# BrAD-seq: Breath Adapter Directional sequencing: a streamlined, ultra-simple and fast library preparation protocol for strand specific mRNA library construction.

Brad Thomas Townsley, Michael F Covington, Yasunori Ichihashi, Kristina Zumstein and Neelima Roy Sinha

# BrAD-seq: <u>Br</u>eath <u>A</u>dapter <u>D</u>irectional sequencing: a streamlined, ultra-simple and fast library preparation protocol for strand specific mRNA library construction.

Brad Townsley[1], Michael F. Covington[1], Yasunori Ichihashi[1,2,], Kristina Zumstein and Neelima Sinha[1,*]

University of California at Davis, Department of Plant Biology, Davis, California, United States

[2] Present address: RIKEN Center for Sustainable Resource Science, Yokohama, Kanagawa, 230-0045 Japan

**\* Correspondence:** Corresponding Author, University of California at Davis, Department of Plant Biology, One Shields Avenue, Davis, California, 95616, United States.

nrsinha@ucdavis.edu

## ABSTRACT

Next Generation Sequencing (NGS) is driving rapid advancement in biological understanding and RNA-sequencing (RNA-seq) has become an indispensable tool for biology and medicine. There is a growing need for access to these technologies although preparation of NGS libraries remains a bottleneck to wider adoption. Here we report a novel method for the production of strand specific RNA-seq libraries utilizing the terminal breathing of double-stranded cDNA to capture and incorporate a sequencing adapter. Breath Adapter Directional sequencing (BrAD-seq) reduces sample handling and requires far fewer enzymatic steps than most available methods to produce high quality strand-specific RNA-seq libraries. The method we present is optimized for 3-prime Digital Gene Expression (DGE) libraries and can easily extend to full transcript coverage shotgun (SHO) type strand-specific libraries and is modularized to accommodate a diversity of RNA and DNA input materials. BrAD-seq offers a highly streamlined and inexpensive option for RNA-seq libraries.

## INTRODUCTION

Next Generation Sequencing (NGS) technologies have rapidly become foundational tools of genomics research (Koboldt et al., 2013). In particular, RNA-sequencing (RNA-seq) has transformed gene expression analyses and promoted the study of non-model organisms at an unprecedented level of detail with the ability to generate transcriptome assemblies for virtually any species (Sémon, 2014). On the most commonly used Illumina platform the ability to sequence a large number of biological samples requires the creation of libraries from nucleic acid samples with specified sequence "adapters" at the termini of the molecules. There are a variety of methods available to generate adapter-added libraries from nucleic acid samples from a variety of source materials, however the process still remains technically challenging, laborious, and expensive, thereby limiting widespread access to the technology.

37 Here we present a novel and efficient method for constructing strand specific RNA-seq libraries
38 in a simple, rapid, and inexpensive modular format. The method is optimized to create strand
39 specific 3-prime Digital Gene Expression (DGE – providing readout from the 3' end of the
40 mRNA) and can be adapted for strand-specific non-DGE shotgun type (SHO) and more
41 conventional non-strand specific (CNV) RNA-seq libraries, in addition to utilizing a variety of
42 DNA source materials. 3-prime DGE libraries are often preferred for gene expression studies
43 because a single mRNA yields approximately 1 sequence read reducing potential sources of bias.

44 Strand specific RNA-seq requires the directional addition of unique 5-prime and 3-prime adapter
45 sequences during preparation of the cDNA libraries. This is accomplished in a number of ways
46 among the various NGS library preparation protocols. These include, the ligation of a known
47 sequence to the 5-prime portion of mRNA molecules prior to cDNA synthesis (Lister et al.,
48 2008), removal of the template RNA strand followed by randomly primed 2$^{nd}$ strand synthesis
49 (Armour et al., 2009), labeling of first or second strand cDNA molecules with dUTP for
50 enzymatic degradation prior to enrichment (Parkhomchuk et al., 2009) and the use of terminal
51 transferases to add defined nucleotides to the cDNA molecules (Zhu et al., 2001;Tang et al.,
52 2010), with each method having advantages and shortcomings (Regev et al., 2012). Our method
53 for directional NGS library construction considerably simplifies and accelerates the library
54 construction process. Only around 10 milligrams of cytoplasmically dense plant tissue such as
55 Shoot Apical Meristem (SAM) or leaf primordia (slightly larger amounts for mature tissue), are
56 required for RNA-seq library production, and an individual worker can readily complete the
57 procedure starting from tissue in a single day.

58 We utilize an aspect of nucleic acid chemistry that has not been exploited in available methods to
59 generate strand specific libraries. Double stranded nucleic acids undergo a phenomenon called
60 "breathing" where the individual strands will momentarily separate to expose the bases (von
61 Hippel et al., 2013). This process happens at a higher rate at the ends of double stranded nucleic
62 acids (von Hippel et al., 2013). We exploit this transient terminal breathing to incorporate an
63 adapter oligonucleotide that includes the Illumina TruSeq PE1 sequence specifically at the 5-
64 prime terminus of the RNA-cDNA duplex. Breath capture allows for streamlined strand-specific
65 library protocols not requiring prior second strand synthesis or removal of template RNA,
66 allowing construction of either 3-prime DGE or shotgun (SHO) type strand specific libraries.

67 From these basic strand specific modules we further developed additional compatible modules to
68 accommodate a variety of nucleic acid species as input materials - single-stranded RNA, double-
69 stranded DNA and single-stranded DNA. This provides a general purpose platform for creation
70 of libraries for gene expression studies, genomic DNA libraries as well as from the products of
71 amplification of minute samples such as DNA obtained in Chromatin Immunoprecipitation
72 (ChIP) experiments and RNA from Laser Capture Microdissected (LCM) tissue samples. The
73 use of common modules in this platform minimizes the number of individual reagents required to
74 generate any number of library types, as well as standardizes the handling and manipulation
75 steps, reducing the learning curve and minimizing the potential for human error.

## MATERIALS AND METHODS

77 A schematic diagram of the reaction steps for strand-specific library synthesis is shown in Figure
78 1. Brief protocol for non-strand specific "conventional" (CNV) RNA-seq libraries can be found
79 in Supplementary methods 1. Detailed directions for strand specific DGE RNA-seq as well as

80  strand specific SHO RNA-seq and non-strand CNV RNA-seq and DNA-seq protocol variants
81  can be found in Supplementary methods 2.  All oligonucleotides used in this study were ordered
82  from Life Technologies (Thermo Fisher Scientific) at 50 nanomole scale, desalted with no
83  additional purification.

## Plant material

85  Tomato seeds (*S. lycopersicum* cv M82: LA3475) were provided by the Tomato Genetics
86  Resource Center, University of California, Davis. After sterilization (50% bleach for one minute
87  followed by rinse with water), seeds were placed onto water-soaked paper towels in Phytatrays
88  (Sigma) in the dark for three days at room temperature to allow germination. The germinated
89  seeds within Phytatrays were placed into a growth chamber at 22°C with 70% relative humidity
90  and a photoperiod of 16 h light/8 h dark for another four days. Seedlings were then transplanted
91  into Sunshine Mix soil (Sun Gro). After growing in soil for 11 days, P5 leaf primordia (the leaf
92  sample) and SAM (consisting of the SAM and 4 younger leaf primordia) were dissected
93  carefully using razor blades and harvested into RNase-free tubes.

## mRNA isolation
95  Tissues were processed and lysed as described previously by Kumar et al. (Kumar et al., 2012)
96  using zircon beads and Lysate Binding Buffer containing Sodium dodecyl sulfate in place of
97  Lithium dodecyl sulfate.  mRNA was isolated from 200 µl of lysate per sample.  1 µl of 12.5 µM
98  of 5-prime biotinylated polyT oligonucleotide containing a 5-prime 20 nucleotide arbitrary
99  spacer sequence followed by 20 thymidine nucleotides (5'-bio-
100 ACAGGACATTCGTCGCTTCCTTTTTTTTTTTTTTTTTTTT-3') was added to each lysate
101 sample, mixed by pipetting several times and allowed to stand for 10 minutes.  Following
102 incubation, captured mRNAs were isolated from the lysate by the addition of 20 µl of LBB
103 washed Streptavidin-coated magnetic beads (New England BioLabs, Cat. # S1420S).  The bead-
104 lysate mixture was mixed by pipetting and allowed to stand an additional 10 minutes.  Samples
105 were placed on a 96-well magnetic separator (Edge BioSystems, Cat. # 57624) and washed as
106 previously described (Kumar et al., 2012) with the following modifications.  A) Wash volumes
107 of WBA, WBB and LSB were 300 µl each and buffers were chilled on ice prior to use.  B)
108 mRNA elution was done into 16 µl of 10 mM TrisHCl pH 8 containing 1 mM β-
109 mercaptoethanol.

## mRNA fragmentation, 3-prime adapter priming and cDNA synthesis

111 mRNA fragmentation was accomplished using magnesium ions (3 mM) at elevated temperature
112 (Supplementary Figure 1) and was standardized at 90 seconds.  Priming for the cDNA synthesis
113 reaction was carried out in a single reaction mixture for Strand Specific-DGE, Strand Specific-
114 SHO, and non-Strand Specific libraries were fragmented in a reaction containing 1.5 µl 5X RT
115 buffer (Thermo scientific, Cat. # EP0441), 1 µl of priming adapter and 7.5 µl of the sample
116 mRNA in a total reaction volume of 10 µl.  Mixtures were spun down and incubated in a
117 thermocycler.  The following oligonucleotides and thermocycler programs were used for each
118 library type.  Additional details and comments can be found in Supplementary Methods 2.

119 **DGE**: 1 µl of 2 µM oligo L-3ILL-20TV.2  (5'-
120 GTGACTGGAGTTCAGACGTGTGCTCTTCCGATCTTTTTTTTTTTTTTTTTTTTV-3')
121 (25°C 1 second, 94°C 1.5 min, 30°C 1 min, 20°C 4 min, 20°C hold).

122  **SHO**: 1 μl of 5 μM oligo L-3ILL-N8.2 (5'-
123  GTGACTGGAGTTCAGACGTGTGCTCTTCCGATCTNNNNNNNN-3') (25°C 1 second,
124  94°C 1.5 min, 4°C 5 min, 20°C hold).

**cDNA synthesis**

126  cDNA was synthesized by addition of 5 μl of the following reaction mixture to the fragmented
127  and primed mRNA: 1.5 μl 5X Thermo Scientific RT buffer (Thermo scientific, Cat. # EP0441),
128  1.5 μl 0.1M Dithiothreitol (DTT), 1 μl H2O, 0.5 μl 25mM dNTPs (Thermo Scientific, Cat. #
129  R1121), 0.5 μl RevertAid RT enzyme (Thermo scientific, Cat. # EP0441) (total reaction volume
130  15 μl).  The reaction mixture was set up at room temperature and placed in a thermocycler
131  running the following program: (25°C 10 min, 42°C 50 min, 50°C 10 min, 70°C 10 min, 4°C
132  hold).  cDNA was cleaned and size-selected prior to "breath capture" or second strand synthesis
133  by addition of 5 μl 50 mM EDTA pH 8.0 and 30 μl Agencourt AMPure XP beads (Beckman,
134  Cat. # A63881) to each sample and mixed by pipetting.  After 5 minutes, samples were placed on
135  a magnetic tray, supernatant was removed, and pellets were washed twice with 300 μl 80%
136  ethanol without pellet disruption.  Residual ethanol was removed with 20-μl pipette tip and
137  samples were allowed to air-dry until no visible traces of liquid were detectable.

**5-prime duplex breath capture adapter addition (strand specific)**

139  5-prime adapter addition was done by rehydrating the cDNA bound to bead-pellet with 4 μl 10
140  μM pre-annealed 5-prime double stranded adapter oligo at room temperature.  Double stranded
141  5-prime adapter was prepared by making a stock solution containing 10 mM each of oligos
142  5pSense8n and 5pAnti (5pSense8n 5'-CCTACACGACGCTCTTCCGATCTNNNNNNNN-3',
143  5pAnti 5'-AGATCGGAAGAGCGTCGTGTAGG-3') in H$_2$O, dispensing to 100 μL volumes in
144  strip tubes and annealing them in a thermocycler running the following program: (94 C, 1min (94
145  C, 10 sec) X 60 cycles -1 C/cycle, 20 C 1 min, 4 C hold).  Subsequently, 6 μl of the following
146  reaction mixture was added, mixed by pipetting to fully re-suspend the pellet and incubated at
147  room temperature for 15 minutes: 3.5 μl H$_2$O, 1 μl 10X Thermo Pol I reaction buffer (Thermo
148  Scientific, Cat. # EP0041), 1 μl 250 mM MgCl$_2$ (made fresh and stored at -20 C), 0.25 μl 25 mM
149  dNTPs (Thermo Scientific, Cat. # R1121), 0.25 μl Thermo DNA Pol I (Thermo Scientific, Cat. #
150  EP0041) (10 μl total reaction volume).  The pre-enrichment libraries on beads were washed and
151  size-selected using Agencourt AMPure XP beads present from the previous step by adding 10 μl
152  50 mM EDTA pH 8.0 and 30 μl ABR solution (15% PEG 8000, 2.5 M NaCl), mixed thoroughly
153  by pipetting and allowed to stand for 5 minutes prior to placing on the magnetic tray.
154  Supernatant was removed and pellets were washed twice with 300 μl 80% ethanol, without pellet
155  disruption.  Residual ethanol was removed with 20-μl pipette tip and samples were allowed to
156  air-dry until no visible traces of liquid were detectable. Pellets were re-suspended in 22 μl 10mM
157  Tris pH 8.0, allowed to stand 1 minute and place on the magnetic tray.  Supernatant was
158  transferred without beads to fresh strip tubes and stored at -20°C prior to enrichment.

159  For Conventional library steps please see Supplementary Methods 1 or the detailed protocol in
160  Supplementary methods 2.

**PCR enrichment and index sequence addition (strand-specific and non-strand-specific)**

162  The enrichment step was done using full length oligonucleotides containing the full adapter
163  sequence as well as short oligonucleotides complementary to the distal-most portion of the
164  adapter arms to ensure predominantly full-length amplification products.  PCR enrichment was
165  carried out by combining 1 μl of the 2 μM uniquely-indexed ILL-INDEX oligonucleotide (ILL-
166  INDEX 5'-
167  CAAGCAGAAGACGGCATACGAGATxxxxxxxxGTGACTGGAGTTCAGACGTGTGCTCT
168  TCCGAT-3') (Supplementary table 1: Oligonucleotide sequences) with 9 μL of the master mix:
169  4 μl 5X Phusion HF Buffer, 2.6 μl $H_2O$, 1 μl 2 μM PE1 primer (PE1 5'-
170  AATGATACGGCGACCACCGAGATCTACACTCTTTCCCTACACGACGCTCTTCCGATC
171  T-3'), 1 μl 8 μM each S1 + S2 primers (S1 5'-AATGATACGGCGACCACCGA-3',  S2 5'-
172  CAAGCAGAAGACGGCATACGA-3'),  0.2 μl 25mM dNTPs, 0.2 μl Phusion Polymerase
173  (Thermo scientific, Cat. # F-530L) and 10 μl of pre-enrichment cDNA in a total reaction volume
174  of 20 μl. Half of the PCR mix (10 μl) was placed in separate sample tubes stored at -20 C as
175  backup for samples where more cycles of enrichment were needed.  The remaining 10 μl were
176  spun down and placed in a thermocycler using the program: (98 C 30 seconds, (98 C 10 seconds,
177  65 C 30 seconds, 72 C 30 seconds) 11 cycles, 72 C 5 min, 10 C hold).  Samples showing only
178  very faint enrichment were re-amplified with 14 cycles of enrichment from the backup PCR
179  samples.  2 μl of each library sample was run on a 1% agarose gel, with 1 μl of O'GeneRuler 100
180  bp DNA ladder (Thermo Scientific, Cat. # SM1143) for size and quantity reference, at 100 volts
181  for 20 minutes.  The remaining 8 μl of enriched library sample was cleaned and size selected
182  using 12 μl of fresh Agencourt AMPure XP beads and washing twice with 80% ethanol as in
183  previous wash steps.  The libraries were eluted from the pellet with 10 μl 10mM Tris pH 8.0,
184  quantified, and pooled as previously described (Kumar et al., 2012).  50 bp single end
185  sequencing was carried out at the Vincent J. Coates Genomic sequencing Facility at UC
186  Berkeley.

187  **Bioinformatics**

188  Bioinformatics and statistical analysis was carried out using the iPlant Atmosphere cloud service
189  (Goff et al., 2011).  Reads were trimmed to 42bp and quality filtered using FASTX-Toolkit
190  (http://hannonlab.cshl.edu/fastx_toolkit/) and scripts developed by Comai lab, UC Davis
191  (http://comailab.genomecenter.ucdavis.edu/).  Reads were mapped using Bowtie (Langmead et
192  al., 2009) with the parameters specified in Supplementary Table 2.  Read quality analysis was
193  performed using FASTQC (http://www.bioinformatics.bbsrc.ac.uk/projects/fastqc/).  The code
194  that was used to perform each of the bioinformatic steps is available at https://github.com/plant-
195  plasticity/townsley-fips-2015/ and FASTQ files for RNA-seq data used in this study can be
196  downloaded from Dryad data repository (link can only be provided in proof due to Dryad data
197  hosting policies).

198  **RESULTS AND DISCUSSION**

199  To evaluate our strand-specific library preparation method, we prepared Shoot Apical Meristem
200  (SAM) and leaf primordium (Leaf) samples using the new BrAD-seq DGE method and our
201  previously-developed HTR method for a pairwise comparative analysis. In this protocol we add
202  sample-identifying index sequences to the library molecules during the enrichment stage (Meyer
203  and Kircher, 2010).

204

**Library enrichment**

Although as a matter of procedure we don't typically quantify mRNA concentration prior to library synthesis to maintain higher throughput, when beginning experiments with unfamiliar materials it can be of utility to have some idea how many enrichment cycles would be reasonable to try.  To ascertain the relationship between the input mRNA concentration and the number of enrichment cycles chosen, 22 mRNA samples which were used for DGE library synthesis were quantified on a Bioanalyzer using the RNA 6000 Pico kit (Agilent Technologies).  This information was correlated with the number of cycles used for enrichment of each library sample and the concentration of washed libraries (Supplementary figure 2).  The relationship suggests that below about 10ng/μl of mRNA it may be worthwhile to start with about 14 enrichment cycles at the first attempt, although individual preferences in interpretation of gel images and targeted final concentrations for pooling of samples will ultimately be important factors in deciding on the ideal number of enrichment cycles.

**Read Quality**

To avoid inclusion of sequence originating from the 5-prime adapter capture strand, the first 8 bases of DGE libraries was trimmed prior to analysis.  For HTR libraries the percentage of reads mapping was also found to be higher (77.8% vs. 74.1%) when the first 8 bases were trimmed, so for all analyses trimmed FASTQ files were generated for samples prior to the quality filtering step.  The mapping rate improves in trimmed HTR libraries because during cDNA synthesis random primers anneal with mismatches, incorperating non-native sequence into cDNA molecules.

The overall quality scores for the raw DGE libraries was lower than HTR (Supplementary figure 3) due to the inclusion of cDNA inserts containing polyA tracts.  These low complexity sequences cannot be mapped to reference sequences and they are largely removed prior to mapping by quality filtering (Figure 2A and Supplementary Figure 3).

Since a population of strand-specific cDNA molecules highly enriched at the 3-prime of mRNA transcripts should be comprised of a smaller number of unique sequences for each transcript, identical reads from independent cDNA molecules are expected at a higher level than in non-strand-specific and non-DGE libraries.  We do indeed observe higher sequence duplication for DGE and strand specific library types than HTR (Figure 2B).  Non-DGE strand specific libraries have fuller transcript length coverage and show lower sequence duplication than DGE libraries resulting from higher sequence complexity (Supplementary figure 4).  Strand specific tomato SHO libraries made from similarly staged developing tomato leaves and Arabidopsis strand specific libraries(Hsu et al., 2013) downloaded from the Gene Expression Omnibus (Acession: GSE38879) made using a deoxy-Uracil (dU) marked strand specific method(Wang et al., 2011) were also assessed and possess similar rates of duplication to one another (Supplementary figure 4).  To remove differences in sequencing depth between samples as a factor in read duplication counts a random subsample of 1 million reads was used from each FASTQ file for duplication analysis.

244 Additionally, in 3-prime DGE libraries not all poly-A runs are removed by quality filtering.
245 Homonucleotide "A" repeats make up the predominant duplicated sequences in DGE libraries,
246 comprising ~0.3% of quality filtered reads. After quality-filtering, GC content and per base
247 sequence content differ between DGE and HTR (Figure 2C) with lower GC content in strand-
248 specific DGE library reads. Wheras individual base compositions in non-strand-specific libraries
249 (e.g., HTR libraries) should contain roughly equal amounts of G to C and A to T nucleotides,
250 G/C and A/T ratios are unequal for the coding strand of mRNAs. The proportions of each
251 nucleotide in the sense strand of annotated tomato coding sequences is 22.1% G, 18.5% C,
252 29.9% A, 29.4% T. This closely matches the observed proportions in the DGE sequences: 22.5%
253 G, 15.2% C, 28.5% A, 33.8% T (Figure 2D). Quality scores, sequence content and GC
254 distribution show similar performance between SHO and dU library methods (Supplementary
255 figure 5)

**Adapter and rRNA contamination**

257 Adapter contamination was higher in DGE libraries than in HTR (Figure 3A) consisting of ~5%
258 of reads in DGE compared with ~1% of reads in HTR. This may be due to the use of higher
259 PEG concentrations in the the bead washing step in the DGE protocol. This could increase bead
260 binding of small products. Approximately 1% of reads from DGE libraries could be attributed to
261 ribosomal contamination compared with 0.22% to 0.39% in HTR libraries (Figure
262 3B) and approximately 3% in tomato libraries made with a commercial Illumina kit (Kumar et
263 al., 2012). Increased rRNA in DGE compared to HTR is lilely due to single step mRNA
264 isolation compared to two stage mRNA re-isolation in the HTR process.
265
266 **Read mapping**

267 To reliably compare DGE and HTR libraries we created a set of reference sequences consisting
268 of the annotated tomato coding sequence plus an additional downstream portion corresponding to
269 the genome sequence 3-prime to the stop codon. Plant 3-prime untranslated regions (3'-UTRs)
270 are variable in length and average around 200 bp (Mignone et al., 2002) but many 3'-UTRs are
271 not annotated. For the purpose of this study 500 bp of downstream genomic sequence was
272 chosen to encompass most 3'-UTR sequences and appended to the annotated ITAG2.4 coding
273 sequences (ITAGcds+500). An additional mapping reference was generated specifically for
274 DGE libraries consisting of the 3-prime 500 bp of the coding sequence plus an additional 500 bp
275 representing the 3'-UTR (ITAG500+500) to minimize the effect of mis-priming of the 3-prime
276 polyT containing adapter onto any A-rich regions within coding sequences.

277 The proportion of reads mapping one or more times to the plus and minus strands of the
278 ITAGcds+500 reference is higher in DGE (85-87%) than HTR (77-78%) (Figure 3C)
279 demonstrating that a large majority of reads in both methods originate from mRNAs.

**DGE 3-prime selectivity**

281 There is a strong selectivity of the DGE library protocol for the 3-prime portion of mRNA
282 transcripts wheras reads derived from HTR are more evenly distributed across transcripts.
283 (Supplementary Figure 6). Although the ITAG500+500 reference sequences are, on average, 608
284 bp shorter than the ITAGcds+500 reference sequences, more DGE reads map uniquely and

285  strand-specifically to the ITAG500+500 reference (78% to 81%) than the HTR reads mapping
286  uniquely to the ITAGcds+500 reference (73% to 78%).

**Strand specifity**

288  To evaluate strand-specificity of the DGE libraries, reads were mapped to tomato coding
289  sequences only (Figure 3D) to exclude reads mapping to overlapping UTR regions.
290  Approximately 99% of mapped reads in DGE libraries and 50% of mapped reads in HTR
291  libraries localize to the sense strand, indicating a very high degree of strand-specificity for the
292  DGE libraries.  Directional information of the cDNA molecule is preserved because only the
293  cDNA strand of the RNA-cDNA duplex can serve as a template for Pol I.  We have successfully
294  produced libraries using this method with E. coli Pol I, Klenow fragment, and Klenow exo-
295  (Supplementary Figure 1C) indicating the exonuclease activity of Pol I is not required for the
296  process to work efficiently.

297  A large majority of uniquely mapped reads (95%) in the DGE libraries map to a region +/- 500
298  bp of the annotated stop codons of ITAGcds+500 reference (Table 1), whereas HTR libraries
299  show a more even distribution across the transcript (Figure 4A).  The DGE reads localize almost
300  entirely to the 3-prime region of the transcript including downstream of the annotated stop
301  codon, suggesting that only this interval is necessary for mapping DGE reads.  HTR reads by
302  comparison show a more even distribution but still bias toward sequence at the 3-prime of the
303  transcript.  Since not all coding sequences are 1kb or longer the read locations were also scaled
304  to the portion of the coding sequence (Figure 4B).  HTR libraries still show a slight bias for
305  sequences near the 3-prime end of the CDS.  SHO libraries show similar transcript coverage to
306  HTR although SHO coverage shows somewhat higher 5-prime transcript representation
307  (Supplementary figure 7).

308  To ascertain the degree of sequence selection bias introduced by the adapter capture process, 20
309  nucleotides upstream of the first mapped nucleotide for each read was extracted from the FASTA
310  mapping reference for base composition (Figure 4C) and information content (Supplementary
311  Figure 8).  Positions -8 through -1 correspond to the cDNA region annealed to the 8 bp single
312  stranded portion of the adapter responsible for breath capture of the DNA-RNA duplex.
313  Positions -20 through -9 correspond to the "shielded" double stranded portion of the adapter
314  containing the Illumina TruSeq PE1 sequence.  Despite the presence of the shielding
315  oligonucleotide, the positions approaching the -9 map location corresponding to the last few
316  bases of the adapter show some sequence bias near the end of the double stranded region
317  (Supplementary figure 9). This suggests that duplex breathing of the adapter at the capturing end
318  transiently exposes the first few internal bases, allowing for increased interaction with cDNA
319  sequences with some complementarity.  While the degree and range of this sequence selection
320  bias is significantly improved over earlier versions of this protocol utilizing un-shielded single
321  stranded adapters, it may still be further improved by converting the first base of the random
322  8mer into an extended double-stranded shield region.  Retention of the template mRNA strand
323  prevents access to the interior portions of the cDNA. This restricts the interactions of the adapter
324  to the terminal portion of the cDNA, which provides control of library size through mRNA
325  fragmentation and limits the effects of sequence specific secondary structures.  Increasing
326  Magnesium concentration in the breath capture reaction to 20 mM improves library yield
327  (Supplementary Figure 1B) potentially through increased strength of base-pair interactions
328  between the cDNA strand and the capturing adapter.  The strand specificity of the DGE libraries

329    also allows for unambiguous assignment of the transcript of origin for genes in which the
330    terminator regions overlap (Supplementary Figure 8).

**Detection of gene expression**

332    Reads were analyzed from equally-sized subsets of pre-quality-filtered reads (Table 2).  The
333    number of transcripts with mapped reads is reduced in both DGE and HTR libraries when
334    excluding non-uniquely-mapped reads.  The limited span of the transcript incorporated into DGE
335    libraries, in combination with retaining only uniquely mapped reads and strand specificity may
336    reduce the false detection of transcripts where genomic locations of transcripts overlap and
337    where coding sequences are highly conserved.

338    Correlation between replicates is higher for DGE than HTR samples (Figure 5 and
339    Supplementary Table 3).  R-squared values for all pairwise comparisons of Log2-transformed
340    expression showed higher correlation between DGE (SAM 0.96, Leaf 0.95) replicates than HTR
341    (SAM 0.91, Leaf 0.93).  These values are also similar for DGE and *Arabidopsis* dU libraries
342    (0.96) as well as between HTR and SHO (0.92).  Variation between DGE and HTR experimental
343    samples was also assessed using multidimensional scaling (MDS) (Figure 6A).  Both DGE and
344    HTR samples cluster by tissue type although distance between SAM and Leaf clusters is greater
345    along dimension 2 for DGE libraries suggesting a high power of discrimination between tissues
346    by gene expression.  Differential gene expression calls between DGE and HTR show a high
347    degree of overlap (Table 3).  We found very strong correlation ($r_s = 0.92$) between the $\log_2$ fold-
348    change of genes that are differentially regulated (FDR $< 0.05$) in SAM vs. leaf samples for both
349    library preparation methods. The correlation remains very strong when considering genes
350    differentially regulated for only the DGE method ($r_s = 0.87$; orange in Figure 6B) or only the
351    HTR method ($r_s = 0.87$; blue in Figure 6B).

352    To compare within and across method differential expression results, we divided the samples
353    into ten groups of two replicates. The ten sample groups were: 2 HTR leaf, 2 HTR SAM, 3 DGE
354    leaf, and 3 DGE SAM. Within each library preparation method, we performed differential gene
355    expression analysis for all combinations of leaf x SAM. This resulted in 4 comparisons for HTR
356    and 9 for DGE. With these, we were able to calculate Spearman's Ranked Correlation
357    Coefficient for all combinations of leaf-SAM differentially expressed genes within (45 for DGE
358    and 6 for HTR) and between (36 for DGE vs. HTR) each library preparation method
359    (Supplementary figure 10).  We found that although the fold-change of differentially regulated
360    genes is less correlated when comparing between library preparation methods than within, both
361    between- and within-method comparisons show very strongly correlation.

**Cost**

363    We sought to minimize library prep cost and complexity by developing a protocol that uses
364    mostly unmodified oligonucleotides and minimizes handling, steps, and reagents.  The cost of
365    isolating mRNA and making strand-specific libraries with this method is extraordinarily low,
366    with magnetic bead, dNTP, and enzyme costs totaling \$2.96/sample including mRNA isolation
367    or \$1.98 if making libraries from mRNA.  Even allowing for the additional cost of consumables,
368    chemical reagents and an extra 10% volume for reaction master mixes, this method provides a
369    20-40 fold cost reduction over available commercial strand-specific methods (e.g. NEBNext®

370     Ultra™ Directional RNA Library Prep Kit for Illumina® 96 reactions Cat. # E7420L, SureSelect
371     Strand Specific RNA-Seq Library Preparation kit for 96 samples reactions Cat. # G9691A).

372     **Protocol development**

373     We had initially set out to modify a template switching protocol, but ended up making a
374     discovery that would enable us to create arguably the cheapest and fastest RNA-seq protocol to
375     date.  Our original goal was to try to use adapter-encoded index sequences together with barcode
376     sequences within the primary reads to achieve extremely dense multiplexing of samples.  The 5-
377     prime adapters were designed as single-stranded molecules with a partial Illumina PE1 sequence
378     followed by a 9-base-pair sequence (a 6 base pair barcode and 3 terminal guanines) to facilitate
379     base-pairing with non-templated cytosines added to the cDNA by MMLV polymerase. The
380     addition of adapter sequence to the cDNA was done in a second reaction using *E. coli*
381     Polymerase I following a size-selection bead cleanup to avoid "background cDNA" composed of
382     adapter concatamers.

383     Our initial libraries showed a highly heterogeneous enrichment of identical pooled test mRNA
384     dependent on the barcode sequence contained in the adapter (Supplementary figure 11), with
385     significant visible banding due to massive overrepresentation of specific amplicons which vary
386     with the adapter barcode sequence.  Following trimming of the first 9 nucleotides from the
387     Illumina reads, mapping to tomato transcripts, and clustering of samples unexpectedly showed
388     grouping based on barcode sequence and not on sample type (Supplementary figure 12).
389     Additionally, in the first attempt libraries only a small numbers of transcripts accounted for the
390     majority of read counts.

391     Further investigation of these unexpected results showed that, while cDNA libraries that could be
392     sequenced on the Illumina platform were produced, the priming mechanism did not utilize
393     template switching as originally envisioned.  Sequence analysis of the transcript reference
394     sequences located 5-prime to the first mapped nucleotide of the trimmed reads showed an
395     extreme bias in the sequenced tomato transcripts for nucleotides matching the barcode sequence
396     and "G" repeats (Supplementary figures 13, 14 and 15) and further upstream sequences
397     continued to include similarity to the PE1 sequence of the adapter.  This indicated that base-
398     pairing interactions between the terminal portion of the double-stranded cDNA and the barcode-
399     containing portion of the adapter were selecting the transcripts that would be represented in the
400     libraries.

401     Despite the rarity of any particular 9 base pair sequence in a given genome (one instance every
402     3.8e-06 bases), 74% of reads contained a perfect 9 base pair match to the barcode followed by 3
403     "G"s in the pre-trimmed portion of the read (Supplementary figure 15).  This showed that the
404     dominant template for the sequencing reaction was the strand primed from 3-prime end of the
405     adapter using the cDNA as a template.  Consequently, the addition of non-templated "C"s by
406     MMLV reverse transcriptase to the cDNA molecule likely blocked priming on the adapter
407     oligonucleotide forcing the majority of sequenced molecules to originate from the second strand.

408     This suggested that there was a breathing effect in the double stranded template. We redesigned
409     the 5-prime adapters to take advantage of this breath-capture effect and eliminate the sequence
410     biases created by our early adapters.  The portion of the adapter containing the Illumina PE1
411     sequence was shielded by annealing a complementary sequence oligonucleotide and the

412 following 9 bases were replaced with variable length extensions of random mixed-base
413 sequences, with extensions between 6 and 8 nucleotides outperforming shorter and longer
414 variants.  Adapter variants incorporating blocking groups at the 3-prime end of the random
415 nucleotide extension performed extremely poorly indicating that priming from this strand was
416 essential for library formation using this process.

417 **CONCLUSION**

418 We have developed a rapid and inexpensive method for making strand-specific 3-prime DGE
419 RNA-seq libraries from tissue in a multiplexed format. The entire process can be completed in a
420 single working day.  To our knowledge this is the first library construction process to utilize the
421 terminal breathing of nucleic acid duplexes to selectively and directionally add adapter
422 sequences.  We have further developed the process to include modules allowing the creation of a
423 variety of library types.  We have also used the core DGE method on a number of species in
424 addition to *S. lycopersicum* including *C. pentagona, S. pennellii, S. pimpinellifolium, S. neorickii*
425 and *N. tobacum.*  To date we have successfully used our DGE protocol to study differential gene
426 expression in a number of studies relating to development and abiotic stress with good results.
427 We have added and adapted modules to this core protocol for our own purposes and we provide
428 those modules as well so that others can also use this protocol as the basis for a universal RNA
429 and DNA-seq library protocol family.  In the hope of helping to democratize NGS sequencing
430 technologies we offer an inexpensive and easily implemented protocol for the preparation of
431 NGS libraries.

443 **REFERENCES**This work used the Vincent J. Coates Genomics Sequencing Laboratory at UC Berkeley,
444     supported by NIH S10 Instrumentation Grants S10RR029668 and S10RR027303.
445 Armour, C.D., Castle, J.C., Chen, R., Babak, T., Loerch, P., Jackson, S., Shah, J.K., Dey, J., Rohl, C.A.,
446     Johnson, J.M., and Raymond, C.K. (2009). Digital transcriptome profiling using selective hexamer
447     priming for cDNA synthesis. *Nature Methods* 6**,** 647-U635.
448 Goff, S.A., Vaughn, M., Mckay, S., Lyons, E., Stapleton, A.E., Gessler, D., Matasci, N., Wang, L., Hanlon,
449     M., Lenards, A., Muir, A., Merchant, N., Lowry, S., Mock, S., Helmke, M., Kubach, A., Narro, M.,
450     Hopkins, N., Micklos, D., Hilgert, U., Gonzales, M., Jordan, C., Skidmore, E., Dooley, R., Cazes, J.,
451     Mclay, R., Lu, Z., Pasternak, S., Koesterke, L., Piel, W.H., Grene, R., Noutsos, C., Gendler, K., Feng,
452     X., Tang, C., Lent, M., Kim, S.-J., Kvilekval, K., Manjunath, B.S., Tannen, V., Stamatakis, A.,
453     Sanderson, M., Welch, S.M., Cranston, K.A., Soltis, P., Soltis, D., O'meara, B., Ane, C., Brutnell, T.,
454     Kleibenstein, D.J., White, J.W., Leebens-Mack, J., Donoghue, M.J., Spalding, E.P., Vision, T.J.,

455         Myers, C.R., Lowenthal, D., Enquist, B.J., Boyle, B., Akoglu, A., Andrews, G., Ram, S., Ware, D.,
456         Stein, L., and Stanzione, D. (2011). The iPlant collaborative: cyberinfrastructure for plant biology.
457         *Frontiers in Plant Science* 2.
458 Hsu, P.Y., Devisetty, U.K., and Harmer, S.L. (2013). Accurate timekeeping is controlled by a cycling
459         activator in Arabidopsis. *Elife* 2.
460 Koboldt, D.C., Steinberg, K.M., Larson, D.E., Wilson, R.K., and Mardis, E.R. (2013). The Next-Generation
461         Sequencing Revolution and Its Impact on Genomics. *Cell* 155**,** 27-38.
462 Kumar, R., Ichihashi, Y., Kimura, S., Chitwood, D.H., Headland, L.R., Peng, J., Maloof, J.N., and Sinha, N.R.
463         (2012). A high-throughput method for Illumina RNA-Seq library preparation. *Frontiers in Plant*
464         *Science* 3.
465 Langmead, B., Trapnell, C., Pop, M., and Salzberg, S.L. (2009). Ultrafast and memory-efficient alignment
466         of short DNA sequences to the human genome. *Genome Biology* 10.
467 Lister, R., O'malley, R.C., Tonti-Filippini, J., Gregory, B.D., Berry, C.C., Millar, A.H., and Ecker, J.R. (2008).
468         Highly integrated single-base resolution maps of the epigenome in Arabidopsis. *Cell* 133**,** 523-
469         536.
470 Mignone, F., Gissi, C., Liuni, S., and Pesole, G. (2002). Untranslated regions of mRNAs. *Genome biology* 3**,**
471         REVIEWS0004-REVIEWS0004.
472 Parkhomchuk, D., Borodina, T., Amstislavskiy, V., Banaru, M., Hallen, L., Krobitsch, S., Lehrach, H., and
473         Soldatov, A. (2009). Transcriptome analysis by strand-specific sequencing of complementary
474         DNA. *Nucleic Acids Research* 37.
475 Regev, A., Levin, J.Z., and Yassour, M. (2012). Comprehensive comparative analysis of strand-specific
476         RNA sequencing methods. *ArrayExpress Archive*.
477 Sémon, S.P.a.M. (2014). Transcriptomics of developing embryos and organs: A raising tool for evo–devo.
478         *Journal of Experimental Zoology*.
479 Tang, F., Barbacioru, C., Nordman, E., Li, B., Xu, N., Bashkirov, V.I., Lao, K., and Surani, M.A. (2010). RNA-
480         Seq analysis to capture the transcriptome landscape of a single cell. *Nature Protocols* 5**,** 516-535.
481 Von Hippel, P.H., Johnson, N.P., and Marcus, A.H. (2013). Fifty Years of DNA "Breathing": Reflections on
482         Old and New Approaches. *Biopolymers* 99**,** 923-954.
483 Wang, L., Si, Y., Dedow, L.K., Shao, Y., Liu, P., and Brutnell, T.P. (2011). A Low-Cost Library Construction
484         Protocol and Data Analysis Pipeline for Illumina-Based Strand-Specific Multiplex RNA-Seq. *Plos*
485         *One* 6.
486 Zhu, Y.Y., Machleder, E.M., Chenchik, A., Li, R., and Siebert, P.D. (2001). Reverse transcriptase template
487         switching: A SMART (TM) approach for full-length cDNA library construction. *Biotechniques* 30**,**
488         892-897.

489

## Figure legends

491 **Figure 1** | Schematic diagram of strand-specific library synthesis mechanism. mRNAs are
492 fragmented by heat and magnesium (1) and primed for cDNA synthesis by an adapter-containing
493 oligonucleotide (2 and 3).  Size selection and cleanup removes unincorperated oligonucleotides
494 and small cDNA fragments (4).  Transient duplex breathing at the terminus of the RNA-cDNA
495 hybrid (5) facilitates interaction with the single-stranded portion of the 5-prime capturing adapter
496 (6) and *E. coli* DNA Polymerase I catalyses its incorporation into a complete library molecule
497 (7).

498 **Figure 2 |** Analysis of library quality and characteristics.   Percentage of reads passing all quality
499 filtering steps (A). Sequence duplication levels for DGE and HTR (B). GC content of reads in
500 DGE and HTR (C). The average GC content is lower and the distribution broader in DGE than
501 HTR.  The composition of individual nucleotides differs between the strand-specific DGE and
502 non-strand-specific HTR libraries (D). Sequence bias is more evident in the HTR libraries in the
503 first several positions of the trimmed quality-filtered reads.  Error bars reflect standard deviation
504 among samples separated by tissue and method (A) or by method (B,C)

505 **Figure 3 |** Read mapping and strand specifity.  Fraction of reads coming from adapter (A) and
506 ribosomal RNA (B) contamination.  Reads mapping to either strand of ITAGcds+500 reference
507 (C). Coding sequence mapped reads belonging to plus strand (D).

508 **Figure 4** | Transcript coverage and cDNA sequence selection bias.  Localization of DGE and
509 HTR reads within the mapping reference (A), DGE reads mapped to 1.5KB window localize
510 near the annotated stop codon.  Base frequencies for transcript nucleotides upstream of mapped
511 reads (B).

512 **Figure 5** | Log2-transformed expression correlations for representative sample pairs for each
513 sample DGE and HTR using a representative pair of samples for each.  Mean R-squared values
514 for all DGE and HTR.

515 **Figure 6 |** Multi Dimensional Scaling (MDS) plot for DGE and HTR SAM and Leaf samples
516 (A). SAM vs. Leaf Log$_2$ fold change comparison between DGE and HTR  (B).

517 **Table 1** | DGE read mapping location in ITAGcds+500 reference with respect to the stop codon.

| Fraction of mapped reads | Region of reference sequence |
|---|---|
| >50% | -60 to +120 |
| >75% | -150 to +200 |
| >85% | -250 to +250 |
| >95% | -500 to +500 |

518

519 **Table 2** | Transcript detection for pre-quality-filtered subsets of 6.5M reads each for DGE and HTR. Non-
520 uniquely mapping reads mapping to both strands of ITAGcds+500 reference (Blue), uniquely mapping
521 reads mapping to both strands of ITAGcds+500 (Purple) and uniquely mapping reads mapping to sense
522 strand of ITAG500+500 reference (Red).

523

| | | | non-uniquely mapping | | | Uniquely mapping | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | Mapping to both strands | | | | | | Maping to sense strand | | |
| | | | ITAGcds+500 | | | | | | ITAG500+500 | | |
| Combined sample | Innitial reads | Passing QF | Mapped | Percent mapped | Transcripts detected | Reads mapping | Percent QF reads mapped | Transcripts with hits | Reads mapping | Percent QF reads mapped | Transcripts with hits |
| DGE-SAM | 6,500,000 | 5,255,791 | 4,449,163 | 85 | 23,348 | 4,252,370 | 81 | 21,618 | 4,113,253 | 78 | 20,922 |
| DGE-Leaf | 6,500,000 | 5,230,179 | 4,442,859 | 85 | 23,395 | 4,232,606 | 81 | 21,574 | 4,117,670 | 79 | 20,893 |
| HTR-SAM | 6,500,000 | 5,745,924 | 4,508,993 | 78 | 24,931 | 4,355,096 | 76 | 22,999 | | | |
| HTR-Leaf | 6,500,000 | 5,741,410 | 4,447,320 | 77 | 24,526 | 4,280,954 | 75 | 22,627 | | | |

524
525 **Table 3 |** Differential gene expression calls for DGE and HTR library samples.

| FDR 0.05 | DGE Total | HTR total | DGE only | Both | HTR only |
|---|---|---|---|---|---|

| | | | | | |
|---|---|---|---|---|---|
| Up (S vs. L) | 2534 | 1386 | 1630 | 904 | 482 |
| Down (S vs. L) | 3014 | 1751 | 1615 | 1399 | 352 |
| FDR 0.01 | | | | | |
| Up (S vs. L) | 1766 | 722 | 1251 | 515 | 207 |
| Up (S vs. L) | 2376 | 1128 | 1413 | 963 | 165 |

526

## Supplementary figures

528 **Supplementary figure 1 |** RNA fragmentation by 3 mM magnesium at 94 degrees at increasing time
529 intervals (A). Effect on library output of MgCl concentration in breath capture reaction using *E. coli*
530 Polymerase I (B). Breath capture reaction is successfully facilitated by *E. coli* polymerase I (2.5 U),
531 Klenow fragment (1.25 U) and Klenow exo- (1.25 U) (C).  Lanes shown in C are 4,2 and 2 technical
532 replicates respectively. Breath capture reactions (B and C) were carried out at room temperature for 15
533 minutes.

534 **Supplementary figure 2** | RNA starting amounts vs library amplification, number cycles used and
535 concentration of washed libraries prior to pooling.

536 **Supplementary figure 3 |** Pre and post quality filtering PHRED scores for DGE and HTR libraries used
537 in this study.

538 **Supplementary figure 4** | Sequence duplication rates per million quality filtered reads.  High throughput
539 HTR 23.12 % (Black dashed), DGE 66.15% (Orange solid),  Shotgun (SHO) 53.63% (Yellow solid),
540 deoxy-Uracil marked (dU) 48.28% (Blue solid)

541 **Supplementary figure 5** | FastQC analytics on filtered read information for additional strand specific
542 library methods, Shotgun (SHO) (A,C,E) and deoxy-Uracil marked (dU) (B,D,F).  Quality scores (A,B),
543 Base composition (C,D), Percentage GC content (E,F).

544 **Supplementary figure 6 |** Genomic mapping location of uniquely mapped reads in DGE and HTR.  DGE
545 reads show predominant localization to 3-prime of transcripts.

546 **Supplementary figure 7** | Transcript coverage trace for SHO libraries

547 **Supplementary figure 8 |** Discrimination of read origin.  DGE reads can be positively assigned to their
548 transcript of origin when transcripts overlap or are in close proximity by strand specificity of the reads.

549 **Supplementary figure 9 |** Sequence logos displaying information content for 20 bases upstream of
550 mapped reads.

551 **Supplementary figure 10 |** Pairwise comparisons of differential gene expression shows higher
552 correlation within each method than between methods.

553 **Supplementary figure 11 |** Heterogeneous amplification from identical mRNA samples by single-
554 stranded adapters containing barcode sequences near the 3-prime end.

555 **Supplementary figure 12 |** Hierarchical clustering of library samples made with single stranded barcode
556 containing adapters shows grouping only by barcode sequence.

557 **Supplementary figure 13 |** Overrepresentation of reads mapping to positions containing Guanine repeats.

558 **Supplementary figure 14 |** Highly uneven distribution of mapping locations in libraries made with
559 prototype adapters.

560 **Supplementary figure 15 |** Sequence information content for reads upstream of the first mapping
561 nucleotide for the trimmed reads.

562 **Supplementary tables**

563 **Supplementary table 1** | Oligonucleotides

564       *See attached supplementary table 1 spreadsheet

565 **Supplementary table 2 |** User defined parameters for scripts used in this study.

| | |
|---|---|
| fastx_trimmer | -f 9 -Q 33 |
| trimFastqQuality.py | 20 35 |
| read_N_remover.py | |
| adapterEffectRemover.py | 41 |
| Bowtie: | |
| non-strand specific, non-uniquely mapped | -a --best --strata -v 1 -p 4 --sam --tryhard |
| non-strand specific, uniquely mapped | -a --best --strata -m 1 -v 1 -p 4 --sam --tryhard |
| strand specific, uniquely mapped | -a --best --strata --norc -m 1 -v 1 -p 4 --sam --tryhard |
| | |

566

567 **Supplementary table 3 |** R-squared values for all pairwise replicate sample comparisons log2
568 normalized read counts.

| Mean | HTR leaf | L_HTR_A 4 | L_HTR_B 5 | L_HTR_C 6 | L_HTR_D 7 |
|---|---|---|---|---|---|
| 0.9277 | L_HTR_A 4 | | | | |
| | L_HTR_B 5 | 0.9283 | | | |
| | L_HTR_C 6 | 0.9307 | 0.9201 | | |
| | L_HTR_D 7 | 0.9320 | 0.9229 | 0.9324 | |

| Mean | HTR SAM | S_HTR_A 5 | S_HTR_A 6 | S_HTR_B 6 | S_HTR_B 7 |
|---|---|---|---|---|---|
| 0.9064 | S_HTR_A 5 | | | | |
| | S_HTR_A 6 | 0.9111 | | | |
| | S_HTR_B 6 | 0.9086 | 0.9293 | | |
| | S_HTR_B | 0.9209 | 0.8887 | 0.8797 | |

| 7 | | | | |
|---|---|---|---|---|

| Mean | DGE Leaf | L_DGE_A 4 | L_DGE_B 5 | L_DGE_C 6 | L_DGE_D 7 | L_DGE_E 7 | L_DGE_E 8 | L_DGE_F 1 |
|---|---|---|---|---|---|---|---|---|
| 0.9523 | L_DGE_A 4 | | | | | | | |
| | L_DGE_B 5 | 0.9599 | | | | | | |
| | L_DGE_C 6 | 0.9605 | 0.9600 | | | | | |
| | L_DGE_D 7 | 0.9614 | 0.9603 | 0.9607 | | | | |
| | L_DGE_E 7 | 0.9572 | 0.9553 | 0.9519 | 0.9581 | | | |
| | L_DGE_E 8 | 0.9541 | 0.9541 | 0.9516 | 0.9578 | 0.9657 | | |
| | L_DGE_F 1 | 0.9360 | 0.9328 | 0.9347 | 0.9405 | 0.9429 | 0.9433 | |

| Mean | DGE SAM | S_DGE_A 5 | S_DGE_A 6 | S_DGE_B 7 | S_DGE_C 7 | S_DGE_D 8 | S_DGE_E 1 | S_DGE_F 2 |
|---|---|---|---|---|---|---|---|---|
| 0.9582 | S_DGE_A 5 | | | | | | | |
| | S_DGE_A 6 | 0.9591 | | | | | | |
| | S_DGE_B 7 | 0.9564 | 0.9599 | | | | | |
| | S_DGE_C 7 | 0.9567 | 0.9565 | 0.9518 | | | | |
| | S_DGE_D 8 | 0.9588 | 0.9594 | 0.9524 | 0.9608 | | | |
| | S_DGE_E 1 | 0.9564 | 0.9590 | 0.9522 | 0.9615 | 0.9631 | | |
| | S_DGE_F 2 | 0.9575 | 0.9614 | 0.9557 | 0.9594 | 0.9617 | 0.9623 | |

| Mean | dU | dU_1 | dU_2 | dU_3 |
|---|---|---|---|---|
| 0.9564 | dU_1 | | | |
| | A2 | 0.9564 | | |
| | A3 | 0.9582 | 0.9545 | |

| Mean | SHO | SHO_1 | SHO_2 | SHO_3 |
|---|---|---|---|---|
| 0.9225 2 | SHO_1 | | | |
| | SHO_2 | 0.926926 | | |

| SHO_3 | 0.920931 | 0.919703 | ███████ |

Mg++  Mg++

AAAAAAAAAAAAAAA

NNNNNNNN

V

**1. Heat RNA fragmentation**

Mg++  Mg++

AAAAAAAAAAAAAAA

NNNNNNNN

V

**2. 3-prime adapter priming**

AAAAAAAAAAAAAAA

NNNNNNNN

V

**3. cDNA synthesis**

AAAAAAAAAAAAAAA

NNNNN

V

**4. Bead cleanup and size selection**

AAAAAAAAAAAAAAA

V

**5. Duplex terminal breathing**

AAAAAAAAAAAAAAA

V

**6. 5-prime adapter breath capture**

NNNNNNNN

AAAAAAAAAAAAAAA

V

**7. Adapter incorporation**

NNNNNNNN

AAAAAAAAAAAAAAA

V

NNNNNNNN

AAAAAAAAAAAAAAA

V

NNNNNNNN

V

Figure 2.TIF

Figure 3.TIF

Figure 4.TIF

Figure 5.TIF

Figure 6.TIF



**Leading MDS**

**A**

**Log$_2$ fold change (HTR method)**

**B**

DE in both = 0.92
DE in DGE= 0.87
DE in HTR= 0.87