# BMC Research Notes

Short Report

# Intragenic tandem repeats in *Daphnia magna*: structure, function and distribution

Isabelle Colson, Louis Du Pasquier and Dieter Ebert*

Address: Basel University, Zoological Institute, Vesalgasse 1, CH-4051 Basel, Switzerland

Email: Isabelle Colson - i.colson@bangor.ac.uk; Louis Du Pasquier - dupasquier@dial.eunet.ch; Dieter Ebert* - dieter.ebert@unibas.ch

* Corresponding author

## Abstract

**Background:** Expressed sequence tag (EST) databases provide a valuable source of genetic data in organisms whose genome sequence information is not yet compiled. We used a published EST database for the waterflea *Daphnia magna* (Crustacea:Cladocera) to isolate variable number of tandem repeat (VNTR) markers for linkage mapping, Quantitative Trait Loci (QTL), and functional studies.

**Findings:** Seventy-four polymorphic markers were isolated and characterised. Analyses of repeat structure, putative gene function and polymorphism indicated that intragenic tandem repeats are not distributed randomly in the mRNA sequences; instead, dinucleotides are more frequent in non-coding regions, whereas trinucleotides (and longer motifs involving multiple-of-three nucleotide repeats) are preferentially situated in coding regions. We also observed differential distribution of repeat motifs across putative genetic functions. This indicates differential selective constraints and possible functional significance of VNTR polymorphism in at least some genes.

**Conclusion:** Databases of VNTR markers situated in genes whose putative function can be inferred from homology searches will be a valuable resource for the genetic study of functional variation and selection.

## Background

Waterfleas of the genus *Daphnia* (Crustacea:Cladocera) are small planktonic crustaceans found in standing freshwater bodies around the world. They have a long history as model organisms for evolutionary, ecological and ecotoxicological research. Recently, the genus has been the focus of a major sequencing effort, and the full genome sequence of *Daphnia pulex* is now available [1]. Genomic resources are steadily being developed for another species of the genus, *D. magna*. In particular, a database of around 12,000 expressed sequence tags (EST) is currently available [1,2], providing a useful resource to isolate polymor-

phic genetic markers in this species. Developing genetic markers from transcribed sequences offers specific advantages compared to traditional methods of screening enriched genomic libraries. Apart from the lower cost and higher speed of development, EST-derived genetic markers have a higher probability of being functionally significant and of being located in gene-rich regions [3-5]. This makes them highly useful markers for QTL mapping of ecologically-relevant phenotypes and for the study of selection in natural populations. Although it could be thought that functional constraints might limit polymorphism levels in genic repeated sequences, comparative

studies have reported both lower [6] and higher [4] levels of polymorphism in genic microsatellites as compared with genomic microsatellites. Polymorphism of transcribed repeated sequences can have direct phenotypic consequences both in terms of protein function [7] and in terms of regulating gene expression [8]; it has also been hypothesised to play an important role in evolvability and phenotypic adaptation [9-12]. Here, we report the development of 74 polymorphic VNTR markers from the EST database and explore their patterns of polymorphism in relation to repeat sequence structure and putative gene function.

## Methods

The Tandem Repeat Finder (TRF) software [13] was used to recover tandemly repeated sequences from the EST database [1]. Sequences containing only mononucleotides repeats were discarded, and redundant sequences were merged using the CAP assembly software [14]. The 346 single sequences obtained from 531 ESTs were translated using the expasy "Translate" software [15]. The amino acid sequence of the longest open reading frame (ORF) was then blasted against protein databases [16] and used in InterProScan searches [17] in order to identify functional domains and transmembrane regions. E-values < 0.0001 were accepted as significant homology in the blast searches. When translation did not produce an obvious candidate ORF, blastX searches were carried out from the nucleotide sequence. Putative function was inferred from the identity of homologous sequences and from the presence of functional domains. Six broad functional categories were defined: 1. Proteins involved in metabolism, including energy metabolism and protein synthesis (MET); 2. Proteins involved in signalling pathways and regulation of gene expression (SIG); 3. Surface or integumental proteins (SUR); 4. Proteins involved in defense (pathogens and stress) (DEF); 5. Other proteins with known function (OTH), regrouping proteins involved in development, transport and cell structure, functions that were represented by only a few loci; 6. Proteins of unknown function (UNK), regrouping loci with non-annotated homologous sequences and loci with no significant homologous sequence in Genbank. The position of the tandem repeat in the mRNA sequence (ORF, 5'UTR or 3'UTR) was determined using the gene prediction software FGenesh [18,19]. Primers were designed for 218 loci, using the "Primer 3" software [20]. DNA from 18 individuals representing six populations from Europe and North-America (UK, Germany, Belgium, Finland, Hungary and Canada) was extracted with E.Z.N.A tissue DNA mini kit (Peqlab, Germany) and used in PCR reactions. Depending on the locus, we performed either standard or hot start PCR. Standard PCR reactions were carried out in 12.5 µl reactions containing 1× PCR reaction buffer (Sigma Aldrich), 1.5 or 3.5 mM $MgCl_2$ depending on the locus, 200 µM of each dNTP, 0.2 µM of each primer (with the forward primer fluorescently labelled) and 0.5 unit Taq polymerase (Sigma Aldrich). An initial denaturation step of 4 minutes at 94°C was followed by 35 cycles of 94°C for 30 seconds, 53°C for 30 seconds, and 72°C for 30 seconds, followed by a final extension step of 72°C for 4 minutes. Hotstart PCR was performed with thermo-start PCR master mix (ABGene, Epsom, UK) with 1.5 mM or 3.5 mM of $MgCl_2$ depending on the locus, and 0.2 µM of each primer (with the forward primer fluorescently labelled). PCR conditions were as described above, except for an initial incubation at 94°C for 15 minutes. Primer sequences and PCR conditions for polymorphic VNTR loci are described in Additional File 1. PCR products were run on an ABI 310 automated sequencer (Applied Biosystems, Foster City, USA) and analysed with the Genemapper software (Applied Biosystems, Foster City, USA). Polymorphism (number of alleles) was assessed at 106 loci, which consistently amplified DNA from all individuals and for which no more than 2 alleles per individual were present. Furthermore, the 106 loci were blasted against the genome to check for any potential gene duplication.

We analysed contingency tables using the $\chi^2$ test when sample sizes were large enough (less than 20% of cells containing less than 5 cases). Otherwise, Yates' correction was employed [21]. We conducted nonparametric correlations using Spearman rank correlation factor rho, performed with SPSS 15.0.
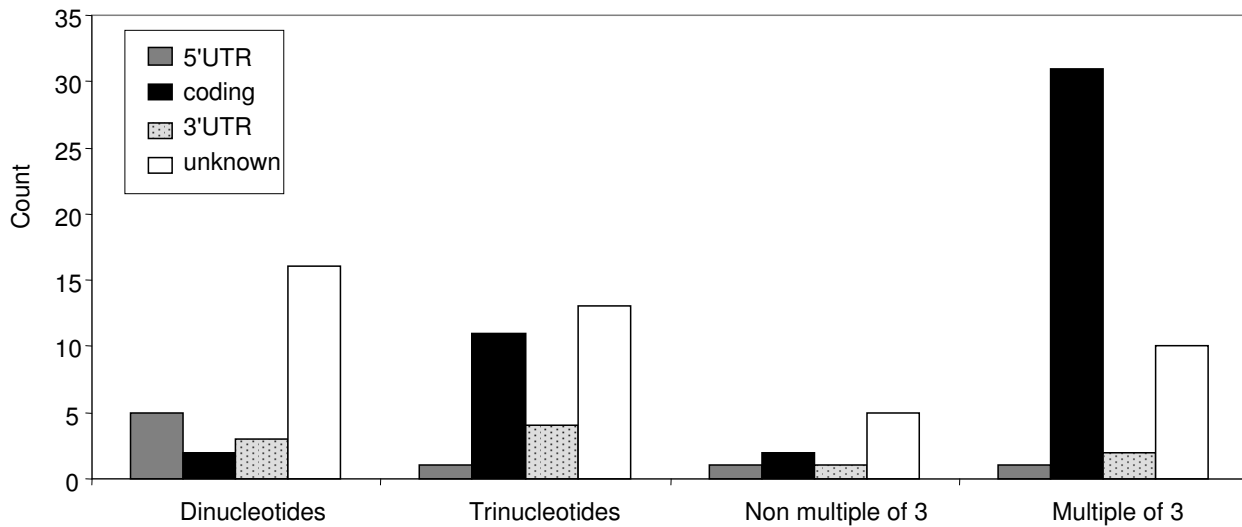
## Results and Discussion

**Figure 1**
**Distribution of repeat structures and localisation in the mRNA**
FGenesh software. "5'UTR": 5' untranslated region; "coding": protein coding region; 3'UTR: 3' untranslated region; "unknown": loci for which the location of the repeated motif could not be ascertained. Total number of loci = 106.
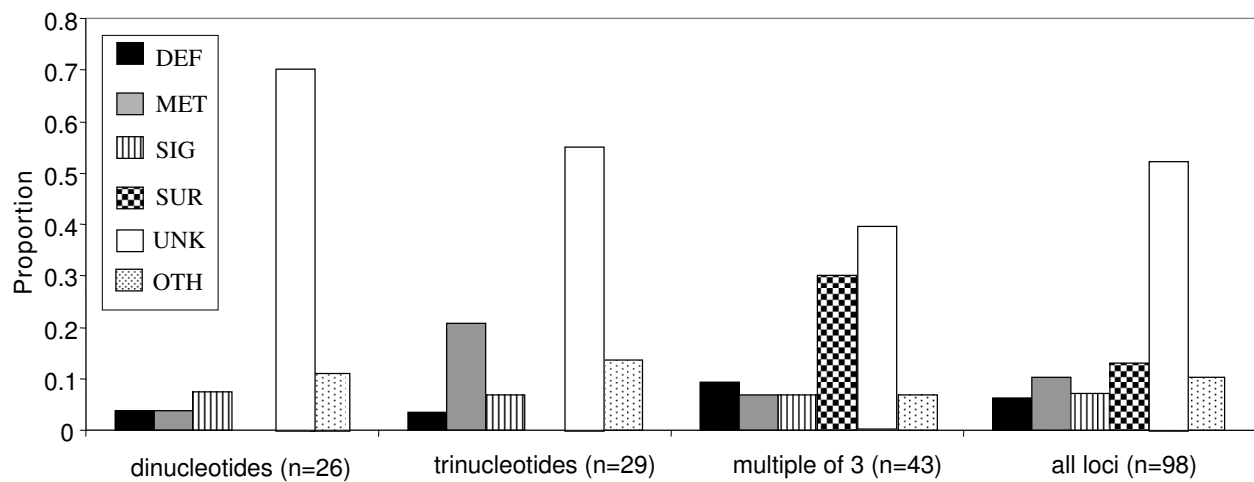


**Figure 2**
**Functional distribution of EST loci containing tandem repeats**
analysis due to their low number. Total number of loci = 98.

*Saccharomyces cerevisiae*

*D. magna*

$\chi$ = 3.048, df = 3, p > 0.05), repeat localisation in mRNA ($\chi^2$ = 0.87, df = 2, p > 0.05, with Yates correction), and putative protein function ($\chi^2$ = 1.88, df = 4, p > 0.05, with Yates correction). However, putative defense genes showed a higher proportion of polymorphic loci (5 out of 6) than other functional categories. The number of alleles was significantly positively correlated with the number of repeats (Spearman rank correlation coefficient: 0.303, p < 0.01 between total number of repeats and number of alleles; 0.218, p < 0.05 between number of perfect repeats and number of alleles).

Homology searches against the           genome identified 49 EST with partial homology of fragments longer than 100 base pairs. Eleven of these had multiple partial homologues situated in distinct genomic locations in the           genome, indicating some degree of sequence duplication and paralogy in the           genome. However, most of the homologue sequences (39) were only partial and not encompassing the whole amplicon (with either one or both primer sequences missing). We found that only ten loci had a           homolog encompassing the whole amplicon, i.e. including both primer sequences in the same scaffold, allowing for the presence of introns. In all cases, only one complete homologue was identified. From this analysis and in view of the genotyping results we are confident that, although gene duplication seems to be a common feature in the           genome (at least in           ), our genetic markers represent single loci.

The EST database allowed fast, cheap *in silico*
*Daphnia magna*

morphic loci (28/39, 72%) than loci coding for intracellular proteins (17/31, 55%). These trends can tentatively be interpreted as repeated sequences playing a role in the evolutionary dynamics of host-pathogen relationships (see [7] for a discussion of this topic in pathogens). However, more data and much more targeted analyses, which fall outside the scope of this report, are needed to further explore this possibility.

We observed a significant positive correlation between polymorphism and number of repeats amongst our loci, as previously observed in genomic microsatellites [26]. However, interruption of the length of perfect repeat array did not correlate with lower polymorphism, as is the case in genomic microsatellites [27]. This discrepancy could be explained by differences in mutational and selective constraints in intragenic and genomic microsatellites, in particular in relation to third codon position redundancy. Also, our dataset includes loci with longer repeat structures ("minisatellites") for which replication slippage might not be the primary mutational process.

EST databases are increasingly being used as a resource to develop VNTR markers, which are likely to be very informative in genome screens for functionally relevant polymorphism. The 74 VNTR markers for *D. magna* described here will be useful in producing the first genetic linkage map in the species (increasing marker density in gene rich areas), and in QTL mapping of evolutionary and ecologically relevant traits. To illustrate the potential of our VNTR markers, 34 of the 74 polymorphic markers described here were found to distinguish between two European clones used to develop recombinant lines for mapping purposes (unpublished data). The availability of markers with potentially functionally relevant polymorphism,

## Competing interests

## Authors' contributions

## References

1. **WFleabase**  [http://wfleabase.org]
2. Watanabe H, Tatarazako N, Oda S, Nishide H, Uchiyama I, Morita M, Igushi T: **Analysis of expressed sequence tags of the water flea Daphnia magna.** *Genome* **48:**
3. Vasemägi A, Nilsson J, Primmer CR: **Expressed sequence tag-linked microsatellites as a source of gene-associated polymorphisms for detecting signatures of divergent selection in Atlantic salmon (Salmo salar L.).** *Mol Biol Evol*
4. Coulibaly I, Gharbi K, Danzmann RG, Yao J, Rexroad CE III: **Characterization and comparison of microsatellites derived from repeat-enriched libraries and expressed sequence tags.** *Anim Genet* **36:**

**Development of polymorphic expressed sequence tag-derived microsatellites for the extension of the genetic linkage map of the black tiger shrimp (Penaus monodon).** **37:** **Genic microsatellite markers in plants: features and applications.** *Trends Biotechnol*

**Intragenic tandem repeats generate functional variability.** **37:** **Some microsatellites may act as novel polymorphic cis**

**341:** Caburet S, Cocquet J, Vaiman D, Veitia RA: **evolutionary "agility".** *BioEssays*

*Trends Genet* **Tandemly repeated DNA: why should anyone care?** *Mutat Res* **598:** Mularoni L, Veitia RA, Mar Albà M: **Highly constrained proteins contain an unexpectedly large number of amino acid tandem repeats.** **89:** **Tandem repeat finder: a program to analyze DNA sequences.** **27:** Huang X, Madan A: 1999, 868-877.
15. [http://www.expasy.ch/tools/dna.html]
16. Altschul SF, Warren G, Miller W, Myers EW, Lipman DJ: **Basic local alignment search tool.** *J Mol Biol*

*Genome Res* **10:** **Fgenesh** berry.phtml?topic=fgenesh&group=programs&sybgroup=gfind]
20. Rozen S, Skaletsky HJ: In *Bioinformatics Methods and Protocols: Methods in Molecular Biology*

**independence  [Computer  software].**

Andersen SO, Højrup P, Roepstorff P: *Insect Biochem Molec* 1994, 153-176.
23. Ijima M, Hashimoto T, Matsuda Y, Nagai T, Yamano Y, Ichi T, Osaki T, Kawabata S-I:

### Additional file 1
*PCR conditions for polymorphic VNTR loci. Label: F-primer fluorescent label used. [MgCl₂]: Concentration of MgCl₂ used in the PCR buffer. An asterisk following MgCl₂ concentration indicates hot start PCR.*
Click here for file
[http://www.biomedcentral.com/content/supplementary/1756-0500-2-206-S1.DOC]

### Additional file 2
*Description of the 74 polymorphic VNTR loci. EST locus: locus containing the VNTR; Size: size of the repeated motif; Sequence: repeat consensus sequence. N: number of repeats, [the number of perfect repeatsm with exact consensus sequence, is shown in brackets]; A: number of alleles. *: total number of alleles in loci with more than one repeated motif. He: Gene diversity.*

### Additional file 3
*Results of homology searches and putative functions of EST loci transmembrane regions. Function categories: DEF: defense; MET: metabolism; OTH: other; SIG: signaling and gene expression regulation; SUR: surface and integumental proteins; UNK: unknown function. No hit: no significant homolog found (cut-off E-value 0.0001).*

0500-2-206-S3.DOC]

### Additional file 4
*Distribution of allele sizes of polymorphic loci in each sampled location.*

*FEBS J*

**Ecology, Epidemiology, and Evolution of Parasitism in Daphnia.**

Qi W, Nong G, Preston JF, Ben-Ami F, Ebert D: 2009 in press.

26. Schloetterer C: 2000, 365-371.
27. Santibáñez-Koref M, Gangeswaran R, Hancock J: **rates in Mammalian genomes.** 2001, **18:**2119-2123.