

The Nature, Extent, and Consequences of Genetic Variation in the *opa* Repeats of *Notch* in *Drosophila*

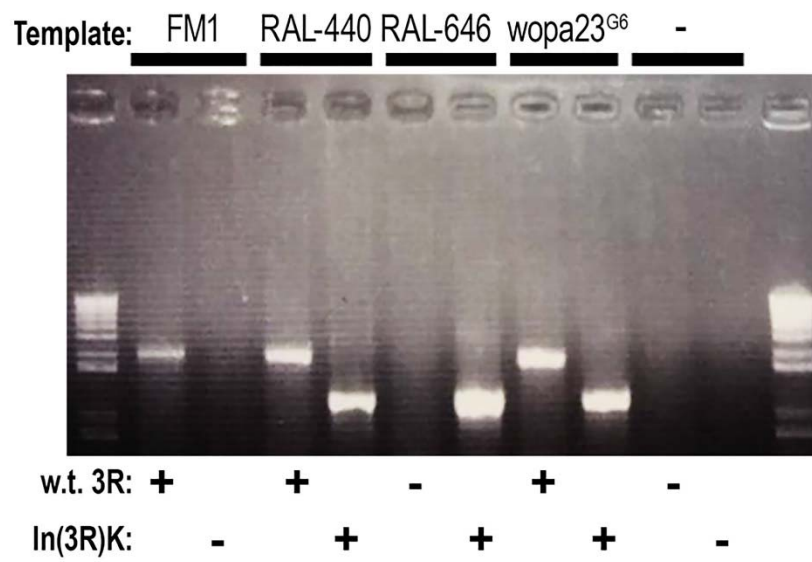
Clinton Rice^{*†}, Danielle Beekman^{*†}, Liping Liu^{*} and Albert Erives^{*1}

^{*}Department of Biology, University of Iowa, Iowa City, IA, 52242-1324, USA

[†]These authors contributed equally to this work.

¹Corresponding author: Department of Biology, University of Iowa, Iowa City, IA, 52242-1324, USA
E-mail: albert-erives@uiowa.edu

DOI: [10.1534/g3.115.021659](https://doi.org/10.1534/g3.115.021659)



File S1 Figure showing genotyping results for the Kodani inversion.

Supplementary File 2. Deconstruction of the original DGRP-646 (“Line646”) sequencing reads and alignment

1. We retrieved all original Illumina sequencing reads matching the exact CDS for the start (“PQ₉...”) or end (either “...Q₉L” or “...Q₅LGGLE”) of the Notch “opa” repeats underlying the “Line646” assembly (SRR835083). After confirming that the sequencing reads mapped to Notch CDS, we aligned and trimmed the sequences to focus on the opa repeat regions as shown below.
2. Aligning to reference (asterisked), we count a total of 67 mismatches from eight mismatched reads (10/13 reads). Furthermore, 3 of 4 codons spanning 5'-CAT (His) encode 5'-CAR (Gln). Number of mismatches (highlighted in fuchsia) relative to the reference sequence are given. Dashes indicate the region deleted in the inferred true alignment (explained in step #4 below).

REF.	*	<u>CCCCAGCAACAGCAGCAGCAGCAGCAACAGCAACAGCAGCAACATCAGCAGCAGCAACAGCAGCAGCAACAGCAGCAGCAGCAACAGCAGCAGCAACTC</u>
2766069	1	<u>CCCCAGCAACAGCAGCAGCAGCAGCAGCAGCA</u> -----
2766082	2	<u>CCCCAGCAACAGCAGCAGCAGCAGCAGCAGCAACAGCAGCAGC</u> -----
2766083	6	<u>CCCCAGCAACAGCAGCAGCAGCAGCAGCAGCAGCAACAGCAGCAGCAAAAGCAGCAGCAGCAACA</u>
2766075	25	----- <u>CCACAAAAGCACAAGAACAAGACACCACAACACCAC</u> <u>CAGCAACAGCAGCAGCAGCAACAGCAGCAGCAACTC</u>
2766076	9	----- <u>CAACACCCCCCGCAAAACAGCAGCAGCAACAGCAA</u> <u>CAGCAGCAACAGCAGCAGCAACTC</u>
2766089	2	----- <u>GCAGCAGCAGCAGCAACAGCAGCAGCAACAGCAGCAGCAGCAACAGCAGCAGCAACTC</u>
2766072	10	----- <u>GCAGCAGCAACGCCACCGGCACCACAACACCAGCAGCAGCAACAGCAGCAGCAACTC</u>
2766087	1	----- <u>GCAGCAGCAGCAACAGCAGCAGCAACAGCAGCAGCAGCAACAGCAGCAGCAACTC</u>
2766074	10	----- <u>CCCCACCAACAGCACCAACCAGCAGCAGCAGCAACAGCAGCAGCAACTC</u>
2766092	0	----- <u>GCAGCAACAGCAGCAGCAGCAACAGCAGCAGCAACTC</u>
2766078	1	----- <u>CCGCAGCAGCAGCAACAGCAGCAGCAACTC</u>
2766088	0	----- <u>GCAGCAGCAGCAACAGCAGCAGCAACTC</u>
2766090	0	----- <u>CAGCAGCAACAGCAGCAGCAACTC</u>

3. Removing reads with >2 mismatches in the opa repeat region reveals the origin of the DGRP “Line646” consensus of a single non-synonymous substitution (H→Q), two additional synonymous substitutions, and no indels. Three of the four mismatches are present in at least three sequences, suggesting that these are polymorphisms.

REF.	*	<u>CCCCAGCAACAGCAGCAGCAGCAGCAACAGCAACAGCAGCAACATCAGCAGCAGCAACAGCAGCAGCAACAGCAGCAGCAGCAACAGCAGCAGCAACTC</u>
2766069	1	<u>CCCCAGCAACAGCAGCAGCAGCAGCAGCAGCA</u> -----
2766082	2	<u>CCCCAGCAACAGCAGCAGCAGCAGCAGCAGCAACAGCAGCAGC</u> -----
2766089	2	----- <u>GCAGCAGCAGCAGCAACAGCAGCAGCAACAGCAGCAGCAGCAACAGCAGCAGCAACTC</u>
2766087	1	----- <u>GCAGCAGCAGCAACAGCAGCAGCAACAGCAGCAGCAGCAACAGCAGCAGCAACTC</u>
2766092	0	----- <u>GCAGCAACAGCAGCAGCAGCAACAGCAGCAGCAACTC</u>
2766078	1	----- <u>CCGCAGCAGCAGCAACAGCAGCAGCAACTC</u>
2766088	0	----- <u>GCAGCAGCAGCAACAGCAGCAGCAACTC</u>
2766090	0	----- <u>CAGCAGCAACAGCAGCAGCAACTC</u>
“Line646”	3	<u>CCCCAGCAACAGCAGCAGCAGCAGCAGCAACAGCAGCAGCAGCAGCAGCAACAGCAGCAGCAACAGCAGCAGCAACAGCAGCAGCAACTC</u> < DGRP 1.0/2.0

Rice et al. (2015)

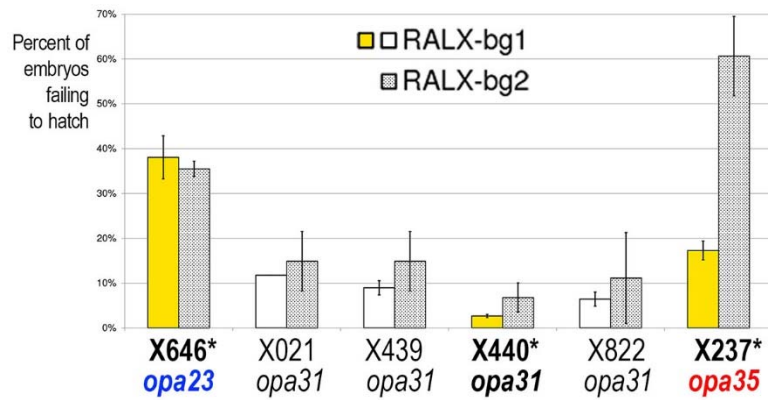
4. However, after aligning to the Sanger sequence for RAL-646, which we re-sequenced and found to contain a large gap, we find that mismatches are reduced to 56 differences from only 6 reads (previously was 67 differences from 10/13 reads). This represents a substantial improvement in alignment quality for a greater number of the original Illumina reads underlying DGRP assembly. With our new alignment, a total of seven reads have a perfect consensus.

REF.	1	<u>CCCCAGCAACAGCAGCAGCAGCAGCAACAGCAACAGCAGCAACATCAGCAGCAGCAACAGCAGCAGCAACAGCAGCAGCAACAGCAGCAGCAACTC</u>
Sanger-646 *		<u>CCCCAGCAACAGCAGCAGCAGCAGCAGCAGCAACAGCAGCAGCAACAGCAGCAGCAGCAACAGCAGCAGCAACTC</u>
2766069	0	<u>CCCCAGCAACAGCAGCAGCAGCAGCAGCAGCAACAGCAGCAGCAACAGCAGCAGCAGCAACAGCAGCAGCAACTC</u> -----CA
2766082	0	<u>CCCCAGCAACAGCAGCAGCAGCAGCAGCAGCAACAGCAGCAGCAACAGCAGCAGCAGCAACAGCAGCAGCAACTC</u> -----CAACAGCAGCAGC
2766083	1	<u>CCCCAGCAACAGCAGCAGCAGCAGCAGCAGCAACAGCAGCAGCAACAGCAGCAGCAACAGCAGCAGCAACTC</u> -----CAACAGCAGCAGCAAAAGCAGCAGCAGCAACA
2766075	24	<u>CCACAAAAGCACAAGAACAAGACCAGCAACAGCAGCAGCAACAGCAGCAGCAACAGCAGCAGCAACTC</u> -----CAACACCAGCAGCAACAGCAGCAGCAACAGCAGCAGCAACTC
2766076	10	<u>CAACACCCCCAGCAGCAGCAGCAGCAGCAGCAACAGCAGCAGCAACAGCAGCAGCAACAGCAGCAGCAACTC</u> -----AAACAGCAGCAGCAACAGCAAAGCAGCAGCAACAGCAGCAGCAACTC
2766089	0	<u>GCAGCAGCAGCAGCAGCAGCAGCAACAGCAGCAGCAACAGCAGCAGCAACAGCAGCAGCAACTC</u> -----CAACAGCAGCAGCAACAGCAGCAGCAACAGCAGCAGCAACTC
2766072	8	<u>GCAGCAGCAACGCAGCAGCAGCAGCAGCAGCAACAGCAGCAGCAACAGCAGCAGCAACAGCAGCAGCAACTC</u> -----CACCGGCAGCAGCAACAGCAGCAGCAACAGCAGCAGCAACTC
2766087	0	<u>GCAGCAGCAGCAGCAGCAGCAGCAACAGCAGCAGCAACAGCAGCAGCAACAGCAGCAGCAACTC</u> -----CAACAGCAGCAGCAACAGCAGCAGCAACAGCAGCAGCAACTC
2766074	11	<u>CCCCACCAGCAGCAGCAGCAGCAGCAGCAACAGCAGCAGCAACAGCAGCAGCAACTC</u> -----CAACAGCAGCAACAGCAGCAGCAACAGCAGCAGCAACTC
2766092	0	<u>GCAGCAACAGCAGCAGCAGCAGCAACAGCAGCAGCAACTC</u> -----GCAGCAACAGCAGCAGCAGCAACAGCAGCAGCAACTC
2766078	1	<u>CCGCAGCAGCAGCAACAGCAGCAGCAACTC</u> -----CCGCAGCAGCAGCAACAGCAGCAGCAACTC
2766088	0	<u>GCAGCAGCAGCAACAGCAGCAGCAACTC</u> -----GCAGCAGCAGCAACAGCAGCAGCAACTC
2766090	0	<u>CAGCAGCAACAGCAGCAGCAACTC</u> -----CAGCAGCAACAGCAGCAGCAACTC

5. Thus, using the seven reads without mismatches to each other, which represent a majority of reads (7/13), we can reproduce the exact Sanger RAL-646 genotype. We thus conclude that *opa23* reads were sequenced in the original RAL646 (DGRP 1.0) but were misassembled because large gaps are too prohibitive to search.

Sanger-646 *		<u>CCCCAGCAACAGCAGCAGCAGCAGCAGCAGCAACAGCAGCAGCAACAGCAGCAGCAGCAACAGCAGCAGCAACTC</u>
2766069	0	<u>CCCCAGCAACAGCAGCAGCAGCAGCAGCAGCAACAGCAGCAGCAACAGCAGCAGCAGCAACAGCAGCAGCAACTC</u>
2766082	0	<u>CCCCAGCAACAGCAGCAGCAGCAGCAGCAGCAACAGCAGCAGCAACAGCAGCAGCAGCAACAGCAGCAGCAACTC</u>
2766089	0	<u>GCAGCAGCAGCAGCAACAGCAGCAGCAACAGCAGCAGCAGCAACAGCAGCAGCAACTC</u> -----GCAGCAGCAGCAACAGCAGCAGCAACAGCAGCAGCAACTC
2766087	0	<u>GCAGCAGCAGCAACAGCAGCAGCAACAGCAGCAGCAACAGCAGCAGCAACTC</u> -----GCAGCAGCAGCAACAGCAGCAGCAACAGCAGCAGCAACTC
2766092	0	<u>GCAGCAACAGCAGCAGCAGCAGCAACAGCAGCAGCAACTC</u> -----GCAGCAACAGCAGCAGCAGCAACAGCAGCAGCAACTC
2766088	0	<u>GCAGCAGCAGCAACAGCAGCAGCAACTC</u> -----GCAGCAGCAGCAACAGCAGCAGCAACTC
2766090	0	<u>CAGCAGCAACAGCAGCAGCAACTC</u> -----CAGCAGCAACAGCAGCAGCAACTC

6. **EXTRA:** Illumina sequencing reads are 75 nucleotides long and a single read could have encompassed the 75 base pairs exactly underlying the "PQ₂₃L" coding sequence. We searched for Q₂₃ reads or the central Q₂₁, Q₁₉, Q₁₇, or Q₁₅ coding sequences to see if we could find reads that span the center poly-Q coding sequence without otherwise being anchored on the sides by non-polyQ encoding sequence. We did not find any such reads mapping to the *Notch* locus.



All 12 lines available by request from Erives Lab.

* These RALX-bg1 lines are also available from BDSC (one each of *opa23*, *opa31*, and *opa35* in bg1).

File S3 Figure showing the separate embryonic assay results for the RALX-bg1 and RALX-bg2 series.