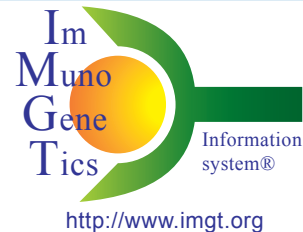


IMGT/Automat: the strategy for the annotation of human and mouse cDNA nucleotide sequences of IG and TR

G eraldine Folch, Joumana Jabado-Michaloud, Fatena Bellahcene, Laetitia Regnier, V eronique Giudicelli and Marie-Paule Lefranc



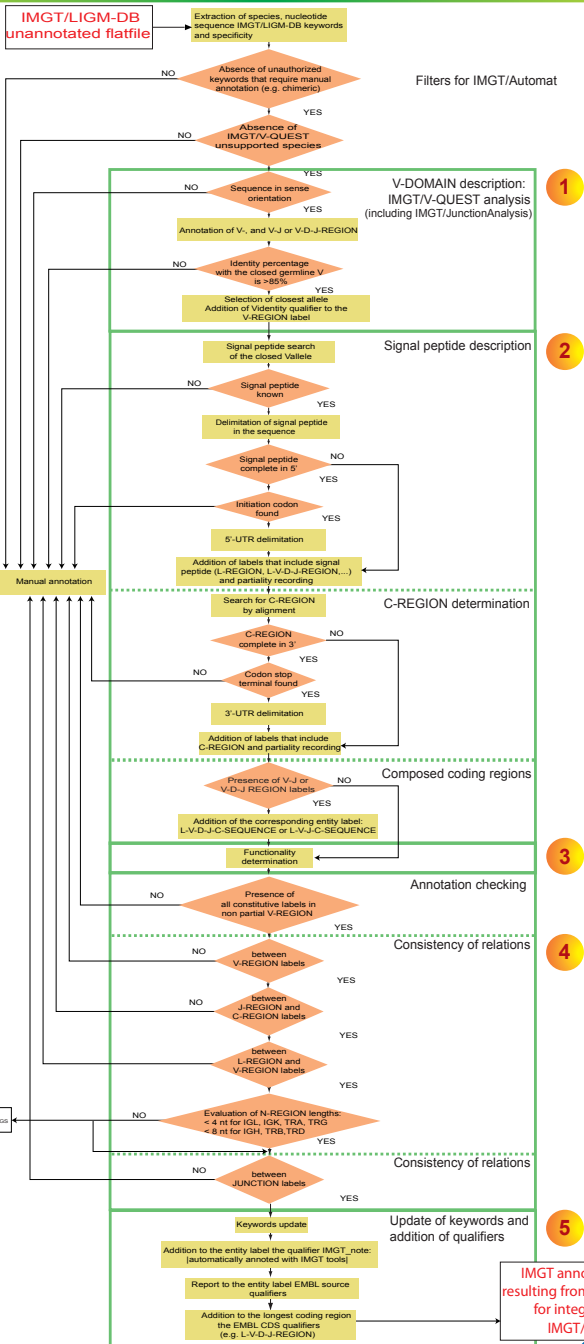
IMGT , the international ImMunoGeneTics information system , Laboratoire d'ImmunoG en tique Mol culaire LIGM, Universit  Montpellier 2, Institut de G n tique Humaine IGH, UPR CNRS 1142, 141 rue de la Cardonille, F-34396 Montpellier cedex 05, France

<http://www.imgt.org>

The cDNA sequences of immunoglobulins (IG) and T cell receptors (TR) represent more than one half of the sequences in the IMGT  nucleotide database IMGT/LIGM-DB [1] and 75% of them are from human and mouse. A few cDNA are germline but the great majority results from a V-D-J or V-J gene rearrangement, spliced to a C gene. The IG and TR genes have been studied extensively in IMGT  (<http://www.imgt.org>) [2], which allowed to set up their nomenclature and the corresponding germline reference sequences. These standardized reference directories sets (one for each group of each locus) and the IMGT-ONTOLOGY axioms and derived concepts [3] are the key elements indispensable to perform the annotation of IG and TR cDNA sequences. A Java program, IMGT/Automat [4], was developed by IMGT , to automatically annotate the IG and TR cDNA sequences and to produce a totally automatic and complete annotation. More than 9,000 human and mouse cDNA have already been successfully automatically annotated. The quality of the cDNA automatic annotation is equivalent to the quality of the annotation achieved by a human expert. The IMGT  strategy is currently the only way, in the field of immunogenetics, to guarantee the annotation quality and the management of an always increasing number of IG and TR cDNA nucleotide sequences.

IMGT/Automat includes five main tasks:
In a first step IMGT/Automat implements IMGT/V-QUEST [5]. The description of the V-D-J junction is performed by the IMGT/JunctionAnalysis [6] tool. In a second step, IMGT/Automat delimits the signal peptide, the constant region and the composed coding regions (for example: L-V-D-J-C-REGION). In a third step, the functionality of the sequence (a concept of identification) is defined. The fourth step corresponds to a thorough annotation checking. In a fifth and final step, keywords are updated and qualifiers on biological origin and methodology used (concepts of obtention) are integrated, and the annotated flat file is generated.

IMGT/Automat main tasks



1 V-DOMAIN description: IMGT/V-QUEST Analysis (including IMGT/JunctionAnalysis)

V-DOMAIN description (V-J-REGION and V-D-J-REGION) is performed by IMGT/V-QUEST analysis. Detailed analysis of JUNCTION is performed by the integrated IMGT/JunctionAnalysis tool

Alignment for V-GENE and allele identification

```

BC024289          <-----FR1-IGHJ3-21*01
AB019439          gaggtgcagctggtggagctgtgggg...ggctgtgcagga
M99658  IGHV3-21*02          -----a-----t-t-ac-
M99675  IGHV3-48*01          -----t-----t-ac-
AB019438  IGHV3-48*02          -----t-----t-ac-
AJ879484  IGHV3-h*01 (P)      -----t-----t-a-
    
```

Alignment for J-GENE and allele identification

```

BC024289          tctccgccagctaacctcctactgtgacttcgatctctggggg
J00256  IGHJ2*01             tctccgccagctaacctcctactgtgacttcgatctctggggg
M25625  IGHJ4*03             -----a-----t-t-ac-
    
```

Results of IMGT/JunctionAnalysis

Maximum number of accepted mutations in 3V-REGION = 2, D-REGION = 4, 5J-REGION = 2

Input	V-NAME	3V-REGION	N1	D-REGION	N2
BC024289	IGHV3-21*01	tgtgggagaga t		tctgggagcta ... acttt	

Input	5J-REGION	J name	D name	Vmut	Dmut	Jmut	Ngc
BC024289	ctactgtgacttcgatctctgg	IGHJ2*01	IGHJ3-10*01	0	4	0	3/7

- Identification** of the sequence (chain type for ex: IG-Heavy)
- Classification** of the V, D, J genes and alleles
- Description** of the IG and TR specific constitutive motifs
- Delimitation** of the framework regions (FR-IMGT) and complementarity determining regions (CDR-IMGT)
- Numbering** of the codons

2 Signal peptide, C-REGION and composed coding regions

Signal peptide, C-REGION and composed coding regions description is performed using the L-V-J-C-SEQUENCE and L-V-D-J-C-SEQUENCE prototypes.

- Identified regions:** L-V-D-J-C-SEQUENCE, L-REGION, V-D-J-REGION, L-V-D-J-C-REGION, V-D-J-C-REGION, C-REGION
- Composed coding regions for:**
 - Entity: L-V-D-J-C-SEQUENCE
 - Entity: L-REGION, L-V-D-REGION, L-V-D-J-REGION, L-V-D-J-C-REGION, V-D-J-C-REGION, C-REGION
- Composed coding regions for:**
 - Entity: L-V-J-C-SEQUENCE
 - Entity: L-V-J-REGION, L-V-J-C-REGION, V-J-C-REGION, C-REGION

3 Functionality determination

The functionality of the sequence is defined according to the biological rules of the IMGT Scientific chart. The sequence is **PRODUCTIVE** if the coding region has an open reading frame, with no stop codon and no defect described in the initiation codon, splicing sites and/or regulatory elements, and an in-frame JUNCTION. The sequence is **UNPRODUCTIVE** if the JUNCTION is out-of-frame and/or the presence of stop codon(s) and/or frameshift mutation(s), and/or a defect described in the splicing sites and/or the regulatory element(s), and/or unusual features (TRANSLOCATED, GENE FUSION...) and/or changes of conserved

4 Annotation checking

Annotation checking comprises several steps (see figure), for examples:
Presence of all constitutive labels by comparison with the prototype (e.g. I-REGION, V-REGION, D-REGION...)
Consistency of relations between labels (e.g. L-REGION adjacent_in_its_3_prime_with V-REGION, FR1-IMGT is_included_with_same_5_prime_in V-REGION)

5 Annotated IMGT FLAT-FILE resulting from IMGT/Automat

ID	BC024289	IMGT/LIGM annotation : automatic, mRNA, HUM, 1630 BP	FT V-D-J-REGION	121_486
XX	BC024289		FT	translatability="EVLVESGGGKPKGGLRSCAASGFTTSSYMMWRRQAP
XX			FT	GKGLVWSSSSSSSYTYADSWKGRFTRSDNANKSLYQMSLRAEDTAYVYC
XX			FT	ASGLRLSITVYFDLWGRGLTVVSS
XX			FT	121_416
XX			FT	allele="IGHJ3/21*01"
XX			FT	igene="IGHJ3/21*01"
XX			FT	Viterbi="95.31%
XX			FT	(296288 nt)"
XX			FT	CDS_language="1.16"
XX			FT	putative_ImMTC="3 sides"
XX			FT	translatability="EVLVESGGGKPKGGLRSCAASGFTTSSYMMWRRQAP
XX			FT	GKGLVWSSSSSSSYTYADSWKGRFTRSDNANKSLYQMSLRAEDTAYVYC
XX			FT	ASGLRLSITVYFDLWGRGLTVVSS
XX			FT	121_195
XX			FT	AA_IMGT="AA 66 to 65, AA 60 to 61 missing"
XX			FT	translatability="EVLVESGGGKPKGGLRSCAAS"
XX			FT	184_186
XX			FT	translatability="EVLVESGGGKPKGGLRSCAAS"
XX			FT	186_199
XX			FT	AA_IMGT="AA 27 to 38, AA 31, 32, 33, 34 are missing"
XX			FT	220_270
XX			FT	AA_IMGT="AA 66 to 65, AA 60 to 61 missing"
XX			FT	220_270
XX			FT	AA_IMGT="AA 66 to 65, AA 60 to 61 missing"
XX			FT	220_270
XX			FT	translatability="EVLVESGGGKPKGGLRSCAAS"
XX			FT	295_439
XX			FT	AA_IMGT="AA 66 to 104, AA 103 to 101 missing"
XX			FT	translatability="YADSWKGRFTRSDNANKSLYQMSLRAEDTAYVYC"
XX			FT	400_400
XX			FT	400_433
XX			FT	AA_IMGT="AA 105 to 117 including 112, 111, 111"
XX			FT	translatability="EVLVESGGGKPKGGLRSCAAS"
XX			FT	400_496
XX			FT	JUNCTION
XX			FT	in_frame
XX			FT	translatability="CARDLRLTYSYFDLWGRGLTVVSS"
XX			FT	3V-REGION
XX			FT	400_416
XX			FT	N1-REGION
XX			FT	codon_start=2
XX			FT	133_423
XX			FT	D-REGION
XX			FT	allele="IGHJ3/10*01"
XX			FT	igene="IGHJ3/10*01"
XX			FT	133_423
XX			FT	codon_start=3
XX			FT	N2-REGION
XX			FT	400_434
XX			FT	translatability="I"
XX			FT	435_496
XX			FT	J-REGION
XX			FT	435_496
XX			FT	allele="IGHJ3/10*01"
XX			FT	igene="IGHJ3/10*01"
XX			FT	435_496
XX			FT	putative_ImMTC="3 sides"
XX			FT	Viterbi="100.00%
XX			FT	(8353 nt)"
XX			FT	translatability="YFDLWGRGLTVVSS"
XX			FT	454_456
XX			FT	J-TRP
XX			FT	FR4-IMGT
XX			FT	454_456
XX			FT	AA_IMGT="AA 118 to 128"
XX			FT	translatability="YFDLWGRGLTVVSS"
XX			FT	487_1476
XX			FT	C-REGION
XX			FT	allele="IGHJ1/2"
XX			FT	igene="IGHJ1/2"
XX			FT	translatability="ASTKPKPSPVLPAPSKSSTGGTAALGCLVDVYFPEVPSVW
XX			FT	NSCALTSYFTHYDFLWGRGLTVVSSSSTGGTAALGCLVDVYFPEVPSVW
XX			FT	KYKPKKSKDKHTKPCQAPAEKLVSSYFPKPKDKITLSRTPKDFLMSRTPEVLCQVDVSH
XX			FT	EDPEVFWYDGVVDDVHNAKPKPKEQYSSYTRVSLVTLKQVSLKQVSTKPKK
XX			FT	SKMALKAPKTSKSGKAPKPKPVDFYLPSSDLTQVSLVTLKQVYFVSDA
XX			FT	VEVSEVSPKNTKPKTTFVLSLVDGSLFPVSKLTVKSRDGGWVFCSSVMHEAL
XX			FT	NIPTKQSLSPK"
XX			FT	STOP-CODON
XX			FT	1477_1479