

SUPPORTING INFORMATION

Gene Model Annotations for *Drosophila melanogaster*: The Rule-Benders

Madeline A. Crosby^{*1}, L. Sian Gramates^{*}, Gilberto dos Santos^{*}, Beverley B. Matthews^{*}, Susan E. St. Pierre^{*}, Pinglei Zhou^{*}, Andrew J. Schroeder^{*}, Kathleen Falls^{*}, David B. Emmert^{*}, Susan M. Russo^{*}, William M. Gelbart^{*}, and the FlyBase Consortium

^{*}Department of Molecular and Cellular Biology, Harvard University, Cambridge, MA 02138

¹Corresponding author: FlyBase, Biological Laboratories, Harvard University, 16 Divinity Avenue, Cambridge, MA 02138, 617.495.9925, crosby@morgan.harvard.edu

DOI: 10.1534/g3.115.018937

Figures and tables in Supporting Information (this file)

Figure S1. Polycistronic locus with monocistronic, dicistronic and tricistronic alternative transcripts.

Figure S2. Conservation of protein sequence beyond stop-codon readthroughs.

Table S1: Standardized comments used by FlyBase for flagging exceptional transcripts

Table S2: GenBank flags used in transcript and protein RefSeq entries

Table S3: Genes annotated with a non-AUG translation start in release 6.04

Supporting Information submitted as separate files:

File S1 Complete listing of polycistronic loci.

File S2 Complete listing of genes that share exons with other genes.

File S3 Complete listing of genes with multiphasic exons.

File S4 Listing of all annotated introns with non-canonical splices.

File S5 Complete listing of genes annotated with a stop-codon readthrough.

Available for download as Excel files at www.g3journal.org/lookup/suppl/doi:10.1534/g3.115.018937/-/DC1

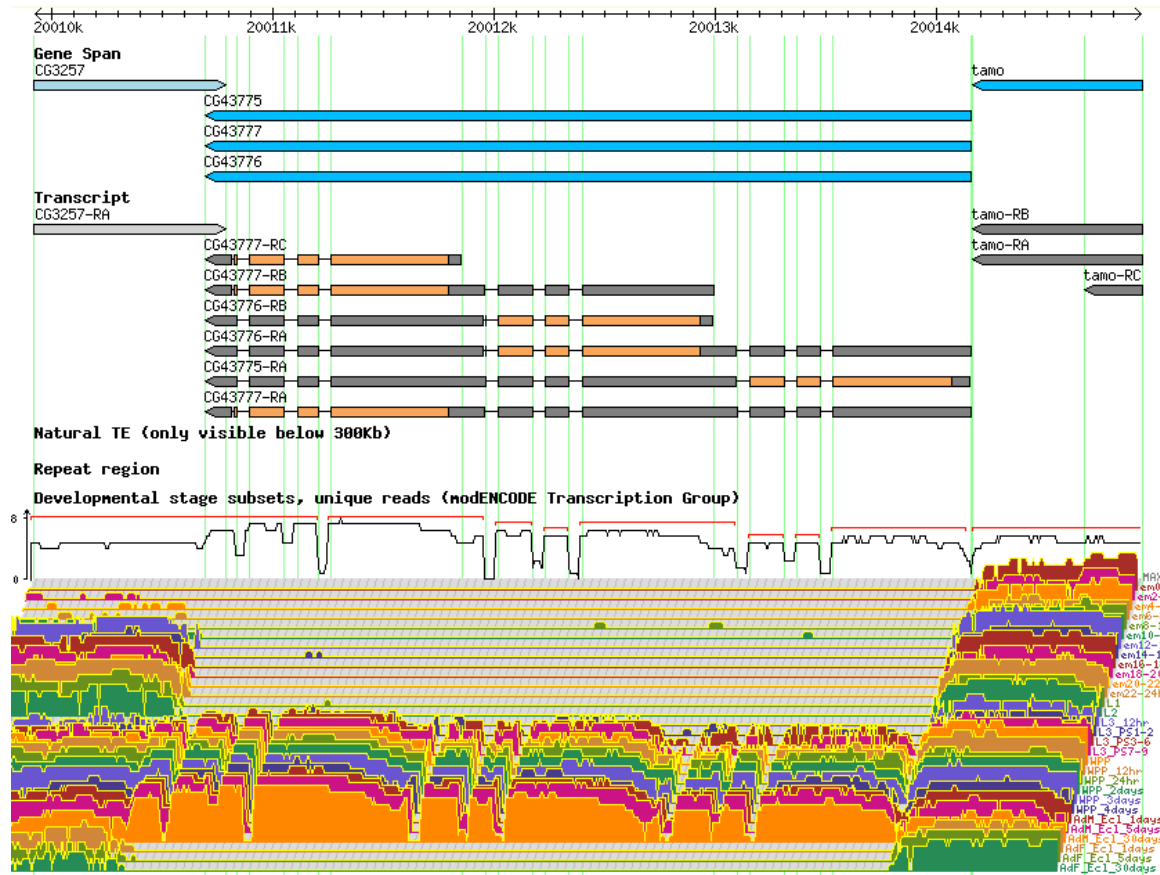


Figure S1 Polycistronic locus with monocistronic, dicistronic and tricistronic alternative transcripts. All annotated introns are supported by cDNA and RNA-Seq junction data; downstream transcription start sites are supported by RAMPAGE TSS data. A GBrowse view showing (top to bottom) the gene extents and the gene models; unstranded RNA-Seq coverage data corresponding to a developmental series (early embryos, top, to adults, bottom). More information on data presented in GBrowse may be found at http://flybase.org/wiki/FlyBase:GBrowse_Tracks#General.

CLUSTAL 2.1 multiple sequence alignment

```

Dmel_trn-PB      IRETIKGLWGNLSALGRKEREYQKTFCEDEYMSRQHHPHPCSLGIHSTFPNTYAPHH--P 58
Dpse_trn        IREYLKGLWGSALGRKEREYQKTFCEDEYMARHQHHPHPCSLGIHSTFPNTYAPHH--P 58
Dvir_trn        IREYLKGGVWGNLSALGRKEREYQKTFCEDEYMSRLQHHPHPCSLGIHSTFPNTYAPHQATA 60
Dgri_trn        IREYLKGLWGNLSALGRKEREYQKTFCEDEYMSRLQHHPHPCSLGIHSTFPNTYAPHQTTA 60
Dmel_caps-PD    IREMLKG---HSALGRKEREYQKTFSEDEYMSRPP-PGGGG-VHPAAGG-----YP 46
Dpse_caps       IREMLKG---HSALGRKEREYQKTFSEDEYMRTPP-PGCGG-VHPAAA-----YP 45
Dgri_caps       IRELKFG---HSALGRKEREYQKTFSEDEYMRTPPAPGCVGGVHPGSG-----YP 47
Dvir_caps       IREMLKG---HSALGRKEREYQKTFSEDEYMRTPP-PGCGG-VHPASG-----YP 45
* * : * *      * * * * * * * * * * * * * * * * * * * * * * * * * * * * * *
Dmel_trn-PB      GAT---HHYGMCP-MPVNDLGAII-DPQKFKQLVVP-----TATMIS-EKK 98
Dpse_trn        SAPG---HHYGMCTGMPINDLNAAGDPQKFKQLVVP-----TGLSMN-EKK 102
Dvir_trn        SAP---HHHYGMCP-MPINDLNAV--DGQHKFKQLVVP-----ITATLMH-EKK 102
Dgri_trn        AAAAAHHHYGMCP-MPINDLNSG--DGQHKLQLQVPA-----MSATLLHTEQK 107
Dmel_caps-PD    YIAG-----NSRMIPVTELXLEAPPPQLRGRGG-----GGGASTAS--GA 85
Dpse_caps       CSQSQ---YMGSRPIPVTLELXLEAPPPQMRGRGGIASTTTTTTSSGSGSSGSAAP 101
Dgri_caps       CPSSYNTTYQLGSRPIPVTLELXLEVPPPPQLRGRGG-----APTSH-----GAS 91
Dvir_caps       CPSGINSSQYLGSRSIPVTELXLEAPPPQLRGRGG-----APSSHNGPMASSTSGTP 98
. : : : *      . : : : *      . : : : *      . : : : *
Dmel_trn-PB      LNNNKALVSQGAIDDSASFVLMKATMGRDQVH-----QNPQ-----136
Dpse_trn        LNNNKALASQAAIDDSASFVLMKATLARDHLQ-----QHPHPHQHPQHPQ--QHQS 155
Dvir_trn        LNNNKSLA--GSVDSASFVLMKATMGRERQQQ-----AQLQQQHQLQQQQL--QQQHP 153
Dgri_trn        LNNNKALA--GSIDDSANFVLMKATMNRDQHQ-----QQQQQLQQQQQLPHQQHP 160
Dmel_caps-PD    VQQLQVPSAVDQAS--NSFAQLSHIHYMTNNGQ-----QQAQQQSTSKMHHSQ 133
Dpse_caps       LQQLQVPSVDHAAAASSFAQLSHIHYMTN-----QTIASPMTPNQMHHSQ 148
Dgri_caps       LQQLQVPSAIDA--HAPFAQLSHIHYMTNPLTASSAASAAAASTTATTPRSHHSQ 149
Dvir_caps       LQQLQ-VPSAVDAS-HAPFAQLSHIHYMTNPS-----AATTTATTPRSHHSQ 146
. : : : *      . : : : *      . : : : *      . : : : *
Dmel_trn-PB      ---LNHYTKPQFLSATAVGDSCYS---YADVPMVHGAPLGGP--NQQLRLTQEHFK 186
Dpse_trn        VQSKLNHYTKPQFLSATAAVGDSCYS---YADVPMVHAAPLGAAPLQQLRLTHEHFK 211
Dvir_trn        HQSKLNHYTKPQFLAATAVADSCYS---YADVPMVHAP-----QQLRVTHEHFK 200
Dgri_trn        QQSKLNHYTKPQFLAATAVADSCYS---YADVPLVHAP-----QQLRITHEHFK 207
Dmel_caps-PD    QDMRLLACNGGKTLNA-TSLPRHRP---PVVQESTLSHYSQP---LANGIRLTQDHFN 184
Dpse_caps       QDMRLLANGGGKTLNA-ASLPRHR---MQESTLSHYSOPLA---LANGIRLTQDHFN 198
Dgri_caps       QDMRLLANG--KALGLNASLPRHMAGRQCGVQESTLSHYSQVPG---IRLTQDHFN 200
Dvir_caps       QDMRLLANGGGKVLGLNASLPRHMAGRQCGVQESTLSHYSOPLAN--GHAGIRLTQDHFN 204
* * : * *      * * : * *      . : : *      . : : *      . : : * * *
Dmel_trn-PB      QR-----ELYDQEMGS-EILDHNYIYSNTHYSMPLE 216
Dpse_trn        HRVAGT-----GEHYDNEVNS-EILDPNYIYSNAHYSMPLE 246
Dvir_trn        QREQR-----PRDFEADINGEMDPNYIYSNAHYSMPLE 236
Dgri_trn        QRE-----RDFDENPLGEMDPNYIYSNAHYSMPLE 239
Dmel_caps-PD    HNQQ-----SHNQHYGG-VYAKPCDAMSEPGYIHNNSHYSPLD 223
Dpse_caps       HNSHGHHGHHGHHGHHGPHSLGHPHSHTTHGGGVYAKACDAMTEPGYIHNNSHYSPLD 258
Dgri_caps       HN-----GGVYAKPCDAMTEPGYIHNNSHYSPLD 230
Dvir_caps       HNGIGG-----GAGVYAKPCDAMAEAGYIHNNSHYSPLD 239
. : : * * * * *      . : : * * * * *
Dmel_trn-PB      QLGRSKTPTPPPMPALPLRNLG-----CATTGRRSFQKSAQKQQQNNNTLRQFTHX 270
Dpse_trn        QMGRSKTPTPPPMPALPLRNLG-----CATTGRRSFQKSA-----NNNTLRQFTHX 294
Dvir_trn        QMGRNKTPTPPVPALPLRNLG-----CATTGRRSFQK-TPAHNN-NTSTMRQFTHX 288
Dgri_trn        QMGRSKTPTPPPLPALPLRNLG-----CATTGRRSLQHRPTATN---TLRQFTHX 288
Dmel_caps-PD    HDLPP-SPTPTPPPMPALPLRNGVMALIHGNTTGRSFSN-----NNNVSTLSNNNH 275
Dpse_caps       HDLPP-SPTPTPPPMPALPLRNGVMALIHGNTTGRSFSNSNSNS-NNNVATLSNNNH 315
Dgri_caps       HDMPPTPTPPPMPALPLRNLGMLVHGNTTGRSFSNSNSNSNNNNVATLSNNNH 289
Dvir_caps       HDMPPTPTPPPMPALPLRNGMMLVHGNTTGRSFSN-----NNNVATLSNNNH 292
. : : * * * * *      . : : * * * * *
Dmel_trn-PB      -----SSTYRRRQLSIYA-- 283
Dpse_trn        -----STVDYRRRQLSIYA-- 308
Dvir_trn        -----SSEHYRRRQLSIYA-- 302
Dgri_trn        -----SSEHYRRRQLSIFA-- 302
Dmel_caps-PD    ---GIGGGV-VAVGGTVGNNGSLRRYH-- 300
Dpse_caps       ---GGGGGGGGGAILASSNNGSLRRYH-- 341
Dgri_caps       ---GLVLANG--CSG-INNNNGSLRRYH-- 311
Dvir_caps       GVGGGAVLVNGNYNTNGSNNNNGSLRRYH-- 322
. : : *      . : : *

```

Figure S2 Conservation of protein sequence beyond stop-codon readthroughs. Clustal alignment (<http://www.ebi.ac.uk/Tools/msa/clustalw2/>) of the carboxy end of *tartan* gene proteins from several *Drosophila* species (if annotated with a short read-through extension) and the predicted carboxy end of *caps* gene proteins (if annotated with a long read-through extension) from several *Drosophila* species. The *caps* readthrough (denoted by 'X') is at alignment position 82; the *trn* readthrough is at alignment position 300. Note regions of low complexity, variable conservation and variable length interspersed with regions of protein sequence conservation.

Table S1 Standardized comments used by FlyBase for flagging exceptional transcripts

FlyBase transcript-associated standardized comments
Polycistronic transcript.
Dicistronic transcript.
Dicistronic transcript; alternative monocistronic transcript(s) exist.
Monocistronic transcript; alternative dicistronic transcript(s) exist.
Unconventional splice site postulated (<i>[splice donor-acceptor specified]</i>).
Unconventional splice site invoked (<i>[splice donor-acceptor specified]</i>); sequence altered due to transposon insertion; this splice may not occur in vivo.
Unconventional splice site(s) invoked due to gap in genomic sequence; this splice does not occur in vivo.
Unconventional splice site(s) invoked (<i>[splice donor-acceptor]</i>); within a dynamic region of nested TEs that may differ in different strains; this splice may not occur in vivo.
Unconventional splice site postulated: transcript subject to HAC1-type intron splice site recognition and cleavage.
Trans-spliced.
Unconventional translation start postulated (<i>[codon specified]</i> encoding Met).
Stop-codon suppression (UGA as Sec) postulated; reflected in aa sequence of predicted polypeptide.
Stop-codon suppression postulated (<i>[codon specified]</i>); reflected in aa sequence of predicted polypeptide.
Double stop-codon suppression postulated (<i>[codons specified]</i>); reflected in aa sequence of predicted polypeptide.
Translational frameshifting postulated: +1 frameshift reflected in aa sequence of predicted polypeptide.
TAA stop codon is completed by the addition of 3' A residues to the mRNA.
Start codon not determined.

Table S2 GenBank flags used in transcript and protein RefSeq entries

GenBank flags for exceptional cases
gene /exception="dicistronic gene" [transcript RefSeq entries only]
CDS /exception="nonconsensus splice site"
CDS /trans_splicing
CDS /note="non-AUG ([<i>codon specified</i>]) translation initiation"
CDS /transl_except=(pos:x..y,aa:Met)
CDS /transl_except=(pos:x..y,aa:Sec)
CDS /transl_except=(pos:x..y,aa:OTHER)
CDS /ribosomal_slippage
CDS /transl_table=5
CDS /transl_except=(pos:x,aa:TERM)
CDS /note="TAA stop codon is completed by the addition of 3' A"
CDS /note="start codon not determined"

Table S3 Genes annotated with a non-AUG translation start in release 6.04

Gene	Start codon	Reference (FlyBase Reference ID)
Eip74EF	CUG	Burtis <i>et al</i> 1990 (FBrf0051390), Boyd and Thummel 1993 (FBrf0064374)
cpo	CUG	Bellen <i>et al</i> 1992 (FBrf0056119)
ewg	CUG	de Simone and White 1993 (FBrf0059052)
Eip78C	CUG	Stone and Thummel 1993 (FBrf0064719)
Syn	CUG	Klagges <i>et al</i> 1996 (FBrf0087510)
att-ORFA	CUG	Madigan <i>et al</i> 1996 (FBrf0089733)
Fmr1	CUG	Beerman and Jongens 2011 (FBrf0213401)
Trpy	CUG	FlyBase analysis
CG4629	CUG	FlyBase analysis
CG11076	CUG	FlyBase analysis
CG16890	CUG	FlyBase analysis
Cha	GUG	Sugihara <i>et al.</i> 1990 (FBrf0052176)
Akt1	GUG	FlyBase analysis
Klp54D	ACG	Andjelkovic <i>et al.</i> 1995 (FBrf0079853)
NAT1	AUU	Takahashi <i>et al</i> 2005 (FBrf0184018)
Gsc	AUU	FlyBase analysis
CG11836	AUU	FlyBase analysis
Wnk	AUU	FlyBase analysis
Sh	AUU	FlyBase analysis
CG14989	UUG	FlyBase analysis
Jwa	UUG	FlyBase analysis
sol	UUG	FlyBase analysis
CG43778	UUG	FlyBase analysis
CG2162	AUC	FlyBase analysis
CG43921	AUC	FlyBase analysis
CG30334	AUC	FlyBase analysis
CG43273	AUC	FlyBase analysis