

Title

Matthias Grenié¹, Jean-François Gôût², Michael Lynch²

1 Département de Biologie, École Normale Supérieure de Lyon, Lyon, France

2 Biology Department, Indiana University, IN, United States of America

Abstract

Author Summary

Introduction

Since Ohno first hypothesized the influence of Whole-Genome Duplications (WGD), scientists kept showing that a broad number of organisms experienced several rounds of WGD : yeast, insects, Angiosperma, Vertebrates, Salmonids, and many others. WGD are evolutionary event when the genome of a given individual is duplicated, meaning that the whole genome is in two copies, duplicated pairs of genes are called *paralogs*. WGD can also occur when two closely related species hybridize and form a fertile descendant.

WGDs may be involved in many evolutive radiations as they provide the raw material to explore new evolutionary landscapes. A recent study on the horseshoe crab genome underlined that WGDs may be a more common phenomenon than stated. They showed that at least a WGD occurred in this conserved lineage, thus WGDs may not be evolutionary drivers.

Still, it seems that WGDs are widespread along the tree of life and since Ohno numerous models have been developed to explain the retention rate we observe between duplicate genes (reviewed in). Gene Balance Hypothesis is for example a well explored hypothesis in the litterature, according to this hypothesis dosage-sensitive gene are more retained because of stoichiometry issues, it would thus explain the over retention of transcription factors and multicomplexes proteins.

To understand the consequences of WGDs we have been studying various *Paramecium* species. *Paramecium* are a Ciliates group (See position on phylogenetic tree). Aury *et al.* showed that at least three round of WGDs occured in the *Paramecium* genus, two of which occured in the *aurelia* complex. It is a cryptic species complex of 13 species who are reproductively isolated. Studying the evolutionary fate of duplicate genes amond this complex is of great interest.

Recently, several studies unraveled the link between gene expression in gene retention after WGD in *Paramecium*. Genes that are highly expressed are more likely to be retained than genes with a low expression. The COSTEX model proposed by the authors states that gene evolution is more constrained as gene expression is high. It explains well why highly expressed genes are more retained and why they evolve more slowly than other genes.

Having established this link between gene retention and expression level of the gene, it is normal to try to get better understanding of *Paramecium* regulatory sequences. The compactness of *Paramecium* genomes makes the study of regulatory elements easier than in other eukaryotic species. Indeed, with an average of 300nt long intergenic regions in *P. tetraurelia* the genome is similar to the one of yeast.

Various strategies to study regulatory elements evolution have been developed ; the idea being that regulatory elements can be conserved along evolution, some tools can detect these more conserved regions and thus output putative motifs. Phylogenetic footprinting is one of the approaches developed when trying to detect motifs among several species, the motif finding process is weighted by the phylogenetical relationships of the given sequences, *id est*, if a motif is found in two closely related species it will have less weight than if it is detected between two distant species.

In this paper we develop a phylogenetic footprinting workflow to study the *cis*-regulatory elements among three *aurelia* species : *P. biaurelia*, *P. sexaurelia*, *P. tetraurelia* and a more basal species *P.*

caudatum. The three *aurelia* experienced two specific rounds of WGD compared to *P. caudatum*, thus when studying orthology and paralogy groups between genes, they can contain up to 13 different genes : 1 from *P. caudatum* and 4 from each *aurelia* species. Using those groups we use our workflow to identify conserved motifs in the upstream regions in each group.

Materials and Methods

developed a whole pipeline (show simple pipeline graph)

- Families upstream sequences extraction
- CDSs extraction and alignment
- CDSs phylogenetic tree
- BigFoot identification (explanation of phylogenetic score and alignment score)
- MEME research
- Comparison MEME and BigFoot
- Identification of given motif in species genome
- Correlation between motifs and expression levels

Used TranslatorX, PhyML, BigFoot, MEME.

We set up a pipeline to make our analyse (Fig.), the code is available at.

Genomes and Annotation

We used annotation and sequence from our previous analysis for *P.* and *P.* etc. (see .) and the additional sequence of.

Gene families

Looking at the phylogenetic tree of the *aurelia* species, two WGDs occurred at the root and affected three of our species. We have established some gene families from WGD2 using comparison, each family contains a set of orthologous genes between the four *aurelia* species studied, and eventually, the paralogous gene found in each species ; at maximum the families contain 13 different genes. We identified 5781 families.

Upstream sequences extraction

We considered only families with at least 4 genes. Then, using our assembly and annotation of each species genome, we extracted upstream regions from 15nt with a cut-off at 250nt of all genes of the family. If the upstream region of a gene overlapped with another gene we discarded the gene, if the upstream region was less than 15nt long we also discarded the gene. Considering discarded genes, we kept only families with at least 4 members in our datasets. Under these conditions, we extracted genes from 5008 families.

Coding Sequences extraction and alignment

Phylogenetic footprinting requires a phylogenetic tree when detecting motifs, to weigh the phylogenetic signal of given motifs. A motif conserved between two close species will have less importance than a motif conserved in two distant species. Because we are focusing on the conservation of upstream sequences we chose not to use them to avoid circularity. Instead, corresponding coding sequences (CDS) were extracted and used to model phylogenetic trees for each family. We preferred to have a gene tree over species tree, to avoid eventual inconsistencies because of gene conversion (ref. needed).

CDSs in each family were aligned using TranslatorX (ref. needed) a protein-guided alignment software. The Maximum Likelihood (ML) tree was then computed using PhyML (ref. needed) default parameters.

Phylogenetic footprinting

We used a phylogenetic footprinting software BigFoot (ref. needed) to detect highly conserved motifs in upstream sequences. We used 10000 burn-in cycles and 20000 cycles with a sampling rate of 1000 for the Hidden Monte-Carlo Markov Chain (HMCMC) process. BigFoot aligns the given sequences with gaps and tries to identify conserved and non-conserved regions, the stochastic process of the HMCMC let BigFoot refines its alignment. BigFoot assigns to each nucleotide an alignment score, which represents the confidence in the alignment, and a prediction score, measuring the phylogenetic signal of the given nucleotide, *i.e.*, the more conserved the nucleotide is, according to the phylogenetic tree, the better the score.

Using these scores we detected motifs of at least 6 nucleotides long, alignment score over 0.8 and phylogenetic score over 0.9. Because of the known biological nature of Transcription Factor Binding Sites, we allowed for a gap in motifs of 2 nucleotide, so that the scores could drop under the thresholds in those gaps.

Among the previous 5008 families, we identified 1060 motifs in 735 families.

Comparison with MEME

To confirm our phylogenetic footprinting findings, used MEME (ref. needed) to analyze the motifs in each family. MEME use a statistical process to find motifs. Usually, among a set of given

Pipeline

Description of the full pipeline

Challenges

Species tree or gene tree ?

Using a phylogenetic footprinting program means we have to use a phylogenetic tree and depending on the phylogenetic tree we are using, the evolutionary signal used in the footprinting is not identical.

The species tree (see figure.) gives us the relation between all considered, accounting for the various splits between species along with WGDs. The problem is that, not all gene families follow this tree. In particular, gene conversion, is heavily involved in shifting genes tree away from species tree : paralogs gene recombine and ...

Motif detection strategy

BigFoot does not output directly identified motifs, instead it produces two files with an alignment of the sequences and associated phylogenetic and alignment scores. The phylogenetic scores is computed, etc. The alignment scores, etc.

Transcription factor binding sites are known to be generally conserved but degenerate on certain positions. For example, ... showed that this motif was conservedNN... with two highly variable positions (denoted by N; meaning A; T; C; G using IUPAC notation). Thus, to seek biologically relevant motifs, we had to take into account that in the middle of motifs, the phylogenetic score could drop on several positions. (Show a phylogenetic score profile?)

To answer this problem we use a sliding window method of 8 nucleotides in our analysis : for each family, we looked at the scores of 8 nucleotides at the time and slide along the sequences. If the window contained at least 6 bases with scores above our thresholds, we would retain this motif. Then, from this particular region we would try to extend the sequence by adding adjacent nucleotides with good scores.

Measure motif relevance

MEME was shown to have a very high False Positive Rate of the discovery (Ref. needed). That is why many studies combine multiple motif detection tools (Ref. needed).

In our case, MEME is of particular interest as it obtains motifs using a totally different method from BigFoot.

To assess the relevance of motifs found from BigFoot's outputs, we compared them to a well-known motif finding program : MEME. For each family we identified overlapping motifs between MEME and BigFoot. We computed an overlapping index as follows : $0 \leq \frac{\text{nucleotidesincommon}}{\text{sizeofthesmallestmotif}} \leq 1$, if this index was over 0.9 we would then considered the motif as relevant.

Results

From our orthologs and paralogs gene families we had 5751 families, only 5008 families were analyzed according to our conditions. Among these 5008 families BigFoot found 811 motifs in 608 families. From these 811 motifs 117 matche

Perspectives

Conservation among species. Major results is that. Divergent resolution of WGD \rightarrow divergence in motifs ?

After identifying motifs should relate presence/absence of motifs to duplicated genes fate.

Motifs detection should take phylogeny into account for comparative analysis. Not the same value.

Need to improve pipeline for degenerate motifs, for the moment, extract only exact motif in the genome. Need to measure diversity among detected motifs -¿ clustering tools, and suppress redundancy Implement other motif finding tools to validate results.

Materials and Methods

Acknowledgments

Figure Legends

Tables