

1 **Supplementary file 2 - Sub-assemblies of the *Echis coloratus* venom gland transcriptome**

2 In an attempt to determine the minimum required amount of sequencing to fully sequence and assemble the venom gland
 3 transcriptome of *Echis coloratus*, sub-sets of RNA-seq reads were extracted and assembled (Table S2). Paired venom
 4 gland reads were first interleaved using the `shuffleSequences.pl` perl script (part of the Velvet *de novo* assembly
 5 program [1]) so that each read pair was maintained during sub-sampling. Using the commands `head` and `tail`, 3 sub-
 6 sets (designated as “head”, “middle” and “tail”) of either 2, 4, 8 or 10 million reads were taken from an RNA-seq dataset
 7 containing 44,678,609 paired-end reads. These data were assembled using Trinity [2,3], with parameters set to run as a
 8 single-end read dataset (as there is only one .fastq input file), but with the added command-line parameter `--`
 9 `run_as_paired` to indicate that the data contains paired-end data.

10

11 Supplementary Table S2. Assembly metrics for sub-assemblies of the venom gland transcriptome of *Echis*
 12 *coloratus*.

| Sub-sample | Sample size (million reads) | Total number of contigs | Number of contigs ≥ 300 nt | Total length (nt) | Max. contig size (nt) | Contig N50 (nt) |
|----------------------------|-----------------------------|-------------------------|---------------------------------|-------------------|-----------------------|-----------------|
| H E A D | 2 | 24,585 | 14,744 | 10,302,850 | 7,474 | 808 |
| | 4 | 34,990 | 22,184 | 17,605,771 | 7,860 | 1,023 |
| | 8 | 45,207 | 30,121 | 27,542,537 | 11,824 | 1,293 |
| | 10 | 48,349 | 32,660 | 31,623,176 | 11,824 | 1,420 |
| M I D D L E | 2 | 23,915 | 14,229 | 10,036,594 | 7,840 | 837 |
| | 4 | 34,383 | 21,736 | 17,282,856 | 8,970 | 1,027 |
| | 8 | 44,759 | 29,946 | 27,116,697 | 11,738 | 1,279 |
| | 10 | 47,832 | 32,451 | 30,985,872 | 11,752 | 1,387 |
| T A I L | 2 | 24,170 | 14,513 | 10,059,952 | 8,547 | 810 |
| | 4 | 34,735 | 21,994 | 17,315,514 | 8,165 | 1,004 |
| | 8 | 44,956 | 29,988 | 27,283,356 | 11,803 | 1,284 |
| | 10 | 48,116 | 32,535 | 31,022,314 | 11,805 | 1,382 |

13

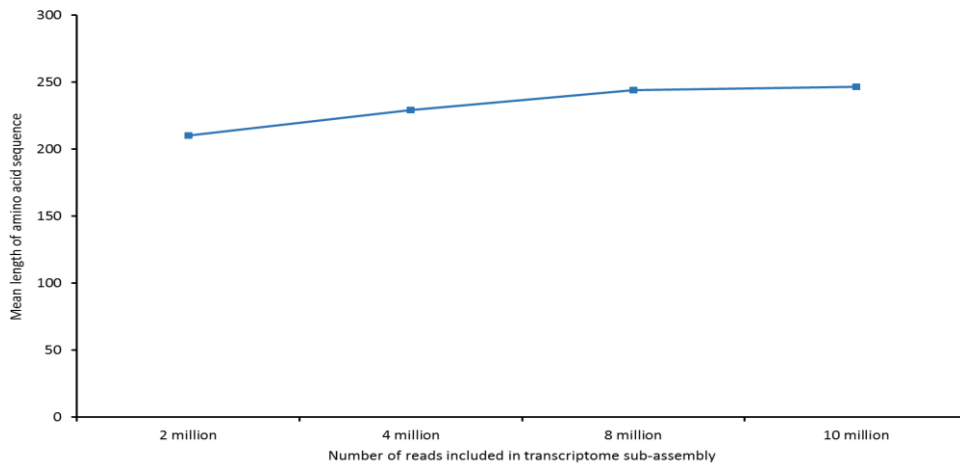
14

15 Local blast surveys were then carried out using BLAST+ version 2.2.27 [4] to identify previously characterised putative
 16 toxin genes in *E. coloratus*. The majority of transcripts encoding putative toxin genes appear to be present in venom gland
 17 transcriptome assemblies generated from only 2 million paired-end reads (here presence is defined as the transcript being
 18 found in all three (Head/Middle/Tail) sub-assemblies) (Table S3).

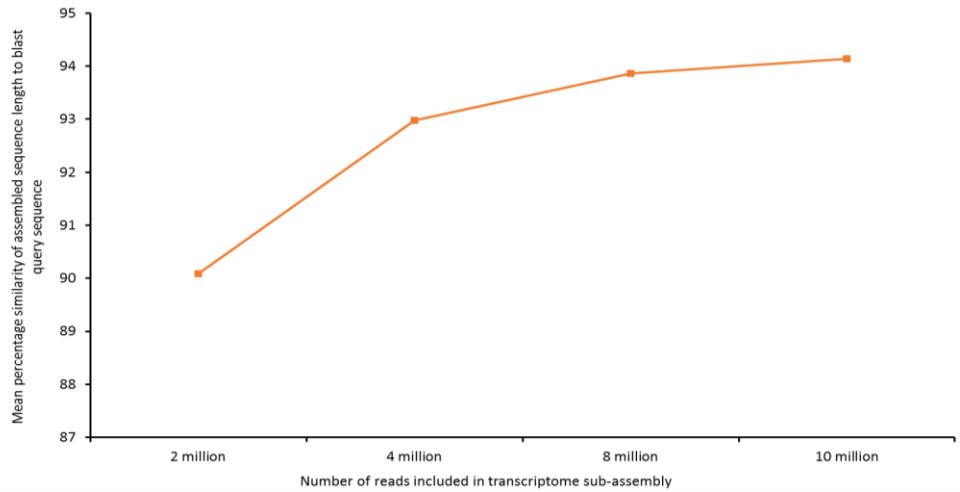
20 Supplementary Table S3. Presence/absence of putative toxin transcripts in sub-assemblies of the venom gland
 21 transcriptome of *Echis coloratus*. Detected transcripts are shaded, transcripts not found are shaded grey. H, head; M,
 22 middle; T, tail.

| Gene | Sub-sample size and position | | | | | | | | | | | |
|------------------------|------------------------------|---|---|-----------------|---|---|-----------------|---|---|------------------|---|---|
| | 2 million reads | | | 4 million reads | | | 8 million reads | | | 10 million reads | | |
| | H | M | T | H | M | T | H | M | T | H | M | T |
| 3ftx-a | + | - | + | + | + | + | + | + | + | + | + | + |
| 3ftx-b | + | + | + | + | + | + | + | + | + | + | + | + |
| ache - transcript 1 | - | + | - | + | + | + | + | + | + | + | + | + |
| complement c3 | + | + | + | + | + | + | + | + | + | + | + | + |
| crisp-b | + | + | + | + | + | + | + | + | + | + | + | + |
| crotamine-like | + | + | + | + | + | + | + | + | + | + | + | + |
| c-type lectins a-k | + | + | + | + | + | + | + | + | + | + | + | + |
| cystatin e/m | + | + | + | + | + | + | + | + | + | + | + | + |
| dpp 3 | + | + | + | + | + | + | + | + | + | + | + | + |
| dpp 4 | + | + | + | + | + | + | + | + | + | + | + | + |
| esp-e1 | + | + | + | + | + | + | + | + | + | + | + | + |
| ficolin | + | + | + | + | + | + | + | + | + | + | + | + |
| kallikrein | + | + | + | + | + | + | + | + | + | + | + | + |
| kunitz 1 | + | + | + | + | + | + | + | + | + | + | + | + |
| kunitz 2 | + | + | + | + | + | + | + | + | + | + | + | + |
| lao-a | + | + | + | + | + | + | + | + | + | + | + | + |
| lao-b1 | + | + | + | + | + | + | + | + | + | + | + | + |
| lao-b2 | + | + | + | + | + | + | + | + | + | + | + | + |
| lipa-a | + | + | + | + | + | + | + | + | + | + | + | + |
| lipa-b | - | - | - | + | + | + | + | + | + | + | + | + |
| ngf | + | + | + | + | + | + | + | + | + | + | + | + |
| PLA ₂ IIA-c | + | + | + | + | + | + | + | + | + | + | + | + |
| PLA ₂ IIA-d | + | + | + | + | + | + | + | + | + | + | + | + |
| PLA ₂ IIA-e | + | + | + | + | + | + | + | + | + | + | + | + |
| PLA ₂ IIE | - | - | - | - | + | - | - | + | + | - | - | + |
| plb | + | + | + | + | + | + | + | + | + | + | + | + |
| renin | + | + | + | + | + | + | + | + | + | + | + | + |
| serine protease a-f | + | + | + | + | + | + | + | + | + | + | + | + |
| svmp-a | + | + | + | + | + | + | + | + | + | + | + | + |
| svmp-b | + | + | + | + | + | + | + | + | + | + | + | + |
| svmp-c | - | - | + | - | + | + | + | + | + | + | + | + |
| svmp-d | - | - | - | - | + | + | + | + | + | + | + | + |
| svmp-e | + | + | + | + | + | + | + | + | + | + | + | + |
| svmp-f | + | + | + | + | + | + | + | + | + | + | + | + |
| svmp-g | + | + | + | + | + | + | + | + | + | + | + | + |
| svmp-i | + | + | + | + | + | + | + | + | + | + | + | + |
| svmp-j | + | + | + | + | + | + | + | + | + | + | + | + |
| svmp-k | + | + | + | + | + | + | + | + | + | + | + | + |
| svmp-m | + | - | + | + | + | + | + | + | + | + | + | + |
| svmp-n | + | + | + | + | + | + | + | + | + | + | + | + |
| svmp-p | + | - | + | + | + | + | + | + | + | + | + | + |
| svmp-q | + | + | + | - | - | + | + | + | + | + | + | + |
| svmp-t | + | + | + | + | + | + | + | + | + | + | + | + |
| vegf-a | - | - | - | + | - | + | + | + | + | + | + | + |
| vegf-c | - | - | - | - | - | - | + | - | + | + | + | + |
| vegf-f | + | + | + | + | + | + | + | + | + | + | + | + |
| waprin | + | - | + | + | + | + | + | + | + | + | + | + |

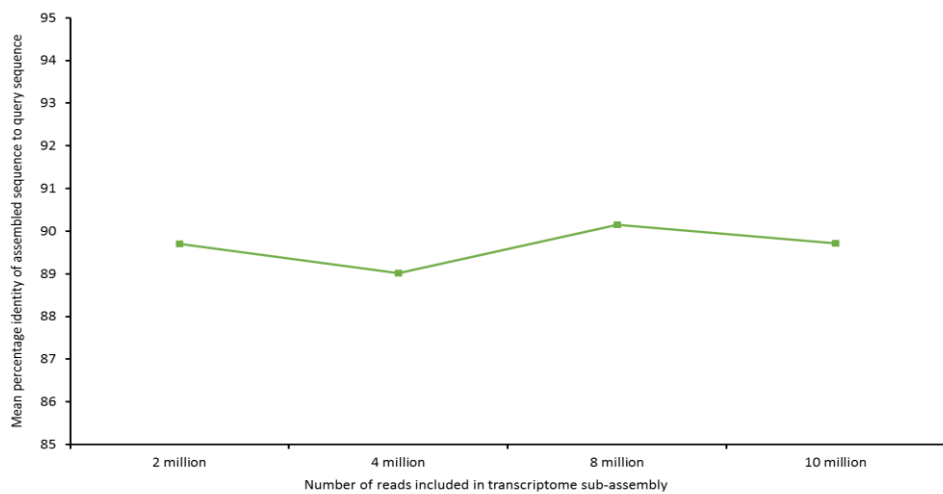
24 As the number of reads used for assembly increases the mean length of the amino acid sequence encoded by the assembled
25 transcript also increases, although there is only a 36 amino acid increase between 2 million and 10 million reads (Figure
26 S1).



27



28



29

30 Supplementary figure S1. Analysis of sequence assembly quality based on local blast surveys using previously
31 characterised amino acid sequences from *Echis coloratus* venom gland. Top - mean length of amino acid sequence
32 matches in sub-assemblies, Middle - mean percentage length of query sequence covered by assembled sequence. Bottom
33 - mean percentage similarity of assembled sequence to query sequence in sub-assemblies.

34

35 However, the number of contigs ≥ 300 bp roughly doubles (Table S1), meaning considerably fewer contigs which are
36 likely to be unplaced paired reads are present in the transcriptome assembly. To gain insight into how this increase in
37 length relates to the quality of the assembled toxin transcript sequences, the percentage of the query sequence covered by
38 the newly assembled sequence was calculated. Again there is only a minor improvement of 4% following an increase
39 from 2 million reads to 10 million (Figure S1). The mean percentage similarity between assembled sequence and query
40 sequence appears to be more variable across the sub-assemblies, with no apparent consistent improvement as the number
41 of reads increases (Figure S1). As the query sequences used for local BLAST searches were obtained from an assembly
42 of multiple *E. coloratus* venom gland datasets in order to represent an overabundance of sequencing, and the sub-
43 assemblies were assembled from a different set of venom gland reads, it should be expected that not all blast alignments
44 will have a 100% match between query and subject due to variation between individuals. However, a lower % identity
45 would indicate that either sequencing errors were incorporated into the assembly or there has been a misassembly, both
46 likely due to a reduced depth of sequencing coverage.

47

48 **Conclusion**

49 Around 8 million reads appears to be sufficient sequencing depth to capture all putative toxin-encoding transcripts to a
50 suitable assembly quality. The Illumina HiSeq2500 sequencing platform can currently produce 300-400 million 100nt
51 paired-end reads in “high output” mode, or 200-300 million 150nt paired-end reads in “rapid run” mode. With this in
52 mind, and 8 million paired-end reads assumed to be the minimum sequencing depth required to fully capture all putative
53 toxin transcripts, it is possible to sequence ~40 venom gland libraries on one sequencing lane of the Illumina HiSeq2500
54 (in “high output” mode).

55

56 **Supplementary references**

- 57 1. Zerbino DR, Birney E. 2008. Velvet: Algorithms for de novo short read assembly using de bruijn graphs. *Genome Res*
58 **18**: 821-829.
- 59 2. Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, Adiconis X, Fan L, Raychowdhury R, Zeng Q.
60 2011. Full-length transcriptome assembly from RNA-seq data without a reference genome. *Nat Biotechnol* **29**: 644-652.
- 61 3. Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, Bowden J, Couger MB, Eccles D, Li B, Lieber M. 2013.
62 De novo transcript sequence reconstruction from RNA-seq using the trinity platform for reference generation and analysis.
63 *Nat Protocols* **8**: 1494-1512.
- 64 4. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL. 2009. BLAST+: Architecture
65 and applications. *BMC Bioinformatics* **10**: 421.