

# The transcriptome of *Anguillicola crassus* sampled by pyrosequencing

Emanuel G Heitlinger<sup>\*1,2</sup> Stephen Bridgett<sup>3</sup> Anna Montazam<sup>3</sup> Horst Taraschewski<sup>1</sup> and Mark Blaxter<sup>2</sup>

<sup>1</sup>Department of Ecology and Parasitology, Zoological Institute 1, University of Karlsruhe, Kornblumenstrasse 13, Karlsruhe, Germany

<sup>2</sup>Institute of Evolutionary Biology, The Ashworth laboratories, The University of Edinburgh, King's Buildings Campus, Edinburgh, UK <sup>3</sup>The GenePool Sequencing Service, The Ashworth laboratories, The University of Edinburgh, King's Buildings Campus, Edinburgh, UK

Email: Emanuel G Heitlinger<sup>\*</sup> - emanuelheitlinger@gmail.com; Stephen Bridgett - sbridgett@staffmail.ed.ac.uk; Anna Montazam - Anna.Montazam@ed.ac.uk; Horst Taraschewski - dc20@rz.uni-karlsruhe.de; Mark Blaxter - mark.blaxter@ed.ac.uk;

<sup>\*</sup>Corresponding author

## Abstract

---

**Background:** *Anguillicola crassus* has been introduced from Asia, where it parasitises the Japanese eel *Anguilla japonica*, to Europe 30 years ago. Here it parasitises the endangered, commercially exploited European eel *Anguilla anguilla*, but differences in the parasite's phenotype and host relations (natural versus novel host) are not understood yet. Phylogenetics places *A. crassus* at a key position for the emergence of parasitism, basal to one of the major clades of parasitic nematodes.

**Results:** After extensive screening of 756,363 raw pyrosequencing reads, we assembled 353,055 into 11,372 contigs spanning 6,575,121 bases and additionally obtained 21,153 singletons and lower quality contigs spanning 6,157,974 bases. We obtained annotations for roughly 55% of the contigs and roughly 30% of the tentatively unique genes (TUGs) confirming the high quality of especially the contigs. We identified 5,112 high quality single nucleotide polymorphisms (SNPs) and suggest 199 of them as most suitable markers for population-genetic studies. The correlation between different analyses provided further insights and confirmed biologically relevant expectations: we found an overabundance of predicted signal peptide cleavage sites in sequence conserved in Nematoda and novel in *A. crassus*, correlations between coding polymorphism and differential expression and between evolutionary conservation and presence of orthologs with lethal

RNAi-phenotypes in *C. elegans*. GO-term analysis identified an enrichment of peptidases and subunits of the respiratory chain for transcripts under positive selection. Enzymes for energy metabolism were also found enriched in genes differentially expressed between European and Asian *A. crassus*.

**Conclusions:** The transcriptome of *A. crassus* is a basis for molecular research on this important species. It furthermore has the potential to provide unique insights into the evolution of parasitism in the Spirurina. We identified energy metabolism as a candidate phenotype for differences between European and Asian worms due to modification or even divergent evolution of gene expression.

---

## Background

The nematode *Anguillicola crassus* Kuwahara, Niimi et Itagaki, 1974 is a native parasite of the Japanese eel *Anguilla japonica* [1]. Adults localise to the swim bladder where they feed on blood [2]. Larvae are transmitted via crustacean intermediate hosts [3]. Originally endemic to East-Asian populations of the Japanese eel (*Anguilla japonica*), *A. crassus* has attracted interest due to recent anthropogenic expansion of its geographic and host ranges to Europe and the European eel (*Anguilla anguilla*). Recorded for the first time in 1982 in North-West Germany [4], where it was most likely introduced through live-eel trade [5,6], *A. crassus* has spread rapidly through populations of its newly acquired host [7]. At the present day it is found in all *An. anguilla* populations except those in Iceland [8]. *A. crassus* can be regarded as a model for invasive parasite introduction and spread [9].

In its colonised host prevalence and mean intensity of infection are higher than in *An. japonica* [10,11], which is accompanied by a larger body mass of adult worms, an earlier onset of reproduction and a larger egg output [12]. These modifications of the lifecycle as well as the virulence of *A. crassus* in its new host have been attributed to an inadequate immune response in *An. anguilla* [13]. Only *An. japonica* is capable of killing histotropic larvae of the parasite after vaccination [14] or under high infection pressure [15]. Accordingly mainly *An. anguilla* is affected by pathology, including thickening and inflammation of the swim bladder wall [16].

The genus *Anguillicola* is placed in the nematode suborder Spirurina (clade III *sensu* [17]) [18,19]. The Spirurina are exclusively parasitic and include important human pathogens (the causative agents of

filariasis and ascariasis) as well as prominent veterinary parasites. Molecular phylogenetic analyses place *Anguillicola* in a clade of spirurine nematodes (Spirurina B of [20]) that have an freshwater or marine intermediate host, but infect a wide range of carnivorous definitive hosts. Spirurina B is sister to the main Spirurina C, including the agents of filariasis and ascariasis), and thus *A. crassus* may be used as an outgroup taxon to understand the evolution of parasitic phenotypes in these species.

Recent advances in sequencing technology (often termed Next Generation Sequencing; NGS), provide the opportunity for rapid and cost-effective generation of genome-scale data. The Roche 454 platform [21] offers longer reads than other NGS technologies, and thus is suited to de novo assembly of genome-scale data in previously understudied species. Roche 454 data has particular application in transcriptomics [22]. The difference in the biology of *A. crassus* in *An. japonica* (coevolved) and *An. anguilla* (recently captured) eel hosts likely results from an interaction between different host and parasite responses, underpinned by definitive differences in host genetics, and possible genetic differentiation between the invading European and endemic Asian parasites. As part of a programme to understand the invasiveness of *A. crassus* in *An. anguilla*, we are investigating differences in gene expression and genetic distinction between invading European and endemic Asian *A. crassus* exposed to the two different host species. Here we report on the generation of a reference transcriptome for *A. crassus* based on Roche 454 data, and explore patterns of gene expression and diversity.

## Methods

### Nematode samples, RNA extraction, cDNA synthesis and Sequencing

*A. crassus* from *An. japonica* were sampled from Kao-Ping river and an adjacent aquaculture in Taiwan as described in [15]. Worms from *An. anguilla* were sampled in Sniardwy Lake, Poland (53.751959N, 21.730957E) and from the Linkenheimer Altrhein, Germany (49.0262N, 8.310556E). After determination of the sex of adult nematodes, they were stored in RNA-later (Quiagen, Hilden, Germany) until extraction of RNA. RNA was extracted from individual adult male and female nematodes and from a population of L2 larvae (Table 1). RNA was reverse transcribed and amplified into cDNA using the MINT-cDNA synthesis kit (Evrogen, Moscow, Russia). For host contamination screening a liver-sample from an uninfected *An. japonica* was also processed. Emulsion PCR was performed for each cDNA library according to the manufacturer's protocols (Roche/454 Life Sciences), and sequenced on a Roche 454 Genome Sequencer FLX. Raw sequencing reads are archived under study-accession number SRP010313 in the NCBI Sequence Read Archive (SRA; <http://www.ncbi.nlm.nih.gov/Traces/sra>) [23].

All samples were sequenced using the FLX Titanium chemistry, except for the Taiwanese female sample T1, which was sequenced using FLX standard chemistry, to generate between 99,000 and 209,000 raw reads. For the L2 larval library, which had a larger number of non-*A. crassus*, non-*Anguilla* reads, we confirmed that these data were not laboratory contaminants by screening Roche 454 data produced on the same run in independent sequencing lanes.

### **Trimming, quality control and assembly**

Raw sequences were extracted in fasta-format (with the corresponding qualities files) using sffinfo (Roche/454) and screened for adapter sequences of the MINT-amplification-kit using cross-match [24] (with parameters -minscore 20 -minmatch 10). Seqclean [25] was used to identify and remove poly-A-tails, low quality, repetitive and short (<100 base) sequences. All reads were compared to a set of screening databases using BLAST (expect value cutoff  $E < 1e-5$ , low complexity filtering turned off: -F F). The databases used were (a) a host sequence database comprising an assembly of the *An. japonica* Roche 454 data, a unpublished assembly of *An. anguilla* Sanger dideoxy sequenced expressed sequence tags (made available to us by Gordon Cramb, University of St Andrews) and transcripts from EelBase [26] a publicly available transcriptome database for the European eel; (b) a database of ribosomal RNA (rRNA) sequences from eel species derived from our Roche 454 data and EMBL-Bank; and (c) a database of rRNA sequences identified in our *A. crassus* data by comparing the reads to known nematode rRNAs from EMBL-Bank. This last database notably also contained xenobiont rRNA sequences. Reads with matches to one of these databases over more than 80% of their length and with greater than 95% identity were removed from the dataset. Screening and trimming information was written back into sff-format using sffile(Roche 454). The filtered and trimmed data were assembled using the combined assembly approach [22]: two assemblies were generated, one using Newbler v2.6 [21] (with parameters -cdna -urt), the other using Mira v3.2.1 [27] (with parameters -job=denovo,est,accurate,454). The resulting two assemblies were combined into one using Cap3 [28] at default settings and contigs were labeled by whether they derived from both assemblies or one assembly only (for a detailed analysis of the assembly categories see the supporting Methods file).

### **Post-assembly classification and taxonomic assignment of contigs**

After assembly contigs were assessed a second time for host and other contamination by comparing them (using BLAST) to the three databases defined above, and also to nembase4, a nematode transcriptome database derived from whole genome sequencing and EST assemblies [29,30]. For each contig, the

highest-scoring match was recorded as long as it spanned more than 50% of the contig. We also compared the contigs to the NCBI non-redundant nucleotide (NCBI-nt) and protein (NCBI-nr) databases, recording the taxonomy of all best matches with expect values better than 1e-05. TUGs with a best hit to non-Metazoans and to Chordata within Metazoa were additionally excluded from further analysis.

### Protein prediction and annotation

Protein translations were predicted from the contigs using prot4EST (version 3.0b) [31]. Proteins were predicted either by joining single high scoring segment pairs (HSPs) from a BLAST search of uniref100 [32], or by ESTscan [33], using as training data the *Brugia malayi* complete proteome back-translated using a codon usage table derived from the BLAST HSPs, or, if the first two methods failed, simply the longest ORF in the contig. For contigs where the protein prediction required insertion or deletion of bases in the original sequence, we also imputed an edited sequence for each affected contig. Annotations with Gene Ontology (GO), Enzyme Commission (EC) and Kyoto Encyclopedia of Genes and Genomes (KEGG) terms were inferred for these proteins using Annot8r (version 1.1.1) [34], using the annotated sequences available in uniref100 [32]. Up to 10 annotations based on a BLAST similarity bitscore cut-off of 55 were obtained for each annotation set. The complete *B. malayi* proteome (as present in uniref100) and the complete *C. elegans* proteome (as present in wormbase v.220) were also annotated in the same way. SignalP V4.0 [35] was used to predict signal peptide cleavage sites and signal anchor signatures for the *A. crassus*-transcriptome and similarly again for the proteomes of the two model-worms. Additionally InterProScan [36] (command line utility iprscan version 4.6 with options -cli -format raw -iprlookup -seqtype p -goterms) was used to obtain domain based annotations for the high credibility assembly (highCA) derived contigs.

We recorded the presence of a lethal RNAi-phenotype in the *C. elegans* ortholog of each TUG using the biomart-interface [37] to wormbase v. 220 through the R-package biomaRt [38].

### Single nucleotide polymorphism analysis

We mapped the raw reads against the the complete set of contigs, replacing imputed sequences for originals where relevant, using ssaha2 [39] (with parameters -kmer 13 -skip 3 -seeds 6 -score 100 -cmatch 10 -ckmer 6 -output sam -best 1). From the ssaha2 output, pileup-files were produced using samtools [40], discarding reads mapping to multiple regions. VarScan [41] (pileup2snp) was used with default parameters on pileup-files to output lists of single nucleotide polymorphisms (SNPs) and their locations. For enrichment

analysis of GO-terms we used the R-package GOstats [42].

Using Samtools [40] (mpileup -u) and Vcftools [43] (view -gcv) we genotyped individual libraries for the list of previously found overall SNPs. Genotype-calls were accepted at a phred-scaled genotype quality threshold of 10. In addition to the relative heterozygosity (number of homozygous sites/number of heterozygous sites) we used the R package Rhh [44] to calculate internal relatedness [45], homozygosity by loci [46] and standardised heterozygosity [47] from these data.

We confirmed the significance of heterozygote-heterozygote correlation by analysing the mean and 95% confidence intervals from 1000 bootstrap replicates estimated for all measurements.

### **Gene-expression analysis**

Read-counts were obtained from the bam-files generated also for genotyping using the R-package Rsamtools [48]. Counts to off target data and lowCA contigs were disregarded. Furthermore contigs with less than 32 reads over all libraries were excluded from analysis, to avoid inference based on too low overall expression values. Because very low coverage from library E1 and L2 leading highly variable normalised data, we excluded these libraries from analysis.

The statistic of Audic and Claverie [49] as implemented in ideg6 [50] was used to contrast single libraries. Differential expression between libraries from different sex of worms was accepted for genes differing between all female libraries E2, T1 and T2 versus the male (M) library ( $p < 0.01$ ) but not within any of the female libraries at the same threshold. Differential expression between libraries from European and Asian origin was accepted for genes differing between libraries E2 versus T1 and T2 ( $p < 0.01$ ) but not between T1 versus T2.

### **Over-representation analyses**

Prior to analysis of GO-term over-representation (based on dn/ds or expression values) we used the R-package annotationDbi [51] to obtain a full list of associations (also with higher-level terms) from Annot8r-annotations. We then used the R-package topGO [52] to traverse the annotation-graph and analyse each node in the annotation for over-representation of the associated term in the focal gene-set compared to a appropriate universal gene-set (all contigs with dn/ds values or all contigs analysed for gene-expression) with the “classic” method and Fisher’s exact test. From the resulting tables we removed uninformative terms, for which an ancestral term already was already in the table and no additional counts supported overrepresentation.

We used Mann-Whitney u-tests to test the influence of factors on dn/ds values, when multiple contrasts between groups (factors) were investigated we used Nemenyi-Damico-Wolfe-Dunn tests. For overrepresentation of one group (factor) in other groups (factors) we used Fisher’s exact test.

## General coding methods

The bulk of analysis (unless otherwise cited) presented in this paper was carried out in R [53] using custom scripts. We used a method provided in the R-packages Sweave [54] and Weaver [55] for “reproducible research” combining R and L<sup>A</sup>T<sub>E</sub>Xcode in a single file. All intermediate data files needed to compile the present manuscript from data-sources are provided upon request. For visualisation we used the R-packages ggplot2 [56] and VennDiagram [57].

## Results

### Sampling *A. crassus*

One female worm and one male worm were sampled from an aquaculture with height infection loads in Taiwan. An additional female worm was sampled from a stream with low infection pressure adjacent to the aquaculture. All these worms were parasitising endemic *An. japonica*. A female worm and pool of L2 larval stages were sampled from *An. anguilla* in the river Rhine, one female worm from a lake in Poland. All adult worms were filled with large amounts of host-blood, therefore we anticipated abundant host-contamination in sequencing data and decided to sequence a liver sample of an uninfected *An. japonica* for screening.

### Sequencing, trimming and pre-assembly screening

A total of 756,363 raw sequencing reads were generated for *A. crassus* (Table 1). These were trimmed for base call quality, and filtered by length to give 585,949 high-quality reads (spanning 169,863,104 bases). In the eel dataset from 159,370 raw reads 135,072 were assembled after basic quality screening.

We then screened the *A. crassus* reads for contamination by host (30,071 reads matched previously sequenced eel genes or our own *An. japonica* 454 transcriptome, which had been assembled into 10,639 mRNA contigs. 181,783 reads matched large or small subunit nuclear or mitochondrial ribosomal RNA sequences of *A. crassus* (Table 1). In addition to fish mRNAs, we identified (and removed) 5,286 reads in the library derived from the L2 nematodes that had significant similarity to cercozoan (likely parasite) ribosomal RNA genes (Table 1).

## Assembly and taxonomic classification

We assembled the remaining 353,055 reads (spanning 100,491,819 bases) using the combined assembler strategy [22] and Roche 454 GSAssembler (version 2.6) and MIRA (version 3.21) [27]. From this we derived 13,851 contigs that were supported by both assembly algorithms, 3,745 contigs only supported by one of the assembly algorithms and 22,591 singletons that were not assembled by either approach (Table 2). When scored by matches to known genes, the contigs supported by both assemblers are of the highest credibility, and this set is thus termed the high credibility assembly (highCA). Those with evidence from only one assembler and the singletons are of lower credibility (lowCA). These datasets are the most parsimonious (having the smallest size) for their quality (covering the largest amount of sequence in reference transcriptomes). In the highCA parsimony and low redundancy is prioritized, while in the complete assembly (highCA plus lowCA) completeness is prioritized. The 40187 sequences (contig consensus and singletons) in the complete assembly are referred to below as tentatively unique genes (TUGs).

We screened the complete assembly for residual host contamination, and identified 3,441 TUGs that had higher, significant similarity to eel (and chordate) sequences (our and publicly available 454 ESTs and EMBLBank Chordata proteins) than to nematode sequences [30].

Given our prior identification of cercozoan ribosomal RNAs, we also screened the complete assembly for contamination with other transcriptomes.

1,153 TUGs were found mapping to Eukaryota outside of the kingdoms Metazoa, Fungi and Viridiplantae. These hits included a wide range of protists ranging from Apicomplexa (mainly Sarcocystidae, 28 hits and Cryptosporidiidae 10 hits) over Bacillariophyta (diatoms, mainly Phaeodactylaceae, 41 hits) and Phaeophyceae (brown algae, mainly Ectocarpaceae, 180 hits) and Stramenopiles (Albuginaceae, 63 hits) to Kinetoplastida (Trypanosomatidae, 26 hits) and Heterolobosea (Vahlkampfiidae, 38 hits).

Additionally we found 298 TUGs with hits to fungi (e.g. Ajellomycetaceae, 53 hits) and 585 TUGs with hits to plants.

Hits outside the Eukaryota were mainly to Bacteria (825 hits) and within those mostly to members of the Proteobacteria (484 hits). No hits were found to Wolbachia or related Bacteria known as symbionts of nematodes and arthropods. 9 TUGs were hitting sequence from Viruses and 8 from Archaea.

We excluded all TUGs with best hits outside Metazoa and our assembly thus has 32,525 TUGs, spanning 12,733,095 bases (of which 11,372 are highCA-derived, and span 6,575,121 bases) that are likely to derive from of *A. crassus*.



## Protein prediction

For 32,418 TUGs a protein was predicted using prot4EST [31] (Table 2). The full open reading frame was obtained in 353 TUGs, while for 2,683 the 5' end and for 8,283 the 3' end was complete. In 13,383 TUGs the corrected sequence with the imputed ORF was slightly changed compared to the raw sequence.

## Annotation

We obtained basic annotations with orthologous sequences from *C. elegans* for 9,556 TUGs, from *B. malayi* for 9,664 TUGs, from nempep [29,30] for 11,620 TUGs and with uniprot proteins for 11,115 TUGs.

We used Annot8r [34] to assign gene ontology (GO) terms for 6,511 TUGs, Enzyme Commission (EC) numbers for 2,460 TUGs and Kyoto Encyclopedia of Genes and Genomes (KEGG) pathway annotations for 3,846 TUGs (Table 2). Additionally 5,125 highCA derived contigs were annotated with GO terms through InterProScan [36]. Nearly one third (6,989) of the *A. crassus* TUGs were annotated with at least one identifier, and 1,831 had GO, EC and KEGG annotations (Figure 1).

We compared our *A. crassus* GO annotations for high-level GO-slim terms to the annotations (obtained the same way) for the complete proteome of the filarial nematode *B. malayi* and the complete proteome of *C. elegans* (Figure 2).

Correlation shows the occurrence of terms for the partial transcriptome of *A. crassus* to be more similar to the proteome of *B. malayi* (0.95; Spearman correlation coefficient) than to the proteome of *C. elegans* (0.9). Also the two model-nematode compared to each other (0.91) are less similar in the occurrence of terms than the two parasites.

We inferred presence of signal peptide cleavage sites in the predicted protein sequence using SignalP [35].

We predicted 920 signal peptide cleavage sites and 65 signal peptides with a transmembrane signature.

Again these predictions are more similar to predictions using the same methods for the proteome *B. malayi* (742 signal peptide cleavage sites and 41 with transmembrane anchor) than for the proteome of *C. elegans* (4,273 signal peptide cleavage sites and 154 with transmembrane anchor).

We inferred the presence of a lethal RNAi phenotype in the orthologous annotation of *C. elegans*. For 259 TUGs a non-lethal phenotype was inferred, for 6,029 TUGs a lethal phenotype.

## Evolutionary conservation

*A. crassus* TUGs were classified as conserved, conserved in Metazoa, conserved in Nematoda, conserved in Spirurina or novel to *A. crassus* by comparing them to public databases and using two BLAST bit-score

cutoffs to define relatedness (Table 3).

Roughly a third and a quarter of the highCA derived contigs were categorized as conserved across kingdoms at a bitscore threshold of 50 and 80, respectively. Roughly half or 3/5 of the these contigs were identified as novel in *A. crassus*.

The remaining highCA contigs spread across intermediate relatedness-levels. More sequences were categorised as novel at the phylum level (Nematoda) compared to kingdom and clade III level and the number of contigs at intermediate relatedness-levels was roughly consistent for the two bitscore thresholds. The latter points about intermediate conservation levels were also true, when all TUGs were analysed. The numbers of TUGs categorised at these intermediate levels roughly doubled. In contrast, the proportion of additional conserved lowCA TUGs was small compared to additional TUGs categorised as novel in *A. crassus*, mirroring the higher amount of erroneous sequence.

Proteins predicted to be novel to Nematoda and novel in *A. crassus* were significantly enriched in signal peptide annotation compared to conserved proteins, proteins novel in Metazoa and novel in clade III (Fisher's exact test  $p < 0.001$  ; Figure 3).

The proportion of lethal RNAi phenotypes was significantly higher for orthologs of TUGs conserved across the kingdom level (97.23%) than for orthologs of TUGs not conserved (94.59%) across kingdoms ( $p < 0.001$ , Fisher's exact test).

### Identification of single nucleotide polymorphisms

We called single nucleotide polymorphisms (SNPs) on the 1,100,522 bases of the TUGs that had coverage of more than 8-fold available using VARSscan [41]. We excluded SNPs predicted to have more than 2 alleles or that mapped to an undetermined (N) base in the reference, and retained 10,496 SNPs. The ratio of transitions (ti; 6,908) to transversion (tv; 3,588) in this set was 1.93. Using the prot4EST predictions and the corrected sequences, 7,189 of the SNPs were predicted to be inside an ORF, with 2,322 at codon first positions, 1,832 at second positions and 3,035 at third positions. As expected ti/tv inside ORFs (2.39) was higher than outside ORFs (1.25). The ratio of synonymous polymorphisms per synonymous site to non-synonymous polymorphisms per non-synonymous site (dn/ds) was 0.45. We filtered these SNPs to exclude those that might be associated with analytical bias. As Roche 454 sequences have well-known systematic errors associated with homopolymeric nucleotide sequences [58], we analysed the effect of exclusion of SNPs in, or close to, homopolymer regions. We observed changes in ti/tv and in dn/ds when SNPs were discarded using different size thresholds for homopolymer runs and proximity thresholds

(Figure 4). Based on this we decided to exclude SNPs with a homopolymer-run as long as or longer than 4 bases inside a window of 11 bases (5 to bases to the right, 5 to the left) around the SNP. We also observed a relationship between TUG dn/ds and TUG coverage, associated with the presence of sites with low abundance minority alleles (less than 7% of the allele calls), suggesting that some of these may be errors. Removing low abundance minority allele SNPs from the set removed this effect (Figure 5). Our filtered SNP dataset includes 5,113 SNPs. We retained 4.65 SNPs per kb of contig sequence, with 7.95 synonymous SNPs per 1000 synonymous bases and 2.44 non-synonymous SNPs per 1000 non-synonymous bases. A mean dn/ds of 0.244 was calculated for the 765 TUGs (765 highCA-derived contigs) containing at least one synonymous SNP.

### **Polymorphisms associated with biological processes**

We consolidated our annotation and polymorphism analyses by examining correlations between nonsynonymous variability and particular classifications.

Signal peptide containing proteins have been shown to have higher rates of evolution than cytosolic proteins in a number of nematode species. In *A. crassus*, TUGs predicted to contain signal peptide cleavage sites in SignalP showed non-significant a trend towards higher dn/ds values than TUGs without signal peptide cleavage sites ( $p = 0.22$ ; two sided Mann-Whitney-test).

Positive selection can be inferred from dn/ds analyses, and we defined TUGs with a dn/ds higher than 0.5 as positively selected. We identified over-represented GO ontology terms associated with these putatively positively selected genes (Table 4; Additional Figures 1). Within the molecular function category, “peptidase activity” was the most significantly overrepresented term and had twelve TUGs supporting the overrepresentation. The highlighted twelve peptidases annotated with twelve unique orthologs in *C. elegans* and *B. malayi*.

Other overrepresented terms abundant over categories pointed to subunits of the respiratory chain e.g. “heme-copper terminal oxidase activity” and “cytochrome-c oxidase activity” in molecular function and “mitochondrion” in cellular compartment (Table 4 and Additional Figures 1).

At both bitscore thresholds contigs novel in clade III and novel in *A. crassus* had a significantly higher dn/ds than other contigs (Figure 6, novel.in.metazoa - novel.in.Ac, 0.005 and 0.015; novel.in.nematoda - novel.in.Ac, 0.005 and 0.002; novel.in.nematoda - novel.in.clade3, 0.207 and 0.045; comparison, p-value from bitscore of 50 and p-value from bitscore of 80, Nemenyi-Damico-Wolfe-Dunn test, given only for significant comparisons).

Orthologs of *C. elegans* transcripts with lethal RNAi phenotype are expected to evolve under stronger selective constraints. Indeed the values of  $dn/ds$  showed a non-significant trend towards lower values in TUGs with orthologs with a lethal phenotype compared to a non-lethal phenotypes ( $p=0.815$ , two-sided U-test).

### SNP markers for single worms

We used Samtools [40] and Vcftools [43] to call genotypes in single worms (adult sequencing libraries). This resulted in 199 informative sites in 152 contigs, where two alleles were found in at least one assured genotype at least in one of the worms.

Internal relatedness [45], homozygosity by loci [46] and standardised heterozygosity [47] were all highlighting the Taiwanese worm from aquaculture (sample T1) as the most and the European worm from Poland (sample E2) as the least heterozygous individual. The other worms had intermediate values between these two extremes.

We confirmed the genome-wide significance of these estimates using heterozygosity-heterozygosity correlation [44]. These tests confirmed the representativeness of the 199 SNP-markers for the whole genome in population genetic studies ( $\mu = 0.78$ ,  $ci_l=0.444$ ;  $\mu = 0.86$  and  $ci_l = 0.596$ ;  $\mu = 0.87$  and  $ci_l = 0.632$ ; mean and lower bound of 95% confidence intervals from 1000 bootstrap replicates for internal relatedness, homozygosity by loci and standardised heterozygosity). Using a higher number of genotyped individuals these markers would allow to assess the amount of inbreeding in populations of *A. crassus*.

### Differential expression

We also analysed gene-expression inferred from mapping. Of the 353,055 reads 252,388 (71.49%) mapped uniquely (with their best hit) to the fullest assembly (including the all assembled contigs as a “filter” later removing screened out sequences for analysis). The number of reads mapping is given for each library Table 1, to get unbiased estimates of expression we removed also all contigs with a coverage lower than 32 reads overall and thus analysed 658 contigs using ideg6 [50] for normalisation the statistic of Audic and Claverie [49] for detection of differences.

54 contigs showed an expression predominantly in the male library, 56 contigs in the female library. 56 contigs were primarily expressed in the libraries from Taiwan, 22 contigs in the European library.

Overrepresentation of GO-terms differentially expressed between the male and female libraries highlighted especially ribosomal proteins oxidoreductases and collagen processing enzymes as enriched

(Table 6a and Additional Figures 1). Ribosomal proteins were all overexpressed in the male library, oxidoreductases and collagen processing enzymes were overexpressed female libraries. Overrepresentation of GO-terms differentially expressed between libraries from worms of European and Asian origin highlighted catalytic activity especially related to metabolism (Table 6b; Additional Figures 1). Acyltransferase contigs were all upregulated in the European libraries. However, the expression patterns for other contigs connected to metabolism did not show concerted up or down-regulation (eg. for “steroid biosynthetic process” 2 contigs were downregulated in the European library, 3 contigs upregulated). Enrichment of signal-positives was not found in any category of overexpressed genes. Differentially expressed genes also showed no pattern of enrichment in conservation categories and no enrichment of *C. elegans* orthologs with lethal/non-lethal RNAi-phenotypes. Significantly elevated dn/ds was found for contigs differentially expressed according to worm-origin (Fisher’s exact test  $p=0.007$ ; also both up- or downregulated were significant). Contigs overexpressed in the female libraries showed elevated levels of dn/ds (Fisher’s exact test  $p=0.041$ ). In contrast male overexpressed genes showed decreased levels of dn/ds (Fisher’s exact test  $p=0.014$ ).

## Discussion

We have generated a de novo transcriptome for *A. crassus* an important invasive parasite that threatens wild stocks of the European eel *An. anguilla*. These data enable a broad spectrum of molecular research on this ecologically important and evolutionary interesting parasite. As *A. crassus* lives in close association with its host, we have used exhaustive filtering to attempt to remove all host-derived, and host-associated organism-derived contamination from the data. To do this we have also generated a transcriptome dataset from the definitive host *An. japonica*. The non-nematode, non-eel data identified, particularly in the L2 sample, showed highest identity to flagellate protists, which may have been parasitising the eel (or the nematode). Encapsulated objects observed in eel swim bladder walls [15] could be due solely to immune attrition of *A. crassus* larvae or to other coinfections.

A second examination of sequence origin was performed after assembly, employing higher stringency cutoffs. Similar taxonomic screening was used in a garter snake transcriptome project [59], and an analysis of lake sturgeon tested and rejected hypotheses of horizontal gene-transfer when xenobiont sequences were identified [60]. A custom pipeline for transcriptome assembly from pyrosequencing reads [61] proposed the use of EST3 [62] to infer sequence origin based simply on nucleotide frequency. We were not able to use this approach successfully, probably due to the fact that xenobiont sequences in our data set derive from

multiple sources with different GC content and codon usage.

Compared to other NGS transcriptome sequencing projects [63], the combined assembly approach generated a smaller number of contigs that had lower redundancy and higher completeness. Projects using the mira assembler often report substantially greater numbers of contigs for datasets of similar size (see e.g. [64]), comparable to the mira sub-assembly in our approach. The use of oligo(dT) to capture mRNAs probably explains the bias towards 3' end completeness and a relative lack of true initiation codons in our protein prediction. This bias is near-ubiquitous in deep transcriptome sequencing projects (e.g. [65]).

\*\*\*START: I moved the chunk below up (from a paragraph you did not correct last time) to have annotation before SNPs

We were able to obtain high-quality annotations for a large set of TUGs: For roughly 30% of the complete assembly and over 50% of our highCA assembly BLAST-based annotations could be obtained. 45% of the contigs in the highCA assembly were additionally decorated with domain-based annotations through InterProScan [36].

Comparison with complete protein sequence from the genomes of *B. malayi* and *C. elegans* showed a remarkable degree of agreement regarding the occurrence of terms in the two parasitic worms. This agreement was higher than with the free living nematode *C. elegans* and even the two genome-sequencing-derived proteomes showed less agreement with each other than the filarial parasite with our dataset. This implies that our transcriptome is truly a representative partial genome [66] of a parasitic nematode.

Analysis of conservation identified more sequence novel in Nematode than in the eukaryote kingdom or in clade III this is in agreement with prevalence of genic novelty in the Nematoda [67]. Furthermore the basal position of *A. crassus* in clade III could be leading to most novelty in the clade not being shared with *A. crassus*.

TUGs predicted to be novel in the phylum Nematoda and novel to *A. crassus* contained the highest proportion of signal-positives. This confirms observations made in a study on *Nippostrongylus brasiliensis* [68], where signal positives were reported as less conserved. Interestingly enrichment of signal sequence bearing TUGs in our dataset was constrained to sequences novel in nematodes and *A. crassus* (i.e. not to the level of clade III). This may be explained, with two different hypotheses involving the basal position of *A. crassus*: first the signal positives shared with all nematodes could be conserved molecules not excreted by parasites. A different class of secreted/excreted molecules with prominent role in host parasite interactions would not have arisen early in the evolution of parasitism in clade III - or be too fast-evolving

- and thus be detected as specific to deeper sub-clades (i.e. to *A. crassus* in our dataset). A second explanation would be, that orthologs of excreted parasite-specific genes could be among those shared with other nematodes and the fewer shared with clade III implying a predisposition to parasitism outside of the Spirurina or even the convergent evolution of secreted molecules in other parasitic nematodes. However analysis of dn/ds (see below) across conservation categories favors the first hypothesis, as it identifies a higher amount of positive selection in TUGs novel to clade III and *A. crassus* than to nematodes.

\*\*\*End: I moved the above chunk up to have annotation before SNPs

We generated transcriptome data from multiple *A. crassus* of Taiwanese and European origin, and identified SNPs both within and between populations. Screening of SNPs in or adjacent to homopolymer regions improved overall measurements of SNP quality. The ratio of transitions to transversions (ti/tv) increased. Such an increase is explained by the removal of “noise” associated with common homopolymer errors [58]. The value of 1.925 (1.25 outside, 2.39 inside ORFs) is in good agreement with the overall ti/tv of humans (2.16 [69]) or *Drosophila* (2.07 [70]). The ratio of non-synonymous SNPs per non-synonymous site to synonymous SNPs per synonymous site (dn/ds) decreased with removal of SNPs adjacent to homopolymer regions from 0.45 to 0.244 after full screening. The most plausible explanation is the removal of error, as unbiased error would lead to a dn/ds of 1. While dn/ds is not unproblematic to interpret within populations [71], the assumption of negative (purifying) selection on most protein-coding genes makes lower mean values seem more plausible. We used a threshold value for the minority allele of 7% for exclusion of SNPs, based on an estimate that approximately 10 haploid equivalents were sampled (5 individual worms plus an negligible contribution from L2 larvae in the L2 library and within the female adult worms). The benefit of this screening was mainly a reduction of non-synonymous SNPs in high coverage contigs, and a removal of the dependence of dn/ds on coverage. Working with an estimate of dn/ds independent of coverage, efforts to control for sampling biased by depth (i.e. coverage; see [72] and [63]) could be avoided.

\*\*\* you corrected up to her last time

Also in comparison with published intra-species values of dn/ds our final estimate of seems plausible: in transcripts from the female reproductive tract of *Drosophila* dn/ds was 0.15 [73] and 0.21 in the male reproductive tract [74] (although for ESTs specific to the male accessory gland were shown to have a higher dn/ds of 0.47). A pyrosequencing study in the parasitic nematode *Ancylostoma caninum* [75] reported dn/ds of 0.3.

When the whole of coding sequences are studied, of which only a small subset of sites can be under positive selection, dn/ds of 0.5 has been suggested as threshold for assuming positive selection [73] instead of the

classical threshold of 1 [76]. The use of this threshold for positive selection led to the identification of over-represented of GO-term highlighting very interesting transcripts: twelve peptidases under positive selection (from 43 with a  $dn/ds$  obtained) meant an enrichment in the category. All twelve have different orthologs in *B. malayi* and *C. elegans* and are conserved across kingdoms. Despite their conservation peptidases are thought to have acquired new and prominent roles in host-parasite interaction compared to free living organisms: in *A. crassus* a trypsin-like proteinase has been identified thought to be utilised by the tissue-dwelling L3 stage to penetrate host tissue and an aspartyl proteinase thought to be a digestive enzyme in adults [2]. The twelve proteinases under positive selection could be the targets of the adaptive immunity developed against *A. crassus* [14, 77], which is often only elicited against subtypes of larvae [78]. Genotyping of individual worms identified a set of 199 SNPs with highest credibility and a high information content for population-genetic studies. The low number of SNPs inferred with assured genotypes reflects both the additional variance in allele contribution introduced by the sampling process of transcripts involved in the generation of transcriptomic data and the stringency of software more targeted at even higher throughput genotyping (VCFtools is rather designed for genomic data from the solexa platform [79]). Nevertheless, levels of genome-wide heterozygosity found for the 5 adult worms examined in our study are in agreement with microsatellite data [80] showing reduced heterozygosity in European populations of *A. crassus*. The polish female worm from our study can be regarded highly inbred, the worm from a wild *An. japonica* in Taiwan highly outbred.

We employed methods to developed for the comparison of cDNA-libraries to make inference about possible differential gene-expression according to experimental groups (origin of sequencing-libraries) [49]. Such approaches are widely used with pyrosequencing-data (e.g. [75]). For the statistically valid comparison of conditions however, the unit of replication would be the individual library and approaches respecting this fact would be desirable. However, we were not able to use the R-packages DESeq [81] or edgeR [82] developed for count data from deep sequencing (but more targeted towards RNA-seq on the solexa-platform) as both repetition and throughput of our present experiment were too low. As a result the differentially expressed genes are by no means significant for the investigated conditions, but just for the specific cDNA-libraries. With these reservations we identified genes differentially expressed between libraries prepared from worms of different sex and worms from different origin.

Genes over-expressed in male *A. crassus* comprise major sperm proteins well known for their high expression in nematode sperm [83]. A surprise was the overexpression of ribosomal proteins in the male library.



That collagen processing enzymes are overexpressed in female worms, filled with developing embryos and larvae, is in line with a complicated regulation and modulation of collagen in nematode larval development [84].

The overexpression acetyl-CoA acetyltransferase in European worms are interesting especially because of the role of these enzymes in fatty-acid  $\beta$ -oxidation in peroxisomes and mitochondria [85]. Together with a change in steroid metabolism and the enrichment of mitochondrially localized enzymes these suggest changes in energy metabolism of *A. crassus* from different origins. Possible explanations would include a change to more or less aerobic processes in worms in Europe due to their bigger size and/or increased availability of nutrients.

Contigs overexpressed in the female libraries showed elevated levels of dn/ds but genes overexpressed in males decreased levels of dn/ds. The first finding is unexpected, as overexpressed in female libraries will also contain contigs related to larval development (such as the collagen modifying enzymes discussed above), these larval transcripts in turn are expected to be under purifying selection because of pleiotropic effects of genes in early development [86]. Also the second finding is in slight contrast to published results for male specific traits and transcripts are often showing hallmarks of positive selection [74,87]. . In *Ancylostoma caninum* however, female-specific transcripts showed an enrichment of “parasitism genes” [75] and a possible explanation would be a similar enrichment of positively selected parasitism related in our dataset. For males the decreased dn/ds can be explained by the by the high number of ribosomal proteins, which are all show very low levels of dn/ds (that these proteins are found differentially expressed remains puzzling though), while single transcripts e.g. major sperm protein (expressed in the male library only) showed elevated dn/ds but did not level the overall effect. But this also has a positive aspect: it is unlikely that correlation of differential expression with positive selection results from mapping artifacts, as all the ribosomal proteins identified overexpressed in males have very low dn/ds.

Genes differential expressed according to worm-origin (in either direction) showed significantly elevated levels of dn/ds. This is interpretable as a correlation between sequence evolution and phenotypic modification in different host-environments or even correlation between sequence evolution and evolution of gene-expression. Thus, whether expression of these genes is modified in different hosts or evolved rapidly in a contemporary divergence between European and Asian populations of *A. crassus*, is in the center of a future research program building on the reference transcriptome presented here. For such an analysis it is important to disentangle the influence of the host and the nematode population in a co-inoculation experiment. Such a project will also use the individual worm as the level of replication for “conditions”

(that is, worm-population and host-species) to allow rigid hypothesis testing. Based on the pilot evaluation presented here differences in these factors are expected overlap with differences in male vs. female worms and the careful cross-examination of the above factors with worm-sex is advised.

Population genetic approaches using the SNP-markers presented here directly or populations genomic approaches choosing to use the SNPS found here as gold-standard in comparison with higher throughput technology, constitute another field of future research on *A. crassus*. The maintenance or loss of variation in European populations in or close to genes under general positive selection will be of major interest in such projects.

## **Conclusions**

The *A. crassus* transcriptome provides a basis of molecular research on this ecologically important species. It further allows insight in the evolution of parasitism complementing the catalogue of available transcriptomic data with a member of the Spirurina phylogenetically distant to so far sequenced parasites in this clade. Differences in energy metabolism between European and Asian *A. crassus* constitute a candidate phenotype relevant for phenotypic modification or contemporary divergent evolution as well as for the long term evolution of parasitism.

## **Competing interests**

The authors declare no competing interests.

## **Authors contributions**

EGH and MB conceived and designed the experiments. EGH carried out bioinformatic analyses. SB assisted in bioinformatic analyses. AM prepared sequencing libraries. HT provided close supervision throughout. EGH and MB interpreted results and prepared the manuscript. All authors have read and approved the final manuscript.

## **Acknowledgements**

This work has been made possible through a grant provided to EGH by Volkswagen Foundation, "Förderinitiative Evolutionsbiologie".

## References

1. Kuwahara A, Niimi H, Itagaki H: **Studies on a nematode parasitic in the air bladder of the eel I. Descriptions of *Anguillicola crassa* sp. n. (Philometridea, Anguillicolidae).** *Japanese Journal for Parasitology* 1974, **23**(5):275–279.
2. Polzer M, Taraschewski H: **Identification and characterization of the proteolytic enzymes in the developmental stages of the eel-pathogenic nematode *Anguillicola crassus*.** *Parasitology Research* 1993, **79**:24–7, [<http://www.ncbi.nlm.nih.gov/pubmed/7682326>].
3. De Charleroy D, Grisez L, Thomas K, Belpaire C, Ollevier F: **The life cycle of *Anguillicola crassus*.** *Diseases of Aquatic Organisms* 1990, **8**(2):77–84.
4. Neumann W: **Schwimblasenparasit *Anguillicola* bei Aalen.** *Fischer und Teichwirt* 1985, :322.
5. Kooops H, Hartmann F: ***Anguillicola*-infestations in Germany and in German eel imports.** *Journal of Applied Ichthyology* 1989, **5**:41–45, [<http://onlinelibrary.wiley.com/doi/10.1111/j.1439-0426.1989.tb00568.x/abstract>].
6. Koie M: **Swimbladder nematodes (*Anguillicola* spp.) and gill monogeneans (*Pseudodactylogyrus* spp.) parasitic on the European eel (*Anguilla anguilla*).** *ICES J. Mar. Sci.* 1991, **47**(3):391–398, [<http://icesjms.oxfordjournals.org/cgi/content/abstract/47/3/391>].
7. Kirk RS: **The impact of *Anguillicola crassus* on European eels.** *Fisheries Management & Ecology* 2003, **10**(6):385–394, [<http://dx.doi.org/10.1111/j.1365-2400.2003.00355.x>].
8. Kristmundsson A, Helgason S: **Parasite communities of eels *Anguilla anguilla* in freshwater and marine habitats in Iceland in comparison with other parasite communities of eels in Europe.** *Folia Parasitologica* 2007, **54**(2):141.
9. Taraschewski H: **Hosts and Parasites as Aliens.** *Journal of Helminthology* 2007, **80**(02):99–128, [<http://journals.cambridge.org/action/displayAbstract?fromPage=online&aid=713884>].
10. Münderle M, Taraschewski G, Klar B, Chang CW, Shiao JC, Shen KN, He JT, Lin SH, Tzeng WN: **Occurrence of *Anguillicola crassus* (Nematoda: Dracunculoidea) in Japanese eels *Anguilla japonica* from a river and an aquaculture unit in SW Taiwan.** *Diseases of Aquatic Organisms* 2006, **71**(2):101–8, [<http://www.ncbi.nlm.nih.gov/pubmed/16956057>].
11. Lefebvre FS, Crivelli AJ: ***Anguillicolosis*: dynamics of the infection over two decades.** *Diseases of Aquatic Organisms* 2004, **62**(3):227–32, [<http://www.ncbi.nlm.nih.gov/pubmed/15672878>].
12. Knopf K, Mahnke M: **Differences in susceptibility of the European eel (*Anguilla anguilla*) and the Japanese eel (*Anguilla japonica*) to the swim-bladder nematode *Anguillicola crassus*.** *Parasitology* 2004, **129**(Pt 4):491–6, [<http://www.ncbi.nlm.nih.gov/pubmed/15521638>].
13. Knopf K: **The swimbladder nematode *Anguillicola crassus* in the European eel *Anguilla anguilla* and the Japanese eel *Anguilla japonica*: differences in susceptibility and immunity between a recently colonized host and the original host.** *Journal of Helminthology* 2006, **80**(2):129–36, [<http://www.ncbi.nlm.nih.gov/pubmed/16768856>].
14. Knopf K, Lucius R: **Vaccination of eels (*Anguilla japonica* and *Anguilla anguilla*) against *Anguillicola crassus* with irradiated L3.** *Parasitology* 2008, **135**(5):633–40, [<http://www.ncbi.nlm.nih.gov/pubmed/18302804>].
15. Heitlinger E, Laetsch D, Weclawski U, Han YS, Taraschewski H: **Massive encapsulation of larval *Anguillicoloides crassus* in the intestinal wall of Japanese eels.** *Parasites and Vectors* 2009, **2**:48, [<http://www.parasitesandvectors.com/content/2/1/48>].
16. Würtz J, Taraschewski H: **Histopathological changes in the swimbladder wall of the European eel *Anguilla anguilla* due to infections with *Anguillicola crassus*.** *Diseases of Aquatic Organisms* 2000, **39**(2):121–34, [<http://www.ncbi.nlm.nih.gov/pubmed/10715817>].
17. Blaxter M, De Ley P, Garey J, X Liu L, Scheldeman P, Vierstraete A, Vanfleteren J, Mackey L, Dorris M, Frisse L, Vida J, Thomas W: **A molecular evolutionary framework for the phylum Nematoda.** *Nature* 1998, **392**(6671):71–75, [<http://dx.doi.org/10.1038/32160>].
18. Nadler SA, Carreno RA, Meja-Madrid H, Ullberg J, C Pagan C, Houston R, Hugot J: **Molecular Phylogeny of Clade III Nematodes Reveals Multiple Origins of Tissue Parasitism.** *Parasitology* 2007, **134**(10):1421–1442, [<http://journals.cambridge.org/action/displayAbstract?fromPage=online&aid=1279744>].

19. Wijová M, Moravec F, Horák A, Lukes J: **Evolutionary relationships of Spirurina (Nematoda: Chromadorea: Rhabditida) with special emphasis on dracunculoid nematodes inferred from SSU rRNA gene sequences.** *International Journal for Parasitology* 2006, **36**(9):1067–75, [<http://www.ncbi.nlm.nih.gov/pubmed/16753171>].
20. Laetsch DR, Heitlinger EG, Taraschewski H, Nadler SA, Blaxter M: **The phylogenetics of Anguillicolidae (Nematoda: Anguillicolidea), swimbladder parasites of eels.** *BMC Evolutionary Biology* under review.
21. Margulies M, Egholm M, Altman WE, Attiya S, Bader JS, Bemben LA, Berka J, Braverman MS, Chen YJ, Chen Z, Dewell SB, Du L, Fierro JM, Gomes XV, Godwin BC, He W, Helgesen S, Ho CH, Ho CH, Irzyk GP, Jando SC, Alenquer ML, Jarvie TP, Jirage KB, Kim JB, Knight JR, Lanza JR, Leamon JH, Lefkowitz SM, Lei M, Li J, Lohman KL, Lu H, Makhijani VB, McDade KE, McKenna MP, Myers EW, Nickerson E, Nobile JR, Plant R, Puc BP, Ronan MT, Roth GT, Sarkis GJ, Simons JF, Simpson JW, Srinivasan M, Tartaro KR, Tomasz A, Vogt KA, Volkmer GA, Wang SH, Wang Y, Weiner MP, Yu P, Begley RF, Rothberg JM: **Genome sequencing in microfabricated high-density picolitre reactors.** *Nature* 2005, **437**:376–380, [<http://dx.doi.org/10.1038/nature03959>].
22. Kumar S, Blaxter ML: **Comparing de novo assemblers for 454 transcriptome data.** *BMC Genomics* 2010, **11**:571, [<http://dx.doi.org/10.1186/1471-2164-11-571>].
23. Sayers EW, Barrett T, Benson DA, Bolton E, Bryant SH, Canese K, Chetvernin V, Church DM, Dicuccio M, Federhen S, Feolo M, Fingerman IM, Geer LY, Helmberg W, Kapustin Y, Krasnov S, Landsman D, Lipman DJ, Lu Z, Madden TL, Madej T, Maglott DR, Marchler-Bauer A, Miller V, Karsch-Mizrachi I, Ostell J, Panchenko A, Phan L, Pruitt KD, Schuler GD, Sequeira E, Sherry ST, Shumway M, Sirotkin K, Slotta D, Souvorov A, Starchenko G, Tatusova TA, Wagner L, Wang Y, Wilbur WJ, Yaschenko E, Ye J: **Database resources of the National Center for Biotechnology Information.** *Nucleic Acids Res.* 2012, **40**:13–25. [[Http://www.ncbi.nlm.nih.gov/pubmed/22140104](http://www.ncbi.nlm.nih.gov/pubmed/22140104)].
24. Green P: *PHRAP documentation.* 1994, [<http://www.phrap.org>].
25. Pertea G, Huang X, Liang F, Antonescu V, Sultana R, Karamycheva S, Lee Y, White J, Cheung F, Parvizi B, Tsai J, Quackenbush J: **TIGR Gene Indices clustering tools (TGICL): a software system for fast clustering of large EST datasets.** *Bioinformatics* 2003, **19**:651–652, [<http://www.ncbi.nlm.nih.gov/pubmed/12651724>].
26. Coppe A, Pujolar JM, Maes GE, Larsen PF, Hansen MM, Bernatchez L, Zane L, Bortoluzzi S: **Sequencing, de novo annotation and analysis of the first Anguilla anguilla transcriptome: EelBase opens new perspectives for the study of the critically endangered European eel.** *BMC Genomics* 2010, **11**:635, [<http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3012609>].
27. Chevreux B, Pfisterer T, Drescher B, Driesel AJ, Muller WE, Wetter T, Suhai S: **Using the miraEST assembler for reliable and automated mRNA transcript assembly and SNP detection in sequenced ESTs.** *Genome Res.* 2004, **14**:1147–1159, [<http://www.ncbi.nlm.nih.gov/pmc/articles/PMC419793>].
28. Huang X, Madan A: **CAP3: A DNA sequence assembly program.** *Genome Res.* 1999, **9**:868–877, [<http://www.ncbi.nlm.nih.gov/pmc/articles/PMC310812>].
29. Parkinson J, Whitton C, Schmid R, Thomson M, Blaxter M: **NEMBASE: a resource for parasitic nematode ESTs.** *Nucl. Acids Res.* 2004, **32**(suppl\_1):D427–430, [[http://nar.oxfordjournals.org/cgi/content/abstract/32/suppl\\_1/D427](http://nar.oxfordjournals.org/cgi/content/abstract/32/suppl_1/D427)].
30. Elsworth B, Wasmuth J, Blaxter M: **NEMBASE4: The nematode transcriptome resource.** *Int. J. Parasitol.* 2011, **41**:881–894, [<http://www.ncbi.nlm.nih.gov/pubmed/21550347>].
31. Wasmuth J, Blaxter M: **prot4EST: Translating Expressed Sequence Tags from neglected genomes.** *BMC Bioinformatics* 2004, **5**:187, [<http://www.biomedcentral.com/1471-2105/5/187>].
32. Bairoch A, Bougueleret L, Altairac S, Amendolia V, Auchincloss A, Argoud-Puy G, Axelsen K, Baratin D, Blatter MC, Boeckmann B, Bolleman J, Bollondi L, Boutet E, Quintaje SB, Breuza L, Bridge A, deCastro E, Ciapina L, Coral D, Coudert E, Cusin I, Delbard G, Dornevil D, Roggli PD, Duvaud S, Estreicher A, Famiglietti L, Feuermann M, Gehant S, Farriol-Mathis N, Ferro S, Gasteiger E, Gateau A, Gerritsen V, Gos A, Gruaz-Gumowski N, Hinz U, Hulo C, Hulo N, James J, Jimenez S, Jungo F, Junker V, Kappler T, Keller G, Lachaize C, Lane-Guermonprez L, Langendijk-Genevaux P, Lara V, Lemerrier P, Le Saux V, Lieberherr D,

- Lima TdeO, Mangold V, Martin X, Masson P, Michoud K, Moinat M, Morgat A, Mottaz A, Paesano S, Pedruzzi I, Phan I, Pilboud S, Pillet V, Poux S, Pozzato M, Redaschi N, Reynaud S, Rivoire C, Roehert B, Schneider M, Sigrist C, Sonesson K, Staehli S, Stutz A, Sundaram S, Tognolli M, Verbregue L, Veuthey AL, Yip L, Zuletta L, Apweiler R, Alam-Faruque Y, Antunes R, Barrell D, Binns D, Bower L, Browne P, Chan WM, Dimmer E, Eberhardt R, Fedotov A, Foulger R, Garavelli J, Golin R, Horne A, Huntley R, Jacobsen J, Kleen M, Kersey P, Laiho K, Leinonen R, Legge D, Lin Q, Magrane M, Martin MJ, O'Donovan C, Orchard S, O'Rourke J, Patient S, Pruess M, Sitnov A, Stanley E, Corbett M, di Martino G, Donnelly M, Luo J, van Rensburg P, Wu C, Arighi C, Arminski L, Barker W, Chen Y, Hu ZZ, Hua HK, Huang H, Mazumder R, McGarvey P, Natale DA, Nikolskaya A, Petrova N, Suzek BE, Vasudevan S, Vinayaka CR, Yeh LS, Zhang J: **The Universal Protein Resource (UniProt)** 2009. *Nucleic Acids Res.* 2009, **37**:D169–174, [<http://www.ncbi.nlm.nih.gov/pubmed/18836194>].
33. Iseli C, Jongeneel CV, Bucher P: **ESTScan: a program for detecting, evaluating, and reconstructing potential coding regions in EST sequences.** *Proc Int Conf Intell Syst Mol Biol* 1999, :138–148, [<http://www.ncbi.nlm.nih.gov/pubmed/10786296>].
  34. Schmid R, M B: **annot8r: GO, EC and KEGG annotation of EST datasets.** *BMC Bioinformatics* 2008, **9**:180, [<http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2324097>].
  35. Petersen TN, Brunak S, von Heijne G, Nielsen H: **SignalP 4.0: discriminating signal peptides from transmembrane regions.** *Nat. Methods* 2011, **8**:785–786, [<http://www.ncbi.nlm.nih.gov/pubmed/21959131>].
  36. Zdobnov EM, Apweiler R: **InterProScan—an integration platform for the signature-recognition methods in InterPro.** *Bioinformatics* 2001, **17**:847–848, [<http://www.ncbi.nlm.nih.gov/pubmed/11590104>].
  37. Kasprzyk A: **BioMart: driving a paradigm change in biological data management.** *Database (Oxford)* 2011, **2011**:bar049, [<http://www.ncbi.nlm.nih.gov/pubmed/22083790>].
  38. Durinck S, Spellman PT, Birney E, Huber W: **Mapping identifiers for the integration of genomic datasets with the R/Bioconductor package biomaRt.** *Nat Protoc* 2009, **4**:1184–1191, [<http://www.ncbi.nlm.nih.gov/pubmed/19617889>].
  39. Ning Z, Cox AJ, Mullikin JC: **SSAHA: a fast search method for large DNA databases.** *Genome Res.* 2001, **11**:1725–1729, [<http://www.ncbi.nlm.nih.gov/pmc/articles/PMC311141>].
  40. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis GR, Durbin R: **The Sequence Alignment/Map format and SAMtools.** *Bioinformatics* 2009, **25**(16):2078–2079, [<http://dx.doi.org/10.1093/bioinformatics/btp352>].
  41. Koboldt DC, Chen K, Wylie T, Larson DE, McLellan MD, Mardis ER, Weinstock GM, Wilson RK, Ding L: **VarScan: variant detection in massively parallel sequencing of individual and pooled samples.** *Bioinformatics* 2009, **25**:2283–2285, [<http://www.ncbi.nlm.nih.gov/pubmed/19542151>].
  42. Falcon S, Gentleman R: **Using GStats to test gene lists for GO term association.** *Bioinformatics* 2007, **23**:257–258, [<http://www.ncbi.nlm.nih.gov/pubmed/17098774>].
  43. Danecek P and Auton,†A and Abecasis, G and Albers CA and Banks, E and DePristo, MA and Handsaker RE and Lunter G and Marth GT and Sherry ST and McVean GT and Durbin T and the 1000 Genomes Project: **The variant call format and VCFtools.** *Bioinformatics* 2011, **27**:2156–2158, [<http://www.ncbi.nlm.nih.gov/pubmed/21653522>].
  44. Alho JS, Valimaki K, Merila J: **Rhh: an R extension for estimating multilocus heterozygosity and heterozygosity-heterozygosity correlation.** *Mol Ecol Resour* 2010, **10**:720–722, [<http://www.ncbi.nlm.nih.gov/pubmed/21565077>].
  45. Amos W, Wilmer JW, Fullard K, Burg TM, Croxall JP, Bloch D, Coulson T: **The influence of parental relatedness on reproductive success.** *Proc. Biol. Sci.* 2001, **268**:2021–2027, [<http://www.ncbi.nlm.nih.gov/pubmed/11571049>].
  46. Aparicio JM, Ortego J, Cordero PJ: **What should we weigh to estimate heterozygosity, alleles or loci?** *Mol. Ecol.* 2006, **15**:4659–4665, [<http://www.ncbi.nlm.nih.gov/pubmed/17107491>].
  47. ColtMan W, G PJ, A SJ, ton JM P: **Parasite-mediated selection against inbred Soay sheep in a free-living, island population.** *Evolution* 1999, **81**:1259–1267, [<http://www.jstor.org/stable/2640828>].
  48. Morgan M, Pagès H: *Rsamtools: Import aligned BAM file format sequences into R / Bioconductor* [<http://bioconductor.org/packages/release/bioc/html/Rsamtools.html>]. [R package version 1.4.3].

49. Audic S, Claverie JM: **The significance of digital gene expression profiles.** *Genome Res.* 1997, **7**:986–995, [<http://www.ncbi.nlm.nih.gov/pubmed/9331369>].
50. Romualdi C, Bortoluzzi S, D’Alessi F, Danieli GA: **IDEG6: a web tool for detection of differentially expressed genes in multiple tag sampling experiments.** *Physiol. Genomics* 2003, **12**:159–162, [<http://www.ncbi.nlm.nih.gov/pubmed/12429865>].
51. Pages H, Carlson M, Falcon S, Li N: *AnnotationDbi: Annotation Database Interface.* [R package version 1.16.10].
52. Alexa A, Rahnenfuhrer J: *topGO: topGO: Enrichment analysis for Gene Ontology* 2010. [R package version 2.6.0].
53. R Development Core Team: *R: A language and environment for statistical computing.* R Foundation for Statistical Computing, Vienna, Austria 2009, [<http://www.R-project.org>].
54. Leisch F: **Sweave: Dynamic Generation of Statistical Reports Using Literate Data Analysis.** In *Compstat 2002 — Proceedings in Computational Statistics.* Edited by Härdle W, Rönz B, Physica Verlag, Heidelberg 2002:575–580, [<http://www.stat.uni-muenchen.de/~leisch/Sweave>]. [ISBN 3-7908-1517-9].
55. Falcon S: **Caching code chunks in dynamic documents.** *Computational Statistics* 2009, **24**(2):255–261, [<http://www.springerlink.com/content/55411257n1473414>].
56. Wickham H: *ggplot2: elegant graphics for data analysis.* Springer New York 2009, [<http://had.co.nz/ggplot2/book>].
57. Chen H, Boutros PC: **VennDiagram: a package for the generation of highly-customizable Venn and Euler diagrams in R.** *BMC Bioinformatics* 2011, **12**:35, [<http://www.ncbi.nlm.nih.gov/pubmed/21269502>].
58. Balzer S, Malde K, Jonassen I: **Systematic exploration of error sources in pyrosequencing flowgram data.** *Bioinformatics* 2011, **27**:i304–309, [<http://www.ncbi.nlm.nih.gov/pubmed/21685085>].
59. Schwartz TS, Tae H, Yang Y, Mockaitis K, Van Hemert JL, Proulx SR, Choi JH, Bronikowski AM: **A garter snake transcriptome: pyrosequencing, de novo assembly, and sex-specific differences.** *BMC Genomics* 2010, **11**:694, [<http://www.ncbi.nlm.nih.gov/pubmed/21138572>].
60. Hale MC, Jackson JR, Dewoody JA: **Discovery and evaluation of candidate sex-determining genes and xenobiotics in the gonads of lake sturgeon (Acipenser fulvescens).** *Genetica* 2010, **138**:745–756, [<http://www.ncbi.nlm.nih.gov/pubmed/20386959>].
61. Papanicolaou A, Stierli R, Ffrench-Constant RH, Heckel DG: **Next generation transcriptomes for next generation genomes using est2assembly.** *BMC Bioinformatics* 2009, **10**:447, [<http://www.ncbi.nlm.nih.gov/pubmed/20034392>].
62. Emmersen J, Rudd S, Mewes HW, Tetko IV: **Separation of sequences from host-pathogen interface using triplet nucleotide frequencies.** *Fungal Genet. Biol.* 2007, **44**:231–241, [<http://dx.doi.org/10.1016/j.fgb.2006.11.010>].
63. O’Neil ST, Dzurisin JD, Carmichael RD, Lobo NF, Emrich SJ, Hellmann JJ: **Population-level transcriptome sequencing of nonmodel organisms *Erynnis propertius* and *Papilio zelicaon*.** *BMC Genomics* 2010, **11**:310, [<http://www.ncbi.nlm.nih.gov/pubmed/20478048>].
64. Gregory R, Darby AC, Irving H, Coulibaly MB, Hughes M, Koekemoer LL, Coetzee M, Ranson H, Hemingway J, Hall N, Wondji CS: **A De Novo Expression Profiling of *Anopheles funestus*, Malaria Vector in Africa, Using 454 Pyrosequencing.** *PLoS ONE* 2011, **6**:e17418, [<http://www.ncbi.nlm.nih.gov/pubmed/21364769>].
65. Kunstner A, Wolf JB, Backstrom N, Whitney O, Balakrishnan CN, Day L, Edwards SV, Janes DE, Schlinger BA, Wilson RK, Jarvis ED, Warren WC, Ellegren H: **Comparative genomics based on massive parallel transcriptome sequencing reveals patterns of substitution and selection across 10 bird species.** *Mol. Ecol.* 2010, **19 Suppl 1**:266–276, [<http://www.ncbi.nlm.nih.gov/pubmed/20331785>].
66. Parkinson J, Anthony A, Wasmuth J, Schmid R, Hedley A, Blaxter M: **PartiGene—constructing partial genomes.** *Bioinformatics* 2004, **20**(9):1398–1404, [<http://bioinformatics.oxfordjournals.org/cgi/content/abstract/20/9/1398>].

67. Wasmuth J, Schmid R, Hedley A, Blaxter M: **On the Extent and Origins of Genic Novelty in the Phylum Nematoda.** *PLoS Neglected Tropical Diseases* 2008, **2**(7):e258, [<http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2432500>].
68. Marcus Y, Parkinson J, Fernandez C, Daub J, Selkirk M, Blaxter M, Maizels R: **Signal sequence analysis of expressed sequence tags from the nematode *Nippostrongylus brasiliensis* and the evolution of secreted proteins in parasites.** *Genome Biology* 2004, **5**(6):R39, [<http://genomebiology.com/2004/5/6/R39>].
69. Yang H, Chen X, Wong WH: **Completely phased genome sequencing through chromosome sorting.** *Proc. Natl. Acad. Sci. U.S.A.* 2011, **108**:12–17, [<http://www.ncbi.nlm.nih.gov/pubmed/21169219>].
70. Adey A, Morrison H, Asan X, Xun X, Kitzman J, Turner E, Stackhouse B, MacKenzie A, Caruccio N, Zhang X, Shendure J: **Rapid, low-input, low-bias construction of shotgun fragment libraries by high-density in vitro transposition.** *Genome Biol.* 2010, **11**(12):R119, [<http://genomebiology.com/content/11/12/R119>].
71. Kryazhimskiy S, Plotkin JB: **The population genetics of dN/dS.** *PLoS Genet.* 2008, **4**:e1000304, [<http://www.ncbi.nlm.nih.gov/pubmed/19081788>].
72. Novaes E, Drost DR, Farmerie WG, Pappas GJ, Grattapaglia D, Sederoff RR, Kirst M: **High-throughput gene and SNP discovery in *Eucalyptus grandis*, an uncharacterized genome.** *BMC Genomics* 2008, **9**:312, [<http://www.ncbi.nlm.nih.gov/pubmed/18590545>].
73. Swanson WJ, Wong A, Wolfner MF, Aquadro CF: **Evolutionary expressed sequence tag analysis of *Drosophila* female reproductive tracts identifies genes subjected to positive selection.** *Genetics* 2004, **168**:1457–1465, [<http://www.ncbi.nlm.nih.gov/pubmed/15579698>].
74. Swanson WJ, Clark AG, Waldrip-Dail HM, Wolfner MF, Aquadro CF: **Evolutionary EST analysis identifies rapidly evolving male reproductive proteins in *Drosophila*.** *Proc. Natl. Acad. Sci. U.S.A.* 2001, **98**:7375–7379, [<http://www.ncbi.nlm.nih.gov/pubmed/11404480>].
75. Wang Z, Abubucker S, Martin J, Wilson RK, Hawdon J, Mitreva M: **Characterizing *Ancylostoma caninum* transcriptome and exploring nematode parasitic adaptation.** *BMC Genomics* 2010, **11**:307, [<http://www.ncbi.nlm.nih.gov/pubmed/20470405>].
76. Miyata T, Yasunaga T: **Molecular evolution of mRNA: a method for estimating evolutionary rates of synonymous and amino acid substitutions from homologous nucleotide sequences and its application.** *J. Mol. Evol.* 1980, **16**:23–36, [<http://www.ncbi.nlm.nih.gov/pubmed/6449605>].
77. Knopf K, Madriles Helm A, Lucius R, Bleiss W, Taraschewski H: **Migratory response of European eel (*Anguilla anguilla*) phagocytes to the eel swimbladder nematode *Anguillicola crassus*.** *Parasitology Research* 2008, **102**(6):1311–6, [<http://www.ncbi.nlm.nih.gov/pubmed/18311570>].
78. Molnár K: **Formation of parasitic nodules in the swimbladder and intestinal walls of the eel *Anguilla anguilla* due to infections with larval stages of *Anguillicola crassus*.** *Diseases of Aquatic Organisms* 1994, **20**(3):163–170.
79. Danecek P and others: **The variant call format and VCFtools.** *Bioinformatics* 2011, **27**:2156–2158, [<http://www.ncbi.nlm.nih.gov/pubmed/21653522>].
80. Wielgoss S, Taraschewski H, Meyer A, Wirth T: **Population structure of the parasitic nematode *Anguillicola crassus*, an invader of declining North Atlantic eel stocks.** *Molecular Ecology* 2008, **17**(15):3478–95, [<http://www.ncbi.nlm.nih.gov/pubmed/18727770>].
81. Anders S, Huber W: **Differential expression analysis for sequence count data.** *Genome Biol.* 2010, **11**:R106, [<http://www.ncbi.nlm.nih.gov/pubmed/20979621>].
82. Robinson MD, McCarthy DJ, Smyth GK: **edgeR: a Bioconductor package for differential expression analysis of digital gene expression data.** *Bioinformatics* 2010, **26**:139–140, [<http://www.ncbi.nlm.nih.gov/pubmed/19910308>].
83. Scott AL: **Nematode sperm.** *Parasitol. Today (Regul. Ed.)* 1996, **12**:425–430, [<http://www.ncbi.nlm.nih.gov/pubmed/15275275>].
84. Johnstone IL: **Cuticle collagen genes. Expression in *Caenorhabditis elegans*.** *Trends Genet.* 2000, **16**:21–27, [<http://www.ncbi.nlm.nih.gov/pubmed/10637627>].
85. Middleton B: **The oxoacyl-coenzyme A thiolases of animal tissues.** *Biochem. J.* 1973, **132**:717–730, [<http://www.ncbi.nlm.nih.gov/pubmed/4721607>].

86. Cutter AD, Ward S: **Sexual and temporal dynamics of molecular evolution in *C. elegans* development.** *Mol. Biol. Evol.* 2005, **22**:178–188, [<http://www.ncbi.nlm.nih.gov/pubmed/15371532>].
87. Eberhard WG: **Evolutionary conflicts of interest: are female sexual decisions different?** *Am. Nat.* 2005, **165 Suppl 5**:19–25, [<http://www.ncbi.nlm.nih.gov/pubmed/15795858>].



## Figures

### Figure 1 - Number of contigs annotated with different methods

Number of annotations obtained for Gene Ontology (GO), Enzyme Commission (EC) and Kyoto Encyclopedia of Genes and Genomes (KEGG) terms through Annot8r [34] for all TUGs (a) and for highCA derived contigs (b). The latter includes additional domain-based annotations obtained with InterProScan [36].

### Figure 2 - Comparing high level GO-slim annotations

For Gene Ontology (GO) categories molecular function, cellular compartment and biological process the number of terms in high level GO-slim categories is given as obtained through Annot8r [34].

### Figure 3 - Enrichment of Signal-positives for categories of evolutionary conservations

Proportions of SignalP-predictions for each category of evolutionary conservation. Generally - across bit-score thresholds - TUGS novel in nematodes and in *A. crassus* have the highest proportion of signal-positives.

### Figure 4 - Changes in ti/tv and dn/ds due to exclusion of homopolymer-runs

When SNPs in or adjacent to homopolymeric regions are removed changes in ti/tv and dn/ds are observed: as the overall number of SNPs is reduced both ratios change to more plausible values. Note the reversed axis for dn/ds to plot these lower values to the right. For homopolymer length  $> 3$  a linear trend for the total number of SNPs and the two measurements is observed. A width of 11 for the screening window provides most plausible values (suggesting specificity) while still incorporating a high number of SNPs (sensitivity).

### Figure 5 - SNP calling and SNP categories

Overabundance of SNPs at codon-position two (a) and of non-synonymous SNPs (c) for low percentages of the minority allele. (b) Significant positive correlation of coverage and dn/ds before removing these SNPs at a threshold of 7% ( $p < 0.001$ ,  $R^2 = 0.015$ ) and (d) absence of such a correlation afterwards ( $R^2 < 0.001$ ,  $p = 0.192$ ).

### **Figure 6 - Positive selection and evolutionary conservation**

Box-plots for  $dn/ds$  in TUGs according to different categories of evolutionary conservation. Significant comparisons are novel.in.metazoa - novel.in.Ac (0.005 and 0.015), novel.in.nematoda - novel.in.Ac (0.005 and 0.002), novel.in.nematoda - novel.in.clade3 (0.207 and 0.045; p-value for bitscore of 50 and 80, Nemenyi-Damico-Wolfe-Dunn test).

## Tables

**Table 1 - Sampling, trimming and pre-assembly screening, library statistics**

For libraries two sequencing libraries from European eels (E1 and E2) one from L2-larvae (L2), one from male (M) and two from Eels in Taiwan (T1 and T2) the following statistics are given. life.st = lifecycle stage: f for female m for male. source.p = source population: R for Rhine, P for Poland, C for cultured, W for wild. raw.reads = raw number of sequencing reads obtained. lowqal = number of reads discarded due to low quality or length in *Seqclean* [25]. AcrRNA = number of reads hitting *A. crassus*-rRNA (screened). eelmRNA = number of reads hitting eel transcriptome-sequences (screened). eelrRNA = number of reads hitting eel-rRNA genes (screened). Cercozoa = number of reads hitting cercozoan rRNA (screened). valid = number of reads valid after screening (assembled). valid.span = number of bases valid (assembled). mapping.unique = number of reads mapping uniquely to the assembly. mapping.Ac = number of reads mapping to the part of the assembly considered *A. crassus* origin (see post-assembly screening). mapping.MN = number of reads mapping to the highCA-derived part of the assembly (and also *A. crassus* origin). over.32 = number of reads mapping to contigs with overall coverage of more than 32 reads (considered in gene-expression analysis).

library	E1	E2	L2	M	T1	T2
life.st	adult f	adult f	L2 larvae	adult m	adult f	adult f
source.p	Europe R	Europe P	Europe R	Asia C	Asia C	Asia W
raw.reads	209325	111746	112718	106726	99482	116366
lowqal	92744	10903	15653	15484	7947	27683
AcrRNA	76403	11213	30654	31351	24929	7233
eelmRNA	4835	3613	1220	1187	7475	11741
eelrRNA	13112	69	1603	418	514	38
Cercozoa	0	0	5286	0	0	0
valid	22231	85948	58302	58286	58617	69671
valid.span	7167338	24046225	16661548	17424408	14443123	20749177
mapping.unique	12023	65398	39690	36782	42529	55966
mapping.Ac	8359	61073	12917	31673	37306	50445
mapping.MN	5883	48009	8475	18998	28970	41963
over.32	3595	34115	1602	10543	21413	22909

**Table 2 - Assembly classification and contig statistics**

Summary statistics for contigs from different assembly-categories given in columns as highCA = high credibility assembly; lowCA = low credibility assembly, combined = complete assembly.

Rows indicate summary statistics: total.contigs = numbers of total contigs, fish.contigs = number of contigs hitting eel-mRNA or Chordata in NCBI-nr or NCBI-nt (screened out), xeno.contigs = number of contigs with best hit (NCBI-nr and NCBI-nt) to non-eukaryote (screened out), remaining.contigs = number

of contigs remaining after this screening, remaining.span = total length of remaining contigs, non.u.cov = non-unique mean base coverage of contigs, cov = unique mean base coverage of contigs, p4e.“X” = number protein predictions derived in p4e, where “X” describes the method of prediction (see Methods), full.3p = number of contigs complete at 3’, full.5p = number of contigs complete at 5’, GO = number of contigs with GO-annotation, KEGG = number of contigs with KEGG-annotation, EC = number of contigs with EC-annotation, nem.blast = number of contigs with BLAST-hit to nematode in nr, any.blast = number of contigs with BLAST-hit to nematode or non-nematode (eukaryote non chordate) sequence in NCBI-nr.

	lowCA	highCA	combined
total.contigs	26336	13851	40187
rRNA.contigs	829	59	888
fish.contigs	2419	1022	3441
xeno.contigs	1935	1398	3333
remaining.contigs	21153	11372	32525
remaining.span	6157974	6575121	12733095
non.u.cov	14.665	10.979	12.840
cov	2.443	6.838	4.624
p4e.BLAST-similarity	4357	5664	10021
p4e.ESTScan	8324	3597	11921
p4e.LongestORF	8352	2085	10437
p4e.no-prediction	93	14	107
full.3p	5909	2714	8623
full.5p	1484	1270	2754
full.l	104	185	289
GO	2636	3875	6511
EC	967	1493	2460
KEGG	1609	2237	3846
IPR	0	7557	7557
nem.blast	4869	5821	10690
any.blast	5107	6008	11115

**Table 3 - Evolutionary conservation and novelty**

The kingdom Metazoa (novel.in.metazoa), the phylum Nematoda(novel.in.nematoda) and clade III (Spirurina; novel.in.spirurina) were assessed for occurrences of BLAST-hits at two different bitscore thresholds (50 = bit.50 and 80 = bit.80). TUGs without any hit at a given threshold were categorized as novel in *A. crassus* (novel.in.Ac). Both novelty and conservation can be derived from this (numbers for conservation would be the cumulative sum of lower-level novelty).

	conserved	novel.in.metazoa	novel.in.nematoda	novel.in.clade3	novel.in.Ac
bit.50.all	5604	1715	2173	1485	21548
bit.80.all	3506	1383	2015	1525	24096
bit.50.highCA	3479	876	1010	601	5406
bit.80.highCA	2457	833	1084	716	6282

**Table 4 - Over-representation of GO-terms in positively selected**

Significantly ( $p < 0.05$ ) over-represented GO-terms in contigs putatively under positive selection. Horizontal lines separate categories of the GO-ontology. First category is molecular function, second biological process, last cellular compartment. P-values (p.value) for over-representation are given along with the number of positively selected contigs (Significant;  $dn/ds > 0.5$ ) and the number of contigs with this annotation for which a  $dn/ds$  was obtained (Annotated) and the description of the GO-term (Term). For a graph of induced GO-terms see also Additional Figures 1.

GO.ID	Term	Annotated	Significant	Expected	p.value
GO:0008233	peptidase activity	43	13	6.08	0.0034
GO:0015179	L-amino acid transmembrane transporter activity	2	2	0.28	0.0198
GO:0043021	ribonucleoprotein binding	6	3	0.85	0.0396
GO:0042277	peptide binding	10	4	1.41	0.0397
GO:0070011	peptidase activity, acting on L-amino acid peptides	35	9	4.95	0.0442
GO:0004175	endopeptidase activity	25	7	3.54	0.0488
GO:0042594	response to starvation	15	7	2.13	0.0022
GO:0009083	branched chain family amino acid catabolic process	3	3	0.43	0.0027
GO:0006914	autophagy	12	6	1.70	0.0031
GO:0009063	cellular amino acid catabolic process	10	5	1.42	0.0071
GO:0009267	cellular response to starvation	7	4	0.99	0.0093
GO:0006520	cellular amino acid metabolic process	44	12	6.24	0.0128
GO:0006915	apoptosis	78	18	11.06	0.0147
GO:0009308	amine metabolic process	57	14	8.08	0.0189
GO:0005997	xylulose metabolic process	2	2	0.28	0.0199
GO:0006739	NADP metabolic process	2	2	0.28	0.0199
GO:0007616	long-term memory	2	2	0.28	0.0199
GO:0009744	response to sucrose stimulus	2	2	0.28	0.0199
GO:0010172	embryonic body morphogenesis	2	2	0.28	0.0199
GO:0015807	L-amino acid transport	2	2	0.28	0.0199
GO:0050885	neuromuscular process controlling balance	2	2	0.28	0.0199
GO:0007281	germ cell development	17	6	2.41	0.0226
GO:0090068	positive regulation of cell cycle process	17	6	2.41	0.0226

GO:0042981	regulation of apoptosis	64	15	9.07	0.0232
GO:0051329	interphase of mitotic cell cycle	23	7	3.26	0.0320
GO:0044106	cellular amine metabolic process	55	13	7.80	0.0325
GO:0031571	mitotic cell cycle G1/S transition DNA damage checkpoint	14	5	1.98	0.0355
GO:0010564	regulation of cell cycle process	34	9	4.82	0.0377
GO:0006401	RNA catabolic process	6	3	0.85	0.0398
GO:0010638	positive regulation of organelle organization	6	3	0.85	0.0398
GO:0009056	catabolic process	149	28	21.12	0.0398
GO:0008219	cell death	93	19	13.18	0.0441
GO:0007154	cell communication	144	27	20.41	0.0455
GO:0051726	regulation of cell cycle	52	12	7.37	0.0474
GO:0030330	DNA damage response, signal transduction by p53 class mediator	15	5	2.13	0.0475
GO:0033238	regulation of cellular amine metabolic process	15	5	2.13	0.0475
GO:0030532	small nuclear ribonucleoprotein complex	7	4	0.99	0.0093
GO:0005739	mitochondrion	137	28	19.38	0.0113
GO:0005682	U5 snRNP	2	2	0.28	0.0198
GO:0015030	Cajal body	2	2	0.28	0.0198
GO:0046540	U4/U6 x U5 tri-snRNP complex	2	2	0.28	0.0198
GO:0016607	nuclear speck	6	3	0.85	0.0396

**Table 5 - Measurements of multi-locus heterozygosity for single worms**

Genotyping for a set of 199 SNPs, different measurements were obtained to assess genome-wide heterozygosity. Measurements for relative heterozygosity (rel.het; number of homozygous sites/ number of heterozygous sites), internal relatedness (int.rel; [45]), homozygosity by loci (ho.loci; [46]) and standardized heterozygosity (std.het; [47]) are given with the number of SNPs informative for this library (inform.snp). All these measurements are pointing to sample T1 (Taiwanese worm from aquaculture) as the most heterozygous and sample E2 (the European worm from Poland) as the least heterozygous individual. Heterozygote-heterozygote correlation [44] confirmed the genome-wide significance of these markers.

	rel.het	int.rel	ho.loci	std.het	inform.snps
T2	0.45	-0.73	0.59	1.00	121.00
T1	0.93	-0.95	0.34	1.62	136.00
M	0.37	-0.73	0.66	0.84	92.00
E1	0.38	-0.83	0.60	0.91	65.00
E2	0.18	-0.35	0.82	0.50	140.00

**Table 6 - Over-representation of GO-terms differentially expressed**

Significantly ( $p < 0.05$ ) over-represented GO-terms in contigs differentially expressed between male and female worms (a) or between European and Asian origin (b). Horizontal lines separate categories of the GO-ontology. First category is molecular function, second biological process, last cellular compartment. P-values (p.value) for over-representation are given along with the number of differentially expressed contigs (Significant) and the number of contigs with this annotation analysed (Annotated) and the description of the GO-term (Term). For a graph of induced GO-terms see also Additional Figures 1.

a)

b)

## **Additional Files**

### *Additional text*

The additional text describes the assembly process and evaluation of assembly quality in further detail. This text also contains figures and tables.

### *Additional tables*

Additional table 1 a lists all data computed on the contig level, including sequences (raw, coding, imputed, protein) additional table 1 b lists only the metadata not including sequences.

Additional table 2 lists high quality SNPs.

Additional tables 3 list contigs differentially expressed between male and female worms (a) and European and Asian worms (b). Normalised counts and the logarithm of fold changes are given.

### *Additional figures*

Additional Figures 1: subgraphs of the GO-ontology categories induced by the top 10 terms identified as enriched in different sets of genes. Boxes indicate the 10 most significant terms. Box colour represents the relative significance, ranging from dark red (most significant) to light yellow (least significant). In each node the category-identifier, a (eventually truncated) description of the term, the significance for enrichment and the number of DE / total number of annotated gene is given. Black arrows indicate a is “is-a” relationship. GO-ontology category and the set of genes analysed for the enrichment are indicated in each figure.