# Guidelines for RNA-seq projects: applications and opportunities in non-model decapod crustacean species

Tuan Viet Nguyen[*†a], Hyungtaek Jung[†b], Guiomar Rotllant[c], David Hurwood[d], Peter Mather[d], Tomer Ventura[a]

[a] GeneEcology Research Centre, Faculty of Science, Health, Education and Engineering, University of the Sunshine Coast, 4 Locked Bag, Maroochydore, QLD 4558, Australia

[b] Centre for Tropical Crops and Biocommodities, Institute for Future Environment, Queensland University of Technology, GPO Box 2434, Brisbane QLD 4001, Australia

[c] Institut de Ciències del Mar (CSIC), Passaig Marítim de la Barceloneta 37, 08003 Barcelona, Spain

[d] School of Earth, Environment and Biological Sciences, Science and Engineering Faculty, Queensland University of Technology, GPO Box 2434, Brisbane QLD 4001, Australia

[*]: Corresponding author: Tuan Viet Nguyen (tnguyen@usc.edu.au)

[†]: Equal contribution

32

## Abstract (150-200 words)

34 Next Generation Sequencing (NGS) has dramatically changed the way biological research is being

35 conducted in the post-genomic era and they have only been utilized widely over the recent decade

36 for studies of non-model decapod crustacean species, predominantly by sequencing the

37 transcriptome of various tissues across different life stages. NGS can now provide a rapid, cost-

38 effective solution for discovery of genetic markers crucial in many applications that would previously

39 have otherwise taken years to develop. Sequencing of the entire transcriptome (referred to as RNA

40 sequencing; RNA-seq) is one of the most popular NGS tools. RNA-seq studies of non-model species in

41 crustacean taxa however, have faced some problems, including a lack of "good" experimental study

42 design, a relative paucity of gene annotations, combined with limited knowledge of genomic

43 technologies and analyses. The aim of the current review is to assist crustacean biologists to develop

44 a better appreciation of the applications and scope of RNA-seq analysis, understand the basic

45 requirements for optimal RNA-seq studies and provide an overview of each step from RNA-seq

46 experimental design to bioinformatics approaches to data analysis. Insights that have resulted from

47 RNA-seq studies across a wide range of non-model decapod species are also summarized.

48

# INTRODUCTION

Next generation sequencing (NGS) technologies have rapidly transitioned bioscience into the post-genomic era, resulting in easier, cheaper, and faster DNA sequencing. Application of advanced NGS platforms has allowed multiple techniques to be developed that address biological challenges. These include; RNA-sequencing (RNA-seq), whole-exome sequencing, chromatin immunoprecipitation sequencing (ChIP-seq), microRNA sequencing (miRNA-seq), restriction assisted DNA sequencing (RAD-seq), and small RNA sequencing. Among these, RNA-seq is a technique that has revolutionized gene expression studies and marker discovery (single sequence repeats [SSRs]/microsatellites and single nucleotide polymorphisms [SNPs]) (Das et al. 2016; Lister et al. 2009; Marguerat and Bähler 2010; Mykles et al. 2016; Ozsolak and Milos 2011; Wang et al. 2009; Wilhelm and Landry 2009). The RNA-seq platform is based on the analysis of the transcriptome - a small portion of the whole genome that is transcribed from chromosomal DNA into RNA molecules - a dynamic set of elements that change depending on developmental stages or physiological conditions. Also, by analysing the sequenced transcriptome, genetic polymorphisms including SNPs and SSRs can be mined and analysed with ease (Jaramillo et al. 2016; Jin et al. 2013; Jung et al. 2011; Jung et al. 2016; Lv et al. 2014; Meng et al. 2015; Nguyen et al. 2016).

While RNA-seq techniques have had a major impact on model species (which in this review is defined as a species with a well-characterized genome, e.g. *Daphnia pulex*), the application of RNA-seq approaches in non-model decapod crustacean taxa is still limited by the small size of the research community and the subsequent bottleneck of bioinformatics analysis capabilities. Many NGS analytical tools are available and by default, are developed for model species, making it difficult for researchers investigating non-model organisms to navigate through and identify appropriate tools. Designing and evaluating a pipeline for transcriptomics projects in non-model species therefore, can be considered a crucial step prior to project initiation. While the transcriptome can encompass many categories of different types of RNA (mi-RNA, small nuclear RNA, non-coding RNA etc.), this review will

76  focus mainly on mRNA sequencing using Second Generation Sequencing (SGS) technology – we intend

77  to use the same classification proposed by Schadt et al. (2010), that defined Sanger sequencing as First

78  Generation, "wash-and-scan" sequencing technology as Second Generation, and single molecule real

79  time sequencing as Third Generation. Under this classification scheme, SGS includes a number of

80  platforms, notably Illumina, Solid, Ion Torrent/Ion Proton, Roche 454; whereas PacBio and Oxford

81  Nanopore are classified as Third Generation Sequencing (TGS) Technology. Here we will focus

82  primarily on different strategies to initiate a transcriptome study, briefly addressing several platforms

83  that currently are available, as well as recommending a number of experimental designs,

84  bioinformatics software for *de novo* assembly and specific data analyses for decapod crustacean

85  species. Finally, we review recent biological insights gained from application of SGS in crustacean

86  transcriptomics and highlight opportunities as well as challenges for applied RNA-seq in the future.


87  ## OVERVIEW OF RNA-SEQ TECHNOLOGY

---

88  ### PRE-SEQUENCING

89  New sequencing technologies and new sequencing chemistries are being developed rapidly. The

90  arrival of SGS, and more recently TGS, has completely changed the way researchers approach

91  unanswered phenomena in basic, applied, and clinical research. Each sequencing platform is based on

92  different proprietary chemistries and technologies and each has unique strengths and weaknesses.

93  Details on sequencing chemistry have been summarized elsewhere (Goodwin et al. 2016; Koboldt et

94  al. 2013; Metzker 2010; Reuter et al. 2015). Currently, Illumina is the most widely utilized SGS for RNA-

95  seq, since the platform enables deep coverage of the transcriptome and provides long, low-error reads

96  that are suitable for mapping to reference genomes and transcriptome assemblies (Goodwin et al.

97  2016; Metzker 2010; Niedringhaus et al. 2011). Performance benchmarking of many SGS platforms

98  has been conducted for several years (Finseth and Harrison 2014; Glenn 2011; Goodwin et al. 2016;

99  Lahens et al. 2017; Lam et al. 2012; Liu et al. 2012) and an online archive of sequencing platforms is

100      available on the market and can be found at https://allseq.com/knowledgebank/. Given the popularity

101      of Illumina Sequencers in general, it tends to be the technology most widely applied in crustacean

102      transcriptome projects (Havird and Santos 2016).

103      In brief, RNA-seq includes the use of an SGS platform to generate a huge amount of sequence data.

104      Due to technical constraints of the approach (most SGS platforms can only generate short to medium

105      length reads, approximately 50-300 bp), RNA transcripts must be fragmented into shorter sequences.

106      In the absence of a reference genome, short reads are then reconstructed to make a reference

107      transcriptome, referred to as a *de novo* assembly. Following this, raw reads can be realigned (or

108      mapped) to the previously generated reference sequence and counted, thus providing a digital

109      measurement of specific transcript abundances that can facilitate biological interpretation. Where key

110      genes are targeted (based on either high differential expression or previously identified in the

111      literature), they can be validated by replicating samples across a range of experimental conditions (eg.

112      in different tissue types, at different life history stages, between sexes, etc.). A popular approach for

113      validation includes quantitative real-time PCR (qRT-PCR) where relative transcript abundance can be

114      assessed under more strictly controlled conditions. Most RNA-seq strategies that utilize SGS can be

115      summarized by a basic workflow (Figure 1).

116 *Experimental design*

117 Designing an RNA-seq experiment requires a solid biological understanding of the taxa under

118 investigation and the question(s) to be addressed. Poor or inappropriate decisions at this stage can

119 result in a large amount of unusable data. A good experimental design for every NGS–based

120 experiment therefore, is a basic requirement that cannot be over-emphasized.

121 In general, several factors must be considered prior to the initiation of any well-designed sequencing

122 project. Essentially, an appropriate experimental design is a balance between the level of biological

123 versus technical replication (Figure 2) and the resulting depth of coverage for each tissue type, life

124 stage, sex etc., within a framework of time and financial constraints. It is advisable that researchers

125 without much prior experience should seek suggestions from professional service providers including

126 bioinformaticians and biostatisticians, as well as the sequencing provider. This review highlights some

127 of the pitfalls to be aware of and sets the scene for appropriate study design. Several studies provide

128 direction on how to design a statistically valid RNA-seq experiment (Auer and Doerge 2010; Conesa et

129 al. 2016; Fang and Cui 2011; Yang and Wei 2015). In general, a comprehensive transcriptome requires

130 multiple tissues from multiple developmental stages while gene expression studies require samples

131 that represent contrasting treatments (e.g. male vs female, control vs hormone treated, salinity vs

132 freshwater acclimation, or different developmental/life history stages).

133 *Biological and technical replicates*

134 In the NGS context, technical replication refers to multiple libraries from the same biological sample

135 (i.e. the technical steps are performed separately) (Figure 2). While potentially increasing the depth

136 of reads, any variation recorded among technical replicates will also help identify inconsistencies

137 associated with sampling techniques, PCR biases or sequencing errors. In some rare cases, sample

138 collection, storage or processing can be a source of technical variance owing to the relative instability

139 of RNA. It is advisable to employ several randomization techniques during sequencing, for example,

140 multiplexing (mixing of different libraries, each tagged using a different barcode), splitting technical

141    repeats between multiple lanes, or randomization of different libraries in the same lane (an excellent

142    review on statistical randomization for RNA-seq can be found elsewhere (Auer and Doerge 2010)).

143    Ultimately, this type of replication provides some measure of the quality and/or reliability of the

144    analysis.

145    Biological replication alternatively, relates to different biological samples (e.g. same tissue type but

146    from different individuals) that are processed separately (Figure 2). Biological replication is desirable

147    since it quantifies natural variation among individuals within the experimental cohort. Furthermore,

148    increasing the sample size (number of biological replicates) not only increases sequencing depth, but

149    also provides greater statistical power to detect differences among treatments where they may exist.

150    Nevertheless, with a very large sample size, accommodating both technical and biological variation

151    can become very costly and may also result in a complex assay to analyse. When sequencing

152    individuals from a population with large levels of genetic variation, for example when dealing with

153    wild-caught individuals, the more biological replicates, the more likely it is to capture genuine

154    differential expression among groups. In general, most SGS experiments conducted on crustaceans

155    tend to be under-replicated and while there is no gold standard for this matter, it is currently

156    acceptable for RNA-seq experiments to consist of a minimum of three biological samples to provide

157    adequate statistical power; a number of published studies have shown that the power to detect

158    differential expressed genes improves from two samples to five samples per treatment (Dillies et al.

159    2013; Kvam et al. 2012). Similarly, other studies have proposed that sequencing fewer reads and

160    including more biological replicates is an effective strategy to increase statistical power and accuracy

161    in large-scale differential expression RNA-seq studies (Liu et al. 2014). More recently, results suggest

162    that at least six biological replicates may be needed in more sophisticated RNA-seq experiments and

163    up to 12 replicates per experimental group (Schurch et al. 2016). However, for samples that are very

164    different from each other in terms of transcription level (for example, differential expression profiles

165    between brain versus ovary), less replication may also be acceptable. It is also important to highlight

166    the fact that replicates in an RNA-seq-based study are required for publication in some journals (e.g.,

167    refer to section 2.6.7 at https://www.frontiersin.org/about/author-guidelines). To conclude, we

168    would recommend maximizing biological replicates to include at least three samples for each

169    experimental condition in every non-model decapod crustacean RNA-seq study.

170    ***Choice of sequencing platforms***

171    There are several sequencing methods that researchers can choose from, including single-end (SE)

172    /paired-end (PE) reads, strand-specific, or non-strand-specific library preparation. The decision on

173    which is selected will be based on the desired outcome of the study but will also depend on budget

174    constraints. For experiments on crustacean species in general, PE sequencing is recommended to

175    obtain a reliable *de novo* assembly where no reference genome is readily available. Long read

176    sequencing (e.g. PacBio, Nanopore sequencing),  proven to be suitable for enhancing continuity of *de*

177    *novo* transcriptome assembly, is currently relatively expensive and its application has been described

178    elsewhere (Cartolano et al. 2016; Chen et al. 2017; Kuo et al. 2017). Illumina short read sequencing

179    however, is by far the most widely used platform for transcriptome sequencing in crustaceans due to

180    its cost-effectiveness (unit price per nucleotide), fast sequencing times and higher raw read accuracy.

181    Another consideration is to choose whether stranded sequencing will be needed. In brief, a stranded-

182    specific RNA-seq can retain the gene orientation (sense or antisense transcript). A number of studies

183    have attempted to compare between stranded vs non-stranded approaches and most have shown

184    that a stranded RNA-seq approach is more advantageous due to better assembly of unannotated

185    genes, ability to detect genes on the antisense strand as well as improved continuity of transcripts.

186    New *de novo* assembly programs like Trinity (Grabherr et al. 2011) have a special mode for strand-

187    specific data analysis that has proven to be more effective than non-stranded data (Levin et al. 2010;

188    Parkhomchuk et al. 2009; Sultan et al. 2012; Zhao et al. 2015; Zhong et al. 2011). We therefore

189    recommend strand-specific RNA-seq if possible for non-model decapod crustacean studies (Havird

190    and Santos 2016).

191 ***Depth of sequence (number of reads)***

192 The amount of sequencing needed for a given sample is determined by the aims of the experiment,

193 the number of transcribed transcripts and the nature of the species' RNA samples (this is due to the

194 fact that crustacean genomes can be quite complex compared to other invertebrates). To our

195 knowledge, there has been no attempt to investigate the depth required for effective RNA-seq studies

196 in crustaceans. A study of chicken RNA-seq data revealed that approximately 30 million reads

197 (Illumina-75 bp PE) covered all annotated genes, while 10 million reads detected only ~80% (Wang et

198 al. 2011). Whereas RNA-seq samples from six different phyla (Annelida, Arthropoda, Chordata,

199 Cnidaria, Ctenophora and Mollusca) has suggested that approximately 20 million reads for tissue

200 samples and 30 million for whole-animal samples were required to provide a good balance between

201 total coverage and noise (Francis et al. 2013). Based on these data, it is acceptable that 20 million PE

202 reads per sample for a diploid crustacean organism is a reasonable target to aim for, although there

203 is no specific benchmark for all sequencing experiments.

204 It is also important to note that in order to detect transcripts with low expression, a deeper sequencing

205 strategy may be needed. In the guidelines for the ENCODE project (https://www.encodeproject.org/),

206 an experiment to evaluate similarity between two transcriptional profiles, requires 30 million PE reads

207 that must be mapped to the genome or known transcriptome. Guidelines to detect novel elements or

208 quantification of known transcript isoforms requires deeper sequencing (Refer to the whole guideline

209 at https://www.encodeproject.org/about/experiment-guidelines/#guideline). Another tool, Scotty,

210 can be used to assist in the design phase of RNA-Seq experiments (Busby et al. 2013). This program

211 can confirm if the design applied has sufficient statistical power to detect differentially expressed

212 genes (DEGs) at the predetermined level required. The program is freely available online at

213 http://bioinformatics.bc.edu/marthlab/scotty/scotty.php/.

214 An interim conclusion to be drawn from the above sections is that there exists a trade-off between

215 depth of reads per sample and the number of samples (which include technical and biological repeats).

216    The technology employed and financial limitations usually dictate a fine balance between these

217    factors.

218    ***Tissues RNA extraction and cDNA library preparation***

219    Library preparation is a crucial step prior to sequencing. It consists of a number of stages including

220    RNA extraction, proper storage of RNA, quality checking of RNA, mRNA isolation and finally cDNA

221    library generation.

222    In brief, extraction of total RNA from target tissue can be undertaken immediately on-site or samples

223    can be stored in RNA-later® solution for later extraction. It is important to note that RNA is extremely

224    fragile and degrades readily if stored under inappropriate conditions. Additionally, ribonucleases

225    (RNases) which enzymatically degrade RNA pose a constant threat of contamination and degradation

226    of purified RNA. Traditionally, RNA can be stored at −20°C, −80°C (most desirable) or in liquid nitrogen

227    (-196°C) to provide protection. RNA storage solutions that include chelating agents which inhibit

228    RNase activity, can be used, although these might interfere with reverse transcription and should thus

229    be removed prior to these steps. To our knowledge, there is no crustacean specific RNA extraction kit

230    available on the market, however several commercial kits for RNA extraction are still usable for

231    crustaceans, in addition to in-house (modified) versions of RNA extraction methods that use Beta-

232    mercaptoethanol or phenol-based compounds, with the latter being more popular in recent

233    publications. A detailed review on the effect of RNA extraction methods on RNA-seq can be found

234    elsewhere (Sultan et al. 2014). RNA can then be assessed for quality and quantity using a Nanodrop®

235    spectrophotometer or BioAnalyzer®. It is important to note that RNA integrity number (RIN) that has

236    been used as a standardized metric of RNA quality for vertebrate species, is not usually valid for

237    crustacean samples with non-typical RNA profiles. RIN is calculated based on the ratio between 18S

238    ribosomal RNA (rRNA) and 28S rRNA band intensities, which are usually very conserved across

239    eukaryotes. However, the 28S rRNA of arthropods tends to break down into two subunits, preventing

240    a reliable RIN value calculation (Macharia et al. 2015; McCarthy et al. 2015; Winnebeck et al. 2010).

241   RNA can be stored and shipped in ambient conditions after desiccation with RNA-stable solution

242   (Seelenfreund et al. 2014). An important consideration when it comes to RNA extraction in crustacean

243   species is tissues with high pigment content (i.e. eyestalk). For these tissue, extra caution is suggested

244   to avoid extracting pigment contamination that will affect the quality of library preparation. Currently,

245   there is no threshold for deciding if a sample is too degraded for whole-transcriptome analysis. In

246   most cases however, sequencing facilities provide users with specific guidelines and technical notes

247   recommended for producing the best results. Moreover, depending on each sequencing platform,

248   different cDNA library preparation protocols may be required.


249      **POST-SEQUENCING**

250   Post-sequencing analyses include quality checking of raw sequences, trimming, *de novo* assembly of

251   trimmed reads, read mapping and quantification, DEG assessment and finally biological interpretation

252   (Figure 1).


253   **Quality control for SGS data**

254   Current SGS runs generate millions, or even hundreds of millions of read sequences. Technologies

255   advancement reduces the error rate; however, every platform still produces read errors that require

256   the application of a quality control program post-sequencing. Read errors, while relatively negligible

257   in number compared with the massive dataset generated, still pose a hurdle for downstream analysis.

258   For instance, errors in base-calling cause improper connection of nodes in *de novo* assembly (thus

259   expanding running time and increase memory needed to store the nodes). In addition, incorrect SNP

260   detection can result from an inability to differentiate between a true polymorphism and a sequencing

261   error (Kelley et al. 2010). Several quality control tools have been developed for NGS data (most

262   popular tools are summarized in Table 1).

263   In general, quality control of raw reads from NGS sequencers can be completed in a few simple steps.

264   Raw      read      statistics      can      then      be      checked      with      FASTQC      software

265  (https://www.bioinformatics.babraham.ac.uk/projects/fastqc/). A variety of parameters can be used

266  to trim the data. The most important is the PHRED quality score (a base-calling score ranking system

267  that allows users to judge the confidence of a nucleotide presumed to be correctly called (Ewing et al.

268  1998)). Some other considerations include reads average length, total number of base pairs and

269  adapters contamination. In addition, reads generated on the Illumina platform are considered to have

270  a relatively higher error rate towards the 3'-end of the read (Schirmer et al. 2015), so if a drop in

271  quality is detected, it is acceptable to trim off a portion of the read from that end. Some commonly

272  used criteria for trimming reads include; minimum read length, minimum quality score, and

273  homopolymer trimming. Read duplication is also a factor to consider during the quality control step

274  in NGS projects. In brief, read duplications are identical reads that map to the same genomic location

275  (effects of PCR amplification bias, excess computational resources, and errors). Raw reads may also

276  need to be cleaned from artificially introduced sequences - PCR primers or sequencing adapters; these

277  are usually addressed in most quality control packages. In a benchmarking study, it was shown that

278  trimming applied in every sequencing project will improve not only quality of the results, but also

279  reduce analysis duration (Del Fabbro et al. 2013). In general, normal quality trimming with a PHRED

280  score ranging from 20 to 30 is normal for most RNA-seq experiments, while a PHRED score threshold

281  of 30 or above is usually required for variant calling experiments (Ledergerber and Dessimoz 2011).

282  However, in one particular study, the authors highlighted that although strict trimming is usually

283  applied, in some cases a more gentle trimming (PHRED score <2 or <5) might be more optimal

284  (MacManes 2014). This is due to the fact that short and low expressed transcripts suffer from heavy

285  negative bias when using harsh trimming (MacManes 2014). Therefore, lowering the PHRED score

286  threshold in the quality control step can result in a greater transcript discovery rate. As a conclusion,

287  we suggest gentle trimming initially as suggested in the above study. A list of some popular software

288  packages for NGS quality control can be found in Table 1.

289

**_De novo_ assembly for non-model species and transcript clustering**

For non-model decapod species, it is often difficult to align RNA-seq data to a reference genome from relatively recently diverged organisms (currently there are very few crustacean reference genomes available – see section 3.2). An alternative strategy therefore, is to construct a _de novo_ assembly (new assembly) from high-quality reads. The primary aim is to extend the short reads from the sequencer into longer continuous sequences (contigs) that reflect the mRNAs transcribed in the cell without any chimeric/fusion events. A number of _de novo_ transcriptome assemblers have been developed (initially they were simply modified genome assemblers), including the Velvet/Oases pipeline (Schulz et al. 2012; Zerbino and Birney 2008), SOAPdenovo (https://soap.genomics.org.cn/soapdenovo.html) and Trans-Abyss (Robertson et al. 2010). More recently, the Trinity software (Grabherr et al. 2011) has become available, developed specifically for _de novo_ transcriptome assembly from short-read RNA-seq data. Since reads from SGS are short in length compared with pyrosequencing output (Liu et al. 2012), transcriptome _de novo_ assemblers often employ a De Bruijn graph algorithm instead of the traditional **O**verlap **L**ayout **C**onsensus (OLC). This minimizes the amount of memory required to handle numerous parallel calculations. Further information on graph algorithms can be found elsewhere (Li et al. 2012; Miller et al. 2010).

Most _de novo_ assemblers are freely distributed but usually required operating using command line, which deters many biologists without programming skills. To overcome this issue, bioinformatics platforms such as Galaxy (https://www.usegalaxy.org/) and CyVerse (https://www.cyverse.org/) embed command line packages into user-friendly interfaces. Yet, there is a limited flexibility in utilizing these tools. Learning how to use command line programming can be time consuming and potentially is out of reach for many non-model biology researchers and this can slow the pace at which NGS studies are performed on these species. To address this, users can use commercial products (usually with a "point-and click" user-friendly interface) that are available on the market. A summary of some

314    notable *de novo* assemblers can be found in Table 2. A performance comparison of commonly used

315    *de novo* transcriptome assemblers can be found elsewhere (Amin et al. 2014; Finseth and Harrison

316    2014; Ghangal et al. 2013; Surget-Groba and Montoya-Burgos 2010; Zhao et al. 2011).

317    One significant challenge associated with *de novo* assembly is the lack of software to identify the

318    assembly that is most accurate. To address this challenge, **S**equence **C**omparative **A**nalysis using

319    **N**etworks (SCAN) was created (Misner et al. 2013). SCAN uses a reference dataset (from a related

320    genome) to identify the most accurate *de novo* assembly and to classify "good" transcripts in these

321    assemblies (Misner et al. 2013). A similar program was generated for this purpose, named DETONATE

322    (an abbreviation of *DE novo* **T**ranscript**O**me r**N**a-seq **A**ssembly with or without the **T**ruth **E**valuation)

323    (Xie et al. 2014). This program combines multiple factors into a single evaluation score that then can

324    be    used    to    select    the    best    assembler.    The    software    is    distributed    freely    at

325    [https://deweylab.biostat.wisc.edu/detonate](https://deweylab.biostat.wisc.edu/detonate). Another approach is to employ the CEGMA pipeline

326    (**C**ore **E**ukaryotic **G**enes **M**apping **A**pproach) (Parra et al. 2007) or BUSCO (**B**enchmarking **U**niversal

327    **S**ingle-**C**opy **O**rthologs) (Simão et al. 2015). These programs scan the *de novo* assembly against a

328    dataset of core eukaryotic genes that are well conserved across several eukaryotic taxa, to calculate

329    the coverage of protein-coding genes, thus estimating the degree of completeness of the

330    reconstruction and the full-length complement of transcript sequences comprising the *de novo*

331    transcriptome assembly. As a concluding remark, benchmarking assemblies are an option that can be

332    trialled, but the practice is still in development.

333    **Transcriptome mapping**

334    Following *de novo* assembly, reads can be aligned against the *de novo* assembly (mapping). The

335    mapping step can serve two purposes: i) a remapping step can be used to assess the assembly quality

336    and ii) the alignment can then be quantified; gene expression levels can be inferred from the total

337    counts of reads aligned to each contig. Furthermore, mapping also enables variant calling for

338    transcripts of interest.

339 Stringent parameters may result in a small subset of reads mapped, while less stringent settings

340 reduce mapped read specificity. To gain a balance between sensitivity and specificity, trials with

341 different parameters can be performed. Popular aligners for RNA-seq include: Bowtie 1 (Langmead et

342 al. 2009) and Bowtie 2 (Langmead and Salzberg 2012), BWA (Li and Durbin 2009), GSnap (Wu and Nacu

343 2010), and commercial programs including CLC Genomics Workbench®, DNA-STAR® or Partek

344 Genomics®. A detailed list of available aligners can be found at

345 https://www.ebi.ac.uk/~nf/hts_mappers/ (Fonseca et al. 2012). Comparisons of different aligners

346 usually takes into consideration running time, accuracy, as well as the sensitivity of mapped reads

347 (Baruzzo et al. 2017; Grant et al. 2011; Hatem et al. 2013; Li and Homer 2010). Critically, for non-model

348 organisms where no genome sequence is available, it is hard to define which are the best mapping

349 parameters to apply. This is due to the occurrence of isoforms and splice variants that cannot be

350 accurately determined without access to a reference genome. Reads can be mapped randomly to

351 shared exons between splice variants, biasing the resulting count and confounding the biological

352 interpretation.

**Quantifying transcript level and analysis of differential gene expression**

354 To quantify gene expression, RNA-seq reads need to be aligned to a reference genome from model

355 organisms or to the transcriptome sequences reconstructed using *de novo* assembly strategies for

356 organisms without reference genome sequences. The number of mapped reads is calculated based on

357 the outcome of the alignment and can be used to estimate the relative expression level of individual

358 genes. Following this, statistical methods are applied to test for significant differences among

359 experimental groups. The data however, first needs to be normalized since there are inherent

360 differences in total reads per sample, resulting in over-represented long transcripts. With rapid

361 development of RNA-seq technology, there are now numerous tools available to estimate gene

362 expression levels, which vary in their efficiency. Popular RNA-seq quantification (reads counting) tools

363 include: RSEM (Li and Dewey 2011), eXpress (Roberts and Pachter 2012), HTSeq (Anders et al. 2015),

364 Salmon (Patro et al. 2017) and kallisto (Bray et al. 2016). Several studies have also been conducted to

365    compare the pros and cons of each tool (Chandramohan et al. 2013; Li and Homer 2010; Teng et al.

366    2016).

367    DEG analysis programs perform statistical tests to determine if fold change results under different

368    experimental conditions are significant (e.g. among tissues types, life stages etc.). Many programs

369    have been developed for DEG analysis (a brief summary of popular DEG tools can be found in Table 3)

370    and several comparative assessments are available (Khang and Lau 2015; Kvam et al. 2012; Rajkumar

371    et al. 2015; Robles et al. 2012; Soneson and Delorenzi 2013; Zhang et al. 2014). Much like assembly

372    and mapping, there is no guarantee as to which tool is the best, or which parameters will result in the

373    highest accuracy or robustness of the results generated (Zhang et al. 2014). Most DEG call methods

374    are designed to address analysis of RNA-seq experiments that have biological replicates. There are a

375    few tools however, that can handle non-replicated experiments (e.g. GFOLD (Feng et al. 2012), EdgeR

376    (Robinson et al. 2010), NOISeq (Tarazona et al. 2011). A recent study recommended using EdgeR

377    (Robinson et al. 2010) or DESeq2 (Love et al. 2014) for experiments with less than 12 replicates per

378    group, while they suggest studies with more than 12 replicates should use DESeq2 for the statistical

379    analysis (Schurch et al. 2016). An alternative strategy is to employ several software packages and then

380    compare the outcome of each approach, highlighting not only the similarity, but also differences

381    among these analyses. Fold change is an important parameter to consider, but will depend on the

382    number of reads that are assigned to a specific transcript. If the depth is low, yet with high fold change

383    between groups, it should be considered as noise. For example, 10 X 100 base reads mapped onto a

384    1 Kb transcript per sample in one group (giving an average depth of 1) compared to 1 read on average

385    per sample in the other group is a 10-fold change, yet the coverage is very low and should be validated

386    using additional samples via qPCR.

387    **Annotation of transcripts**

388    After all reads have been assembled *de novo* into contigs, the next step is to annotate all the contigs

389    based on the most up-to-date database (i.e. identify homology to previously characterized genes). The

390    most common way to annotate a large number of transcripts is the Basic Local Alignment Search Tool

391    (BLAST). As the number of contigs in every *de novo* assembly can be thousands to a few hundred

392    thousand sequences, usage of an automated search tool, in particular BLAST+ (Camacho et al. 2009),

393    is essential. For non-model species, many candidate protein databases are available including the non-

394    redundant protein database (nr), UniProtKB/Swiss-Prot database, and the Reference Sequence

395    database (RefSeq). RefSeq (nucleotide and protein) and UniProt/Swiss-Prot (protein) consist of

396    curated, well annotated sequences, whereas the nr database includes both curated and non-curated

397    databases. For most crustacean RNA-seq experiments, the nr database is considered to be the best

398    choice due to the fact that very few crustacean genes have been properly annotated to date, a

399    problem that has been highlighted (Clark and Greenwood 2016; Das and Mykles 2016).

400    After transcripts have been scanned against the protein database and assigned annotations, there is

401    a variety of downstream packages that can further analyse a contig, including Gene Ontology (GO)

402    term analysis, functional enrichment analysis, protein domain analysis (PFAM domain search -

403    pfam.xfam.org), and Kyoto Encyclopedia of Genes and Genomes (KEGG) pathway maps (Kanehisa and

404    Goto 2000). Each entry in the sequence database can be classified into a number of biologically

405    relevant terms. In GO analysis, most genes can be assigned to one out of three basic ontologies:

406    cellular component, biological process or molecular function. When comparing samples from two

407    groups, differentially represented GO terms can help define the mechanism via which the groups differ

408    from one another. Similarly, when using KEGG, contigs can be assigned with components of specific

409    pathways. Differential expression can allow a detailed assessment of the changes in pathways

410    between studied groups. PFAM shows domains within the open reading frames of contigs that enables

411    characterization of the protein function, based on the architecture of domains in a polypeptide chain.

412    A number of software packages are capable of extracting vast numbers of GO terms from public

413    databases including Blast2GO (Conesa et al. 2005), DAVID - https://david.abcc.ncifcrf.gov/ (Huang da

414    et al. 2009) or ermineJ (Lee et al. 2005). Among these programs, Blast2GO stands out as an easy to

415    use, point-and-click program that has become very popular in the last few years.

416     In addition to Blast2GO, databases like PFAM (Finn et al. 2016), eggNOG (Huerta-Cepas et al. 2015)

417     and InterProScan (Jones et al. 2014) can be employed to predict the function of unknown proteins.

418     The Trinity RNA-seq package, Trinotate ([https://trinotate.github.io/](https://trinotate.github.io/)), uses UniProt, eggNOG and GO

419     Pathway databases for annotating novel sequences and these have been widely used in recent years

420     (Das and Mykles 2016; Das et al. 2016).


421     **Validation of RNA-seq results**

422     Validation is a very important step in every RNA-seq study. There is generally a very high correlation

423     between RNA-seq and qPCR results with respect to relative gene expression. A significant point is that

424     testing the same RNA samples used in the NGS platform for validation with techniques like qPCR or

425     digital PCR only validates the sequencing accuracy result. Therefore, additional, independent

426     biological replicates should be included to properly validate the biological interpretation from the

427     RNA-seq experiment. Essentially, validation post sequencing is now mandatory for publication. An

428     approach has been proposed to set the minimum acceptable standard for qPCR validation (Fang and

429     Cui 2011) that takes a number of factors into consideration, including the number of genes tested and

430     the number of isoform transcripts detected in the transcriptome. Nevertheless, up-to-date, qPCR

431     techniques offer the easiest way to validate data in a transcriptomics study. One important note for

432     researchers who are unfamiliar with the technology is that some RNA-seq pipelines allow RNA-seq

433     analysis at the gene level (Trinity/RSEM for instance (Haas et al. 2013)). However, there is a deeper

434     level of transcripts component (in which transcripts can be isoforms resulting from alternative splicing

435     events or a single nucleotide variation). Therefore, researchers should design primers that are not

436     included in these regions to avoid unreliable qPCR results between biological replicates. As a

437     concluding remark for this section, several studies have compared RNA-seq results to qPCR data, and

438     have found excellent correlations between the approaches (Everaert et al. 2017; Rajkumar et al. 2015;

439     Wu et al. 2014).


440

# PROMISES AND CHALLENGES OF RNA-SEQ BASED STUDIES IN CRUSTACEANS

**Studied topics**

SGS has revolutionized biological science, shifting it toward the post-genomics era. Transcriptomics studies in crustaceans include either:

1. Sequencing and annotation of the transcriptome of one (or several) tissue/s, or a whole individual of a particular taxa in a specific developmental stage or under specific experimental conditions.
2. Applying RNA-seq to identify DEGs among different physiological conditions, treatments, developmental stages and/or tissues.
3. Identification of novel transcripts – enzymes, receptors, hormones, neuropeptides.
4. Screening for variant mutations - SNPs, SSRs and/or microsatellites.
5. A combination of the above.

To date, several RNA-seq projects have been initiated on a variety of crustacean species. In Table 5, we have summarized several RNA-seq based studies on crustacean taxa that have been conducted over the last few years based on the following categories: *Aquatic toxicology*, *Reproduction & sexual differentiation*, *Disease resistance & immunology*, *Developmental biology*, and *Physiology*. This is however, by no means an exhaustive list as hundreds of applied RNA-seq studies have been undertaken in recent years, rather the list here illustrates several model comparative RNA-seq approaches.

**Ongoing challenges for applied RNA-seq studies of crustaceans**

– *Experimental limitation:* A good experimental design will have a major impact on data outcomes; it can prevent wasted resources and help avoid the generation of unpublishable results. The balance between sequencing cost and experimental design constraints is a major issue that has been highlighted in many review articles. Due to budgetary limitations, there will always be an incentive to cut costs by sequencing with higher depth but with little or no biological replication. Furthermore, where depth is added, a large number of reads will be also mapped to the already well-covered regions, while if additional replication was available, greater statistical power can be achieved resulting in better biological inference. To resolve this problem, optimal guidelines for the design of RNA-seq experiments are needed and should be applied accordingly. In parallel, biological replicates (at least 3 or greater) are required for an RNA-seq study to reach a basic publishable level.  As another recommendation for best practice, is undertaking a pilot-sequencing project where a high number of libraries are run on one lane initially. This can be valuable in assessing the feasibility of the larger experiment, as well as providing a good indicator for how to address trade-offs between obtaining high quality output vs cost. Finally, although RNA-seq methods are becoming more robust and reliable and sometimes qPCR validations are proven to be unnecessary, a section for qPCR validation of selected genes/transcripts of interest may be beneficial to reveal the biological insights if the study has limited replications. Therefore, we recommend that for reliable biological interpretation and validation of RNA-seq analysis, the candidate genes themselves are tested for expression, rather than choosing random genes or genes showing high expression levels.

– *In silico annotation and functional annotation:* Annotation of RNA-seq data is based loosely on BLAST searches. In fact, many BLAST results produce "hypothetical", "predicted", "uncharacterized", or "low-quality" assignments. This highlights the fact that gene databases for non-model species currently, are very limited. To add another layer of complexity, *Daphnia pulex*, the model species currently available for crustaceans, has a large number of genes that currently remain

486    unannotated. Furthermore, when compared with Decapoda, it is very remotely related and in many

487    cases, shares higher similarity with insects than with other crustacean taxa. Further downstream

488    annotation is also a constraint for crustacean RNA-seq studies, as specific GO classification and KEGG

489    pathways are still not available for these taxa. As a result, drawing biological interpretations from

490    predicted results can be problematic. Moreover, similarities in structure do not necessarily correlate

491    with equivalent functionalities. It is crucial therefore, to highlight that *in silico* prediction is only

492    speculative and functional annotation is very important to validate any biological interpretations (in

493    particular for novel genes). RNAi technology (gene silencing) is now the go-to method for gene

494    functional studies in decapod crustaceans and it has been already applied in some cases (Sagi et al.

495    2013). Gene editing technologies, for example CRISPR/Cas9 technology, have emerged recently and

496    hold great potential for functional annotation in decapod crustacean species (Mykles and Hui 2015).

497    Employing RNAi and/or CRISPR-Cas9 in RNA-seq studies would be extremely helpful to highlight key

498    genes and resolve functional roles of novel genes for crustacean species.

499    – *Combining transcriptomics/RNA-seq with other OMICS techniques:* In parallel with advances in

500    RNA-seq technologies, other OMICS technologies including genomics, proteomics, metagenomics

501    phylogenomics and phenomics have also developed rapidly. This highlights a challenge for RNA-seq

502    studies, to make use of other OMICS approaches and to utilize them to create a multilayer outcome.

503    One key reason why decapod crustacean genomes are not yet available is that they are often very

504    large and complex which makes them hard to resolve. Nevertheless, draft genomes of a few

505    crustacean species have been made publicly available recently including draft genomes for some

506    decapods including: *N. denticulata* (Kenny et al. 2014), *P. vannamei* (Yu et al. 2015), *E. sinensis* (Song

507    et al. 2016), *P. hawaiensis* (Kao et al. 2016), *P. monodon and M. japonicus (Yuan et al. 2017), and P.

508    virginalis (Gutekunst et al. 2018)*. Utilizing these new genomic resources will allow better gene

509    annotation and functional annotation of crustacean gene pathways. There is no doubt that in the near

510    future, when the cost barrier for sequencing is essentially overcome, coupled with improved

sequencing technologies, combining RNA-seq approaches with integrated OMICS will enable

researchers to answer the most complex of biological questions.

## CONCLUSIONS

To conclude, RNA-seq offers great promise for crustacean studies. It is a very powerful tool that can

lead to developing a better understanding of underlying pathways and mechanisms that form the

basis of many scientific questions. The guidelines offered here for future RNA-seq studies of

crustaceans are an attempt to assist biologists who are not familiar with the complex and diverse array

of bioinformatics software that are currently available. It is also important however, to highlight the

gap between *in silico* prediction from RNA-seq analysis and *in vivo* results. This may be explained in

general, by limitations on experimental designs in the past, the lack of annotation databases for

crustacean species, as well as the need for question-driven research. In the future, we also suggest

that RNA-seq should be integrated with other OMICs technologies to increase data output as well as

improving biological insights.

534 **Species abbreviation**

535 *P. trituberculatus: Portunus trituberculatus*

536 *S. henanense: Sinopotamon henanense*

537 *P. vannamei: Penaeus vannamei*

538 *P. monodon: Penaeus monodon*

539 *M. japonicus: Marsupenaeus japonicus*

540 *P. virginalis: Procambarus virginalis or*
541 *Procambarus fallax forma virginalis*

542 *S. olivacea: Scylla olivacea*

543 *S. paramamosain: Scylla paramamosain*

544 *E. sinensis: Eriocheir sinensis*

545 *N. norvegicus: Nephrops norvegicus*

546 *F. merguiensis: Fenneropenaeus merguiensis*

547 *P. clarkii: Procambrarus clarkii*

548 *M. rosenbergii: Macrobrachium rosenbergii*

549 *S. verreauxi : Sagmariasus verreauxi*

550 *N. denticulata: Neocaridina denticulata*

551 *P. hawaiensis: Parhyale hawaiensis*

552 *E. carinicauda: Exopalaemon carinicauda*

553 *M. olfersi: Macrobrachium olfersi*

554 *P. elegans: Palaemon elegans*

555 *P. australiensis: Paratya australiensis*

556 *E. carinicauda: Exopalaemon carinicauda*

557 *H. rubra: Halocaridina rubra*

558

## List of Figures

560

561 Figure 1. **A simple flowchart of an RNA-seq/transcriptomics study.**

562 Figure 2. **Technical samples versus biological samples.**

## List of Tables

564

565 Table 1. **Commonly used quality control software for NGS data**. *+: Yes, -: No*

566 Table 2. **Summary of commonly used de novo transcriptome assemblers.**

567 *Platforms -*  *Linux,*  *Windows,*  *MacOS. License - C: Commercial product, F: Free.*

568 Table 3. **Popular read aligners for RNA-seq.** *+: Yes, -: No*

569 Table 4. **Widely used DEG analysis tools.**

570 Table 5. **Summary of recent RNA-seq studies in non-model decapod crustacean species.**

571

572

573    **References**

574    Amin, S., P. Prentis, E. Gilding & A. Pavasovic, 2014. Assembly and annotation of a non-model
575        gastropod (Nerita melanotragus) transcriptome: a comparison of *De novo* assemblers. BMC
576        Research Notes 7(1):488 doi:10.1186/1756-0500-7-488.
577    Anders, S., P. T. Pyl & W. Huber, 2015. HTSeq--a Python framework to work with high-throughput
578        sequencing data. Bioinformatics 31(2):166-9 doi:10.1093/bioinformatics/btu638.
579    Auer, P. L. & R. W. Doerge, 2010. Statistical design and analysis of RNA sequencing data. Genetics
580        185(2):405-416 doi:10.1534/genetics.110.114983.
581    Baruzzo, G., K. E. Hayer, E. J. Kim, B. Di Camillo, G. A. FitzGerald & G. R. Grant, 2017. Simulation-
582        based comprehensive benchmarking of RNA-seq aligners. Nature Methods 14(2):135-139
583        doi:10.1038/nmeth.4106.
584    Bray, N. L., H. Pimentel, P. Melsted & L. Pachter, 2016. Near-optimal probabilistic RNA-seq
585        quantification. Nature Biotechnology 34(5):525-527 doi:10.1038/nbt.3519.
586    Busby, M. A., C. Stewart, C. A. Miller, K. R. Grzeda & G. T. Marth, 2013. Scotty: a web tool for
587        designing RNA-Seq experiments to measure differential gene expression. Bioinformatics
588        29(5):656-657 doi:10.1093/bioinformatics/btt015.
589    Camacho, C., G. Coulouris, V. Avagyan, N. Ma, J. Papadopoulos, K. Bealer & T. L. Madden, 2009.
590        BLAST+: architecture and applications. BMC Bioinformatics 10:421 doi:10.1186/1471-2105-
591        10-421.
592    Cartolano, M., B. Huettel, B. Hartwig, R. Reinhardt & K. Schneeberger, 2016. cDNA library
593        enrichment of full length transcripts for SMRT long read sequencing. PLoS ONE
594        11(6):e0157779 doi:10.1371/journal.pone.0157779.
595    Chandramohan, R., P. Y. Wu, J. H. Phan & M. D. Wang, 2013. Systematic assessment of RNA-Seq
596        quantification tools using simulated sequence data. ACM Conference on Bioinformatics,
597        Computational Biology and Biomedicine 2013 doi:10.1145/2506583.2506648.
598    Chen, S.-Y., F. Deng, X. Jia, C. Li & S.-J. Lai, 2017. A transcriptome atlas of rabbit revealed by PacBio
599        single-molecule long-read sequencing. Scientific Reports 7(1):7648 doi:10.1038/s41598-017-
600        08138-z.
601    Clark, K. F. & S. J. Greenwood, 2016. Next-Generation Sequencing and the crustacean immune
602        system: The need for alternatives in immune gene annotation. Integrative and Comparative
603        Biology 56(6):1113-1130 doi:10.1093/icb/icw023.
604    Conesa, A., S. Gotz, J. M. Garcia-Gomez, J. Terol, M. Talon & M. Robles, 2005. Blast2GO: a universal
605        tool for annotation, visualization and analysis in functional genomics research.
606        Bioinformatics 21(18):3674-6 doi:10.1093/bioinformatics/bti610.
607    Conesa, A., P. Madrigal, S. Tarazona, D. Gomez-Cabrero, A. Cervera, A. McPherson, M. W. Szcześniak,
608        D. J. Gaffney, L. L. Elo, X. Zhang & A. Mortazavi, 2016. A survey of best practices for RNA-seq
609        data analysis. Genome Biology 17(1):13 doi:10.1186/s13059-016-0881-8.
610    Das, S. & D. L. Mykles, 2016. A comparison of resources for the annotation of a *de novo* assembled
611        transcriptome in the molting gland (Y-Organ) of the Blackback Land Crab, *Gecarcinus*
612        *lateralis*. Integrative and Comparative Biology 56(6):1103-1112 doi:10.1093/icb/icw107.
613    Das, S., S. Shyamal & D. S. Durica, 2016. Analysis of annotation and differential expression methods
614        used in RNA-seq Studies in crustacean systems. Integrative and Comparative Biology
615        56(6):1067-1079 doi:10.1093/icb/icw117.
616    Del Fabbro, C., S. Scalabrin, M. Morgante & F. M. Giorgi, 2013. An extensive evaluation of read
617        trimming effects on Illumina NGS data analysis. PLoS ONE 8(12):e85024
618        doi:10.1371/journal.pone.0085024.
619    Dillies, M.-A., A. Rau, J. Aubert, C. Hennequet-Antier, M. Jeanmougin, N. Servant, C. Keime, G. Marot,
620        D. Castel, J. Estelle, G. Guernec, B. Jagla, L. Jouneau, D. Laloë, C. Le Gall, B. Schaëffer, S. Le
621        Crom, M. Guedj & F. Jaffrézic, 2013. A comprehensive evaluation of normalization methods

622  for Illumina high-throughput RNA sequencing data analysis. Briefings in Bioinformatics
623      14(6):671-683 doi:10.1093/bib/bbs046.
624  Everaert, C., M. Luypaert, J. L. V. Maag, Q. X. Cheng, M. E. Dinger, J. Hellemans & P. Mestdagh, 2017.
625      Benchmarking of RNA-sequencing analysis workflows using whole-transcriptome RT-qPCR
626      expression data. Scientific Reports 7:1559 doi:10.1038/s41598-017-01617-3.
627  Ewing, B., L. Hillier, M. C. Wendl & P. Green, 1998. Base-calling of automated sequencer traces using
628      Phred. I. Accuracy assessment. Genome research 8(3):175-185.
629  Fang, Z. & X. Cui, 2011. Design and validation issues in RNA-seq experiments. Briefings in
630      Bioinformatics doi:10.1093/bib/bbr004.
631  Feng, J., C. A. Meyer, Q. Wang, J. S. Liu, X. Shirley Liu & Y. Zhang, 2012. GFOLD: a generalized fold
632      change for ranking differentially expressed genes from RNA-seq data. Bioinformatics
633      28(21):2782-8 doi:10.1093/bioinformatics/bts515.
634  Finn, R. D., P. Coggill, R. Y. Eberhardt, S. R. Eddy, J. Mistry, A. L. Mitchell, S. C. Potter, M. Punta, M.
635      Qureshi, A. Sangrador-Vegas, G. A. Salazar, J. Tate & A. Bateman, 2016. The Pfam protein
636      families database: towards a more sustainable future. Nucleic Acids Research 44(D1):D279-
637      D285 doi:10.1093/nar/gkv1344.
638  Finseth, F. R. & R. G. Harrison, 2014. A comparison of Next Generation Sequencing technologies for
639      transcriptome assembly and utility for RNA-Seq in a non-model bird. PLoS ONE
640      9(10):e108550 doi:10.1371/journal.pone.0108550.
641  Fonseca, N. A., J. Rung, A. Brazma & J. C. Marioni, 2012. Tools for mapping high-throughput
642      sequencing data. Bioinformatics doi:10.1093/bioinformatics/bts605.
643  Francis, W. R., L. M. Christianson, R. Kiko, M. L. Powers, N. C. Shaner & S. H. Haddock, 2013. A
644      comparison across non-model animals suggests an optimal sequencing depth for *de novo*
645      transcriptome assembly. BMC Genomics 14(1):167.
646  Ghangal, R., S. Chaudhary, M. Jain, R. S. Purty & P. Chand Sharma, 2013. Optimization of *De Novo*
647      short read assembly of Seabuckthorn *Hippophae rhamnoide*s L. transcriptome. PLoS ONE
648      8(8):e72516 doi:10.1371/journal.pone.0072516.
649  Glenn, T. C., 2011. Field guide to next-generation DNA sequencers. Molecular Ecology Resources
650      11(5):759-769 doi:10.1111/j.1755-0998.2011.03024.x.
651  Goodwin, S., J. D. McPherson & W. R. McCombie, 2016. Coming of age: ten years of Next Generation
652      Sequencing technologies. Nature Reviews Genetics 17(6):333-351 doi:10.1038/nrg.2016.49.
653  Grabherr, M., B. Haas, M. Yassour, J. Levin, D. Thompson, I. Amit, X. Adiconis, L. Fan, R.
654      Raychowdhury, Q. Zeng, Z. Chen, E. Mauceli, N. Hacohen, A. Gnirke, N. Rhind, F. di Palma, B.
655      Birren, C. Nusbaum, K. Lindblad-Toh, N. Friedman & A. Regev, 2011. Full-length
656      transcriptome assembly from RNA-Seq data without a reference genome. Nature
657      Biotechnology 29(7):644-652 doi:10.1038/nbt.1883.
658  Grant, G. R., M. H. Farkas, A. D. Pizarro, N. F. Lahens, J. Schug, B. P. Brunk, C. J. Stoeckert, J. B.
659      Hogenesch & E. A. Pierce, 2011. Comparative analysis of RNA-Seq alignment algorithms and
660      the RNA-Seq unified mapper (RUM). Bioinformatics 27(18):2518-2528
661      doi:10.1093/bioinformatics/btr427.
662  Gutekunst, J., R. Andriantsoa, C. Falckenhayn, K. Hanna, W. Stein, J. Rasamy & F. Lyko, 2018. Clonal
663      genome evolution and rapid invasive spread of the marbled crayfish. Nature Ecology &
664      Evolution 2(3):567-573 doi:10.1038/s41559-018-0467-9.
665  Haas, B. J., A. Papanicolaou, M. Yassour, M. Grabherr, P. D. Blood, J. Bowden, M. B. Couger, D.
666      Eccles, B. Li, M. Lieber, M. D. MacManes, M. Ott, J. Orvis, N. Pochet, F. Strozzi, N. Weeks, R.
667      Westerman, T. William, C. N. Dewey, R. Henschel, R. D. LeDuc, N. Friedman & A. Regev,
668      2013. *De novo* transcript sequence reconstruction from RNA-Seq: reference generation and
669      analysis with Trinity. Nature Protocols 8(8):10.1038/nprot.2013.084
670      doi:10.1038/nprot.2013.084.
671  Hatem, A., D. Bozdağ, A. E. Toland & Ü. V. Çatalyürek, 2013. Benchmarking short sequence mapping
672      tools. BMC Bioinformatics 14(1):184 doi:10.1186/1471-2105-14-184.

673 Havird, J. C. & S. R. Santos, 2016. Here we are, but where do we go? A systematic review of
674     crustacean transcriptomic studies from 2014–2015. Integrative and Comparative Biology
675     56(6):1055-1066 doi:10.1093/icb/icw061.

676 Huang da, W., B. T. Sherman & R. A. Lempicki, 2009. Systematic and integrative analysis of large
677     gene lists using DAVID bioinformatics resources. Nature Protocols 4(1):44-57
678     doi:10.1038/nprot.2008.211.

679 Huerta-Cepas, J., D. Szklarczyk, K. Forslund, H. Cook, D. Heller, M. C. Walter, T. Rattei, D. R. Mende, S.
680     Sunagawa & M. Kuhn, 2015. eggNOG 4.5: a hierarchical orthology framework with improved
681     functional annotations for eukaryotic, prokaryotic and viral sequences. Nucleic Acids
682     Research 44(D1):D286-D293 doi:10.1093/nar/gkv1248.

683 Jaramillo, M. L., F. Guzman, C. L. Paese, R. Margis, E. M. Nazari, D. Ammar & Y. M. R. Müller, 2016.
684     Exploring developmental gene toolkit and associated pathways in a potential new model
685     crustacean using transcriptomic analysis. Development Genes and Evolution 226(5):325-337
686     doi:10.1007/s00427-016-0551-6.

687 Jin, S., H. Fu, Q. Zhou, S. Sun, S. Jiang, Y. Xiong, Y. Gong, H. Qiao & W. Zhang, 2013. Transcriptome
688     analysis of androgenic gland for discovery of novel genes from the oriental river prawn,
689     *Macrobrachium nipponense*, using Illumina Hiseq 2000. PloS one 8(10):e76840
690     doi:10.1371/journal.pone.0076840.

691 Jones, P., D. Binns, H.-Y. Chang, M. Fraser, W. Li, C. McAnulla, H. McWilliam, J. Maslen, A. Mitchell, G.
692     Nuka, S. Pesseat, A. F. Quinn, A. Sangrador-Vegas, M. Scheremetjew, S.-Y. Yong, R. Lopez &
693     S. Hunter, 2014. InterProScan 5: Genome-scale protein function classification. Bioinformatics
694     30(9):1236-1240 doi:10.1093/bioinformatics/btu031.

695 Jung, H., R. E. Lyons, H. Dinh, D. A. Hurwood, S. McWilliam & P. B. Mather, 2011. Transcriptomics of
696     a Giant Freshwater Prawn (*Macrobrachium rosenbergii*): *De novo* assembly, annotation and
697     marker discovery. PLoS ONE 6(12):e27938 doi:10.1371/journal.pone.0027938.

698 Jung, H., B.-H. Yoon, W.-J. Kim, D.-W. Kim, D. Hurwood, R. Lyons, K. Salin, H.-S. Kim, I. Baek, V. Chand
699     & P. Mather, 2016. Optimizing hybrid *de novo* transcriptome assembly and extending
700     genomic resources for Giant Freshwater Prawns (*Macrobrachium rosenbergii*): The
701     identification of genes and markers associated with reproduction. International Journal of
702     Molecular Sciences 17(5):690 doi:10.3390/ijms17050690.

703 Kanehisa, M. & S. Goto, 2000. KEGG: Kyoto Encyclopedia of Genes and Genomes. Nucleic Acids
704     Research 28(1):27-30.

705 Kao, D., A. G. Lai, E. Stamataki, S. Rosic, N. Konstantinides, E. Jarvis, A. Di Donfrancesco, N.
706     Pouchkina-Stancheva, M. Sémon, M. Grillo, H. Bruce, S. Kumar, I. Siwanowicz, A. Le, A.
707     Lemire, M. B. Eisen, C. Extavour, W. E. Browne, C. Wolff, M. Averof, N. H. Patel, P. Sarkies, A.
708     Pavlopoulos & A. Aboobaker, 2016. The genome of the crustacean *Parhyale hawaiensis*, a
709     model for animal development, regeneration, immunity and lignocellulose digestion. eLife
710     5:e20062 doi:10.7554/eLife.20062.

711 Kelley, D. R., M. C. Schatz & S. L. Salzberg, 2010. Quake: quality-aware detection and correction of
712     sequencing errors. Genome Biology 11(11):R116 doi:10.1186/gb-2010-11-11-r116.

713 Kenny, N. J., Y. W. Sin, X. Shen, Q. Zhe, W. Wang, T. F. Chan, S. S. Tobe, S. M. Shimeld, K. H. Chu & J.
714     H. Hui, 2014. Genomic sequence and experimental tractability of a new decapod shrimp
715     model, *Neocaridina denticulata*. Marine Drugs 12(3):1419-37 doi:10.3390/md12031419.

716 Khang, T. F. & C. Y. Lau, 2015. Getting the most out of RNA-seq data analysis. PeerJ 3:e1360
717     doi:10.7717/peerj.1360.

718 Koboldt, Daniel C., Karyn M. Steinberg, David E. Larson, Richard K. Wilson & E. R. Mardis, 2013. The
719     next-generation sequencing revolution and its impact on genomics. Cell 155(1):27-38
720     doi:10.1016/j.cell.2013.09.006.

721 Kuo, R. I., E. Tseng, L. Eory, I. R. Paton, A. L. Archibald & D. W. Burt, 2017. Normalized long read RNA
722     sequencing in chicken reveals transcriptome complexity similar to human. BMC Genomics
723     18(1):323 doi:10.1186/s12864-017-3691-9.

724 Kvam, V. M., P. Liu & Y. Si, 2012. A comparison of statistical methods for detecting differentially
725     expressed genes from RNA-seq data. American Journal of Botany 99(2):248-256
726     doi:10.3732/ajb.1100340.
727 Lahens, N. F., E. Ricciotti, O. Smirnova, E. Toorens, E. J. Kim, G. Baruzzo, K. E. Hayer, T. Ganguly, J.
728     Schug & G. R. Grant, 2017. A comparison of Illumina and Ion Torrent sequencing platforms in
729     the context of differential gene expression. BMC Genomics 18(1):602 doi:10.1186/s12864-
730     017-4011-0.
731 Lam, H. Y. K., M. J. Clark, R. Chen, R. Chen, G. Natsoulis, M. O'Huallachain, F. E. Dewey & L. Habegger,
732     2012. Performance comparison of whole-genome sequencing platforms. Nature
733     Biotechnology 30 doi:10.1038/nbt.2065.
734 Langmead, B. & S. L. Salzberg, 2012. Fast gapped-read alignment with Bowtie 2. Nature Methods
735     9(4):357-359 doi:10.1038/nmeth.1923.
736 Langmead, B., C. Trapnell, M. Pop & S. L. Salzberg, 2009. Ultrafast and memory-efficient alignment of
737     short DNA sequences to the human genome. Genome Biology 10(3):R25 doi:10.1186/gb-
738     2009-10-3-r25.
739 Ledergerber, C. & C. Dessimoz, 2011. Base-calling for next-generation sequencing platforms.
740     Briefings in Bioinformatics 12(5):489-497 doi:10.1093/bib/bbq077.
741 Lee, H. K., W. Braynen, K. Keshav & P. Pavlidis, 2005. ErmineJ: tool for functional analysis of gene
742     expression data sets. BMC Bioinformatics 6(1):269.
743 Levin, J. Z., M. Yassour, X. Adiconis, C. Nusbaum, D. A. Thompson, N. Friedman, A. Gnirke & A. Regev,
744     2010. Comprehensive comparative analysis of strand-specific RNA sequencing methods.
745     Nature Methods 7(9):709-15 doi:10.1038/nmeth.1491.
746 Li, B. & C. N. Dewey, 2011. RSEM: accurate transcript quantification from RNA-Seq data with or
747     without a reference genome. BMC Bioinformatics 12:323 doi:10.1186/1471-2105-12-323.
748 Li, H. & R. Durbin, 2009. Fast and accurate short read alignment with Burrows–Wheeler transform.
749     Bioinformatics 25(14):1754-1760 doi:10.1093/bioinformatics/btp324.
750 Li, H. & N. Homer, 2010. A survey of sequence alignment algorithms for next-generation sequencing.
751     Briefings in Bioinformatics 11(5):473-483 doi:10.1093/bib/bbq015.
752 Li, Z., Y. Chen, D. Mu, J. Yuan, Y. Shi, H. Zhang, J. Gan, N. Li, X. Hu, B. Liu, B. Yang & W. Fan, 2012.
753     Comparison of the two major classes of assembly algorithms: overlap-layout-consensus and
754     de-Bruijn-graph. Briefings in Functional Genomics 11(1):25-37 doi:10.1093/bfgp/elr035.
755 Lister, R., B. D. Gregory & J. R. Ecker, 2009. Next is now: new technologies for sequencing of
756     genomes, transcriptomes, and beyond. Current Opinion in Plant Biology 12(2):107-118
757     doi:10.1016/j.pbi.2008.11.004.
758 Liu, L., Y. Li, S. Li, N. Hu, Y. He, R. Pong, D. Lin, L. Lu & M. Law, 2012. Comparison of Next-Generation
759     Sequencing systems. Journal of Biomedicine and Biotechnology 2012:11
760     doi:10.1155/2012/251364.
761 Liu, Y., J. Zhou & K. P. White, 2014. RNA-seq differential expression studies: more sequence or more
762     replication? Bioinformatics 30(3):301-304 doi:10.1093/bioinformatics/btt688.
763 Love, M. I., W. Huber & S. Anders, 2014. Moderated estimation of fold change and dispersion for
764     RNA-seq data with DESeq2. Genome Biology 15(12):550 doi:10.1186/s13059-014-0550-8.
765 Lv, J., P. Liu, B. Gao, Y. Wang, Z. Wang, P. Chen & J. Li, 2014. Transcriptome analysis of the *Portunus
766     trituberculatus*: de novo assembly, growth-related gene identification and marker discovery.
767     PLoS one 9(4):e94055.
768 Macharia, R. W., F. L. Ombura & E. O. Aroko, 2015. Insects RNA profiling reveals absence of hidden
769     break in 28S Ribosomal RNA molecule of Onion Thrips, *Thrips tabaci*. Journal of Nucleic Acids
770     2015:8 doi:10.1155/2015/965294.
771 MacManes, M. D., 2014. On the optimal trimming of high-throughput mRNA sequence data.
772     Frontiers in Genetics 5:13 doi:10.3389/fgene.2014.00013.
773 Marguerat, S. & J. Bähler, 2010. RNA-seq: from technology to biology. Cellular and Molecular Life
774     Sciences 67(4):569-579 doi:10.1007/s00018-009-0180-6.

775    McCarthy, S. D., M. M. Dugon & A. M. Power, 2015. 'Degraded' RNA profiles in Arthropoda and
776         beyond. PeerJ 3:e1436 doi:10.7717/peerj.1436.
777    Meng, X.-l., P. Liu, F.-l. Jia, J. Li & B.-Q. Gao, 2015. *De novo* transcriptome analysis of *Portunus*
778         *trituberculatus* ovary and testis by RNA-Seq: Identification of genes involved in gonadal
779         development. PLoS ONE 10(6):e0128659 doi:10.1371/journal.pone.0128659.
780    Metzker, M. L., 2010. Sequencing technologies—the next generation. Nature Reviews Genetics
781         11(1):31-46 doi:10.1038/nrg2626.
782    Miller, J. R., S. Koren & G. Sutton, 2010. Assembly algorithms for Next-Generation Sequencing data.
783         Genomics 95(6):315-327 doi:10.1016/j.ygeno.2010.03.001.
784    Misner, I., C. Bicep, P. Lopez, S. Halary, E. Bapteste & C. E. Lane, 2013. Sequence comparative
785         analysis using networks: Software for evaluating *de novo* transcript assembly from Next-
786         Generation Sequencing. Molecular Biology and Evolution 30(8):1975-1986
787         doi:10.1093/molbev/mst087.
788    Mykles, D. L., K. G. Burnett, D. S. Durica, B. L. Joyce, F. M. McCarthy, C. J. Schmidt & J. H. Stillman,
789         2016. Resources and recommendations for using transcriptomics to address grand
790         challenges in comparative biology. Integrative and Comparative Biology 56(6):1183-1191
791         doi:10.1093/icb/icw083.
792    Mykles, D. L. & J. H. Hui, 2015. *Neocaridina denticulata*: A decapod crustacean model for Functional
793         Genomics. Integrative and Comparative Biology 55(5):891-7 doi:10.1093/icb/icv050.
794    Nguyen, C., T. G. Nguyen, L. Van Nguyen, H. Q. Pham, T. H. Nguyen, H. T. Pham, H. T. Nguyen, T. T.
795         Ha, T. H. Dau & H. T. Vu, 2016. *De novo* assembly and transcriptome characterization of
796         major growth-related genes in various tissues of *Penaeus monodon*. Aquaculture 464:545-
797         553 doi:10.1016/j.aquaculture.2016.08.003.
798    Niedringhaus, T. P., D. Milanova, M. B. Kerby, M. P. Snyder & A. E. Barron, 2011. Landscape of Next-
799         Generation Sequencing technologies. Analytical Chemistry 83(12):4327-4341
800         doi:10.1021/ac2010857.
801    Ozsolak, F. & P. M. Milos, 2011. RNA sequencing: advances, challenges and opportunities. Nature
802         Reviews Genetics 12(2):87-98.
803    Parkhomchuk, D., T. Borodina, V. Amstislavskiy, M. Banaru, L. Hallen, S. Krobitsch, H. Lehrach & A.
804         Soldatov, 2009. Transcriptome analysis by strand-specific sequencing of complementary
805         DNA. Nucleic Acids Research 37(18):e123-e123 doi:10.1093/nar/gkp596.
806    Parra, G., K. Bradnam & I. Korf, 2007. CEGMA: a pipeline to accurately annotate core genes in
807         eukaryotic genomes. Bioinformatics 23(9):1061-1067 doi:10.1093/bioinformatics/btm071.
808    Patro, R., G. Duggal, M. I. Love, R. A. Irizarry & C. Kingsford, 2017. Salmon provides fast and bias-
809         aware quantification of transcript expression. Nature Methods doi:10.1038/nmeth.4197.
810    Rajkumar, A. P., P. Qvist, R. Lazarus, F. Lescai, J. Ju, M. Nyegaard, O. Mors, A. D. Børglum, Q. Li & J. H.
811         Christensen, 2015. Experimental validation of methods for differential gene expression
812         analysis and sample pooling in RNA-seq. BMC Genomics 16(1):548 doi:10.1186/s12864-015-
813         1767-y.
814    Reuter, Jason A., D. V. Spacek & Michael P. Snyder, 2015. High-Throughput Sequencing Technologies.
815         Molecular Cell 58(4):586-597 doi:10.1016/j.molcel.2015.05.004.
816    Roberts, A. & L. Pachter, 2012. Streaming fragment assignment for real-time analysis of sequencing
817         experiments. Nature Methods 10:71 doi:10.1038/nmeth.2251.
818    Robertson, G., J. Schein, R. Chiu, R. Corbett, M. Field, S. D. Jackman, K. Mungall, S. Lee, H. M. Okada,
819         J. Q. Qian, M. Griffith, A. Raymond, N. Thiessen, T. Cezard, Y. S. Butterfield, R. Newsome, S.
820         K. Chan, R. She, R. Varhol, B. Kamoh, A.-L. Prabhu, A. Tam, Y. Zhao, R. A. Moore, M. Hirst, M.
821         A. Marra, S. J. M. Jones, P. A. Hoodless & I. Birol, 2010. *De novo* assembly and analysis of
822         RNA-seq data. Nature Methods 7(11):909-912 doi:10.1038/nmeth.1517.
823    Robinson, M. D., D. J. McCarthy & G. K. Smyth, 2010. edgeR: a Bioconductor package for differential
824         expression analysis of digital gene expression data. Bioinformatics 26(1):139-140
825         doi:10.1093/bioinformatics/btp616.

826 Robles, J. A., S. E. Qureshi, S. J. Stephen, S. R. Wilson, C. J. Burden & J. M. Taylor, 2012. Efficient
827      experimental design and analysis strategies for the detection of differential expression using
828      RNA-Sequencing. BMC Genomics 13(1):484 doi:10.1186/1471-2164-13-484.
829 Sagi, A., R. Manor & T. Ventura, 2013. Gene silencing in Crustaceans: From basic research to
830      biotechnologies. Genes 4(4):620 doi:10.3390/genes4040620.
831 Schadt, E. E., S. Turner & A. Kasarskis, 2010. A window into third-generation sequencing. Human
832      Molecular Genetics 19(R2):R227-R240 doi:10.1093/hmg/ddq416.
833 Schirmer, M., U. Z. Ijaz, R. D'Amore, N. Hall, W. T. Sloan & C. Quince, 2015. Insight into biases and
834      sequencing errors for amplicon sequencing with the Illumina MiSeq platform. Nucleic Acids
835      Research 43(6):e37-e37 doi:10.1093/nar/gku1341.
836 Schulz, M. H., D. R. Zerbino, M. Vingron & E. Birney, 2012. Oases: robust *de novo* RNA-seq assembly
837      across the dynamic range of expression levels. Bioinformatics 28(8):1086-1092
838      doi:10.1093/bioinformatics/bts094.
839 Schurch, N. J., P. Schofield, M. Gierliński, C. Cole, A. Sherstnev, V. Singh, N. Wrobel, K. Gharbi, G. G.
840      Simpson, T. Owen-Hughes, M. Blaxter & G. J. Barton, 2016. How many biological replicates
841      are needed in an RNA-seq experiment and which differential expression tool should you
842      use? RNA 22(6):839-851 doi:10.1261/rna.053959.115.
843 Seelenfreund, E., W. A. Robinson, C. M. Amato, A.-C. Tan, J. Kim & S. E. Robinson, 2014. Long term
844      storage of dry versus frozen RNA for next generation molecular studies. PLoS ONE
845      9(11):e111827 doi:10.1371/journal.pone.0111827.
846 Simão, F. A., R. M. Waterhouse, P. Ioannidis, E. V. Kriventseva & E. M. Zdobnov, 2015. BUSCO:
847      Assessing genome assembly and annotation completeness with single-copy orthologs.
848      Bioinformatics 31(19):3210-3212 doi:10.1093/bioinformatics/btv351.
849 Soneson, C. & M. Delorenzi, 2013. A comparison of methods for differential expression analysis of
850      RNA-seq data. BMC Bioinformatics 14(1):91 doi:10.1186/1471-2105-14-91.
851 Song, L., C. Bian, Y. Luo, L. Wang, X. You, J. Li, Y. Qiu, X. Ma, Z. Zhu, L. Ma, Z. Wang, Y. Lei, J. Qiang, H.
852      Li, J. Yu, A. Wong, J. Xu, Q. Shi & P. Xu, 2016. Draft genome of the Chinese mitten crab,
853      *Eriocheir sinensis*. GigaScience 5:5 doi:10.1186/s13742-016-0112-y.
854 Sultan, M., V. Amstislavskiy, T. Risch, M. Schuette, S. Dökel, M. Ralser, D. Balzereit, H. Lehrach & M.-
855      L. Yaspo, 2014. Influence of RNA extraction methods and library selection schemes on RNA-
856      seq data. BMC Genomics 15(1):675 doi:10.1186/1471-2164-15-675.
857 Sultan, M., S. Dökel, V. Amstislavskiy, D. Wuttig, H. Sültmann, H. Lehrach & M.-L. Yaspo, 2012. A
858      simple strand-specific RNA-Seq library preparation protocol combining the Illumina TruSeq
859      RNA and the dUTP methods. Biochemical and Biophysical Research Communications
860      422(4):643-646 doi:10.1016/j.bbrc.2012.05.043.
861 Surget-Groba, Y. & J. I. Montoya-Burgos, 2010. Optimization of *de novo* transcriptome assembly
862      from next-generation sequencing data. Genome Research 20(10):1432-1440
863      doi:10.1101/gr.103846.109.
864 Tarazona, S., F. García-Alcalde, J. Dopazo, A. Ferrer & A. Conesa, 2011. Differential expression in
865      RNA-seq: A matter of depth. Genome Research 21(12):2213-2223
866      doi:10.1101/gr.124321.111.
867 Teng, M., M. I. Love, C. A. Davis, S. Djebali, A. Dobin, B. R. Graveley, S. Li, C. E. Mason, S. Olson, D.
868      Pervouchine, C. A. Sloan, X. Wei, L. Zhan & R. A. Irizarry, 2016. A benchmark for RNA-seq
869      quantification pipelines. Genome Biology 17(1):74 doi:10.1186/s13059-016-0940-1.
870 Wang, Y., N. Ghaffari, C. Johnson, U. Braga-Neto, H. Wang, R. Chen & H. Zhou, 2011. Evaluation of
871      the coverage and depth of transcriptome by RNA-Seq in chickens. BMC Bioinformatics 12
872      doi:10.1186/1471-2105-12-S10-S5.
873 Wang, Z., M. Gerstein & M. Snyder, 2009. RNA-Seq: a revolutionary tool for transcriptomics. Nature
874      Reviews Genetics 10(1):57-63 doi:10.1038/nrg2484.
875 Wilhelm, B. T. & J.-R. Landry, 2009. RNA-Seq—quantitative measurement of expression through
876      massively parallel RNA-sequencing. Methods 48(3):249-257.

877     Winnebeck, E. C., C. D. Millar & G. R. Warman, 2010. Why Does Insect RNA Look Degraded? Journal
878          of Insect Science 10:159 doi:10.1673/031.010.14119.
879     Wu, A. R., N. F. Neff, T. Kalisky, P. Dalerba, B. Treutlein, M. E. Rothenberg, F. M. Mburu, G. L.
880          Mantalas, S. Sim, M. F. Clarke & S. R. Quake, 2014. Quantitative assessment of single-cell
881          RNA-sequencing methods. Nature Methods 11(1):41-46 doi:10.1038/nmeth.2694.
882     Wu, T. D. & S. Nacu, 2010. Fast and SNP-tolerant detection of complex variants and splicing in short
883          reads. Bioinformatics 26(7):873-881 doi:10.1093/bioinformatics/btq057.
884     Xie, Y., G. Wu, J. Tang, R. Luo, J. Patterson, S. Liu, W. Huang, G. He, S. Gu, S. Li, X. Zhou, T.-W. Lam, Y.
885          Li, X. Xu, G. K.-S. Wong & J. Wang, 2014. SOAPdenovo-Trans: *De novo* transcriptome
886          assembly with short RNA-Seq reads. Bioinformatics doi:10.1093/bioinformatics/btu077.
887     Yang, C. & H. Wei, 2015. Designing Microarray and RNA-Seq experiments for greater systems biology
888          discovery in modern plant genomics. Molecular Plant 8(2):196-206
889          doi:10.1016/j.molp.2014.11.012.
890     Yu, Y., X. Zhang, J. Yuan, F. Li, X. Chen, Y. Zhao, L. Huang, H. Zheng & J. Xiang, 2015. Genome survey
891          and high-density genetic map construction provide genomic and genetic resources for the
892          Pacific White Shrimp *Litopenaeus vannamei*. Scientific Reports 5:15612
893          doi:10.1038/srep15612.
894     Yuan, J., X. Zhang, C. Liu, Y. Yu, J. Wei, F. Li & J. Xiang, 2017. Genomic resources and comparative
895          analyses of two economical penaeid shrimp species, *Marsupenaeus japonicus* and *Penaeus
896          monodon*. Marine Genomics doi:10.1016/j.margen.2017.12.006.
897     Zerbino, D. R. & E. Birney, 2008. Velvet: Algorithms for *de novo* short read assembly using de Bruijn
898          graphs. Genome Research 18(5):821-829 doi:10.1101/gr.074492.107.
899     Zhang, Z. H., D. J. Jhaveri, V. M. Marshall, D. C. Bauer, J. Edson, R. K. Narayanan, G. J. Robinson, A. E.
900          Lundberg, P. F. Bartlett, N. R. Wray & Q.-Y. Zhao, 2014. A comparative study of techniques
901          for differential expression analysis on RNA-Seq data. PLoS ONE 9(8):e103207
902          doi:10.1371/journal.pone.0103207.
903     Zhao, Q.-Y., Y. Wang, Y.-M. Kong, D. Luo, X. Li & P. Hao, 2011. Optimizing *de novo* transcriptome
904          assembly from short-read RNA-Seq data: a comparative study. BMC Bioinformatics 12
905          doi:10.1186/1471-2105-12-S14-S2.
906     Zhao, S., Y. Zhang, W. Gordon, J. Quan, H. Xi, S. Du, D. von Schack & B. Zhang, 2015. Comparison of
907          stranded and non-stranded RNA-seq transcriptome profiling and investigation of gene
908          overlap. BMC Genomics 16(1):675 doi:10.1186/s12864-015-1876-7.
909     Zhong, S., J. G. Joung, Y. Zheng, Y. R. Chen, B. Liu, Y. Shao, J. Z. Xiang, Z. Fei & J. J. Giovannoni, 2011.
910          High-throughput illumina strand-specific RNA sequencing library preparation. Cold Spring
911          Harbor Protocols 2011(8):940-9 doi:10.1101/pdb.prot5652.

912