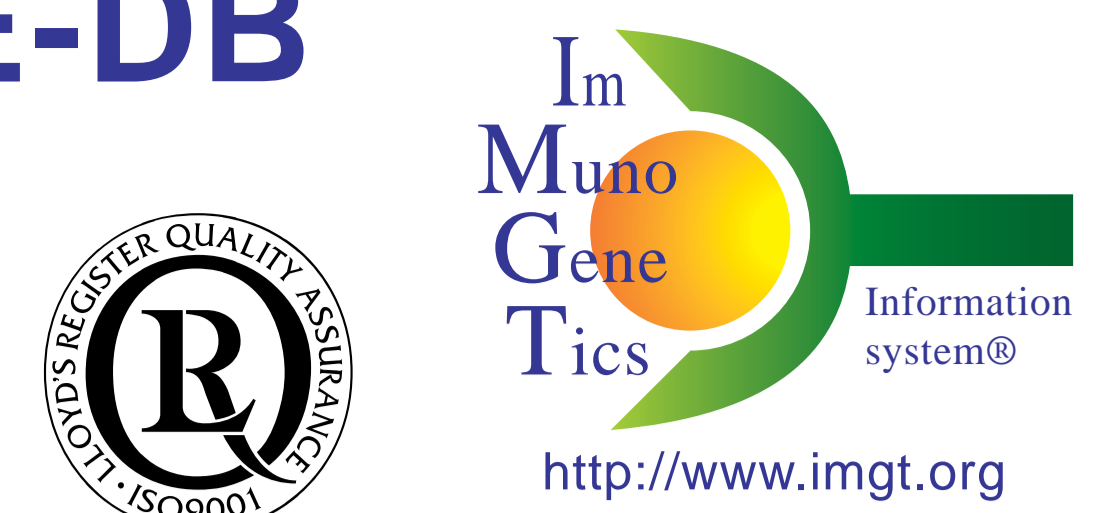


IMGT® biocuration of IG and TR in IMGT/LIGM-DB and IMGT/GENE-DB

Géraldine Folch*, Joumana Michaloud*, Marine Peralta, Mélanie Arrivet, Imène Chentli, Mélissa Cambon, Pascal Bento, Souphatta Sasorith, Typhaine Paysan-Lafosse, Patrice Duroux, Véronique Giudicelli, Sofia Kossida* and Marie-Paule Lefranc*

*Equal contribution

Université Montpellier and CNRS, Laboratoire d'ImmunoGénétique Moléculaire (LIGM), Institut de Génétique Humaine (IGH), UPR CNRS 1142, Montpellier (France)



http://www.imgt.org

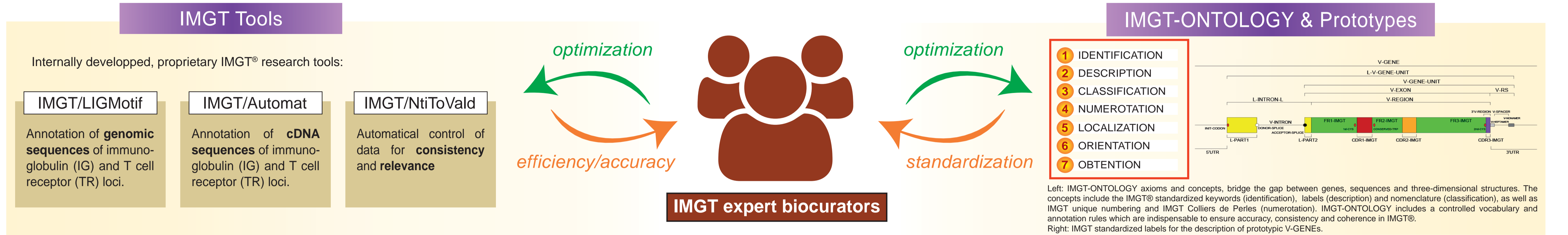
IMGT®, the international ImMunoGeneTics information system, <http://www.imgt.org>, is the global reference in immunogenetics and immunoinformatics [1]. By managing the extreme diversity and complexity of the antigen receptors of the adaptive immune response, the immunoglobulins (IG) or antibodies and the T cell receptors (TR) [2,3] (2.10¹² different specificities per individual), IMGT® is at the origin of immunoinformatics, a science at the interface between immunogenetics and bioinformatics [4]. IMGT® is based on the concepts of IMGT-ONTOLOGY [5] and these concepts are used for expert annotation and standardized knowledge in IMGT/LIGM-DB, the IMGT® database of IG and TR nucleotide sequences from human and other vertebrate species and in IMGT/GENE-DB, the IMGT® gene and allele database. The IMGT/LIGM-DB biocuration pipeline of IG and TR sequences includes IMGT/LIGMotif, for the analysis of large genomic DNA sequences, and IMGT/Automat, for the automatic annotation of rearranged cDNA sequences. Analysis results are checked for consistency, both manually and by using IMGT® tools (IMGT/NiToVald, IMGT/V-QUEST, IMGT/BLAST, etc.). The annotated sequences are integrated in IMGT/LIGM-DB and include the sequence identification (IMGT® keywords), the gene and allele classification (IMGT® nomenclature), the constitutive and specific motif description (IMGT® labels in capital letters, no plural), the translation of the coding regions (IMGT® unique numbering) [4,5]. For genomic IMGT/LIGM-DB sequences containing either an IG or TR variable (V), diversity (D) or joining (J) gene in germline configuration or a constant (C) gene, the gene and allele information is entered in IMGT/GENE-DB. In parallel, the IMGT® Répertoire is updated (Locus representations, Gene tables and Protein displays (for new genes), Alignments of alleles (for new and/or confirmatory alleles)) and the IMGT® reference directory [1,4] is completed (sequences used for gene and allele comparison and assignment in IMGT® tools (IMGT/V-QUEST, IMGT/HighV-QUEST for next generation sequencing (NGS), IMGT/DomainGapAlign) and databases (IMGT/2Dstructure-DB, IMGT/3Dstructure-DB). An IMGT/GENE-DB entry also provides information on the rearranged cDNA and gDNA entries (with links to IMGT/LIGM-DB) and on the three-dimensional structures (with links to IMGT/3Dstructure-DB). IMGT/GENE-DB is the official repository of IG and TR genes and alleles. IMGT® gene names were approved by HGNC and endorsed by WHO-IUIS, the World Health Organization (WHO)-International Union of Immunological Societies (IUIS) Nomenclature Subcommittee for IG and TR. Reciprocal links exist between IMGT/GENE-DB and HGNC, NCBI and Vega. The definition of antibodies published by the WHO International Nonproprietary Name (INN) Programme is based on the IMGT® concepts [6], and allows easy retrieval via IMGT/mAb-DB query [1,4]. The IMGT® standardized annotation has allowed to bridge the gaps for IG or antibodies and TR between fundamental and medical research, veterinary research, repertoire analysis, biotechnology related to antibody engineering, diagnostics and therapeutical approaches.

[1] Lefranc M-P et al. Nucleic Acids Res 43:413-422 (2015) PMID: 25378316,
[2] Lefranc M-P, Lefranc G. The Immunoglobulin FactsBook (2001),

[3] Lefranc M-P, Lefranc G. The T cell receptor FactsBook (2001),
[4] Lefranc M-P. Front Immunol 5:22 (2014) PMID: 24600447,

[5] Giudicelli V, Lefranc, M-P. Front Genet 3:79 (2012) PMID: 22654892,
[6] Lefranc M-P. mAbs 3(1):1-2 (2011) PMID: 21099347.

IMGT® Expert Biocuration Pipeline



IMGT/LIGM-DB

177 115 sequences
351 species

IMGT/LIGM-DB includes all germline (non-rearranged) and rearranged IG and TR genomic DNA and complementary DNA sequences published in generalist databases. IMGT/LIGM-DB allows searches from the Web interface according to biological and immunogenetic criteria. For a given entry, nine types of display are available including the IMGT flat file, the translation of the coding regions and the analysis by the IMGT/V-QUEST tool. The annotations hugely enhance the quality and the accuracy of the distributed detailed information.

IMGT/GENE-DB

3 927 genes
5630 alleles
24 species

IMGT/GENE-DB Query Page allows the search of IG/TR genes according to IMGT-ONTOLOGY's seven axioms. IMGT/GENE-DB entry displays accurate gene data related to genome (gene localization), allelic polymorphisms (number of alleles, IMGT reference sequences, functionality, etc.) gene expression (known cDNAs) and protein structures (IMGT Colliers de Perles, IMGT/3Dstructure-DB). It provides internal links to the IMGT sequence databases and the IMGT Web resources as well as external links to genome and generalist sequence databases.

Web Resources

1 IMGT flat file

1 IDENTIFICATION: Keywords

genomic-DNA=MoleculeType
germline=ConfigurationType
regular=StructureType
functional=Functionality
Homo sapiens=Taxon
Ig-Heavy-Mu=ChainType
variable=GeneType

2 DESCRIPTION: Labels

V-GENE=Entity
V-REGION=CoreRegion
FR1-IMGT=SubRegion

3 CLASSIFICATION: Nomenclature

IGHV=Group
IGHV3=Subgroup
IGHV3-66=Gene
IGHV3-66*04=Allele

4 NUMERATION

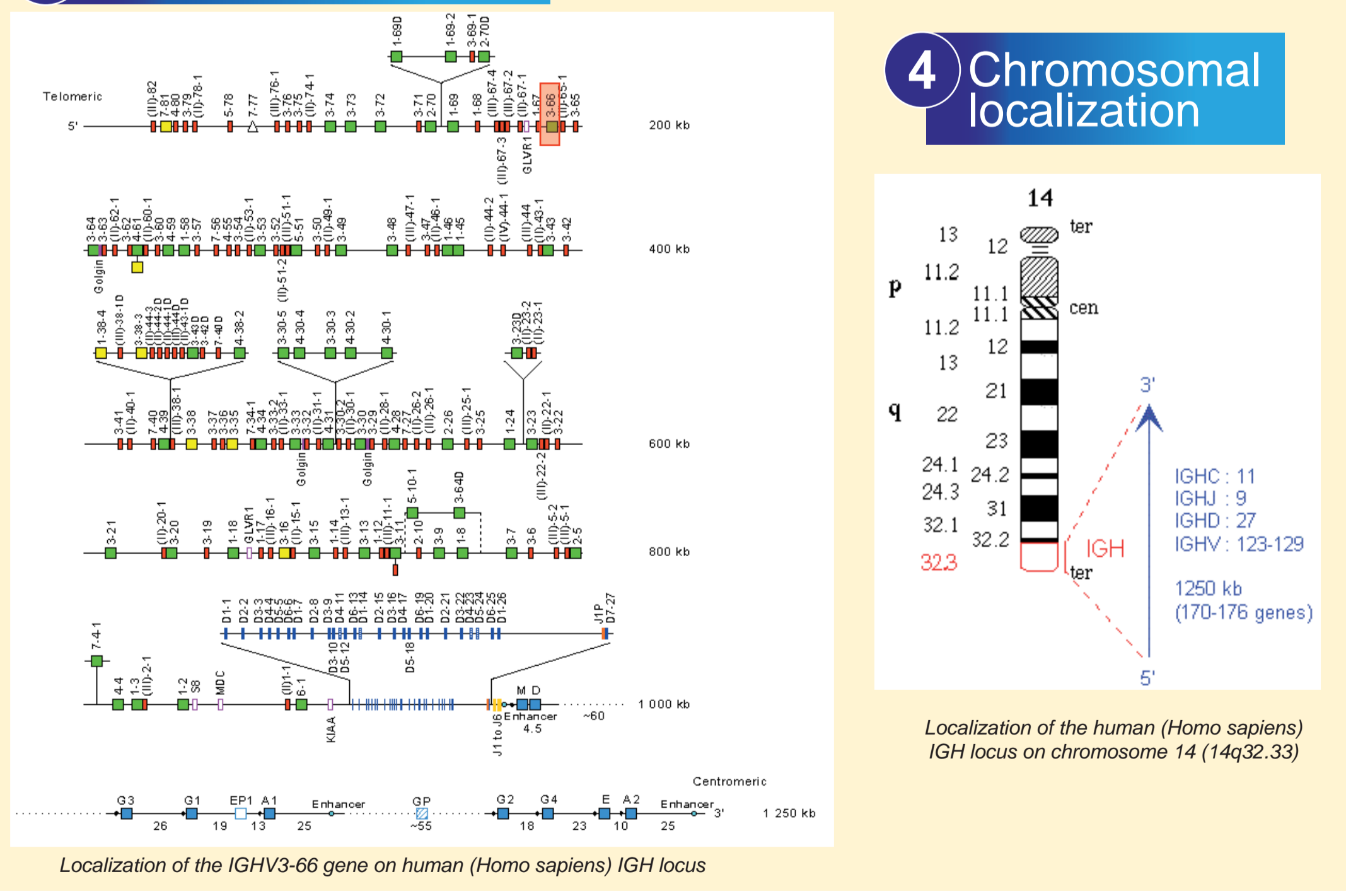
(8.7.2)-V-REGION CDR lengths
1 to 26. AA.10 is missing-AA IMGT numbering

```
ID X70208; SV 1; linear; gDNA; STD; HUM; 830 BP.  
XX X70208; X68844;  
XX  
DT 15-MAY-1995 (Rel. 1995S2-1, arrived in LIGM-DB)  
DT 31-MAR-2015 (Rel. 2015J4-2, Last updated, Version 19)  
XX  
DE H.sapiens DNA for Igh heavy chain (MTGLA)  
XX  
KW antigen receptor; Immunoglobulin superfamily (IgSF); immunoglobulin (IG);  
KW Ig-Heavy; Ig-Heavy-Mu; variable; lymphoma; IMGT reference sequence;  
KW regular; gDNA; germline; functional; V-gene.  
CC IMGT/LIGM-DB annotation level: by annotators  
FH Key Location/Qualifiers  
FH  
FT V-GENE 1..830  
FT /IMGT_allele="IGHV3-66*04"  
FT /IMGT_gene="IGHV3-66"  
FT /ID_sref="taxon:9606"  
FT /mol_type="genomic DNA"  
FT /organism="Homo sapiens"  
FT 1..291  
FT 5'UTR 292..337  
FT L-PART1 /translation="MEFGLSWFLVALLK"  
FT 292..324  
FT INTRON-SPICE 337..339  
FT 338..438  
FT V-INTRON /codon_start=3  
FT ACCEPTOR-SPICE 439..440  
FT 439..742  
FT V-EXON /translation="VQEVQLVSGGGLVQPQGLLRLSCAASGFTVSSNYSWVR  
FT QAPQGLVSWVSYSGGTYVADYVGRFTISRNKNTLYLQMSLRADYVYCA*  
FT YCA*  
FT L-PART2 439..449  
FT /codon_start=3  
FT /translation="VQC*"  
FT V-REGION 450..742  
FT /translation="EVQLVSGGGLVQPQGLLRLSCAASGFTVSSNYSWVRQAP  
FT GKLEWVSYSGGTYVADYVGRFTISRNKNTLYLQMSLRADYVYCA*  
FT R*  
FT /IGV_lengths="(8,7,2)"  
FT /IMGT_allele="IGHV3-66*04"  
FT /IMGT_gene="IGHV3-66"  
FT 450..524  
FT FR1-IMGT /translation="VQA...to 26; AA.10 missing"  
FT 5'UTR 741..830  
SQ Sequence 830 BP; 203 A; 292 C; 237 G; 188 T; 0 other;  
cctaaatgaa taccaggca cactcaata atataaatt atatttctt gaatgtagg 60  
ataatacca atctctccc agggacctt catctgact agccaccgct cttctctag 120  
ctgtgatta ctgtgtaga caccactga gggagccca ttgtgccc agacacac 180  
ctctctcga ggaatctg ggaatcag gcccggggc cttcaggag 240  
830
```

2 Gene table

IMGT gene name	IMGT allele name	IGT	Chromosomal localization	a	b	Positions in the locus	IMGT/LIGM-DB reference sequence	Positions in the reference sequence (from IMGT V-REGION)	Accession numbers	IMGT/LIGM-DB sequences from the literature (from the end of V-REGION)	Positions in the reference sequence (from the end of V-REGION)
IGHV3	IGHV3-66*01	F	14q32.33	-	-	150000-160000	IGHV3-66	150000-160000	U08182	U08182	150000-160000
	IGHV3-66*02	F	14q32.33	-	-	150000-160000	IGHV3-66	150000-160000	U08182	U08182	150000-160000
	IGHV3-66*03	F	14q32.33	-	-	150000-160000	IGHV3-66	150000-160000	U08182	U08182	150000-160000
	IGHV3-66*04	F	14q32.33	-	-	150000-160000	IGHV3-66	150000-160000	U08182	U08182	150000-160000

3 Locus representation

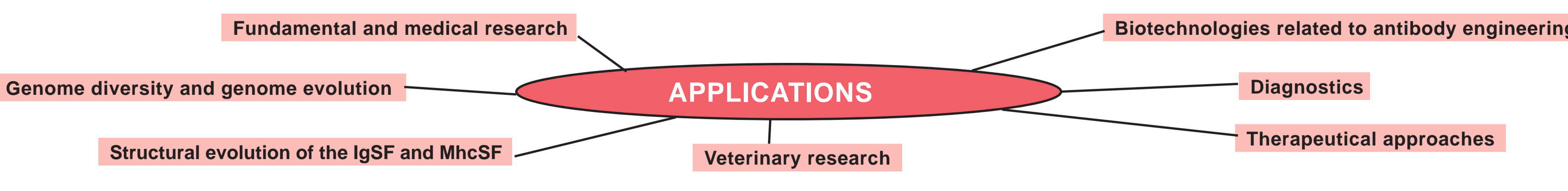


5 External links



INTEROPERABILITY

WEB INTERFACE



IMGT® director: Sofia Kossida (Sofia.Kossida@igh.cnrs.fr)
IMGT® founder and executive director emeritus: Marie-Paule Lefranc (Marie-Paule.Lefranc@igh.cnrs.fr)
Bioinformatics manager: Véronique Giudicelli (Veronique.Giudicelli@igh.cnrs.fr)
Computer manager: Patrice Duroux (Patrice.Duroux@igh.cnrs.fr)

