

Overview:

I was given a 50k base long sequence from the *D. erecta* and told to annotate all the proteins in the sequence. First I ran my fosmid through RepeatMasker to eliminate all repeats and then I ran the masked sequence through Genscan, Nscan, and comparing it to the homologous region on *D. melanogaster* that has already been annotated. Through these processes I was able to find all the genes and pseudogenes in my fosmid and these are shown below in Overview-Figure 1. I found 3 functional genes and 2 pseudogenes. Through the entire process I also discovered which tools were helpful and which were not. Nscan predicted 2 functional genes with accurate exon/intron borders while Genscan predicted 6 functional genes of which only 2 were functional and real genes. I found that once I had the homologous *D. melanogaster* region, blasting the mRNA of the genes in the homologous region against my fosmid proved to be the most useful technique in finding functional genes and finding accurate exon/intron borders. Sixpack was also helpful in finding the exon/intron borders because it showed all 6 reading frames for my fosmid.



Overview-Figure 1: Genes 1,2, and 3 are 3 genes homologous to 3 genes in *D. melanogaster*. Pseudogene 1 and 2 are 2 pseudogenes homologous to 2 real genes in *D. melanogaster*.

Genes:

Gene 1 is homologous to CG7133(NP_649379.1) in *D. melanogaster*. The gene only has 1 isoform. There is only 1 exon and its sequence is on the minus strand running from 1274-288bp on the fosmid. My starting point for this gene was with Genscan's predicted gene 1. Results for Genscan's gene 1 are in the figure labeled "Gene 1-Figure 1" below.

<u>Gn</u>	<u>Ex</u>	<u>Type</u>	<u>S</u>	<u>.Begin</u>	<u>...End</u>	<u>.Len</u>	<u>Fr</u>	<u>Ph</u>	<u>I/Ac</u>	<u>Do/T</u>	<u>CodRg</u>	<u>P....</u>	<u>Tscr..</u>
1.02	PlyA	-		188	183	6							1.05
1.01	Sngl	-		1274	285	990	2	0	42	43	586	0.985	46.37
1.00	Prom	-		2814	2775	40							-7.76

```
>Dere3_dna|GENSCAN_predicted_peptide_1|329_aa
MSNVYKDHYHVLGLARNASDSEIREAFRRSLQYHPDKNENGAGEFLKINDAYRVLIDHH
KRASYDRRLSFRDLEAIIIPSENASGQLSELRIKTS PGNFHKKLKVAVVIGGVLVGTYYA
YRVFQKSPPIIPVPQPITPPAIP TQELSELHPGYLWTLISGLVTLRSKRILSLGKLATGA
NIRVSPSTL KAPLSSAAEVVAKTVIQGGVGVSTATSSSSLATPANVVAKMAKTLFKGSR
SGLYTASKTVVVPSTEILHWSKTA AKCTLATLKKGTSYARSPVPSGSVLLSSI PSALRAF
AGTFRPALVKSKNIVVKKASASWIQKRL
```

Gene 1-Figure 1: Results for Genscan's gene 1

I extracted bp 1-3000 from my fosmid using extractseq and blasted it against *D. melanogaster* in Flybase. I only got 1 hit that had a reasonable e-score. The results are in the figure labeled "Gene 1-Figure 2" below.

>gnl|dmel|3L type=chromosome_arm; loc=3L:1..24543557; ID=3L; dbxref=GB:AE014296; MD5=ec7148cae3daabbd2a226eaa6e85d7c2; length=24543557; release=r5.17; species=Dmel; Length = 24543557

HSP # = 1 , Score = 270.093 bits (136) , Expect = 2.05781e-70
 Identities = 315 / 378 (83.3%) , Positives = 315 / 378 (83.3%) , Gaps = 15 / 378 (4%)
 Strand = Plus / Plus

Genome View		Subject FASTA	
Query:	593	CCACGTTTCGCTGGCGTTGCTAAGGACGAACTGGAAGTTGCAGTTGAACTACGCCCTCCTG	652
Subject:	22066752		22066811
Query:	653	GCCCTTGGATTACTGTTTTGGCTACGACTTCGCGAGCCGAAAGATAAGGGTCTTTCAGAG	712
Subject:	22066812		22066871
Query:	713	TACTAGGAGAACTCGAATATTTGCTCCTGTAGCCAGTTTTCCGAGACTCAGTATTCTTT	772
Subject:	22066872		22066931
Query:	773	TCGATCTCAATGTCACTAGCCCGGAGATCAATGTCCATAGGTATCCGGGATGCAGTTCGC	832
Subject:	22066932		22066991
Query:	833	TTAGTTCCTGCGTTGGAATAGCTGGAGGAGTGATTGGCTGGGAACTGGGATTATTGGCG	892
Subject:	22066992		22067036
Query:	893	GAGATTTCTGGAATACCCGGTACGCCACATAGGTGCCAACCCAGCACGCCTCCGATTACGA	952
Subject:	22067037		22067096
Query:	953	CAGCCACTTTGAGCTTTT 970	
Subject:			22067114

Gene 1-Figure 2: Results from nblast of bp 1-3000 from my fosmid and *D. melanogaster* genome.

The region shown in Gene 1-Figure 2 showed a homology between bp 593-970 of my extracted sequence and a region in the gene CG7133 of *D. melanogaster*. I then did a pblast of Genscan's predicted peptide sequence for gene 1 against *D. melanogaster*. Once again I only got 1 hit that had a reasonable e-score. The results are shown below in the figure labeled "Gene 1-Figure 3."

>gnl|dmel|FBpp0078138 type=protein; loc=3L:complement(22066375..22067436); ID=FBpp0078138; name=CG7133-PA; parent=FBgn0037150, FBtr0078485; dbxref=FlyBase:FBpp0078138, FlyBase_Annotation_IDs:CG7133-PA, GB_protein:AAF51805.2, REFSEQ:NP_649379, GB_protein:AAF51805; MD5=6cdf10d2532767019459ec3c5a423d60; length=353; release=r5.17; species=Dmel; Length = 353

HSP # = 1 , Score = 280.796 bits (717) , Expect = 7.83597e-76
 Identities = 161 / 298 (54%) , Positives = 191 / 298 (64.1%) , Gaps = 31 / 298 (10.4%)

Subject FASTA	
Query:	1 MSNVYKDHVHVLGLARNASDSEIREAFRRLSLQYHPDKNENGAGEFLKINDAYRVLIDHH 60
Subject:	1 MS+VY+DHY VLGL RNA+DSEI++AFRRLSLQYHPDKNE+GA EFL+IN+A+RVLIDH 60
Query:	61 KRASYDRRLSFRDLEAIIIPSENASGQSEL-----RIIKTSPGNFHKKLVAVVIGGVL 114
Subject:	61 +RA YD D+EAIIP+ENA+GQL EL +T P +F +KLKVA IGG+L 120
Query:	115 VGTYYAVRVFQKXXXXXXXXXXXXXXXXXTOELSELHPGYLWTLISGLVTLRSKRILSLG 174
Subject:	121 VGTYYVYRVFQK-----PPSIPVFRPIPTOELSDSHLGSWTLSSGLLALRSKRILGLG 175
Query:	175 KLATGANIRVSPSTLKAPLSSAAEVVAKTVIQGPGVGXXXXXXXXXXXXXXXXNVVAK---- 230
Subject:	176 KLA AN RVSP TL AP SAAA+VVAKTVI+G VG NV K 235
Query:	231 -----MAKTLFKGSRSLGYTASKTIVVVPST----EILHWSKTA--AKCTLATL 272
Subject:	236 +AKTL +GSR+G Y+A KTV + +L+W+ T K T AT+ 293

Gene 1-Figure 3: Results of a pblast of Genscan's predicted peptide sequence for gene 1 against *D. melanogaster*.

The predicted peptide sequence from Genscan was homologous to CG7133 from *D. melanogaster*. Since both the nblast and pblast for Genscan's predicted gene 1 matched up with CG7133 I expected to find a homologous gene to CG7133 in my fosmid. I decided to run a nblast 2 of the mRNA sequence from CG7133 against my fosmid to see where in my fosmid the homologous sequence was. The results are below in the figure labeled "Gene 1-Figure 4."

```

Score = 520 bits (281), Expect = 2e-149
Identities = 580/715 (81%), Gaps = 57/715 (7%)
Strand=Plus/Minus

Query 184   AAAATGAGCGATGTCTACGAAGATCACTACCAGGTTCTGGGCTTACCGAGAAATGCCACC 243
          ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
Sbjct 1277   AAAATGAGCAACGTCTACAAAGACCCTACCATGTTCTGGGCTTGGCGAGAAACGCCAGC 121

Query 244   GACAGTGAGATTAAG-GATGCTTTTCGGCGGCTGTCCCTGCAATATCATCCCGACAAAA 302
          ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
Sbjct 1217   GACAGTGAGA-TCAGAGAAGCTTTTCGGCGGTTGTCCCTGCAATATCATCCCGACAAAA 115

Query 303   CGAGGATGGAGCGAAGGAGTTCCTTAGAATCAACGAGGCCCATCGCGTCTGATTGACCA 362
          ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
Sbjct 1158   CGAGAATGGAGCGGGGAGTTCCTTAAAATCAACGACGCCTACCGCGTGCTAATTGACCA 109

Query 363   TCAGAGAAGGGCCTT-GTACGATTG-CTGC-T-TCCAGTCCATGGACGTT-GAAGCCATT 417
          ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
Sbjct 1098   TCATAAAAGGGC-TTCGTACGATCGTC-GCCTGTC--GTTTAGGGAC-TTAGAAGCCATC 104

Query 418   ATTCCCG-CTGAGAACGCTAATGGCCAACCTGCCTGAATTGGGAAATCCATTCTTCCCAAT 476
          ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
Sbjct 1043   ATTCC-GTCCGAGAACGCTAGTGGCCAACCTGTCTGAATTACGAA-TC-AT-CA----AA- 993

Query 477   GCCACCCGAAACGCCGCTGCTAGTTTTTCGCGAAAAGCTCAAAGTTGCTGCCTTCATCGG 536
          ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
Sbjct 992   ---AC---AT-CGCC---TGGTAATTTCCACAAAAGCTCAAAGTGGCTGTCTGTAATCGG 943

Query 537   AGGCTTGCTGGTTGGTACATACGTGGGGTACCGGGTATTCAGAAACCGCCCACTCTA- 595
          ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
Sbjct 942   AGGCGTGTCTGGTTGGCACCTATGTGGCGTACCGGGTATTCAGAAATCTCCGCCAA-TAA 884

Query 596   TCCAGTTCCTCCCGCCCAAT---TCC---A---A---C----GCAGGAACCTAAGCGAT-TCG 639
          ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
Sbjct 883   TCCAGTTCCTCCAGCAATCACTCCTCCAGCTATTCCAACGCAGGAACCTAAGCGAAT-G 825

Query 640   CATCTCGGATCCCTATGGACATTG-TCCTCCGGCTATTGGCATTGAGATCGAAAAGAAT 698
          ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
Sbjct 824   CATCCCGGATACCTATGGACATTGATC-TCCGGGCTAGTGACATTGAGATCGAAAAGAAT 766

Query 699   ACTGGGACTCGGCAAACTGGCGCCATGG-GCAAATACTCGAGTTTCTCCTAGGACTCTGA 757
          ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
Sbjct 765   ACTGAGTCTCGGAAAACCTGGCTACA-GGAGCAAATATTGAGTTTCTCCTAGTACTCTGA 707

Query 758   ACGCGCCCTTTTCTTCGGCTGCCAAAGTCGTGGCCAAAACAGTAATCCGAGG-CCAAAGA 816
          ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
Sbjct 706   AAGCACCTTATCTTCGGCTGCCGAAGTCGTAGCCAAAACAGTAATCCAAGGGCCAG-GA 648

Query 817   GCCGTAGTTCCTTCTGCAACTTCCAGTTCGTTCCTTAGCATCGGCTGCGAACGTGG 871
          ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
Sbjct 647   GGCGTAGTTCAACTGCAACTTCCAGTTCGTTCCTTAGCAACGCCAGCGAACGTGG 593

```

Gene 1-Figure 4: Results from a nblast 2 of the mRNA sequence from CG7133 against my fosmid

I ran my entire masked fosmid through sixpack to get all the reading frames and then used the results from Gene 1-Figure 4 to see where the first start codon was and where the first stop codon was on the same reading frame. That is where I got the exon borders at 1274-288bp on my fosmid. I then ran gene 1 with its coordinates for the exon through Gene checker which verified the gene. I ran a pblast 2 of the peptide sequence from Gene checker for gene 1 against the peptide sequence for CG7133. The results are shown below in the figure labeled "Gene 1-Figure 5".

```
>lcl|42903 FBpp0078138 type=protein; loc=3L:complement(22066375..22067436);
ID=FBpp0078138; name=CG7133-PA; parent=FBgn0037150,FBtr0078485;
dbxref=FlyBase:FBpp0078138,FlyBase Annotation IDs:CG7133-PA,GB_protein:AAF51805.2,REFSEQ:NP_649379,GB_protein:
MD5=6cdf10d2532767019459ec3c5a423d60; length=353; release=r5.17;
species=Dmel;
Length=353
```

```
Score = 283 bits (725), Expect = 3e-81, Method: Compositional matrix adjust.
Identities = 186/338 (55%), Positives = 223/338 (65%), Gaps = 39/338 (11%)
```

```
Query 1 MSNVYKDHYHVLGLARNASDSEIREAFRRSLQYHPDKNENGAGEFLKINDAYRVLIDHH 60
MS+VY+DHY VLGL RNA+DSEI++AFRRSLQYHPDKNE+GA EFL+IN+A+RVLIDH
Sbjct 1 MSDVYEDHYQVLGLPRNATDSEIKDAFRRSLQYHPDKNEDGAKEFLRINEAHRVLIDHQ 60

Query 61 KRASYDRRLSFRDLEAIIPSENASGQLSEL-----RIIKTSPGNFHKKLVAVVIGGVL 114
+RA YD D+EAIIP+ENA+QQL EL +T P +F +KLVVA IGG+L
Sbjct 61 RRALYDCCCFQSM DVEAIIPAENANGQLPELGNPFFPMPPETPPASFREKLVAAPIGGLL 120

Query 115 VGTYYVAYRVFQKSPPIIPVQPITPPAIPTQELSELHPGYLWTLISGLVTLRSKRILSLG 174
VGTYYV YRVFQK P P IPTQELS+ H G LWTL SGL+ LRSKRIL LG
Sbjct 121 VGTYYVGYRVFQKP-----PPSIPVPRPIPTQELSDSHLGSWTLSSGLLALRSKRILGLG 175

Query 175 KLATGANIRVSPSTLKAPLSSAAEVVAKTVIQGPGVGSTATSSSSSLATPANVVAK---- 230
KLA AN RVSP TL AP SSAA+VVAKTVI+G VGS+ATSSSSLA+ ANV K
Sbjct 176 KLAPWANTRVSPRTLNAFFSSAAKVVAKTVIRGQRAVGSSATSSSSSLASAAANVAVKSLPS 235

Query 231 -----MAKTLFKGSRSGLYTASKTVVVPST---EILHWSKT-----AAKCTLAT 271
+AKTL +GSR+G Y+A KTV + +L+W+ T A T+A
Sbjct 236 KASVNSATETVAKTLQSGSRAGPYSALKTVWSSAVSYLRSLLNWATTPKWKGATPATIAG 295

Query 272 LKKG-----TSYARSPVPSGVLSSIPSALRAFAGTF 304
L +S A++ + + S I S LRA AG F
Sbjct 296 LVHNSRPVWSSAAKNTLAALMYACSKISSYLRLALAGRF 333
```

Gene 1-Figure 5: Results from a blast 2 of the peptide sequence from Gene checker for gene 1 against the peptide sequence for CG7133.

The protein sequence is not incredibly well conserved as shown by the 55% for identities, 65% for positives, and that the peptide sequence for gene 1 from the fosmid is 329 a.a. and the peptide sequence for CG7133 is 353 a.a. I wanted to find out why there was 24 a.a. difference in peptide length so I ran a clustal2w of the coding sequence for CG7133 and my gene 1. The results are shown below in the figure labeled “Gene 1-Figure 6”.

FBtr0078485 ATGAGCGATGTCTACGAAGATCACTACCAGTTCTGGGCTTACCGAGAAATGCCACCGAC 60
 EMBOSS_001 ATGAGCAACGTCTACAAAGACCACTACCATGTTCTGGGCTTGGCGAGAAACGCCAGCGAC 60
 ***** * ***** ** * ***** * ***** * ***** * ***** *

FBtr0078485 AGTGAGATTAAGGATGCTTTTCGGCGGCTGTCCCTGCAATATCATCCCGACAAAAACGAG 120
 EMBOSS_001 AGTGAGATCAGAGAAGCTTTTCGGCGGTTGTCCCTGCAATATCATCCCGACAAAAACGAG 120
 ***** * ** * ***** * ***** * ***** * ***** *

FBtr0078485 GATGGAGCGAAGGAGTTCCTTAGAATCAACGAGGCCATCGCGTCTGATTGACCATCAG 180
 EMBOSS_001 AATGGAGCGGGGAGTTCCTTAAAATCAACGACGCCTACCGCGTCTAATTGACCATCAT 180
 ***** * ***** * ***** * ***** * ***** * ***** *

FBtr0078485 AGAAGGGCCTTGTACGATTGCTGCTTCCAGTCCATGGACGTTGAAGCCATTATCCCGCT 240
 EMBOSS_001 AAAAGGGCTTCGTACGATCGTCGCCTGTCGTTTAGGGACTTAGAAGCCATCATTCCGTCC 240
 * ***** * ***** * ** * ** * ***** * ***** * ***** *

FBtr0078485 GAGAACGCTAATGGCCAACCTGCCTGAATTGGGAAATCCATTCTTCCAATGCCACCCGAA 300
 EMBOSS_001 GAGAACGCTAGTGGCCAACCTGCTGAATTACGAA--TCAT-----CAAA 282
 ***** * ***** * ***** * ** * ** * ***** * ***** *

FBtr0078485 ACGCCGCTGCTAGTTTTTCGCGAAAAGCTCAAAGTTGCTGCCTTCATCGGAGGCTTGCTG 360
 EMBOSS_001 ACATCGCTGGTAATTTCCACAAAAAGCTCAAAGTGGCTGTCGTAATCGGAGGCGTGCTG 342
 ** ***** * ** * * ***** * ***** * ***** * *****

FBtr0078485 GTTGGTACATACGTGGGGTACCGGGTATTCCAGAAACCGCCCATCTATCCAGTTCCC 420
 EMBOSS_001 GTTGGCACCTATGTGGCGTACCGGGTATTCCAGAAATCTCCGCCAATAATCCAGTTCCC 402
 ***** * ** * ** * ***** * ***** * ***** * *****

FBtr0078485 CGGCCAAT-----TCCAACGCAGGAATAAGCGATTGCGATCTCGGATCC 465
 EMBOSS_001 CAGCCAATCACTCCTCCAGCTATTCGAACGCAGGAATAAGCGAACTGCGATCCCGGATAC 462
 * ***** * ***** * ***** * ***** * ***** *

FBtr0078485 CTATGGACATTGCTCCTCCGGGCTATTGGCATTGAGATCGAAAAGAATACTGGGACTCGGC 525
 EMBOSS_001 CTATGGACATTGATCTCCGGGCTAGTGACATTGAGATCGAAAAGAATACTGAGTCTCGGA 522
 ***** * ***** * ***** * ***** * ***** * *****

FBtr0078485 AAACGGCGCCATGGGCAAATACTCGAGTTTCTCCTAGGACTCTGAACGGCCCTTTTCT 585
 EMBOSS_001 AAACGGCTACAGGAGCAAATATTGAGTTTCTCCTAGTACTCTGAAAGCACCCCTTATCT 582
 ***** * ** * ***** * ***** * ***** * ***** *

FBtr0078485 TCGGCTGCCAAAGTCGTGGCCAAAACAGTAATCCGAGGCCAAAGAGCCGTAGTTCTTCT 645
 EMBOSS_001 TCGGCTGCCGAAGTCGTAGCCAAAACAGTAATCCAAGGCCAGGAGGCGTAGTTCAACT 642
 ***** * ***** * ***** * ***** * ***** * *****


```

FBtr0078485      GCAACTTCCAGTTCGTCCTTAGCATCGGCTGCGAACGTGGCTGTGAAATCCCTACCAAGC 705
EMBOSS_001      GCAACTTCCAGTTCGTCCTTAGCAACGCCAGCGAACGTGGTTGCCAAA----- 690
*****
***** ** * ***** ** ***

FBtr0078485      AAAGCTTCAGTGAATTCGCTACGGAAACAGTCGCTAAAAACACTTCCAGGGATCACGA 765
EMBOSS_001      -----ATGGCGAAAAACACTTTTCAAGGGATCCCGA 720
*****
* ** ***** * ***** **

FBtr0078485      GCTGGACCATATTCAGCTTTGAAAACAGTT-TGGTCCTCAGCGTTTCATACTGCGCTC 824
EMBOSS_001      TCTGGACTATACAGCTTCCAAAACAGTTGTGGTACCGAG---TACCGAGATTTTACAT 777
***** ** * ***** ***** ** * ** * * *

FBtr0078485      TCTTCTAAATTGGGCAACTACACCGAAATGGGGGAAGGCAACACCAGCAACCATTGCTGG 884
EMBOSS_001      TGGTCCAAAACGG-----CAGCAAAATGTACCCTTGTACCTTGAAGAAGGGTAC-AT 829
* ** ** * ** * ** * ** * ** * ** * **

FBtr0078485      TTTGGTTCATAACTCCAGACCAGTTTGGTCATCTGCAGCAAAAAATACCCTTGTGCCTT 944
EMBOSS_001      CCTATGCTAGGAGTCCAGTCCCTTCGGGTTCTGTGCTGCTTTCAAGCATACCATCGGCTT 889
* * * ***** ** * ** * ** * ** * ** * **

FBtr0078485      GATGTATGCATGCTCCAAAATACTTCGTATTTGCGCGCTTTGGCCGGGCGCTTTATTAG 1004
EMBOSS_001      TACGCGCTTTTCCCGGAACCTTTTCGTCCAGCTTTAGTAAAAATCCAAGAACATTGTTGTGA 949
* * ** * ** * ** * ** * ** * * * * * *

FBtr0078485      TCCACCTTTAGACGCCTCCAAAGACTACAATCTTTTAAAGGAAGCTTGAAAAATTGA 1062
EMBOSS_001      AGAAAGCCAGCGCCAGTTGGATACAAAAAAACGCCTATAG----- 990
* * * * * * * * * * * *

```

Gene 1-Figure 6: Results from a clustalw2 of the coding sequence for CG7133 and my gene 1. The highlighted hyphens are regions in CG7133 that are not present in my gene. The highlighted “A” of CG7133 and the highlighted “TAG” of my gene show a nonsense mutation.

The results of the clustalw2 show that there are several regions of base pairs in the sequence of CG7133 that are not present in my gene 1. Also, there is a nonsense mutation where there was a point mutation in my gene from an “A” to a “T” resulting in a premature stop codon in my gene 1. Since my gene 1 protein was not very conserved with the CG7133 protein I thought that maybe it may be a pseudogene in my fosmid. First I looked at the function of the CG7133 in *D. melanogaster* and found that it functions in unfolded protein binding and heat shock protein binding. This function seems important so *D. erecta* probably has this functioning gene. I ran a nblast 2 of the mRNA from CG7133 against the entire *D. erecta* genome to find out if maybe the ortholog of CG7133 was found somewhere else in *D. erecta* but I did not get any good hits. Based on the importance of the gene and the fact that in flies pseudogenes are rare, I concluded that this gene is probably a real gene and that though there is not very high conservation between the proteins, the protein coded for by my gene 1 is still functional.

Gene 2 is homologous to CG7130(NP_649380.1) in *D. melanogaster*. The gene only has 1 isoform. There is only 1 exon and its sequence is on the minus strand running from 3823-3443bp on the fosmid. My starting point for this gene was with Nscan’s predicted gene 1. Results for Nscan’s gene 1 are in the figure labeled “Gene 2-Figure 1” below.

<u>Exon</u>	Strand	Begin	End	Length
-------------	--------	-------	-----	--------

1	-	3439	3868	429
---	---	------	------	-----

127 aa

MGKDYYKILGIERNASSEEVKKGYYRRMALRYHPDKNDHPQAEHFREVVA
AFEVLSKDKEKRETYDKYGEEGLRCDDEPATFAQPTSDMLPFMCAVGGTVL
FAFAAYKTFQFFNRKKEATDGDGSSSD

Transcript:

ATGGGTAAGGATTACTACAAGATTCTGGGCATCGAGAGAAATGCGTCCAG
CGAAGAAGTGAAGAAAGGATACCGCCGGATGGCTCTCCGCTACCATCCAG
ACAAGAACGACCATCCGCAGGCTGAGGAGCACTTCAGGGAGGTGGTGGCC
GCCTTCGAAGTGCTCTCCGACAAGGAAAAGCGCGAGACATACGACAAGTA
CGGCGAGGAGGGCCTCAGGTGTGATGACGAGCCGGCGACCTTCGCCCAGC
CCACGTCAGACATGCTCCCCTTTATGTGCGCCGTCGGAGGAACTGTGCTC
TTTGCATTCGCCGCCTATAAGACCTTCCAGTTTTTCAACCGGAAAAGGA
GGCTACCGACGGCGATGGATCGTCCTCGGAC

Gene 2-Figure 1: Results for Nscan's gene 1

I extracted bp 3000-4000 from my fosmid using extractseq and blasted it against *D. melanogaster* in Flybase. I only got 1 hit that had a reasonable e-score. The results are in the figure labeled "Gene 2-Figure 2" below.

>gnl|dmel|3L type=chromosome_arm; loc=3L:1..24543557; ID=3L; dbxref=GB:AE014296; MD5=ec7148cae3daabbd2a226eaa6e85d7c2; length=24543557; release=r5.17; species=Dmel; Length = 24543557

HSP # = 1 , Score = 511.94 bits (258) , Expect = 3.94308e-144
 Identities = 353 / 384 (91.9%) , Positives = 353 / 384 (91.9%) , Gaps = 3 / 384 (0.8%)
 Strand = Plus / Minus

Genome View		Subject FASTA	
Query: 1	ATGGGTAAGGATTACTACAAGATTCTGGGCATCGAGAGAAATGCGTCCAGCGAAGAAGTG	60	
Subject: 22068734		ATGGGTAAGGATTACTACAAGATTCTGGGCATCGAGAGAAATGCGTCCAGCGAAGAAGTG	22068675
Query: 61	AAGAAAGGATACCGCCGGATGGCTCTCCGCTACCATCCAGACAAGAACGACCATCCGCAG	120	
Subject: 22068674		AAGAAAGGATACCGCCGGATGGCTCTCCGCTACCATCCAGACAAGAACGACCATCCGCAG	22068615
Query: 121	GCTGAGGAGCACTTCAGGGAGGTGGTGGCCGCTTCGAAGTGTCTCCGACAAGGAAAAG	180	
Subject: 22068614		GCTGAGGAGCACTTCAGGGAGGTGGTGGCCGCTTCGAAGTGTCTCCGACAAGGAAAAG	22068555
Query: 181	CGCGAGACATACGACAAGTACGCGAGGAGGGCCTCAGGTGTGATGACGAGC---CGGCG	237	
Subject: 22068554		CGCGAGATATACGACCAGCAGCGAGGAGGGTCTCAAATGTGATGACGAGCCTGCTGCG	22068495
Query: 238	ACCTTCGCCACGCCACGTCAGACATGCTCCCTTTATGTGCGCCGTCGGAGGAACTGTG	297	
Subject: 22068494		ACCTTCGCCACGCCACGTCAGACATGCTCCCTTCAATGTGCGCCGTCGGAGGAACTGTG	22068435
Query: 298	CTCTTTGCATTGCGCCCTATAAGACCTTCCAGTTTTC AACCGGAAAAAGGAGGCTACC	357	
Subject: 22068434		CTCTTTGCATTGCGCCCTATAAGACATTCAGTTTTC AACCGGAAAAAGGAGGCTACC	22068375
Query: 358	GACGGCGATGGATCGTCCTCGGAC	381	
Subject: 22068351		CACGGCGATGGATCGTCCTCGGAC	22068351

Gene 2-Figure 2: Results from nblast of bp 3000-4000 from my fosmid and *D. melanogaster* genome.

The region shown in Gene 2-Figure 2 showed a homology between bp 1-381 of my extracted sequence and a region in the gene CG7130 of *D. melanogaster*. I then did a pblast of Nscan's predicted peptide sequence for gene 1 against *D. melanogaster*. Once again I only got 1 hit that had a reasonable e-score. The results are shown below in the figure labeled "Gene 2-Figure 3."

>gnl|dmel|FBpp0078137 type=protein; loc=3L:complement(22068348..22068734); ID=FBpp0078137; name=CG7130-parent=FBgn0037151, FBtr0078484; dbxref=FlyBase:FBpp0078137, FlyBase_Annotation_IDs:CG7130-PA, GB_protein:AAF51806.1, REFSEQ:NP_649380, GB_protein:AAF51806; MD5=fcb3086a827f52f13ec0072806552ab; length=128

HSP # = 1 , Score = 242.662 bits (618) , Expect = 4.15323e-65
 Identities = 118 / 128 (92.2%) , Positives = 122 / 128 (95.3%) , Gaps = 1 / 128 (0.8%)

Subject FASTA	
Query: 1	MGKDYYKILGIERNASSEEVKKGYYRRMALRYHPDKNDHPQAEHFREVVAAFEVLSDEK 60
Subject: 1	MGKDYYKILGIERNASSE+VKKGYRRMALRYHPDKNDHPQAE FREVVAAFEVL DKEK 60
Query: 61	RETYDKYGEEGLRCDDEP-ATFAQPTSDMLPFMCAVGGTVLFAFAAYKTFQFFNRKKEAT 119
Subject: 61	REIYDQHGEGLKCDDEPAATFAQPTDMLPFMCAVGGTVLFAFAAYKTFQFFNRKKEAT 120
Query: 120	DGDGSSSD 127
Subject: 121	HGDGSSSD 128

Gene 2-Figure 3: Results of a pblast of Nscan's gene 1 predicted peptide sequence for against *D. melanogaster*.

The predicted peptide sequence from Nscan was homologous to CG7130 from *D. melanogaster*. Since both the nblast and pblast for Nscan's predicted gene 1 matched up with CG7130 I expected to find a homologous gene to CG7130 in my fosmid. I decided to run a nblast 2 of the mRNA sequence from CG7130 against my fosmid to see where in my fosmid the homologous sequence was. The results are below in the figure labeled "Gene 2-Figure 4."

```
>lcl|59535 Dere3_dna range=fosmid17:1-50000 5'pad=0 3'pad=0 strand=+ repeatMasking=none
Length=50000
```

```
Score = 756 bits (838), Expect = 0.0
Identities = 534/607 (87%), Gaps = 9/607 (1%)
Strand=Plus/Minus

Query 94 CAGAACAAGGAGAATACCATTTCCACCACAATGGGTAAGGATTACTACAAGATTCTGGGCA 153
|
Sbjct 3852 CAGAACAAGGAGAATCCCATTTCCACCACAATGGGTAAGGATTACTACAAGATTCTGGGCA 3793

Query 154 TCGAGAGGAATGCGTCCAGCGAAGACGTCAGAAGGGATACCGCCGGATGGCTCTCCGCT 213
|
Sbjct 3792 TCGAGAGAAATGCGTCCAGCGAAGAAGTGAAGAAAGGATACCGCCGGATGGCTCTCCGCT 3733

Query 214 ACCATCCGACACAAGAACGACCATCCGCAGGCCGAGGAGCAGTTTAGGGAGGTGGTGGCCG 273
|
Sbjct 3732 ACCATCCAGACAAGAACGACCATCCGCAGGCTGAGGAGCACTTCAGGGAGGTGGTGGCCG 3673

Query 274 CCTTCGAAGTGCTCTTTGATAAGGAAAAGCGCGAGATATACGACCAGCACGGCGAGGAGG 333
|
Sbjct 3672 CCTTCGAAGTGCTCTCCGACAAGGAAAAGCGCGAGACATACGACAAGTACGGCGAGGAGG 3613

Query 334 GTCTCAAATGTGATGACGAGCCTGCTGCGACCTTCGCCAGCCCCAGCCAGACATGCTCC 393
|
Sbjct 3612 GCCTCAGGTGTGATGACGAGC---CGGCGACCTTCGCCAGCCCCAGTCAGACATGCTCC 3556

Query 394 CCTTCATGTGCGCCGTCGGAGGAACCGTGCTCTTTGCGTTCCGCCCTACAAGACATTCC 453
|
Sbjct 3555 CCTTTATGTGCGCCGTCGGAGGAACGTGCTCTTTGCATTCCGCCCTATAAGACCTTCC 3496

Query 454 AGTTCTTCAACCGGAAAAAAGAGGCTACCCACGGCGATGGATCCTCCTCGGACTGAGCTA 513
|
Sbjct 3495 AGTTTTTCAACCGGAAAAAAGAGGCTACCGACGGCGATGGATCGTCTCCTCGGACTGAGCTA 3436

Query 514 AGGATCCAAGGGCTTGTGAAGCAATTCGGGTACCTAGCGTTCTTCGCTGAATAGTCTT 573
|
Sbjct 3435 ACGATCGGAGGGCTTGGTGAAGCAATACCGGGGATCTAGCGTCTTCACTGAATAGTCTT 3376

Query 574 TAAGATTAAATTTATAGGAACTTAATTATTGACTGTTTATCTAATGAATCCTGCGTACTT 633
|
Sbjct 3375 TAAGATTAAATTTATAGGAACTTGTACATTGACTGTTGATCTCATTATTATGTG-TAGTT 3317

Query 634 ATTGATTAAATTTATTTATTTATTTAGTAAGATAAAATAAAAATTATGGAATTCGTCGGACT 693
|
Sbjct 3316 ACTGAT--ATTGCTGGTTTTATTTAGTAA---AAACAAAAGTTGTGGAATGCGGCAACT 3262

Query 694 GTTGCTC 700
|
Sbjct 3261 CTTGCTC 3255
```

Gene 2-Figure 4: Results from a nblast 2 of the mRNA sequence from CG7130 against my fosmid

I ran my entire masked fosmid through sixpack to get all the reading frames and then used the results from Gene 2-Figure 4 to see where the first start codon was and where the first stop codon was on the same reading frame. That is where I got the exon borders at 3823-3443bp on my fosmid. I then ran gene 2 with its coordinates for the exon through Gene checker which verified the gene. I ran a pblast 2 of the peptide sequence from Gene checker for gene 2 against the peptide sequence for CG7130. The results are shown below in the figure labeled "Gene 2-Figure 5".

```
>lcl|41915 FBpp0078137 type=protein; loc=3L:complement(22068348..22068734);
ID=FBpp0078137; name=CG7130-PA; parent=FBgn0037151,FBtr0078484;
dbxref=FlyBase:FBpp0078137,FlyBase Annotation IDs:CG7130-PA,GB protein:AAF51806.1,REFSEQ:N
MD5=fcb3086a827f52f13ec00728065552ab; length=128; release=r5.17;
species=Dmel;
Length=128
```

```
Score = 242 bits (618), Expect = 1e-69, Method: Compositional matrix adjust.
Identities = 118/128 (92%), Positives = 122/128 (95%), Gaps = 1/128 (0%)
```

```
Query 1 MGKDYYKILGIERNASSEEVKKGYYRRMALRYHPDKNDHPQAEHFREVVAAFEVLSKKEK 60
Sbjct 1 MGKDYYKILGIERNASSE+VKKGYRRMALRYHPDKNDHPQAE FREVVAAFEVL DKEK 60

Query 61 RETYDKYGEEGLRCDDEP-ATFAQPTSDMLPFMCAVGGTVLFAFAAYKTFQFFNRKKEAT 119
Sbjct 61 REIYDQHGEGLKCDDEPAATFAQPTDMLPFMCAVGGTVLFAFAAYKTFQFFNRKKEAT 120

Query 120 DGDGSSSD 127
Sbjct 121 HGDGSSSD 128
```

Gene 2-Figure 5: Results from a pblast 2 of the peptide sequence from Gene checker for gene 2 against the peptide sequence for CG7130.

The protein sequence is pretty well conserved as shown by the 92% for identities, 95% for positives, and that the peptide sequence for gene 1 from the fosmid is 127 a.a. and the peptide sequence for CG7130 is 128 a.a. I can conclude from this that gene 2 is a real functional gene whose ortholog is CG7130 from *D. melanogaster*.

Gene 3 is homologous to RpLP0 (NP_524211.1) in *D. melanogaster*. The gene only has 1 isoform. There are 2 exons whose sequences are on the plus strand running from 4450-4503 and 4579-5475bp on the fosmid. My starting point for this gene was with Nscan's predicted gene 2. Results for Nscan's gene 2 are in the figure labeled "Gene 3-Figure 1" below.

Exon Strand Begin End Length

```
1 + 4341 4503 162
2 + 4578 5478 900
```

317 aa

```
MVRENKAAWKAQYFIKVVLELDFEFPKCFIVGADNVGSKQMQNIIRTSLRGL
AVVLMGKNTMMRKAIRGHLENNPQLEKLLPHIKGNVGFVFTKGDLAEVRD
KLLESKVRAPARPGAIAPLHVIIPAQNTGLGPEKTSFFQALSIPTKISKG
TIEIINDVPILKPGDKVGASEATLLNMLNISPFSYGLIVSQVYDSGSIFS
PEILDIKPEDLRAKFQQGVANLAAVCLSVGYPTIASAPHSIANGFKNLLA
IAATTEVEFKEATTIKEYIKDPSKFAAAASVSAAAPAAGGAAEKKEEAKKV
ESESEEDDDMGFGLFD
```

Transcript:

ATGGTTAGGGAGAACAAGGCAGCATGGAAGGCTCAGTACTTCATCAAGGT
TGTGGAAGTGTTCGATGAGTTCCTCAAGTGCTTCATCGTGGGCGCCGACA
ACGTTGGCTCCAAGCAGATGCAGAACATCCGTACCAGCCTGCGTGGACTG
GCCGTGCTGCTTATGGGCAAGAACACCATGATGCGCAAGGCCATCCGCGG
TCATCTGGAGAACAACCCGCAGCTGGAGAAGCTGCTGCCCCACATCAAGG
GTAACGTGGGCTTCGTTTTACCAAGGGCGATCTCGCCGAGGTGCGTGAC
AAGCTGTTGGAGTCCAAGGTGCGCGCCCCCGCCGTCCCGGCGCTATTGC
CCCTCTGCACGTCATCATCCCGGCCCAGAACACCGGCTTGGGACCCGAGA
AGACCAGTTTCTTCCAGGCCCTGTCCATCCCGACCAAGATTTCCAAGGGA
ACAATTGAAATCATCAACGATGTGCCCATCCTGAAGCCCGGCGACAAGGT
CGGCGCCTCCGAGGCAACGCTGCTCAACATGTTGAACATCTCGCCCTTCT
CGTACGGTTTGATCGTCAGCCAGGTGTACGACTCCGGCTCGATCTTTTCG
CCTGAGATTCTGGACATTAAGCCCGAGGATCTGCGCGCCAAGTTCCAGCA
GGGAGTGGCCAACCTGGCCGCCGTTTGTGTTGTCTGTGGGCTACCCACCA
TTGCCTCGGCCCCGCACAGCATTGCCAACGGATTCAAGAACCTGCTGGCC
ATTGCTGCCACCACCGAGGTGGAGTTCAAGGAGGCGACCACCATCAAGGA
GTACATCAAGGACCCAGCAAGTTCGCCGCCGCTGCCTCGGTTTCGGCTG
CCCCCGCCGCCGGCGGAGCTGCCGAGAAGAAGGAGGAGGCCAAGAAAGTC
GAGTCCGAGTCCGAGGAGGAGGACGATGATATGGGCTTCGGTCTGTTCGA
C

Gene 3-Figure 1: Results for Nscan's gene 2

I extracted bp 4000-6000 from my fosmid using extractseq and blasted it against *D. melanogaster* in Flybase. I only got 1 hit that had a reasonable e-score. The results are in the figure labeled "Gene 3-Figure 2" below.

>gnl|dmel|3L type=chromosome_arm;loc=3L:1..24543557;ID=3L;dbxref=GB:AE014296;
MD5=ec7148cae3daabbd2a226eaa6e85d7c2;length=24543557;release=r5.17;species=Dmel;
Length = 24543557

HSP # = 1 , Score = 2030.42 bits (1024) , Expect = 0
Identities = 1348 / 1459 (92.4%) , Positives = 1348 / 1459 (92.4%) , Gaps = 19 / 1459 (1.3%)
Strand = Plus / Plus

Genome View

Subject FASTA

Query:	215	TCGCTATCGATGTGGTCACACTTGCTTCCGGCGCCAACTTCCCTCTTTCCGTTCTGTGAG	274
Subject:		TCGCCATCGAAGCGGTACACTGGGTGCCGCCCAACTTCACTCTTTCCGTTCTGTGAG	22069211
22069152			
Query:	275	CGAAAACCGAAAAAGTCTGTGCTTTGGTAAGTATTATTAAGAGCGAAAAAGATGTTGCAT	334
Subject:		CGAAAACCGAAAAAGTCTGTGCTTTGGTAAGTGTGCTAAAAAGTTCGGAATAATGTTGCAT	22069271
22069212			
Query:	335	CCCGAGCTTTTTTGGGTGAATAACTGTTGCATGGCGCTGGCCAGTACCGACTAATCGAG	394
Subject:		CCCGAGCATTTTCGGGTACATAACTGTTCCACGGCGGTGGTCCAGCAAAGACTAATCGTT	22069331
22069272			
Query:	395	ATCACATCTTCGCGAGTTCTTAAATTCACCCGACGAGTCCCTAATACAAAATCAAATGG	454
Subject:		ATCACGCCTTTTCGCGAGTTCTTAAATTCACCCGACGAGTCCCTAATACACAATTAATGG	22069391
22069332			
Query:	455	TTAGGGAGAAACAAGGCAGCATGGAAGGCTCAGTACTTCATCAAGGTTGTGGTAAGTATAG	514
Subject:		TTAGGGAGAAACAAGGCAGCGTGAAGGCTCAGTACTTCATCAAGGTTGTGGTAAGTATAG	22069451
22069392			
Query:	515	AACCG-----CCCTCACTAGCTCGCCCTGGCTTATGCTCTTAACTAATCCTCGCT	565
Subject:		AACCTTATAGAATTCGCTCACTAGCTGGCGCTGGCTTATGCTGTTAACTGATCC-----	22069506
22069452			

Query: 566	AATTCCTCCTCCAGGAACTGTTTCGATGAGTTCCCAAGTGCTTCATCGTGGGCGCCGACA	625
Subject: 22069507	-----CTCCTCCAGGAACTGTTTCGATGAGTTCCCAAGTGCTTCATCGTGGGCGCCGACA	22069561
Query: 626	ACGTTGGCTCCAAGCAGATGCAGAACATCCGTACCAGCCTGCGTGGACTGGCCGTCGTGC	685
Subject: 22069562	ACGTGGGCTCCAAGCAGATGCAGAACATCCGTACCAGCCTGCGTGGACTGGCCGTCGTGC	22069621
Query: 686	TTATGGGCAAGAACACCATGATGCGCAAGGCCATCCGCGGTCATCTGGAGAACAACCCGC	745
Subject: 22069622	TTATGGGCAAGAACACCATGATGCGCAAGGCCATCCGCGGTCATCTGGAGAACAACCCGC	22069681
Query: 746	AGCTGGAGAAGCTGCTGCCCCACATCAAGGGTAACGTGGGCTTCGTTTTACCAAGGGCG	805
Subject: 22069682	AGCTGGAGAAGCTGCTACCCACATCAAGGGCAACGTGGGATTCGTGTTACCAAGGGCG	22069741
Query: 806	ATCTCGCCGAGGTGCGTGACAAGCTGTTGGAGTCCAAGGTGCGCGCCCCGCCGTC	865
Subject: 22069742	ATCTCGCCGAGGTGCGCGACAAGCTGCTGGAGTCCAAGGTGCGCGCCCCGCCGTC	22069801
Query: 866	GCGCTATTGCCCTCTGCACGTCATCATCCCGGCCAGAACACCGGCTTGGACCCGAGA	925
Subject: 22069802	GCGCTATTGCCCTCTGCACGTCATCATCCCGGCCAGAACACCGGCTTGGACCCGAGA	22069861
Query: 926	AGACCAGTTTCTTCCAGGCCCTGTCCATCCCGACCAAGATTTCCAAGGGAACAATTGAAA	985
Subject: 22069862	AGACCAGTTTCTTCCAGGCCCTGTCCATCCCGACCAAAAATTTCCAAGGGAACAATTGAAA	22069921
Query: 986	TCATCAACGATGTGCCATCCTGAAGCCCGGCGACAAGGTCGGCGCCTCCGAGGCAACGC	1045
Subject: 22069922	TCATCAACGATGTGCCATCCTGAAGCCTGGCGACAAGGTCGGCGCCTCCGAGGCGACAC	22069981
Query: 1046	TGCTCAACATGTTGAACATCTCGCCCTTCTCGTACGGTTTGATCGTCAGCCAGGTGTACG	1105
Subject: 22069982	TGCTCAACATGTTGAACATCTCGCCCTTCTCGTACGGTCTGATTGTCAACCAGGTCTACG	22070041

Score = 1700 bits (920), Expect = 0.0
 Identities = 1023/1074 (95%), Gaps = 2/1074 (0%)
 Strand=Plus/Plus

```

Query 189  GGAACTGTTTCGATGAGTTCCCAAAGTGCTTCATCGTGGGCGCCGACAACGTGGGCTCCAA 248
          |||
Sbjct 4578  GGAACTGTTTCGATGAGTTCCCAAAGTGCTTCATCGTGGGCGCCGACAACGTGGGCTCCAA 4637

Query 249  GCAGATGCAGAACATCCGTACCAGCCTGCGTGGACTGGCCGTGCTGCTTATGGGCAAGAA 308
          |||
Sbjct 4638  GCAGATGCAGAACATCCGTACCAGCCTGCGTGGACTGGCCGTGCTGCTTATGGGCAAGAA 4697

Query 309  CACCATGATGCGCAAGGCCATCCGCGGTTCATCTGGAGAACAACCCGCAGCTGGAGAAGCT 368
          |||
Sbjct 4698  CACCATGATGCGCAAGGCCATCCGCGGTTCATCTGGAGAACAACCCGCAGCTGGAGAAGCT 4757

Query 369  GCTACCCACATCAAGGGCAACGTGGGATTCGTGTTCAACAAGGGCGATCTCGCCGAGGT 428
          |||
Sbjct 4758  GCTACCCACATCAAGGGTAACGTGGGCTTCGTTTTCAACAAGGGCGATCTCGCCGAGGT 4817

Query 429  GCGCGACAAGCTGCTGGAGTCCAAGGTGCGCGCCCCCGCCCGTCCCGGCGCTATTGCCCC 488
          |||
Sbjct 4818  GCGTGACAAGCTGTTGGAGTCCAAGGTGCGCGCCCCCGCCCGTCCCGGCGCTATTGCCCC 4877

Query 489  TCTGCACGTCATCATCCCGGCGCAGAACACCGGCTTGGGACCCGAGAAGACCAGTTTCTT 548
          |||
Sbjct 4878  TCTGCACGTCATCATCCCGGCCCCAGAACACCGGCTTGGGACCCGAGAAGACCAGTTTCTT 4937

Query 549  CCAGGCCCTGTCCATCCCGACCAAAATTTCCAAGGGAACAATTGAAATCATCAACGATGT 608
          |||
Sbjct 4938  CCAGGCCCTGTCCATCCCGACCAAGATTTCCAAGGGAACAATTGAAATCATCAACGATGT 4997

Query 609  GCCCATCCTGAAGCCTGGCGACAAGGTGCGCGCCTCCGAGGCGACACTGCTCAACATGTT 668
          |||
Sbjct 4998  GCCCATCCTGAAGCCCGGCGACAAGGTGCGCGCCTCCGAGGCAACGCTGCTCAACATGTT 5057

Query 669  GAACATCTCGCCCTTCTCGTACGGTCTGATTGTCAACCAGGTCTACGACTCCGGCTCGAT 728
          |||
Sbjct 5058  GAACATCTCGCCCTTCTCGTACGGTTTGATCGTCAGCCAGGTGTACGACTCCGGCTCGAT 5117
  
```

Gene 3-Figure 4: Results from a nblast 2 of the mRNA sequence from CG7133 against my fosmid. The 3 exons from RpLP0 have homologous regions in the fosmid.

I ran my entire masked fosmid through sixpack to get all the reading frames. This gene was a lot trickier to match up to my fosmid than the other 2 genes. Nscan predicted only 2 exons though there are 3 in RpLP0. I decided to use the data from Gene 3-Figure 4 to find the exon/intron borders. When I used the first exon from RpLP0 to find the first starting codon in my gene 3 and then used the second and third exons to find the other 2 exons I got a peptide sequence that was much shorter than the one for RpLP0. I then looked at where translation starts in RpLP0 by running its mRNA through sixpack and looking for the first start codon. Translation starts at the bp 135 in the middle of exon 2. The first exon is non-coding in *D. melanogaster*, but when I looked at the homologous sequence for exon 1 in my fosmid I got a start codon where there was not one in *D. melanogaster*. When I compared exon 1 of RpLP0 with its homologous region in my fosmid I found that there was a mutation at bp 4224 in my fosmid. There is an “A” at the homologous position in *D. melanogaster* and a “T” in the fosmid. This mutation causes a premature stop codon. I have highlighted the mutation below in figure Gene 3-Figure 5 for clarity.

Score = 115 bits (62), Expect = 1e-27
 Identities = 85/96 (88%), Gaps = 2/96 (2%)
 Strand=Plus/Plus

```

Query 1  GGT-ATCTTATTTCGCCATCGAAGCGGTCACACTGGGTGCCGCGCCAACCTTCACTCTTTC 59
          |||
Sbjct 4204  GGTAATTTTAA-TCGCTATCGATGTTGGTTCACACTTGCTTCCGCGCCAACCTTCCCTCTTTC 4262

Query 60  CGTTCTGTGAGCGAAAACCGAAAAGTCTGTGCTTTG 95
          |||
Sbjct 4263  CGTTCTGTGAGCGAAAACCGAAAAGTCTGTGCTTTG 4298
  
```

Gene 3-Figure 5: Exon 1 from RpLP0 and its homologous region in the fosmid. The highlighted section shows the mutation from “A” in RpLP0 to “T” in the fosmid. This causes a premature start codon.

I thought that maybe this gene is a pseudogene then. I ran the mRNA of RpLP0 against the entire genome of *D. erecta* to see if the gene is present any where else in *D. erecta*, but I did not have any hits. I then read about the function of the RpLP0 gene and found that it is involved in translation, DNA repair, translational elongation, and ribosome biogenesis. The gene serves an important function. I concluded that a functional, real ortholog of RpLP0 has to exist in *D. erecta* based on the idea that the gene has an important function, is not found anywhere else in the genome, and that the transcript and peptide sequence of the predicted gene 2 from Nscan matches so closely to RpLP0. I decided to just use the coding sequence for the RpLP0 protein and find where it was homologous in my fosmid. I found the start codon in my fosmid to start at bp 4450 and for the exon to end at bp 4503 by looking at the homologous region of the end of exon 2 in RpLP0 and looking for a “GT” where 95% of introns start. I found the exon borders for exon 2 of gene 3 on the fosmid to be at 4579-5475 bp. I had to look at the region in the fosmid homologous to the start and end of exon 3 in RpLP0. In order to find the start of the exon I looked for an “AG” where 95% of introns end. I looked for the first stop codon in the reading frame to find the end of the third exon. I then ran gene 3 with its coordinates for the 2 exon through Gene checker which verified the gene. I ran a pblast 2 of the peptide sequence from Gene checker for gene 3 against the peptide sequence for RpLP0. The results are shown below in the figure labeled “Gene 3-Figure 6”.

name=RpLP0-PA; parent=FBgn0000100, FBtr0078481; dbxref=FlyBase:FBpp0078134, FlyBase_Annotation_IDs:CG7490-PA, GB_protein:AAF51807.1, REFSEQ:NP_524211, GB_protein:AAF51807; MD5=86e1796e988a2ee9e406941fb4905ecb; length=317; release=r5.17; species=Dmel; Length = 317

HSP # = 1 , Score = 550.821 bits (1418) , Expect = 3.80774e-157
Identities = 270 / 271 (99.6%) , Positives = 271 / 271 (100%)

Subject FASTA

```

Query: 1 MVRENKAAWKAQYFIKVVELFDEFKCFIVGADNVGSKQMQRNIRTSRGLAVVLMGKNTM 60
Subject: 1 MVRENKAAWKAQYFIKVVELFDEFKCFIVGADNVGSKQMQRNIRTSRGLAVVLMGKNTM 60

Query: 61 MRKAI RGHLENNPQLEKLLPHIKGNVGFVFTKGD LAEVRDKLLESKVRAPARPGA IAPLH 120
Subject: 61 MRKAI RGHLENNPQLEKLLPHIKGNVGFVFTKGD LAEVRDKLLESKVRAPARPGA IAPLH 120

Query: 121 VIIPAQNTGLGPEKTSFFQALS IPTKISKGTIEI INDPVILKPGDKVGASEATLLNMLNI 180
Subject: 121 VIIPAQNTGLGPEKTSFFQALS IPTKISKGTIEI INDPVILKPGDKVGASEATLLNMLNI 180

Query: 181 SPFSYGLIVSQVYDSGSIFSPEILDIKPEDLRAKFQQGVANLAAVCLSVGYPTIASAPHS 240
Subject: 181 SPFSYGLIVSQVYDSGSIFSPEILDIKPEDLRAKFQQGVANLAAVCLSVGYPTIASAPHS 240

Query: 241 IANGFNLLAIAATTEVEFK EATTIKEYIKD 271
Subject: 241 IANGFNLLAIAATTEVEFK EATTIKEYIKD 271

```

Gene 3-Figure 6: Results from a pblast 2 of the peptide sequence from Gene checker for gene 1 against the peptide sequence for CG7133.

The protein sequence is incredibly well conserved as shown by the 99.6% for identities, 100% for positives, and that the peptide sequence for both RpLP0 and gene 3 in the fosmid are both 271 a.a in length. This data supports the claim even more that his gene is a real gene.

Pseudogene 1 is homologous to Sfp79B in *D. melanogaster*. I based my conclusion that this gene is a pseudogene after analyzing it in several ways. First I noticed that neither Nscan nor Genscan provided a predicted gene that was an ortholog to Sfp79B which made me first think that it may be a pseudogene. I looked at the function of Sfp79B in *D. melanogaster*. The function is unknown though the protein stands for the seminal fluid

I ran my fosmid through sixpack and used the results of the nblast 2 to find where the start codon was for the ortholog of msopa in the fosmid. The start codon was at bp 11752 and the stop codon was at bp 11944. I ran a pblast 2 of the peptide sequence for msopa against the peptide sequence for the ortholog in the fosmid and only had the first 24 amino acids with any kind of similarity. Also, the length of the peptide sequence for msopa is 83 a.a. and the length of the peptide sequence in the ortholog is only 63 a.a. Both of these results just provided more evidence that the ortholog is a pseudogene. I decided to run a clustal2w of the coding sequence for the translation of msopa against the coding sequence for the translation of the ortholog. The results are shown below in the figure labeled, "Pseudogene 2-Figure 2."

```

FBtr0078482      ATGAACTTCATACAGATCGCCGTGCTGTTTCGTCTGGTTCGCAGTGGCCTT
Dere3_dna        ATGAACTTCCTACAGATCGCCTTGCTGGTGGTCTAGTGGCAGTGGCCTT
***** * ***** * ***** * ***** * ***** * *****

FBtr0078482      GGCCAGACCACAGGAAGATCC---GGCAAATCTGCCAGCTCCAGAGGCAGCAGCACC
Dere3_dna        GGCCAGAGCACAGGATGATCCACCGACAGATCTGCCAGCTCCAGACGCAACAAAACCACC
***** * ***** * ***** * * * ***** * ***** * * * *

FBtr0078482      ACCAGCAGCAGCAGCAGCACCACCAGCAGCAGCAGCAGC-----ACCACCAGCACCACCA
Dere3_dna        AGCAGCAGCAGCAGCTGGTGTCTCCAGCTGGTGTCCCGGGTAAAAATAACCAAAATGTCAA
* ***** * * ***** * * * * * * * * * * * * * * * * * *

FBtr0078482      GCACCACCAGCTGCAGCACCT---CAA
Dere3_dna        TCACAAC--GTTGTGACCATGGATAA
*** * * * * * * *

```

Pseudogene 2-Figure 2: Results of clustal2w for the coding sequence for the translation of msopa against the coding sequence for the translation of the ortholog.

There seem to be a lot of mutations in the 2 sequences. The highlighted sequences show where a "C" in msopa is mutated to a "T" in the ortholog causing a nonsense mutation. This is the reason that the protein coded for by the ortholog is 20 a.a. shorter. I concluded that the ortholog for msopa in the fosmid is a pseudogene because the sequences for the transcripts are only somewhat similar for the first 135 or so base pairs, the peptide sequences are only somewhat similar for the first 24 a.a., the ortholog codes for a protein that is 20 a.a. shorter than the msopa protein, and the ortholog codes for a protein that is too short to be a real functioning protein.

A chart has also been added that shows all the exon borders for the 3 real genes.

Gene	BP of exon 1	BP of exon 2
Gene1	1274-288	
Gene2	3823-3443	
Gene3	4450-4503	4579-5475

Genes-Figure-1: exon borders for the 3 genes in the fosmid

Clustal Analysis:

FBgn0105459 -----TCACGCACCTCGGC 14
FBgn0240193 NNNNNNNNNNNNGGCGGGAGCGCGCACCTTGGACTCCAACAGCTTGTGCGCACCTCGGC 60
EMBOSS_001 -----
FBgn0100484 -----TTTCGCTCACTGAACG 16

FBgn0105459 GAGATCGCCCTTGGTGAAAACGAAGCCCACGTTACCCTTGAT--GTGGGGCAGCAGCTTC 72
FBgn0240193 GAGATCGCCCTTGGTGAAACACGAAGCCCACGTTACCCTTGAT--GTGGGGCAGTAGCTTC 118
EMBOSS_001 -----TGGTGAAACACGAATCCCACGTTGCCCTTGAT--GTGGGGTAGCAGCTTC 47
FBgn0100484 GAAAGAGGGTAAGTTGGCGCCGAAATAGCGATGCGGAAAATATGTGAAATATCGATATT 76
* * * * *

FBgn0105459 TCCAGCTGCGGGTTGTTCTCCAGATGACCGCGGATGGCCTTGCGCATCATGGTGTCTTG 132
FBgn0240193 TCCAGCTGCGGGTTGTTCTCCAGATGACCGCGGATGGCCTTGCGCATCATGGTGTCTTG 178
EMBOSS_001 TCCAGCTGCGGGTTGTTCTCCAGATGACCGCGGATGGCCTTGCGCATCATGGTGTCTTG 107
FBgn0100484 TCCTCATGTCAAAAAAT-TCCGGAT-ATCGATTGATTTTGGAAAATTGTTCCGATACTT 134
*** ** * * * * * * * * *

FBgn0105459 CCCATAAGCACGACGGCCA-GTCC-ACGCAGGCTGGTACGGATGTTCTGCATCTGCTTGG 190
FBgn0240193 CCCATAAGCACGACGGCAA-GTCC-ACGCAGGCTGGTACGGATGTTCTGCATCTGCTTGG 236
EMBOSS_001 CCCATAAGCACGACGGCCA-GTCC-ACGCAGGCTGGTACGGATGTTCTGCATCTGCTTGG 165
FBgn0100484 CCAGGTTATAAAAAATAATGTCTTGGGCAGACGAGCCTGAACATTAATATGATTCTG 194
** * * * * * * * * *

FBgn0105459 AGCCAACGTTGTGCGGCGCCACGATGAA-GCACTTGGGGAACCTCATCGAACAGTTCCTGG 249
FBgn0240193 AGCCAACGTTGTGCGGCGCCACGATGAA-GCACTTGGGGAACCTCATCGAACAGTTCCTGG 295
EMBOSS_001 AGCCCACGTTGTGCGGCGCCACGATGAA-GCACTTGGGGAACCTCATCGAACAGTTCCTGG 224
FBgn0100484 AGAAAAAATGTCTGTGTCATCAAAATTGATTGTTATTTTACGCTTATTGT-CATAAACTGG 253
** * * * * * * * * *

FBgn0105459 AGGAGGAATTAGCGAGGATTAGTTAAGAGCATAAGCCAGG-GGCGAGCTAGTG-AGGGCG 307
FBgn0240193 AGGACGGATTAGCGAGGATTAGTTGGAAGCATATGCCAGG-GGCGAGGCAGGGGAGAGCG 354
EMBOSS_001 AGGAGGGATCAGTTAA---CAGCATAAGCCAGGCGCCAGCTAGTGAGCGAATTCTATAAG 281
FBgn0100484 CTC-GAAATTAGAAACCAATGATGCCAGACAGGCGTGAATCAGCGATGCTATAAAATTTAG 312
* * * * * * * * *

FBgn0105459 GTTCTATACTTACCACAACCTTGATGAAGTACTGAGCCTTCCATGCTGCCTTGTCTCCC 367
FBgn0240193 GTTCTATACTTACCACAACCTTGATGAAGTACTGAGCCTTCCATGCTGCCTTGTCTCCC 414
EMBOSS_001 GTTCTATACTTACCACAACCTTGATGAAGTACTGAGCCTTCCACGCTGCCTTGTCTCCC 341
FBgn0100484 AGATGATACTTTTTCGA-----TGAACGAGTATTTTCTGCCAAATGGGCAATATTTTTCG 367
***** * * * * * * * * *

FBgn0105459 TAACCATTTTIGATTTT-GTATTAGGGACTCGTCGGGTGAATTTAAGAAGCTGCGGAAGATG 426
FBgn0240193 TAACCATTTTIGATCGT-GTATTAGGGACTCGTCGGATGAATTTAAGAAGCTGCAAAAGACG 473
EMBOSS_001 TAACCATTTTAAATTGT-GTATTAGGGACTCGTCGGGTGAATTTAAGAAGCTGCGAAAGGCG 400
FBgn0100484 ATACCAATAGTATAACCAGACTAAGAATGTAAAGAATGTTTAAATAAATTTTCTGTTGTA 427

**** * ** * * * * * * * * * * * * *

FBgn0105459 TGATCTCGATTAGTCGGTACTGGGCCAGCGCCATGCAACAGTTATTCACCCAAAAAGCT 486
FBgn0240193 TCATCTCAATTAGTCAGCACTTGGCCAACACCAATTCAACAGTTATTGGACCGAAAATCCT 533
EMBOSS_001 TGATAACGATTAGTCTTTGCTGGACCACCGCCGTGGAACAGTTATGTACCCGAAAATGCT 460
FBgn0100484 ATTTATTCTTAATTACTTTTTGTACCAAGATTCCAATTTTACTTTCCATGCGATACTACT 487

* * * * * * * * * * * * * * * *

FBgn0105459 CGGGATGCAACATCTTTTCGCTCTTTTAATAATACTTACCAAAGCACAGACT-TTTCGGT 545
FBgn0240193 CGGGATGCAACAATTTTTCGGGCTTTTAGCATTACTTACCAAAGCACAGACT-TTTCGGT 592
EMBOSS_001 CGGGATGCAACATTATTCCGAACTTTTCAGCAACTTACCAAAGCACAGACT-TTTCGGT 519
FBgn0100484 TTCACCCGAATTGAAGCGCTCTCCGTTTCGGTTGAGTGATAGGTGCACAGATTATTTTGAA 547

** * * * * * * * * * * * * * * *

FBgn0105459 TTTTCGCTCACAGAACGGAAAGAGGGGAAGTTGGCGCCGGAAGCAAGTGTGACCACATCGAT 605
FBgn0240193 TTTTCGCTCACAGAACGGAAAGAGGGGAAGTTGGCGCCGGAAGCCAGTGTGACCGCATCGAT 652
EMBOSS_001 TTTTCGCTCACAGAACGGAAAGAGTGAAGTTGGCGCCGCGCACCCAGTGTGACCGCTTCGAT 579
FBgn0100484 GTTCAAATCCTTTTTTGTACAAGGACTTTAAGTACTTAATCGAGCA-GTCTGAATTTAT 606

*** * * * * * * * * * * * * * * *

FBgn0105459 AGCGATAAAATTACCATATGGTCTAAAAAAATACCAATTTGAAAAGCGACATTAATAT 665
FBgn0240193 TACGATAAAATTACCGCATGGTCTAAAAA--ATACCAATTCGATAAAAAGAAATTATTAT 710
EMBOSS_001 GGCGAATAAGATACCGCACGGTCTGAAAAA-ATACCAAACCGGTAAGCGATATGAATAT 638
FBgn0100484 TTCTTTAGAAATTC---AAAGTTTGACACAACACTACTATTTCGTTTTATATTGCATTA AAC 663

* * * * * * * * * * * * * * * *

FBgn0105459 TTT--ATCTTTGAAATATATAATTCCTTAAAAAATACATGAAAAAATAATTTAACAAAT 723
FBgn0240193 TTT--TTCTTTGAAATATATAATTCCTTGTGAAATACTTGAATAAAAAAAC----- 759
EMBOSS_001 TTTCTATTTTCGAAATGTGTAATTCCTTTTAAAGTACTCAAAAAAGATAGATTTA----- 693
FBgn0100484 AAA-----CAAACTAAATATTTTCTACAAATTTAAATATAGAAACAAATAAAC----- 713

** * * * * * * * * * * * * * * *

FBgn0105459 TATAAATTATTAATAAATGCCAAATGGGAGCTTAACAGTGCAGCTACACAGCTGCTGCGC 783
FBgn0240193 -ATGAATTATTAATTAATG-----GAAGCCTAA-AGTGCC-TTACA----- 797
EMBOSS_001 -ACAGATCATTAAATTAATA-----ATAGGCTAATAGTGCTCCTATACAGCTGCTACGA 745
FBgn0100484 -GTAAATCTGCATTGAACATAA-----ATATATTTATTTTATTATTTCATAATCTGTGTGA 768

** * * * * * * * * * * * * * * *

```

FBgn0105459 TAGACAGAATTTCTTCTGTAAAC--GTGTTAACAGTATTCTGGAGGTAGGGTCACACTTG 841
FBgn0240193 -----ATTACTTCTGTAAATA-GTGTTAACAGTATTTGAAGGTAGGGTTACTCTTG 848
EMBOSS_001 TCGTCAGAATCTCTTCTGTAACTAACAGTGCTTCTGGTGGTAAGGTAGGGTCACTCTGG 805
FBgn0100484 TATTT---TTGGGTTTTAATAAATTATTTAAAAAAGCAATTGTTTCAAATGGTCCAT 825
          *   * * *   * * * * * * * * * *

FBgn0105459 T-ACTGGCGGGCAATAAATTTATCGATTGGGATGTTAGTCT-GTTTATTTGCTAAATTTTC 899
FBgn0240193 TTACTGGCGGGCAATAAATTTATCGATTATC-IGTTAGTTTCAGTTTATTTGTTATATTGT 907
EMBOSS_001 CCACTGGCGGCCGATAAATTTATCGATTGAC-----TGTTTATTTGTAATAATGT 854
FBgn0100484 CCAATGGCGCTGTCTCTTTACTGGTCTTATCGATGCATCGATAACTTCATAACTTTCT 885
          * * * * * * * * * * * * * * * *

FBgn0105459 TTC--ATTGGCGAATCAAGTTACCAAAATATAC-----ACAAGCACCAACAAG 945
FBgn0240193 TTTCAATCAGCGAATTGAACTACCGAAATAT-----CACCAAG 945
EMBOSS_001 TTTAAATCAGCGAATCGAATCAGCAATATATACTCGCATTACCACATAAGCACCAACAAG 914
FBgn0100484 CTGGAAATCAAAACACAAGTTGAAATAAGATTCT-----GCAATAAATAAA 931
          * * * * * * * * * * * * * *

FBgn0105459 TTAATGCGCTGCCTCACATACTATAG-----CA 973
FBgn0240193 TTGAGGCGCTGCCTCACATAGCATA-----CA 973
EMBOSS_001 TCGAGACGCTGCCTCACATTCCATACACGGCGCAACACCGCTAGACGACCATCTCAGCCA 974
FBgn0100484 TTATTGTACCTGAGCACCCGATTAAGAGAAAAGGAA-----CA 969
          * * * * * * * * * * * *

FBgn0105459 GAAC---AAGGAG-AATCCCATTTCACCACA 1000
FBgn0240193 GAAC---AAGGAG-AATCCCATTTCACCACA 1000
EMBOSS_001 GAAC---AAGGAG-AATACCATTTCACCACA 1001
FBgn0100484 AAGCGITAGGGAGTAACAGTTCGTTATCACC 1000
          * * * * * * * * * * * *

```

Clustal-Figure-2: Results from a clustal2w of the first 1k of bp upstream from the start codon for gene 2(top sequence) and 3 of its orthologs in *D. melanogaster*(third sequence), *D. yakuba*(second sequence), and *D. ananassae*(fourth sequence).

I found the TATA box which is highlighted in the figure. Though there are deletions in the other sequences, if you look 1 bp upstream, there is a “T” which still makes the TATA box functional. I could not find the initiators or any DEP’s.

Repeats:

The fosmid was checked for repeats using RepeatMasker in order to eliminate any repetitious elements before the fosmid was checked for gene features. RepeatMasker generated 2 tables. The first table below labeled, “Repeats-Figure-1,” shows the total amount of each kind of repeat in the fosmid. The second table below labeled, “Repeats-Figure-2,” shows the base pair location of every repeat found in the fosmid.

	number of elements*	length occupied	percentage of sequence
Retroelements	0	0 bp	0.00 %
SINEs:	0	0 bp	0.00 %
Penelope	0	0 bp	0.00 %
LINEs:	0	0 bp	0.00 %
CRE/SLACS	0	0 bp	0.00 %
L2/CR1/Rex	0	0 bp	0.00 %
R1/LOA/Jockey	0	0 bp	0.00 %
R2/R4/NeSL	0	0 bp	0.00 %
RTE/Bov-B	0	0 bp	0.00 %
L1/CIN4	0	0 bp	0.00 %
LTR elements:	0	0 bp	0.00 %
BEL/Pao	0	0 bp	0.00 %
Ty1/Copia	0	0 bp	0.00 %
Gypsy/DIRS1	0	0 bp	0.00 %
Retroviral	0	0 bp	0.00 %
DNA transposons	10	1402 bp	2.80 %
hobo-Activator	0	0 bp	0.00 %
Tc1-IS630-Pogo	0	0 bp	0.00 %
En-Spm	0	0 bp	0.00 %
MuDR-IS905	0	0 bp	0.00 %
PiggyBac	0	0 bp	0.00 %
Tourist/Harbinger	0	0 bp	0.00 %
Other (Mirage, P-element, Transib)	3	692 bp	1.38 %
Rolling-circles	0	0 bp	0.00 %
Unclassified:	0	0 bp	0.00 %
Total interspersed repeats:		1402 bp	2.80 %
Small RNA:	0	0 bp	0.00 %
Satellites:	0	0 bp	0.00 %
Simple repeats:	4	212 bp	0.42 %
Low complexity:	5	302 bp	0.60 %
bases masked:		1916 bp (3.83 %)	

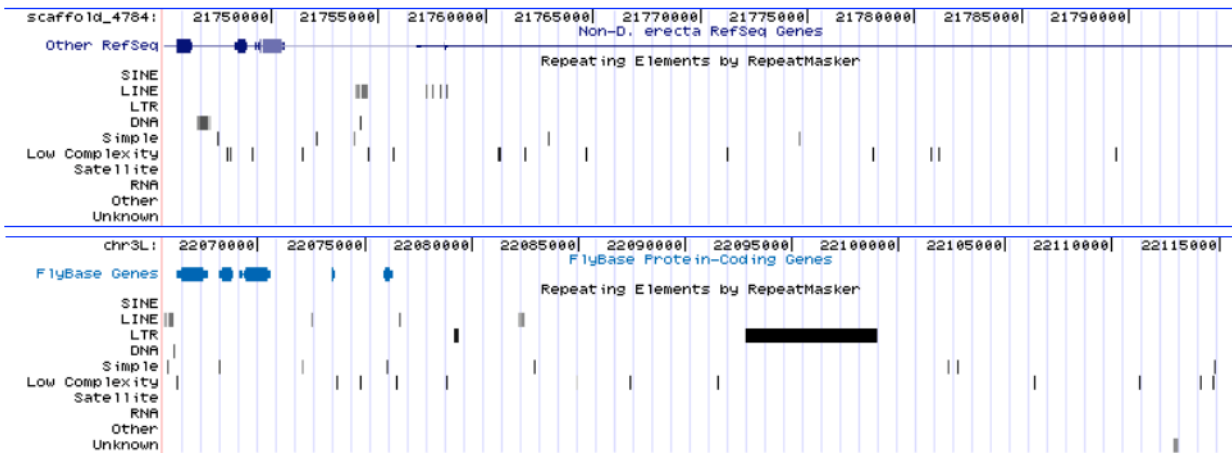
Repeats-Figure-1: total amount of each kind of repeat in the fosmid

+ score	% div.	% del.	% ins.	query sequence	-position in query-			C matching + repeat	repeat class/family	-position in repeat--			
					begin	end	(left)			(left) begin	end	begin (left)	linkage id/graphic
± 293	19.9	0.0	7.2	Dere3_dna	1572	1690	(48310)	C PROTOP_A	DNA/P	(296)	809	699	1 *
± 769	11.1	4.4	0.0	Dere3_dna	1662	1796	(48204)	C PROTOP	DNA/P	(452)	4028	3888	2 *
± 1665	11.8	1.9	0.4	Dere3_dna	1789	2059	(47941)	C PROTOP	DNA/P	(794)	3686	3412	2
± 226	19.1	9.7	10.3	Dere3_dna	2055	2229	(47771)	C PROTOP_B	DNA/P	(282)	871	698	3 *
± 495	7.3	0.0	0.0	Dere3_dna	2487	2554	(47446)	+ (TATG)n	Simple_repeat	3	70	(0)	4
± 22	75.4	0.0	0.0	Dere3_dna	3053	3109	(46891)	+ AT_rich	Low_complexity	1	57	(0)	5
± 35	82.1	0.0	0.0	Dere3_dna	4083	4166	(45834)	+ AT_rich	Low_complexity	1	84	(0)	6
± 26	57.7	0.0	0.0	Dere3_dna	6427	6452	(43548)	+ AT_rich	Low_complexity	1	26	(0)	7
± 201	7.4	0.0	0.0	Dere3_dna	7107	7133	(42867)	+ (CATCG)n	Simple_repeat	1	27	(0)	8
± 221	19.5	2.5	3.8	Dere3_dna	8851	8930	(41070)	+ (CTAA)n	Simple_repeat	2	80	(0)	9
± 405	11.1	0.0	1.6	Dere3_dna	8952	9015	(40985)	+ Helitron-1_DYak	DNA/Helitron	49	111	(10735)	10
± 503	14.3	12.4	0.8	Dere3_dna	9018	9130	(40870)	+ DNAREP1_DM	DNA/Helitron	469	594	(0)	11
± 276	5.9	0.0	0.0	Dere3_dna	9162	9195	(40805)	C PROTOP_A	DNA/P	(564)	541	508	1 *
± 1439	14.3	2.2	5.7	Dere3_dna	9194	9504	(40496)	+ DNAREP1_DM	DNA/Helitron	293	593	(1)	12
± 254	21.4	0.0	0.0	Dere3_dna	12181	12236	(37764)	+ DNAREP1_DM	DNA/Helitron	539	594	(0)	13
± 262	19.6	0.0	0.0	Dere3_dna	12522	12577	(37423)	+ DNAREP1_DM	DNA/Helitron	539	594	(0)	14
± 312	16.1	0.0	0.0	Dere3_dna	12849	12904	(37096)	+ DNAREP1_DM	DNA/Helitron	539	594	(0)	15
± 312	16.1	0.0	0.0	Dere3_dna	13176	13231	(36769)	+ DNAREP1_DM	DNA/Helitron	539	594	(0)	16
± 62	85.9	0.0	2.8	Dere3_dna	15606	15714	(34286)	+ AT_rich	Low_complexity	1	106	(0)	17
± 225	10.8	0.0	0.0	Dere3_dna	17882	17918	(32082)	+ (CA)n	Simple_repeat	2	38	(0)	18
± 26	53.9	0.0	0.0	Dere3_dna	44347	44372	(5628)	+ AT_rich	Low_complexity	1	26	(0)	19

Repeats-Figure-2: shows the base pair location of every repeat found in the fosmid

Synteny:

This section discusses the synteny between the fosmid and the homologous region in *D. melanogaster*. I looked to see if all the genes in the fosmid were from the same region of the *D. melanogaster* genome. I also looked for a similarity of repetitious elements. UCSC Genome Browser was used to get the area of *D. erecta* that the fosmid came from and to get the homologous area of *D. melanogaster* in order to compare the location of the genes and repetitious elements. Below in the figure labeled “Synteny-Figure-1” are the 2 areas from *D. erecta* and *D. melanogaster* that were compared.



Synteny-

Figure-1: 2 homologous areas from *D. erecta*(top) and *D. melanogaster*(bottom) from UCSC genome browser

The genes are generally from about the same region of the genome. *D. melanogaster* of course has the genes, Sfp79B and msopa, which are not present in *D. erecta* because they are pseudogenes. The simple and low complexity repeats seem to be the most conserved between the two genomes. *D. erecta* has a DNA repeat element

right after gene 1. This repeat seems to have caused the other 2 genes that come after gene 1 to shift to the right making those 2 genes farther away from gene 1 than the 2 homologous genes in *D. melanogaster* are from the first homologous gene in *D. melanogaster*. *D. melanogaster* also has the long LTR(long terminal repeat elements) that could be a transposon or retroposon.

Tools Used:

RepeatMasker- <http://www.repeatmasker.org/>

Extractseq- <http://gander.wustl.edu/cgi-bin/emboss>

Sixpack- <http://gander.wustl.edu/cgi-bin/emboss>

Flybase- <http://flybase.org/>

UCSC Genome Browser- <http://genome.ucsc.edu/>

Clustal2w- <http://www.ebi.ac.uk/Tools/clustalw2/index.html>

NCBI blast- <http://blast.ncbi.nlm.nih.gov/Blast.cgi>

References