



Gene prediction in the immunoglobulin loci

Vikram Sirupurapu, Yana Safonova and Pavel Pevzner

Genome Res. published online May 11, 2022

Access the most recent version at doi:[10.1101/gr.276676.122](https://doi.org/10.1101/gr.276676.122)

P<P Published online May 11, 2022 in advance of the print journal.

Accepted Manuscript Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.

Creative Commons License This article is distributed exclusively by Cold Spring Harbor Laboratory Press for the first six months after the full-issue publication date (see <https://genome.cshlp.org/site/misc/terms.xhtml>). After six months, it is available under a Creative Commons License (Attribution-NonCommercial 4.0 International), as described at <http://creativecommons.org/licenses/by-nc/4.0/>.

Email Alerting Service Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or [click here](#).

Advance online articles have been peer reviewed and accepted for publication but have not yet appeared in the paper journal (edited, typeset versions may be posted when available prior to final publication). Advance online articles are citable and establish publication priority; they are indexed by PubMed from initial publication. Citations to Advance online articles must include the digital object identifier (DOIs) and date of initial publication.

To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Published by Cold Spring Harbor Laboratory Press

Gene prediction in the immunoglobulin loci

Vikram Sirupurapu¹, Yana Safonova^{1,2}, and Pavel A. Pevzner^{1,*}

¹Computer Science and Engineering Department, University of California San Diego, San Diego,
USA

²The Department of Computer Science, Johns Hopkins University, Baltimore, USA

* Corresponding author: ppezvner@ucsd.edu

Abstract

The V(D)J recombination process rearranges the variable (V), diversity (D), and joining (J) genes in the immunoglobulin loci to generate antibody repertoires. Annotation of these loci across various species and predicting the V, D, and J genes (IG genes) is critical for studies of the adaptive immune system. However, since the standard gene finding algorithms are not suitable for predicting IG genes, they have been semi-manually annotated in very few species. We developed the IGDetective algorithm for predicting IG genes and applied it to species with the assembled IG loci. IGDetective generated the first large collection of IG genes across many species and enabled their evolutionary analysis, including the analysis of the “bat IG diversity” hypothesis. This analysis revealed extremely conserved V genes in evolutionary distant species indicating that these genes may be subjected to the same selective pressure, e.g., pressure driven by common pathogens. IGDetective also revealed extremely diverged V genes and a new family of evolutionary conserved V genes in bats with unusual non-canonical cysteines. Moreover, in difference from all other previously reported antibodies, these cysteines are located within complementarity-determining regions. Since cysteines form disulfide bonds, we hypothesize that these cysteine-rich V genes might generate antibodies with non-canonical conformations and could potentially form a unique part of the immune repertoire in bats. We also

23 analyzed the diversity landscape of the recombination signal sequences and revealed their features
24 that trigger the high/low usage of the IG genes.

25 **Introduction**

26 *Antibodies* (or *immunoglobulins*) are the key components of the immune system of jawed vertebrates
27 that provide adaptive immune response by recognizing and neutralizing *antigens*. Antibodies are not
28 encoded in the germline genome but rather result from somatic *VDJ recombinations* (Tonegawa,
29 1983). This process affects an immunoglobulin (IG) loci containing the families of the *variable* (V),
30 *diversity* (D), and *joining* (J) genes (referred to as *IG genes*) by selecting one V, one D gene, and one J
31 gene, and concatenating them together to generate one of the antibody chains. Antibodies are further
32 diversified by *somatic hypermutations* (Dudley et al. 2005).

33 The diversity of the IG loci is driven by the variety of antigens: different species encounter different
34 antigens and develop their unique ways to fight them through mutations in IG genes. As a result, the
35 IG loci have rapidly evolved independently in different species and resulted in a diverse collection of
36 IG genes that remain largely unknown since both sequencing highly-repetitive IG loci and predicting
37 IG genes in these loci are challenging tasks (Das et al., 2012, Pettinello and Dooley, 2014).

38 Mammalian genomes have three IG loci: *heavy chain (IGH)*, *kappa light chain (IGK)*, and *lambda*
39 *light chain (IGL)* as well as four *T-cell antigen receptor (TCR)* loci (TRA, TRB, TRG, and TRD). In
40 this work, we mainly focus on the IGH locus. The V, D, and J genes in the IGH locus (and the
41 fragments of the IGH locus containing these genes) are also referred to as IGHV, IGHD, and IGHJ,
42 respectively.

43 *Diversity of immunoglobulin genes*. Studies of adaptive immune responses across various vertebrate
44 species open new therapeutics opportunities (Muyldermans and Smider, 2016). For example, studies
45 of single-chain camelid antibodies led to the development of *nanobodies* that are able to diagnose
46 cancer (Rashidian et al., 2015; Keyaerts et al., 2016), while studies of ultralong cow antibodies

47 revealed that they recognize various HIV strains (Sok et al., 2017). Understanding the diversity of the
48 adaptive immune systems across various vertebrates can also contribute to analyzing the spread of
49 newly emerged zoonotic pathogens. It is particularly important for bat species that possess a unique
50 immune system capable of neutralizing viruses that are often lethal to other mammals, such as rabies,
51 Ebola, SARS-CoV, MERS-CoV, and SARS-CoV-2 (Teeling et al., 2018).

52 Schountz et al., 2017 formulated a “bat IG diversity” hypothesis that argues that, since bats have a
53 greater combinatorial diversity of IG genes than other mammals, their large naive antibody repertoires
54 do not need substantial affinity maturation to successfully neutralize antigens. Indeed, while the
55 human IGH locus has only 55 functional V genes, the little brown bat is estimated to have at least 200
56 V genes (Bratsch et al., 2011), suggesting that bats may be better equipped for clonal selection of B
57 cells responding to viral antigens. However, the IG genes in bats remain poorly characterized,
58 moreover, the only support for the “bat IG diversity” hypothesis comes from a probabilistic model
59 (Bratsch et al., 2011) rather than an annotated IG loci in a well-assembled bat genome.

60 *Annotation of immunoglobulin genes.* Annotation of the IG loci (i.e., predicting IG genes) is a
61 prerequisite for most follow-up immunogenomics studies. Since assembly and annotation of the IG
62 loci for novel species is complicated by their highly-repetitive structure (Matsuda et al., 1998; Watson
63 et al., 2013; Rodriguez et al., 2020), the sequences of the IG loci are only known for a few species.
64 However, recent advancements in long-read sequencing enabled the first contiguous assemblies of
65 highly-repetitive genomic regions and the Vertebrate Genome Project (VGP) now aims to generate
66 high-quality reference genomes for all vertebrate species (Rhie et al., 2020). So far, the VGP has
67 generated well-assembled genomes of nearly 150 vertebrate species.

68 Although many IG loci have been assembled in the last two years, their automated annotation remains
69 an open problem. Previous studies of IG genes combined a time-consuming experimental approach
70 with semi-manual computational analysis, such as in the studies of the platypus (Gambon-Deza et al.,
71 2009), the cow (Ma et al., 2016), and the ferret (Wong et al., 2020). These studies used the human IG
72 genes for a similarity-based detection of IG genes in the novel species and may have missed diverged

73 IG genes. For example, since most (if not all) V genes in mammalian species were predicted based on
74 their similarities with human V and J genes, it remains unclear whether there exist still unknown
75 families of V genes that are highly diverged from the human IG genes.

76 *Annotation of recombination signal sequences.* De novo prediction of IG genes relies on identifying
77 conserved *recombination signal sequences (RSSs)* that flank IG genes and enable the VDJ
78 recombination. Since short RSSs have many spurious RSS-like occurrences in the genome (*cryptic*
79 *RSSs*), their prediction and downstream IG gene annotation is a challenging problem. It is further
80 complicated by the fact that some cryptic RSSs are implicated in unusual genomic rearrangements
81 outside the IG and TCR loci (Messier et al., 2003) and sometimes play an important role in antibody
82 generation. For example, cryptic RSSs flanking the human *LAIR1* gene participate in the off-target
83 VDJ recombination and generates a new type of antibodies (where *LAIR1* represents an additional
84 domain of an unusual VDJ region) that broadly neutralize *Plasmodium* parasites (Tan et al., 2016).
85 Although Teng et al., 2016 analyzed the impact of off-target VDJ recombinations on lymphocyte
86 genomes, no attempts to analyze the landscape of cryptic RSSs within the IG loci have been made yet.

87 The problem of identifying RSSs in a genomic sequence was considered by Merelli et al., 2010
88 (RSSsite tool) and Olivieri et al., 2013 (VgeneExtractor tool). VgeneExtractor further used its RSS
89 predictions for detecting V genes in various species. However, this method does not account for the
90 variations in the RSSs within species and does not report D and J genes. Safonova and Pevzner, 2020
91 recently benchmarked RSSsite and demonstrated that it results in a high false-positive rate. Even
92 though they developed a more accurate SEARCH-D algorithm for detecting RSSs of D genes, the
93 problem of detecting RSSs for all types of IG genes and the follow-up evolutionary analysis of IG
94 genes across multiple species remains open.

95 *Role of cysteines in immunoglobulins.* Cysteines are structurally important amino acids that form
96 *disulfide bonds* in immunoglobulins (Frangione et al., 1969). Conserved cysteines of human
97 immunoglobulins are referred to as *canonical* cysteines (Tonegawa, 1983). Non-canonical cysteines
98 are common in some biomedically important species, such as llamas and cows (de Los Rios, 2015;

99 Prabakaran and Chowdhury, 2020). The current consensus is that the additional disulfide bonds in
100 immunoglobulins increase their stability and enrich the complexity of antigen-binding site topology
101 (Wu et al., 2012). Since cysteine patterns represent a structurally important feature of antibodies, we
102 analyzed non-canonical cysteines in the newly identified IG genes.

103 *IGDetective tool*. We describe the IGDetective tool for predicting IG genes, apply it for annotating IG
104 loci in the recently assembled mammalian genomes, generate the largest set of IG genes to date, and
105 reveal surprising diversity of IG genes across multiple species, such as a new family of unusual
106 cysteine-rich V genes in bats.

107 Since all previous attempts to annotate V (J) genes in newly sequenced species relied on their
108 similarities with known human IG genes, it remains unclear whether there exist V (J) genes that are
109 not similar to the human ones. In addition to a similarity-based search for IG genes
110 (IterativeIGDetective mode), IGDetective has a BlindIGDetective mode for annotating IG genes in
111 the absence of any prior knowledge about the sequences of IG genes in other species. We benchmark
112 BlindIGDetective on well-annotated human, mouse, and cow genomes, demonstrate that it
113 automatically derives nearly all known IG genes in these species in a blind fashion, apply it to newly-
114 sequenced genomes to reveal highly-diverged IG genes, and reveal new families of highly diverged V
115 genes in bats that evade the similarity-based approach to predicting IG genes.

116 **Results**

117 **From predicting RSSs to predicting IG genes.** Sequences of human IG genes were inferred as the
118 result of painstaking manual analysis at the dawn of the DNA sequencing era (Li et al., 2002).
119 Consequently, sequences of all non-human IG genes were inferred based on similarities with human
120 IG genes (Sitnikova and Su, 1998). IterativeIGDetective automates this approach and iteratively
121 extends it by predicting more and more distant IG genes at each iteration. In contrast,
122 BlindIGDetective is designed to predict IG genes even if they share no similarities with currently
123 known IG genes.

124 Given a newly assembled genome, both IterativeIGDetective and BlindIGDetective start from
125 predicting RSSs in this genome. Each RSS consists of a conserved *heptamer* followed by a non-
126 conserved spacer (12 nt long in D genes and 23 nt long in V and J genes in the IGH locus), and a
127 conserved *nonamer* (**Figure 1A**). Each V (J) gene has an RSS at the 3' (5') end and each D gene is
128 flanked by RSSs on both ends. We henceforth refer to the 3' (5') flanking signal of a V (J) gene as
129 RSSV (RSSJ) and the 5' (3') flanking signals of a D gene as $RSSD_{left}$ ($RSSD_{right}$).

130 Since RSS motifs are highly conserved across all mammals, the human RSS motif (**Figure 1A**) can be
131 used for predicting RSS motifs in other mammalian species. IGDetective forms a *motif profile* from
132 all known RSSs in a reference genome and uses this profile to evaluate the *likelihood ratio* of an
133 arbitrary string from a target genome to decide whether this string represents a putative RSS (see
134 Methods). For each genomic position, it computes the likelihood ratio that there is an RSS flanking
135 this position and classifies a position as a candidate RSS if this ratio exceeds a *likelihood threshold*
136 (see Methods).

137 **IterativeIGDetective pipeline.** IterativeIGDetective predicts IG genes in the *target* genome by
138 leveraging the knowledge of IG genes in well-studied *reference* genomes (human, mouse, and cow).
139 Although the highly-repetitive IGH loci in a few other genomes have been semi-manually assembled
140 and annotated (e.g., IG loci in pig (Eguchi-Ogawa et al. 2010), goat (Du et al. 2018), and rabbit (Gertz
141 et al. 2013)), it remains unclear how accurate these short-read assemblies are since it is difficult to
142 assemble the IG loci even from long reads (Bankevich and Pevzner, 2020), let alone from short reads.
143 In the absence of accurate IG annotation tools, it is also unclear what are the false positive/false
144 negative rates of the manually predicted IG genes in these assemblies.

145 **Figure 1A** illustrates the IterativeIGDetective pipeline with emphasis on detecting V genes. It starts
146 by identifying a contig (or multiple contigs) containing the IGH locus in the target species and finding
147 candidate RSSs in this locus based on their similarity to the known RSSs in the reference species.
148 Afterward, it analyzes the genomic region flanked by the found RSSs to predict the IG genes
149 themselves. In the case of V and J genes (that are longer and more conserved than D genes), it

150 classifies a region preceding/following the identified RSS as a novel V/J gene if its similarity with a
151 known V/J gene exceeds a *similarity threshold* determined by percent identity and shared *k*-mers (See
152 Supplemental Method “Identification of candidate V and J genes” for parameters specifying similarity
153 thresholds).

154 IterativeIGDetective extends the *non-iterative mode* described above (that ends with the similarity-
155 based identification of V/J genes), by the *iterative mode* for identifying novel V/J genes whose
156 similarity with known V/J genes does not exceed the similarity threshold. The iterative mode extends
157 the set of known V/J genes by the newly identified V/J genes and uses this extended set to iteratively
158 repeat the non-iterative mode until no novel V/J genes are found (**Figure 1A**).

159 **The challenge of identifying highly-divergent IG genes.** Although IterativeIGDetective identifies
160 many candidate IG genes in target species, it is not capable of detecting distant IG genes that
161 significantly deviate from all canonical human IG genes, e.g., IG genes from a distant family that has
162 not been discovered yet. Reducing the similarity threshold in IterativeIGDetective increases its false
163 positive rate and does not necessarily lead to identifying distant IG genes.

164 We hypothesize that, similarly to IG genes in known families, highly-diverged IG genes from a novel
165 family should (i) display some degree of pairwise similarity to all other IG genes from the same
166 family, and (ii) be located within a relatively short region of the genome. Based on these two “IG
167 family” criteria, BlindIGDetective identifies novel IG genes which do not resemble known human IG
168 genes by constructing the *similarity graph* described below.

169 **Similarity graph.** Given a position *s* in a genome, a parameter *direction* (downstream “-” or upstream
170 “+”) and an integer *segment-length* (*L*), we define the *s*-fragment as the segment of length *L* either
171 downstream or upstream from this position depending on the parameter *direction*. We define the
172 *coding length* of an *s*-fragment as the length of the longest out of three *open reading frames* ending at
173 position *s* (in the case *direction*=”-”) or starting at position *s* (in the case *direction*=”+”). For
174 simplicity, below we assume that *direction*=”-”.

175 Given a set S of positions in a genome, parameters *direction*, and parameter *segment-length*, we
176 define an edge-weighted *similarity graph* (referred to as the *S-graph*) as follows. Each vertex in this
177 graph corresponds to a position in the set S and each edge connects *similar s*-fragments, where
178 similarity is established based on percent identity between *s*-fragments. All isolated vertices are
179 removed from the similarity graph. The edge-weight of an edge (s,s') is defined as the *percent identity*
180 between the *s*-fragment and the *s'*-fragment. Given a parameter *span* (default value 0.5 Mb), vertices
181 in the same connected component of the similarity graph are classified as *co-located* if they are
182 separated by less than *span* nucleotides in the genome.

183 Given a vertex v in a connected component, its *clump* is defined as the set of all vertices co-located
184 with v after removing all vertices whose percent identity with any other single vertex in the clump
185 does not exceed the similarity threshold (default value 70%). A vertex with a maximum-size clump in
186 a given connected component is called the *center* of this component (ties are resolved randomly). Two
187 clumps are *linked* if the distance between their center vertices does not exceed a distance threshold
188 (default value 1 Mbp). BlindIGDetective constructs *clusters* using single linkage clustering of linked
189 clumps and analyzes the constructed *large clusters* (of size larger than the default value *smallSize*=3)
190 as putative IG genes within a single IG locus. Clusters are further analyzed as candidates for new IG
191 families as they satisfy the “IG family” criteria specified above.

192 **BlindIGDetective pipeline.** Given a position-set S , parameter *direction*, and parameter *segment-*
193 *length*, BlindIGDetective constructs the *S-graph* (**Figure 1B**). Below, we limit attention to the case
194 when S is the set of starting positions of RSSVs predicted by IGDetective, *direction*= “-” and
195 *segment-length*= 350 bp. Since this setting models the search for V genes, we refer to the resulting *S-*
196 *graph* (*s*-fragments) as the *V-graph* (*v*-fragments). The default parameters *direction* and *segment-*
197 *length* will need to be modified for D-graphs and J-graphs since they depend on the position of the
198 gene (5' or 3' end) with respect to the RSS and the typical length of the gene. Since IG genes in
199 known genomes are located within relatively short regions (e.g., human IGHV genes are located
200 within 850 kbp long IGHV locus), a clump of co-located vertices within a connected component of

201 the V-graph may reveal a family of V genes in a newly sequenced genome. To generate such clumps
202 for each connected component, BlindIGDetective forms a clump of the center vertex in this
203 component, removes vertices of the constructed clump from this connected component, and iterates
204 until the component has no vertices left. It further combines clumps into clusters as described above.

205 We benchmark BlindIGDetective on the entire human genome and demonstrate that it reveals the vast
206 majority of human V (J) genes, for example, it finds 52 out of 57 known human IGHV genes without
207 any prior information about their sequences, while limiting possible false positives (with respect to
208 known human V genes) to 7. Some of these false positives may represent novel candidate V genes in
209 the human genome (see section “BlindIGDetective reveals novel candidate V genes in the human
210 genome”). Although BlindIGDetective missed a small number (5) of 57 human IGHV genes with
211 “weak” RSSs, these genes can be easily identified by slightly reducing the likelihood threshold with
212 follow-up similarity search against 52 identified IGHVs. In addition to finding IGHV genes,
213 BlindIGDetective finds V genes from other IG and TR loci and even *orphan* V genes on
214 Chromosomes 15 and 16 (**Supplemental Table S1**) that resulted from segmental duplications of the
215 human IGH locus (Nagaoka et al., 1994). We thus argue that applying this approach to any
216 mammalian species would reveal most V genes in this species, including highly-diverged V genes
217 missed by IterativeIGDetective as well as unusual genes that are affected by off-target VDJ
218 recombinations.

219 **[FIGURE 1]**

220 **Figure 1. IterativeIGDetective (A) and BlindIGDetective (B) pipelines.** (A) IterativeIGDetective iteratively
221 extends the set of identified IGHV genes. “Known V genes” box represents known V genes in a reference
222 genome. “RSSV motif” box represents a profile formed by the reference RSSVs for human V genes. (A1) After
223 identifying a contig containing the IGH locus in the target genome, IterativeIGDetective identifies candidate
224 RSSs for V genes in this contig based on similarities with the human RSS motif. A region preceding a true
225 positive RSS represents a V gene while a region preceding a false positive RSS does not. (A2) A region
226 preceding a candidate RSS is classified as a human-like V gene if its similarity with a known human V gene

227 exceeds a similarity threshold. (A3) Target-like V genes in the target genome are identified based on similarities
228 with human-like V genes detected in step (A2). (A3*) Target-like V genes are iteratively identified based on
229 previously detected target-like V genes, until no new genes are identified. (B) BlindIGDetective constructs the
230 V-graph, analyzes connected components in this graph, finds clumps of co-localized fragments in each
231 connected component of this graph, and combines the found clumps into clusters that represent candidate
232 families of V genes. (B1) Candidate RSSVs in the entire target genome are identified based on similarities with
233 the human RSSV motif for V genes formed by the reference human RSSVs (represented by the “RSSV motif”
234 box). (B2) The V-graph is constructed on the vertex-set of all RSSVs. Two vertices (RSSVs) in the V-graph are
235 connected by an edge if fragments preceding these RSSVs are similar. True (false) positive RSSVs form large
236 (small) connected components in the V-graph. (B3) Each connected component in the V-graph is partitioned
237 into clumps of co-located genes. (B4) Non-trivial clumps (containing multiple RSSs from the same connected
238 component and clustered within a short region of the genome) represent putative V genes within a putative IG
239 loci. Note that a vertex is not included in a clump if it is not similar to all other vertices in this clump (like light
240 green vertex in the rightmost clump). At the final step (not shown), BlindIGDetective combines the identified
241 clumps into clusters to reveal IG genes.

242 **Datasets.** We extracted known IG genes from the IMGT database (Lefranc et al., 2015) for the three
243 reference genomes and mapped them to the IGH locus of the same species as described in
244 **Supplemental Method “Annotating IG genes in reference genomes”, Supplemental Table S2.**
245 **Table 1** presents information about the mapped IG genes in the reference species that we refer to as
246 *canonical* IG genes. We also generated a *combined set* of IG genes by combining human, mouse, and
247 cow IG genes, thus adding one more “reference” to the human, mouse, and cow references. We refer
248 to the profile of all RSSs in this set (for each type of IG gene) as the *combined RSS* and denote it as
249 *RSS**.

250 We selected twenty target mammalian species for prediction of IG genes: three great ape species
251 (Kronenberg et al., 2018) and seventeen species assembled by the VGP consortium (Rhie et al., 2020)
252 that represent a wide range of biological orders (**Figure 2A, Supplemental Table S3**). For each target
253 species, we identified contigs containing fragments of the IGH loci (*IGH-contigs*) as described in the

254 **Supplemental Method “Identifying putative IGH-contigs in a genome assembly”, Supplemental**
255 **Table S4.**

256 **Benchmarking IterativeIGDetective.** We benchmark IterativeIGDetective on the IG loci of the
257 human, mouse, and cow genomes by assuming that one of them represents a reference genome and
258 one of the remaining ones represents a target genome. Since the IG genes and RSSs in these species
259 are known, we evaluate IterativeIGDetective based on (i) maximizing the number of the known RSSs
260 and IG genes predicted by IterativeIGDetective and (ii) minimizing the number of false RSSs and IG
261 genes predicted by IGDetective. To this end, we attempt to maximize the True Positive Rate (TPR)
262 while minimizing the False Discovery Rate (FDR).

263 We tabulated heptamer and nonamer likelihood thresholds (see Methods) (**Supplemental Table S5**)
264 and visualized the percentage of RSSs and heptamers/nonamers passing the human reference L_{min}
265 likelihood thresholds (**Supplemental Figure S1**). The vast majority of the heptamers and nonamers in
266 the genome have very low likelihood ratios.

267 Given the identified likelihood thresholds for each reference (human, mouse, cow, or combined), we
268 first launched IterativeIGDetective with these thresholds on the same species by considering them as
269 target species. We then tabulated the number of detected candidate signals in the human IGH locus (in
270 four cases that represent identifications of putative human RSSs using RSS profiles in human, cow,
271 mouse and combined) and computed true positives (candidate signals representing canonical signals),
272 false positives (candidate signals that are not canonical signals), and false negatives (canonical signals
273 that are not candidate signals) in **Supplemental Table S6**. Afterward, we tabulated the RSS detection
274 statistics of IGDetective on all combinations of four references and four targets in **Supplemental**
275 **Method “Extended benchmarking of IterativeIGDetective”, Supplemental Table S7**. Finally, we
276 extracted the V, D, and J genes from the candidate RSSs determined by launching
277 IterativeIGDetective in non-iterative mode (see section “Identification of candidate IG genes” in
278 Methods) with the RSS profile based on the combined reference. **Table 1** provides information about
279 the results of IterativeIGDetective on the reference species and using the combined RSS* profile.

V genes						
species	# of canonical genes	# of predicted genes	# of true positive genes	FDR	TPR	F1
human	70	57	50	0.12	0.71	0.78
cow	36	23	21	0.08	0.58	0.71
mouse	154	97	92	0.05	0.59	0.73
combined	260	177	162	0.08	0.62	0.74
D genes						
species	# of canonical genes	# of predicted genes	# of true positive genes	FDR	TPR	F1
human	25	17	15	0.11	0.6	0.71
cow	14	12	9	0.30	0.64	0.66
mouse	15	9	8	0.11	0.53	0.66
combined	54	38	32	0.17	0.59	0.68
J genes						
species	# of canonical genes	# of predicted genes	# of true positive genes	FDR	TPR	F1
human	6	6	5	0.16	0.83	0.83
cow	12	4	4	0	0.33	0.5
mouse	4	2	2	0	0.5	0.66
combined	22	12	11	0.08	0.5	0.64

280 **Table 1. Information about IG genes predicted by IterativeIGDetective based on the *combined* RSS***
281 **profile.** Rows refer to the species in which IterativeIGDetective predicts IG genes based on the RSS* profile.
282 “FDR”, “TPR” and “F1” columns represent false discovery rate, true positive rate, and F1 score, respectively
283 (see Methods). Since the same candidate D gene could potentially be reported twice on both the forward and
284 reverse strands, such a D gene is considered a true positive if either reported D gene’s start and end index
285 matches a reference gene’s start and end index. Some of the true positive predictions represent pseudogenes that
286 either have an in-frame stop codon or do not participate in VDJ recombination. We classify a detected gene as a
287 true positive if (i) its end index is the same as the corresponding reference gene’s end index, and (ii) its start
288 index is within 3 nucleotides towards the 5’ direction of the corresponding reference gene’s start index. This
289 ensures that the gene is predicted with an offset of at most +1 amino acid.

290 **Detecting IG genes in target genomes.** We applied IterativeIGDetective using the combined RSS*
291 profile to IGH-contigs of all target species (**Supplemental Method “Analysis of RSSs in reference
292 and target species”, Supplemental Table S8**). Since the identified IGH-contigs are usually longer
293 than the IGH loci, the predicted RSSs may include many false positives. For example, the number of
294 predicted RSSV candidates for a single species varies from 69 to 7027 with the median value 995
295 (**Supplemental Table S3**). However, further similarity-based filtering (described in Supplemental
296 Method “Identification of candidate V and J genes”) of regions flanking these candidate RSSVs
297 greatly reduces the number of false positive predictions, resulting in 3–64 V genes per species (the

298 median value is 34) (**Supplemental Table S3**). In total, IterativeIGDetective found 1021 candidate V
299 genes across twenty target species, including 50 target-like V genes (**Supplemental Table S9**). 34 out
300 of 50 target-like V genes share at least 80% percent identity with other V genes identified at the
301 previous iteration (**Supplemental Table S9**). After filtering candidates with stop codons in the open
302 reading frame, the number of candidate V genes was reduced to 581.

303 The number of predicted RSSJs varies from 54 to 3911 and from 4 to 21 (with the median value 8)
304 after similarity-based filtering, resulting in a total of 174 candidate J genes. After applying the
305 additional filters based on the conservation of the tryptophan-encoding TGG codon in the candidate J
306 genes (**Supplemental Method “Comparative analysis of IGHJ gene candidates”, Supplemental**
307 **Figure S2**), the number of candidate J genes was reduced to 60.

308

309 The number of predicted RSSDs varies from 1 to 17 with a median value of 4. IterativeIGDetective
310 identified a total of 137 candidate D genes which were extracted as short regions flanked by the
311 predicted RSSDs (without any filtering). After redefining the boundaries (see **Supplemental Method**
312 **“Computing boundaries of the IGH loci using predicted IG genes”**) of the IGH loci, we discarded
313 45 candidate D genes located outside these loci (to minimize the number of false positive D genes),
314 resulting in a set of 92 candidate D genes.

315

316 For each target species, we also found positions of *constant* (C) immunoglobulin genes in its
317 assembly by aligning highly-conserved human IGHC genes using Bowtie2 (Langmead and Salzberg,
318 2012). The number of IGHC genes per species varies from 2 to 19 with the median value 7, resulting
319 in 149 IGHC gene candidates.

320

321 **The IGH loci widely vary in length across mammalian species.** We analyzed the positions of the
322 candidate V, D, J, and C genes within the assembled genome in order to identify the boundaries of the
323 IGH loci and their lengths assuming the standard V→D→J→C ordering (see Methods). Long repeats

324 within the IGHV locus often break its assembly into multiple contigs, with one of the contigs
325 containing the first *i* V genes (referred to as the *IGH-start*) and another contig (referred to as an *IGH-*
326 *end*) containing all (or some of) the remaining V genes as well as D, J, and C genes (**Figure 2B**). See
327 **Supplemental Method “Unusual IGH locus in the sloth genome”, Supplemental Figure S3.**

328 The sloth and the spear-nosed bat have the longest IGH loci among analyzed species (6.7 and 4.6
329 Mbp, respectively), while aquatic animals (vaquita, blue whale, platypus, sea lion, and dolphin) have
330 the shortest IGH loci, varying from 311 kbp for the vaquita to 607 kbp for the dolphin.

331 Our estimate of the length of the platypus IGH locus (457 kbp) is higher than the previous estimate by
332 Gambon-Deza et al., 2009 (271 kbp). The analysis of the IGH-start contig (containing V genes only)
333 and IGH-end contig (containing V, D, J and C genes) in platypus revealed an unusual feature. It
334 turned out that, the IGH-end contig contains the entire IGH locus in platypus, while the IGH-start
335 contig contains V genes from the *TCR μ* locus, a unique T-cell receptor locus found only in marsupials
336 and monotremes (Miller, 2010). Since previous studies demonstrated that V genes from this locus are
337 more similar to immunoglobulin V genes than TCR V genes (Miller, 2010), this finding illustrates
338 that IGDetective is capable of detecting unusual TCR genes. Nevertheless, to limit analysis to the IGH
339 loci only, we discarded V genes from the platypus *TCR μ* locus from further analyses.

340 **[FIGURE 2]**

341 **Figure 2. Information about the IGH loci in twenty target mammalian species.** (A) The phylogenetic tree
342 formed by twenty target and three reference species. The tree was subsampled from the Tree of Life constructed
343 in Hedges et al., 2015. Each species is shown by its common name and a VGP identifier specified in the
344 parenthesis if available. Each species is encoded by a unique color (left vertical color panel) and a color
345 representing its order (right vertical color panel). The list of orders is shown in the upper left corner. Here and
346 below visualization was performed using the BioRender and Iroki (Moore et al., 2020) tools. (B) Information
347 about the IGH loci of twenty target species. Each line corresponds to a target species and shows fragments of
348 the IGH locus with positions of candidate V (blue), D (orange), J (green), and C (red) genes. For six out of
349 twenty species, the IGH-end covers less than 80% of the predicted IGH locus length (lynx, blue whale, mastiff

350 bat, horseshoe bat, chimpanzee, and grey squirrel). We showed both the IGH-start and the IGH-end for these six
351 species and only the IGH-end for the remaining species. For a better visualization, all IGH loci are shown as
352 having the same length that does not reflect their real lengths (the bar plots next to the IGH loci show the
353 predicted lengths). The map on the right shows the counts of the productive V, D, J, and C genes identified in
354 the IGH locus. A C gene is classified as productive if its translated regions (defined by the closest human C
355 gene) does not contain stop codons. Non-zero counts are shown in green. The green triangles indicate partially
356 found and (likely) partially missing D and J genes in the sloth IGH locus. Although IterativeIGDetective did not
357 identify any J genes in two species (chimpanzee and spear-nosed bat) within the boundaries of the IGH loci, it
358 found a highly conserved candidate J gene in a short contig in the chimpanzee assembly (denoted as 1* in the
359 map of the right). IterativeIGDetective did not identify any candidate J genes in the spear-nosed bat, presumably
360 because all its RSSJs did not pass the likelihood threshold.

361 **Comparative analysis reveals highly-similar V genes in evolutionary distant species.** We
362 combined the candidate V genes (referred to simply as V genes) across all target species with known
363 V genes in reference species and constructed a phylogenetic tree on their amino acid sequences using
364 Clustal Omega (Sievers et al., 2011). **Figure 3A** shows the computed tree where leaves (representing
365 V genes) are colored according to the species they belong to and the order of the species. “Cutting”
366 this tree by a horizontal line at a *height threshold* results in subtrees formed by clusters of similar V
367 genes. We classified a cluster as *large* if it contains more than 5 V genes and as *multi-species* (multi-
368 order) if it includes V genes from multiple species (orders). Since large multi-species clusters are the
369 main focus of the comparative analysis, we selected the height threshold 1.12 maximizing the number
370 of these clusters. The resulting 219 clusters include 43 large clusters, 25 large multi-species clusters,
371 and 7 large multi-order clusters (**Supplemental Figure S4**).

372 For each large multi-species cluster formed by V genes g_1, \dots, g_n from species s_1, \dots, s_m , we computed its
373 *gene distance* and *species distance*. The gene distance is computed as the $\max\{GeneDist(g_i, g_j)\}$ for
374 all pairs of genes g_1, \dots, g_n , where $GeneDist(x, y)$ represents the fraction of non-matching positions in
375 the alignment between genes x and y . The species distance is computed as $\max\{SpeciesDist(s_i, s_j)\}$ for
376 all pairs of species s_1, \dots, s_m , where $SpeciesDist(x, y)$ is the distance between species x and y in the tree

377 shown in **Figure 2A**. **Figure 3B** illustrates correlation between species distances and gene distances
 378 for 25 large multi-species clusters (Pearson's correlation $r=0.77$, $P\text{-value}=7.15\times 10^{-6}$). Seven multi-
 379 order clusters (referred to as C1-C7) are represented in blue in **Figure 3B**, and, for all of them but one
 380 (cluster C7 shown as a blue asterisk), the gene distance is either well-estimated or under-estimated by
 381 the regression line. The unusual C7 cluster is formed by V genes from six target species (chimp,
 382 gorilla, orangutan, marmoset, red squirrel, and horseshoe bat) that are similar to human genes IGHV3-
 383 30, IGHV3-30-3, IGHV3-30-5, and IGHV3-33 (**Figure 3C**). This cluster reveals a surprising
 384 conservation of V genes across distant species, e.g., the amino acid sequence of the gene HS_bat_58
 385 in the horseshoe bat has even fewer differences (8) with the human V genes in this cluster than some
 386 chimpanzee V genes (9). We thus conjecture that V genes in this cluster are subjected to the same
 387 selective pressure, e.g., driven by common pathogens that are faced by the species in this cluster.

388

[FIGURE 3]

389 **Figure 3. Comparative analysis of mammalian IGHV genes.** (A) A phylogenetic tree of IGHV genes in
 390 twenty target (581 V genes) and three reference (310 V genes) species. Edges corresponding to clusters C1–C7
 391 described in (B) are shown in blue. The scale is shown on the left. The upper and lower horizontal bars show
 392 colors of species and their orders, respectively. List of species and their colors are specified below the tree
 393 (species from the same order are shown by a colored vertical bar on the left). (B) The plot on the left shows the
 394 species distance (x -axis) and the gene distance (y -axis) for 25 large multi-species clusters of V genes. Red and
 395 blue dots correspond to single-order and multi-order clusters, respectively. The linear regression line is shown in
 396 grey. The Pearson's correlation (r) and P -value (P) are shown on the top of the plot. Each of seven multi-order
 397 clusters C1–C7 is represented as a pie-chart on the right. An inner (outer) wedge in each pie-chart corresponds
 398 to a species (an order) and the wedge size is proportional to the number of V genes it contains. (C) Multiple
 399 alignment of 21 V genes from the cluster C7. Four human V genes from this cluster are shown on the top. Non-
 400 human genes are denoted according to the short names of species, and "HS_bat" refers to the horseshoe bat. A
 401 position in a non-human V gene is shown as "." if the amino acid at this position matches the corresponding
 402 amino acid in one of four human V genes and by the corresponding amino acid otherwise. Red rectangles show
 403 positions of CDR1 and CDR2 according to the IMGT notation. Green bars show positions of two conserved
 404 cysteines (one located close to the start of CDR1 and another located close to the end of the V gene).

405 **A new family of cysteine-rich IGHV genes.** 254 out of 309 canonical human V genes (82%) listed
406 in the IMGT database with productive amino acid sequences (including allelic variants) have two
407 *canonical cysteines* located at conserved positions (**Figure 3C**). We classify a V gene as *cysteine-rich*
408 if it contains four or more cysteines and analyze cysteine-rich V genes among all V genes shown in
409 **Figure 3A**. There are no human cysteine-rich V genes and only 2 (1 mouse (cow) cysteine-rich V
410 genes. It remains unclear whether the mouse and cow cysteine-rich V genes represent pseudogenes
411 rather than functional genes, e.g., the cysteine-rich V gene in cow does not contribute to antibody
412 repertoires (Safonova et al., 2022).

413 715 out of 891 V genes (80%) from both reference and target species have 2 cysteines and only 61
414 (7%) are cysteine-rich. Cysteine-rich V genes are grouped together in the phylogenetic tree (shown in
415 dark green in **Figure 4A**) and appear only in 2 out of 25 identified large multi-species clusters,
416 including a multi-order cluster C4 (**Figure 3B**) and a single-order cluster that we refer to as C*.
417 Cluster C4 contains 27 V genes from six species: dolphin, blue whale, sloth, horseshoe bat, spear-
418 nosed bat, and mastiff bat (**Figure 4B**). Cluster C* consists of 12 V genes from grey and red squirrels
419 (**Figure 4B**). 25 out of 27 V genes in cluster C4 are cysteine-rich and 11 out of 12 V genes in cluster
420 C* are cysteine-rich.

421 **Figure 4C** shows that, in addition to two canonical cysteines at conserved positions, most V genes
422 from clusters C4 and C* have two other cysteines, also at conserved positions (one cysteine in CDR1
423 and another in CDR2). While several antibodies with a single cysteine in either CDR1 or CDR2 were
424 reported before (Wu et al, 2012; Prabakaran and Chowdhury, 2020), antibodies with cysteines in both
425 CDR1 and CDR2 have not been reported yet. Since cysteines form disulfide bonds, we hypothesize
426 that cysteine-rich V genes might generate unusual antibodies with non-canonical conformations and
427 could potentially form a unique part of bat immunity against the great variety of viruses they host.

428 The cysteine-rich cluster C4 includes a sloth V gene (denoted as “sloth_38” in **Figure 4C**) with the
429 unusual 6 aa long suffix CVLLCE classified as the beginning of CDR3. The vast majority of known
430 V genes have the conserved CAR suffix and thus contribute to at most three first amino acids of

431 CDR3s. Two known exceptions from this rule are the cattle IGHV1-7 gene that contributes to
 432 ultralong antibodies in combination with ultralong D gene IGHD8-2 (Wang et al., 2013; Safonova et
 433 al., 2022) and the platypus IGHV1-20 gene with unknown functions (Gambon-Deza et al., 2009). We
 434 hypothesize that, similar to the cattle IGHV1-7 gene contributing to generation of ultralong
 435 antibodies, both sloth mChoDid1_38 and platypus IGHV1-20 V genes with long CDR3 prefixes may
 436 generate (alone or in combination with D genes) antibodies with non-canonical structural features.
 437 Below we analyze unusual candidate D genes in these species and perform comparative analysis of all
 438 detected D genes.

439 **[FIGURE 4]**

440 **Figure 4. Two novel cysteine-rich clusters of mammalian IGHV genes.** (A) A phylogenetic tree of IGHV
 441 genes colored according to the number of cysteines in their amino acid sequences. (B) Only 2 out of 25 large
 442 multi-species clusters of V genes contain cysteine-rich V genes. The description of pie plots is provided in the
 443 caption to **Figure 3B**. The number within the parenthesis next to the species name indicates the number of V
 444 genes from this species. (C) Multiple alignment of genes from two clusters shown in (B). Three non-cysteine-
 445 rich V genes are marked with a grey circle on the left. The gene on top of each cluster is chosen as the sequence
 446 of a V gene with the minimum average distance from other genes in the cluster. Green (purple) bars show
 447 positions of canonical (non-canonical) cysteines. Some proteins contain cysteines (that are shown in purple)
 448 outside these positions. “HS_bat”, “M_bat”, and “SN_bat” in the top alignment refer to V genes of the horse-
 449 shoe bat, the mastiff bat, and spear-nosed bat, respectively. “Red_sq” and “grey_sq” in the bottom alignment
 450 refer to V genes of the red squirrel and the grey squirrel, respectively.

451 **Comparative analysis of mammalian D genes.** IGDetective identified 92 candidate D genes in
 452 target species (**Figure 5A**). For the comparative analysis, we combined these D gene candidates with
 453 81 known D genes of three reference species (27 human, 31 mouse, and 23 cow D genes), resulting in
 454 a set of 173 D genes. For the sake of simplicity, we refer to both reference and candidate D genes as
 455 simply D genes.

456 **Figure 5B** shows that only five D genes are shared among two or more species. Since, in difference
 457 from V and J genes, D genes are short and very diverse, it remains unclear whether there exist specific

458 features of D genes that are shared among nearly all mammalian species. To reveal such features, we
459 searched for common substrings in all 173 D genes. We translated each D gene into three reading
460 frames and extracted all its 4-mers in the amino acid alphabet. In total, we collected 738 4-mers, with
461 128 of them appearing in multiple species. The maximum number of species represented by a single
462 4-mer is 6.

463 We constructed the *Hamming graph* on 128 shared 4-mers by connecting two 4-mers by an edge if
464 they differ in a single amino acid (Safonova et al., 2015). It turned out that all connected components
465 in this graph are small (less than 5 vertices) with exception of two components consisting of 43 and
466 42 4-mers and covering 68% of all 4-mers appearing in multiple species (**Figure 5C**). These two
467 components cover all highly abundant 4-mers (4-mers that are present in 3–6 species). Below we
468 focus on the first component since the second component represents the same substrings of D genes as
469 the first component but translated in a different reading frame.

470 The 4-mers in the first component represents 16 out 23 analyzed species and are mostly formed by
471 amino acids G, S, and Y. We refer to D genes that encode these 4-mers as *G/S/Y-rich* D genes. A half
472 of all possible single-nucleotide mutations of cysteine-encoding codons (TGT and TGC) result in
473 codons GGT, AGT, TAT, and TCT encoding amino acids G, S, S, and Y, respectively. These three
474 amino acids are extremely frequent in the longest cattle D gene IGHD8-2 (**Figure 5C**) where they
475 play a special *cysteine-triggering* role in ultralong cattle antibodies: somatic hypermutations create
476 new cysteines from these amino acids, forming new disulfide bonds, and reshaping the resulting
477 antibody (Wang et al., 2013). We conjecture that the *G/S/Y-rich* D genes may play a similar cysteine-
478 triggering role as the IGHD8-2 gene in cows.

479 **Unusual cysteine-rich D genes in the platypus genome.** Platypus and sloth genomes have two V
480 genes with long non-canonical suffixes that represent the beginnings of CDR3s: the platypus gene
481 IGHV1-20 with suffix LAAELLYCR and the sloth gene “sloth_38” with suffix CVLLCE
482 (**Supplemental Figure S5**). The only other known V gene with a long non-canonical suffix is the cow

483 gene IGHV1-7 that plays a special role in generating ultralong CDR3s by recombining with the
484 longest known D gene (IGHD8-2) of length 148 nt (Wang et al., 2013, Safonova et al., 2022). This
485 long D gene is highly unusual: all but three amino acids in its translation are either cysteines or
486 cysteine-triggering amino acids G, S, and Y (**Figure 5C**). We thus searched for D genes with similar
487 properties in the platypus and sloth IGH loci. However, IGDetective, which uses rather stringent
488 parameters for RSS search, did not report any unusual (long or C/G/S/Y-rich) D genes. In contrast to
489 IGDetective, SEARCH-D (Safonova and Pevzner, 2020) uses relaxed parameters for finding RSSs at
490 the cost of reporting more false positive D gene candidates.

491 We thus launched SEARCH-D on the platypus (457 kbp long) and sloth (6.7 Mbp long) IGH loci.
492 SEARCH-D reported 45 and 76 D gene candidates (simply referred to as D genes) for the platypus
493 and the sloth, respectively (including 6 platypus and 1 sloth D gene candidates reported by
494 IGDetective). Since the candidate D genes in the sloth are scattered through the entire IGH locus, we
495 were unable to identify the location of the sloth IGHD locus. (**Supplemental Figure S6**). In contrast,
496 29 out of 45 candidate D genes in platypus form a dense 60 kbp long cluster pointing to a previously
497 unknown location of the IGHD locus (**Figure 5D,E**).

498 Similarly to other mammalian IGHD loci (Safonova and Pevzner, 2020), the identified 60 kb long
499 fragment harboring 29 candidate D genes in platypus is a tandem repeat (**Figure 5D,E**), reinforcing
500 the conclusion that this region indeed represents the IGHD locus. We classify a D gene as *cysteine-*
501 *rich* if it contains at least two cysteines in one of its reading frames (only 3 out of 25 human D genes
502 are cysteine-rich). Clustering 29 candidate D genes in platypus revealed 4 groups of similar D genes
503 with percent identity $\geq 70\%$ (referred as D1–D4) that include many cysteine-rich D genes: 12 out of 15
504 D genes in these groups are cysteine-rich and all genes in these groups, similarly to the cow D gene
505 IGH8-2, have many cysteine-triggering amino acids (**Figure 5E**).

506 Even though these 15 candidate D genes have rather diverged RSSs (**Figure 5F**), the high level of
507 their sequence conservation indicates that they are likely functional. We assume that, similarly to

508 human IGHD2 genes, these D genes can be responsible for generating antibodies with a disulfide
 509 bond inside CDR3s (Prabakaran and Chowdhury, 2020). The high number of these D genes suggests
 510 that the fraction of such antibodies is likely higher in the platypus repertoires as compared to the
 511 human repertoires. This finding agrees with a study by Johansson et al., 2002 (that reported an
 512 unusually high percentage of cysteine-rich antibodies in platypus antibody repertoires) and extends it
 513 by revealing the germline D genes contributing to the cysteine-rich antibodies. Further analysis of
 514 platypus Rep-Seq data will help to determine if these cysteine-rich D genes are recombined with
 515 IGHV1-20 (with unusual suffix LAAELLYCR) and shed light on their role in antibody repertoires.

516

[FIGURE 5]

517 **Figure 5. Comparative analysis of D genes.** (A) The distribution of the counts and lengths of D genes for 20
 518 target species. (B) D genes shared among two or more reference and target species. Green cells show species
 519 containing the corresponding D gene candidates. Species from left to right: chimp, gorilla, human, orangutan,
 520 mastiff bat, horseshoe bat, and otter. (C) The largest connected component of the Hamming graph on amino
 521 acid 4-mers of D genes. The component is shown by the subgraph of the Hamming graph (left subpanel) and the
 522 amino acid content at each position of the 4-mer (right subpanel). Vertices of the Hamming graph are colored
 523 according to the number of species they represent: from 2 (pale green) to 6 (dark green). The amino acid
 524 sequence of the G/S/Y-rich cow D gene IGHD8-2 is shown on the bottom of the right subpanel. Panels D–F
 525 illustrate the analysis of D genes in the platypus genome. (D) Positions of D genes detected by SEARCH-D in
 526 the platypus IGH locus. D genes are colored according to their lengths: 50 nt or less (purple), from 51 to 100 nt
 527 (green), and from 51 to 150 nt (orange). (E) The dotplot on the left shows the alignment of the $\cong 60$ kbp long
 528 platypus IGHD locus against itself. Positions and sequences of genes from four D gene families with two
 529 cysteines are shown on the right. (F) Motif logos of $RSSD_{left}$ heptamer (L7), $RSSD_{left}$ nonamer (L9), $RSSD_{right}$
 530 heptamer (R7), $RSSD_{right}$ nonamer (R9) for families D1–D4. Positions that do not match nucleotides in the
 531 consensus RSSs computed using the combined references are highlighted in grey. Consensus RSSs for the
 532 combined reference are shown in **Supplemental Figure S7**.

533 **Benchmarking BlindIGDetective.** BlindIGDetective constructed the human V-graph from 28394
 534 sites in the human genome that passed the RSSV likelihood threshold. A connected component is

535 classified as either small (of size at most *smallSize*), large (of size larger *smallSize* but smaller than
536 *giantSize*), or giant (of size at least *giantSize*) for default values *smallSize*=3 and *giantSize*=500. The
537 vast majority of candidate RSSs in the human genome are false RSSs that often represent isolated
538 vertices or vertices of giant components that likely originated from spurious RSSs within repeated
539 regions. Indeed, the human V-graph contains 6942 isolated vertices and one giant component on
540 20768 vertices. The vast majority of vertices in the giant component represent spurious repeats that
541 we have ignored in further analysis. **Supplemental Table S10** provides information about the V-
542 graphs for three reference species (human, mouse, and cow) and two selected target species (the
543 spear-nosed bat and the horseshoe bat).

544 For each vertex in the V-graph, we define the *percent identity*, *coding length*, *annotation index*, and
545 *conservation index* (see **Supplemental Methods “Analyzing connected components in the
546 similarity graph” and “Speeding-up BlindIGDetective”** for details). The conservation index of a
547 vertex is defined as the percent identity between its *v*-fragment and the closest *predicted gene*
548 (predicted genes could either be canonical genes from reference species or candidate genes detected
549 by IterativeIGDetective). A vertex is *annotated* if its conservation is at least $PI_{\text{annotation}}$ (default value =
550 90%). For annotating human (cow, mouse) vertices, we define conservation with respect to human
551 (cow, mouse) canonical V genes. However, for annotating the target species, we use either a
552 conservation threshold of $PI_{\text{annotation}} = 90\%$ with respect to IG genes in this species predicted by
553 IterativeIGDetective or a conservation threshold of $PI_{\text{annotation}} = 80\%$ with respect to human canonical
554 V genes.

555 A vertex in the V-graph is classified as *accordant* if its coding length exceeds the *minimum coding
556 length* threshold (default value *minCL*=200 bp). Since clusters with short coding lengths are likely
557 formed by spurious RSS (for reference species, all V genes, except one, have length exceeding 208
558 bp), below we focus on *accordant* clusters (clumps) defined as clusters (clumps) with coding lengths
559 exceeding the *minCL* threshold. We note that *accordant* clusters (clumps) may contain both *accordant*
560 and non-*accordant* vertices.

561 We launched BlindIGDetective with the RSS profile corresponding to a 23 nt spacer. This setting is
562 aimed at finding V genes in IGH, IGL, TRA, TRB, TRD or TRG loci (that all have RSSs with 23 nt
563 spacer) but not the IGK locus (since RSSs in this locus have a 12 nt rather than a 23 nt long spacer).
564 However, the IGK locus can be easily identified by simply changing a spacer length from 23 to 12 nt
565 in BlindIGDetective.

566 An annotated vertex in a cluster is said to be part of a specific V gene locus (IGH, IGL, TRA, TRB, or
567 TRG) defined by its most similar (as measured by percent identity) annotated V gene, where the set of
568 annotated V genes can be sourced either from a canonical set of V genes (if available) or predictions
569 made by IterativeIGDetective. A cluster is classified as *annotated* by a specific V gene locus if all its
570 annotated vertices are assigned to this locus.

571 **BlindIGDetective reveals novel candidate V genes in the human genome.** BlindIGDetective
572 identified 79 clumps from 179 connected components in the human V-graph on 684 vertices (after
573 removing isolated vertices and the giant component). It further removed clumps with very small spans
574 (below 1 kb) as such clumps are typically formed by multiple candidate (likely spurious) RSSs
575 located within a short region. Afterward, it combined the remaining clumps into clusters as described
576 in the subsection “The similarity graph”, resulting in 11 (8) large (accordant) clusters in the human
577 genome (**Table 2**). 7 out of 8 accordant clusters revealed seven known loci of human IG genes (IGH
578 and IGL loci are represented by the largest clusters of size 59 and 31, respectively).

579 **Table 2** shows that the cluster representing human IGH (IGL) genes is formed by 5 (5) clumps,
580 including 2 (1) unannotated clumps (**Supplemental Table S11**). There are only 2 (7) unannotated
581 fragments contained in IGH (IGL) annotated clumps, including 1 (3) accordant *v*-fragments.
582 Moreover, unannotated fragments within annotated clumps still show high median conservation of
583 83% for both IGH and IGL clusters.

584 Since 2 (7) unannotated fragments in IGH (IGL) annotated clumps have large median conservation,
585 we launched an IgBLAST search on them and revealed significant hits with E-values of at most $4e-76$

586 (1e-60) and percent identities of at least 84 (75) % against human V genes IGHV3-48 and IGHV1-46
587 (for 2 unannotated fragments in IGHV clumps) and IGLV1-44, IGLV3-22, IGLV3-21, IGLV3-9,
588 IGLV3-31 and IGLV7-46 genes (for 5 unannotated fragments in IGLV clumps). Alignment of 1 (3)
589 out of 2 (7) unannotated and accordant IGH (IGL) ν -fragments from annotated clusters revealed that
590 they align with their closest human V gene in a reading frame containing a stop codon. We therefore
591 suggest that these accordant unannotated fragments could represent previously undiscovered V genes
592 (or pseudogenes) that may be affected by RAG proteins during VDJ recombination.

593 The unannotated clumps representing human IGH (IGL) genes accounted for 5 (2) unannotated ν -
594 fragments in the IGH (IGL) cluster and included one accordant ν -fragment in the IGH locus (coding
595 length 597 nt) with a low percent identity with known human genes (under 62%). Unannotated
596 clumps retained a low median conservation under 62% in both IGH and IGL clusters. A similar
597 IgBLAST search of these the 5 (2) unannotated ν -fragments in the IGH cluster revealed hits for only
598 three ν -fragments against human IGHV genes IGHV3-7, IGHV3-11, and IGHV3-21 with low E-
599 values of 8×10^{-13} , 2×10^{-18} , and 7×10^{-20} and percent identities of 61%, 62%, and 63%, respectively.
600 Although these three ν -fragments share some (albeit low) similarity with known human V genes, they
601 are missing in the IMGT database. Moreover, all these ν -fragments are located in the IGH locus
602 within a short distance from the canonical IGHV3 gene (at distance 25 kbp, 11 kbp, and 61 kbp,
603 respectively). We therefore suggest that they represent distant IGHV3 genes missed by earlier
604 methods for annotating V genes.

605 The remaining 2 (2) ν -fragments in the IGH (IGL) loci had high E-values exceeding 0.42.
606 **Supplemental Figure S8** shows two alignments between the two pairs of these ν -fragments (that
607 extends through the entire sequence) and suggests these 2+2 ν -fragments could represent
608 undiscovered genes (or pseudogenes) that are not similar to known human V genes and therefore
609 would not have been discovered through V gene finding methods reliant on similarity with previously
610 identified genes.

611 See **Supplemental Method “BlindIGDetective results on cow and mouse genomes”** and
 612 **Supplemental Table 12** for details of BlindIGDetective benchmarking on other reference species.

Human														
cluster ID	cluster size	# clumps / [L.C]	PI	coding length (nt)	cluster density (%)	cluster span (Mb)	center vertex		AI/AI80	conservation wrt species genes		conservation wrt human genes		locus
							chr	coordinate (Mb)		min	med	min	med	
H0	59	5/27	91	342	29	0.93	14	106.34	0.88/0.92	56	100	56	100	IGH
H1	31	5/21	89	348	18	0.85	22	22.43	0.71/0.84	55	100	55	100	IGL
H2	8	4/2	100	201	14	0.58	15	21.72	0.5/0.5	54	83	54	83	IGH*15
H3	8	3/4	99	345	57	1.93	16	33.83	0.75/0.75	53	100	53	100	IGH*16
H4	6	1/6	95	435	100	0.03	7	38.36	0.83/1	88	100	88	100	TRG
H5	6	2/4	83	360	47	0.18	14	21.90	1/1	100	100	100	100	TRA
H6	4	2/2	88	330	33	0.12	7	142.50	0.75/1	88	100	88	100	TRB
H7	4	2/2	87	252	33	0.28	17	3.22	0/0	54	56	54	56	
Spear-nosed bat														
cluster ID	cluster size	# clumps / [L.C]	PI	coding length (nt)	cluster density (%)	cluster span (Mb)	center vertex		AI/AI80	conservation wrt species genes		conservation wrt human genes		locus
							contig	coordinate (Mb)		min	med	min	med	
PD0	130	12/38	94	357	19	1.672	S13	59.11	0/0.67	53	56	54	82	IGL
PD1	56	5/27	90	186	21	0.834	S16	0.97	0.38/0.46	54	86	55	79	IGH
PD2	50	4/27	94	315	30	0.762	S16	2.33	0.66/0.64	68	95	65	81	IGH
PD3	35	9/8	96	327	16	1.401	S3	145.30	0/0.43	53	55	73	80	TRA
PD4	22	3/8	84	192	19	0.387	MS7	0.38	0.86/0.5	72	100	62	80	IGH
PD5	14	3/8	84	195	30	0.163	MS2	0.08	0.36/0.5	53	88	55	81	IGH
PD6	14	3/5	92	177	23	0.127	MS12	0.06	0.79/0.5	76	99	67	80	IGH
PD7	13	3/8	89	126	15	0.206	MS33	0.06	0.92/0.46	84	96	69	72	IGH
PD8	12	4/5	90	369	23	0.052	MS7	0.01	0/0.67	55	56	56	83	IGL
PD9	10	1/10	93	378	49	0.167	S11	85.66	0/0.1	55	57	72	76	TRB
PD10	9	2/5	79	318	31	0.15	S2	133.59	0.89/0.44	85	100	66	79	IGH
PD11	7	2/5	81	177	52	0.283	S5	85.09	0.71/0.5	54	100	54	83	IGH
PD12	4	1/4	83	210	100	0.048	S3	80.22	0.75/1	90	91	84	87	IGH
PD13	7	2/4	99	222	43	0.629	S14	10.95	0/0	52	53	53	54	
PD14	5	1/5	87	2286	100	0.002	S4	208.06	0/0	51	51	52	53	

613 **Table 2. Information about large clusters derived from the human and spear-nosed bat genomes.**

614 BlindIGDetective constructed 14 large clusters (8 accordant clusters) in the human genome (only accordant
 615 clusters are shown). “L.C” represents the size of the largest clump in the cluster. Annotation index and
 616 annotation80 index are abbreviated to “AI ” and “AI80”. Annotated clusters are highlighted in blue. Clusters are
 617 ordered in the decreasing orders of their sizes. Conservation is shown with respect to predicted V genes from the
 618 same species as well as canonical human V genes, with annotation index (annotation80 index) defined with
 619 respect to the former (latter). Predicted V genes are the canonical V genes for human and are the candidate V
 620 genes predicted by IterativeIGDetective for spear-nosed bat. The locus column classifies the annotated clusters

621 as one of the families of human IG or TCR genes described in **Supplemental Table S1** (highlighted in blue).
622 All annotated human clusters, except for H2, have coding length greater than 315 nt, consistent with the range of
623 coding lengths in known V genes. IGH orphans on Chromosomes 15 and 16 are noted as IGH*15 and IGH*16.
624 Human cluster H2 with coding length 201 nt (locus IGH*15) has shorter coding length because they contain
625 many short pseudogenes. The table also shows all 13 large annotated spear-nosed bat clusters, 7 of which are
626 accordant. Only accordant unannotated spear-nosed bat clusters are shown. In the spear-nosed bat's "contig"
627 column, "mPhyDis1_scaffold_<N>" in the VGP assembly version "mPhyDis1.pri.cur.20200504" is shortened
628 to "MS<N>".

629 **BlindIGDetective reveals highly diverged candidate V genes in the spear-nosed bat genome.**

630 Bratsch et al., 2011 and Schountz et al., 2017 formulated the "bat IG diversity" hypothesis stating that
631 bat's ability to carry many disease-causing pathogens (while themselves being unaffected) could be
632 linked to their large and diverse immunoglobulin gene-set. However, since assembly of IG loci in any
633 species is challenging (Bankevich and Pevzner, 2020), accurately assembled (let alone, annotated) IG
634 loci in bats remained unavailable until recently. In fact, the only support for the claim that bats have a
635 very large number of IG genes comes from a probabilistic model rather than an annotated IG loci in
636 an assembled bat genome: Bratsch et al., 2011 used this model to predict $\cong 240$ V genes in the little
637 brown bat without having access to its genome. IterativeIGDetective results do not support the "bat
638 IG diversity" hypothesis": it reported from only 9, 10, 32 and 63 IGHV genes across four bat species.
639 We thus applied BlindIGDetective to reveal divergent V genes that IterativeIGDetective may have
640 missed. We focused on the spear-nosed bat (only 34 IGHV genes reported by IterativeIGDetective) as
641 the bat species with the longest IGH locus (4.6 Mbp).

642 BlindIGDetective constructed 72 clusters in the spear-nosed bat genome, including 17 (29) accordant
643 (large) clusters (**Table 2**). Vertices in these clusters were annotated using the candidate IGHV genes
644 predicted by IterativeIGDetective and canonical human genes to annotate V genes from IGL, TRA,
645 TRB, and TRG loci. A total of 13 large, annotated clusters were observed in the spear-nosed bat
646 genome.

647 BlindIGDetective identified 9 large IGH clusters (3 of which are accordant) encompassing 27 clumps
648 and 189 ν -fragments. It also identified 2 large IGL clusters with 142 ν -fragments, 1 TRA cluster with
649 35 ν -fragments, and 1 TRB cluster with 10 ν -fragments (all these clusters are accordant). In addition
650 to the large clusters, it identified a small IGH cluster containing 2 ν -fragments. No small clusters were
651 detected for IGL, TRA, and TRB loci. For each ν -fragment, the opening reading frame and the start
652 position of the gene were computed using alignment to the closest human germline V gene. A ν -
653 fragment is classified as *productive* if it has an open reading frame that begins at its start position and
654 terminates at its end position. 147 IGH, 95 IGL, 15 TRA, and 1 TRB annotated ν -fragments were
655 found within the identified clusters; 29 IGH, 58 IGL, 6 TRA, and 0 TRB annotated ν -fragments were
656 productive and thus classified as V gene candidates.

657 BlindIGDetective also identified 44, 47, 20, and 9 unannotated ν -fragments within the IGH, IGL,
658 TRA, and TRB clusters, respectively (including 1 IGH, 20 IGL, 4 TRA, and 7 TRB productive ν -
659 fragments). The alignments of these productive unannotated ν -fragments against the translated human
660 IG and TCR genes revealed high percent identity (in amino acids) varying from 49% to 77% with the
661 average value 68%. To analyze the origin of unannotated candidate V genes, we combined them with
662 productive annotated gene candidates and constructed phylogenetic trees using Clustal Omega
663 (Sievers et al., 2011) for IGH, IGL, and TRA loci. Annotated and unannotated V gene candidates are
664 interspersed in trees for IGH and IGL loci, indicating that unannotated V gene candidates likely
665 represent highly diverged members of the canonical V gene families (**Supplemental Figure S9**). In
666 the TRA locus, unannotated V gene candidates form an independent subtree suggesting that they
667 represent an unknown V gene family.

668 **Prediction of V genes in bats does not support the “bat IG diversity” hypothesis.** Our
669 benchmarking of BlindIGDetective revealed nearly all IG genes identified by IterativeIGDetective
670 and more. However, the “bat IG diversity” hypothesis was not supported by our analysis of the spear-
671 nosed bat genome as we identified a much smaller number of IGHV genes than 200+ V genes
672 predicted in Bratsch et al., 2011 using a probabilistic model applied to the little brown bat. There are

673 three possible explanations for a discrepancy between our analysis and the “bat IG diversity”
674 hypothesis: (i) spear-nosed bats have relatively few genes as compared to little brown bats; (ii)
675 probabilistic model in Bratsch et al., 2011 does not adequately approximates the number of IGHV
676 genes; and (iii) many IGHV genes in spear-nosed bats have “diverged” RSSs that do not pass the
677 default RSS likelihood threshold. **Supplemental Method “Applying BlindIGDetective to the**
678 **horseshoe bat genome”** and **Supplemental Table S13** present BlindIGDetective results on the
679 horseshoe bat and also does not support the “bat IG diversity” hypothesis.

680 **Variations in RSSs trigger high/low usage of human D genes.** In addition to analyzing variations in
681 IG genes, we also analyzed variations in RSSs and their effect on antibody repertoires.

682 The usage of a gene in an antibody repertoire is defined as the percentage of antibodies formed by
683 VDJ recombinations in this repertoire that involve this gene. The IG genes have a highly non-uniform
684 usage that may vary by orders of magnitude, e.g., while the most widely used human D gene (IGHD3-
685 10) is used in ~15% of all human antibodies, some human D genes hardly ever contribute to
686 formation of human antibodies (Safonova and Pevzner, 2019). Since the usage of IG genes is likely
687 affected by the sequence of their RSSs, we analyzed the associations between the gene usage of IG
688 genes and their RSSs. **Supplemental Method “Clustering nonamers in RSSVs”** and **Supplemental**
689 **Figures S10, S11** illustrate that RSSs can be partitioned into subgroups of highly similar signals
690 within the set of all RSSs. By revealing these subgroups of similar RSSs we can shed light on the
691 correlations between RSSs and the usage of the genes they flank.

692 Below we focus on analyzing correlations between the usage of D genes and their RSSs by analyzing
693 immunosequencing datasets containing rearranged VDJ sequences from 24 donors from the study by
694 Levin et al., 2017. For each such VDJ sequence, we used the IgScout tool (Safonova and Pevzner,
695 2019) to identify the D gene that contributed to this sequence. The *individual usage* of a gene for a
696 single dataset is defined as the percentage of total VDJ recombinations derived from this gene in this
697 dataset. Although the individual usages vary, they have a rather low variance across various
698 individuals in a given species (Safonova and Pevzner, 2019). The *usage of an IG gene* is defined as

699 the average of individual usages across all 24 datasets generated in Levin et al., 2017. The *usage of an*
700 *RSS* is defined as the usage of the gene that this signal flanks.

701 Below we analyze N D-genes in a reference species and consider N pairs ($RSSD_{left}$, $RSSD_{right}$)
702 flanking these genes (similar analysis has been conducted for V and J genes). We refer to heptamers
703 (nonamers) in this pair as $l7$ and $r7$ ($l9$ and $r9$) and consider 16-mers $l7l9$ and $r7r9$ as well as a 32-mer
704 $l9l7r7r9$, resulting in various signal types. For each signal type, we computed the $N \times N$ matrix of
705 Hamming distances between all pairs of signals and performed k -means clustering on N -dimensional
706 vectors formed by rows of this matrix.

707 We launched the k -means clustering algorithm (20 runs with 200 iterations for each run) for various
708 cluster numbers and selected an optimal number of clusters based on the elbow method (Yuan and
709 Yang, 2019) for all signal types. We determined three $l7$ and three $l9$ clusters as well as two $r7$ and
710 three $r9$ clusters. Each $l9$ ($r9$) cluster can be decomposed into groups of D genes, where each group is
711 a subset of an $l7$ ($l9$) cluster and no pair of groups can be merged. For three $l9$ clusters, the relative
712 sizes of the largest groups with respect to the cluster size are 100%, 100%, and 70% (**Supplemental**
713 **Table S14**). For three $r9$ clusters, the relative sizes of the largest groups are 53%, 75%, and 100%.
714 We hypothesize that clusters $l7$ and $l9$ (as well $r7$ and $r9$) and the high level of overlap between them
715 can trigger variations in usage patterns of D genes.

716 **Figure 6B** illustrates that 12 out of the total 25 human D genes dominate the vast majority (93%) of
717 the usage. We integrate the usage statistics of the D genes into the clustering process. The distribution
718 of usage between the clusters described earlier is recorded for all combinations of the RSSD signals in
719 the **Supplemental Method “Cluster-based usage of RSSDs”, Supplemental Table S15**. The
720 Kruskal-Wallis test on the usage of the signals belonging to each cluster revealed statistically
721 significant associations for both $RSSD_{left}$ and $RSSD_{right}$ heptamer clusters (P-values = 0.042 and
722 0.007, respectively). We also found statistically significant associations for clusters $l9l7r7r9$, $l7l9$ and
723 $r7r9$ (P-values = 0.02, 0.009, 0.028, respectively). These clusters and their significant usage
724 associations are shown in **Figure 6A,B**.

725 Analysis of clusters revealed motifs of RSSs associated with high-usage D genes. Among the
 726 heptamers, the highest used $RSSD_{left}$ cluster (accounting for over 80 % of usage) has the general
 727 pattern **NACTGTG**, where ‘N’ stands for an arbitrary nucleotide. The most used $RSSD_{right}$ heptamer
 728 cluster (accounting for over 94% of usage) has each heptamer equal to **CACAGTG**. The highest used
 729 nonamer cluster for $RSSD_{left}$ follows a pattern **GGTTTNNNN**, whereas the $RSSD_{right}$ nonamer cluster
 730 follows a pattern **NNNAAAACN**.

731 Analysis of V and J genes did not reveal any associations between RSS clusters and usage patterns
 732 (**Supplemental Method “Finding association between RSSs and usages of V and J genes”**,
 733 **Supplemental Tables S16–17, Supplemental Figures S12–16**).

734 [FIGURE 6]

735 **Figure 6. Clustering and distribution of human RSSDs.** (A) Cluster visualization of *I7*, *r7*, *I9I7r7r9*, *I7I9* and
 736 *r7r9* signals. We shall henceforth refer to the red, blue, green and yellow clusters as clusters 1, 2, 3, and 4,
 737 respectively. The consensus of a cluster is noted as the legend label. PC1 and PC2 refer to the first 2 principal
 738 components of the clustering performed on the signals, as described in subsection “Variations in RSSs trigger
 739 high/low usage of human D genes” in Results. (B) Usage of D genes with respect to clusters on *I7*, *r7*, *I9I7r7r9*,
 740 *I7I9* and *r7r9* signals. The p-value of correlation is depicted on the top right of each panel - P* (P**) represents
 741 a p-value less than 0.05 (0.01). (B:Bottom-right) usage of human D genes. Each of 12 highly used human D
 742 genes (with usage at least 2 %) is represented by a single bar. All remaining low-usage human D genes are
 743 represented by a single bar showing their combined usage equal to 7%.

744 Discussion

745 **IGDetective algorithm.** We benchmarked IterativeIGDetective on three well-annotated IGH loci
 746 (human, mouse, and cow) and demonstrated that it accurately predicts the known V, D, and J genes in
 747 one of these species based on information about RSSs in other species. This observation justifies our
 748 comparative immunogenomics approach to annotating the IGH loci in newly sequenced species.
 749 Although the IterativeIGDetective analysis in this paper is limited to the IGH loci, we plan to extend
 750 it to other IG and TCR loci in a follow-up study. In addition to the three reference species, we applied

751 IterativeIGDetective to twenty mammalian species with the recently assembled IGH loci and
752 predicted 581, 92, and 60 new putative IGHV, IGHD, and IGHJ genes, respectively.

753 In addition to IterativeIGDetective, we developed BlindIGDetective algorithm for predicting novel
754 genes which have diverged from currently known IG genes. BlindIGDetective detects most IG genes
755 in the absence of any information about IG genes in other species and thus opens a possibility to
756 identify highly divergent IG genes. Application of BlindIGDetective to detect V genes from three
757 reference species and two bat species (spear-nosed bat and horseshoe bat) revealed that
758 BlindIGDetective was sensitive not only to canonical V genes from the IGH locus but also to other V
759 gene loci driven by RAG proteins, such as V genes in the IGL, TRA, TRB, and TRG loci. Moreover,
760 BlindIGDetective identified multiple highly-divergent candidate V genes (or pseudogenes) in various
761 species.

762 We plan to combine IterativeIGDetective and BlindIGDetective to extend this analysis to
763 immunological model organisms such as rabbits and llamas as well as non-mammalian vertebrates
764 with the goal to construct a comprehensive database of predicted IG genes across multiple species.

765 **The diversity of IG genes.** We applied IterativeIGDetective to twenty mammalian species with
766 poorly-studied IGH loci, performed comparative analysis of the detected IGHV genes, and identified
767 a highly conservative cluster that covers highly divergent species ranging from primates to bats. We
768 hypothesize that V genes in this cluster are subjected to selective pressure driven by common
769 pathogens or the genetic organization of IGH loci. Further investigation of this conservative cluster
770 will require repertoire sequencing (Rep-Seq) data.

771 In addition to revealing the diversity of IG genes, we also studied the diversity of RSSs, identified
772 clusters of similar RSSDs in humans, revealed associations between these clusters and the usage of
773 the IGHD gene they flank, and found the RSSDs motifs triggering high usage of human D genes.

774 **Unusual cysteine-rich V genes.** We revealed a new family of unusual cysteine-rich V genes in bats
775 and other species that have cysteines in both CDR1 and CDR2. We hypothesize that cysteine-rich V

776 genes might generate unusual antibodies with non-canonical conformations and could potentially
777 form a unique part of bat immunity against the great variety of viruses they host. Further investigation
778 of antibodies derived from these V genes would require both Rep-Seq data and protein structures to
779 shed light on the functions and to estimate the therapeutic potential of such antibodies.

780 **G/S/Y-rich motifs in D genes.** We demonstrated that, despite being highly diverse, D genes in
781 various mammalian species share the G/S/Y-rich motifs that are formed by cysteine-triggering
782 codons. In cows, the G/S/Y-rich IGHD8-2 gene plays a special cysteine-triggering role in ultralong
783 cattle antibodies where somatic hypermutations create new cysteines from these three amino acids,
784 forming new disulfide bonds and reshaping the resulting antibody (Wang et al., 2013). We conjecture
785 that found G/S/Y motifs may play a similar cysteine-triggering role in D genes of many mammals.

786 **Limitations and future developments.** In the last three decades, gene prediction algorithms have
787 evolved from relatively simple statistical tests to sophisticated machine learning approaches (Mathé et
788 al., 2002). However, since even the state-of-the-art gene prediction algorithms generate some false
789 positive genes, they require validation using complementary experimental approaches, such as
790 transcriptome sequencing (Allen et al., 2004) and mass spectrometry (Tanner et al., 2007). Likewise,
791 since IGDetective is merely the first step toward uncovering the diversity of IG genes across
792 vertebrate species, it has to be complemented by complementary experimental approaches, such as
793 antibody repertoire sequencing. Our next goal is to modify IGDetective for working with both
794 genome assemblies and Rep-Seq data.

795 IterativeIGDetective is currently based on constructing the profile matrices of known RSSs, detecting
796 novel RSSs using these profile matrices, and follow-up analysis of genomic sequences flanked by
797 these RSSs. In the future, we plan to develop a Hidden Markov Model for joint analysis of RSSs and
798 the flanking genes. We also plan to complement the existing analysis (based on only three reference
799 genomes without re-training) by bootstrapping when the original model is trained on the references
800 only but the set of references is later extended (based on the reliable predictions of IG genes in new
801 species) for further re-training.

802 IterativeIGDetective failed to identify any J genes in the spear-nosed bat genome, presumably because
 803 all its RSSJ motifs did not pass the default likelihood threshold. Since our current goal is to focus on
 804 highly conserved RSSs and minimize the false positive rate, we did not recompute all IG gene
 805 predictions in the spear-nosed bat with a lower likelihood threshold. In the future, we plan to develop
 806 a version of IterativeIGDetective that iteratively relaxes the RSS search parameters.

807 BlindIGDetective currently starts from identifying highly-conserved RSS in the entire genome and
 808 thus misses all IG genes flanked by less conservative RSSs. Lowering the RSS likelihood threshold
 809 leads to explosion of false RSS thus making BlindIGDetective prohibitively slow. We plan to modify
 810 BlindIGDetective (with the goal of identifying the missed IG genes with less conservative RSSs) by
 811 first identifying the IG-contigs, enriching the set of putative RSSs in these contigs by adding less
 812 conservative RSSs, and combining BlindIGDetective and IterativeIGDetective in a single pipeline.

813 **Methods**

814 **Gene nomenclature.** Followed guidelines of the IMGT nomenclature
 815 (<https://www.imgt.org/IMGTScientificChart/Nomenclature/IMGTnomenclature.html>), we have not
 816 italicized genes of immunoglobulin (IG) and T-cell receptor (TCR) genes.

817 **Profiles and likelihood ratio of strings.** Given a set of k -mers (k -nucleotide-long strings), their
 818 *profile matrix* is a $4 \times k$ matrix *Profile* defined below. The j^{th} column in the profile matrix represents
 819 the j^{th} position in the k -mer and the 4 rows represent the 4 nucleotide bases (A, C, G, and T).
 820 $Profile(i,j)$ represents the frequency of occurrence of the i^{th} base at the j^{th} position in the input k -mers
 821 adjusted for pseudocounts as described in Compeau and Pevzner, 2018. The *consensus string*
 822 (referred to as $consensus=consensus(Profile)$) is a k -mer generated by taking a nucleotide with the
 823 highest frequency at each position of the profile (ties are broken arbitrarily). As opposed to numerous
 824 V and D genes, there are few J genes in the reference species, leading to profile matrices with higher
 825 entropies (when compared to V or D gene profiles) due to pseudocounts.

826 Given a k -mer $S=s_1\dots s_k$ and a $4\times k$ profile $Profile$, the probability that $Profile$ generates a string S is
 827 defined as $prob(S|Profile) = \prod_{j=1,k} Profile(s_j,j)$. Since $prob(S|Profile)$ is maximized when the string S
 828 is a consensus of $Profile$, we define the *likelihood ratio* $L(S)$ of S as

$$829 \quad L(S) = \frac{prob(S|Profile)}{prob(consensus(Profile)|Profile)},$$

830 Given a profile $Profile$ and a *likelihood ratio threshold* L_{min} , a k -mer S is classified as a *candidate* for
 831 $Profile$ if its likelihood ratio $L(S)$ exceeds L_{min} .

832 **RSS detection algorithm.** Given a set of RSSs for each of a reference species (human, mouse, cow,
 833 or combined), we compute their profile matrix (with pseudocounts equal to 1) for heptamers and
 834 nonamers across all signals (RSSV, $RSSD_{left}$, $RSSD_{right}$, and RSSJ). Given a *profile-pair*
 835 $(Profile, Profile')$ of 4×7 and 4×9 stochastic matrices and the associated thresholds L_{min} and L'_{min} for a
 836 heptamer and a nonamer, a *string-pair* (*heptamer*, *nonamer*) is classified as a *candidate RSS* for a
 837 given signal type if both *heptamer* and *nonamer* are candidate strings.

838 We say that a candidate $RSSD_{left}$ and a candidate $RSSD_{right}$ together form a *paired candidate RSSD* if
 839 the $RSSD_{right}$ is located within at most $MaxLength_D$ positions to the right from the $RSSD_{left}$. The
 840 condition is important for D genes since they are flanked by $RSSD_{left}$ and $RSSD_{right}$ and since the vast
 841 majority of known D genes are short (maximum lengths are 37 nt and 29 nt in human and mouse
 842 genomes, respectively). Since the longest currently known D gene is the 148 nt long D gene in cows,
 843 we set the default value $MaxLength_D=150$.

844 **Selection of the likelihood ratio threshold.** We select the L_{min} thresholds for heptamers or nonamers
 845 based on the reference genome chosen. We launch IterativeIGDetective on the reference species,
 846 selecting the heptamer L_{min} and nonamer L'_{min} thresholds for all signal types (RSSV, $RSSD_{left}$,
 847 $RSSD_{right}$, and RSSJ) by performing a grid search within a range (0,0.8] in steps of 0.005. We
 848 compute the TPR (the fraction of known RSSs in the target species detected by IGDetective), the FDR
 849 (the fraction of erroneous RSSs among all RSSs reported by IGDetective), and the F1 score defined
 850 as:

851
$$F1 = \frac{2TPR(1-FDR)}{TPR+(1-FDR)}$$

852 For any given reference, the heptamer L_{min} and nonamer L'_{min} thresholds (determined through grid
853 search) that maximize the F1 score are selected as that references' heptamer L_{min} and nonamer L'_{min} .
854 More details on parameter selection are described in **Supplemental Method “Parameter selection
855 for identifying candidate RSSs”**.

856 **Identification of candidate IG genes.** Details of procedures for identification of candidate V, D, and
857 J genes, including selection of parameters of alignments are described in **Supplemental Methods
858 “Identification of candidate V and J genes”, “Identification of candidate D genes”, “Iterative
859 extension of the set of candidate IG genes”, “Parameter selection for identifying candidate IG
860 genes”, “Aligning candidate IG genes”, and “Alternate similarity thresholds for detection of V
861 genes”**.

862 **Software Availability.** IGDetective is available as Supplemental Code and at
863 github.com/Immunotools/IgDetective. Sequences of identified IGHV, IGHD, and IGHJ genes are
864 available as **Supplemental Tables S19–21** and at github.com/Immunotools/IgDetective_results.

865 **Competing interest statement.** The authors declare no competing financial interests.

866 **Acknowledgments.** This work was supported by the National Science Foundation EAGER award
867 (no. 2032783). All authors contributed to the development of the IGDetective algorithms. VS and YS
868 implemented the IGDetective algorithm. VS and YS analyzed the data and performed computational
869 experiments. All authors conceived the study and wrote the manuscript.

870 **References**

871 Allen, J.E., Pertea, M., Salzberg, S.L. 2004. Computational gene prediction using multiple sources of evidence.
872 *Genome Res.* **14**: 142-148

- 873 Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990. Basic local alignment search tool. *J Mol Biol*
874 **215**(3):403-10.
- 875 Avnir Y, Tallarico AS, Zhu Q, Bennett AS, Connelly G, Sheehan J, Sui J, Fahmy A, Huang CY, Cadwell G, et
876 al., 2014. Molecular signatures of hemagglutinin stem-directed heterosubtypic human neutralizing antibodies
877 against influenza A viruses. *PLoS Pathog* **10**(5):e1004103.
- 878 Avnir Y, Watson CT, Glanville J, Peterson EC, Tallarico AS, Bennett AS, Qin K, Fu Y, Huang CY, Beigel JH,
879 et al., 2016. IGHV1-69 polymorphism modulates anti-influenza antibody repertoires, correlates with IGHV
880 utilization shifts and varies by ethnicity. *Scientific reports* **6**(1):1-3.
- 881 Bankevich, A., Pevzner, P. 2020. mosaicFlye: Resolving long mosaic repeats using long error-prone reads.
882 *bioRxiv* doi: <https://doi.org/10.1101/2020.01.15.908285>.
- 883 Bratsch S, Wertz N, Chaloner K, Kunz TH, Butler JE. 2011. The little brown bat, *M. lucifugus*, displays a
884 highly diverse VH, DH and JH repertoire but little evidence of somatic hypermutation. *Developmental &*
885 *Comparative Immunology* **35**(4):421-30.
- 886 Compeau, P. C.A., Pevzner, P.A. 2018. *Bioinformatics Algorithms: An Active Learning Approach* (3rd edition).
887 Active Learning Publishers
- 888 Das S, Hirano M, Tako R, McCallister C, Nikolaidis N. 2012. Evolutionary genomics of immunoglobulin-
889 encoding loci in vertebrates. *Current Genomics* **13**(2):95-102.
- 890 Dudley DD, Chaudhuri, Bassing CH, Alt FW. 2005. Mechanism and control of V(D)J recombination versus
891 class switch recombination: similarities and differences. *Advances in Immunology* **86**:43-112.
- 892 Du L, Wang S, Zhu Y, Zhao H, Basit A, Yu X, Li Q, Sun X. 2018. Immunoglobulin heavy chain variable region
893 analysis in dairy goats. *Immunobiology* **223**:599-607
- 894 Eguchi-Ogawa T, Wertz N, Sun Z, Piumi F, Uenishi H, Wells K, Chardon P, Tobin GJ, Butler JE. 2010.
895 Antibody repertoire development in fetal and neonatal piglets. XI. The relationship of variable heavy chain gene
896 usage and the genomic organization of the variable heavy chain locus. *J. Immunol* **184**: 3734-3742.
- 897 Frangione B., Milstein C., Pink J.R. 1969. Structural studies of immunoglobulin G. *Nature* **221**:145-148.

- 898 Gambon-Deza F, Sánchez-Espinel C, Magadan-Mompo S. 2009. The immunoglobulin heavy chain locus in the
899 platypus (*Ornithorhynchus anatinus*). *Molecular immunology*. **46**(13):2515-23.
- 900 Gertz EM, Schäffer AA, Agarwala R, Bonnet-Garnier A, Rogel-Gaillard C, Hayes H, Mage RG. 2013.
901 Accuracy and coverage assessment of *Oryctolagus cuniculus* (rabbit) genes encoding immunoglobulins in the
902 whole genome sequence assembly (OryCun2.0) and localization of the IGH locus to chromosome 20.
903 *Immunogenetics* **65**:749-762.
- 904 Hedges SB, Marin J, Suleski M, Paymer M, Kumar S. 2015. Tree of life reveals clock-like speciation and
905 diversification. *Molecular Biology and Evolution* **32**(4):835-45.
- 906 Johansson J, Aveskogh M, Munday B, Hellman L. 2002. Heavy chain V region diversity in the duck-billed
907 platypus (*Ornithorhynchus anatinus*): long and highly variable complementarity-determining region 3
908 compensates for limited germline diversity. *The Journal of Immunology* **168**(10):5155-62.
- 909 Keyaerts M, Xavier C, Heemskerk J, Devoogdt N, Everaert H, Ackaert C, Vanhoeij M, Duhoux FP, Gevaert T,
910 Simon P, et al., 2016. Phase I study of ⁶⁸Ga-HER2-nanobody for PET/CT assessment of HER2 expression in
911 breast carcinoma. *Journal of Nuclear Medicine* **57**(1):27-33.
- 912 Kim SI, Noh J, Kim S, Choi Y, Yoo DK, Lee Y, Lee H, Jung J, Kang CK, Song KH, et al., 2021. Stereotypic
913 neutralizing VH antibodies against SARS-CoV-2 spike protein receptor binding domain in patients with
914 COVID-19 and healthy individuals. *Science Translational Medicine* **13**(578).
- 915 Kronenberg ZN, Fiddes IT, Gordon D, Murali S, Cantsilieris S, Meyerson OS, Underwood JG, Nelson BJ,
916 Chaisson MJP, Dougherty ML, et al., 2018. High-resolution comparative analysis of great ape genomes. *Science*
917 **360**:eaar6343.
- 918 Krumsiek J, Arnold R, Rattei T. 2007. Gepard: a rapid and sensitive tool for creating dotplots on genome scale.
919 *Bioinformatics* **23**(8):1026-8.
- 920 Kruskal WH and Wallis WW. 1952. Use of Ranks in One-Criterion Variance Analysis. *Journal of the American*
921 *Statistical Association* **47**:583-621.
- 922 Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nature Methods* **9**:357.

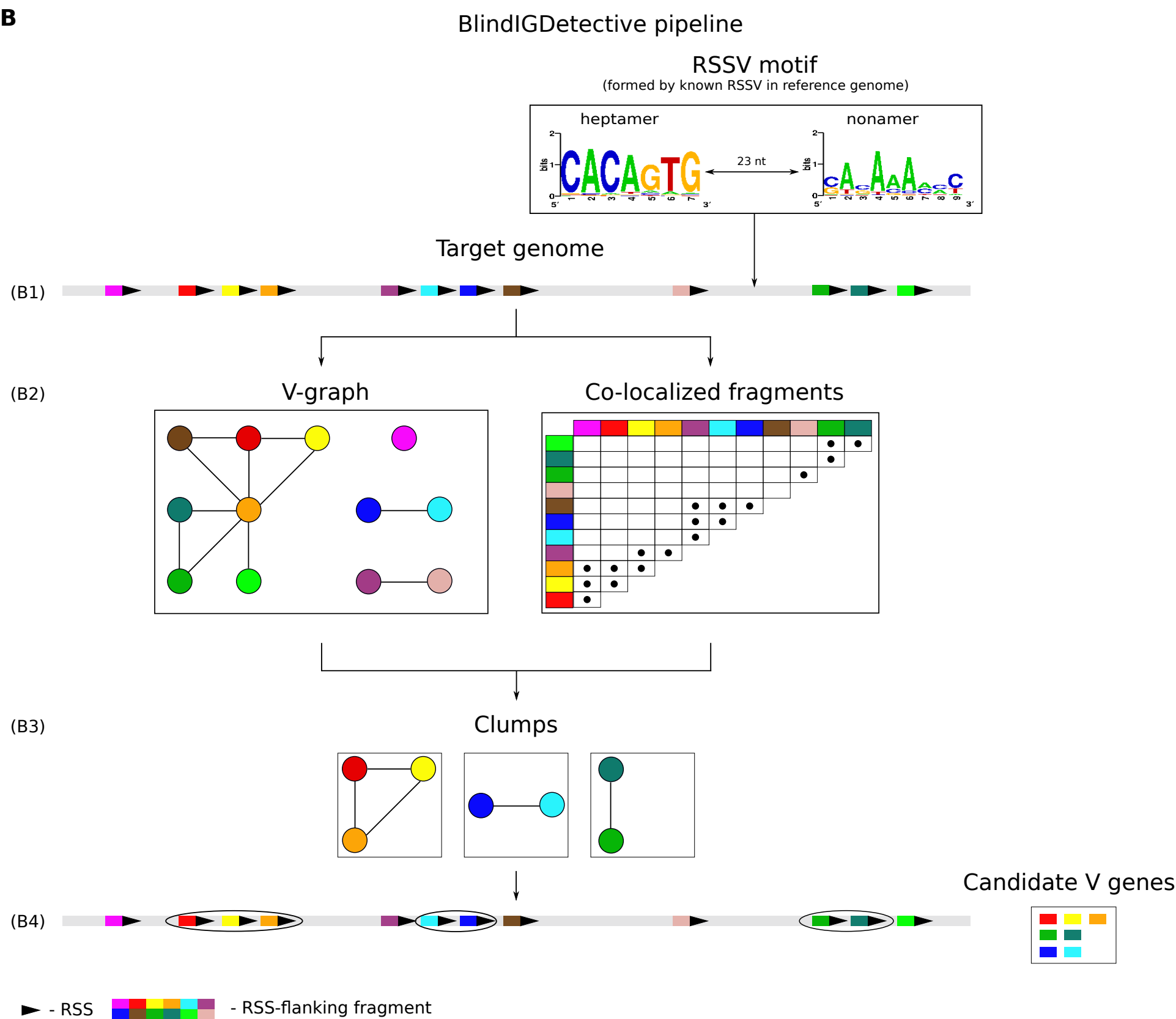
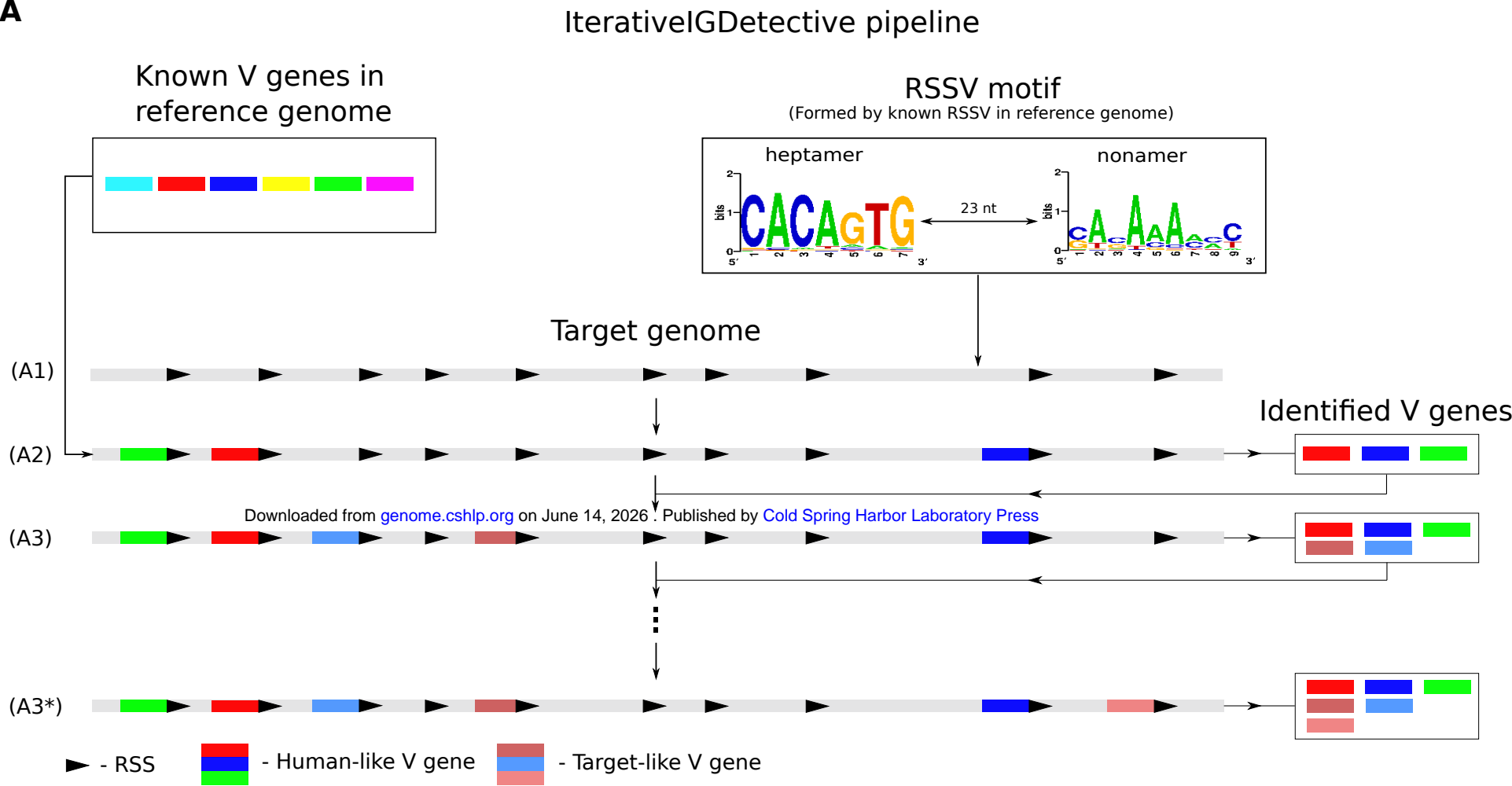
- 923 Lane J, Duroux P, Lefranc MP. 2010. From IMGT-ONTOLOGY to IMGT/LIGMotif: the IMGT® standardized
924 approach for immunoglobulin and T cell receptor gene identification and description in large genomic
925 sequences. *BMC Bioinformatics* **11**(1):1-6.
- 926 Lefranc M-P, Giudicelli V, Duroux P, Jabado-Michaloud J, Folch G, Aouinti S, Carillon E, Duvergey H, Houles
927 A, Paysan-Lafosse T, et al., 2015. IMGT®, the international ImMunoGeneTics information system® 25 years
928 on. *Nucleic Acids Res* **43**(Database issue):D413-22.
- 929 Levin M, Levander F, Palmason R, Greiff L, Ohlin M. 2017. Antibody-encoding repertoires of bone marrow
930 and peripheral blood—a focus on IgE. *J Allergy Clin Immunol* **139**:1026–30.
- 931 Li, H., Cui, X., Pramanik, S., Chimge, N.O. 2002. Genetic diversity of the human immunoglobulin heavy chain
932 VH region. *Immunological Reviews* **190**, 53-68
- 933 Lingwood D, McTamney PM, Yassine HM, Whittle JR, Guo X, Boyington JC, Wei CJ, Nabel GJ. 2012.
934 Structural and genetic basis for development of broadly neutralizing influenza antibodies. *Nature*
935 **489**(7417):566-70.
- 936 De Los Rios M, Criscitiello MF, Smider VV. 2015. Structural and genetic diversity in antibody repertoires from
937 diverse species. *Current Opinion in Structural Biology* **33**:27-41.
- 938 Ma L, Qin T, Chu D, Cheng X, Wang J, Wang X, Wang P, Han H, Ren L, Aitken R, et al., 2016. Internal
939 duplications of DH, JH, and C region genes create an unusual IgH gene locus in cattle. *The Journal of*
940 *Immunology* **196**(10):4358-66.
- 941 Mathé, M., Sagot, M.F., Schiex, T., Rouzé, P. 2002. Current methods of gene prediction, their strengths and
942 weaknesses, *Nucleic Acids Research* **30**, 4103–4117
- 943 Matsuda F, Ishii K, Bourvagnet P, Kuma KI, Hayashida H, Miyata T, Honjo T. 1998. The complete nucleotide
944 sequence of the human immunoglobulin heavy chain variable region locus. *The Journal of Experimental*
945 *Medicine* **188**:2151-62.
- 946 Merelli I, Guffanti A, Fabbri M, Cocito A, Furia L, Grazini U, Bonnal RJ, Milanesi L, McBlane F. 2010.
947 RSSsite: a reference database and prediction tool for the identification of cryptic Recombination Signal
948 Sequences in human and murine genomes. *Nucleic Acids Res.* **38** (suppl. 2): W262-7.

- 949 Messier, T. L., O'Neill, J. P., Hou, S.-M., Nicklas, J. A. & Finette, B. A. 2003. *In vivo* transposition mediated by
950 V(D)J recombinase in human T lymphocytes. *EMBO Journal* **22**, 1381–1388
- 951 Miller RD. 2010. Those other mammals: the immunoglobulins and T cell receptors of marsupials and
952 monotremes. *Seminars in Immunology* **22**, 3-9
- 953 Moore RM, Harrison AO, McAllister SM, Polson SW, Wommack KE. 2020. Iroki: automatic customization and
954 visualization of phylogenetic trees. *PeerJ*. **8**: e8584.
- 955 Muyldermans S, Smider VV. 2016. Distinct antibody species: structural differences creating therapeutic
956 opportunities. *Current Opinion in Immunology* **40**:7-13.
- 957 Nagaoka, H., Ozawa, K., Matsuda, F., Hayashida, H., Matsumura, R., Haino, M., Shin, E. K., Fukita, Y., Imai,
958 T., Anand, R., et al., 1994. Recent translocation of variable and diversity segments of the human
959 immunoglobulin heavy chain from chromosome 14 to chromosomes 15 and 16. *Genomics* **22**: 189-197
- 960 Nagawa F, Ishiguro KI, Tsuboi A, Yoshida T, Ishikawa A, Takemori T, Otsuka AJ, Sakano H. 1998. Footprint
961 analysis of the RAG protein recombination signal sequence complex for V(D)J type recombination. *Mol Cell*
962 *Biol.* **18**: 655–663.
- 963 Olivieri D, Faro J, von Haefen B, Sánchez-Espinel C, Gambón-Deza F. 2013. An automated algorithm for
964 extracting functional immunologic V-genes from genomes in jawed vertebrates. *Immunogenetics* **65**(9):691-702.
- 965 Olivieri, D.N., Gambón-Deza, F. 2019. Iterative Variable Gene Discovery from Whole Genome Sequencing
966 with a Bootstrapped Multiresolution Algorithm. *Computational and Mathematical Methods in Medicine*
967 **3780245**.
- 968 Prabakaran P, Chowdhury PS. 2020. Landscape of non-canonical cysteines in human VH repertoire revealed by
969 immunogenetic analysis. *Cell Reports* **31**(13):107831.
- 970 Pettinello R, Dooley H. 2014. The immunoglobulins of cold-blooded vertebrates. *Biomolecules* **4**:1045-69.
- 971 Qiu X, Ma F, Zhao M, Cao Y, Shipp L, Liu A, Dutta A, Singh A, Braikia FZ, De S, et al., 2020. Altered 3D
972 chromatin structure permits inversional recombination at the IgH locus. *Science Advances* **6**:eaaz8850.

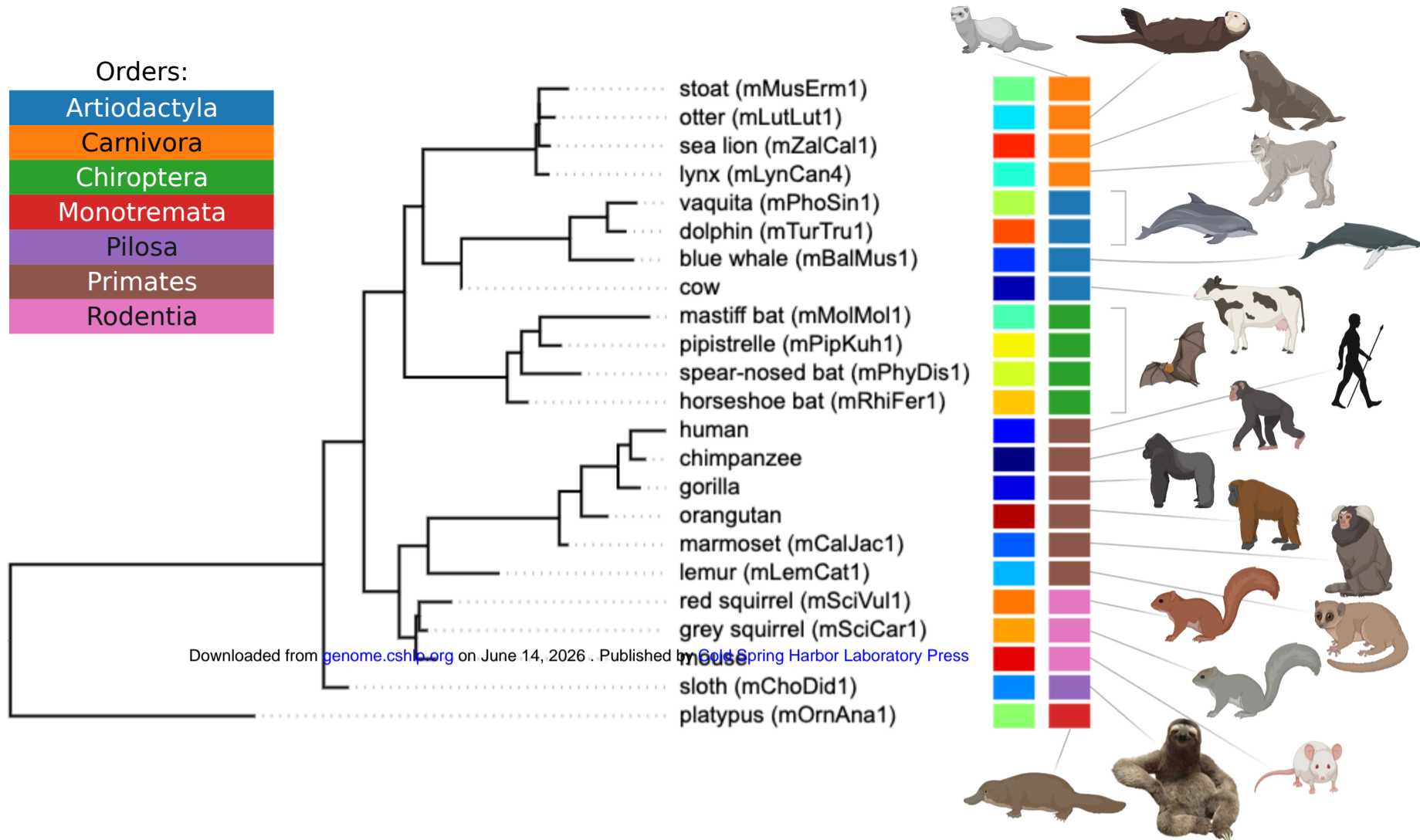
- 973 Rashidian M, Keliher EJ, Bilate AM, Duarte JN, Wojtkiewicz GR, Jacobsen JT, Cragolini J, Swee LK, Victora
974 GD, Weissleder R, et al., 2015. Noninvasive imaging of immune responses. *Proceedings of the National*
975 *Academy of Sciences* **112**:6146-51.
- 976 Rhie A, McCarthy SA, Fedrigo O, Damas J, Formenti G, Koren S, Uliano-Silva M, Chow W, Fungtammasan A,
977 Kim J, et al., 2021. Towards complete and error-free genome assemblies of all vertebrate species. *Nature*
978 **592**(7856):737-46.
- 979 Rodriguez OL, Gibson WS, Parks T, Emery M, Powell J, Strahl M, Deikus G, Auckland K, Eichler EE,
980 Marasco WA, et al., 2020. A Novel Framework for Characterizing Genomic Haplotype Diversity in the Human
981 Immunoglobulin Heavy Chain Locus. *Front Immunol* **11**:2136.
- 982 Safonova Y, Bonissone S, Kurpilyansky E, Starostina E, Lapidus A, Stinson J, DePalatis L, Sandoval W, Lill J,
983 Pevzner PA. 2015. IgRepertoireConstructor: a novel algorithm for antibody repertoire construction and
984 immunoproteogenomics analysis. *Bioinformatics* **31**(12):i53-61.
- 985 Safonova Y, Pevzner PA. 2020. V(DD)J recombination is an important and evolutionary conserved mechanism
986 for generating antibodies with unusually long CDR3s. *Genome Res.* **30**:1547–58.
- 987 Safonova Y, Shin SB, Kramer L, Reecy J, Watson CT, Smith TP, Pevzner PA. 2022. Variations in antibody
988 repertoires correlate with vaccine responses. *Genome Res.* **32**:791-804.
- 989 Schountz T, Baker ML, Butler J, Munster V. 2017. Immunological control of viral infections in bats and the
990 emergence of viruses highly pathogenic to humans. *Frontiers in Immunology* **8**:1098.
- 991 Sievers F, Wilm A, Dineen D, Gibson TJ, Karplus K, Li W, Lopez R, McWilliam H, Remmert M, Söding J, et
992 al., 2011. Fast, scalable generation of high quality protein multiple sequence alignments using Clustal Omega.
993 *Molecular Systems Biology* **7**(1):539.
- 994 Sitnikova, T. Su, C. 1998. Coevolution of immunoglobulin heavy- and light-chain variable-region gene families.
995 *Mol. Biol. Evol.* **15**(6):617–625.
- 996 Sok D, Le KM, Vadnais M, Saye-Francisco KL, Jardine JG, Torres JL, Berndsen ZT, Kong L, Stanfield R, Ruiz
997 J, et al., 2017. Rapid elicitation of broadly neutralizing antibodies to HIV by immunization in cows. *Nature*
998 **548**:108-11.

- 999 Tan J, Pieper K, Piccoli L, Abdi A, Foglierini M, Geiger R, Tully CM, Jarrossay D, Ndungu FM, Wambua J, et
1000 al., 2016. A LAIR1 insertion generates broadly reactive antibodies against malaria variant antigens. *Nature*
1001 **529**:105-9.
- 1002 Tanner, S., Shen, Z., Ng, J., Florea, L., Guigó, R., Briggs, S.P., Bafna, V. 2007. Improving gene annotation
1003 using peptide mass spectrometry. *Genome Res.* **17**(2): 231–239.
- 1004 Teeling EC, Vernes SC, Dávalos LM, Ray DA, Gilbert MT, Myers E, Bat1K Consortium. 2018. Bat biology,
1005 genomes, and the Bat1K Project: to generate chromosome-level genomes for all living bat species. *Annual*
1006 *Review of Animal Biosciences* **6**(1):23-46.
- 1007 Teng G, Maman Y, Resch W, Kim M, Yamane A, Qian J, Kieffer-Kwon KR, Mandal M, Ji Y, Meffre E, et al.,
1008 2015. RAG represents a widespread threat to the lymphocyte genome. *Cell* **162**(4):751-65.
- 1009 Tonegawa S. 1983. Somatic generation of antibody diversity. *Nature* **302**:575-581.
- 1010 Wang F, Ekiert DC, Ahmad I, Yu W, Zhang Y, Bazirgan O, Torkamani A, Raudsepp T, Mwangi W, Criscitiello
1011 MF, et al., 2013. Reshaping antibody diversity. *Cell* **153**:1379-93.
- 1012 Watson CT, Steinberg KM, Huddleston J, Warren RL, Malig M, Schein J, Willsey AJ, Joy JB, Scott JK, Graves
1013 TA, et al., 2013. Complete haplotype sequence of the human immunoglobulin heavy-chain variable, diversity,
1014 and joining genes and characterization of allelic and copy-number variation. *The American Journal of Human*
1015 *Genetics* **92**:530-46.
- 1016 Wong J, Tai CM, Hurt AC, Tan HX, Kent SJ, Wheatley AK. 2020. Sequencing B cell receptors from ferrets
1017 (*Mustela putorius furo*). *PLoS One* **15**:e0233794.
- 1018 Wu L, Oficjalska K, Lambert M, Fennell BJ, Darmanin-Sheehan A, Shúilleabháin DN, Autin B, Cummins E,
1019 Tchistiakova L, Bloom L, et al., 2012. Fundamental characteristics of the immunoglobulin VH repertoire of
1020 chickens in comparison with those of humans, mice, and camelids. *Journal of Immunology* **188**(1):322-33.
- 1021 Ye J, Ma N, Madden TL, Ostell JM. 2013. IgBLAST: an immunoglobulin variable domain sequence analysis
1022 tool. *Nucleic Acids Res.* **41**(Web Server issue): W34-W40.
- 1023 Yuan M, Liu H, Wu NC, Lee CC, Zhu X, Zhao F, Huang D, Yu W, Hua Y, Tien H., et al., 2020. Structural
1024 basis of a shared antibody response to SARS-CoV-2. *Science* **369**(6507):1119-23.

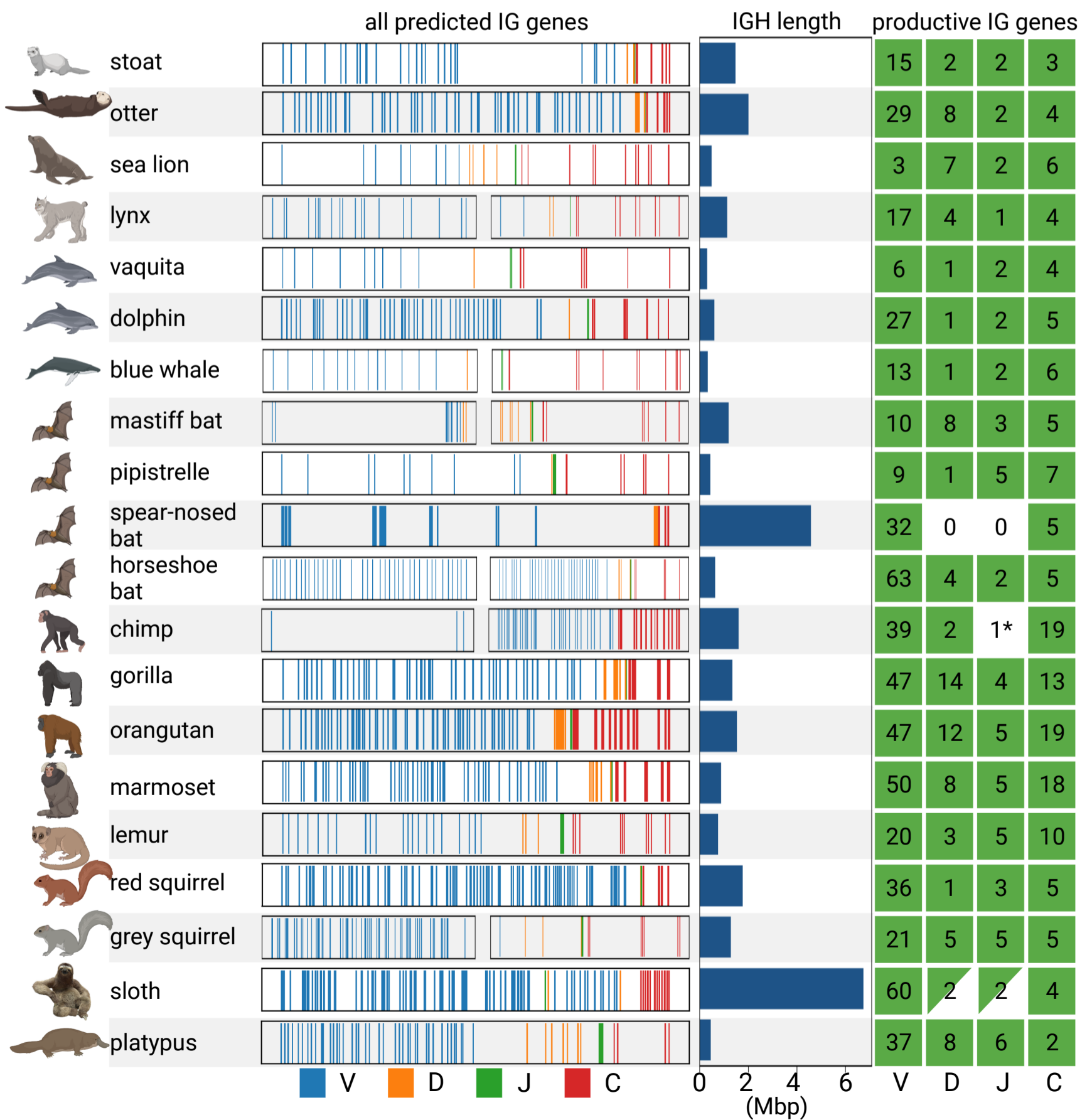
- 1025 Yuan, C. Yang, H. 2019. Research on K-Value Selection method of k-means clustering algorithm. *J -*
1026 *Multidisciplinary Scientific Journal* 2(2), 226-235.

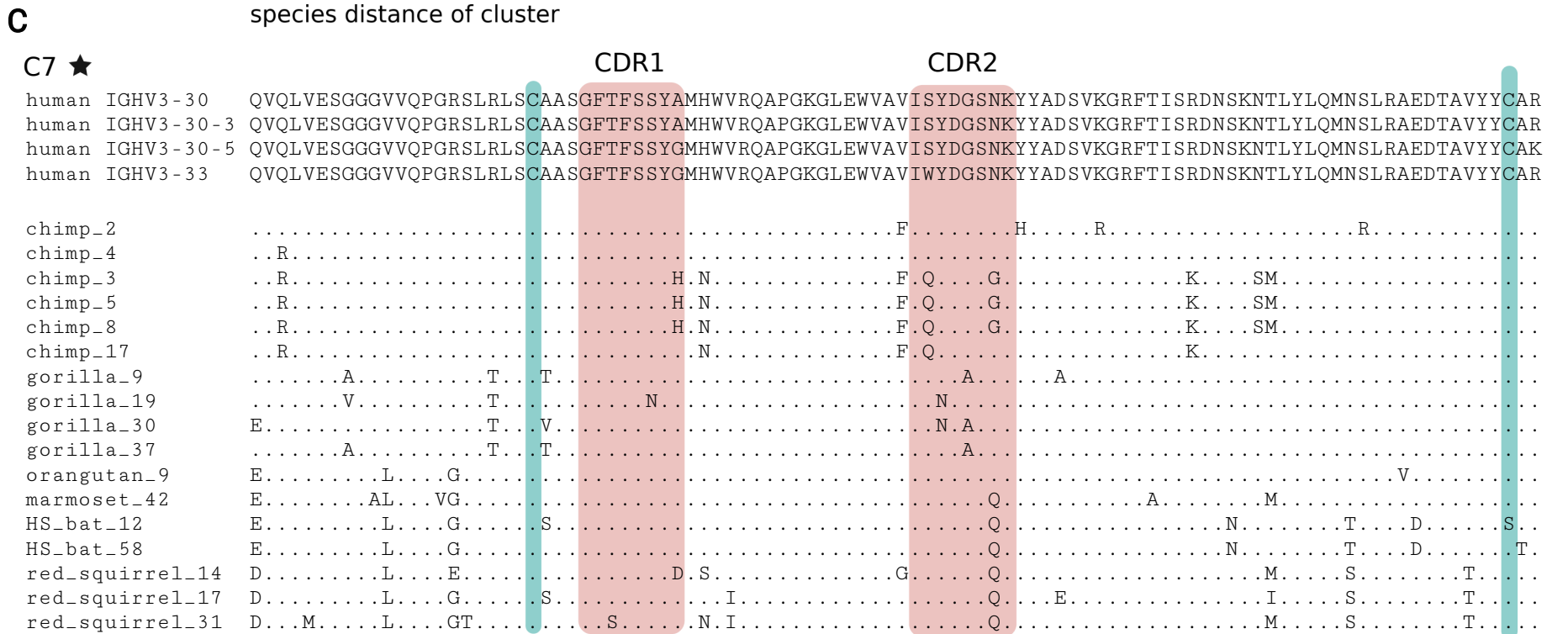
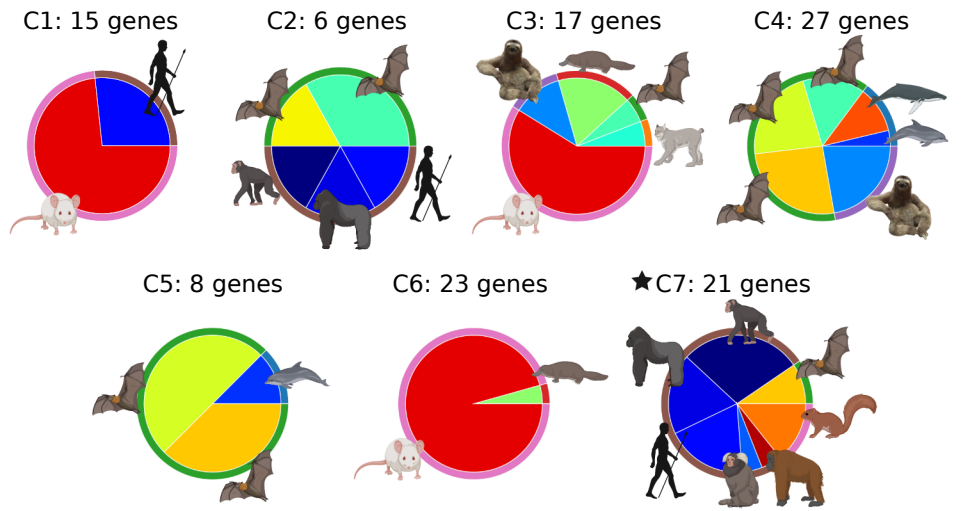
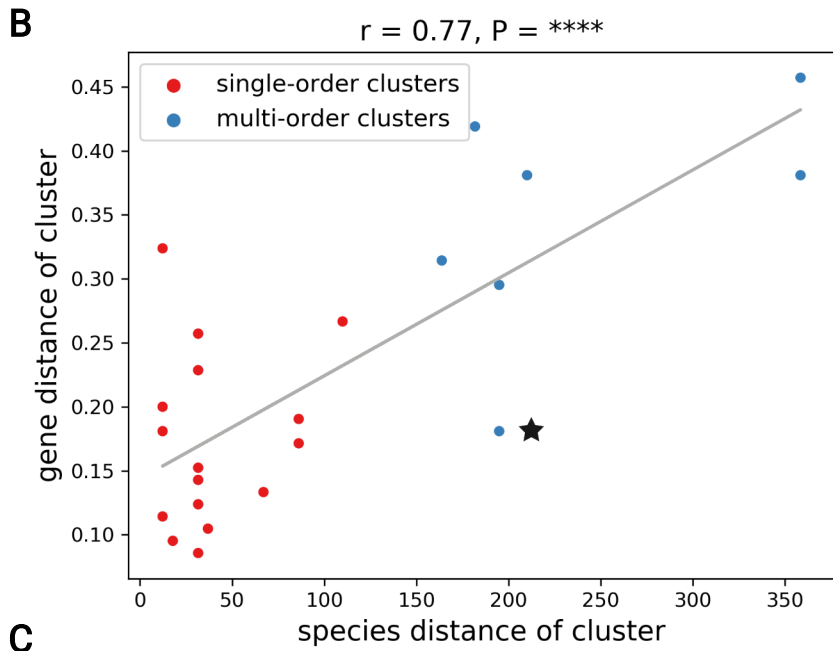
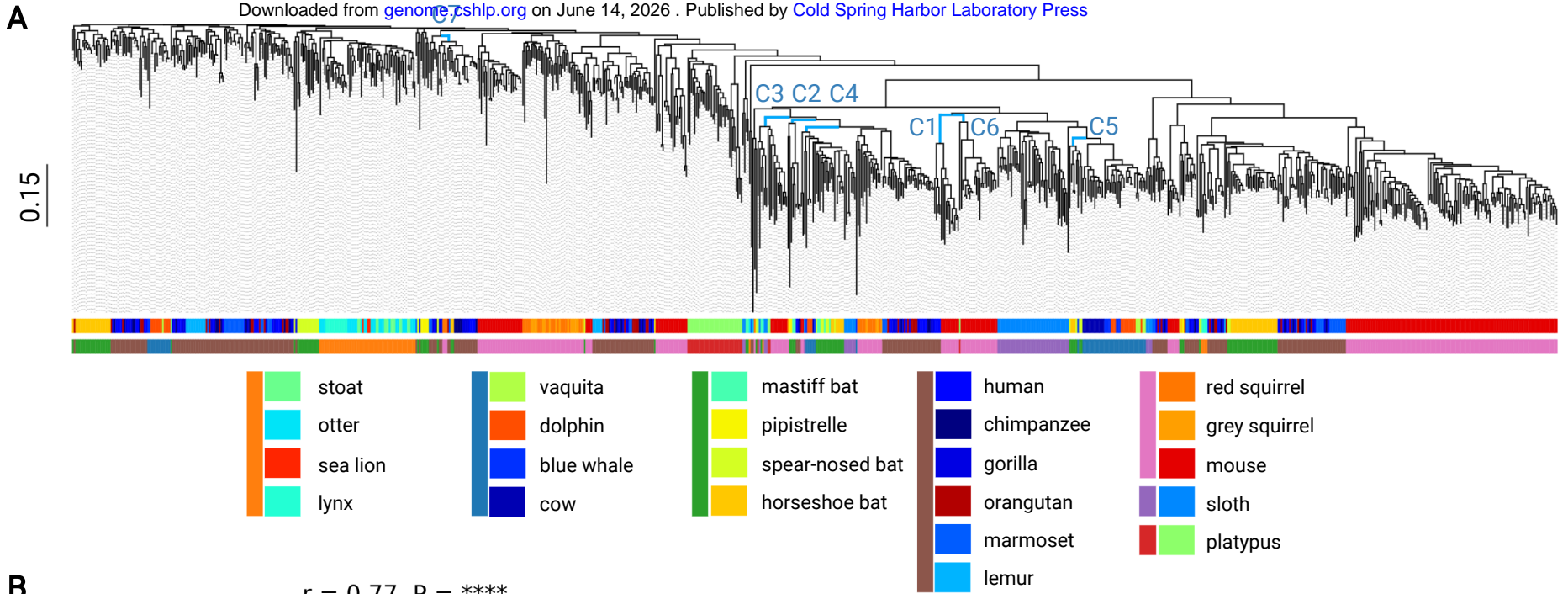


A

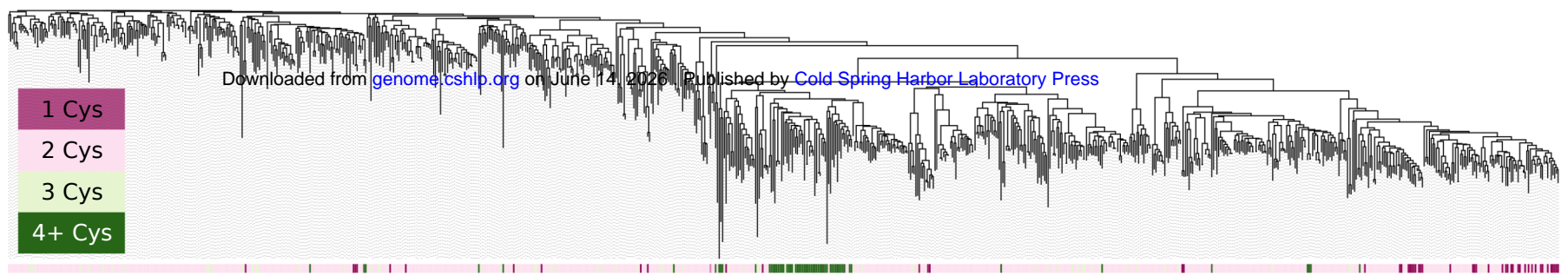


B

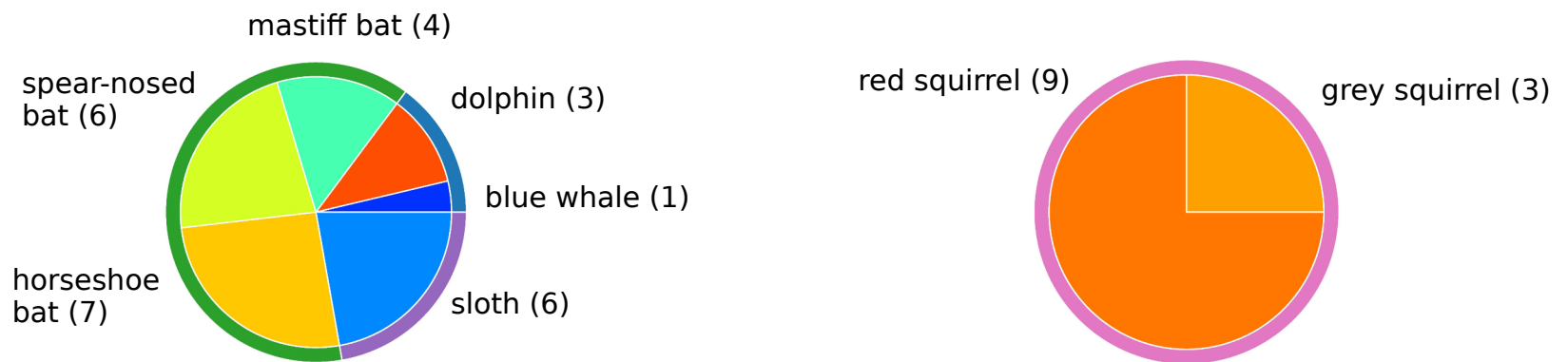




A



B



C

	CDR1	CDR2
HS_bat_59	QVQLQESGPGLVQPSQTLSLTCAVSGFSIT-SGYC	WHWIRQPPGKGLEWIGRICYNGNTYYSPSLKSRSSISRDTSKNQFSLQLSSVTTEDTAVYYCAR---
bluewhale_7	E...QW...LK...T.P...Y.S...SGY...N...T...M.N...ND.T...I...I.H...T...H...MAP...T...K...
sloth_14S...KL...S...T.Y...E.C...Q...Q...I.H...T...H...MAP...T...K...I.H...T...H...MAP...T...K...
sloth_18	...K...S...K...S...SS.Y...A...S...Q...C.Y...S...T...E...T...I.H...T...H...MAP...T...K...
sloth_25	...S...KL...S...T.Y...E.C...Q...Q...I.H...T...E...H...MAP...T...K...I.H...T...E...H...MAP...T...K...
sloth_38	...K.P...D...K...S...T.S...C.L...H...RA...E.S...H...T...M.P...VLLCEH...T...M.P...VLLCE
sloth_39	...K...S...K...S...E...T.K...D...S...Q...Y...S...T...A...T...S...T...A...T...
sloth_51	...K...S...K...S...T.Y...D...Q...E...D.S...VT...A...T...VT...A...T...
M_Bat_1	...L...K...T...S...K...V...AR...N...T...AR...N...T...
M_Bat_3	...K...T...YL...S...TP.K...YLRY...FG.S...F...T...M...H...YLRY...FG.S...F...T...M...H...
M_Bat_6	...K...T...Y.S...W.S...K...M.C.G.D...T...M.C.G.D...T...
M_Bat_8	.L...D...K.K...T...S...D...K...W.H...S...I...A.H.M...N...L...W.H...S...I...A.H.M...N...L...
SN_bat_5	E...K...T...Y...E.G.C...L.R...L.Y.SSS.S...HT...MA...L.Y.SSS.S...HT...MA...
SN_bat_9	...W.TQ.LK...F...Y.Y.I...N...C.S...T...V...D.T...Y.Y.I...N...C.S...T...V...D.T...
SN_bat_22	...R...T...A.AG.Y...AV...D...T...R...A.AG.Y...AV...D...T...R...
SN_bat_25	...K...P...S...Y...G...S...Y...S...TT...R...I...S...Y...S...TT...R...I...
SN_bat_29	...K...T...T...S...K...VS...D...T...R...G...T...S...K...VS...D...T...R...G...
SN_bat_31	...K...I.S...P.S.S...S.H.I.VQQ.KF.SD.S...G...HT...P...P...I.GA...S.H.I.VQQ.KF.SD.S...G...HT...P...P...I.GA...
HS_bat_8	...K...T...S...Q...V.DD.SIA.NSA...T.VT...K...P...S...Q...V.DD.SIA.NSA...T.VT...K...P...
HS_bat_15	...K...T...S...Q...V.DD.SIA.NSA...T.VT...K...P...S...Q...V.DD.SIA.NSA...T.VT...K...P...
HS_bat_21	...K...T...L...D.AI...D.S.A.N.A...K...L...D.AI...D.S.A.N.A...K...
HS_bat_22	...K...T...R.T.D...E.S.N...T...R.T.D...E.S.N...T...
HS_bat_24	...K...T...L...D.AI...D.S.A.N.A...K...L...D.AI...D.S.A.N.A...K...
HS_bat_41	...K...T...D...D.AI...D.SIA.NSA...T...K...D...D.AI...D.SIA.NSA...T...K...
dolphin_3	..H...K.Q...T...S.V...R.D...M.I...A.A.N...HT...S.P...Y...I...T...S.V...R.D...M.I...A.A.N...HT...S.P...Y...I...T...
dolphin_12	..H...R.K.Q...T...S.V...R.D...T.I...A.A.N...HT...S.P...Y...I...T...S.V...R.D...T.I...A.A.N...HT...S.P...Y...I...T...
dolphin_19	..H...K.Q...T...S.V...R.D...M.I...A.A.N...HT...S.P...Y...I...T...S.V...R.D...M.I...A.A.N...HT...S.P...Y...I...T...
red_sq_9	QVQLQESGPGLVKPSQSLSLTCAVSGYSIS-SGYC	WSWIRQPPGKGLEWIGIICSGGSTYYSPSLKSRASISRDTSKNQFSLQLSSLTQTATYYCAR
grey_sq_3T...P...F...H...N...V.V...T...P...F...H...N...V.V...
grey_sq_10T...FT...H...N...G.N...N.V...H...T...FT...H...N...G.N...N.V...H...
grey_sq_17	DL..LK.VS...I..F..T..EF..TT.Y...N..C.LT...L.YYY...N...G...I..F..T..EF..TT.Y...N..C.LT...L.YYY...N...G...
red_sq_1ET...T...F..TT...H.M.E...H.P...C...ET...T...F..TT...H.M.E...H.P...C...
red_sq_10T...A...H...N...V...T...A...H...N...V...
red_sq_11T...F..TT..A.H...NY..G.N...L...W.N.G...T...F..TT..A.H...NY..G.N...L...W.N.G...
red_sq_16	...A...T.I...S...K...VC...N.NN...V...W.N.G...	...A...T.I...S...K...VC...N.NN...V...W.N.G...
red_sq_20	..K..W.T...ET...Y.FT.TT.Y...N...T...C.YD...N...	..K..W.T...ET...Y.FT.TT.Y...N...T...C.YD...N...
red_sq_26R...ETM...S...R...ETM...S...
red_sq_36	.M...C...T...HL...C.YDS.N...V...S...IF...V...	.M...C...T...HL...C.YDS.N...V...S...IF...V...
red_sq_34T...T...S...TTI...S...C.NY...G...Y...A...T...T...S...TTI...S...C.NY...G...Y...A...

