

# MOLECULAR ECOLOGY RESOURCES

## Reference-free transcriptome assembly in non-model animals from next generation sequencing data

Journal:	<i>Molecular Ecology Resources</i>
Manuscript ID:	MER-12-0016.R1
Manuscript Type:	Resource Article
Date Submitted by the Author:	06-Mar-2012
Complete List of Authors:	Cahais, Vincent; Université Montpellier 2, CNRS UMR 5554 Gayral, Philippe; CNRS, Institut des Sciences de l'Évolution; Université Montpellier 2, CNRS UMR 5554 Tsagkogeorga, Georgia; Université Montpellier 2, CNRS UMR 5554 Melo-Ferreira, José; CIBIO, Centro de Investigaç�o em Biodiversidade e Recursos Genéticos, Ballenghien, Marion; Université Montpellier 2, CNRS UMR 5554 Weinert, Lucy; Université Montpellier 2, CNRS UMR 5554 Chiari, Ylenia; University of Montpellier 2, Institut des Sciences de l'Évolution; Belkhir, Khalid; Université Montpellier 2, CNRS UMR 5554 Ranwez, Vincent; Université Montpellier 2, CNRS UMR 5554 Galtier, Nicolas; Université Montpellier 2, CNRS UMR 5554
Keywords:	Bioinformatics/Phyloinformatics, Transcriptomics, Population Genetics - Empirical, Invertebrates

1  
2  
3 **Reference-free transcriptome assembly in non-model**  
4 **animals from next generation sequencing data**  
5  
6

7 Cahais V.<sup>1</sup>, Gayral P.<sup>1</sup>, Tsagkogeorga G.<sup>1,2</sup>, Melo-Ferreira J.<sup>3</sup>, Ballenghien M.<sup>1</sup>, Weinert L.<sup>1,4</sup>,  
8 Chiari Y.<sup>1,3</sup>, Belkhir K.<sup>1</sup>, Ranwez V.<sup>1</sup>, Galtier N<sup>1\*</sup>.  
9

10  
11 <sup>1</sup> Université Montpellier 2, CNRS UMR 5554, Institut des Sciences de l'Evolution de  
12 Montpellier, Place E. Bataillon, 34095 Montpellier, France;

13  
14 <sup>2</sup> School of Biological & Chemical Sciences, Queen Mary University of London, Mile End  
15 Road, London, E1 4NS

16  
17 <sup>3</sup> CIBIO, Centro de Investigação em Biodiversidade e Recursos Genéticos,  
18 Universidade do Porto, Campus Agrário de Vairão, 4485-661 Vairão, Portugal  
19

20 <sup>4</sup> Medical Research Council (MRC), Centre for Outbreak Analysis and Modelling, Imperial  
21 College Faculty of Medicine, London, UK.  
22

23 **Corresponding author:**

24 Nicolas Galtier

25 CNRS UMR5554 – Institut des Sciences de l'Evolution

26 Place E. Bataillon – CC64 – 34095 Montpellier, France

27 Phone: (+33) 467 14 48 18 Fax: (+33) 467 14 36 10

28 nicolas.galtier@univ-montp2.fr  
29

30 **Key words:** NGS, transcriptomics, paralogues, alleles

31 **Running title:** *de novo* NGS-based transcriptome assembly

32

33 **Abstract**

34

35 Next-generation sequencing technologies offer the opportunity for population genomic study  
36 of non-model organisms sampled in the wild. The transcriptome is a convenient and popular  
37 target for such purposes. However, designing genetic markers from NGS transcriptome data  
38 requires assembling gene coding sequences out of short reads. This is a complex task owing  
39 to gene duplications, genetic polymorphism, alternative splicing, and transcription noise.

40 Typical assembling programs return thousands of predicted contigs, whose connection to the  
41 species true gene content is unclear, and from which SNP definition is uneasy. Here the  
42 transcriptomes of five diverse non-model animal species (hare, turtle, ant, oyster, tunicate)  
43 were assembled from newly-generated 454 and Illumina sequence reads. In two species for  
44 which a reference genome is available, a new procedure was introduced to annotate each  
45 predicted contig as either a full length cDNA, fragment, chimera, allele, paralogue, genomic  
46 sequence or other, based on the number of, and overlap between, BLAST hits to the  
47 appropriate reference. Analyses showed that (i) the highest-quality assemblies are obtained  
48 when 454 and Illumina data are combined, (ii) typical *de novo* assemblies include a majority  
49 of irrelevant cDNA predictions, and (iii) assemblies can be appropriately cleaned by filtering  
50 contigs based on length and coverage. We conclude that robust, reference-free assembly of  
51 thousands of genes from transcriptomic next-generation sequence data is possible, opening  
52 promising perspectives for transcriptome-based population genomics in animals. A Galaxy  
53 pipeline implementing our best-performing assembling strategy is provided.

54