

1  
2  
3  
4 Data Note for: *GigaScience*

Corresponding author:

5  
6 Steven J. Castle  
7 USDA ARS  
8 U.S.-ALARC  
9 21881 N. Cardon Lane  
10 Maricopa, Az 85138  
11 Tel: 520-316-6338 Fax: 520-316-6300  
12 steven.castle@ars.usda.gov  
13  
14  
15  
16  
17  
18  
19  
20  
21

22 ***De novo* transcriptome assemblies of four xylem sap-feeding insects**

23  
24  
25  
26 Erica E. Tassone<sup>1</sup>, Charles C. Cowden<sup>2</sup> and S.J. Castle<sup>2\*</sup>  
27  
28  
29  
30  
31

32 <sup>1</sup> Plant Physiology and Genetics Research Unit, U.S. Arid Land Agricultural Research Center,  
33 USDA ARS, Maricopa, AZ 85138 USA  
34

35 <sup>2</sup> Pest Management and Biocontrol Research Unit, U.S. Arid Land Agricultural Research Center,  
36 USDA ARS, Maricopa, AZ 85138 USA  
37  
38  
39  
40

41 \*Corresponding Author  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

2 **Abstract**

3 **Background:** Spittle bugs and sharpshooters are well-known xylem sap-feeding insects and  
4 vectors of the phytopathogenic bacterium *Xylella fastidiosa* (Wells), a causal agent of Pierce's  
5 disease of grapevines and other crop diseases. Specialized feeding on nutrient-deficient xylem  
6 sap is relatively rare among insect herbivores, and only limited genomic and transcriptomic  
7 information has been generated for xylem-sap feeders. To develop a more comprehensive  
8 understanding of biochemical adaptations and symbiotic relationships that support survival on a  
9 nutritionally austere dietary source, transcriptome assemblies for three sharpshooter species and  
10 one spittlebug species were produced.

11 **Findings:** Trinity-based *de novo* transcriptome assemblies were generated for all four xylem-sap  
12 feeders using raw sequencing data originating from whole-insect preps. Total transcripts for each  
13 species ranged from 91,384 for *Cuernia arida* to 106,998 for *Homalodisca liturata* with transcript  
14 totals for *Graphocephala atropunctata* and the spittlebug *Clastoptera arizonana* falling in  
15 between. The percentage of transcripts comprising complete open reading frames ranged from  
16 60% for *H. liturata* to 82% for *C. arizonana*. BUSCO analyses for each dataset indicated quality  
17 assemblies and a high degree of completeness for all four species.

18 **Conclusions:** These four transcriptomes represent a significant expansion of data for insect  
19 herbivores that feed exclusively on xylem sap, a nutritionally deficient dietary source relative to  
20 other plant tissues and fluids. Comparison of transcriptome data with insect herbivores that  
21 utilize other dietary sources may illuminate fundamental differences in the biochemistry of  
22 dietary specialization.

23 **Keywords:** Transcriptome, RNA-seq, Trinity, Insect herbivory, Insect vector, Diet specialization

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

24 **Data description**

25 **Background**

26 Resource partitioning among herbivorous insects spans a continuum between specialists that feed  
27 on one or a few plant species to generalists that are able to utilize hundreds of species belonging  
28 to multiple plant families. A further element of plant partitioning involves the particular location  
29 on a plant or tissue type from which an insect feeds [1]. The diversity of plant feeding strategies  
30 has evolved along with specialized anatomical features such as mouthparts and digestive  
31 systems, unique enzyme complements for processing plant compounds, and partnerships with  
32 symbiotic microbiota that contribute to nutritional gain of the host insect. The transcriptome  
33 assemblies presented here include four species that feed exclusively on sap from xylem vessels, a  
34 relatively rare form of plant feeding from a source that among plant tissues is the most deficient  
35 in nitrogen and carbon content [2]. Three of the transcriptomes represent sharpshooter species  
36 that are members of the subfamily Cicadellinae (Cicadellidae) that belong to the superfamily  
37 Membracoidea (leafhoppers, sharpshooters, treehoppers). The fourth transcriptome represents a  
38 spittlebug (Clastopteridae) that belongs to the superfamily Cercopoidea (spittlebugs,  
39 froghoppers). All are members of the hemipteran suborder Auchenorrhyncha [3]. Their piercing-  
40 sucking mouthparts tap into xylem vessels from which sap is consumed in copious quantities to  
41 compensate for its low nutritional value. Sharpshooters are recognized for their efficient  
42 assimilation of limited nutrients in xylem sap [4], but putative biochemical mechanisms that  
43 enable specialization on xylem sap are unknown. Also unclear is whether the respective roles in  
44 host nutrition played by the dual primary endosymbionts are consistent among xylem feeders [5].  
45 Comparison of transcriptomes of four xylem-feeding insects will provide additional knowledge

1  
2  
3  
4 46 and insight into the survival of ecological specialists on a nutritionally impoverished dietary  
5  
6 47 source.

## 10 48 **Samples**

11  
12 49 The spittlebug *Clastoptera arizonana* Doering was collected in 2014 from a wild population  
13  
14 50 infesting grapevines in Yavapai County, Arizona and established as a glasshouse colony for eight  
15  
16 51 months prior to sample collection in Maricopa, AZ. Samples of the sharpshooter *Cuerna arida*  
17  
18 52 Oman and Beamer (tribe Proconiini) were collected in 2015 by sweep net from a wild population  
19  
20 53 in mixed vegetation in Cochise County, Arizona. The smoke-tree sharpshooter *Homalodisca*  
21  
22 54 *liturata* Ball (Proconiini) was collected in 2015 from *Euphorbia tirucalli* L. plants in Phoenix,  
23  
24 55 AZ. The blue-green sharpshooter *Graphocephala atropunctata* (Signoret) (Cicadellini) was  
25  
26 56 collected in 2013 from a wild population in Orange County, California and maintained as a  
27  
28 57 glasshouse colony on basil (*Osimum basilicum* L.) until samples were collected in 2015. Live  
29  
30 58 adults of unknown age from all four species were homogenized separately in RNAlater<sup>®</sup>  
31  
32 59 (Ambion/Life Technologies, Carlsbad, CA) and stored at -20°C. Total RNA extractions were  
33  
34 60 performed using a Qiagen RNeasy mini kit followed by quality assessment on an Agilent 2100  
35  
36 61 bioanalyzer. Library generation yielding 2 x 100 bp paired end reads (TruSeq RNA Sample  
37  
38 62 Preparation Kit v2; Illumina Inc., San Diego, USA) and sequencing (Illumina HiSeq2000 or  
39  
40 63 HiSeq2500) were performed at the University of Arizona Genomics Center in Tucson, AZ  
41  
42 64 (<http://uagc.arl.arizona.edu>).

## 52 65 **Data filtering**

53  
54 66 The total number of reads, data quantity, and short read archive (SRA) numbers for each of the  
55  
56 67 four xylem-feeding insects are shown in Table 1. For each data set, raw quality was assessed and  
57  
58 68 filtered using both FastQC and Trimmomatic (v 0.32) using the parameters

1  
2  
3  
4 69 ILLUMINACLIP:TruSeq3-PE.fa:2:30:10 LEADING:10 TRAILING:20

5  
6  
7 70 SLIDINGWINDOW:4:25 MINLEN:36 to remove adaptor sequence and filter by quality score.

8  
9  
10 71  (Table 1)

11  
12  
13 72 **Transcriptome assembly**

14  
15 73 All raw data for each insect transcriptome were run through the following pipeline. Prior to  
16  
17 74 assembly, the three replicate samples were concatenated and read abundance was normalized to  
18  
19 75 50X coverage using the *in silico* normalization tool in Trinity v. 2.0.6 [6] to improve assembly  
20  
21 76 time. Each of the datasets were assembled in Trinity using the default parameters, with the  
22  
23 77 addition of the ‘-jaccard clip’ flag to reduce the generation of transcript fusions from non-strand  
24  
25 78 specific data. Open reading frames (ORFs) were predicted using Transdecoder with all run  
26  
27 79 parameters set to default [6]. The transcriptomes were filtered, sorted, and prepared for NCBI  
28  
29 80 transcriptome shotgun assembly (TSA) submission as previously described [7]. To comply with  
30  
31 81 NCBI TSA submission, all transcripts resulting from endosymbionts or bacteria were removed  
32  
33 82 from the final assembly prior to submission.  
34  
35  
36  
37  
38  
39

40 83 **Annotation**

41  
42 84 Functional annotation for each of the transcriptomes was performed at the peptide level using a  
43  
44 85 custom pipeline [7] that defines protein products and assigns transcript names. Predicted proteins  
45  
46 86 and peptides were analyzed using InterProScan 5 [8], using the ‘-iprlookup’ and ‘-goterms’  
47  
48 87 flags, to search all available databases, including Gene Ontology (GO). Each transcriptome was  
49  
50 88 annotated using BLASTp against the UniProt Swiss Prot database (downloaded 11 February  
51  
52 89 2015). Annie [9], a program that cross-references SwissProt BLAST and InterProScan5 results to  
53  
54 90 extract qualified gene names and products, was used to generate the transcript annotation file.  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

91 The resulting .gff3 and .tbl files were further annotated with functional descriptors in  
92 Transvestigator [10].

93 **Transcriptome Quality and Comparisons**

94 Assembled transcriptome metrics showed a high percentage of reads mapping back to each  
95 transcriptome (Table 2) indicating successful assemblies. TransRate [11] scores ranging from  
96 0.16 to 0.42 were used for quality assessment, and BUSCO v. 1.1.b1 (benchmarking universal  
97 single-copy orthologs) results using the arthropod gene set (downloaded December 19, 2015)  
98 [12] indicated that the four transcriptomes have a moderate to high level of completeness. It  
99 should be noted that both the TransRate value (0.16) and BUSCO results for *H. liturata* suggest  
100 this transcriptome may contain more partial transcripts than the other three assemblies.

101 Each of the assembled transcriptomes was used in a reciprocal tBLASTx search to identify  
102 similarities between the four species and their transcriptome assemblies. The final, filtered  
103 transcriptomes were made into nucleotide BLAST databases using NCBI Blast+ (v 2.2.30)  
104 *makeblastdb* tool and all tBLASTx searches were performed using an e-value cutoff of  $1e^{-3}$ . The  
105 tBLASTx results (Table 3) indicate similarities between the four xylem-feeder transcriptomes,  
106 with the lowest (38%) occurring between members of the two superfamilies (spittlebug and all  
107 sharpshooters) and the highest (84%) between *H. liturata* and *C. arida*, members of the same  
108 subfamily [13].



(Tables 2 and 3)

110 **Availability of supporting data**

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

111 The filtered and annotated transcriptomes have been deposited in GenBank as a TSA under the  
112 accessions and BioProject numbers found in Table 1. Datasets further supporting the results of  
113 this article are available in the *GigaScience* repository, GigaDB [13].

114 **Abbreviations**

115 ALARC: Arid Land Agricultural Research Center; BUSCO: Bench-marking Universal Single-  
116 Copy Orthologs; GO: Gene Ontology; ORF: Open Reading Frame; SRA: Short Read Archive;  
117 TSA: Transcriptome Shotgun Assembly; USDA-ARS: United States Department of Agriculture-  
118 Agricultural Research Service

119 **Competing Interests**

120 The authors declare that they have no competing interests.

121 **Authors Contribution**

122 SJC and CCC conceived and performed the experiments; EET analyzed the data and evaluated  
123 the conclusions; EET, SJC, and CCC wrote the manuscript. All authors approved the final  
124 manuscript.

125 **Acknowledgements**

126 The authors thank Tom Perring and Darcy Reed at UC Riverside for providing samples of *G.*  
127 *atropunctata* for transcriptome evaluations. The authors acknowledge funding from the Arizona  
128 Department of Agriculture Grant No. SCBGP-FB 12-12. Bioinformatic analysis was performed  
129 on computing resources available at ALARC. Mention of trade names or commercial products in  
130 this article is solely for the purpose of providing specific information and does not imply  
131 recommendation or endorsement by the US Department of Agriculture. USDA is an equal  
132 opportunity provider and employer.

## References

1. Schoonhoven LM, Van Loon JJA, Dicke M. Insect-plant biology. Oxford University Press, New York. 2005, 421 pp.
2. Mattson WJ. Herbivory in relation to plant nitrogen content. *Ann. Rev. Ecol. Syst.* 1980;11:119-161.
3. Dietrich CH, Rakitov RA, Holmes JL, Black WC. Phylogeny of the major lineages of Membracoidea (Insecta: Hemiptera: Cicadomorpha) based on 28S rDNA sequences. *Mol. Phylogen. Evol.* 2001;18:293-305.
4. Brodbeck BV, Mizell RF, Andersen PC. Physiological and behavioral adaptations of three species of leafhoppers in response to the dilute nutrient content of xylem fluid. *J. Insect Physiol.* 1993;39:73-81.
5. Moran NA, Tran P, Gerardo NM. Symbiosis and insect diversification: an ancient symbiont of sap-feeding insects from the bacterial phylum *Bacteroidetes*. *Appl. Environ. Microbiol.* 2005;71:8802-8810.
6. Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, Bowden J, et al. *De novo* transcript sequence reconstruction from RNA-seq using the Trinity platform for reference generation and analysis. *Nat. Protoc.* 2013;8:1494-1512.
7. Sim SB, Calla B, Hall B, DeRego T, Geib SM. Reconstructing a comprehensive transcriptome assembly of a white-pupal translocated strain of the pest fruit fly *Bactrocera cucurbitae*. *Gigascience.* 2015;4:14.
8. Jones P, Binns D, Chang HY, Fraser M, Li W, McAnulla C, et al. InterProScan 5: genome-scale protein function classification. *Bioinformatics.* 2014;30:1236-40. Accessed 9 November 2015.
9. Tate R, Hall B, DeRego T. Annie the functional annotator - initial release. ZENODO. 2014. <http://doi.org/10.5281/zenodo.10470>. Accessed 27 November 2015.
10. DeRego T, Hall B, Tate R, Geib S. Transvestigator early release. ZENODO. 2014. <http://doi.org/10.5281/zenodo.10471>. Accessed 27 November 2015.
11. Smith-Unna RD, Bournsnel C, Patro R, Hibberd JM, Kelly S. TransRate: reference free quality assessment of de-novo transcriptome assemblies. *bioRxiv.* 2015. <http://dx.doi.org/10.1101/021626>. Accessed 17 September 2015.
12. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics.* 2015;31:3210-2. doi:10.1093/bioinformatics/btv351. Accessed 4 April 2016.



1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

13. Tassone EE, Cowden CC, and Castle SJ. Supporting data for “*De novo* transcriptome assemblies of four xylem sap-feeding insects.” *GigaScience* Database. 2016. <http://dx.doi.org/10.5524/1002578>.

**Table 1** Accession numbers for sequence reads and assembled transcripts for four species of xylem-feeding insects.

Sample	Reads	Size (Gb)	Short Read Archive	BioSample	BioProject
<i>Homalodisca liturata</i>	18,936,520	18.9	SRX1451710	SAMN04293489	PRJNA303151
			SRX1451711	SAMN04293490	“
			SRX1451712	SAMN04293491	“
<i>Clastoptera arizonana</i>	19,038,998	17.8	SRX1451715	SAMN04293493	PRJNA303152
			SRX1451717	SAMN04293494	“
			SRX1451718	SAMN04293495	“
<i>Cuerna arida</i>	14,667,040	18.3	SRX1451216	SAMN04292971	PRJNA303150
			SRX1451218	SAMN04292972	“
			SRX1451467	SAMN04292973	“
<i>Graphocephala atropunctata</i>	16,868,134	8.2	SRX1411425	SAMN04208332	PRJNA299492
			SRX1411426	SAMN04208333	“
			SRX1411427	SAMN04208334	“

**Table 2** Transcriptome assembly statistics and results of BUSCO analysis for four xylem-feeding insects.

	<i>H. liturata</i>	<i>C. arizonana</i>	<i>C. arida</i>	<i>G. atropunctata</i>
<b>Assembly</b>				
Normalized reads	9,468,260	9,519,499	10,714,375	32,429,458
Total no. transcripts	106,998	93,845	91,384	97,830
Average transcript length and range	954 (224 – 30,062)	1,232 (224 – 29,936)	901 (224 – 20,095)	962 (224 – 17,082)
Total assembled bases (all)	102,317,189	115,686,868	79,785,471	94,141,447
N50 (all)	1,650	2,510	1,560	1,692
% GC	37	31	37	39
% mapping	84	91	88	95
TransRate Score	0.16	0.28	0.25	0.42
<b>BUSCO</b>				
Complete (%)	60	82	68	66
Duplicated (%)	23	42	26	24
Fragmented (%)	23	9.2	17	19
Missing (%)	15	8	14	13

**Table 3** Total percent matches from tBLASTx reciprocal searches. Transcriptome used as query on the left, and nucleotide database tBLASTx against which the query was performed is shown at the top.

Query	Nucleotide Database (% similarity)			
	<i>H. liturata</i>	<i>C. arizonana</i>	<i>C. arida</i>	<i>G. atropunctata</i>
<i>Homalodisca liturata</i>	--	42.94	76.31	56.20
<i>Clastoptera arizonana</i>	40.90	--	38.35	38.58
<i>Cuerna arida</i>	83.90	43.40	--	58.36
<i>Graphocephala atropunctata</i>	56.86	40.30	56.10	--

e-value  $\leq 1^{E-3}$