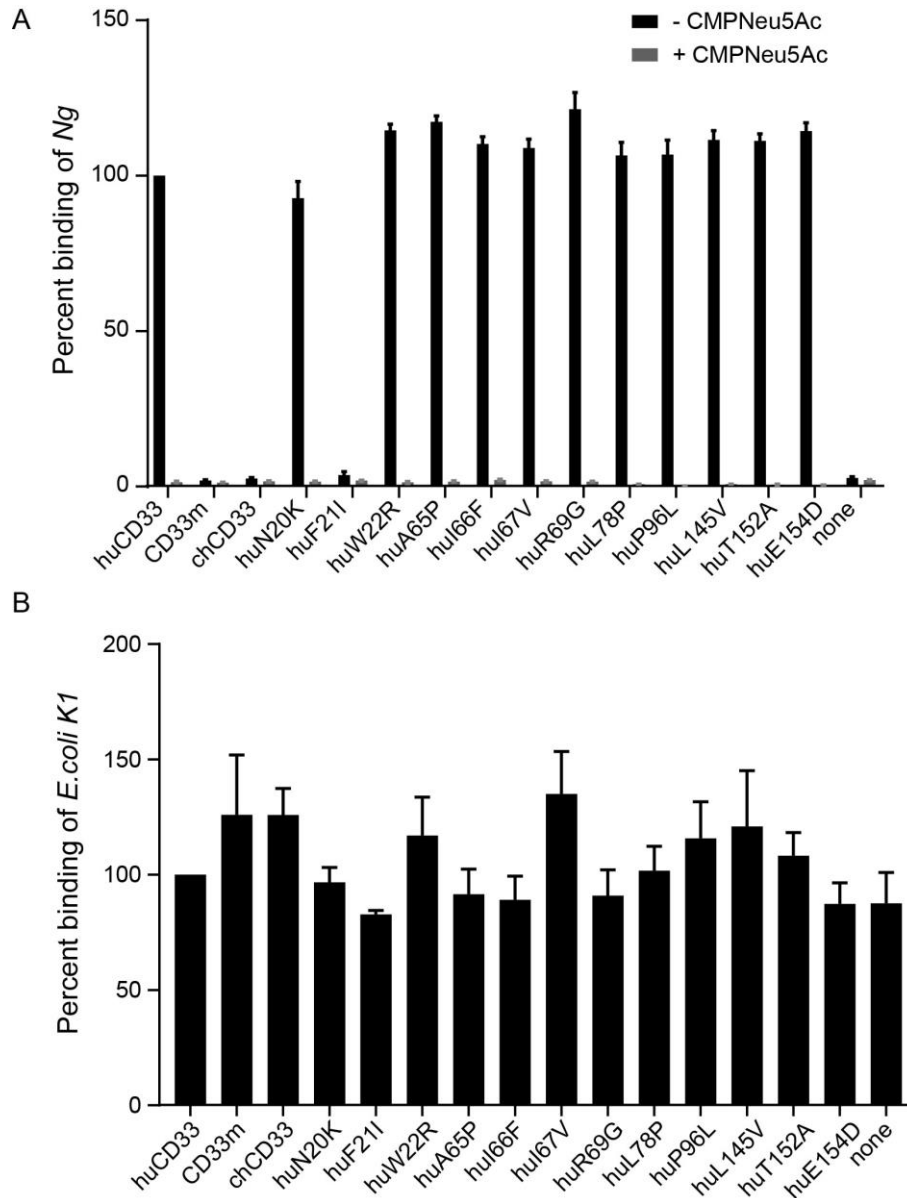
678 **Molecular Simulation.** Starting structures of the V-type domain (residues 18-142) of CD33 were
679 built based on PDB entries 5j0b (A chain) and 6d49 using the graphical interface of YASARA
680 [70]. The two structures differ significantly with respect to the conformation of the C-C' loop

681 (residues 63-71, compare Supplemental Figure S3). A single mutation (G69R) was introduced into
682 5j0b to build CD33(human). The initial 3D models of chCD33 were built by swapping residues:
683 N20K, F21I, W22R, A65P, I67V, R69G (in 6d49), L78P, P96L. An N-glycan core (M3) was
684 attached to Asn100. The side chain of His45 was modeled in two conformations (compare Figure
685 6): ‘down’ (as in PDB entry 6d49) and ‘up’ (as present in PDB entries 5ihb or 5j06 chains A). The
686 systems were solvated in 0.9% NaCl solution (0.15 M) and simulations were performed at 310 K
687 using periodic boundary conditions. The box size was rescaled dynamically to maintain a water
688 density of 0.996 g/ml. Additionally systems were built that contain five molecules of
689 Neu5Gc α OMe or Neu5Ac α OMe distributed in the simulation box which allowed to simulate
690 binding events. Simulations were performed using YASARA with GPU acceleration [71]. In total
691 27 MD trajectories were sampled for huCD33 and 20 for chCD33, most of them covering a
692 microsecond timescale (compare Supplemental Figure S4). Conformational Analysis Tools (CAT,
693 <http://www.md-simulations.de/CAT/>) was used for analysis of trajectory data, general data
694 processing and generation of scientific plots. VMD [72] was used to generate molecular graphics.

695 **Sialoglycan microarray.** The sialoglycan microarray experimental method was adopted from the
696 literature reported earlier [73, 74]. Chemoenzymatically synthesized sialoglycans were quantitated
697 utilizing DMB-HPLC method [75] and 10 mM aqueous stock solutions were prepared. Next, the
698 glycans were diluted to 100 μ M in 300 mM Na-phosphate buffer (pH 8.4) and printed in
699 quadruplets on NHS-functionalized glass slides (PolyAn 3D-NHS; catalog# PO-10400401) using
700 an ArrayIt SpotBot® Extreme instrument. The slides were blocked using 0.05M ethanolamine
701 solution in 0.1 M Tris-HCl (pH 9.0), washed with warm Milli-Q water and dried. Printed slides
702 were fitted in a multi-well microarray hybridization cassette (ArrayIt, CA) and rehydrated using
703 400 μ l of ovalbumin (1% w/v, PBS) for one hour in a humid chamber with gentle shaking. The
704 solution was discarded followed by the addition of 400 μ l solution of the CD33 protein (30 μ g/ml
705 in PBS with 1% w/v ovalbumin) in the individual well. The slides were incubated for 2 h at ambient
706 temperature with gentle shaking followed by washing with PBS-Tween (0.1% v/v) and PBS. The
707 wells were then treated with Cy3-conjugated goat anti-human IgG (1:500 dilution in PBS),
708 incubated for 1h in a dark humid chamber with gentle shaking. After washing and drying, the slides
709 were scanned using a Genepix 4000B scanner (Molecular Devices Corp., Union City, CA) at
710 wavelength 532 nm. Data analysis was performed using the Genepix Pro 7.3 software (Molecular
711 Devices Corp., Union City, CA).

712
713
714

Supplemental Information



715
716

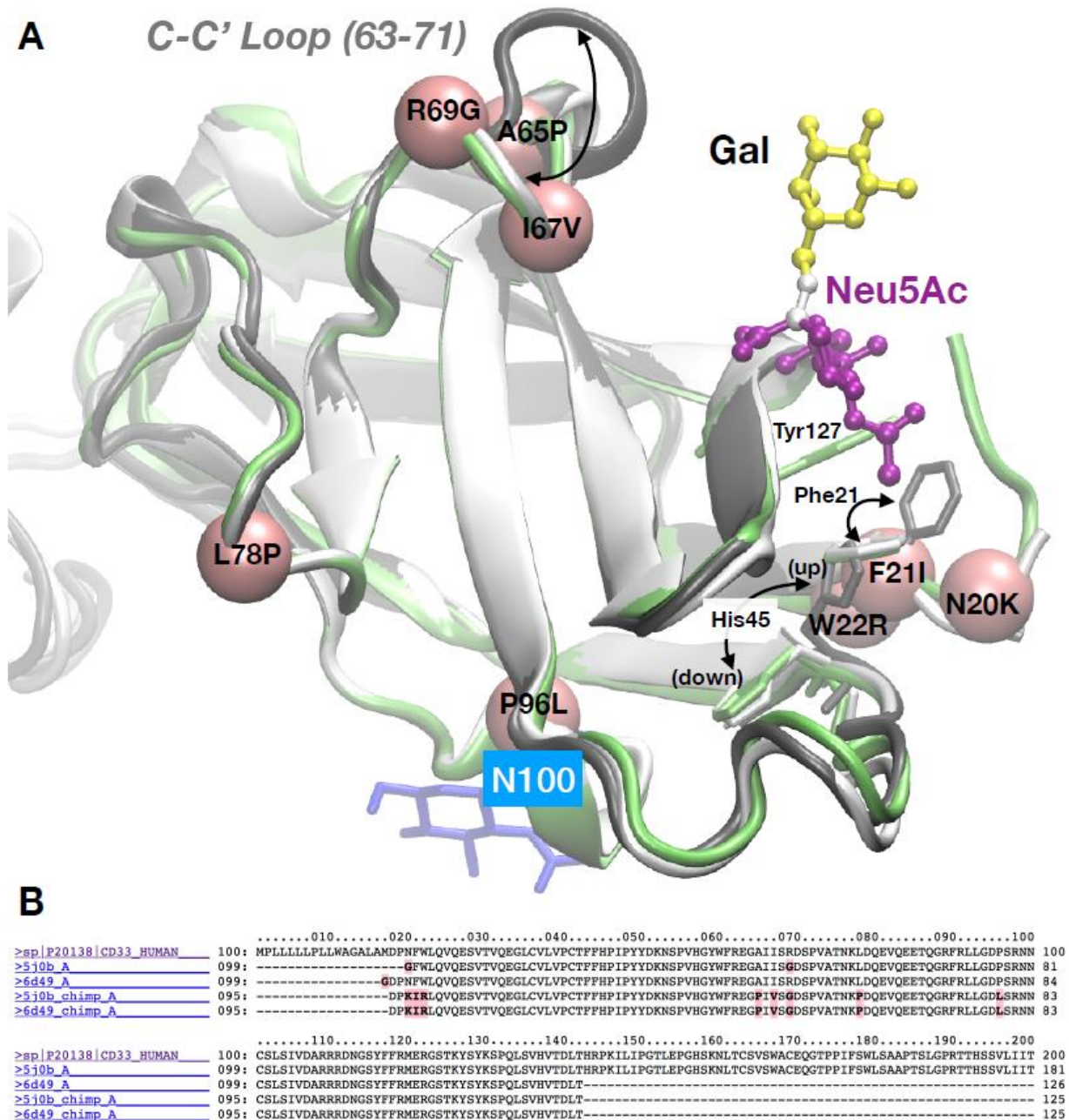
717 **Supplemental Figure S1: While *E. coli* does not bind CD33, human CD33 binding by**
718 ***Neisseria* is Sia-dependent.** (A) Binding of Ng with wildtype or mutant CD33 proteins was
719 determined in the same manner as in Figure 2A. The bacteria for the assay were either grown in
720 presence (+) or absence (-) of exogenous CMP-Neu5Ac as indicated in the legend. All the binding
721 was normalized to wildtype human CD33 binding. Cumulative data from 2 independent
722 experiments, each done in triplet is presented. (B) Binding of *E. coli* K1 was determined using the

723 different CD33 proteins. None of the proteins showed any increased binding to the bacteria relative
724 to no protein (control) containing blank well, indicating that there is no binding of the bacteria
725 with the protein.

Sialic acid	Linkage	huCD33	CD33m	huR>A	huN20K	huF211	huW22R	huA65P	huI66F	chCD33	chK20N	chI21F
Non-Sia		-	-	-	!*	-	-	-	-	-	-	-
Neu5Ac	α 2-3	++	-	-	++	-	+++	-	!*	-	-	++
	α 2-6	+++	-	-	+++	-	+++	-	-	-	-	+++
	α 2-8	-°	-	-	+	-	-°	-	-	-	-	+++
Neu4,5Ac ₂	α 2-3	+	-	-	+	-	+	-	-	-	-	+++
Neu5,9Ac ₂	α 2-3	++	-	-	-/+	-	+	-	-	-	-	+
	α 2-6	-/+	-	-	+	-	-/+	-	-	-	-	+++
Neu5Ac8Me		-	-	-	-	-	-	-	-	-	-	-
Neu5Gc	α 2-3	+	-	-	++	+++	-°	-	!*	+++	-/+	-
	α 2-6	+++	-	-	+++	+++	-	-	-	+++	+++	++
	α 2-8	-	-	-	-	++	-	-	-	+++	++	++
Neu4Ac5Gc	α 2-3	++	-	-	++	+++	-	-	-	++	-	-/+
Neu5Gc9Ac	α 2-3	+	-	-	+	+++	-	-	-	+++	-	-
	α 2-6	++	-	-	+++	+++	-	-	-	+	+	+
Neu5Gc ^{Me}		-	-	-	-	-	-	-	-	-	-	-
Ganglioside type		-	-	-	-	-	-	-	-	-	-	-°

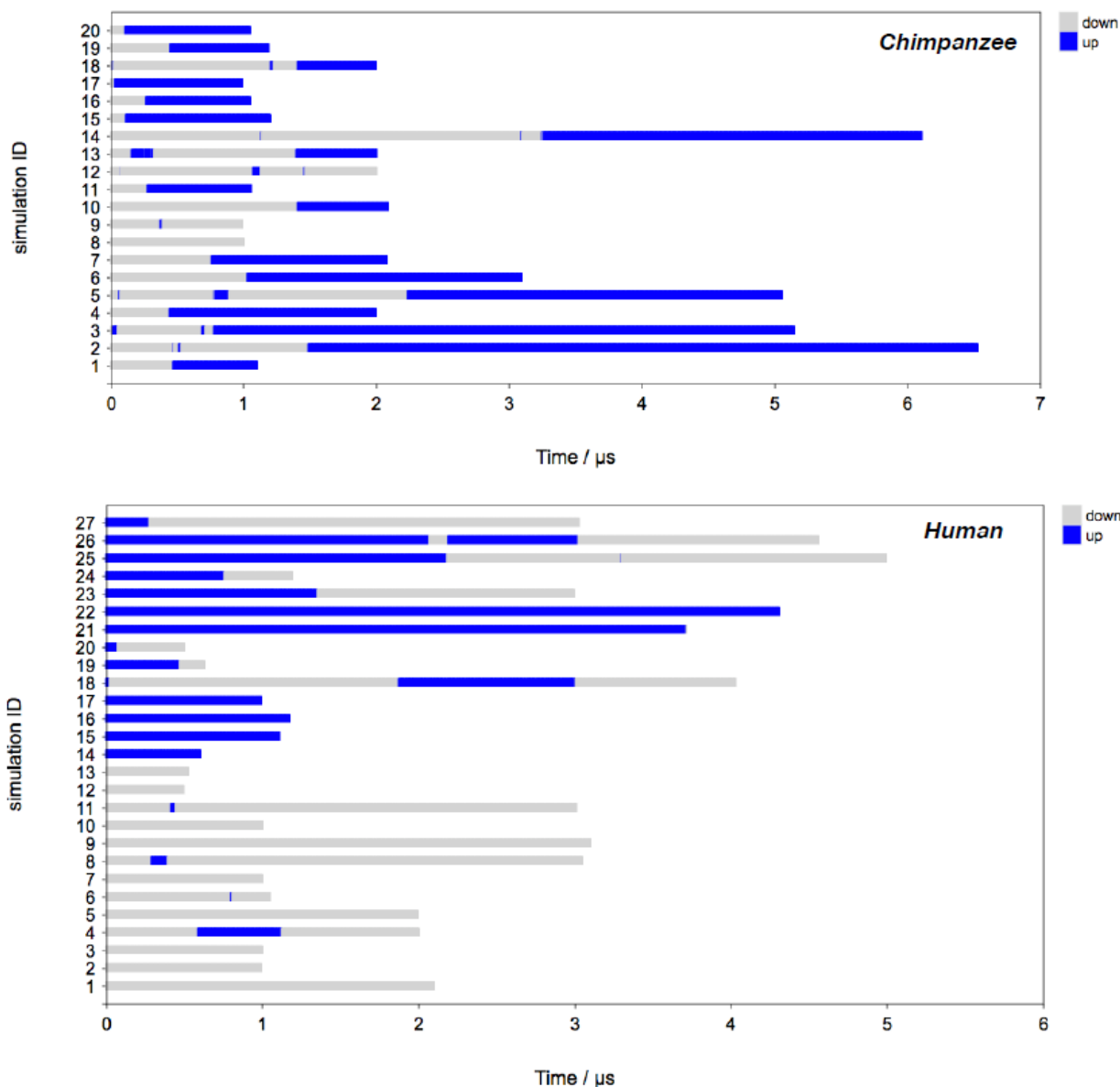
726

727 **Supplemental Figure S2: Summarized result of CD33 sialoglycan binding.** The results of the
728 sialoglycan microarray binding of different wildtype and mutant CD33 proteins presented in
729 Figure 3 are summarized here. A differential sialoglycan binding preference was observed when
730 wildtype and mutant human/chimpanzee CD33 proteins were tested on the microarray. Binding is
731 annotated with a positive (+) symbol and the strength of the binding is indicated by the number of
732 the symbols. +++ indicates a very strong binding. Negative (-) symbol implies non-binding and -
733 +/- indicates very faint interaction. In some cases, only a few sulfated glycans showed strong
734 binding signal (indicated with asterisk). Degree (°) symbols indicate binding with a very few
735 numbers of glycans only. Linkage indicates the nature of the glycosidic bond of the terminal Sia
736 to the underlying glycan.



737

738 **Supplemental Figure S3: Structure and dynamics of CD33.** (A) Examples of x-ray structures
739 of huCD33. PDB entries 5jhb (chain A: dark grey, chain B: white), 6d49 (lime). The dynamics of
740 the C-C' loop and residues Phe21 and His45 are indicated. Positions of mutations present in
741 chimpanzee are labeled on the pink spheres. (B) Amino acid sequences. 1: CD33 human (Uniprot),
742 2-3: PDB entries used for modeling. 3-4: sequences of the chCD33 models.



743

744 **Supplemental Figure S4: MD trajectories of up/down states of Histidine at position 45 (His**

745 **45)**. Molecular dynamics of His45 side chain orientation. Individual MD trajectories of

746 'up'(blue)/'down'(grey) conformational states are shown. For chCD33, it can be observed

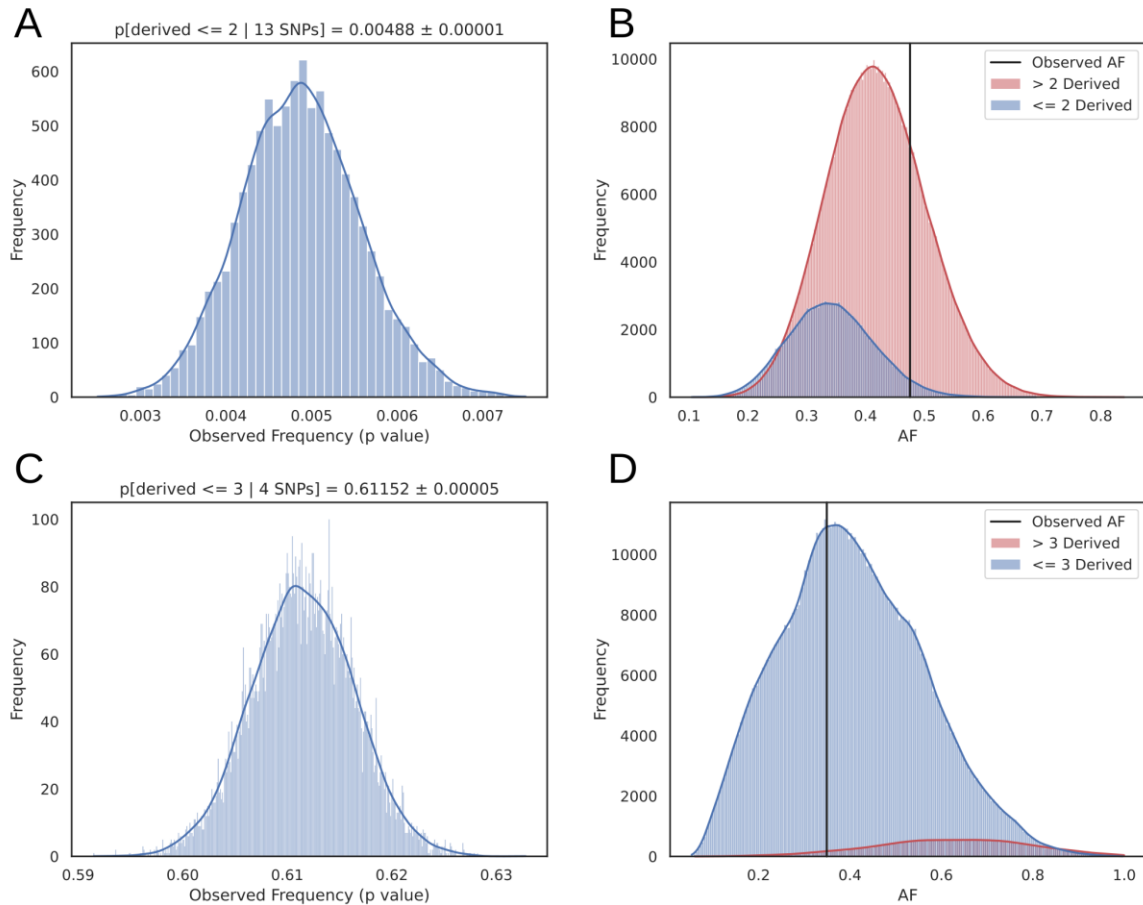
747 reproducibly that simulations started with His45 in a 'down' conformational state undergo a

748 transition to the 'up' conformational state on the microsecond timescale. Therefore, it may be

749 concluded that chCD33 exists mainly with His45 in an 'up' orientation, which would be favorable

750 for binding of Neu5Gc. For huCD33, both conformational states can exist for multiple μs, which

751 explains why huCD33 can bind to Neu5Ac (preferably binds when His45 is ‘down’) and Neu5Gc
752 (preferably binds when His45 is ‘up’).



753

754

755 **Supplemental Figure S5: Resampling analysis of the 5.9 million SNPs from 1000 genomes**

756 **variants.** As an alternative to matching AF directly, the set of filtered SNPs were further restricted
757 to those with a derived population frequency greater than 0.05 resulting in a universe of 5.9 million
758 SNPs. We estimated the probability of observing at most two SNPs derived in either the
759 Neanderthal or Denisovan reference genomes and a mean allele frequency as large as the empirical
760 variants of interest (AF = 0.476). By randomly drawing SNPs, we found that the probability of
761 observing 13 SNPs with such as high global allele frequency and lack of derived alleles in archaic
762 genomes to be highly unusual (p -value = 0.00487 ± 0.00001) (**A**). The low frequency is driven by
763 two factors, as shown in (**B**). Most of the SNPs sampled have more than two archaic-derived SNPs
764 (red curve). Of those with fewer than two archaic-derived SNPs, the overall allele frequency is
765 typically low compared to the target set. With other Siglec SNPs, resampling captures similar
766 properties (**C** and **D**), indicating the LOAD protective set does not represent a random sampling
767 from the genome.

Gene	Associated disease	SNP ID	Description	References
<i>CD33</i>	LOAD	rs12459419, rs3865444	This study for details	[16, 47]
<i>APOE</i>	LOAD, CVD	rs7412, rs429358	Encodes plasma protein APOE, is polymorphic in humans. Three alleles (E2, E3, E4) encode proteins with distinct affinity for lipoprotein particles. The ancestral E4 allele is associated with highest LOAD risk, and increased atherosclerosis and vascular dementia. The derived alleles E2 and E3 seems protective against LOAD, with the lowest risk is in homozygous E2 individuals.	[48, 65, 76]
<i>PICALM</i>	AD	rs3851179 rs10792832	Encodes phosphatidylinositol-binding clathrin assembly protein (PICALM), considered to be one of numerous reproducible risk genes for LOAD.	[49]
<i>SPON1</i>	Dementia	rs2618516	Encodes the developmentally regulated protein F-spondin, reported to be a putative ligand for the amyloid precursor protein (APP).	[50]
<i>TCF7L2</i>	Diabetes	rs7903146	Associated with impaired insulin secretion and enhanced hepatic glucose production.	[52]
<i>ARID5B</i>	AD	rs2588969	Gene encodes a member of AT-rich interaction domain (ARID) family of DNA binding proteins. The encoded protein forms a histone H3K9Me2 demethylase complex with PHD finger protein 2 and regulates the transcription of target genes involved in adipogenesis and liver development.	[51]
<i>PILRA</i>	AD	rs1859788	A cell surface inhibitory receptor that recognizes specific O-glycosylated proteins and expressed on various innate immune cell types including microglia	[53]
<i>CYP3A5</i>	Salt retention and hypertension	rs776746	Cytochrome P450 (CYP) genes are abundant in animal, plant, and bacterial genomes and have evolved to metabolize a variety of diverse compounds.	[54]
<i>PPARG</i>	Diabetes	rs1801282	A nuclear hormone receptor that regulates adipogenesis	[55]
<i>BIN1</i>	AD	rs7561528	Also known as amphiphysin 2, has recently been identified as the most important LOAD risk locus	[20]
<i>SCG2</i>	Hypertension	rs1017448	Secretogranin II (SCG2) associates with hypertension	[56]
<i>CAPN10</i>	Diabetes	rs2975760	CAPN10 encodes a member of the calpain-like cysteine protease family that regulates blood glucose levels.	[57]
<i>AGT</i>	Sodium retention	rs699	Sodium homeostasis links with hypertension	[58]

768 LOAD: Late onset Alzheimer's disease; AD: Alzheimer's disease; CVD: Cardiovascular disease

769

770 **Supplemental Table S1:** Genes affecting cognitive functions in post-reproductive age exhibiting

771 disease-protective alleles uniquely in humans. The corresponding references for each of the genes

772 are mentioned in the table.

Gene	SNP ID	Allele	Global frequency	African	East Asian	European	South Asian	American
<i>CD33</i>	rs12459419	C	0.789	0.949	0.814	0.69	0.84	0.52
		T	0.211	0.051	0.186	0.31	0.16	0.48
	rs3865444	C	0.789	0.949	0.814	0.69	0.84	0.52
		A	0.211	0.051	0.186	0.31	0.16	0.48
<i>APOE</i>	rs7412	C	0.925	0.897	0.9	0.937	0.96	0.96
		T	0.075	0.103	0.1	0.063	0.04	0.04
	rs429358	T	0.849	0.732	0.914	0.845	0.91	0.9
		C	0.151	0.268	0.086	0.155	0.09	0.1
<i>PICALM</i>	rs3851179	T	0.351	0.105	0.407	0.371	0.39	0.39
		C	0.685	0.895	0.593	0.629	0.61	0.61
	rs10792832	A	0.313	0.094	0.409	0.372	0.4	0.39
		G	0.685	0.895	0.593	0.628	0.6	0.61
<i>SPON1</i>	rs2618516	T	0.341	0.259	0.302	0.382	0.52	0.24
		C	0.659	0.741	0.698	0.618	0.48	0.76
<i>TCFLC2</i>	rs7903146	C	0.772	0.74	0.977	0.683	0.7	0.77
		T	0.228	0.26	0.0023	0.317	0.3	0.23
<i>ARID5B</i>	rs2588969	C	0.532	0.472	0.482	0.641	0.63	0.42
		A	0.468	0.528	0.518	0.359	0.37	0.58
<i>PILRA</i>	rs1859788	A	0.341	0.102	0.612	0.321	0.29	0.5
		G	0.659	0.898	0.388	0.679	0.71	0.5
<i>CYP3A5</i>	rs776746	T	0.379	0.82	0.287	0.05	0.33	0.2
		C	0.621	0.18	0.713	0.95	0.67	0.8
<i>PPARG</i>	rs1801282	C	0.93	0.995	0.974	0.88	0.88	0.88
		G	0.07	0.005	0.026	0.12	0.12	0.12
<i>BINI</i>	rs7561528	G	0.8	0.809	0.881	0.683	0.87	0.74
		A	0.2	0.191	0.119	0.317	0.13	0.26
<i>SGC2</i>	rs1017448	C	0.879	0.635	0.963	0.979	0.97	0.95
		T	0.121	0.365	0.037	0.021	0.03	0.05
<i>CAPN10</i>	rs2975760	T	0.882	0.971	0.907	0.841	0.79	0.87
		C	0.118	0.029	0.093	0.159	0.21	0.13
<i>AGT</i>	rs699	A	0.295	0.097	0.147	0.588	0.36	0.36
		G	0.705	0.903	0.853	0.412	0.64	0.64

773

774 **Supplemental Table S2:** Analysis of Gene variants directly or indirectly affecting cognitive
775 function with their human population frequency. The global frequency of the SNPs identified in
776 Supplemental Table S1 was studied across different populations as indicated in the top of the
777 columns.

Amino acid position	Human CD33(M)_EK_Fc Variant	Chimpanzee CD33(M)_EK_Fc Variant	Mutagenesis Primer Pairs_Forward/Reverse_5' > 3'
20	N20K	-	TGGATCCAAAaTTCTGGCTGCAAGTGCAGG TAGCCAGGGCCCCTGCCC
21	F21I	-	GGATCCAAATaTCTGGCTGCAAGTGCAG ATAGCCAGGGCCCCTGCC
22	W22R	-	TCCAAATTTTcGGCTGCAAGTGCAGG TCCATAGCCAGGGCCCCT
65	A65P	-	CCGGAAGGAcCCATTATATC AACCAGTAACCATGAACTG
66	I66F	-	GGAAGGAGCCtTTATATCCAGG CGGAACCAGTAACCATGAAC
67	I67V	-	AGGAGCCATTgTATCCAGGGAC TCCCGGAACCAGTAACCA
69	R69G	-	CATTATATCCgGGGACTCTCCAGTG GCTCCTTCCCGGAACCAG
78	L78P	-	ACAAACAAGCcAGATCAAGAAGTACAGGAG GGCCACTGGAGAGTCCCT
96	P96L	-	CTTGGGGATCiCAGTAGGAACAAC GAGGCGGAATCTGCCCTG
148	L148V	-	GCCCAAATCgTCATCCCTGG CTGTGGGTCAAGTCTGTC
152	T152A	-	CATCCCTGGCgCTCTAGAACC AGGATTTTGGGCCTGTGG
154	E154D	-	GCACTTAGAiCCCGGCCACT CAGGGATGAGGATTTTGGG
21	-	I21F	GGATCCAAAaTCCGGCTGCAAGTG ATAGCCAGGGCCCCTGTG
20	-	K20N	TGGATCCAAAaTCCGGCTGCAAGTGC TAGCCAGGGCCCCTGTGG

778

779 **Supplemental Table S3:** List of the mutagenesis primers used in the study to generate the CD33

780 mutants. Lowercase letters correspond to base change.

781 **Supplemental File S1: List of the glycans used for the sialoglycan microarray.** The complete
782 list of the chemoenzymatically synthesized glycans used to determine the binding profile of
783 different CD33 proteins are presented. The binding intensity of the different proteins (indicated on
784 the top of the columns) towards the corresponding glycan are shown in the heatmap (same heatmap
785 as in Figure 3). The red indicates maximum, and blue indicates minimum binding. R =
786 propylamine linker present in the underlying glycan structure. Gal = galactose, GalNAc = *N*-
787 acetylgalactosamine, Glc = glucose, GlcNAc = *N*-acetyl glucosamine, Fuc = *L*-fucose. The linkage
788 between the monosaccharides is indicated as α - or β - with numbers.

789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833

References

1. Hawkes K., O'Connell J.F., Jones N.G., Alvarez H., and Charnov E.L. (1998). Grandmothering, menopause, and the evolution of human life histories. *Proc Natl Acad Sci U S A.* *95*, 1336-1339.
2. Hawkes K. (2004). Human longevity: the grandmother effect. *Nature.* *428*, 128-129.
3. Hawkes K. (2010). Colloquium paper: how grandmother effects plus individual variation in frailty shape fertility and mortality: guidance from human-chimpanzee comparisons. *Proc Natl Acad Sci U S A.* *107 Suppl 2*, 8977-8984.
4. Johnstone R.A., and Cant M.A. (2010). The evolution of menopause in cetaceans and humans: the role of demography. *Proc Biol Sci.* *277*, 3765-3771.
5. Cant M.A., and Croft D.P. (2019). Life-History Evolution: Grandmothering in Space and Time. *Curr Biol.* *29*, R215-R218.
6. Khan N., Kim S.K., Gagneux P., Dugan L.L., and Varki A. (2020). Maximum reproductive lifespan correlates with CD33rSIGLEC gene number: Implications for NADPH oxidase-derived reactive oxygen species in aging. *FASEB J.* *34*, 1928-1938.
7. Byars S.G., and Voskarides K. (2020). Antagonistic Pleiotropy in Human Disease. *J Mol Evol.* *88*, 12-25.
8. Williams G.C. (1957). Pleiotropy, Natural Selection, and the Evolution of Senescence. *Evolution.* *11*, 398-411.
9. Schwarz F., Springer S.A., Altheide T.K., Varki N.M., Gagneux P., and Varki A. (2016). Human-specific derived alleles of CD33 and other genes protect against postreproductive cognitive decline. *Proc Natl Acad Sci U S A.* *113*, 74-79.
10. Hawkes K. (2016). Genomic evidence for the evolution of human postmenopausal longevity. *Proc Natl Acad Sci U S A.* *113*, 17-18.
11. Läubli H., and Varki A. (2020). Sialic acid-binding immunoglobulin-like lectins (Siglecs) detect self-associated molecular patterns to regulate immune responses. *Cell Mol Life Sci.* *77*, 593-605.
12. Siddiqui S.S., Springer S.A., Verhagen A., Sundaramurthy V., Alisson-Silva F., Jiang W., Ghosh P., and Varki A. (2017). The Alzheimer's disease-protective CD33 splice variant mediates adaptive loss of function via diversion to an intracellular pool. *J Biol Chem.* *292*, 15312-15320.
13. Freeman S.D., Kelm S., Barber E.K., and Crocker P.R. (1995). Characterization of CD33 as a new member of the sialoadhesin family of cellular interaction molecules. *Blood.* *85*, 2005-2012.
14. Bornhöfft K.F., Goldammer T., Rebl A., and Galuska S.P. (2018). Siglecs: A journey through the evolution of sialic acid-binding immunoglobulin-type lectins. *Dev Comp Immunol.* *86*, 219-231.
15. Hernandez-Caselles T., Martinez-Esparza M., Perez-Oliva A.B., Quintanilla-Cecconi A.M., Garcia-Alonso A., Alvarez-Lopez D.M., and Garcia-Penarrubia P. (2006). A study of CD33 (SIGLEC-3) antigen expression and function on activated human T and NK cells: two isoforms of CD33 are generated by alternative splicing. *J Leukoc Biol.* *79*, 46-58.
16. Malik M., Simpson J.F., Parikh I., Wilfred B.R., Fardo D.W., Nelson P.T., and Estus S. (2013). CD33 Alzheimer's risk-altering polymorphism, CD33 expression, and exon 2 splicing. *J Neurosci.* *33*, 13320-13325.

- 834 17. Saha S., Siddiqui S.S., Khan N., Verhagen A., Jiang W., Springer S., Ghosh P., and Varki A.
835 (2019). Controversies about the subcellular localization and mechanisms of action of the
836 Alzheimer's disease-protective CD33 splice variant. *Acta Neuropathol.* *138*, 671-672.
- 837 18. Lamba J.K., Pounds S., Cao X., Downing J.R., Campana D., Ribeiro R.C., Pui C.H., and
838 Rubnitz J.E. (2009). Coding polymorphisms in CD33 and response to gemtuzumab
839 ozogamicin in pediatric patients with AML: a pilot study. *Leukemia.* *23*, 402-404.
- 840 19. Bradshaw E.M., Chibnik L.B., Keenan B.T., Ottoboni L., Raj T., Tang A., Rosenkrantz L.L.,
841 Imboywa S., Lee M., Von Korff A. et al. (2013). CD33 Alzheimer's disease locus: altered
842 monocyte function and amyloid biology. *Nat Neurosci.* *16*, 848-850.
- 843 20. Naj A.C., Jun G., Beecham G.W., Wang L.S., Vardarajan B.N., Buross J., Gallins P.J.,
844 Buxbaum J.D., Jarvik G.P., Crane P.K. et al. (2011). Common variants at MS4A4/MS4A6E,
845 CD2AP, CD33 and EPHA1 are associated with late-onset Alzheimer's disease. *Nat Genet.*
846 *43*, 436-441.
- 847 21. Hollingworth P., Harold D., Sims R., Gerrish A., Lambert J.C., Carrasquillo M.M., Abraham
848 R., Hamshere M.L., Pahwa J.S., Moskvin V. et al. (2011). Common variants at ABCA7,
849 MS4A6A/MS4A4E, EPHA1, CD33 and CD2AP are associated with Alzheimer's disease.
850 *Nat Genet.* *43*, 429-435.
- 851 22. Brinkman-Van der Linden E.C., Angata T., Reynolds S.A., Powell L.D., Hedrick S.M., and
852 Varki A. (2003). CD33/Siglec-3 binding specificity, expression pattern, and consequences
853 of gene deletion in mice. *Mol Cell Biol.* *23*, 4199-4206.
- 854 23. Edler M.K., Sherwood C.C., Meindl R.S., Hopkins W.D., Ely J.J., Erwin J.M., Mufson E.J.,
855 Hof P.R., and Raghanti M.A. (2017). Aged chimpanzees exhibit pathologic hallmarks of
856 Alzheimer's disease. *Neurobiol Aging.* *59*, 107-120.
- 857 24. Edler M.K., Munger E.L., Meindl R.S., Hopkins W.D., Ely J.J., Erwin J.M., Mufson E.J.,
858 Hof P.R., Sherwood C.C., and Raghanti M.A. (2020). Neuron loss associated with age but
859 not Alzheimer's disease pathology in the chimpanzee brain. *Philos Trans R Soc Lond B Biol*
860 *Sci.* *375*, 20190619.
- 861 25. Padler-Karavani V., Hurtado-Ziola N., Chang Y.C., Sonnenburg J.L., Ronaghy A., Yu H.,
862 Verhagen A., Nizet V., Chen X., Varki N. et al. (2014). Rapid evolution of binding
863 specificities and expression patterns of inhibitory CD33-related Siglecs in primates. *FASEB*
864 *J.* *28*, 1280-1293.
- 865 26. Varki A. (2011). Since there are PAMPs and DAMPs, there must be SAMPs? Glycan "self-
866 associated molecular patterns" dampen innate immunity, but pathogens can mimic them.
867 *Glycobiology.* *21*, 1121-1124.
- 868 27. Varki A., and Gagneux P. (2012). Multifarious roles of sialic acids in immunity. *Ann N Y*
869 *Acad Sci.* *1253*, 16-36.
- 870 28. Tortorici M.A., Walls A.C., Lang Y., Wang C., Li Z., Koerhuis D., Boons G.J., Bosch B.J.,
871 Rey F.A., de Groot R.J. et al. (2019). Structural basis for human coronavirus attachment to
872 sialic acid receptors. *Nat Struct Mol Biol.* *26*, 481-489.
- 873 29. Tsai T.Y., Huang M.T., Sung P.S., Peng C.Y., Tao M.H., Yang H.I., Chang W.C., Yang
874 A.S., Yu C.M., Lin Y.P. et al. (2021). SIGLEC-3 (CD33) serves as an immune checkpoint
875 receptor for HBV infection. *J Clin Invest.* *131*, 141965.
- 876 30. Varki A., and Angata T. (2006). Siglecs--the major subfamily of I-type lectins.
877 *Glycobiology.* *16*, 1R-27R.

- 878 31. Consortium G.P., Auton A., Brooks L.D., Durbin R.M., Garrison E.P., Kang H.M., Korbel
879 J.O., Marchini J.L., McCarthy S., McVean G.A. et al. (2015). A global reference for human
880 genetic variation. *Nature*. 526, 68-74.
- 881 32. Prado-Martinez J., Sudmant P.H., Kidd J.M., Li H., Kelley J.L., Lorente-Galdos B.,
882 Veeramah K.R., Woerner A.E., O'Connor T.D., Santpere G. et al. (2013). Great ape genetic
883 diversity and population history. *Nature*. 499, 471-475.
- 884 33. Xue Y., Prado-Martinez J., Sudmant P.H., Narasimhan V., Ayub Q., Szpak M., Frandsen P.,
885 Chen Y., Yngvadottir B., Cooper D.N. et al. (2015). Mountain gorilla genomes reveal the
886 impact of long-term population decline and inbreeding. *Science*. 348, 242-245.
- 887 34. de Manuel M., Kuhlwilm M., Frandsen P., Sousa V.C., Desai T., Prado-Martinez J.,
888 Hernandez-Rodriguez J., Dupanloup I., Lao O., Hallast P. et al. (2016). Chimpanzee genomic
889 diversity reveals ancient admixture with bonobos. *Science*. 354, 477-481.
- 890 35. Green R.E., Krause J., Briggs A.W., Maricic T., Stenzel U., Kircher M., Patterson N., Li H.,
891 Zhai W., Fritz M.H. et al. (2010). A draft sequence of the Neandertal genome. *Science*. 328,
892 710-722.
- 893 36. Prufer K., Racimo F., Patterson N., Jay F., Sankararaman S., Sawyer S., Heinze A., Renaud
894 G., Sudmant P.H., de Filippo C. et al. (2014). The complete genome sequence of a
895 Neanderthal from the Altai Mountains. *Nature*. 505, 43-49.
- 896 37. Edwards J.L., and Apicella M.A. (2004). The molecular mechanisms used by *Neisseria*
897 *gonorrhoeae* to initiate infection differ between men and women. *Clin Microbiol Rev*. 17,
898 965-981.
- 899 38. Landig C.S., Hazel A., Kellman B.P., Fong J.J., Schwarz F., Agarwal S., Varki N., Massari
900 P., Lewis N.E., Ram S. et al. (2019). Evolution of the exclusively human pathogen *Neisseria*
901 *gonorrhoeae*: Human-specific engagement of immunoregulatory Siglecs. *Evol Appl*. 12,
902 337-349.
- 903 39. Apicella M.A., Mandrell R.E., Shero M., Wilson M.E., Griffiss J.M., Brooks G.F., Lammel
904 C., Breen J.F., and Rice P.A. (1990). Modification by sialic acid of *Neisseria gonorrhoeae*
905 lipooligosaccharide epitope expression in human urethral exudates: an immunoelectron
906 microscopic analysis. *J Infect Dis*. 162, 506-512.
- 907 40. Parsons N.J., Patel P.V., Tan E.L., Andrade J.R., Nairn C.A., Goldner M., Cole J.A., and
908 Smith H. (1988). Cytidine 5'-monophospho-N-acetyl neuraminic acid and a low molecular
909 weight factor from human blood cells induce lipopolysaccharide alteration in gonococci
910 when conferring resistance to killing by human serum. *Microb Pathog*. 5, 303-309.
- 911 41. Caugant D.A., and Brynildsrud O.B. (2020). *Neisseria meningitidis*: using genomics to
912 understand diversity, evolution and pathogenesis. *Nat Rev Microbiol*. 18, 84-96.
- 913 42. Seifert H.S. (2019). Location, Location, Location-Commensalism, Damage and Evolution of
914 the Pathogenic *Neisseria*. *J Mol Biol*. 431, 3010-3014.
- 915 43. Carlin A.F., Lewis A.L., Varki A., and Nizet V. (2007). Group B streptococcal capsular sialic
916 acids interact with siglecs (immunoglobulin-like lectins) on human leukocytes. *J Bacteriol*.
917 189, 1231-1237.
- 918 44. Fong J.J., Tsai C.M., Saha S., Nizet V., Varki A., and Bui J.D. (2018). Siglec-7 engagement
919 by GBS β -protein suppresses pyroptotic cell death of natural killer cells. *Proc Natl Acad Sci*
920 *U S A*. 115, 10410-10415.
- 921 45. Munshi M.N. (2017). Cognitive Dysfunction in Older Adults With Diabetes: What a
922 Clinician Needs to Know. *Diabetes Care*. 40, 461-467.

- 923 46. Peila R., Rodriguez B.L., Launer L.J., and Honolulu-Asia A.S. (2002). Type 2 diabetes,
924 APOE gene, and the risk for dementia and related pathologies: The Honolulu-Asia Aging
925 Study. *Diabetes*. *51*, 1256-1262.
- 926 47. Raj T., Ryan K.J., Replogle J.M., Chibnik L.B., Rosenkrantz L., Tang A., Rothamel K.,
927 Stranger B.E., Bennett D.A., Evans D.A. et al. (2014). CD33: increased inclusion of exon 2
928 implicates the Ig V-set domain in Alzheimer's disease susceptibility. *Hum Mol Genet*.
- 929 48. Finch C.E., and Sapolsky R.M. (1999). The evolution of Alzheimer disease, the reproductive
930 schedule, and apoE isoforms. *Neurobiol Aging*. *20*, 407-428.
- 931 49. Xu W., Tan L., and Yu J.T. (2015). The Role of PICALM in Alzheimer's Disease. *Mol*
932 *Neurobiol*. *52*, 399-413.
- 933 50. Liu Z., Dai X., Tao W., Liu H., Li H., Yang C., Zhang J., Li X., Chen Y., Ma C. et al. (2018).
934 APOE influences working memory in non-demented elderly through an interaction with
935 SPON1 rs2618516. *Hum Brain Mapp*. *39*, 2859-2867.
- 936 51. Carrasquillo M.M., Belbin O., Hunter T.A., Ma L., Bisceglia G.D., Zou F., Crook J.E.,
937 Pankratz V.S., Sando S.B., Aasly J.O. et al. (2011). Replication of EPHA1 and CD33
938 associations with late-onset Alzheimer's disease: a multi-centre case-control study. *Mol*
939 *Neurodegener*. *6*, 54.
- 940 52. Helgason A., Pálsson S., Thorleifsson G., Grant S.F., Emilsson V., Gunnarsdottir S.,
941 Adeyemo A., Chen Y., Chen G., Reynisdottir I. et al. (2007). Refining the impact of TCF7L2
942 gene variants on type 2 diabetes and adaptive evolution. *Nat Genet*. *39*, 218-225.
- 943 53. Rathore N., Ramani S.R., Pantua H., Payandeh J., Bhangale T., Wuster A., Kapoor M., Sun
944 Y., Kapadia S.B., Gonzalez L. et al. (2018). Paired Immunoglobulin-like Type 2 Receptor
945 Alpha G78R variant alters ligand binding and confers protection to Alzheimer's disease.
946 *PLoS Genet*. *14*, e1007427.
- 947 54. Thompson E.E., Kuttub-Boulos H., Witonsky D., Yang L., Roe B.A., and Di Rienzo A.
948 (2004). CYP3A variation and the evolution of salt-sensitivity variants. *Am J Hum Genet*. *75*,
949 1059-1069.
- 950 55. Altshuler D., Hirschhorn J.N., Klannemark M., Lindgren C.M., Vohl M.C., Nemesh J., Lane
951 C.R., Schaffner S.F., Bolk S., Brewer C. et al. (2000). The common PPARgamma Pro12Ala
952 polymorphism is associated with decreased risk of type 2 diabetes. *Nat Genet*. *26*, 76-80.
- 953 56. Wen G., Wessel J., Zhou W., Ehret G.B., Rao F., Stridsberg M., Mahata S.K., Gent P.M.,
954 Das M., Cooper R.S. et al. (2007). An ancestral variant of Secretogranin II confers regulation
955 by PHOX2 transcription factors and association with hypertension. *Hum Mol Genet*. *16*,
956 1752-1764.
- 957 57. Vander Molen J., Frisse L.M., Fullerton S.M., Qian Y., Del Bosque-Plata L., Hudson R.R.,
958 and Di Rienzo A. (2005). Population genetics of CAPN10 and GPR35: implications for the
959 evolution of type 2 diabetes variants. *Am J Hum Genet*. *76*, 548-560.
- 960 58. Nakajima T., Wooding S., Sakagami T., Emi M., Tokunaga K., Tamiya G., Ishigami T.,
961 Umemura S., Munkhbat B., Jin F. et al. (2004). Natural selection and population history in
962 the human angiotensinogen gene (AGT): 736 complete AGT sequences in chromosomes
963 from around the world. *Am J Hum Genet*. *74*, 898-916.
- 964 59. Iadecola C. (2013). The pathobiology of vascular dementia. *Neuron*. *80*, 844-866.
- 965 60. Reich D., Green R.E., Kircher M., Krause J., Patterson N., Durand E.Y., Viola B., Briggs
966 A.W., Stenzel U., Johnson P.L. et al. (2010). Genetic history of an archaic hominin group
967 from Denisova Cave in Siberia. *Nature*. *468*, 1053-1060.

- 968 61. Meyer M., Kircher M., Gansauge M.T., Li H., Racimo F., Mallick S., Schraiber J.G., Jay F.,
969 Prufer K., de Filippo C. et al. (2012). A high-coverage genome sequence from an archaic
970 Denisovan individual. *Science*. *338*, 222-226.
- 971 62. Khan N., de Manuel M., Peyregne S., Do R., Prufer K., Marques-Bonet T., Varki N.,
972 Gagneux P., and Varki A. (2020). Multiple Genomic Events Altering Hominin SIGLEC
973 Biology and Innate Immunity Predated the Common Ancestor of Humans and Archaic
974 Hominins. *Genome Biol Evol*. *12*, 1040-1050.
- 975 63. Mendez F.L., Poznik G.D., Castellano S., and Bustamante C.D. (2016). The Divergence of
976 Neandertal and Modern Human Y Chromosomes. *Am J Hum Genet*. *98*, 728-734.
- 977 64. Moon J.M., Aronoff D.M., Capra J.A., Abbot P., and Rokas A. (2018). Examination of
978 Signatures of Recent Positive Selection on Genes Involved in Human Sialic Acid Biology.
979 *G3 (Bethesda)*. *8*, 1315-1325.
- 980 65. Reiman E.M., Arboleda-Velasquez J.F., Quiroz Y.T., Huentelman M.J., Beach T.G., Caselli
981 R.J., Chen Y., Su Y., Myers A.J., Hardy J. et al. (2020). Exceptionally low likelihood of
982 Alzheimer's dementia in APOE2 homozygotes from a 5,000-person neuropathological
983 study. *Nat Commun*. *11*, 667.
- 984 66. Oriá R.B., Patrick P.D., Zhang H., Lorntz B., de Castro Costa C.M., Brito G.A., Barrett L.J.,
985 Lima A.A., and Guerrant R.L. (2005). APOE4 protects the cognitive development in children
986 with heavy diarrhea burdens in Northeast Brazil. *Pediatr Res*. *57*, 310-316.
- 987 67. Angata T., and Varki A. (2000). Cloning, characterization, and phylogenetic analysis of
988 siglec-9, a new member of the CD33-related group of siglecs. Evidence for co-evolution with
989 sialic acid synthesis pathways. *J Biol Chem*. *275*, 22127-22135.
- 990 68. Angata T., Kerr S.C., Greaves D.R., Varki N.M., Crocker P.R., and Varki A. (2002). Cloning
991 and characterization of human Siglec-11. A recently evolved signaling molecule that can
992 interact with SHP-1 and SHP-2 and is expressed by tissue macrophages, including brain
993 microglia. *J Biol Chem*. *277*, 24466-24474.
- 994 69. Kumar S., Stecher G., and Tamura K. (2016). MEGA7: Molecular Evolutionary Genetics
995 Analysis Version 7.0 for Bigger Datasets. *Mol Biol Evol*. *33*, 1870-1874.
- 996 70. Krieger E., and Vriend G. (2014). YASARA View - molecular graphics for all devices - from
997 smartphones to workstations. *Bioinformatics*. *30*, 2981-2982.
- 998 71. Krieger E., and Vriend G. (2015). New ways to boost molecular dynamics simulations. *J*
999 *Comput Chem*. *36*, 996-1007.
- 1000 72. Humphrey W., Dalke A., and Schulten K. (1996). VMD: visual molecular dynamics. *J Mol*
1001 *Graph*. *14*, 33-8, 27.
- 1002 73. Meng C., Sasmal A., Zhang Y., Gao T., Liu C.C., Khan N., Varki A., Wang F., and Cao H.
1003 (2018). Chemoenzymatic Assembly of Mammalian O-Mannose Glycans. *Angew Chem Int*
1004 *Ed Engl*. *57*, 9003-9007.
- 1005 74. Lu N., Ye J., Cheng J., Sasmal A., Liu C.C., Yao W., Yan J., Khan N., Yi W., Varki A. et al.
1006 (2019). Redox-Controlled Site-Specific α 2-6-Sialylation. *J Am Chem Soc*. *141*, 4547-4552.
- 1007 75. Ji Y., Sasmal A., Li W., Oh L., Srivastava S., Hargett A.A., Wasik B.R., Yu H., Diaz S.,
1008 Choudhury B. et al. (2021). Reversible O-Acetyl Migration within the Sialic Acid Side Chain
1009 and Its Influence on Protein Recognition. *ACS Chem Biol*. *16*, 1951-1960.
- 1010 76. Fullerton S.M., Clark A.G., Weiss K.M., Nickerson D.A., Taylor S.L., Stengård J.H.,
1011 Salomaa V., Vartiainen E., Perola M., Boerwinkle E. et al. (2000). Apolipoprotein E
1012 variation at the sequence haplotype level: implications for the origin and maintenance of a
1013 major human polymorphism. *Am J Hum Genet*. *67*, 881-900.