

Complete *de novo* assembly of *Wolbachia* endosymbiont of contemporary *Drosophila simulans* using long-read genome sequencing.

Running Title: *De novo* assembly of *Wolbachia* in *Drosophila simulans*

**Authors:**

Jodie Jacobs<sup>a,b</sup>, Alexandra Lum<sup>a</sup>, Darren D. Lee<sup>a</sup>, Emry Gutierrez<sup>a</sup>, Jonah Dionisio<sup>a</sup>, Camryn Morey<sup>a</sup>, Cade Mirchandani<sup>a,b</sup>, Luke Sylvester<sup>a</sup>, Anne Nakamoto<sup>a,b</sup>, Hailey Loucks<sup>a,b</sup>, Ciara Wanket<sup>b,c</sup>, Ariana Cisneros<sup>a,b</sup>, Alessandro Calicchio<sup>a</sup>, Alexis N. Enstrom<sup>a</sup>, Camille Headrick<sup>a</sup>, Faith Okamoto<sup>a,b</sup>, Harrison Heath<sup>a,b</sup>, Kseniya Malukhina<sup>a</sup>, Petria Russell<sup>a</sup>, Sagorika Nag<sup>a</sup>, Thomas Gillespie<sup>a</sup>, William Sobolewski<sup>a</sup>, Zia Truong<sup>a</sup>, Shelbi L. Russell<sup>#a,b</sup>

<sup>a</sup>Biomolecular Engineering and Bioinformatics Department at the University of California, Santa Cruz

<sup>b</sup>Genomics Institute at the University of California, Santa Cruz

<sup>c</sup>Ecology and Evolutionary Biology Department at the University of California, Santa Cruz

<sup>#</sup>Corresponding Author: Shelbi Russell shelbilrussell@gmail.com

Authorship order was determined based on contributions to data generation and analysis. For authors with equal contributions, ties were resolved alphabetically by first name.

## Abstract

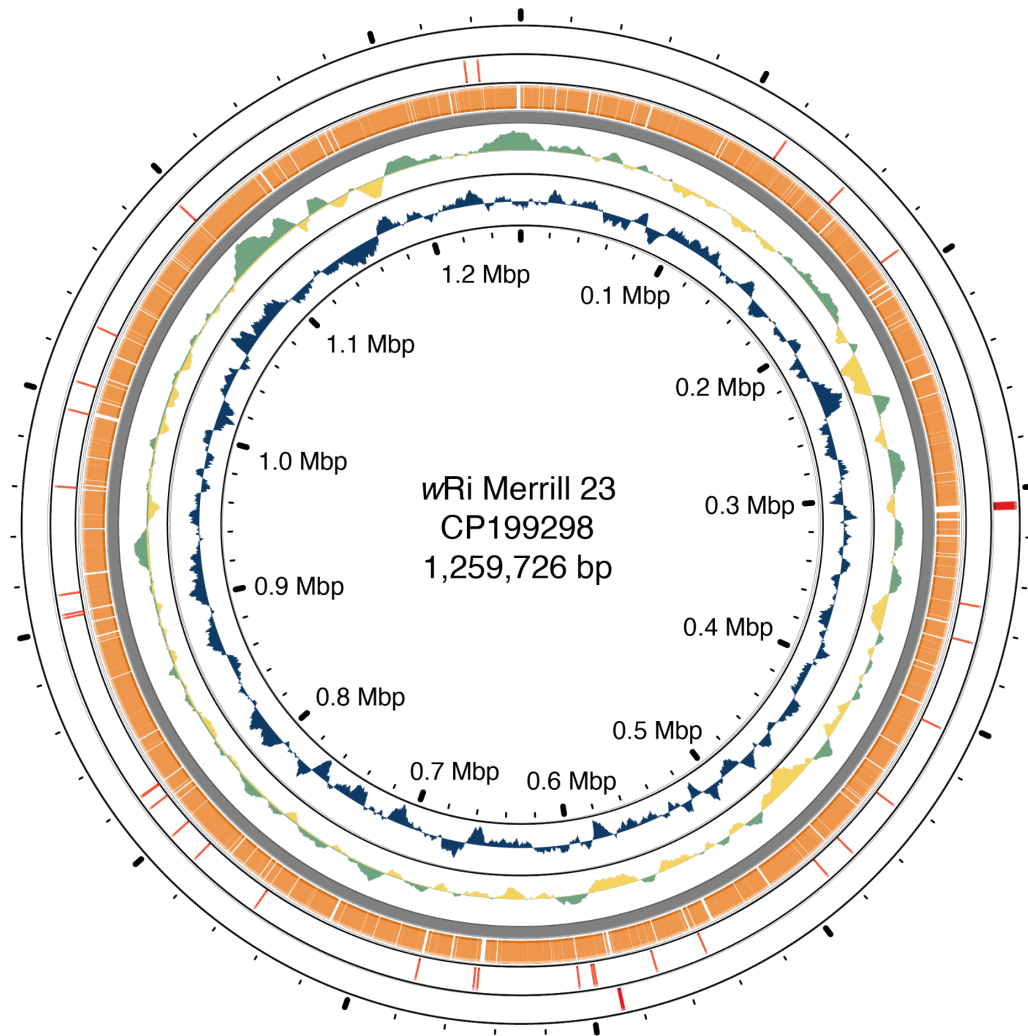
We present a contemporary high-quality, complete *de novo* assembly of *Wolbachia pipientis* (wRi Merrill 23, CP199298), an alphaproteobacterial endosymbiont of *Drosophila simulans*. This assembly was generated using long read sequencing of wRi-infected *D. simulans* embryos collected from the Merrill College at the University of California, Santa Cruz in October 2023.

*Wolbachia pipientis* infects diverse arthropods and nematodes, manipulating host phenotypes through cytoplasmic incompatibility (CI), male killing, and fertility rescue [1,2]. The Riverside strain (wRi) was first identified in California *Drosophila simulans* in the 1980s [3] and rapidly spread statewide due to exceptionally strong CI [4]. Despite its significance in shaping *D. simulans* populations [5], a modern wRi genome has not been assembled and the existing reference genome reflects the wRi present in 1984 [3,6]. Here, we present a complete *de novo* wRi genome assembly from contemporary *D. simulans* collected at the UC Santa Cruz Alan Chadwick Garden, located in Merrill College, in October 2023 (CP199298), providing an updated reference for future studies.

To generate a contemporary wRi genome assembly, we collected wild *D. simulans* flies, established isofemale lines, and performed long-read DNA sequencing of wRi-infected embryos. We established isofemale lines by deploying banana-baited bottles for ~5 days and collecting gravid females onto white food medium. After offspring eclosed, species identity was confirmed by phenotyping males and by PCR using silf-F/R primers to distinguish *D. simulans* from *D. melanogaster* [7] and wsp\_1F/592R primers to confirm wRi identity [8]. We extracted DNA from wRi-infected embryos using the Wizard HMW DNA Extraction Kit (Promega #A2920, Lot: 0000575812) and prepared libraries with the Native Barcoding Kit V14 (SQK-NBD114-24, Lot: NDP1424.10.0010). We sequenced these libraries on the Nanopore MinION Mk1B with a R10 version flow cell (FLO-MIN-114, Lot: 11004365) and MinKNOW v23.07.8 with adaptive sampling (fast model) to deplete *D. simulans* reads (GCF\_016746395.2), yielding 5.4M reads after 20 hours that were subsequently basecalled with Dorado (v0.7.3, hac model). After filtering for host-free reads >3kb, we assembled the wRi genome using Flye [9] following Jacobs and Nakamoto *et al.* (2024) [10], yielding a 1.26 Mb circular assembly with 30x coverage.

To polish the assembly, we generated Illumina short-read whole-genome sequencing data from whole wRi-infected *D. simulans* flies (Merrill 23 stocks). Illumina libraries were prepared using the Tn5 protocol [11] and sequenced on a NovaSeqX Plus. We polished the assembly with Pilon [12] v1.24 using short reads following Jacobs and Nakamoto *et al.* (2024). We assessed the quality of the polished assembly with BUSCO [13] (v5.7.0, rickettsiales\_odb10), which achieved a completeness score of 99.2%, annotated the assembly with Prokka [14] (v1.1.1, kingdom:bacteria) to identify coding sequences (CDS), tRNAs, rRNAs, and ncRNA (Table 1) and calculated and visualized GC content and GC skew with Proksee [15] v1.1.2 (Figure 1). Default parameters were used unless otherwise specified.

Raw sequencing reads and the assembled genome are available under BioProject accession number [PRJNA1312834](https://ncbi.nlm.nih.gov/bioproject/PRJNA1312834). Analysis scripts are available at [https://github.com/jodiejacobs/Jacobs\\_et\\_al\\_2026\\_de\\_novo\\_wRi\\_merrill\\_23\\_assembly](https://github.com/jodiejacobs/Jacobs_et_al_2026_de_novo_wRi_merrill_23_assembly).



**Figure 1. *Wolbachia* wRi genome map.** Concentric circles show (outer to inner): rRNA genes (red), tRNA genes (red), coding sequences (orange), GC skew (green/yellow for high/low), and GC content (blue), with GC metrics plotted as deviations from genome-wide average.

wRi Merrill Annotation summary	
Annotation pipeline	Prokka v1.1.1
Annotation method	kingdom:bacteria
Length (bp)	1,259,726
GC Content	35.22%

Genes (total)	1,283
CDSs (total)	1,246
Genes (RNA)	37
rRNAs	1, 1, 1 (5S, 16S, 23S)
tRNAs	34
ncRNAs	0
Pseudogenes (total)	3

Table 1. Annotation summary statistics.

### Acknowledgements:

The authors acknowledge the University of California Santa Cruz Genomics Institute for providing computational resources, including the Phoenix computational cluster, and support for this project. The authors thank Rion Parsons for his support and the University of California Santa Cruz for use of the Hummingbird computational cluster. The authors thank James Letchinger for the use of his computer for nanopore sequencing. The authors thank the University of California Santa Cruz Baskin Engineering Lab Support team for providing laboratory space and support. Funding for this project was provided by NIH (T32 HG012344) awarded to JJ, CW, HL, AN, CS, and AC, NIH awards to SLR (R00GM135583, R35GM157189), and the NSF-GRFP awarded to AN.

### References

1. Russell SL, Castillo JR. Trends in symbiont-induced host cellular differentiation. *Results Probl Cell Differ*. 2020;69: 137–176.
2. Russell SL, Castillo JR, Sullivan WT. Wolbachia endosymbionts manipulate the self-renewal and differentiation of germline stem cells to reinforce fertility of their fruit fly host. *PLoS Biol*. 2023;21: e3002335.
3. Hoffmann AA, Turelli M, Harshman LG. Factors affecting the distribution of cytoplasmic incompatibility in *Drosophila simulans*. *Genetics*. 1990;126: 933–948.
4. Turelli M, Hoffmann AA. Cytoplasmic incompatibility in *Drosophila simulans*: dynamics and parameter estimates from natural populations. *Genetics*. 1995;140: 1319–1338.
5. Carrington LB, Lipkowitz JR, Hoffmann AA, Turelli M. A re-examination of Wolbachia-

induced cytoplasmic incompatibility in California *Drosophila simulans*. PLoS One. 2011;6: e22565.

6. Klasson L, Westberg J, Sapountzis P, Näslund K, Lutnaes Y, Darby AC, et al. The mosaic genome structure of the *Wolbachia* wRi strain infecting *Drosophila simulans*. Proc Natl Acad Sci U S A. 2009;106: 5725–5730.
7. Faria VG, Sucena É. From nature to the lab: Establishing *Drosophila* resources for evolutionary genetics. Front Ecol Evol. 2017;5. doi:10.3389/fevo.2017.00061
8. Casper-Lindley C, Kimura S, Saxton DS, Essaw Y, Simpson I, Tan V, et al. Rapid fluorescence-based screening for *Wolbachia* endosymbionts in *Drosophila* germ line and somatic tissues. Appl Environ Microbiol. 2011;77: 4788–4794.
9. Kolmogorov M, Yuan J, Lin Y, Pevzner PA. Assembly of long, error-prone reads using repeat graphs. Nat Biotechnol. 2019;37: 540–546.
10. Jacobs J, Nakamoto A, Mastoras M, Loucks H, Mirchandani C, Karim L, et al. Complete de novo assembly of *Wolbachia* endosymbiont of *Drosophila willistoni* using long-read genome sequencing. Sci Rep. 2024;14: 17770.
11. Mirchandani C, Genetti M, Wang P, Pepper-Tunick E, Russell S, Corbett-Detig R. Plate Scale Tn5 based tagmentation library prep protocol v1. 2024. doi:10.17504/protocols.io.4r3l2qmpzpl1y/v1
12. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al. Pilon: an integrated tool for comprehensive microbial variant detection and genome assembly improvement. PLoS One. 2014;9: e112963.
13. Manni M, Berkeley MR, Seppey M, Simão FA, Zdobnov EM. BUSCO update: Novel and streamlined workflows along with broader and deeper phylogenetic coverage for scoring of eukaryotic, prokaryotic, and viral genomes. Mol Biol Evol. 2021;38: 4647–4654.
14. Seemann T. Prokka: rapid prokaryotic genome annotation. Bioinformatics. 2014;30: 2068–2069.
15. Grant JR, Enns E, Marinier E, Mandal A, Herman EK, Chen C-Y, et al. Proksee: in-depth characterization and visualization of bacterial genomes. Nucleic Acids Res. 2023;51: W484–W492.