

1   **Expressed Exome Capture Sequencing (EecSeq): a method for cost-effective exome  
2   sequencing for all organisms with or without genomic resources**

3

4   Jonathan B. Puritz<sup>1,2</sup>

5

6   Katie E Lotterhos<sup>1</sup>

7

8   <sup>1</sup>Department of Marine and Environmental Sciences, Northeastern Marine Science Center, 430  
9   Nahant Rd, Nahant, MA 01908

10

11   <sup>2</sup>Current address: Department of Biological Sciences, University of Rhode Island, 120 Flagg RD,  
12   Kingston, RI 02881

13

14

15

16

17   Running title: EecSeq: exome capture for non-model species

18

19

## 20 Abstract

21 Exome capture is an effective tool for surveying the genome for loci under selection. However,  
22 traditional methods require annotated genomic resources. Here, we present a method for creating  
23 cDNA probes from expressed mRNA, which are then used to enrich and capture genomic DNA  
24 for exon regions. This approach, called “EecSeq”, eliminates the need for costly probe design  
25 and synthesis. We tested EecSeq in the eastern oyster, *Crassostrea virginica*, using a controlled  
26 exposure experiment. Four adult oysters were heat shocked at 36° C for 1 hour along with four  
27 control oysters kept at 14° C. Stranded mRNA libraries were prepared for two individuals from  
28 each treatment and pooled. Half of the combined library was used for probe synthesis and half  
29 was sequenced to evaluate capture efficiency. Genomic DNA was extracted from all individuals,  
30 enriched via captured probes, and sequenced directly. We found that EecSeq had an average  
31 capture sensitivity of 86.8% across all known exons and had over 99.4% sensitivity for exons  
32 with detectable levels of expression in the mRNA library. For all mapped reads, over 47.9%  
33 mapped to exons and 37.0% mapped to expressed targets, which is similar to previously  
34 published exon capture studies. EecSeq displayed relatively even coverage within exons (i.e.  
35 minor "edge effects") and even coverage across exon GC content. We discovered 5,951 SNPs  
36 with a minimum average coverage of 80X, with 3,508 SNPs appearing in exonic regions. We  
37 show that EecSeq provides comparable, if not superior, specificity and capture efficiency  
38 compared to costly, traditional methods.

39

40 **Keywords:** exome capture, population genomics, selection

41

## 42 Introduction

43 The invention of next-generation sequencing has made it possible to obtain massive amounts of  
44 sequence data. These data have given insight into classical problems in evolutionary biology,  
45 including the repeatability of evolution (e.g., Jones *et al.* 2012), the degree of convergent  
46 evolution across distant taxa (e.g., Yeaman *et al.* 2016), and whether selection is driving changes  
47 in existing genetic variation or new mutations (e.g., Reid *et al.* 2016). Despite this rapid progress,  
48 it is still cost prohibitive to sequence dozens or hundreds of full genomes. This limits our ability  
49 to study the genomic basis of local adaptation, which requires large sample sizes for statistical  
50 power (De Mita *et al.* 2013; Lotterhos & Whitlock 2015; Hoban *et al.* 2016). This leads to an  
51 inherent trade-off between sample size and genomic coverage, leading investigators to make  
52 decisions about whether to sequence more individuals (for higher power and precision) versus  
53 more of the genome (for making more accurate statements about the genetic basis of adaptation).

54 Reduced representation library preparation methods offer various kinds of random or targeted  
55 genome reduction, but the available approaches have contrasting advantages and limitations.  
56 RADseq uses restriction enzymes to randomly sample the genome and is appropriate for linkage  
57 mapping and studying neutral processes like gene flow and drift (Puritz *et al.* 2014), but the data  
58 can be limited for understanding the genetic basis of adaptation (Lowry *et al.* 2016, 2017;  
59 Catchen *et al.* 2017; McKinney *et al.* 2017). To focus on coding regions, some investigators have  
60 used RNAseq (De Wit *et al.* 2015); however, only about a dozen individuals can be sequenced  
61 per lane because of log-fold differences in transcript abundance among loci. Additionally, allele-  
62 specific expression limits the confidence in genotypes derived from RNAseq data (Pastinen

63 2010), especially in pooled samples. One increasingly popular option for increasing precision  
64 with larger samples while still maintaining coverage of the entire genome is Pool-seq, which  
65 sequences every individual to very low (1x) coverage and uses the data to calculate allele  
66 frequency of the sample within the pool (Buerkle & Gompert 2013; Schlötterer *et al.* 2014;  
67 Therkildsen & Palumbi 2017). Pool-seq is limited to only estimating allele frequency within  
68 pools, which is a disadvantage because this data cannot be used to understand the fitness of  
69 heterozygotes and some types of statistical analyses would be impossible to perform, such as  
70 haplotype-based analyses (e.g. Fariello *et al.* 2013).

71 To overcome some of these limitations, many investigators have used capture approaches with  
72 biotinylated probes (Jones & Good 2016). Capture approaches have the advantage of enriching  
73 the data for sequences of interest - allowing for individual-level data and a large number of  
74 individuals to be sequenced - but require the investigator to have genomic resources for probe  
75 design and then to purchase the probes from a company. For non-model species, the development  
76 of these resources takes time and a significant amount of bioinformatics expertise. In addition, for  
77 a population-level genomic study with 100s of individuals, probes may cost several tens of  
78 thousands of dollars, depending on how much sequence is captured. Overall, what is needed is a  
79 cost-effective approach to subsample genomes for coding regions, without previously developed  
80 genomic resources. Such an approach would allow for the assessment of rapid adaptation to  
81 environmental disasters such as Deepwater Horizon Oil Spill (Lee *et al.* 2017).

82 Here, we present a novel, cost-effective method of exome capture that synthesizes probes in-situ  
83 from expressed mRNA sequences. Expressed Exome Capture Sequencing (EecSeq) builds upon  
84 existing approaches for in-situ probe synthesis that rely on restriction enzymes to sample the

85 genome or exome (Suchan *et al.* 2016; Schmid *et al.* 2017). To improve capture efficiency, we  
86 developed a novel library preparation procedure that uses standardized procedures to synthesize  
87 cDNA from expressed RNA (without template reduction via restriction digest) and then create  
88 biotinylated probes from cDNA (see **Figure 1** for a conceptual diagram). The EecSeq design  
89 includes custom RNA library adapters that offer several major advantages. The custom adapters  
90 are fully compatible with duplex-specific nuclease normalization, which is included in the  
91 protocol in order to reduce log fold differences in expression - resulting in more even coverage  
92 across high- and low-expressed transcripts. The custom adapters also allow for probe sequencing  
93 - before normalization if differential expression data is desired, or after normalization if probe  
94 abundance data is desired. Moreover, the adapters are easily removed with a single enzymatic  
95 treatment before biotinylation, preventing any interference during hybridization.  
  
96 Our approach is cost-effective and does not require any prior genomic resources, making it a  
97 good choice for studies seeking to understand adaptation in exomes. The approach, however, is  
98 limited in the sense that the probes are designed from expressed RNA, and so investigators  
99 should be careful to choose which tissues and life stages would be relevant. Here, we show  
100 proof-of-concept of the approach in the eastern oyster (*Crassostrea virginica*), and find that the  
101 performance of the approach is comparable, if not superior, to the performance of published  
102 exome capture datasets where probes were designed from sequence data and purchased from a  
103 company.

## 104 Methods

### 105 Experimental overview

106 Expressed exome capture sequencing (EecSeq) is designed with two specific goals: 1) to  
107 eliminate the need for expensive exome capture probe design and synthesis and 2) to focus exon  
108 enrichment of genes that are being expressed relevant to tissue(s) and condition(s) of interest. To  
109 illustrate this conceptually, we exposed adult oysters to a stressor (extreme heat) that would  
110 generate a predictable gene and protein expression profile (expression of heat shock proteins).  
111 Having a predictable coverage profile in the probes allowed us to evaluate whether the genomic  
112 DNA in these exons were captured by the probes. Note, however, that this experiment is not  
113 specifically part of the EecSeq method and that the investigator can choose appropriate tissue(s)  
114 and condition(s) of interest. The steps to probe synthesis and capture are visualized in Figure 1.

115 **Heat shock exposure, tissue collection, and nucleic acid extraction**

116 Eight adult *Crassostrea virginica* individuals were collected and acclimated to a flow-through  
117 seawater system for 24 hours. After acclimation, individuals were randomly assigned to two  
118 treatments, control and heat-shock (HS). HS individuals were placed a small aquaria filled with  
119 36°C filtered seawater for one hour while control individuals were kept in an identical aquarium  
120 filled with 14°C (ambient) filtered seawater. Immediately after the exposure period, all  
121 individuals were shucked and mantle tissue was extracted and frozen in liquid nitrogen in  
122 duplicate. DNA was extracted using the DNeasy kit (Qiagen) and RNA was extracted using TRI  
123 Reagent Solution (Applied Biosystems) using included, standard protocols. DNA was visualized  
124 on an agarose gel and quantified using the Qubit DNA Broad Range kit (Invitrogen). RNA was  
125 visualized on an Agilent BioAnalyzer using the RNA 6000 Nano kit, and was quantified using  
126 the Qubit High Sensitivity Assay Kit (Invitrogen).

127 **Expressed Exome Capture Sequencing**

128 A complete and updated EecSeq protocol can be found at (<https://github.com/jpuritz/EecSeq>).

129 *RNA Adapters*- Custom RNA adapters were used in this protocol. The RNA adapters were similar  
130 to the Illumina TruSeq design, but include the SAI1 restriction site at the 3' end of the "Universal  
131 adapter" and at 5' end of the "Indexed adapter." The presence of this restriction site allows the  
132 Illumina sequence to be removed before hybridization to prevent interference. Note that the  
133 adapters used in this study had an erroneous deletion of a Thymine in position 58 of  
134 "Universal\_SAI1\_Adapter" and in position 8 of all four indexed adapters (the corrected versions  
135 are shown in Table 1, and erroneous version used in this study are shown in Supplemental Table  
136 1). Adapters were annealed in equal parts in a solution of Tris-HCl (pH 8.0), NaCl, and EDTA,  
137 heated to 97.5°C for 2.5 minutes, and then cooled at a rate of 3°C per minute until the solution  
138 reaches a temperature of 21°C.

139 *mRNA Library Preparation and Normalization*- Probes were made from two (of four) control  
140 individuals and two (of four) exposed individuals. The first step for this subset of individuals was  
141 to prepare stranded mRNA libraries using the Kapa Stranded mRNA-Seq Kit (KAPA  
142 Biosystems) with the following modifications: custom adapters were used, 4 micrograms of RNA  
143 per individual were used as starting material, half volume reactions were used for all steps,  
144 adapters were used at a final reaction concentration of 50 nM during ligation, and 12 cycles of  
145 PCR were used for enrichment. Complete libraries were visualized on a BioAnalyzer using the  
146 DNA 1000 kit, quantified using fluorometry, and then 125 ng of each library was taken and  
147 pooled to single library of 500 ng.

148 To reduce the abundance of highly expressed transcripts in our final probe set, complete libraries  
149 were normalized following Illumina's standard protocol for DSN normalization. First, the cDNA

150 library was heat denatured and slowly allowed to reanneal. Next, the library was treated with  
151 duplex-specific nuclease (DSN), which will remove abundant DNA molecules that have properly  
152 annealed. After DSN treatment, the library was SPRI purified and enriched via 12 cycles of PCR.  
153 A subsample of probes was exposed to an additional 12 cycles of PCR to test for PCR artifacts in  
154 probe synthesis. The normalized cDNA library was visualized on a BioAnalyzer using the DNA  
155 1000 kit, quantified with a Qubit DNA Broad Range kit (Invitrogen), and then split into two  
156 equal volume tubes, one to be saved for sequencing and one for probe synthesis. The DNS-  
157 normalized libraries were sequenced on one half lane of HiSeq 4000 by GENEWIZ  
158 ([www.genewiz.com](http://www.genewiz.com)).

159 *Probe Synthesis*-To remove the sequencing adapters, the cDNA library was treated with 100 units  
160 of Sall-HF restriction enzyme (New England Biolabs) in a total volume of 40  $\mu$ l at 37°C for 16  
161 hours. After digestion, the digested library was kept in the same tube, and 4.5  $\mu$ l of 10X Mung  
162 Bean Nuclease Buffer and 5 units of Mung Bean Nuclease (New England Biolabs) were added.  
163 The reaction was then incubated at 30°C for 30 minutes. An SPRI cleanup using AMPure XP  
164 (Agencourt) was completed with an initial ratio of 1.8X. After, visualization of the library on an  
165 Agilent BioAnalyzer, a subsequent SPRI cleanup of 1.5X was completed to remove all digested  
166 adapters. The clean, digested cDNA fragments were then biotin labeled using the DecaLabel  
167 Biotin DNA labeling kit (Thermo Scientific) using the included protocol. The labeling reaction  
168 was then cleaned using a 1.5X SPRI cleanup and fluorometrically quantified. To test the effects  
169 of additional PCR cycles on probe effectiveness, 40 ng of the original, normalized cDNA library  
170 was subjected to an additional 12 cycles of PCR, and then converted to probes as described  
171 above.

172     *Genomic DNA Library Preparation-* Capture was performed on a standard genomic DNA library.  
173     500 ng of genomic DNA from all eight individuals was sheared to a modal peak of 150 base pairs  
174     using a Covaris M220 Focused-ultrasonicator. The sheared DNA was inserted directly into step  
175     2.1 of the KAPA HyperPlus kit with the following modifications: half reaction volumes were  
176     used, and a final adapter:insert molar ratio of 50:1 was used with custom TruSeq-style, barcoded  
177     adapters (note: the adapters contained erroneous mismatches in the barcodes between the top and  
178     bottom oligos; the original oligonucleotide sequences can be found in Supplemental Table 2 and  
179     corrected versions in Supplemental Table 3). After adapter ligation, individuals were pooled into  
180     one single library, and libraries were enriched with 6 cycles of PCR using primers that  
181     complemented the Illumina P5 adapter and Indexed P7 (Supplemental Table 2). The final library  
182     was quantified fluorometrically quantified and analyzed on an Agilent BioAnalyzer.

183     *Hybridization-* Three replicate captures were performed using the set of original probes and the  
184     set of probes with 12 extra cycles of PCR. The hybridization protocol closely followed that of  
185     Suchan *et al.* (2016). 500 ng of probes and 500 ng of genomic DNA library were hybridized  
186     along with blocking oligonucleotides (Table 2) at a final concentration of 20 µM in a solution of  
187     6X SSC, 5 mM EDTA, 0.1% SDS, 2X Denhardt's solution, and 500 ng c<sub>0</sub>t-1 DNA. The  
188     hybridization mixture was incubated at 95°C for 10 minutes, and then 65°C for 48 hours in a  
189     thermocycler. The solution was gently vortexed every few hours.

190     *Exome Capture-* 40 µl of hybridization mixed was added to 200 µl of DynaBeads M-280  
191     Streptavidin beads (Thermo Fisher Scientific). The beads and hybridization mixture were then  
192     incubated for 30 min at room temperature. The mixture was then placed on a magnetic stand  
193     until clear, and the supernatant was removed. This was followed by four bead washes under

194 slightly different conditions. First, the beads were washed with 200 µl 1X SSC and 0.1% SSC  
195 solution, incubated at 65°C for 15 min, placed on the magnet stand, and the supernatant was  
196 removed. Second, the beads were washed with 200 µl 1X SSC and 0.1% SSC solution incubated  
197 at 65°C for 10 minutes, placed on the magnet stand, and the supernatant was removed. Third, the  
198 beads were washed with 200 µl 0.5 SSX and 0.1% SDS solution, incubated at 65°C for 10  
199 minutes, placed on the magnet stand, and the supernatant was removed. Finally, the beads were  
200 washed with 200 µl 0.1X SSC and 0.1% SDS, incubated at 65°C for 10 minutes, placed on the  
201 magnet stand, and the supernatant was removed. Lastly, DNA was eluted from the beads in 22 µl  
202 of molecular grade water heated to 80°C for 10 minutes. The solution was placed on the magnet  
203 and the supernatant was saved. The hybridized fragments were then enriched with 12 cycles of  
204 PCR using the appropriate P5 and P7 PCR primers and cleaned with 1X AMPure XP with a final  
205 elution in 10 mM Tris-HCl (pH 8.0). The six replicate captures, each containing 8 uniquely  
206 barcoded individuals, were sequenced on one half lane (separate from the RNA libraries) on the  
207 HiSeq 4000 platform by GENEWIZ ([www.genewiz.com](http://www.genewiz.com)).

## 208 Bioinformatic Analysis

209 All bioinformatic code, including custom scripts and a script to repeat all analyses, can be found  
210 at (<https://github.com/jpuritz/EecSeq/tree/master/Bioinformatics>)

211 *RNA libraries*- RNA reads were first trimmed for quality and custom adapter sequences were  
212 searched for with Trimmomatic (Bolger *et al.* 2014) as implemented in the dDocent pipeline  
213 (version 2.2.20; Puritz *et al.* 2014). Reads were then aligned to release 3.0 of the *Crassostrea*  
214 *virginica* genome (Accession: GCA\_002022765.4) using the program STAR (Dobin *et al.* 2013).  
215 The genome index was created using NCBI gene annotations for splice junctions. Reads were

216 aligned in a two-step process, first using the splice junctions in the genome index, and then again  
217 using both the splice junctions in the index and additional splice junctions found during the first  
218 alignment. Alignment files from the four libraries were then merged with SAMtools (version 1.4;  
219 Li *et al.* 2009) and filtered for MAPQ > 4, only primary alignments, and reads that were hard/soft  
220 clipped at less than 75 bp. SAMtools (Li *et al.* 2009) and Bedtools (Quinlan 2014) were used to  
221 calculate read and per bp coverage levels for exons, introns, and intergenic regions.

222 *EecSeq Libraries-* Raw reads were first trimmed using the standard methods in the dDocent  
223 pipeline (version 2.2.20; Puritz *et al.* 2014). The DNA adapters contained erroneous mismatches  
224 between the top and bottom oligos in the barcode (original oligonucleotide sequences can be  
225 found in Supplemental Table 2 and corrected versions in Supplemental Table 3). These  
226 differences prevented demultiplexing beyond the capture pool level, and also lead to potentially  
227 erroneous base calls within the first 7 bp of sequencing. To remove these artifacts, the first 7 bp  
228 of every forward read were clipped. Additionally, adapter sequences were searched for in the  
229 paired-end sequences using custom scripts. After trimming, reads were aligned to the reference  
230 genome using BWA (Li & Durbin 2009) with the mismatch parameter lowered from 4 to 3, and  
231 the gap opening penalty lowered from 6 to 5. PCR duplicates were marked using the  
232 *MarkDuplicatesWithMateCigar* module of Picard (<http://broadinstitute.github.io/picard>), and  
233 then SAMtools (Li *et al.* 2009) was used to remove duplicates, secondary alignments, mappings  
234 with a quality score less than ten, and reads with more than 80 bp clipped. SAMtools (Li *et al.*  
235 2009) and Bedtools (Quinlan 2014) were used to calculate read and per bp coverage levels for  
236 exons, introns, and intergenic regions. FreeBayes (Garrison and Marth 2012) was used to call  
237 SNPs.

238 *Calculating Capture Efficiency-* EecSeq is unique amongst exome capture methods because the  
239 probes are not designed directly, i.e. there is no set of *a priori* targets. Additionally, EecSeq is  
240 designed to capture exons that are expressed in the samples used to create probes - not the entire  
241 exome. To compare EecSeq to other capture methods, capture targets were defined as exons that  
242 had more than 35X coverage in the RNAseq (probe) data and confidence intervals were generated  
243 by defining capture targets as 20X RNAseq coverage and 50X RNAseq coverage. We also  
244 calculated a conservative, near-target range of 150 bp on either side of the defined targets. This  
245 range corresponds to the modal DNA fragment length used for the capture libraries with the  
246 expectation that exon probes could capture reads that far from the original target.

## 247 Results

248 *RNA sequencing results-* RNA sequencing, filtering, and mapping statistics can be found in  
249 supplemental Table 3. After filtering, a total of 21,990,025 RNA reads were mapped uniquely to  
250 the eastern oyster genome. Of the total RNA reads, 78% mapped to genic regions of the genome,  
251 and 58% mapped to annotated exon regions. Across all exonic bases in the genome, less than 5%  
252 had more than 50X coverage; however, over 16% had at least 20X coverage and over 45% had at  
253 least 5X coverage (Figure 2).

254 *Exome capture sequencing results-* Six replicate capture pools of the same eight individuals were  
255 sequenced on half a lane of Illumina HiSeq (3 replicates from probes that had been enriched via  
256 12 cycles of PCR and 3 replicates from probes that had been enriched via 24 cycles of PCR). A  
257 summary of exome capture sequencing, filtering, and mapping statistics are shown in Table 2.  
258 On average, there were 47,629,033 raw reads (forward and paired-end) per capture pool and an

259 average of 32,123,268 mapped reads per capture pool after filtration. Across the entire oyster  
260 genome, RNA sequencing coverage and exome sequencing coverage was highly correlated  
261 (Supplemental Figure 1), and across all exon regions total RNA coverage predicted 72.6% of the  
262 variation in exome capture coverage (Figure 3; log-log transformation,  $R^2 = 0.72619$ ,  $p < 0.0001$ ).  
263 Coverage across all exons and expressed exon targets was highly correlated ( $0.984 < r < 0.996$ )  
264 across all replicate captures, and the average capture of pools with standard probes and the  
265 average capture of pools with probes with extra PCR was virtually identical ( $R^2 = 99.1$ ;  $p <$   
266 0.0001).

267

268 *Exome capture efficiency*- Capture sensitivity, or the percentage of targets covered by at least one  
269 read (1X), was high across all replicate pools, regardless of target set (Table 3). Across all  
270 known exons, sensitivity was on average 86.8% across replicate capture pools, and across all  
271 defined target sets, sensitivity was over 99.4%. Increasing the sensitivity threshold from 1X to  
272 10X lowers the sensitivity across all exons but has little effect on sensitivity across defined target  
273 sets (Supplemental Table 4). Sensitivity can also be measured at the per bp level instead of per  
274 exon. The percent of target bases captured is shown as a function of sensitivity threshold (read  
275 depth of capture libraries) in Figure 4.

276 Capture specificity is the percentage of mapped reads that fall within target regions. Across all  
277 exons, capture pools averaged 47.9% reads on target, 6.8% of reads near target (falling within  
278 150 bp of an exon, one modal read length), and 45.3% of reads off-target (more than 150 bp away  
279 from an exon). Across defined expressed exon targets (exons that sequenced to 35x read depth),

280 capture pools averaged 37.1% (C.I. 33.6% - 41.4%) reads on target, 3.55% (C.I. 3.0% - 4.4%) of  
281 reads near target, and 59.38% (C.I. 54.2% - 63.4%) reads off target.

282 For all exons, between the 10th and 90th percentile of exon length (59bp - 517bp), the mean per  
283 basepair coverage averaged  $17.75X \pm 0.06X$  for each capture pool of 8 individuals. When  
284 considering target exons (35X coverage in RNA-derived probes), the mean per basepair coverage  
285 increased to  $61.22X \pm 0.23X$  on average for each capture pool. This breaks down to  
286 approximately 7.66 reads on average per individual per bp within expressed exome targets.

287 Within exons, mean per basepair coverage was evenly distributed across all base pairs with only  
288 slightly lower coverage at the 5' or 3' edges of exons compared to the middle of exons (Figure 5;  
289 Supplemental Figure 3).

290 Mean capture coverage also did not appear to relate to the GC content of the target exon (Figure  
291 6), though it did appear to peak near the mean GC content of 43.57%. To test this, we calculated  
292 the reciprocal of the absolute value of the difference between each exon GC content and the  
293 average GC content, and then tested for a linear relationship to mean coverage. Though we found  
294 this relationship to be significant ( $p > 0.0008$ ), it explained only the 0.0033% of the variance in  
295 coverage, confirming that exon GC content did not affect exon capture in a meaningful way.

296 Coverage did vary significantly between untranslated regions (UTR) within exons and coding  
297 sequence (CDS) within exons (Welch's test  $t = 40.063$ ; degrees of freedom = 135580;  $p <$   
298 0.0001) with a mean coverage for UTR equaling  $11.59X \pm 0.0864$  and a mean coverage for  
299 CDS equaling  $17.71X \pm 0.1261$ . This small but significant coverage difference was also  
300 evident as the percent of target bases greater than a given read depth (Supplemental Figure 2).  
301 This pattern was not surprising, however, because the same pattern was observed for the RNA

302 reads (CDS mean coverage = 13.65X +/- 0.2011; UTR mean coverage = 8.25 +/- 0.1275;  
303 Welch's test  $t = 22.677$ ; degrees of freedom = 129300;  $p < 0.0001$ ), indicating that the probes  
304 also had lower coverage in UTR compared to CDS.

305 *Expressed exon capture*- To visualize the relationship between coverage and an expected  
306 expressed target, we plotted coverage of the six capture pools along two heat shock proteins, Heat  
307 Shock cognate 71 kDa (NCBI Reference Sequence: XM\_022472393.1, Figure 7) and Heat Shock  
308 70 kDa protein 12B-like (NCBI Reference Sequence: XM\_022468697.1; Supplemental Figure 4).  
309 As expected, exons in both genes show elevated coverage that corresponded to the coverage of  
310 the mRNA-derived probes, especially along regions with corresponding CDS with few reads  
311 mapping to intronic or intergenic regions.

312 *SNP Discovery*- A total of 1,011,107 raw SNPs were discovered with 909,792 SNPs having a  
313 quality score higher than 20. A total of 99,169 high quality SNPs were found within known  
314 exons. Of these, 31,579 exome SNPs had at least an average of 16X coverage, 15,760 exome  
315 SNPs had at least an average of 32X coverage, 8,837 exome SNPs had at least an average of 48X  
316 coverage, and 3,508 exome SNPs had at least an average of 80X coverage with an additional  
317 2,443 80X-SNPs found outside of exon regions.

## 318 Discussion

319 Expressed exome capture sequencing (EecSeq) is a novel design for exome capture that uses *in-*  
320 *situ* synthesized biotinylated cDNA probes to enrich for exon sequences, thereby removing the  
321 requirement of *a priori* genomic resources, costly exon probe design, and synthesis. Here, we  
322 showed that EecSeq target enrichment had high levels of sensitivity, with comparable if not

323 superior performance and specificity to traditional methods. EecSeq exon enrichment showed  
324 even coverage levels with exons and across exons with differing levels of GC content. Lastly, we  
325 showed that EecSeq can quickly and cheaply generate thousands of exon SNPs.

326 **Benefits of EecSeq**

327 *Diverse probes-* With EecSeq, cDNA exon probes are constructed *in-situ* from extracted mRNA,  
328 and this allows for the design of a high-diversity probe pool. Traditional sequence capture probes  
329 are typically designed from a single reference genome or individual, and this may limit capture  
330 efficiency on individuals with different SNPs, insertions, or deletions than the reference. While  
331 probes been successfully used to capture sequences in quite divergent species (less than 5%  
332 sequence divergence, Jones & Good 2016), there is evidence that capture success declines as  
333 sequences become less related to the reference. Portik *et al.* (2016) found that for each percent  
334 increase of pairwise divergence, missing data increased 4.76%, sensitivity decreased 4.57%, and  
335 specificity decreased 3.26%. Even with well-designed, commercially available capture kits for  
336 human exon capture, Sulonen *et al.* (2011) found that allele balances for heterozygous variants  
337 tended to have more reference bases than variant bases in the heterozygous variant position  
338 across all methods for probe development. Insertions and deletions (InDels) are arguably an even  
339 larger problem, since these would decrease hybridization with a probe due to a frameshift.

340 *Longer Probes-* Traditional exome capture relies on synthesized RNA or DNA baits. These baits  
341 can be relatively small (60 bp; Bi *et al.* 2012) or range between 95 and 120 bp (Clark *et al.* 2011;  
342 Sulonen *et al.* 2011; Nadeau *et al.* 2012; Chilamakuri *et al.* 2014). In contrast, EecSeq probes  
343 have a modal length of 150 bp but also range up to over 400 bp (data not shown). The longer  
344 length of EecSeq probes likely helps to buffer against divergence between probes and targets.

345 The longer probes may also be the reason why we observed relatively little GC bias in coverage  
346 across exons, and may help explain the uniformity of coverage within exons in EecSeq data.

347 *Cost-* EecSeq provides significant cost and time savings over traditional exome capture and RNA  
348 sequencing (RNAseq). No *a priori* genomic information is necessary for EecSeq, saving  
349 substantial time and money for obtaining these data in non-model organisms. Likewise, the cost  
350 of synthesizing the probes is significantly reduced because probes can be made in-house and do  
351 not have to be designed by a company. On a per sample basis, EecSeq is also significantly  
352 cheaper than RNAseq because (i) commercial DNA library preps are cheaper than those for  
353 mRNA, and (ii) more individuals can be multiplexed on a single lane. For example, the cost of  
354 RNA seq is \$246 per sample (cost estimated using the same RNA kits used with EecSeq and ½  
355 reactions) and assuming that 12 RNAseq libraries can be sequenced in a single lane of Illumina  
356 HiSeq, the cost per sample is (\$1,008; cost of the kit; Kapa Biosystems Stranded mRNA-Seq Kit  
357 with 24 reactions or 48 half-reactions)\*(1/48; the amount used per sample) + \$2700/12 = \$246  
358 per sample). The equivalent cost per sample for EecSeq is \$48.02 per sample (for 96 samples in  
359 one lane of HiSeq; Supplemental Table 6) or \$62.08 per sample if a more conservative  
360 sequencing strategy is used (96 samples sequenced over 1.5 lanes of HiSeq; Supplemental Table  
361 6).

362 *No dependency on restriction sites-* A recently published method, hyRAD-X, (Schmid *et al.*  
363 2017) is similar to EecSeq in that it uses *in-situ* synthesized cDNA probes from expressed mRNA  
364 to capture exome sequences. However, the protocol relies on a restriction digest to fragment  
365 cDNA and ligate on probes. This may result in a reduced template of probes because not all  
366 cDNA fragments will have restriction sites on both ends. To evaluate the possibility that the

367 hyRAD-X would produce a reduced template of probes, we performed crude calculations using  
368 SimRAD in R (Lepais & Weir 2014) on the *C. virginica* exome. Of the 31,383 known mRNA  
369 transcripts in the oyster genome (assuming 1 transcript variant), 29,555 contain at least 2 MseI  
370 cut sites (TTAA). However, there is an SPRI cleanup on the digestion (2X), meaning that at best,  
371 only fragments 100bp and larger are getting through to biotinylation  
372 (<http://www.keatslab.org/blog/pcrpurificationampureandsimple>). SimRAD estimates 220,184 out  
373 of a possible 440,881 fragments. Therefore, at the absolute best hyRAD-X is only sampling  
374  $(29,555/31,383)*(220,184/440,881) = 47\%$  of the exome, though this number may increase  
375 slightly due to transcript variations. Relying on restriction digests may also produce skewed size  
376 distributions in probes which would be magnified in subsequent rounds of PCR. In Schmid *et al.*  
377 (2017), hyRAD-X generated 524 exome SNPs at a minimum of 6X coverage across 27 samples  
378 (compared to the 3,508 exome SNPs discovered at 80X coverage derived from only 8 effective  
379 samples in 6 replicate capture using EecSeq), but they were also studying ancient DNA and so  
380 whether the hyRAD-X protocol results in limited coverage across exons remains to be tested.

381 **Caveats of EecSeq**

382 Despite the demonstrated benefits of EecSeq, there are some potential caveats that should be  
383 considered before employing the method. First, there is no ability to filter out probes that belong  
384 to repetitive sequences, which are often present at high concentrations in large-genome organisms  
385 such as amphibians (Keinath *et al.* 2015) or conifers (De La Torre *et al.* 2014). In one capture  
386 study from designed probes, a small proportion of the probes (unknowingly at the time of probe  
387 development) matched highly repetitive sequences (Syring *et al.* 2016). This resulted in an  
388 inordinate number of reads to these few probe sequences (Syring *et al.* 2016). However, the

389 inclusion of known repetitive sequence blocker in hybridization, such as c<sub>0</sub>t-1 that is used in the  
390 EecSeq protocol, has been shown to nearly double capture efficiency (McCartney-Melstad *et al.*  
391 2016). In general, repetitive elements, short repeats, and low complexity regions are problematic  
392 for all types of probe design and capture.

393 Another caveat of using EecSeq is the need to obtain RNA from relevant samples, although  
394 capture designs or gene expression studies based on transcriptomes face the same challenge.  
395 Note, however, the advantage that EecSeq probes can be made from mRNA pooled from many  
396 individuals, tissues, and conditions of interest. If genes of interest are expressed in tissues that are  
397 difficult to dissect or are in small abundances (such as neurons), then the RNA-based methods  
398 presented here would not be a feasible approach unless pooling multiple extractions.  
399 Additionally, the probes are a limited resource - our results indicate, however, that additional  
400 rounds of PCR on the probes have little effect on capture.

401 **Unique Aspects of EecSeq**

402 Our approach relies on expressed mRNA for probe synthesis and the abundance of particular  
403 mRNAs will vary depending on gene expression. EecSeq includes a normalization step to  
404 decrease the abundance of very common transcripts, but probe pools will still skew towards  
405 highly expressed genes and therefore capture coverage will be higher for those exons. This  
406 aspect of EecSeq can be customized for particular research questions. For projects focused on  
407 total exome capture, pools from multiple individuals, tissue types, and environmental/laboratory  
408 exposures can be constructed to generate a robust probe set. On the contrary, if an investigator is  
409 focused on a subset of genes that are responding to a particular stressor, it is possible to make  
410 probes from organisms exposed that specific condition and then use those probes to capture other

411 individuals. This reduced probe set may also allow for greater multiplexing, but this remains to  
412 be specifically tested. While we have only used mRNA to create probes, there may be  
413 possibilities to capture other types transcribed sequences such as long non-coding RNAs or  
414 possibly even miRNA.

415 Previous work on exome capture probe design has focused on intron/exon boundaries. In  
416 general, it is thought that capture probes that span exon boundaries will result in low coverage of  
417 these regions (Jones & Good 2016) or that certain regions will not be covered at all (Neves *et al.*  
418 2013). Inclusion of too many boundaries may also lower overall capture performance by  
419 increasing off-target capture (Suren *et al.* 2016). EecSeq exome probes are derived from mature  
420 RNA, so some of the probes will span inevitably span exon boundaries. Though exon/intron  
421 boundaries cannot be eliminated in EecSeq, both input mRNA and genomic DNA were  
422 fragmented down to a modal size of 150 base pairs, with the intention of making both smaller  
423 than the average exon size of Eastern Oysters (note that this size is at the lower limit of what is  
424 possible with Illumina sequencing). We found that coverage within exons was fairly uniform,  
425 indicating a lack of "edge effects." We hypothesize that the relative long length of EecSeq probes  
426 (compared to commercially synthesized probes), the near matching length of genomic DNA  
427 fragments, and the length distribution relative to actual exon size helped to ensure uniform exon  
428 coverage.

429 We compared our observed measures of sensitivity and specificity to other recently published  
430 studies in non-model species where probes were designed from bioinformatic resources for the  
431 same species. EecSeq capture efficiency performed as well as or outperformed almost all other  
432 previously published exome capture studies in non-model species (excluding mice and humans;

433 Table 4) with the notable exception of black cottonwood (Zhou & Holliday 2012), a species with  
434 exceptional genomic resources. Note, however, that we analyzed capture efficiency across pools  
435 of 8 individuals, and there could be considerable variability at the individual level that remains to  
436 be quantified.

#### 437 **Conclusions and Future Directions**

438 Here, we have shown that EecSeq effectively targets expressed exons, delivers consistent and  
439 efficient exome enrichment that is comparable to traditional methods of exome capture, and  
440 generates thousands of exome-derived SNPs cost effectively. Additional tests are needed to  
441 examine the efficiency of exome capture across individuals for different species, which should be  
442 coupled with sequencing of EecSeq probes to investigate the effects of probe pool diversity and  
443 sequence divergence between probes and targets on capture. Nonetheless, EecSeq holds  
444 substantial promise as a universally applicable and cost-effective method of exome sequencing  
445 for virtually any macroscopic organism.

#### 446 **Acknowledgements**

447 The authors would like to thank Alan Downey-Wall and Sara Schaal for informative discussions  
448 and input throughout the duration of this project. The authors also thank Nadir Alvarez for  
449 thoughtful comments on a preprint of this manuscript. This work was funded with funds  
450 provided to KEL from Northeastern University and NSF OCE-1635423.

451

452

## 453 References

- 454 Buerkle CA, Gompert Z (2013) Population genomics based on low coverage sequencing: how low should  
455 we go? *Molecular Ecology*, **22**, 3028–3035.
- 456 Bi K, Vanderpool D, Singhal S *et al.* (2012) Transcriptome-based exon capture enables highly cost-  
457 effective comparative genomic data collection at moderate evolutionary scales. *BMC genomics*, **13**,  
458 403.
- 459 Bolger AM, Lohse M, Usadel B (2014) Trimmomatic: a flexible trimmer for Illumina sequence data.  
460 *Bioinformatics*, **30**, 2114–2120.
- 461 Catchen JM, Hohenlohe PA, Bernatchez L *et al.* (2017) Unbroken: RADseq remains a powerful tool for  
462 understanding the genetics of adaptation in natural populations. *Molecular ecology resources*, **17**,  
463 362–365.
- 464 Chilamakuri CSR, Lorenz S, Madoui M-A *et al.* (2014) Performance comparison of four exome capture  
465 systems for deep sequencing. *BMC genomics*, **15**, 449.
- 466 Christmas MJ, Biffin E, Breed MF, Lowe AJ (2017) Targeted capture to assess neutral genomic variation  
467 in the narrow-leaf hopbush across a continental biodiversity refugium. *Scientific reports*, **7**, 41367.
- 468 De Mita S, Thuillet A-C, Gay L *et al.* (2013) Detecting selection along environmental gradients: analysis  
469 of eight methods and their effectiveness for outbreeding and selfing populations. *Molecular ecology*,  
470 **22**, 1383–1399.
- 471 De Wit P, Pespeni MH, Palumbi SR (2015) SNP genotyping and population genomics from expressed  
472 sequences - current advances and future possibilities. *Molecular ecology*, **24**, 2310–2323.
- 473 Dobin A, Davis CA, Schlesinger F *et al.* (2013) STAR: ultrafast universal RNA-seq aligner.  
474 *Bioinformatics*, **29**, 15–21.
- 475 Garrison E, & Marth G (2012) Haplotype-based variant detection from short-read sequencing. *arXiv*,  
476 arXiv:1207.3907
- 477 Fariello MI, Boitard S, Naya H, SanCristobal M, Servin B (2013) Detecting signatures of selection  
478 through haplotype differentiation among hierarchically structured populations. *Genetics*, **193**, 929–  
479 941.
- 480 Hebert FO, Renaut S, Bernatchez, L (2013) Targeted sequence capture and resequencing implies a  
481 predominant role of regulatory regions in the divergence of a sympatric lake whitefish species pair  
482 (*Coregonus clupeaformis*). *Molecular Ecology*, **22**: 4896–4914.
- 483 Hoban S, Kelley JL, Lotterhos KE *et al.* (2016) Finding the Genomic Basis of Local Adaptation: Pitfalls,  
484 Practical Solutions, and Future Directions. *The American naturalist*, **188**, 379–397.
- 485 Jones MR, Good JM (2016) Targeted capture in evolutionary and ecological genomics. *Molecular  
486 ecology*, **25**, 185–202.
- 487 Jones FC, Grabherr MG, Chan YF *et al.* (2012) The genomic basis of adaptive evolution in threespine  
488 sticklebacks. *Nature*, **484**, 55–61.
- 489 Keinath MC, Timoshevskiy VA, Timoshevskaya NY *et al.* (2015) Initial characterization of the large  
490 genome of the salamander *Ambystoma mexicanum* using shotgun and laser capture chromosome  
491 sequencing. *Scientific reports*, **5**, 16413.
- 492 Lee CE, Remfert JL, Opgenorth T *et al.* (2017) Evolutionary responses to crude oil from the Deepwater  
493 Horizon oil spill by the copepod *Eurytemora affinis*. *Evolutionary applications*, **10**, 813–828.
- 494 Lepais O, Weir JT (2014) SimRAD: an R package for simulation-based prediction of the number of loci  
495 expected in RADseq and similar genotyping by sequencing approaches. *Molecular ecology  
496 resources*, **14**, 1314-1321.
- 497 Li H, Durbin R (2009) Fast and accurate short read alignment with Burrows-Wheeler transform.

- 498      *Bioinformatics*, **25**, 1754–1760.  
499    Li H, Handsaker B, Wysoker A *et al.* (2009) The Sequence Alignment/Map format and SAMtools.  
500      *Bioinformatics*, **25**, 2078–2079.  
501    Lotterhos KE, Whitlock MC (2015) The relative power of genome scans to detect local adaptation  
502      depends on sampling design and statistical method. *Molecular ecology*, **24**, 1031–1046.  
503    Lowry DB, Hoban S, Kelley JL *et al.* (2016) Breaking RAD: An evaluation of the utility of restriction site  
504      associated DNA sequencing for genome scans of adaptation. *Molecular ecology resources*.  
505    Lowry DB, Hoban S, Kelley JL *et al.* (2017) Responsible RAD: Striving for best practices in population  
506      genomic studies of adaptation. *Molecular ecology resources*, **17**, 366–369.  
507    McCartney-Melstad E, Mount GG, Bradley Shaffer H (2016) Exon capture optimization in amphibians  
508      with large genomes. *Molecular ecology resources*, **16**, 1084–1094.  
509    McKinney GJ, Larson WA, Seeb LW, Seeb JE (2017) RADseq provides unprecedented insights into  
510      molecular ecology and evolutionary genetics: comment on Breaking RAD by Lowry *et al.* (2016).  
511      *Molecular ecology resources*, **17**, 356–361.  
512    Müller T, Freund F, Wildhagen H, Schmid KJ (2014) Targeted re-sequencing of five Douglas-fir  
513      provenances reveals population structure and putative target genes of positive selection. *Tree genetics  
514      & genomes*, **11**.  
515    Nadeau NJ, Whibley A, Jones RT *et al.* (2012) Genomic islands of divergence in hybridizing Heliconius  
516      butterflies identified by large-scale targeted sequencing. *Philosophical transactions of the Royal  
517      Society of London. Series B, Biological sciences*, **367**, 343–353.  
518    Neves LG, Davis JM, Barbazuk WB, Kirst M (2013) Whole-exome targeted sequencing of the  
519      uncharacterized pine genome. *The Plant journal: for cell and molecular biology*, **75**, 146–156.  
520    Pastinen T (2010) Genome-wide allele-specific analysis: insights into regulatory variation. *Nature  
521      reviews. Genetics*, **11**, 533–538.  
522    Portik DM, Smith LL, Bi K (2016) An evaluation of transcriptome-based exon capture for frog  
523      phylogenomics across multiple scales of divergence (Class: Amphibia, Order: Anura). *Molecular  
524      ecology resources*, **16**, 1069–1083.  
525    Puritz JB, Matz MV, Toonen RJ *et al.* (2014), Demystifying the RAD fad. *Molecular ecology*, **23**: 5937–  
526      5942.  
527    Puritz JB, Hollenbeck CM, Gold JR (2014) dDocent: a RADseq, variant-calling pipeline designed for  
528      population genomics of non-model organisms. *PeerJ*, **2**, e431.  
529    Quinlan AR (2014) BEDTools: The Swiss-Army Tool for Genome Feature Analysis. *Current protocols in  
530      bioinformatics / editorial board, Andreas D. Baxevanis ... [et al.]*, **47**, 11.12.1–34.  
531    Reid NM, Proestou DA, Clark BW *et al.* (2016) The genomic landscape of rapid repeated evolutionary  
532      adaptation to toxic pollution in wild fish. *Science*, **354**, 1305–1308.  
533    Samuels DC, Han L, Li J *et al.* (2013) Finding the lost treasures in exome sequencing data. *Trends in  
534      genetics: TIG*, **29**, 593–599.  
535    Schlötterer C, Tobler R, Kofler R, Nolte V (2014) Sequencing pools of individuals — mining genome-  
536      wide polymorphism data without big funding. *Nature reviews. Genetics*, **15**, 749–763.  
537    Schmid S, Genevest R, Gobet E *et al.* (2017) HyRAD-X, a versatile method combining exome capture and  
538      RAD sequencing to extract genomic information from ancient DNA. *Methods in ecology and  
539      evolution / British Ecological Society*.  
540    Suchan T, Pitteloud C, Gerasimova NS *et al.* (2016) Hybridization Capture Using RAD Probes (hyRAD),  
541      a New Tool for Performing Genomic Analyses on Collection Specimens. *PloS one*, **11**, e0151651.  
542    Sulonen A-M, Ellonen P, Almusa H *et al.* (2011) Comparison of solution-based exome capture methods  
543      for next generation sequencing. *Genome biology*, **12**, R94.  
544    Suren H, Hodgins KA, Yeaman S *et al.* (2016) Exome capture from the spruce and pine giga-genomes.  
545      *Molecular ecology resources*, **16**, 1136–1146.  
546    Syring JV, Tennessen JA, Jennings TN *et al.* (2016) . *Frontiers in plant science*, **7**, 484.

- 547 Therkildsen NO, Palumbi SR (2017) Practical low-coverage genomewide sequencing of hundreds of  
548 individually barcoded samples for population and evolutionary genomics in nonmodel species.  
549 *Molecular ecology resources*, **17**, 194–208.
- 550 Yeaman S, Hodgins KA, Lotterhos KE *et al.* (2016) Convergent local adaptation to climate in distantly  
551 related conifers. *Science*, **353**, 1431–1433.
- 552 Zhou L, Holliday JA (2012) Targeted enrichment of the black cottonwood (*Populus trichocarpa*) gene  
553 space using sequence capture. *BMC genomics*, **13**, 703.

554

## 555 Data Accessibility

556 Raw, demultiplexed sequences are archived at the NCBI Short Read Archive (BioProject:  
557 PRJNA423022). A complete and updated EecSeq protocol can be found at  
558 (<https://github.com/jpuritz/EecSeq>) along with bioinformatic code to repeat all analyses described  
559 in this paper.

## 560 Author Contributions

561 JP conceived the original concept of this work and performed all laboratory and data analysis.  
562 KL contributed all reagents and experimental materials. JP and KL designed the research,  
563 experiments, and data analysis, and wrote the manuscript.

564

565

566

567

568

569

570

571

572

573

574

575

576    **Tables**

577

Oligo Name	Sequence
Universal_SAI1_Adapter	AATGATAACGGGACCACCGAGATCTACACTTTCCCTACACGACGCTTCCGATCTGCGACT*T
Indexed_Adapter_SAI1_I5	P*AGTCGACAGATCGGAAGAGCACACGTCTGAACCTCCAGTCACACAGTGATCTGTATGCCGTCTCTGCTTG
Indexed_Adapter_SAI1_I8	P*AGTCGACAGATCGGAAGAGCACACGTCTGAACCTCCAGTCACACTGAATCTGTATGCCGTCTCTGCTTG
Indexed_Adapter_SAI1_I9	P*AGTCGACAGATCGGAAGAGCACACGTCTGAACCTCCAGTCACGATCTGTATGCCGTCTCTGCTTG
Indexed_Adapter_SAI1_I11	P*AGTCGACAGATCGGAAGAGCACACGTCTGAACCTCCAGTCACGGCTACATCTGTATGCCGTCTCTGCTTG

578    **Table 1. Corrected adapter sequences for mRNA library preparation.**

579    Oligos are listed in a 5' to 3' orientation with "P" indicates a phosphorylation modification to enable  
580    ligation.

581

582

Replicate Capture Pool	Raw Reads	Filtered Reads	Mapped Reads	% Duplicates	Filtered Mapped Reads	% mapping to mitochondrial genome
EC_2	53,493,950	53,118,952	42,403,525	5.8	35,955,539	1.7%
EC_4	44,935,340	44,651,228	35,275,663	6.1	29,519,347	1.6%
EC_7	43,745,614	43,448,296	35,007,184	5.6	29,723,437	2.2%
EC_1	41,402,996	41,145,724	32,668,750	4.5	27,940,717	1.8%
EC_3	56,127,536	55,753,268	44,605,960	5.0	38,103,268	1.9%
EC_12	46,068,760	45,750,394	37,227,067	6.2	31,497,298	2.2%

583 **Table 2. Exome capture sequencing, filtering, and mapping statistics.**  
584 EC\_2, EC\_4, and EC\_7 are the three replicate captures with the original probe pool, and EC\_1, EC\_3, and  
585 EC\_12 are the replicate captures with the probe pool exposed to 12 extra rounds of PCR.  
586

587

Targets	Capture Pool					
	EC_2	EC_4	EC_7	EC_1	EC_3	EC_12
All Exons	88.0%	86.0%	85.8%	86.5%	87.9%	86.4%
20XR Exons	99.5%	99.4%	99.4%	99.4%	99.5%	99.4%
35XR Exons	99.6%	99.6%	99.6%	99.6%	99.6%	99.6%
50XR Exons	99.7%	99.7%	99.7%	99.7%	99.7%	99.7%

588

**Table 3. Exome capture sensitivity with a 1x threshold.**

589

Sensitivity is the percentage of target bp with at least one read mapping successfully. Here, targets are  
590 broken up into subsets: All annotated exons, exons with at least 20X coverage from the RNA library,  
591 exons with at least 35X coverage from the RNA library, and exons with at least 50X coverage from the  
592 RNA library. EC\_2, EC\_4, and EC\_7 are the three replicate captures with the original probe pool, and  
593 EC\_1, EC\_3, and EC\_12 are the replicate captures with the probe pool exposed to 12 extra rounds of PCR.  
594

Reference and species	Num. target genes or exons	Sensitivity % of targeted regions > 10x depth	Specificity % of reads mapping to targeted bases	% of reads mapping near target	% of reads mapping off target	Notes
EecSeq (this study) eastern oyster <i>Crassostrea virginica</i>	71,105 (51,096-110,020)	All exons: 54.7% Expressed Exons: 98.8% ( 97.4% - 99.1%)	All exons: 47.8% Expressed Exons: 37.0% ( 33.6% - 41.4%)	All exons: 28.4% Expressed Exons: 23.6% (22.3% - 25.2%)	All exons: 23.7% Expressed Exons: 39.3% (33.3% - 44.1%)	
(Suren <i>et al.</i> 2016) pine and spruce <i>Picea glauca x engelmannii</i> and <i>Pinus contorta</i>	26824 genes (pine) 28649 genes(spruce)	51% (spruce) and 59% (pine) (all samples, also metrics for 75% of samples)	18.5% (spruce) and 21 % (pine)	37% (spruce) 38% (pine)	44% (spruce) and 41% (pine)	Non-model species, large genomes, near target defined as 500 bp
(Zhou & Holliday 2012) black cottonwood <i>Populus trichocarpa</i>	20.76Mb (5%) of exons, regulatory regions	86.8 % (at 100X coverage about 0-8%)	~93%	On average, approximately 80 base pairs nearest the bait were sequenced at a depth of > 10X	NR	Model species with good genome. Off target defined as > 250bp away.
(Hebert <i>et al.</i> 2013) lake whitefish <i>Coregonus clupeaformis</i>	11,975 nuclear exons, and other genomic markers using 62,438 probes	NR	11.8%	NR	NR	98% of targeted genes (2728) were successfully captured a mean read depth of 31X
(Bi <i>et al.</i> 2012) chipmunk <i>Tamias alpinus</i>	11,975 exons	40.3%	25%	NR	NR	% of exons that were covered by at least one read, > 99%
(Christmas <i>et al.</i> 2017) narrow-leaf hopbush <i>Dodonaea viscosa</i> ssp. <i>angustissima</i>	700 genes	NR	15.7%	NR	NR	Did not account for intron sites
(Syring <i>et al.</i> 2016) whitebark pine <i>Pinus albicaulis</i>	7,849 distinct transcripts	NR	13%	NR	NR	
(Müller <i>et al.</i> 2014) douglas-fir <i>Pseudotsuga menziesii</i>	57,110 exons	90%	32-52% per individual	NR	NR	
(Nadeau <i>et al.</i> 2012) butterflies	BAC loci (3.5 MB; 57,610 baits)	75.6%	33.5%	NR	NR	

595 **Table 4. Comparing specificity and sensitivity across capture methods.**  
596 A summary of sensitivity and specificity of recent exome-capture studies in which probes were  
597 designed from the same species. NR: not reported.  
598

## 599 Figure Legends

### 600 **Figure 1. Conceptual Diagram of Expressed Exome Capture Sequencing.**

601 Upper left panel: The shotgun genomic DNA library that will be captured with probes. Middle left panel: EecSeq  
602 relies on custom RNA adapters that contains a SAI restriction site. Middle upper panel: The adapters are  
603 incorporated into a mRNA library preparation that is normalized with duplex-specific nuclease. Adapters are then  
604 removed with a SAI restriction digest, cDNA probes are subsequently blunted with mung bean nuclease, and  
605 biotinylated via a PCR reaction. Upper right panel: The probes are then hybridized to the shotgun genomic library  
606 with TruSeq style adapters. Exon loci bind to the cDNA probes. Lower panel: Hybridized exon loci and probes are  
607 then captured with magnetic Streptavidin beads. The captured exome fragments are washed several times, eluted,  
608 enriched with PCR, and then sequenced.

### 609 **Figure 2. Distribution of RNA reads across regions of the oyster genome.**

610 Percentage of bases within exons- both coding sequences (CDS) and untranslated exon regions (UTR), intergenic,  
611 and intron regions at various coverage levels.

### 612 **Figure 3. Total DNA and RNA coverage across all exons.**

613 Depth was calculated as the total number of reads overlapping with an exon region. For exome capture depth (y-  
614 axis), reads were summed across all 6 replicate captures. For RNA read depth, reads were summed across all four  
615 libraries. The shape and color of each point was determined by the percentile size of the respective exon (lower 10%  
616 < 59 bp, upper 10% > 517 bp, and the middle 80% was between 57 bp and 517bp). Note that the DNA reads were  
617 sequenced to greater depth than the RNA-derived probes.

### 618 **Figure 4. Per base pair EecSeq capture sensitivity.**

619 To measure EecSeq capture (DNA) sensitivity, capture targets were defined as exons that had more than 35X  
620 coverage in the RNAseq (probe) data. Confidence intervals were generated by defining capture targets between 20X  
621 RNAseq coverage and 50X RNAseq coverage. Near-target mapping were 150 bp on either side of the defined  
622 targets. This range corresponds to the modal DNA fragment length used for the capture libraries with the expectation  
623 that exon probes could capture reads that far from the original target. EC\_2, EC\_4, and EC\_7 are the three replicate  
624 captures with the original probe pool, and EC\_1, EC\_3, and EC\_12 are the replicate captures with the probe pool  
625 exposed to 12 extra rounds of PCR. Depth in this figure is the depth of DNA reads from EecSeq captures.

### 626 **Figure 5. Boxplots of mean per basepair coverage levels plotted across exons size windows.**

627 All annotated exons were broken into 10bp - 30 bp windows depending on overall size and the mean per basepair  
628 coverage per capture was calculated for each window size. The line each box represents the median of mean  
629 coverage values and the box surrounds the 25<sup>th</sup> and 75<sup>th</sup> percentiles. The mean of each bin class is plotted as a black  
630 diamond with standard error bars around it. Outlier points were not plotted. Note that the data for this graph is for  
631 all annotated exons, regardless of expected capture. See Supplemental Figure 3 for a similar plot focused on an  
632 expressed target set.

### 633 **Figure 6. Mean capture depth plotted against exon GC content.**

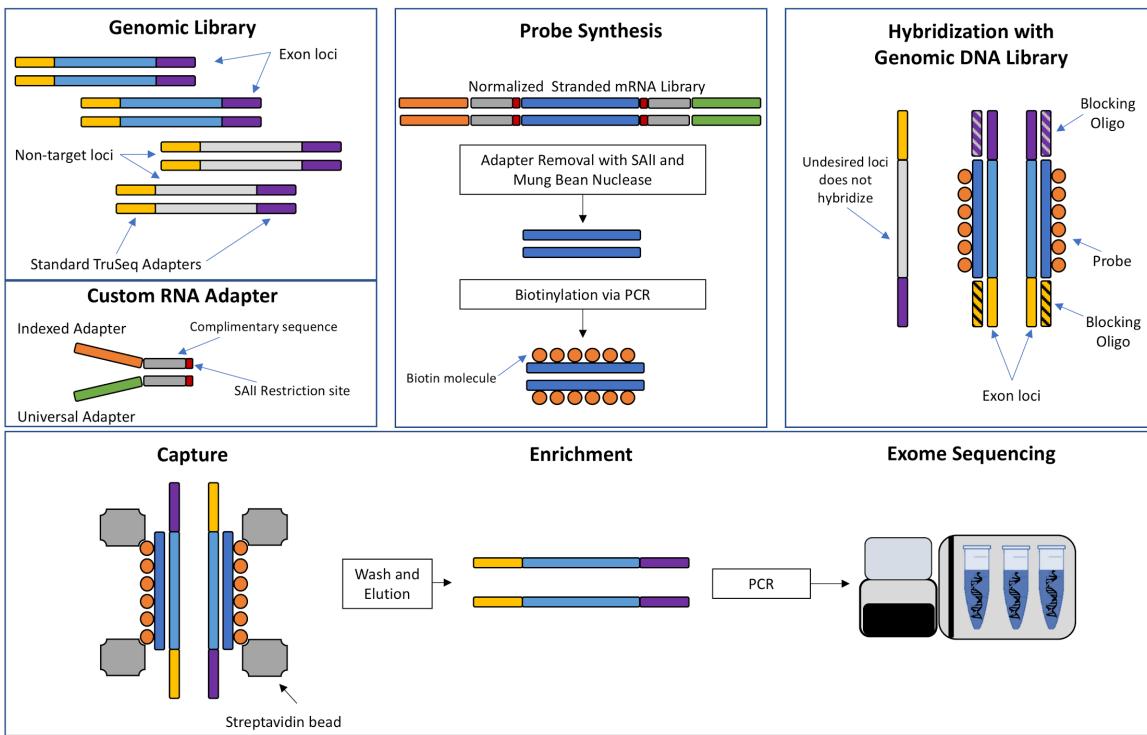
634 Exons were broken up into three size windows: (1) Lower 10%- exons less than 57 bp, (2) Middle 80%- exons  
635 greater than 56 bp and less than 518, (3) Upper 10%- exons greater than 517 bp.

### 636 **Figure 7. EecSeq capture and probe coverage across Heat Shock cognate 71 kDa.**

637 Coverage for each replicate capture pools is plotted along base pairs 32,740,000 to 32,755,000 of reference  
638 Chromosome NC\_035780.1 containing the full gene region of Heat Shock cognate 71 kDa (NCBI Reference  
639 Sequence: XM\_022472393.1), predicted glucose-induced degradation protein 8 homolog (NCBI Reference  
640 Sequence: XM\_022486802), and a partial gene region for rho GTPase-activating protein 39-like (NCBI Reference  
641 Sequence: XM\_022486743.1). Each exome capture pool coverage is plotted in light blue with dashed grey border,  
642 and a rolling 100 bp window average across all pools is plotted in dark blue. Each RNAseq (probe) sample coverage  
643 is plotted in light red with dashed grey border and a rolling 100 bp window average across all pools is plotted in dark  
644 red. Gene regions are marked in purple with exons color coded by gene. Coding sequence (CDS) is marked by a  
645 white bar within exon markers.

646 Figures

647 **Figure 1. Conceptual Diagram of Expressed Exome Capture Sequencing.**

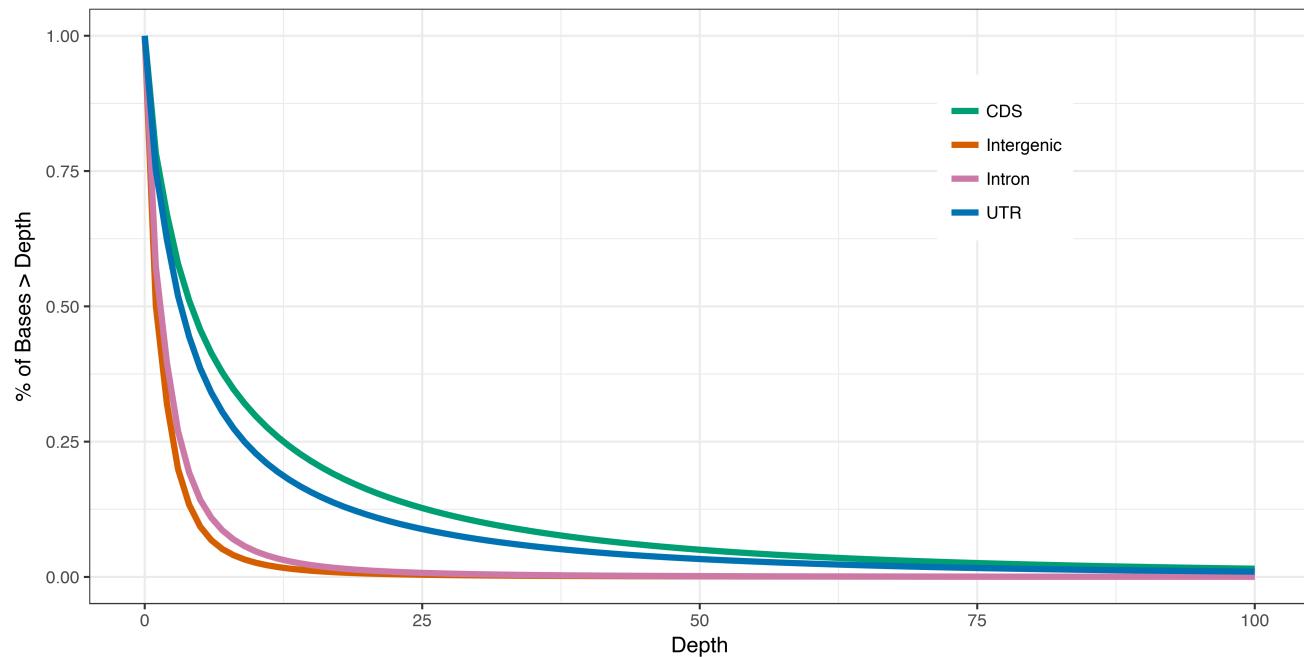


648

649

650

651 **Figure 2. Distribution of RNA reads across regions of the oyster genome.**

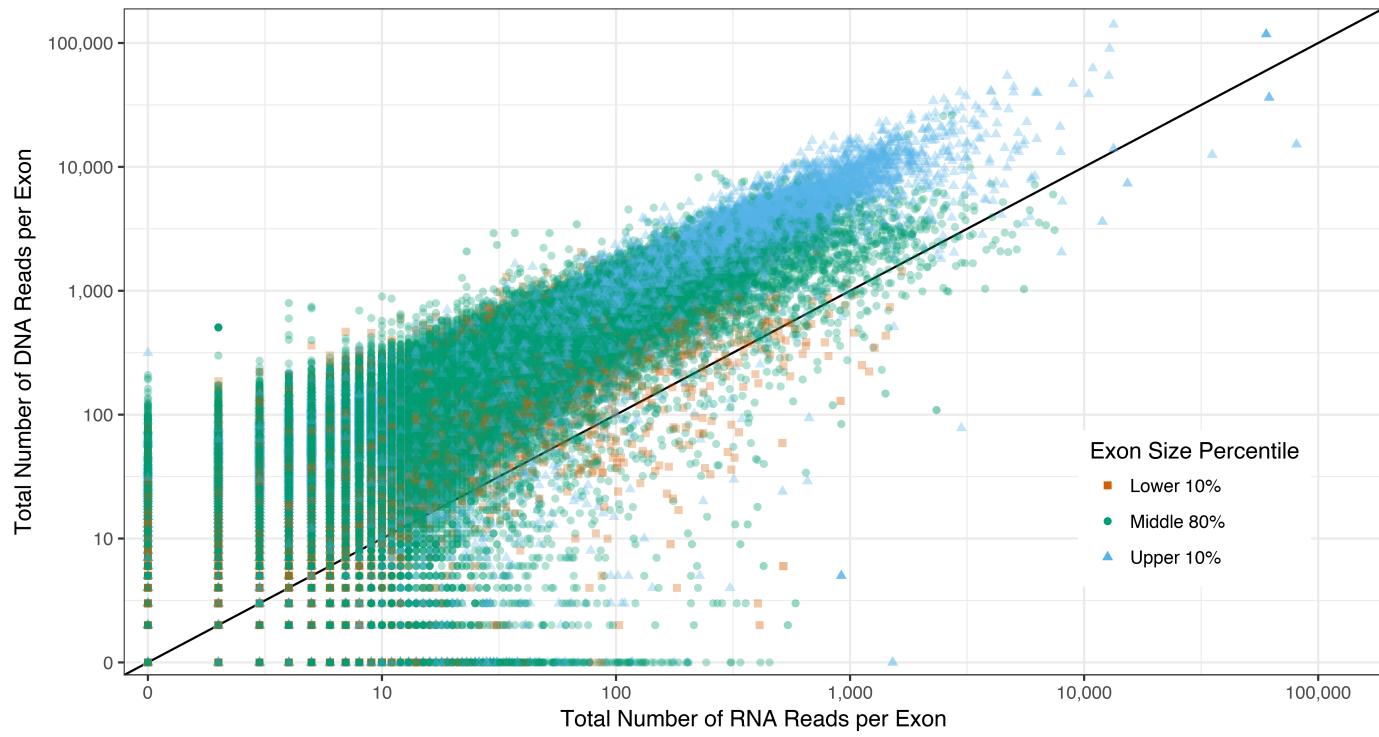


652

653

654

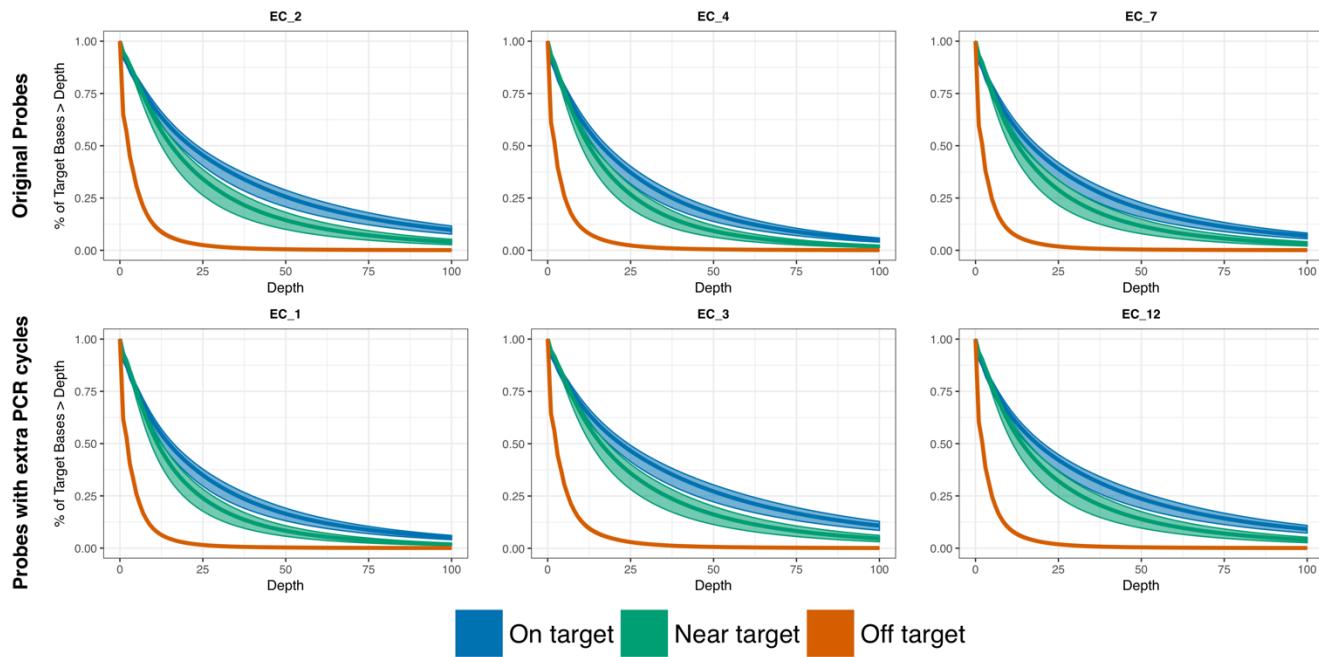
655 **Figure 3. Total DNA and RNA coverage across all exons.**



656

657

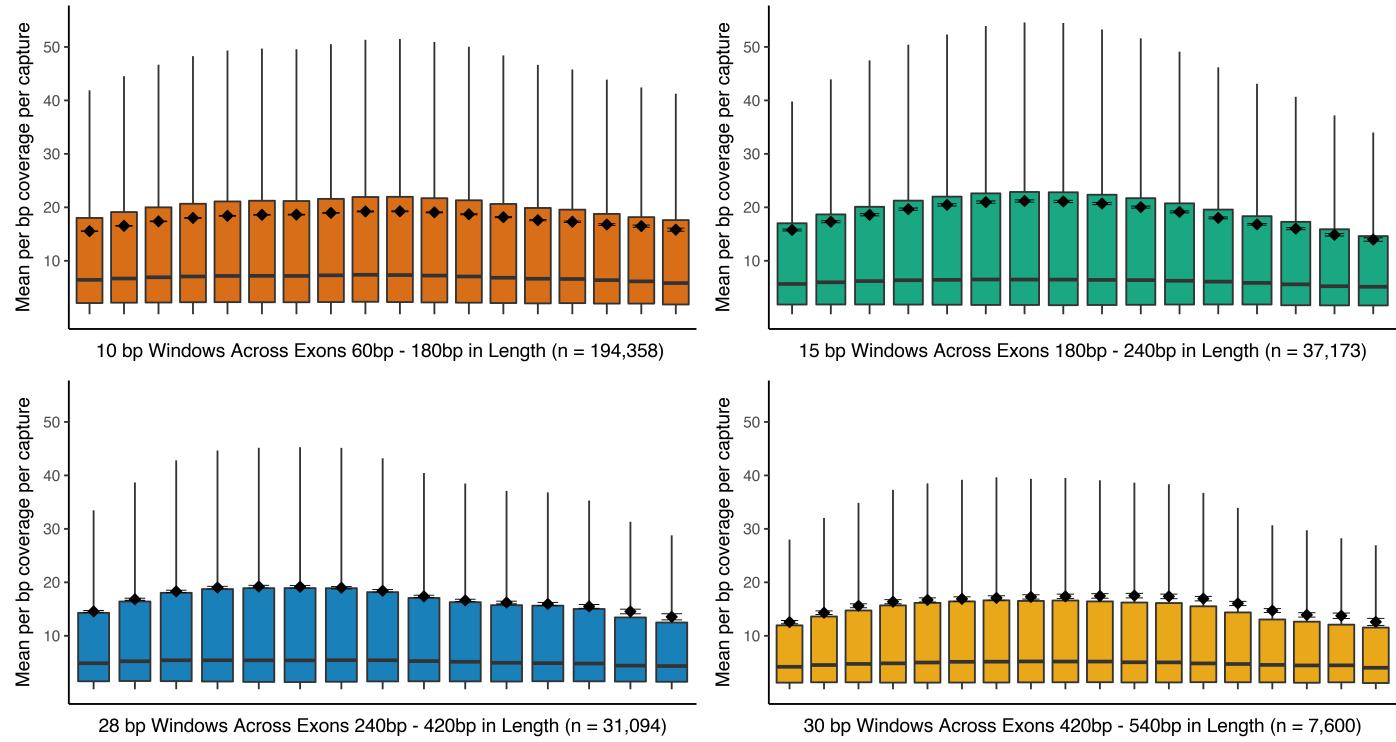
658



659 **Figure 4. Per base pair EecSeq capture sensitivity.**

660

661



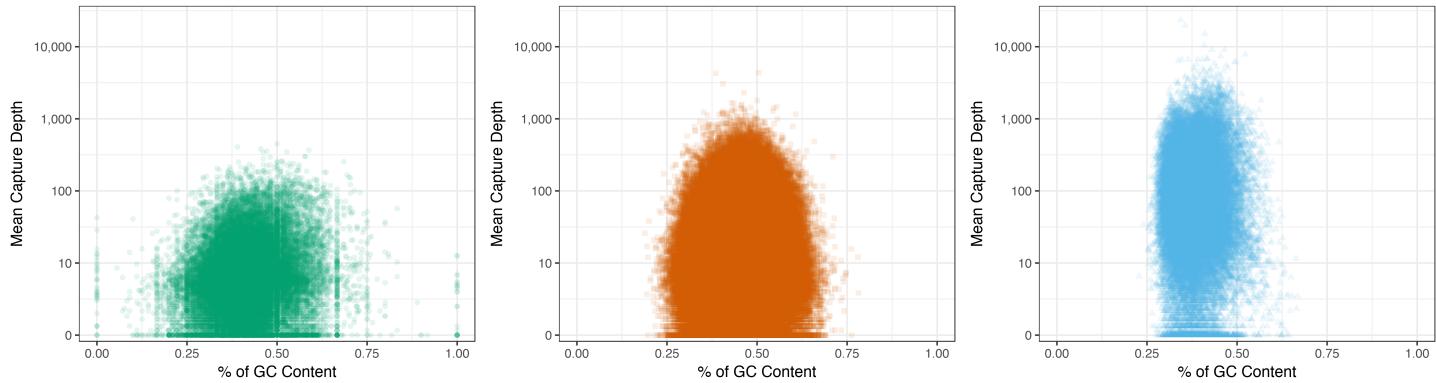
662

663 **Figure 5. Boxplots of mean per basepair coverage levels plotted across exons size windows.**

664

665 **Figure 6. Mean capture depth plotted against exon GC content.**

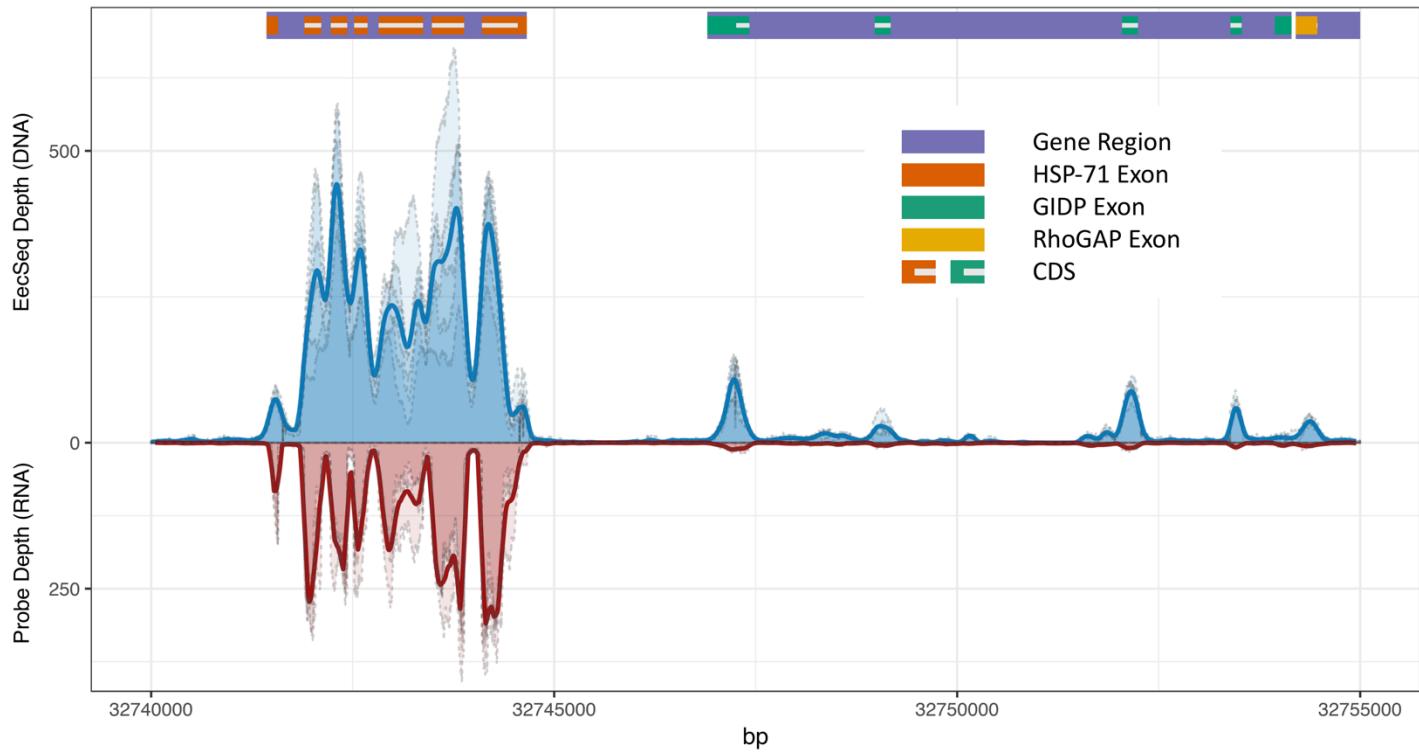
666



667 **Figure 7. EecSeq capture and probe coverage across Heat Shock cognate 71 kDa.**

668

669



670

671

672

673

674

675

676

## **Expressed Exome Capture Sequencing (EecSeq): a method for cost-effective exome sequencing for all organisms with or without genomic resources**

### **Supplemental Material**

Jonathan B. Puritz and Katie E Lotterhos

# Tables

## Supplemental Table 1: Original RNA Adapters

Oligo Name	Sequence
Universal_SAI1_Adapter	AATGATAACGGCACCACCGAGATCTACACTTTCCCTACACGACGCTTCCGATCGTCGACT*T
Indexed_Adapter_SAI1_I5	P*AGTCGACGATCGGAAGAGCACACGTCTGAACCTCCAGTCACACAGTGATCTCGTATGCCGTCTGCTTG
Indexed_Adapter_SAI1_I8	P*AGTCGACGATCGGAAGAGCACACGTCTGAACCTCCAGTCACACTGAATCTCGTATGCCGTCTGCTTG
Indexed_Adapter_SAI1_I9	P*AGTCGACGATCGGAAGAGCACACGTCTGAACCTCCAGTCACGATCTCGTATGCCGTCTGCTTG
Indexed_Adapter_SAI1_I11	P*AGTCGACGATCGGAAGAGCACACGTCTGAACCTCCAGTCACGGCTACATCTCGTATGCCGTCTGCTTG

**Supplemental Table 2: Original DNA Adapters.** Oligos should be paired for adapter formation following 1.1.X pairs with 1.2.X.

DNA_P1.1.1	ACACTTTCCCTACACGACGCTTCCGATCTGCATGG*
DNA_P1.1.2	ACACTTTCCCTACACGACGCTTCCGATCTAACCGAG*T
DNA_P1.1.3	ACACTTTCCCTACACGACGCTTCCGATCTCGATCG*T
DNA_P1.1.4	ACACTTTCCCTACACGACGCTTCCGATCTCGATG*T
DNA_P1.1.5	ACACTTTCCCTACACGACGCTTCCGATCTGCATG*T
DNA_P1.1.6	ACACTTTCCCTACACGACGCTTCCGATCTAACCG*T
DNA_P1.1.7	ACACTTTCCCTACACGACGCTTCCGATCTGGTTGG*T
DNA_P1.1.8	ACACTTTCCCTACACGACGCTTCCGATCTAAGGAG*T
DNA_P1.2.1	PC*CGTACAGATCGGAAGAGCGTCGTAGGGAAAGAGTGT
DNA_P1.2.2	PC*TTGGTAGATCGGAAGAGCGTCGTAGGGAAAGAGTGT
DNA_P1.2.3	PC*GCTAGAGATCGGAAGAGCGTCGTAGGGAAAGAGTGT
DNA_P1.2.4	PC*AGCTAAGATCGGAAGAGCGTCGTAGGGAAAGAGTGT
DNA_P1.2.5	PC*ACGTAAGATCGGAAGAGCGTCGTAGGGAAAGAGTGT
DNA_P1.2.6	PC*GTTGGAGATCGGAAGAGCGTCGTAGGGAAAGAGTGT
DNA_P1.2.7	PC*CCAACAGATCGGAAGAGCGTCGTAGGGAAAGAGTGT
DNA_P1.2.8	PC*TTCCTAGATCGGAAGAGCGTCGTAGGGAAAGAGTGT
DNA_P2.1	P*GATCGGAAGAGCGAGAACAA
DNA_P2.2	GTGACTGGAGTTCACACGTGTGCTTCCGATC*T

**Supplemental Table 3: Correct DNA Adapters.** Oligos should be paired for adapter formation following 1.1.X pairs with 1.2.X.

DNA_P1.1.1	ACACTTTCCCTACACGACGCTTCCGATCTGCATGG*
DNA_P1.1.2	ACACTTTCCCTACACGACGCTTCCGATCTAACCGAG*T
DNA_P1.1.3	ACACTTTCCCTACACGACGCTTCCGATCTCGATCG*T
DNA_P1.1.4	ACACTTTCCCTACACGACGCTTCCGATCTCGATG*T
DNA_P1.1.5	ACACTTTCCCTACACGACGCTTCCGATCTGCATG*T
DNA_P1.1.6	ACACTTTCCCTACACGACGCTTCCGATCTAACCG*T
DNA_P1.1.7	ACACTTTCCCTACACGACGCTTCCGATCTGGTTGG*T
DNA_P1.1.8	ACACTTTCCCTACACGACGCTTCCGATCTAAGGAG*T
DNA_P1.2.1	PC*CATGCAGATCGGAAGAGCGTCGTAGGGAAAGAGTGT
DNA_P1.2.2	PC*TGGTTAGATCGGAAGAGCGTCGTAGGGAAAGAGTGT
DNA_P1.2.3	PC*GATCGAGATCGGAAGAGCGTCGTAGGGAAAGAGTGT
DNA_P1.2.4	PC*ATCGAAGATCGGAAGAGCGTCGTAGGGAAAGAGTGT
DNA_P1.2.5	PC*ATGCAAGATCGGAAGAGCGTCGTAGGGAAAGAGTGT
DNA_P1.2.6	PC*GGTTGAGATCGGAAGAGCGTCGTAGGGAAAGAGTGT
DNA_P1.2.7	PC*CAACCAGATCGGAAGAGCGTCGTAGGGAAAGAGTGT
DNA_P1.2.8	PC*TCCTTAGATCGGAAGAGCGTCGTAGGGAAAGAGTGT
DNA_P2.1	P*GATCGGAAGAGCGAGAACAA
DNA_P2.2	GTGACTGGAGTTCACACGTGTGCTTCCGATC*T

**Supplemental Table 4. RNA sequencing statistics.** Samples 3E and 4E were heat-shocked individuals and Samples 1C and 3C were control individuals.

Sample	Raw Reads	Filtered Reads	Mapped Reads	Filtered Mapped Reads
3E	25259714	23736708	12499756	5682642
4E	24794376	23304626	12619770	5843708
1C	22749260	21494102	12499756	4717835
3C	25485364	23900618	12594813	5745840

### Supplemental Table 5. Exome capture specificity with a 10X threshold.

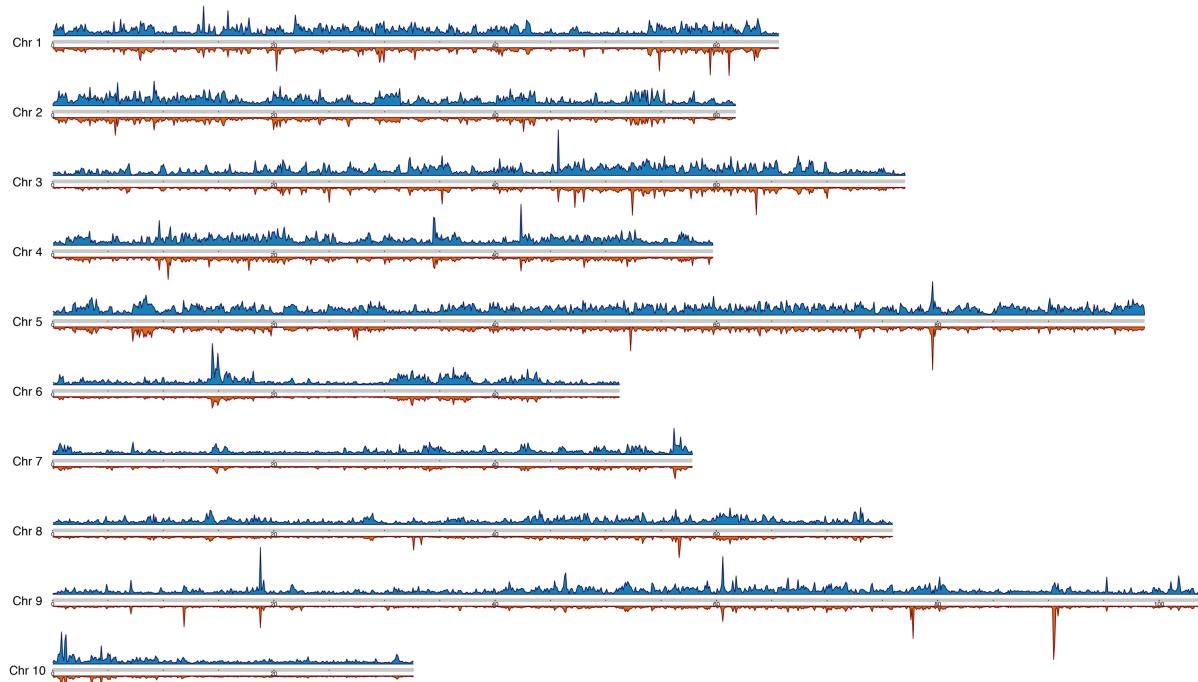
Sensitivity measured here is the percentage of targets with at least 10 reads mapping successfully. Here, targets are broken up into subsets: All annotated exons, exons with at least 20X coverage from the RNA library, exons with at least 35X coverage from the RNA library, and exons with at least 50X coverage from the RNA library. EC\_2, EC\_4, and EC\_7 are the three replicate captures with the original probe pool, and EC\_1, EC\_3, and EC\_12 are the replicate captures with the probe pool exposed to 12 extra rounds of PCR.

10X	Capture Pool					
Targets	EC_2	EC_4	EC_7	EC_1	EC_3	EC_12
All Exons	58.8%	51.9%	52.2%	52.3%	58.7%	54.2%
20XR Exons	98.2%	96.7%	97.1%	96.6%	98.2%	97.6%
35XR Exons	99.1%	98.6%	98.8%	98.6%	99.0%	98.9%
50XR Exons	99.3%	99.0%	99.1%	99.0%	99.2%	99.2%

**Supplemental Table 6. Per sample cost calculations for EecSeq.** The calculations below assume 8 mRNA libraries are used to create probes to capture 96 samples in 8 capture reactions (12 samples per capture). Costs assumes the captured DNA is sequenced in one to one and a half lane(s) of Illumina High Seq 4000. This assumes that coverage levels for 96 samples in one lane would be equivalent to coverage levels seen in six pools of eight samples in half a lane. Cost does not include DNA or RNA extraction. See the github repository for more information on library preps and capture. Based on results in the main paper, this multiplexing strategy would give ~7.66x coverage per individual at exons represented by 35X sequencing depth at RNA-derived probes and would give ~3,508 exome SNPs at 10X coverage. Whether 96 individuals can be sequenced to enough depth in a single lane will depend on the number of megabases represented by the probes, the desired read depth, and the sensitivity and specificity of capture in the focal species. We have included the cost if 96 samples were sequenced over one and half lanes to include this potential variance.

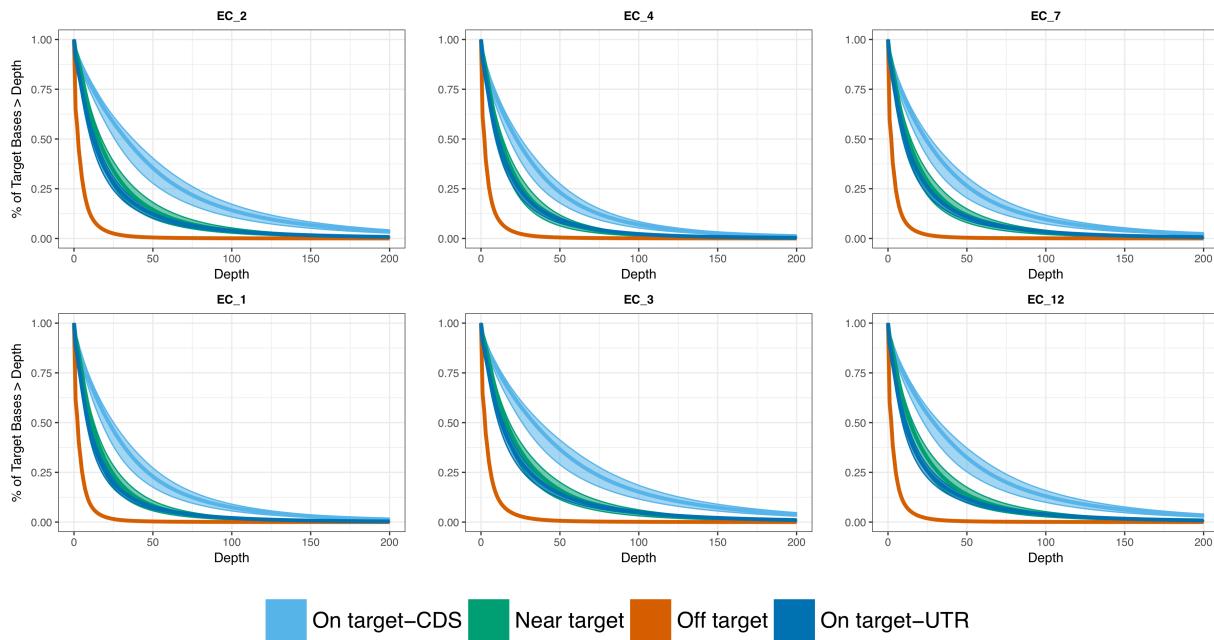
Reagent	Vendor	Price	Total Units	Units Used	Project Cost
Genomic DNA Shearing	Genomic Core Lab	1.50	1	96	\$144.00
Kapa-Stranded mRNA-Seq 24 rxn kit - Illumina	KapaBiosystems KK8420	\$1,008.00	24 rxn	4 (1/2 rxns are used)	\$168.00
Hyper Prep gDNA kit with KAPA Library Amplification Primer Mix (10X) 96 rxn kit	KapaBiosystems KK8504	\$2,496.00	96 rxn	48 (1/2 rxns are used)	\$1,248.00
KAPA Pure Beads (5 mL)	KapaBiosystems KK8000	\$150.00	100	100	\$150.00
Oligos	IDT	\$2,391.85	1000	2	\$4.78
DSN	Evrogen EA001	\$350	50	20	\$140.00
Library Amplification Polymerase	KapaBiosystems KK2611	\$126.00	50	5	\$12.60
SAII-HF Enzyme	NEB R3138S	\$61.00	2000	100	\$3.05
Mung Bean Nuclease	NEB M0250S	\$63.00	1500	50	\$2.10
DecaLabel™ Biotin DNA Labeling Kit	ThermoFisher K0651	\$158.00	10	1	\$15.80
Denhardt's Solution (50X)	ThermoFisher 750018	\$149.00	100	0.0128	\$0.02
Dynabeads™ M-280 Streptavidin	ThermoFisher 11205D	\$496.00	2 mL	0.08	\$19.84
Human Cot-1 DNA (1 mg/ml)	ThermoFisher 15279011	\$230.00	500 µg	4	\$1.84
<b>Total Prep Cost</b>		<b>\$7,680.35</b>			<b>\$1,910.03</b>
				<b>Per Sample</b>	<b>\$19.90</b>
Sequencing on Illumina Hi-Seq 4000	Genomic Core Lab	\$2,700	1	1-1.5	\$2,700
				<b>Per Sample</b>	<b>\$28.13-\$42.18</b>
				<b>Total Per Sample</b>	<b>\$48.02-\$62.08</b>

## Supplemental Figures



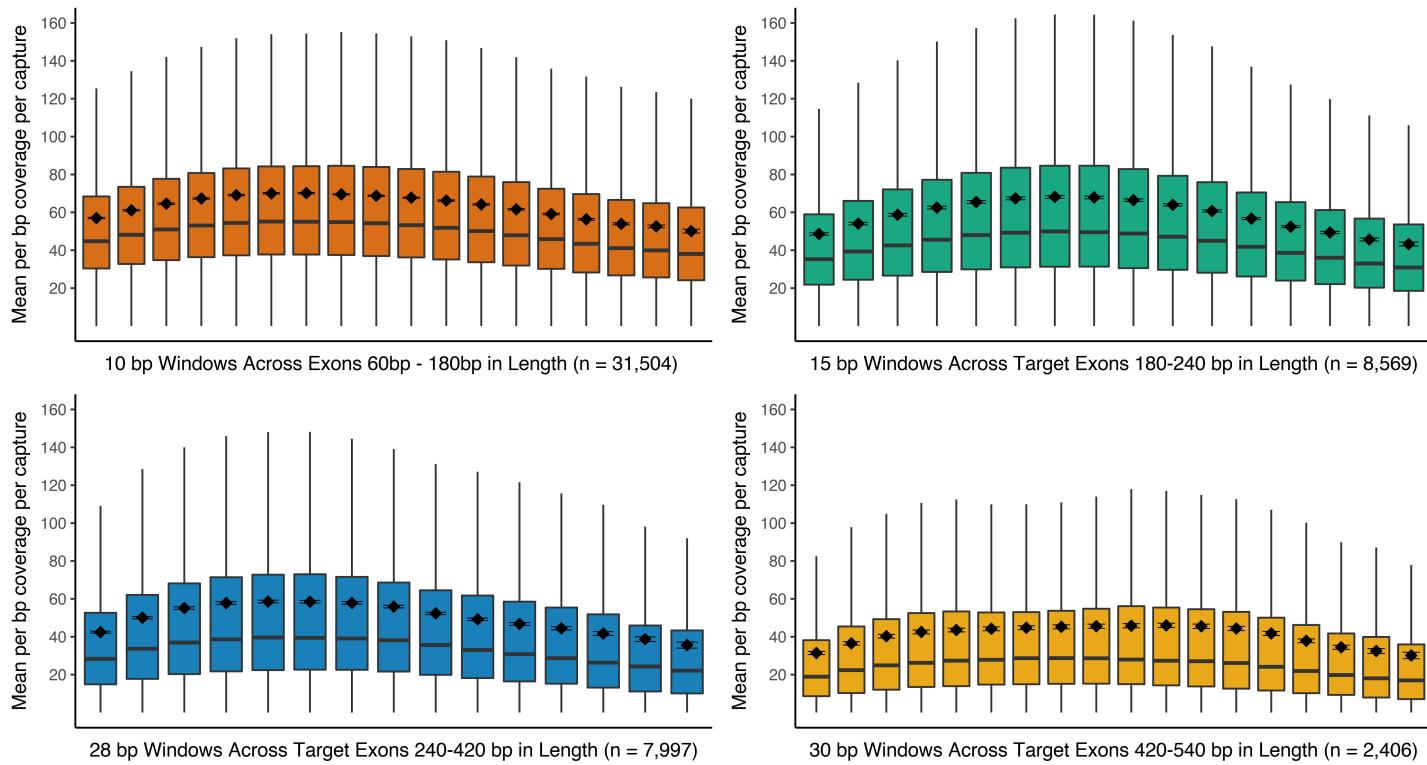
**Supplemental Figure 1. Expressed exome capture reads and RNAseq reads from probes plotted across the eastern oyster genome.**

Read coverage density was plotted in 10,000 bp sliding windows for both total RNA reads (red; below chromosome) and total EecSeq reads (blue; above chromosome) using the karyoploteR package (<https://bioconductor.org/packages/release/bioc/html/karyoploteR.html>).



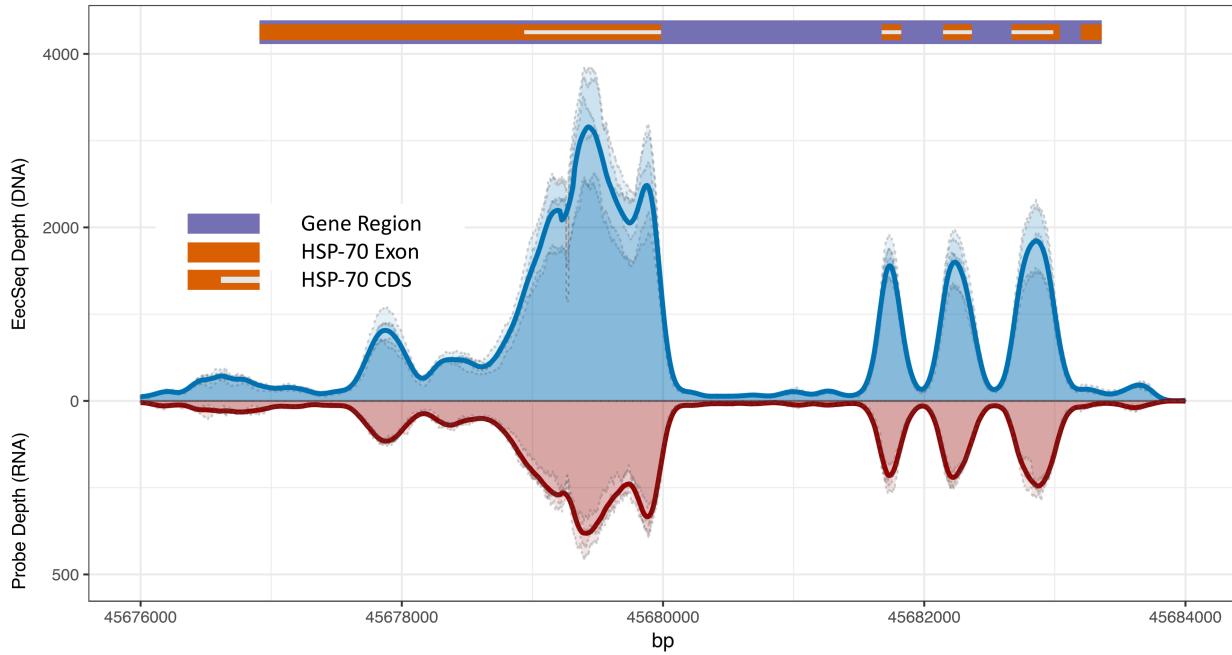
## Supplemental Figure 2. Per base pair sensitivity plot of EecSeq captures including CDS and UTR.

To compare EecSeq to other capture methods, capture targets were defined as exons that had more than 35X coverage in the RNAseq (probe) data and confidence intervals were generated by defining capture targets between 20X RNAseq coverage and 50X RNAseq coverage. Near-target mapping were 150 bp on either side of the defined targets. For this figure, target exons were broken into coding sequence (CDS) and untranslated regions (UTR) for comparisons. This range corresponds to the modal DNA fragment length used for the capture libraries with the expectation that exon probes could capture reads that far from the original target. EC\_2, EC\_4, and EC\_7 are the three replicate captures with the original probe pool, and EC\_1, EC\_3, and EC\_12 are the replicate captures with the probe pool exposed to 12 extra rounds of PCR.



**Supplemental Figure 3. Boxplots of mean per basepair coverage levels plotted across target exons size windows.**

Target exons (those with at least 35X coverage in the RNA data) were broken into 10bp - 30 bp windows depending on overall size and the mean per basepair coverage per capture was calculated for each window size. The line each box represents the median of mean coverage values and the box surrounds the 25<sup>th</sup> and 75<sup>th</sup> percentiles. The mean of each bin class is plotted as a black diamond with standard error bars around it.



#### Supplemental Figure 4. EecSeq capture and probe coverage across Heat Shock cognate 70 kDa.

Coverage for each replicate capture pools is plotted along basepairs 45,760,000 to 45,684,000 of reference Chromosome NC\_035782.1 containing the full gene region of 70 kDa protein 12B-like (NCBI Reference Sequence: XM\_022468697.1). Each exome capture pool coverage is plotted in light blue with dashed grey border and a rolling 100 bp window average across all pools is plotted in dark blue. Each RNAseq (probe) sample coverage is plotted in light red with dashed grey border and a rolling 100 bp window average across all pools is plotted in dark red. Gene regions are marked in purple with exons color coded by gene. Coding sequence (CDS) is marked by a white bar within exon markers.