

IMGT/HighV-QUEST for NGS analysis of IG and TR: statistical analysis of IMGT clonotypes (AA), novel interface and functionalities

Marianne Lèbre, Karthik Kalyan, Patrice Duroux, Véronique Giudicelli, Sofia Kossida, Marie-Paule Lefranc

IMGT®, the international ImMunoGeneTics information system®, Laboratoire d'ImmunoGénétique Moléculaire (LIGM), Institut de Génétique Humaine (IGH), UMR 9002 CNRS-UM, Université de Montpellier (UM), Montpellier (France)

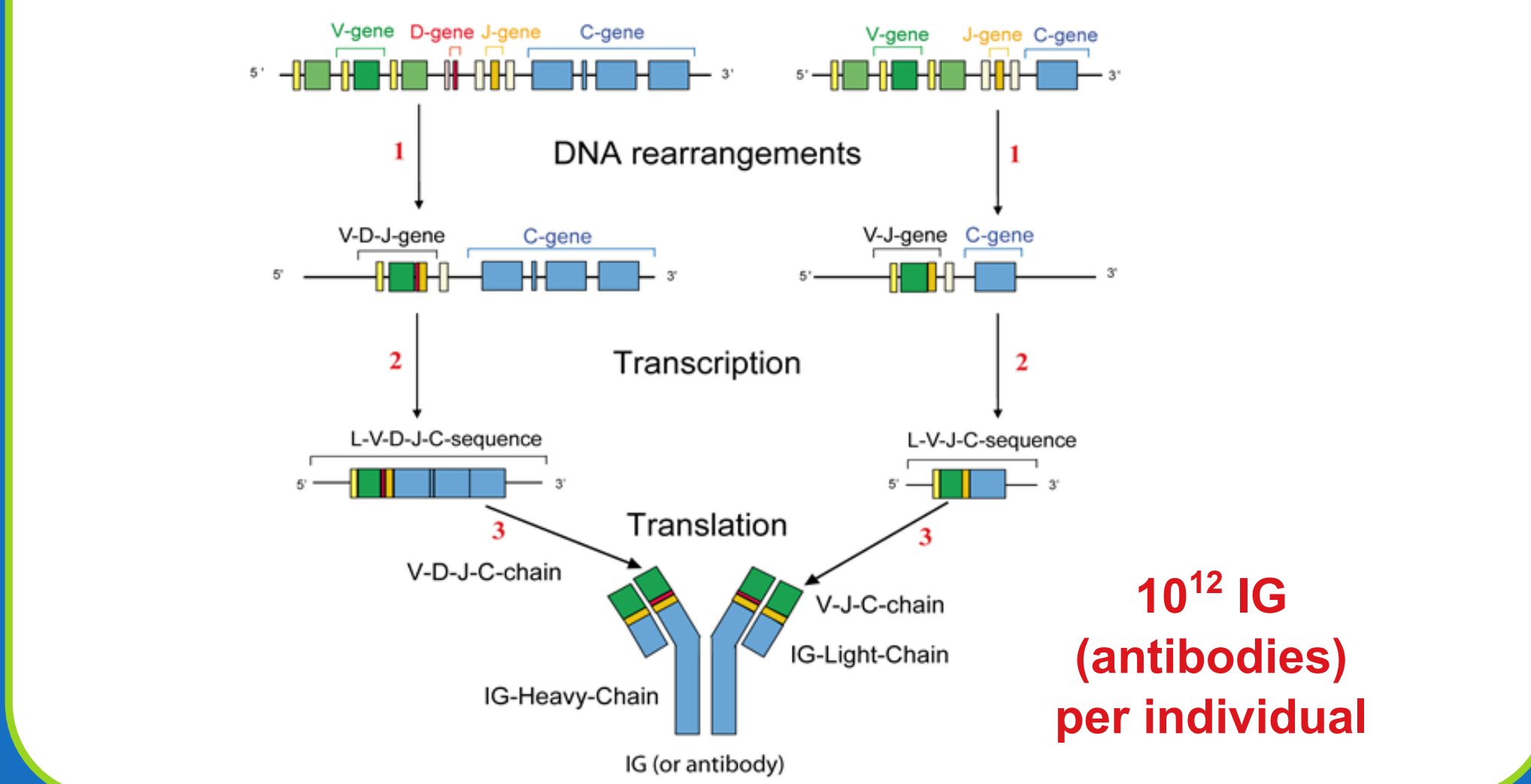


IMGT®, the international ImMunoGeneTics information system®, <http://imgt.org/> [1], is the global reference in immunogenetics and immunoinformatics [2], founded in 1989 by Marie-Paule Lefranc at Montpellier (Université de Montpellier and CNRS). IMGT® is a high-quality integrated knowledge resource specialized in the immunoglobulins (IG) or antibodies, T cell receptors (TR), major histocompatibility (MH) of humans and other vertebrate species. IG and TR are the antigen receptors for the adaptive immune response which characterizes the vertebrates with jaws (*Gnathostomata*) [2]. Their study in normal and pathological conditions is a challenge due to the huge diversity of the variable domain (V-DOMAIN) at the N-terminal end of each chain (10^{12} potential specificities for humans), which results from complex IG and TR synthesis. Since 2010, IMGT® has developed IMGT/HighV-QUEST [3-7], an online portal for the analysis of the IG and TR immune repertoires obtained through the next generation sequencing (NGS) technologies.

[1] Lefranc M.-P. et al. Nucl. Acids Res. 43:D413-422 (2015) PMID: 25378316 [2] Lefranc M.-P. Front. Immunol. 5:22 (2014) PMID: 24600447 [3] Alamyar E. et al. Abstract 60, Poster 27, JOBIM Montpellier (2010) [4] Alamyar E. et al. Immunome Res. 8:1:2 (2012) [5] Li S. et al. Nat. Commun. 4:2333 (2013) PMID: 23995877 [6] Aouinti S. et al. Front. Immunol. 7:339 (2016) PMID: 27667992 [7] Giudicelli V. et al. BMC Immunol. 18(1):35 (2017) PMID: 28651553

IG and TR synthesis

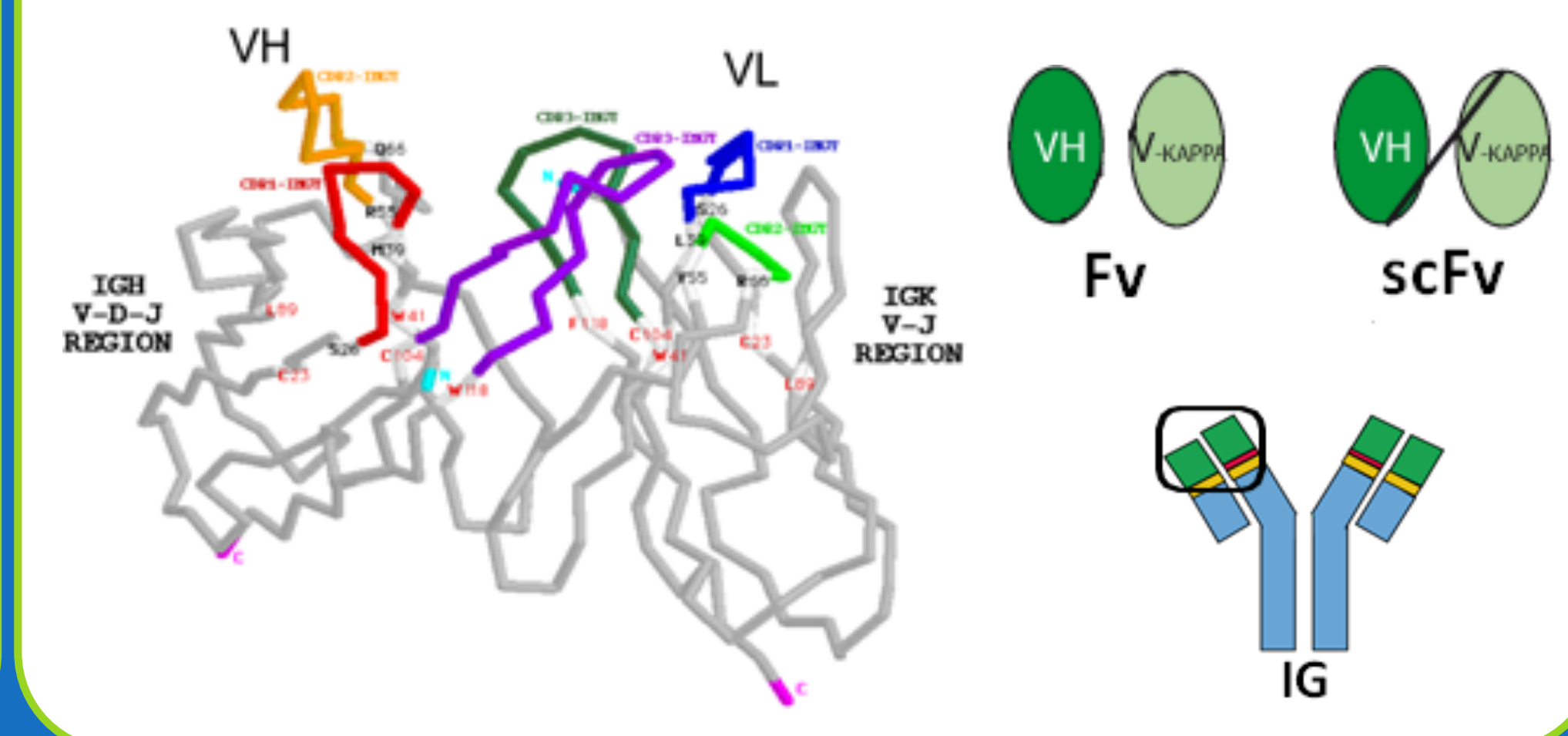
IG and TR are antigen receptors of the adaptive immune response. The synthesis of the V-DOMAIN at the N-terminal end of each IG or TR chain results from genomic DNA rearrangements of variable (V), diversity (D) and joining (J) genes and from junctional diversity.



V-DOMAIN

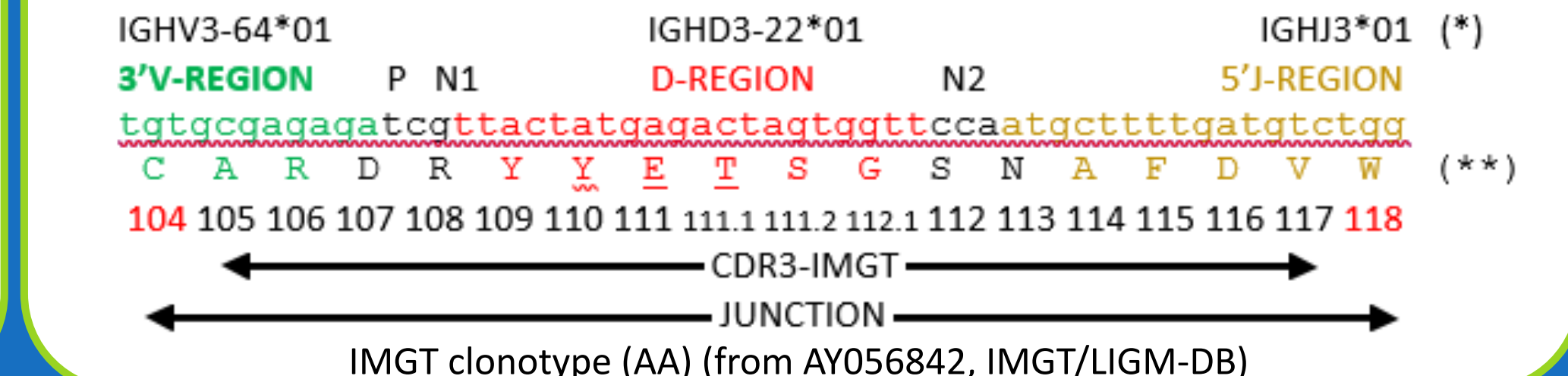
The antigen-binding site is formed by the six complementarity determining regions or CDR-IMGT of two non covalently paired V-DOMAIN (VH/VL for IG, V-ALPHA/V-BETA or V-GAMMA/V-DELTA for TR), defined as Fragment variable (Fv).

An *in vitro* engineered chain made of two V-DOMAIN connected by a linker is a single chain Fragment variable (scFv).



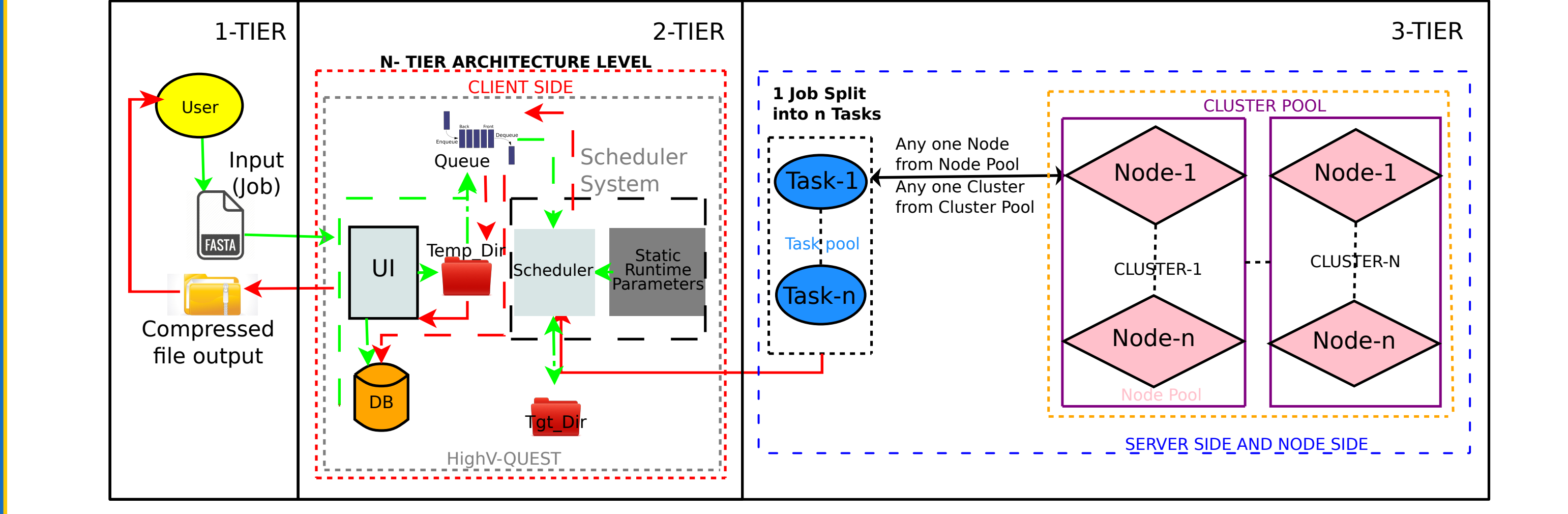
IMGT clonotypes

The 'IMGT clonotype' concept provides a molecular characterization of the IG and TR V-DOMAIN repertoires. IMGT clonotypes are identified by analysis of the nucleotide (nt) sequences of the V-DOMAIN (V-(D)-J-REGION), using IMGT/HighV-QUEST for NGS. An 'IMGT clonotype (AA)' is defined as a unique V-(D)-J rearrangement (with IMGT gene and allele names (*)) and a unique CDR3-IMGT amino acid (AA) sequence (**)



IMGT/HighV-QUEST architecture system

IMGT/HighV-QUEST architecture has moved from a 2-tier to a 3-system: web user interface (UI), database and scheduling- system. The scheduling-system is now a standalone system (shell scripts and cron) which has the possibility to be integrated to an automation tool. The 3-tier architecture enables easier implementation of newly functionalities.



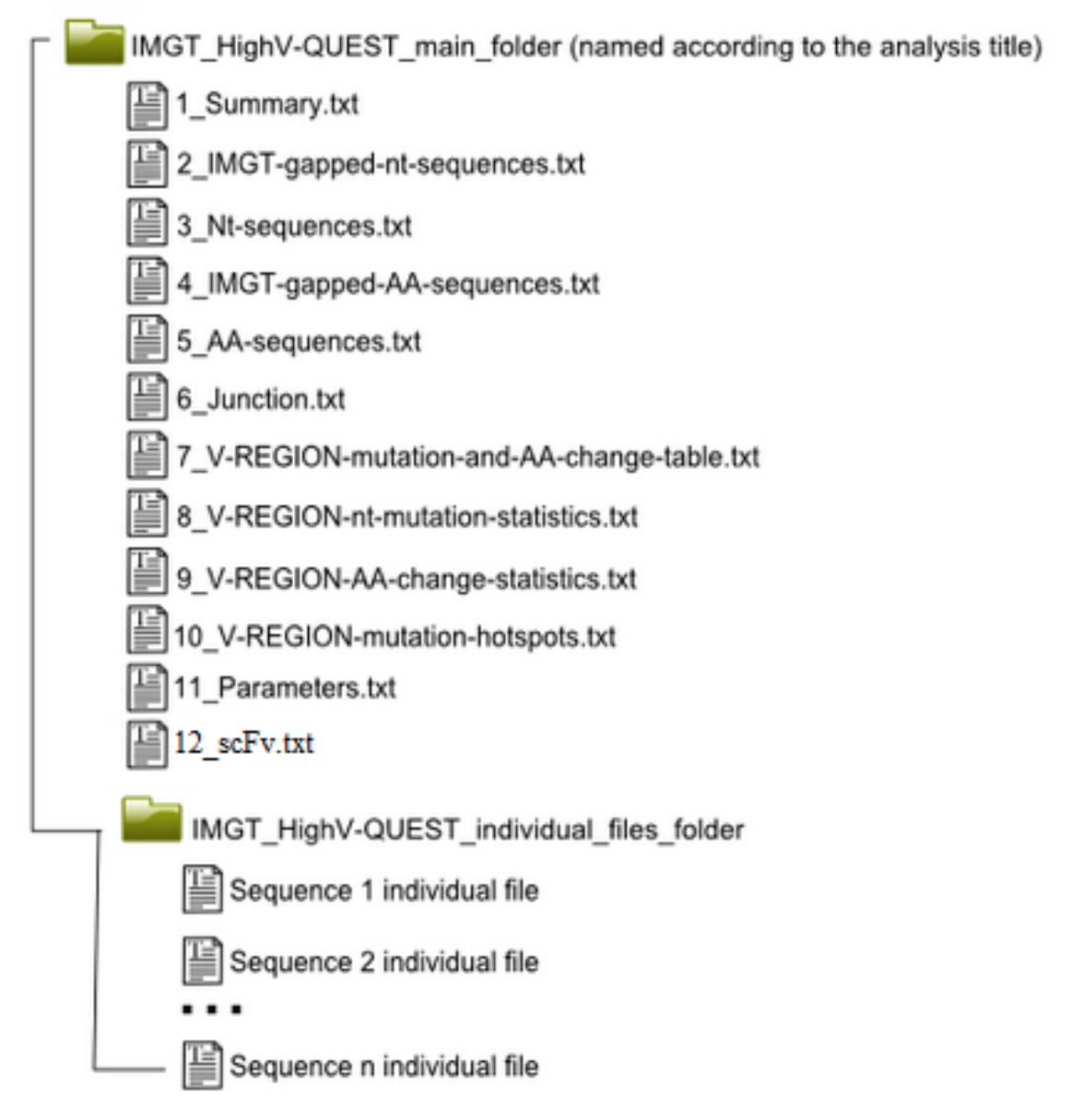
The user can follow the analysis status in the "Analysis history": **queued** (new entry created in the database when data are submitted for a job); **dispatched** (data are sent by the scheduling system for the analysis job when there is enough computational resources available); **completed** (results available and database updated by the scheduling system). A new client UI based on modern web technologies (Bootstrap, Struts2 and Tiles3) has recently been made available. It also integrates a new e-mail update functionality.

IMGT/HighV-QUEST analysis

IMGT/HighV-QUEST analyzes up to 500,000 IG or TR rearranged sequences per run, with the same degree of resolution and high-quality results as IMGT/V-QUEST (same algorithm, same IMGT reference directories). The tool:

- 1) identifies, by default, the insertions/deletions (indels) and correct them
- 2) numbers the user sequences and introduces gaps (IMGT unique numbering)
- 3) identifies the V, D and J GENE and alleles
- 4) characterizes IG somatic hypermutations (nt and AA)
- 5) describes the JUNCTION (IMGT/JunctionAnalysis)
- 6) provides a complete annotation with IMGT labels (IMGT/Automat)
- 7) as an option, analyses scFv (added in 2017 [7]).

IMGT/HighV-QUEST results consist of 11 CSV files (or 12 with the scFv option) provided as an archive file. Files #1 to #10 comprise systematically sequence identification, i.e. the sequence name, the functionality, the names of the closest V-GENE and allele, and files #1 to #6 also include the D and J GENE and alleles. The files #7 to #10 that report the analysis of mutations are used mostly for IG. Files #1 to #10 include one line per submitted sequence, and together may comprise up to 539 columns for a complete results report.



IMGT/HighV-QUEST statistical analysis: IMGT clonotypes (AA) and (nt)

The statistical analysis applies a filter on the IMGT/HighV-QUEST results: only the ones characterized by a V-GENE and allele (single or several alleles), a JUNCTION and a J-GENE and allele (single or several alleles) are filtered-in for statistical analysis. The IMGT/HighV-QUEST statistical analysis, which allows the identification and characterization of the clonotypes [5], may analyse up to one million IMGT/HighV-QUEST results.

A IMGT clonotype (AA and nt) results per locus

- The statistical results are provided in 10 sections (HTML pages):
- IMGT clonotypes (AA) per Nb (1) without or (2) with detailed clonotypes (nt)
- IMGT clonotypes (AA) per V gene (3) without or (4) with detailed clonotypes (nt)
- IMGT clonotypes (AA) per CDR3-IMGT length (AA) (5) without or (6) with detailed clonotypes (nt)
- IMGT clonotypes (AA) by CRD3-IMGT sequence (AA) alphabetical order with detailed clonotypes (nt) (7)
- IMGT clonotype (AA) diversity and expression histograms: per V,(D),J-GENE and per CDR3-IMGT length (8)
- IMGT clonotype (AA) diversity and expression tables: per V, (D), J-GENE and per CDR3-IMGT length (9)
- V gene and allele table: Rearrangements, Nb of sequences and Nb IMGT clonotypes (AA) per V-GENE and allele (10)

ID	Nb	IMGT clonotype (AA) definition	IMGT clonotype (AA) representative sequence	IMGT clonotype (nt)
10000001	1	IGHV3-64*01	IGHV3-64*01	IGHV3-64*01
10000002	1	IGHV3-64*01	IGHV3-64*01	IGHV3-64*01
10000003	1	IGHV3-64*01	IGHV3-64*01	IGHV3-64*01

B IMGT clonotype (AA) results comparison per locus

ID	Nb	IMGT clonotype (AA) definition	IMGT clonotype (AA) representative sequence	IMGT clonotype (nt)
10000001	1	IGHV3-64*01	IGHV3-64*01	IGHV3-64*01
10000002	1	IGHV3-64*01	IGHV3-64*01	IGHV3-64*01
10000003	1	IGHV3-64*01	IGHV3-64*01	IGHV3-64*01

C IMGT/StatClonotype

IMGT/StatClonotype [6] is a tool, downloadable on the IMGT® site, which allows evaluating and exploring, between sets, the significance of pairwise comparison of IMGT clonotype (AA) diversity and expression per V, D and J gene. The IMGT/HighV-QUEST statistical output contains, in 'data' directory, txz file(s) designated as stats_xxx, where 'xxx' is the batch name and the locus type. At least two 'stats_xxx' files are needed to launch a comparative analysis in IMGT/StatClonotype. Integrated in the R package "IMGTStatClonotype", the tool offers a graphical interface to visualize pair wise comparison, per IMGT genes and alleles, of the IMGT clonotype (AA) diversity or expression of any IG or TR immunoprofiles of any species.