# Implementing RNA-Seq data in FlyBase chado

David Emmert

GMOD Community Meeting
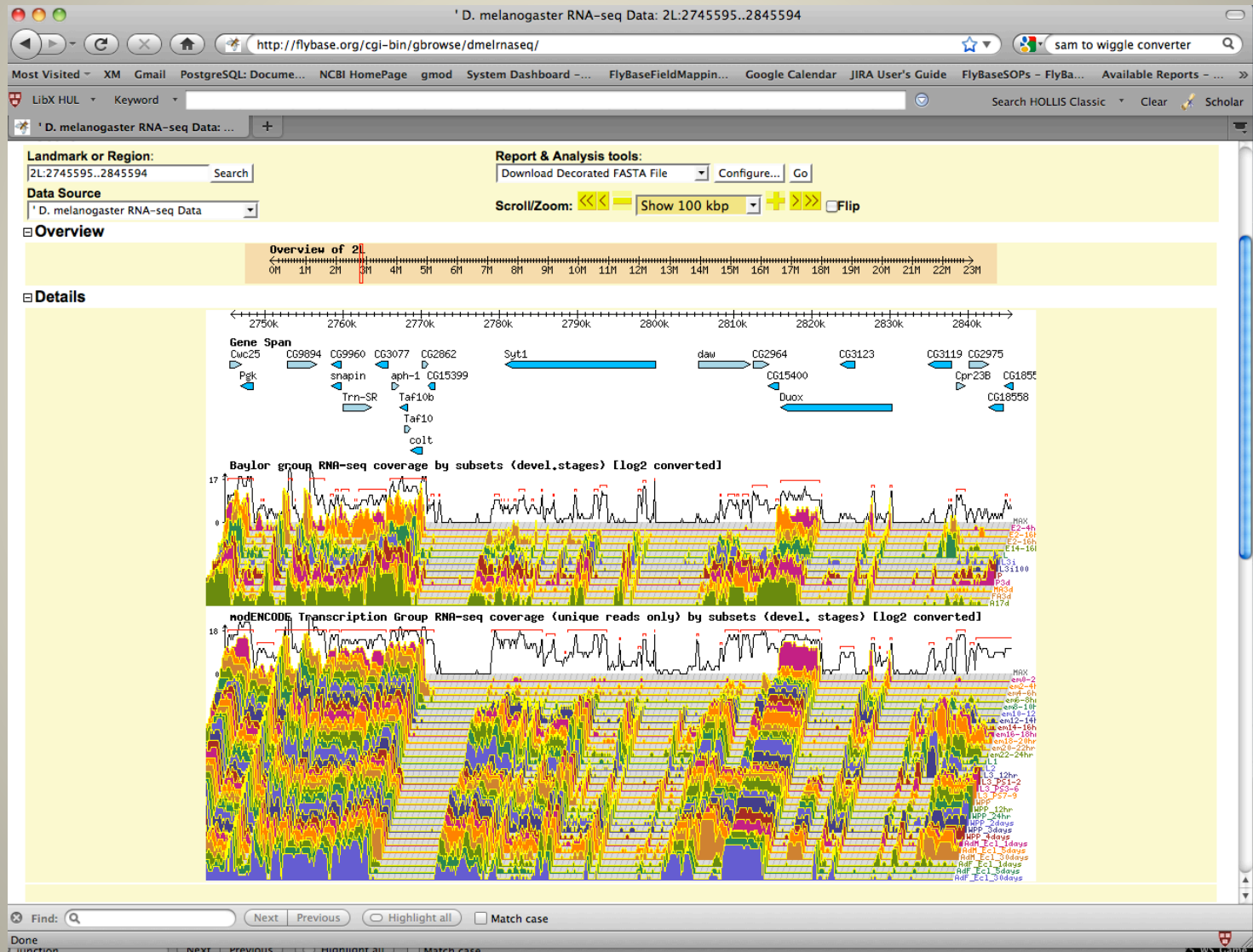
March 6&7, 2011

# RNA-Seq Data

- Data sets:
  - Daines et. al. 2010
    - Solexa/Illumina GAII 65, 75 & 100nt paired-end reads; 12 developmental stages
      - 142.2 million uniquely mapped reads (unstranded)
      - 54,594 unique junctions

  - Graveley et. al. 2010 (modENCODE)
    - Illumina GAII 75 & 76nt single- and paired-end reads; 30 developmental stages
      - 2.25 billion uniquely mapped reads (unstranded)
      - 67,317 unique junctions

# RNA-Seq Coverage Data

- Graphically represented in Gbrowse (but not integrated).

# RNA-Seq Coverage Data

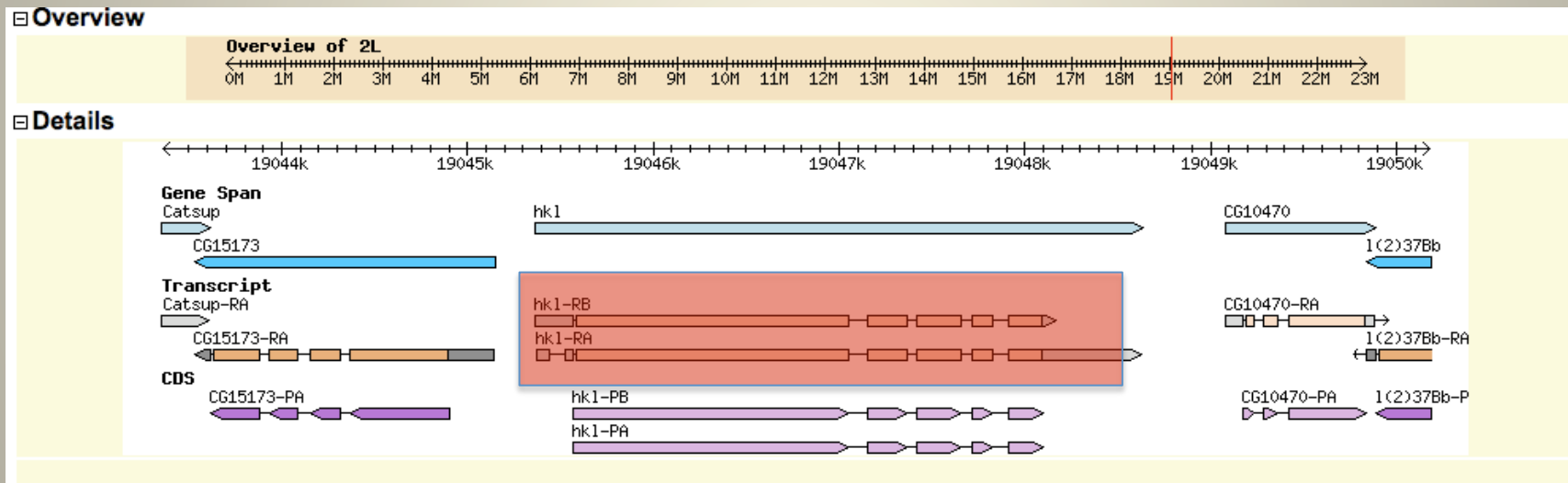- Graphically represented in Gbrowse (but not integrated).

# RNA-Seq Coverage Data

- Integration with genes:

  - Assign expression CV terms to genes based on coverage data.

  - Summarize HT expression data in gene reports.

  - Enable gene search by expression pattern.

  - Enable search for similarly expressed genes.
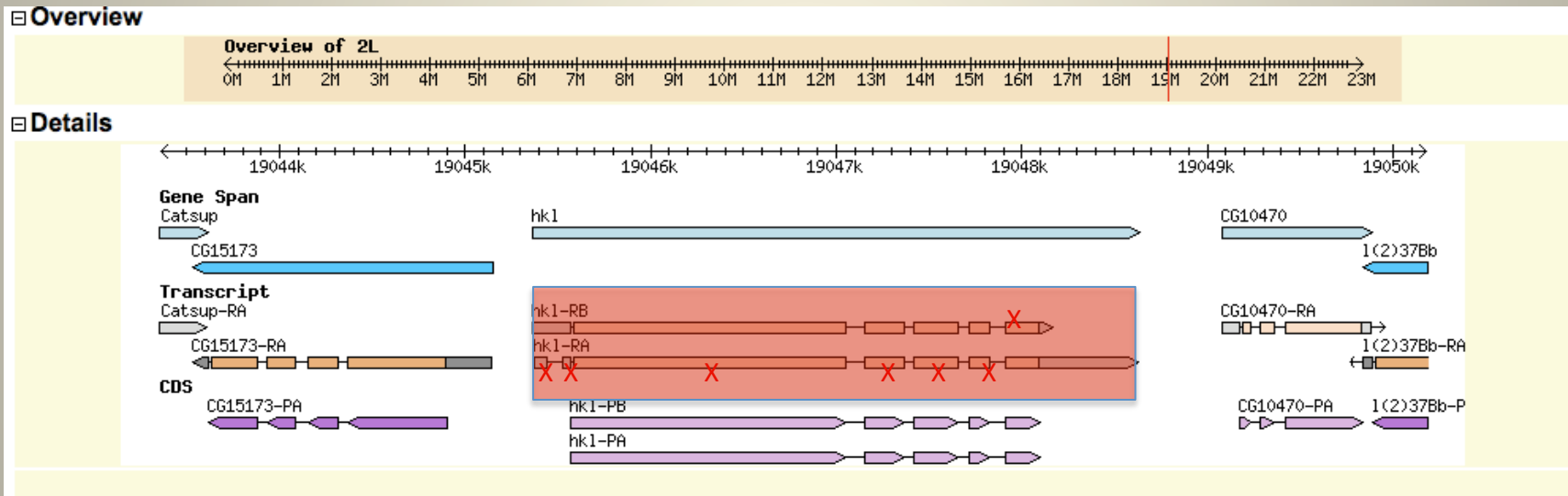
# RNA-Seq Coverage Data

- ## Integration with genes:

1) Determine unique transcribed region for each gene, e.g. *Dmel\hkl*:

# RNA-Seq Coverage Data

- ## Integration with genes:

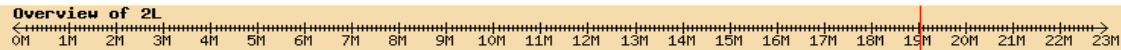1) Determine unique transcribed region for each gene, e.g. *Dmel\hkl*:
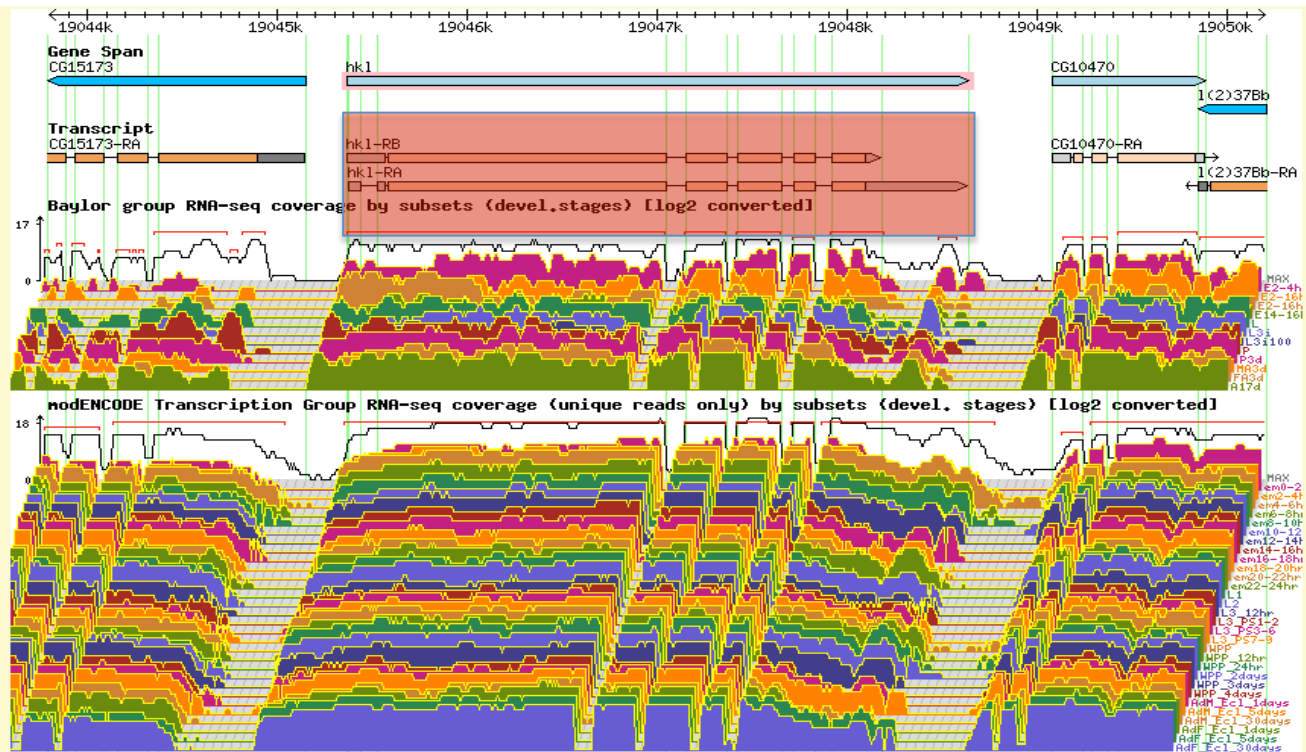
# RNA-Seq Coverage Data

- Integration with genes:

2) Correlate coverage for each stage over unique transcribed region:

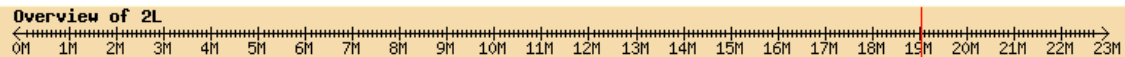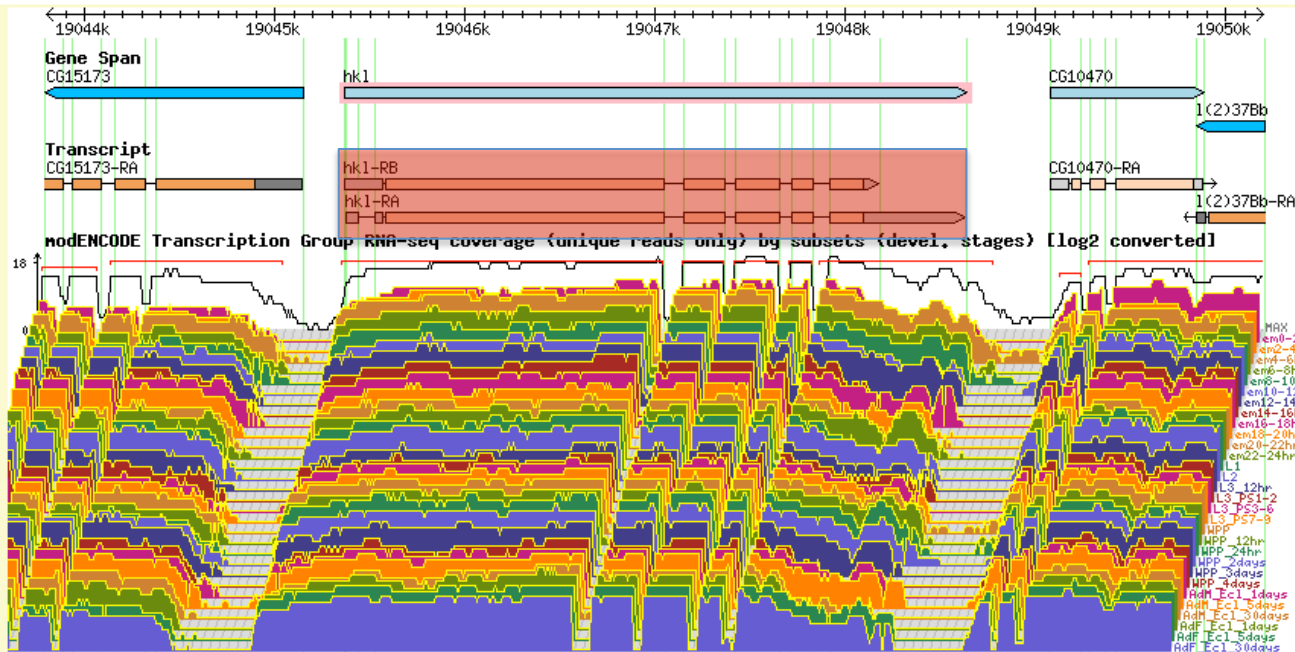# RNA-Seq Coverage Data

- ## Integration with genes:

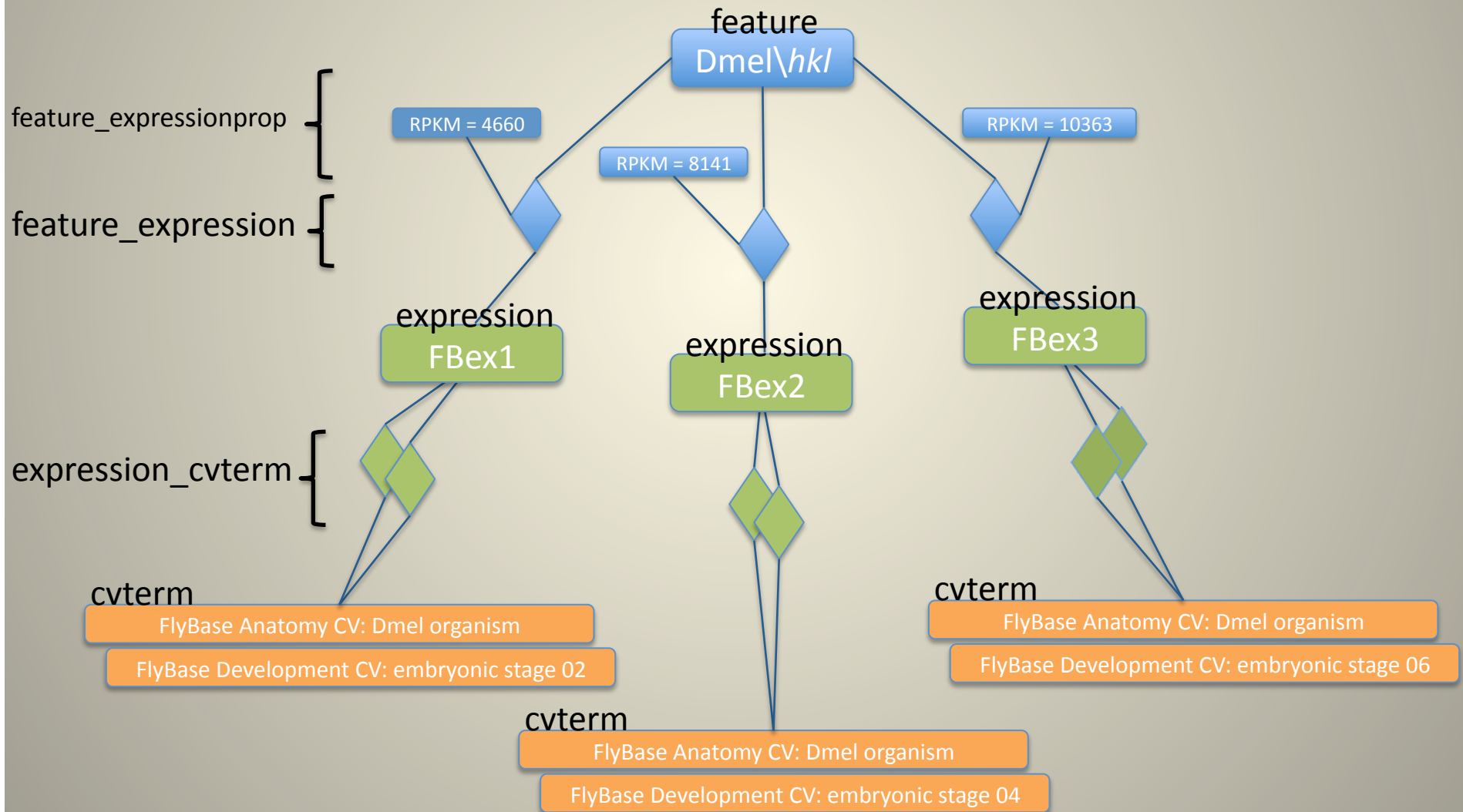3) Calculate $RPKM = \dfrac{total\_exon\_reads}{exon\_length(KB) \times mapped\_reads(millions)}$

*RPKM: **R**eads **p**er **k**ilobase of exon **m**odel per million mapped reads*

# RNA-Seq Coverage Data

- ## Integration with genes:

3) Calculate $\quad RPKM = \dfrac{total\_exon\_reads}{exon\_length(KB) \times mapped\_reads(millions)}$

*RPKM: **R**eads **p**er **k**ilobase of exon **m**odel per million mapped reads*

```
## Coverage for gene:      FBgn0086441 / hkl
## Unique transcribed bases: 2963
```

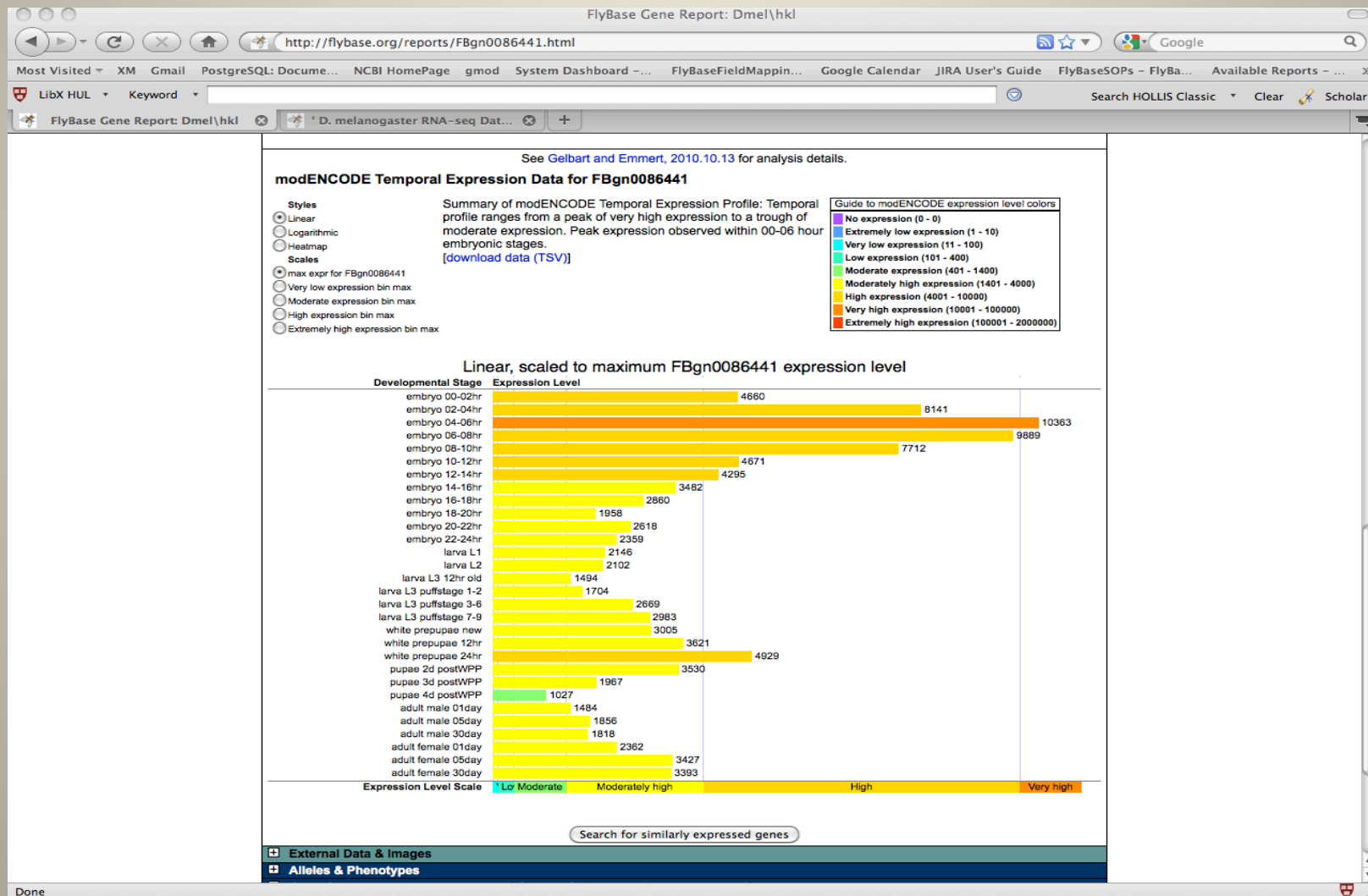| Gene_symbol | Stage | Stage_reads_total | Gene_reads_total | RPKM |
|---|---|---|---|---|
| hkl | embryos0-2hr | 65770867 | 908115 | 4660 |
| hkl | embryos2-4hr | 63321076 | 1527491 | 8141 |
| hkl | embryos4-6hr | 112427066 | 3452237 | 10363 |
| hkl | embryos6-8hr | 72780472 | 2132522 | 9889 |
| hkl | embryos8-10hr | 63545567 | 1452011 | 7712 |
| hkl | embryos10-12hr | 80997587 | 1121070 | 4671 |
| hkl | embryos12-14hr | 97516583 | 1240992 | 4295 |
| hkl | embryos14-16hr | 72245981 | 745353 | 3482 |
| hkl | embryos16-18hr | 79062619 | 670067 | 2860 |
| hkl | embryos18-20hr | 83856061 | 486408 | 1958 |
| hkl | embryos20-22hr | 56404806 | 437496 | 2618 |
| hkl | embryos22-24hr | 79445908 | 555282 | 2359 |
| hkl | L1larvae | 83803085 | 532850 | 2146 |
| hkl | L2larvae | 103442079 | 644133 | 2102 |
| hkl | L3larvae_12hr_post_molt | 55507157 | 245667 | 1494 |
| hkl | L3larvaePS_1-2 | 51235228 | 258629 | 1704 |
| hkl | L3larvaePS_3-6 | 55653242 | 440091 | 2669 |
| hkl | L3larvaePS_7-9 | 66802321 | 590393 | 2983 |
| hkl | white_prepupae | 82817561 | 737369 | 3005 |
| hkl | WPP_12hr | 76325015 | 818933 | 3621 |
| hkl | WPP_24hr | 71693929 | 1047018 | 4929 |
| hkl | pupae_WPP_2d | 85237993 | 891651 | 3530 |
| hkl | pupae_WPP_3d | 88942645 | 518304 | 1967 |
| hkl | pupae_WPP_4d | 77120269 | 234787 | 1027 |
| hkl | adult_male_1d | 77337299 | 340044 | 1484 |
| hkl | adult_male_5d | 95313901 | 524272 | 1856 |
| hkl | adult_male_30d | 67006363 | 361043 | 1818 |
| hkl | adult_female_1d | 88238878 | 617421 | 2362 |
| hkl | adult_female_5d | 65880241 | 669056 | 3427 |
| hkl | adult_female_30d | 66441798 | 668000 | 3393 |

# RNA-Seq Coverage Data

- Chado implementation:

# RNA-Seq Coverage Data

- Reporting & Searching:

# RNA-Seq Coverage Data

- Reporting & Searching:

# RNA-Seq Coverage Data

- Reporting & Searching:

# RNA-Seq Coverage Data

- Reporting & Searching:

# RNA-Seq Junction Data

- Graphically represented in Gbrowse (neither consolidated nor integrated).

# RNA-Seq Junction Data

- Graphically represented in Gbrowse (neither consolidated nor integrated).

# RNA-Seq Junction Data

- Graphically represented in Gbrowse (neither consolidated nor integrated).

# RNA-Seq Junction Data

- Graphically represented in Gbrowse (consolidated but not integrated).

# RNA-Seq Junction Data

- Graphically represented in Gbrowse (consolidated but not integrated).

# RNA-Seq Junction Data

- Consolidation
  - Create persistent exon_junction records in chado.
  - Consolidate junctions by location.
  - Reduces number of records required in DB.
    - E.g., modENCODE set reduces from 1.7M to 67K unique
- Integration
  - Correlate predicted junctions with junctions from gene-model annotations.
  - Correlate junctions with transcripts, genes, cDNAs.

# RNA-Seq Junction Data

- Junctions from all sources:

Total unique RNA-Seq Junctions: 71,082

RNA-Seq
Daines et. al.
Total: 54,594

2,440    3,999    14,597

RNA-Seq
Graveley et. al.
Total: 67,317

46,830

1,325

1,891

3,010

FlyBase gene-model annotation*

* FlyBase release FB2011_02; Dmel r5.34

# RNA-Seq Junction Data

- Correlation of non-matching predicted junctions to annotated gene-model junctions*

RNA-Seq
Daines et. al.
Total: 54,594

2,440

3,999

14,597

RNA-Seq
Graveley et. al.
Total: 67,317

| Type of Correlation | Count |
|---|---|
| Predicted Alternative Junction | 1794 |
| Predicted Exon-Skip Alt. Junction | 606 |
| Predicted Alt. Junction Differing in Splice Donor Site | 7488 |
| Predicted Alt. Junction Differing in Splice Acceptor Site | 5257 |
| Predicted Novel Junction | 5891 |

* FlyBase release FB2011_02; Dmel r5.34

# RNA-Seq Junction Data

- Eventually incorporate junctions from aligned cDNAs & ESTs…



RNA-Seq
Daines et. al.
Total: 54,594

2,440

3,999

14,597

RNA-Seq
Graveley et. al.
Total: 67,317

46,830

1,325

3,010

1,891

FlyBase gene-model annotation*

Aligned cDNAs & ESTs

* FlyBase release FB2011_02; Dmel r5.34

# RNA-Seq Junction Data

- Chado implementation (localization):

| | |
|---|---|
| —— | feature.type = match |
| ···· | feature.type = exon_junction |
| → | featureloc |
| → | feature_relationship |

*part of*     *part of*

Junction

Exonic regions

Genomic Contig

# RNA-Seq Junction Data

- Chado implementation (metadata):

# Acknowledgements

| Bill Gelbart | Andy Schroeder |
|---|---|
| Jim Thurmond | Pinglei Zhou |
| Victor Strelets | Josh Goodman |