

Scientific Project

Master GENIOMHE

2023–2024

Samuel ORTION 

Further development on FTAG Finder, a pipeline to identify
Gene Families and Tandemly Arrayed Genes

Advisors:

Carène RIZZON

Franck SAMSON

Laboratoire de Mathématiques

et Modélisation d'Évry

carene.rizzon@univ-evry.fr

franck.samson@inrae.fr

+33 (0) 1 64 85 35 40

IBGBI

23 Bd. de France

91037 Évry Cedex

keywords: duplicate genes, tandemly arrayed genes, pipeline

Contents

Glossary	6
Acronyms	7
1 Scientific context	9
1.1 Gene duplication mechanisms	9
1.1.1 Whole genome duplication and polyploidisation	9
1.1.2 Unequal crossing-over	9
1.1.3 Retroduplication	11
1.1.4 Transduplication	11
1.1.5 Segment duplication	11
1.2 Fate of duplicate genes in genome evolution	13
1.2.1 Pseudogenization	13
1.2.2 Neofunctionalization	13
1.2.3 Subfunctionalization	13
1.2.4 Functional redundancy	13
1.3 Methods to identify duplicate genes	13
1.3.1 Paralog detection	15
1.3.2 FTAG Finder	15
2 Objectives for the internship	19
2.1 Scientific questions	19
2.2 Extend the existing FTAG Finder Galaxy pipeline	19
2.3 Port FTAG Finder pipeline on a workflow manager	19

List of Figures

1.1	Different types of duplication	10
1.2	Schematic representation of TAG definitions	16

List of Tables

Glossary

allopolyploidisation Polyploidisation with genetic material coming from a diverged species
9

autopolyploidisation Polyploidisation within the same species 9

orthologues Homologous genes whose divergence started at a speciation event 15

polyploidisation Mechanism leading to the acquisition of at least three versions of the same original genome in a species 9

polyspermy Fertilization of an egg by more than one sperm 9

pseudogene A gene like sequence that lost its capacity to transcribe 13

retroduplication Duplication of a gene through retro-transcription of its RNA transcript 11

segment duplication DNA sequences present in multiple locations within a genome that share high level of sequence identity 11

subfunctionalization Fate of a duplicate gene which gets a part of the original gene function, the function being shared among multiple duplicates 13

Acronyms

TAG Tandemly Arrayed Genes 11, 15, 17

WGD Whole Genome Duplication 9

1 Scientific context

It is estimated that between 46% and 65.5% of human genes could be considered as duplicate genes¹ (CORREA et al., 2021). Duplicate genes offers a pool of genetic material available for further experimentation during species evolution.

1.1 Gene duplication mechanisms

Multiple mechanisms may lead to a gene duplication. Their effect ranges from the duplication of the whole genome to the duplication of a fragment of a gene.

1.1.1 Whole genome duplication and polyploidisation

In an event of Whole Genome Duplication (WGD), the entire set of genes present on the chromosomes is duplicated (figure 1.1 (A)). WGD can occur thanks to polyspermy or in case of a non-reduced gamete. Polyploidisation is a mechanism leading to a species with at least three copies of an initial genome. A striking example is probably *Triticum aestivum* (wheat) which is hexaploid. An hexaploid cell have three pairs of homologous chromosomes due to several hybridisation events (GOLOVNINA et al., 2007). We distinguish two kinds of polyploidisations, based on the origin of the duplicate genome: (i) Allopolyploidisation occurs when the supplementary chromosomes come from a divergent species. This is the case for *Triticum aestivum* hybridisation, which consisted in the union of the chromosome set of a *Triticum* species with those of an *Aegilops* species. (ii) Autopolyploidisation consists in the hybridisation or duplication of the whole genome within the same species.

1.1.2 Unequal crossing-over

Another source of gene duplication relies on unequal crossing-over. During cell division, a crossing-over occurs when two chromatids exchange fragments of chromosome. If the cleavage of the two chromatids occurs at different positions, the shared fragments may have different lengths. Homologous recombination of such uneven crossing-over results in the incorporation

¹The estimate vary strongly depending on the criteria in use

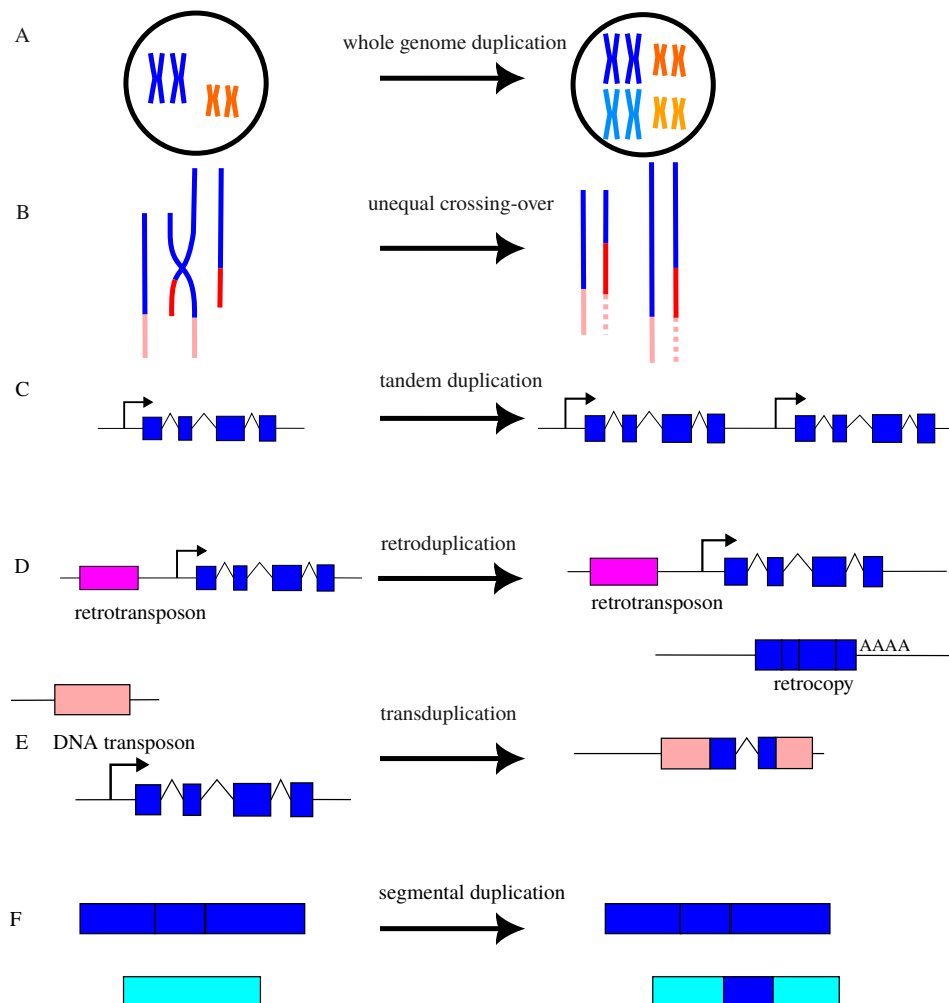


Figure 1.1: Different types of duplication. (A) Whole genome duplication. (B) An unequal crossing-over leads to a duplication of a fragment of a chromosome. (C) In tandem duplication, two (set of) genes are duplicated one after the other. (D) Retrotransposon enables retroduplication: a RNA transcript is reverse transcribed and inserted back without introns and with a polyA tail in the genome. (E) A DNA transposon can acquire a fragment of a gene. (F) Segmental duplication corresponds to long stretches of duplicated sequences with high identity. Adapted from (LALLEMAND et al., 2020) (fig. 1).

of a duplicate region, as depicted in figure 1.1 (B, C). This mechanism leads to the duplication of the whole set of genes present in the fragment. These duplicate genes locate one set after the other: we call them Tandemly Arrayed Genes (TAG). TAG are the kind of gene duplication we will be particularly interested in during this internship.

1.1.3 Retroduplication

Transposable elements play a major role in genome plasticity, and enable gene duplication too. Retrotransposons, or RNA transposons are one type of transposable elements. They share similar structure and replication mechanisms with retroviruses. Retrotransposons replicate in the genome through a mechanism known as “copy-and-paste”. These transposons typically contain a reverse transcriptase gene. This enzyme proceeds in the reverse transcription of an mRNA transcript into its reverse, complementary DNA sequence which can then insert elsewhere in the genome. More generally, retroduplication refers to the duplication of a sequence through reverse transcription of a RNA transcript. Genes duplicated through retroduplication lose their intronic sequences and bring a polyA tail with them in their new locus (figure 1.1 (D)).

1.1.4 Transduplication

DNA transposons are another kind of transposable elements whose transposition mechanism can also lead to gene duplication. This type of transposable element moves in the genome through a mechanism known as “cut-and-paste”. A typical DNA transposon contains a transposase gene. This enzyme recognizes two sites surrounding the donor transposon sequence in the chromosome resulting in a DNA cleavage and an excision of the transposon. The transposase can then insert the transposon at a new genome locus. A transposon may bring a fragment of a gene during its transposition in the new locus (figure 1.1 (E)), leading to the duplication of this fragment.

1.1.5 Segment duplication

Finally, segment duplications, also called *low copy repeats* are long stretches of DNA with high identity score (figure 1.1 (F)). Their exact duplication mechanism remains unclear (LALLEMAND et al., 2020). They may come from an accidental replication, distinct from an uneven cross-over or a double stranded breakage. Transposable elements may well be involved in the mechanism, as a high enrichment of transposable elements has been found at duplicate segments extremities, in *Drosophila* (LALLEMAND et al., 2020).

1.2 Fate of duplicate genes in genome evolution

In his book *Evolution by Gene Duplication*, Susumu OHNO proposed that gene duplication plays a major role in species evolution (OHNO, 1970), because it provides new genetic materials to build on new phenotypes while keeping a backup gene for the previous function. Indeed, duplicate genes may evolve after duplication: they may be inactivated, becoming pseudogenes; they may be deleted or conserved and so, they may acquire new functions.

1.2.1 Pseudogenization

Duplicate genes may be inactivated and become pseudogenes. These pseudogenes keep a gene-like structure, which degrades as and when further genome modifications occur. However, they are no longer expressed.

1.2.2 Neofunctionalization

Duplicate genes may be conserved and gain a new function. For instance, in the set of olfactory receptor genes result from several duplication and deletion events (in *Drosophila*: (NOZAWA and NEI, 2007)), after which the duplicate may specialize in the detection of a particular chemical compound.

1.2.3 Subfunctionalization

Two duplicate genes with the same original function may encounter a subfunctionalization by which each gene conserves only one part of the function.

1.2.4 Functional redundancy

The two gene copies may keep the ancestral function: in this case the organism may increase the quantity of gene product.

1.3 Methods to identify duplicate genes

(LALLEMAND et al., 2020) review the different methods used to detect duplicate genes. These methods depend on the type of duplicate genes they target and vary on computation burden as well as ease of use (LALLEMAND et al., 2020).

1.3.1 Paralog detection

Paralogs are homologous genes derived from a duplication event. We can identify them as homologous genes coming from the same genome, or as homologous genes between different species once we filtered out orthologues (homologous genes derived from a speciation event).

We can use two gene characteristics to assess the homology between two genes: gene structure or sequence similarity. The sequence similarity can be tested with a sequence alignment tool, such as BLAST (ALTSCHUL et al., 1990), Psi-BLAST, and HMMER3 (JOHNSON et al., 2010), or diamond (BUCHFINK et al., 2021), which are heuristic algorithms, which means they may not provide the best results, but do so way faster than exact algorithms, such as the classical Smith and Waterman algorithm (SMITH and WATERMAN, 1981) or its optimized versions PARALIGN (ROGNES, 2001) or SWIMM.

1.3.2 FTAG Finder

Developed in the LaMME laboratory, the FTAG Finder (Families and Tandemly Arrayed Genes Finder) pipeline is a simple pipeline targeting the detection of TAG from the proteome of single species (BOUILLON et al., 2016).

The pipeline proceeds in three steps. First, it estimates the homology links between each pair of genes. Then, it deduces the gene families. Finally, it searches for TAG.

Estimation of homology links between genes

This step consists in establishing a homology relationship between each genes in the proteome. In this step, the typical tool involved is BLAST (Basic Local Alignment Search Tool) (ALTSCHUL et al., 1990) run “all against all” on the proteome.

Several BLAST metrics can be used as an homology measure, such as bitscore, identity percentage, E-value or variations of these. The choice of metrics can affect the results of graph clustering in the following step, and we should therefore chose them carefully (GIBBONS et al., 2015).

Identification of gene families

Based on the homology links between each pair of genes, we construct a undirected weighted graph whose vertices correspond to genes and edges to homology links between them. We apply a graph clustering algorithm on the graph in order to infer the gene families corresponding to densely connected communities of vertices.

FTAG Finder proposes three clustering algorithm alternatives: single linkage, Markov Clustering (VAN DONGEN, 1998) or Walktrap (PONS and LATAPY, 2005).

Detection of TAGs

The final step of FTAG Finder consists in the identification of TAG from the gene families and the positions of genes. For a given chromosome, the tool seeks genes belonging to the same family and located close to each other. The tool allows a maximal number of genes between the homologous genes, with a parameter set by the user. `fig:tag-definitions` is a schematic representation of some possible TAG positioning on a genome associated with their definition in FTAG Finder *Find Tags* step.

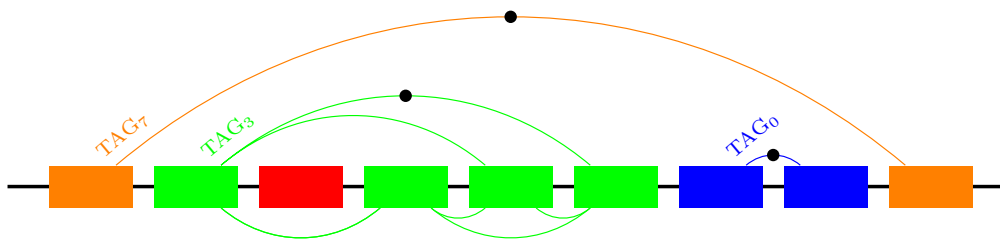


Figure 1.2: Schematic representation of TAG definitions. Several genes are represented on a linear chromosome. The red box represent a singleton gene. Orange boxes represent a TAG with two duplicate genes separated by 7 other genes (TAG₇). Four green boxes constitute a TAG, the gene at the extremities are separated by three genes (TAG₃). The two blue boxes represents a TAG with two genes next to each other TAG₀. The bended edges represents the homology links between each pair of genes of a TAG.

2 Objectives for the internship

2.1 Scientific questions

The underlying question of FTAG Finder is the study of the evolutionary fate of duplicate genes in Eukaryotes.

2.2 Extend the existing FTAG Finder Galaxy pipeline

Galaxy is a web-based platform for running accessible data analysis pipelines, first designed for use in genomics data analysis (GOECKS et al., 2010). Last year, Séanna CHARLES worked on the Galaxy version of the FTAG Finder pipeline during her M1 internship (CHARLES, 2023). I will continue this work.

2.3 Port FTAG Finder pipeline on a workflow manager

Another objective of my internship will be to port FTAG Finder on a workflow manager better suited to larger and more reproducible analysis.

We will have to make a choice for the tool we will use. The two main options being Snakemake and Nextflow. Snakemake is a python powered workflow manager based on rules *à la* GNU Make (KÖSTER and RAHMANN, 2012). Nextflow is a groovy powered workflow manager, which rely on the data flows paradigm (DI TOMMASO et al., 2017). Both are widely used in the bioinformatics community, and their use have been on the rise since they came out in 2012 and 2013 respectively (DJAFFARDJY et al., 2023).

Bibliography

- ALTSCHUL, Stephen F., Warren GISH, Webb MILLER, Eugene W. MYERS, and David J. LIPMAN (Oct. 5, 1990). “Basic Local Alignment Search Tool”. In: *Journal of Molecular Biology* 215.3, pp. 403–410. ISSN: 0022-2836. DOI: 10.1016/S0022-2836(05)80360-2. URL: <https://www.sciencedirect.com/science/article/pii/S0022283605803602> (visited on 04/30/2023).
- BOUILLON, Bérengère, Franck SAMSON, Etienne BIRMELÉ, Loïc PONGER, and Carène RIZZON (2016). *FTAG Finder: Un Outil Simple Pour Déterminer Les Familles de Gènes et Les Gènes Dupliqués En Tandem Sous Galaxy*.
- BUCHFINK, Benjamin, Klaus REUTER, and Hajk-Georg DROST (Apr. 2021). “Sensitive Protein Alignments at Tree-of-Life Scale Using DIAMOND”. In: *Nature Methods* 18.4, pp. 366–368. ISSN: 1548-7105. DOI: 10.1038/s41592-021-01101-x. URL: <https://www.nature.com/articles/s41592-021-01101-x> (visited on 03/28/2024).
- CHARLES, Séanna (2023). *Finalisation du pipeline FTAG (Families and TAG) Finder, un outil de détection des gènes dupliqués sous Galaxy*. Internship Report. Laboratoire de Mathématiques et Modélisation d’Évry.
- CORREA, Margot, Emmanuelle LERAT, Etienne BIRMELÉ, Franck SAMSON, Bérengère BOUILLON, Kévin NORMAND, and Carène RIZZON (May 1, 2021). “The Transposable Element Environment of Human Genes Differs According to Their Duplication Status and Essentiality”. In: *Genome Biology and Evolution* 13.5, evab062. ISSN: 1759-6653. DOI: 10.1093/gbe/evab062. URL: <https://doi.org/10.1093/gbe/evab062> (visited on 09/15/2023).
- DI TOMMASO, Paolo, Maria CHATZOU, Evan W FLODEN, Pablo Prieto BARJA, Emilio PALUMBO, and Cedric NOTREDAME (Apr. 2017). “Nextflow Enables Reproducible Computational Workflows”. In: *Nature Biotechnology* 35.4, pp. 316–319. ISSN: 1087-0156, 1546-1696. DOI: 10.1038/nbt.3820. URL: <https://www.nature.com/articles/nbt.3820> (visited on 03/27/2024).
- DJAFFARDJY, Marine, George MARCHMENT, Clémence SEBE, Raphael BLANCHET, Khalid BELLAJHAME, Alban GAIGNARD, Frédéric LEMOINE, and Sarah COHEN-BOULAKIA (2023). “Developing and Reusing Bioinformatics Data Analysis Pipelines Using Scientific Workflow Systems”. In: *Computational and Structural Biotechnology Journal* 21, p. 2075. DOI: 10.1016/j.csbj

- .2023.03.003. pmid: 36968012. URL: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC10030817/> (visited on 03/26/2024).
- GIBBONS, Theodore R., Stephen M. MOUNT, Endymion D. COOPER, and Charles F. DELWICHE (Dec. 2015). “Evaluation of BLAST-based Edge-Weighting Metrics Used for Homology Inference with the Markov Clustering Algorithm”. In: *BMC Bioinformatics* 16.1, p. 218. ISSN: 1471-2105. DOI: 10.1186/s12859-015-0625-x. URL: <https://bmcbioinformatics.biomedcentral.com/articles/10.1186/s12859-015-0625-x> (visited on 03/19/2024).
- GOECKS, Jeremy, Anton NEKRUTENKO, James TAYLOR, and GALAXY TEAM (2010). “Galaxy: A Comprehensive Approach for Supporting Accessible, Reproducible, and Transparent Computational Research in the Life Sciences”. In: *Genome Biology* 11.8, R86. ISSN: 1474-760X. DOI: 10.1186/gb-2010-11-8-r86. pmid: 20738864.
- GOLOVNINA, K. A., S. A. GLUSHKOV, A. G. BLINOV, V. I. MAYOROV, L. R. ADKISON, and N. P. GONCHAROV (Apr. 1, 2007). “Molecular Phylogeny of the Genus *Triticum* L”. In: *Plant Systematics and Evolution* 264.3, pp. 195–216. ISSN: 1615-6110. DOI: 10.1007/s00606-006-0478-x. URL: <https://doi.org/10.1007/s00606-006-0478-x> (visited on 03/27/2024).
- JOHNSON, L. Steven, Sean R. EDDY, and Elon PORTUGALY (Aug. 18, 2010). “Hidden Markov Model Speed Heuristic and Iterative HMM Search Procedure”. In: *BMC Bioinformatics* 11.1, p. 431. ISSN: 1471-2105. DOI: 10.1186/1471-2105-11-431. URL: <https://doi.org/10.1186/1471-2105-11-431> (visited on 04/09/2024).
- KÖSTER, Johannes and Sven RAHMANN (Oct. 1, 2012). “Snakemake—a Scalable Bioinformatics Workflow Engine”. In: *Bioinformatics (Oxford, England)* 28.19, pp. 2520–2522. ISSN: 1367-4811. DOI: 10.1093/bioinformatics/bts480. pmid: 22908215.
- LALLEMAND, Tanguy, Martin LEDUC, Claudine LANDÈS, Carène RIZZON, and Emmanuelle LERAT (Sept. 4, 2020). “An Overview of Duplicated Gene Detection Methods: Why the Duplication Mechanism Has to Be Accounted for in Their Choice”. In: *Genes* 11.9, p. 1046. ISSN: 2073-4425. DOI: 10.3390/genes11091046. URL: <https://www.mdpi.com/2073-4425/11/9/1046> (visited on 03/19/2024).
- NOZAWA, Masafumi and Masatoshi NEI (Apr. 24, 2007). “Evolutionary Dynamics of Olfactory Receptor Genes in *Drosophila* Species”. In: *Proceedings of the National Academy of Sciences* 104.17, pp. 7122–7127. DOI: 10.1073/pnas.0702133104. URL: <https://www.pnas.org/doi/full/10.1073/pnas.0702133104> (visited on 04/02/2024).
- OHNO, Susumu (1970). *Evolution by Gene Duplication*. Berlin, Heidelberg: Springer Berlin Heidelberg. ISBN: 978-3-642-86661-6. DOI: 10.1007/978-3-642-86659-3. URL: <http://link.springer.com/10.1007/978-3-642-86659-3> (visited on 03/21/2024).
- PONS, Pascal and Matthieu LATAPY (Dec. 12, 2005). *Computing Communities in Large Networks Using Random Walks (Long Version)*. DOI: 10.48550/arXiv.physics/0512106.

arXiv: physics/0512106. URL: <http://arxiv.org/abs/physics/0512106> (visited on 03/30/2024). preprint.

ROGNES, Torbjørn (Apr. 1, 2001). “ParAlign: A Parallel Sequence Alignment Algorithm for Rapid and Sensitive Database Searches”. In: *Nucleic Acids Research* 29.7, pp. 1647–1652. ISSN: 0305-1048. DOI: 10.1093/nar/29.7.1647. URL: <https://doi.org/10.1093/nar/29.7.1647> (visited on 04/09/2024).

SMITH, T. F. and M. S. WATERMAN (Mar. 25, 1981). “Identification of Common Molecular Subsequences”. In: *Journal of Molecular Biology* 147.1, pp. 195–197. ISSN: 0022-2836. DOI: 10.1016/0022-2836(81)90087-5. URL: <https://www.sciencedirect.com/science/article/pii/0022283681900875> (visited on 04/29/2023).

VAN DONGEN, S. (Jan. 1, 1998). “A New Cluster Algorithm for Graphs”. In: R 9814. URL: <https://ir.cwi.nl/pub/4604> (visited on 03/22/2024).

Summary