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

Géraldine Folch, Joumana Jabado-Michaloud, Fatena Bellahcene, Laetitia Regnier, Véronique Giudicelli and Marie-Paule Lefranc

IMGT®, the international ImMunoGeneTics information system®, Laboratoire d'ImmunoGénétique Moléculaire LIGM, Université Montpellier 2, Institut de Génétique Humaine IGH, UPR CNRS 1142, 141 rue de la Cardonille, F-34396 Montpellier cedex 05, France

**Im Muno Gene Tics**
Information system®
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 directory 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.

[1] Giudicelli V. et al. Nucleic Acids Res., 34, D781-784 (2006).
[2] Lefranc M.-P. et al. Nucleic Acids Res., 37, D1006-1012 (2009).
[3] Duroux P. et al. Biochimie, 90, 570-583 (2008).
[4] Giudicelli V. et al. Stud. Health Technol. Inform., 116, 3-8 ( 2005).
[5] Brochet X. et al. Nucleic Acids Res., 36, W503-508 (2008).
[6] Yousfi Monod M. et al. Bioinformatics, 20, i379-385 (2004).

**IMGT/Automat includes five main tasks:**
In a first step IMGT/Automat implements IMGT/V-QUEST [5] . The description of the V-D-J and V-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.

## 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

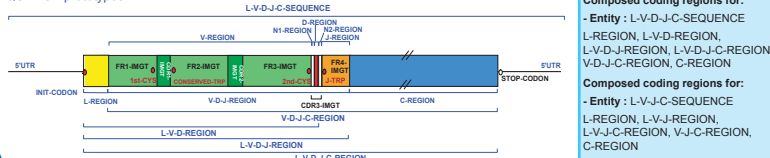**IMGT/V-QUEST analysis provides:**
- **Identification** of the sequence (chain type for ex: IG-Heavy)
- **Classification** of the V, D, J genes and alleles
- **Description** of the IG and TR specific constitutive motifs
- **Delimitation** of the framework regions (FR-IMGT) and complementary determining 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.

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



### IMGT/Automat main tasks