

Supplementary Figures and Table for “A software tool 'CroCo' detects pervasive cross-species contamination in next generation sequencing data”

Authors :

Paul Simion^{1,3}, Khalid Belkhir¹, Clémentine François¹, Julien Veyssier¹, Jochen C. Rink², Michaël Manuel³, Hervé Philippe^{4,5}, Maximilian J. Telford⁶

Affiliations :

1 Institut des Sciences de l'Evolution (ISEM), UMR 5554, CNRS, IRD, EPHE, Université de Montpellier, Montpellier, France

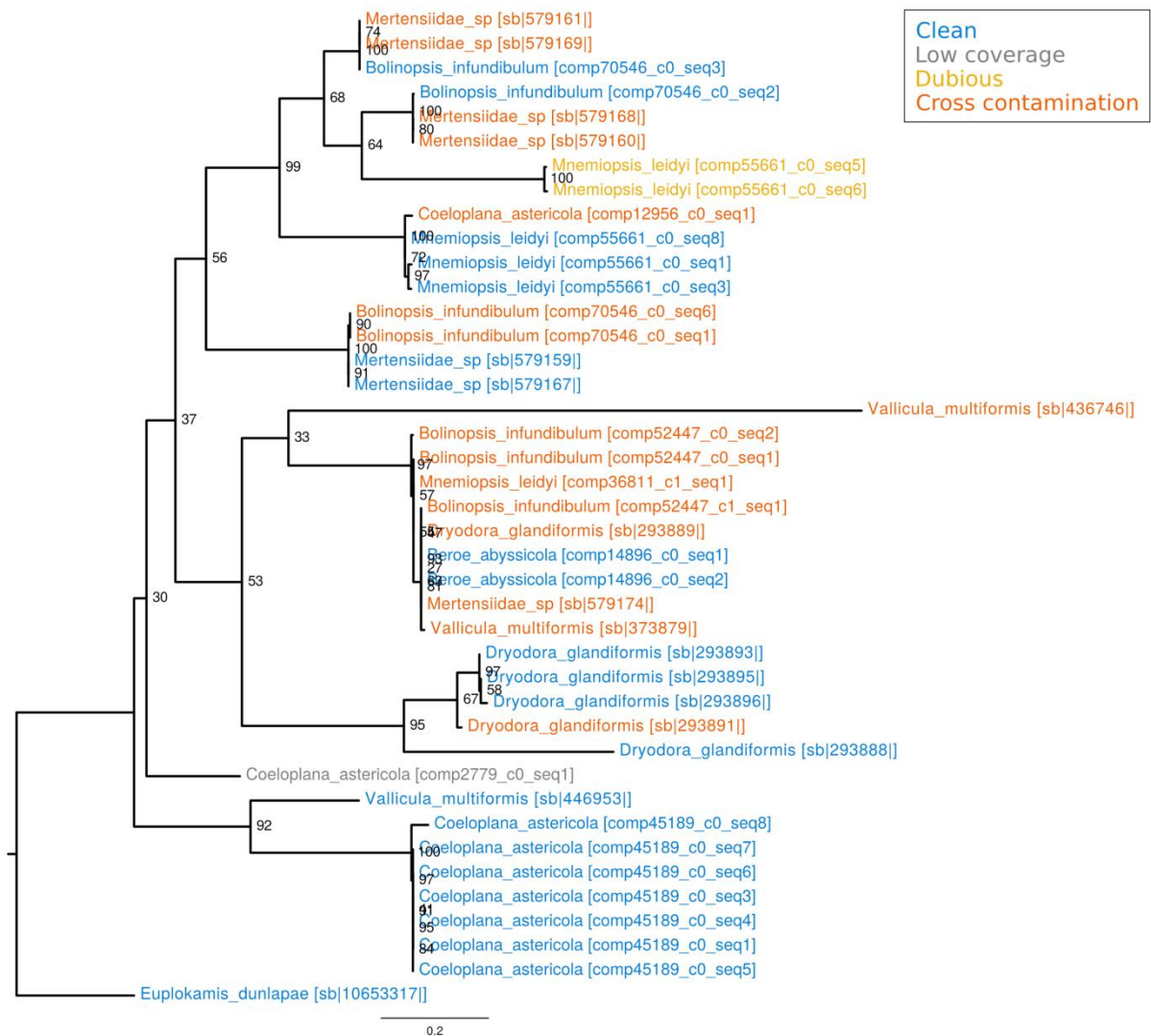
2 Max Plank Institute of Molecular Cell Biology and Genetics, Pfotenhauerstrasse 108, 01307 Dresden, Germany

3 Sorbonne Universités, UPMC Univ Paris 06, CNRS, Evolution Paris-Seine UMR7138, Institut de Biologie Paris-Seine, Case 05, 7 quai St Bernard, 75005 Paris, France

4 Centre de Théorisation et de Modélisation de la Biodiversité, Station d'Ecologie Théorique et Expérimentale, UMR CNRS 5321, Moulis, 09200, France

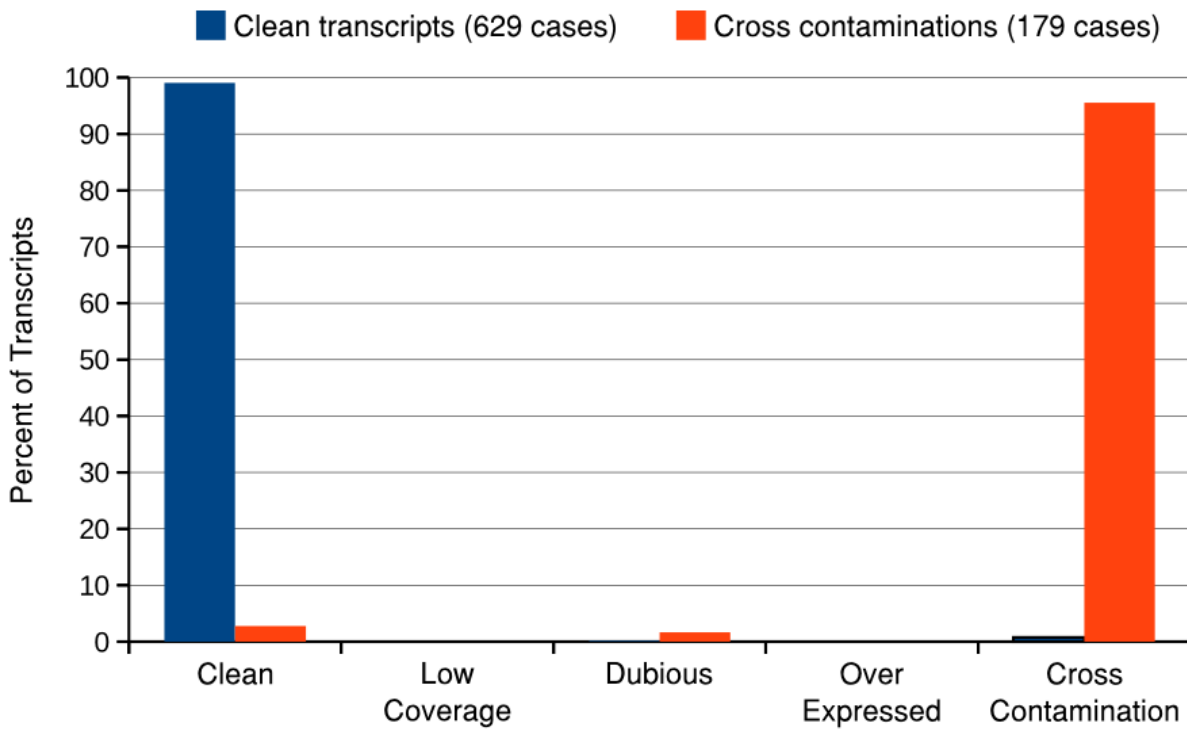
5 Département de Biochimie, Centre Robert-Cedergren, Université de Montréal, Montréal, H3C 3J7 Québec, Canada

6 University College London, Centre for Life's Origins and Evolution, Department of Genetics, Evolution and Environment, Darwin Building, Gower Street, London WC1E 6BT, UK



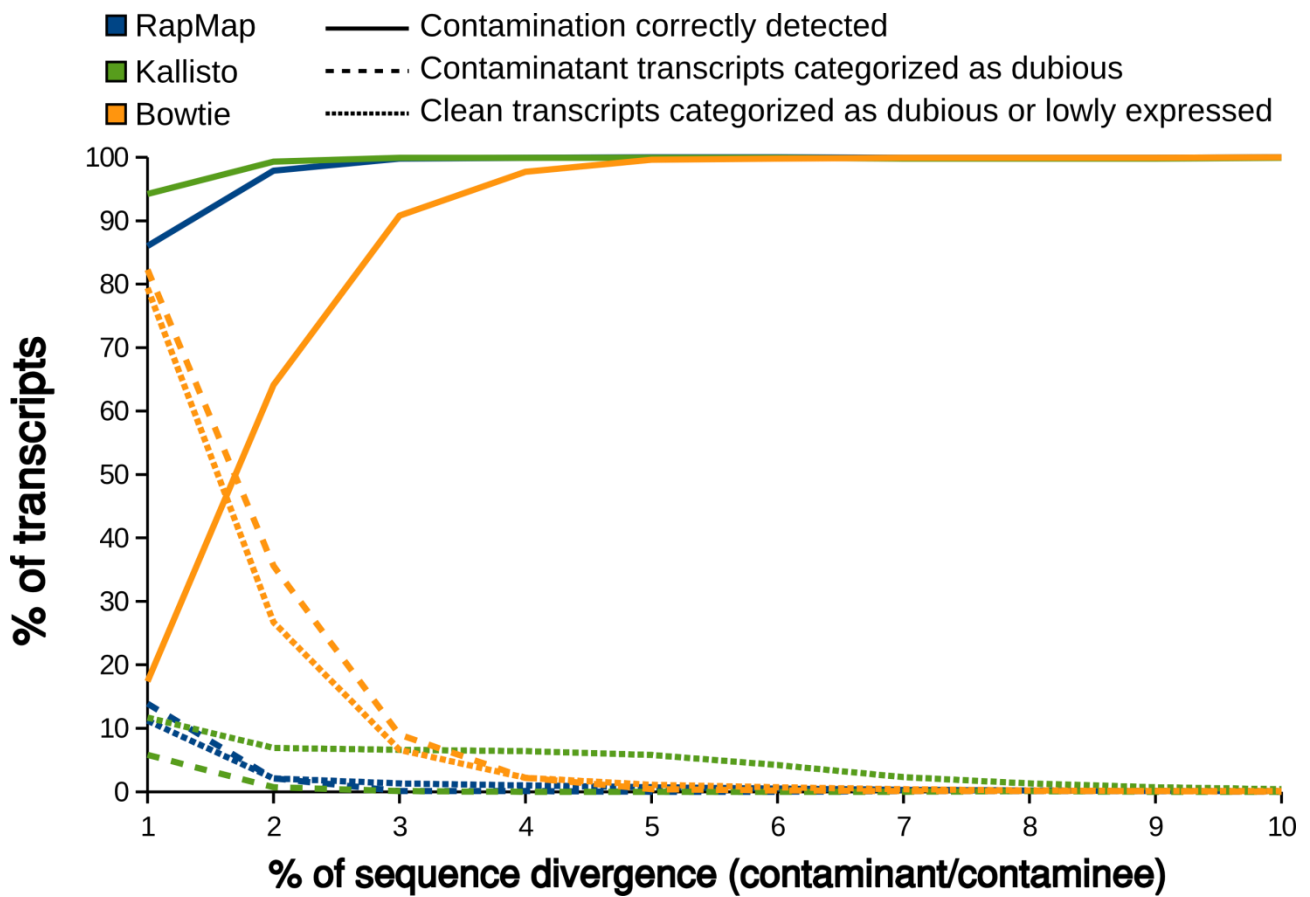
Additional file 1: Figure S1:

Single-gene phylogeny with multiple cross contaminations. Single-gene phylogeny reconstructed from a gene of dataset A belonging to the 14-3-3 gene family (see [Methods section](#) for details), showing at least 11 instances of cross contamination. We used CroCo to categorise transcripts and coloured them accordingly: blue for clean transcripts, grey for low coverage transcripts, orange for dubious transcripts and red for cross contaminations.



Additional file 1: Figure S2:

Comparison between transcript categorization by CroCo and a reference set of manually detected cross contaminations. CroCo categorization into five categories of transcripts previously classified as clean (629 cases, in blue) or as cross contaminations (179 cases, in red) using default parameters.



Additional file 1: Figure S3:

Benchmarking CroCo using simulations. Impact of genetic distance between the contaminant and the contaminee for three different mapping tools on the proportion of cross contamination correctly detected, cross contamination detected as dubious and clean transcripts categorized as anything other than clean. RapMap and Kallisto outperform Bowtie for this task.

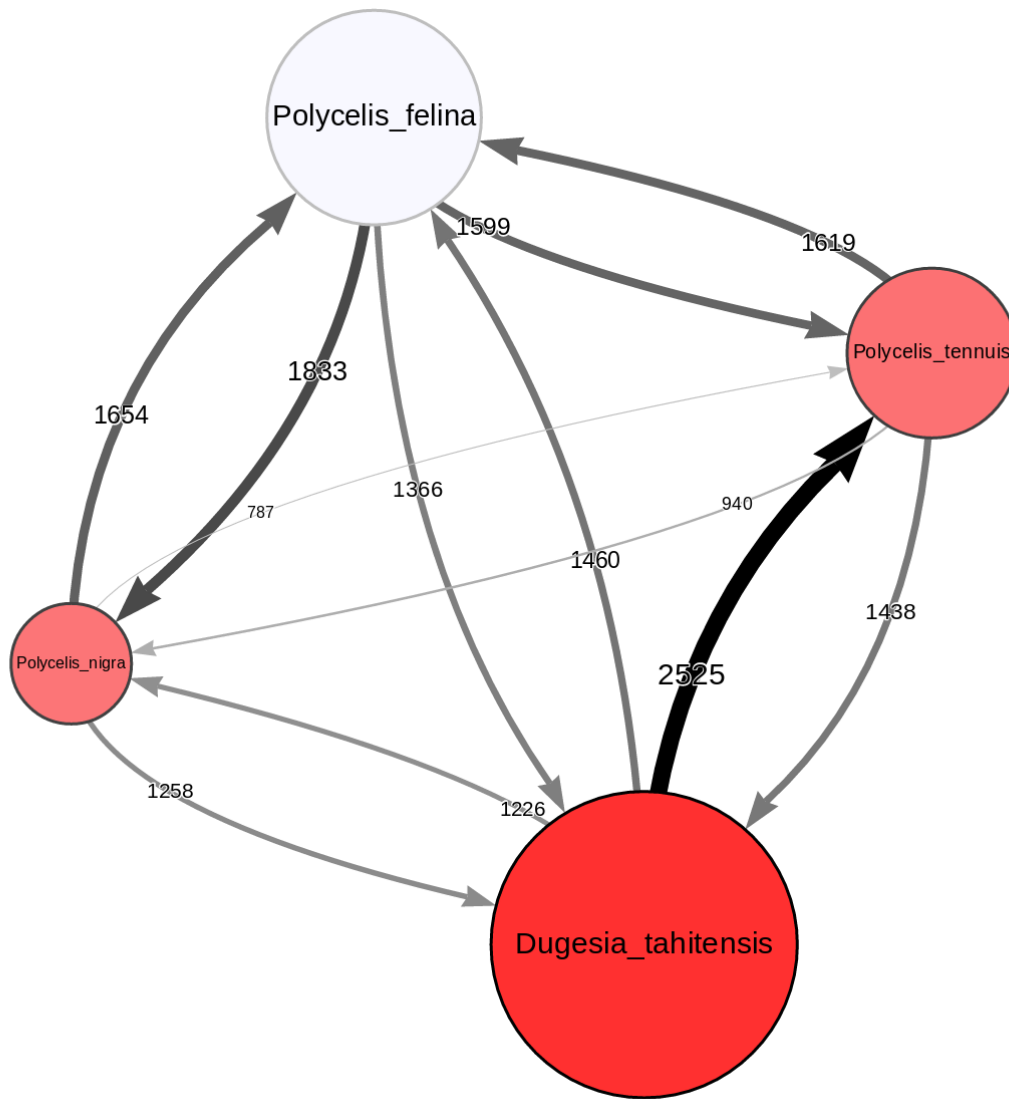
Note for network graph interpretation :

Colors, nodes diameter and arrow sizes in networks are relative to the sampling used and therefore cannot be compared across different sequencing experiments. Example : although *Dryodora glandiformis* looks cleaner in [fig. 2a](#) than *Polycelis nigra* in Additional file 1: Figure S5 based on respective colors, both species have ~6% of their transcriptome that is contaminated.



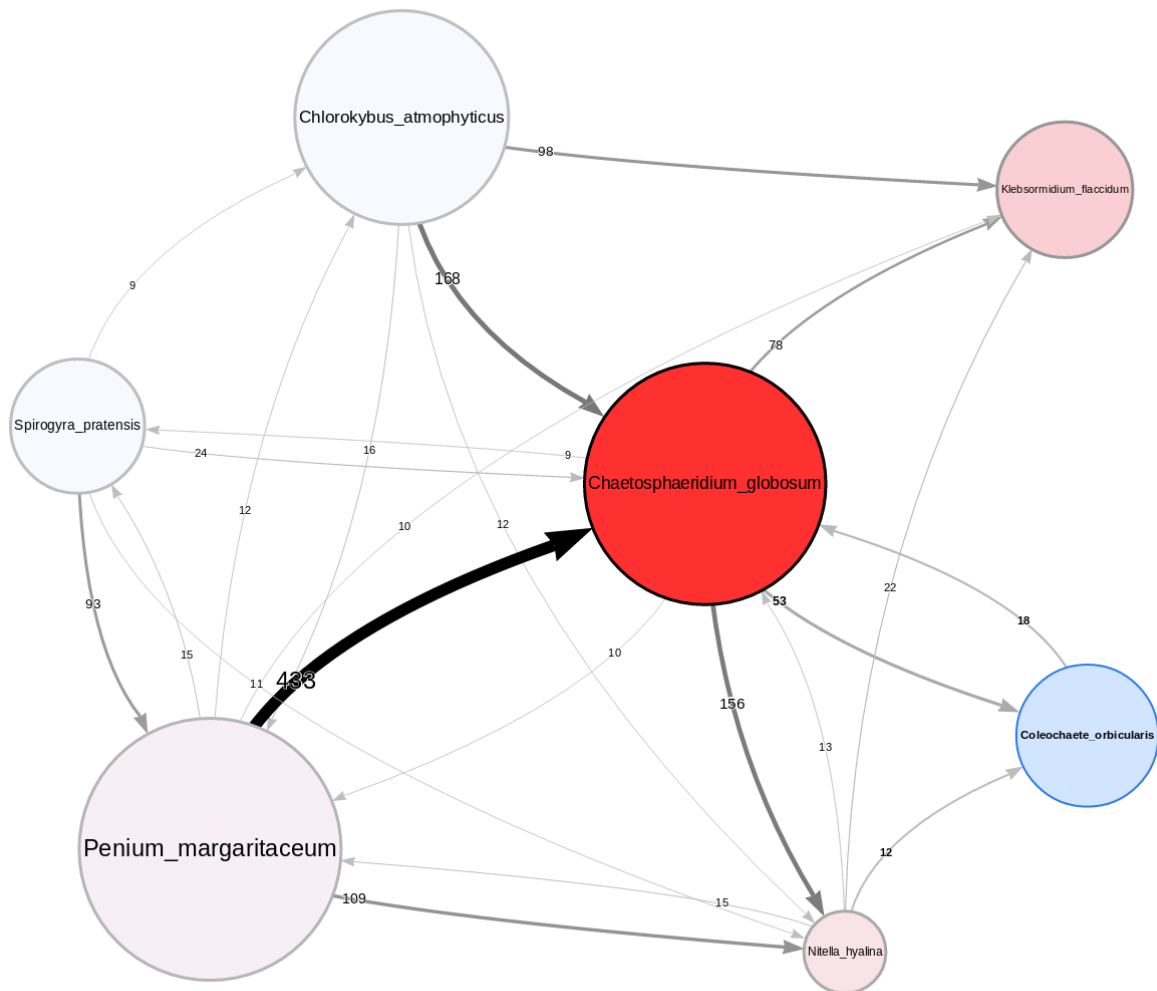
Additional file 1: Figure S4:

Network visualisation of cross contamination patterns in dataset B. Node diameter is proportional to the number of time the sample contaminates another one, node color represent the proportion of its sequence that are contaminated (from white to red), and arrow sizes represent the number of cross contaminations. For clarity, arrows representing less than 2% of the largest cross contamination link are not represented.



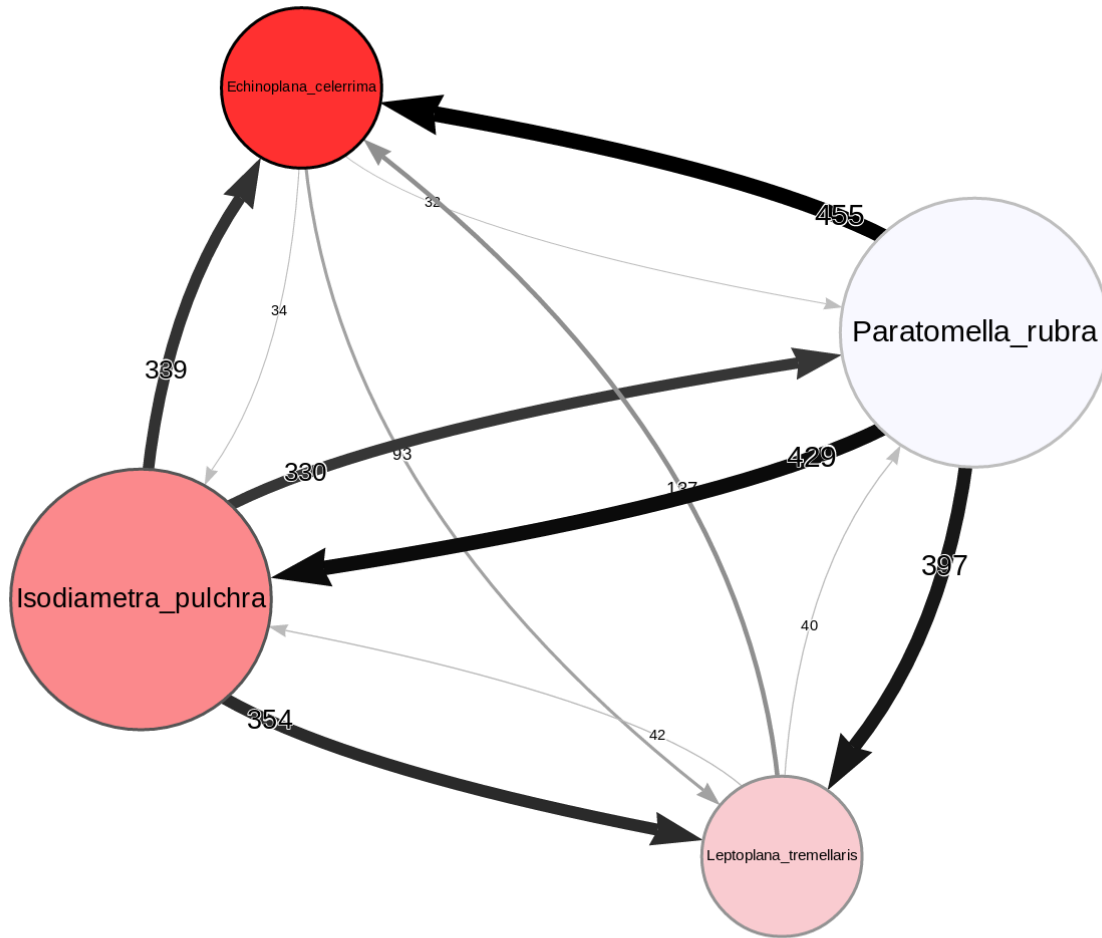
Additional file 1: Figure S5:

Network visualisation of cross contamination patterns in dataset C. Node diameter is proportional to the number of time the sample contaminate another one, node color represent the proportion of its sequence that are contaminated (from white to red), and arrow sizes represent the number of cross contaminations. For clarity, arrows representing less than 2% of the largest cross contamination link are not represented.



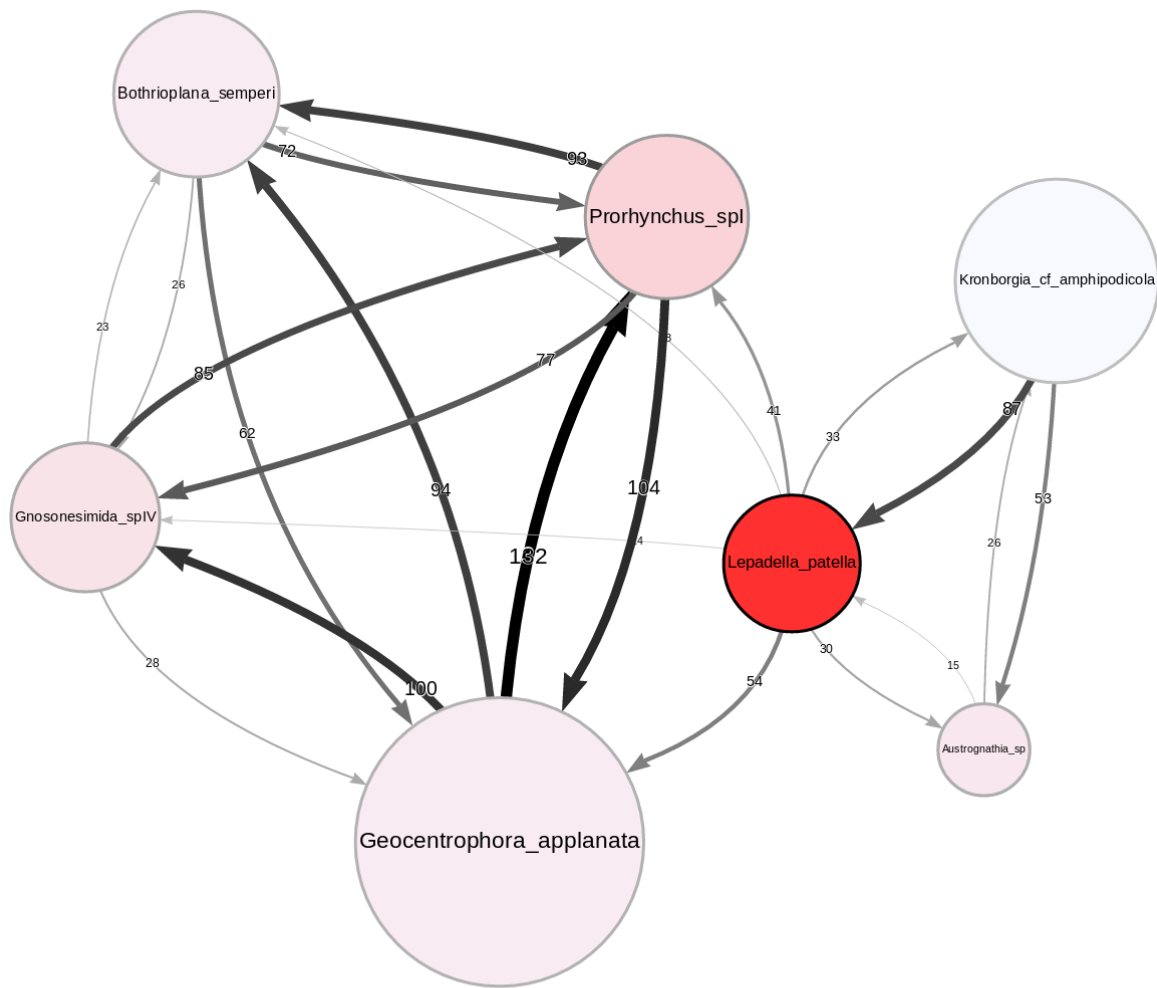
Additional file 1: Figure S6:

Network visualisation of cross contamination patterns in dataset D. Node diameter is proportional to the number of time the sample contaminate another one, node color represent the proportion of its sequence that are contaminated (from white to red), and arrow sizes represent the number of cross contaminations. For clarity, arrows representing less than 2% of the largest cross contamination link are not represented.



Additional file 1: Figure S7:

Network visualisation of cross contamination patterns in dataset E. Node diameter is proportional to the number of time the sample contaminate another one, node color represent the proportion of its sequence that are contaminated (from white to red), and arrow sizes represent the number of cross contaminations. For clarity, arrows representing less than 2% of the largest cross contamination link are not represented.



Additional file 1: Figure S8:

Network visualisation of cross contamination patterns in dataset F. Node diameter is proportional to the number of time the sample contaminate another one, node color represents the proportion of its sequence that are contaminated (from white to red), and arrow sizes represent the number of cross contaminations. For clarity, arrows representing less than 2% of the largest cross contamination link are not represented.

Additional file 1: Table S1:

Datasets from six recent sequencing projects analysed with CroCo. Datasets, species names, taxonomy and accession numbers for sequencing data.

Dataset	Species	Phylum	Accession numbers
A (Moroz <i>et al.</i> 2014)	<i>Beroe abyssicola</i>	Ctenophora	SRR777787
	<i>Bolinopsis infundibulum</i>	Ctenophora	SRR786491
	<i>Coeloplana astericola</i>	Ctenophora	SRR786490
	<i>Dryodora glandiformis</i>	Ctenophora	SRR777788
	<i>Euplokamis dunlapae</i>	Ctenophora	SRR777663
	<i>Mertensiidae</i> sp.	Ctenophora	SRR786492
	<i>Mnemiopsis leidyi</i>	Ctenophora	SRR789900
	<i>Vallicula multiformis</i>	Ctenophora	SRR786489
B (Simion <i>et al.</i> 2017)	<i>Alcyonium palmatum</i>	Cnidaria	SRR3407216
	<i>Antipathes caribbeana</i>	Cnidaria	SRR3407160
	<i>Clathrina coriacea</i>	Porifera	SRR3417192
	<i>Coeloplana meteoris</i>	Ctenophora	SRR3407215
	<i>Grantia compressa</i>	Porifera	SRR3417193
	<i>Lampea pancerina</i>	Ctenophora	SRR3407163
	<i>Leuconia nivea</i>	Porifera	SRR3417190
	<i>Liriope tetraphylla</i>	Cnidaria	SRR3407335
	<i>Lucernariopsis campanulata</i>	Cnidaria	SRR3407219
	<i>Pelagia noctiluca</i>	Cnidaria	SRR3407257
	<i>Plakina jani</i>	Porifera	SRR3417194
	<i>Pleraplysilla spinifera</i>	Porifera	SRR3417588
	<i>Plumapathes pennacea</i>	Cnidaria	SRR3407161
<i>Vallicula multiformis</i>	Ctenophora	SRR3407164	
C (This study)	<i>Dugesia tahitensis</i>	Platyhelminthes	SRR6436040
	<i>Polycelis felina</i>	Platyhelminthes	SRR6388789
	<i>Polycelis nigra</i>	Platyhelminthes	SRR6379017
	<i>Polycelis tenuis</i>	Platyhelminthes	SRR6388485
D (Finet <i>et al.</i> 2010)	<i>Chaetosphaeridium globosum</i>	Streptophyta	SRR064327
	<i>Chlorokybus atmophyticus</i>	Streptophyta	SRR064329
	<i>Coleochaete orbicularis</i>	Streptophyta	SRR036732
	<i>Klebsormidium flaccidum</i>	Streptophyta	SRR064330
	<i>Nitella hyalina</i>	Streptophyta	SRR064326
	<i>Penium margaritaceum</i>	Streptophyta	SRR064328
	<i>Spirogyra pratensis</i>	Streptophyta	SRR036731
E (This study)	<i>Echinoplana celerrima</i>	Platyhelminthes	SRR1796488
	<i>Isodiametra pulchra</i>	Xenacoelomorpha	PRJNA422347
	<i>Leptoplana tremellaris</i>	Platyhelminthes	SRR1797726
	<i>Paratomella rubra</i>	Xenacoelomorpha	PRJNA422367
F (Laumer <i>et al.</i> 2015)	<i>Austrognathia</i> sp.	Gnathostomulida	SRR1976176
	<i>Bothrioplana semperi</i>	Platyhelminthes	SRR1955240
	<i>Geocentrophora applanata</i>	Platyhelminthes	SRR1955490
	<i>Gnosonesimida</i> spIV	Platyhelminthes	SRR1976178
	<i>Kronborgia</i> cf. <i>amphipodicola</i>	Platyhelminthes	SRR1976457
	<i>Lepadella patella</i>	Rotifera	SRR1976570
	<i>Prorhynchus</i> spI	Platyhelminthes	SRR1980634

Additional file 1: Table S2:

Effect of fold difference parameter value on transcripts categorizations. Transcript categories, value of the fold difference parameter, number of transcripts in *Mnemiopsis leidyi* and *Vallicula multiformis*.

	fold difference Parameter Value	<i>Mnemiopsis leidyi</i>	<i>Vallicula multiformis</i>
Clean Transcripts	1.5	117218	5419
	2	114299	5030
	3	98718	4828
Cross Contaminated Transcripts	1.5	2574	52302
	2	2296	47699
	3	2062	33675
Dubious Transcripts	1.5	2229	5384
	2	5426	10376
	3	21241	24602