

PeakSegJoint: fast supervised peak detection via joint segmentation of count data samples

Toby Dylan Hocking
toby.hocking@mail.mcgill.ca
joint work with Guillaume Bourque

June 3, 2015

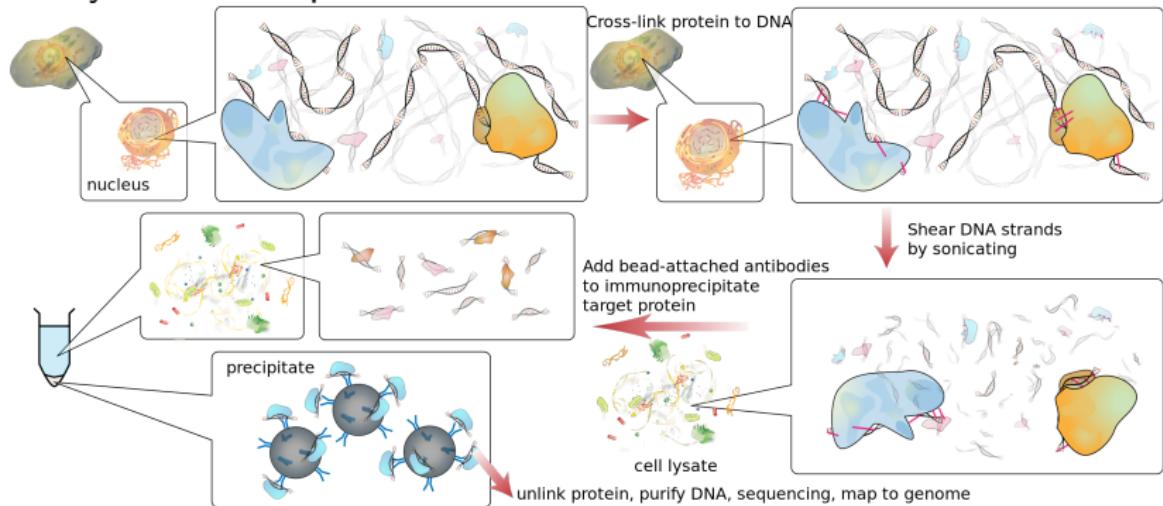
ChIP-seq data and previous work on peak detection

The PeakSegJoint model

Conclusions

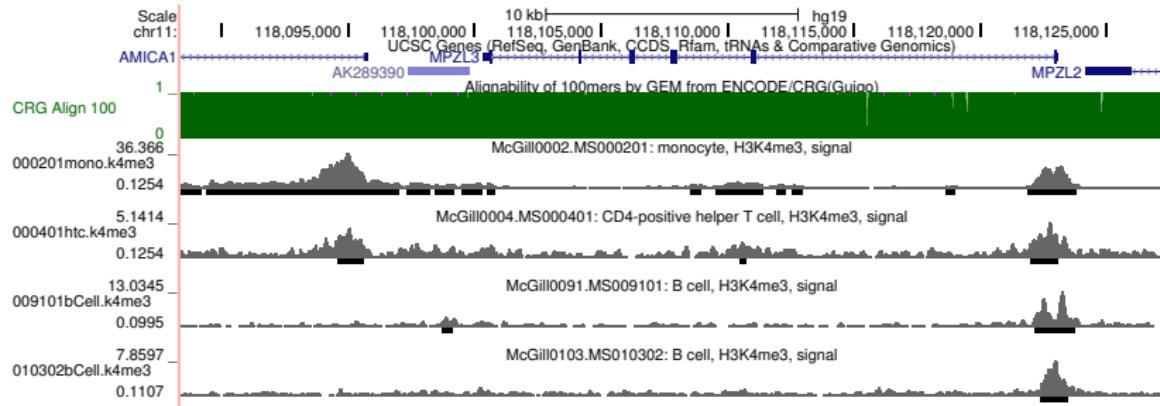
Chromatin immunoprecipitation sequencing (ChIP-seq)

Analysis of DNA-protein interactions.



Source: “ChIP-sequencing,” Wikipedia.

Problem: find peaks in each of several samples



Grey profiles are normalized aligned read count signals.

Black bars are “peaks” called by MACS2 (Zhang et al, 2008):

- ▶ many false positives.
- ▶ overlapping peaks have different start/end positions.

Existing peak detection algorithms

- ▶ Model-based analysis of ChIP-Seq (MACS), Zhang et al, 2008.
- ▶ SICER, Zang et al, 2009.
- ▶ HOMER findPeaks, Heinz et al, 2010.
- ▶ RSEG, Song and Smith, 2011.
- ▶ Histone modifications in cancer (HMCan), Ashoor et al, 2013.
- ▶ ... dozens of others.

Two big questions: how to choose the best...

- ▶ ...algorithm?
- ▶ ...parameters?

How to choose model parameters?

19 parameters for Model-based analysis of ChIP-Seq (MACS), Zhang et al, 2008.

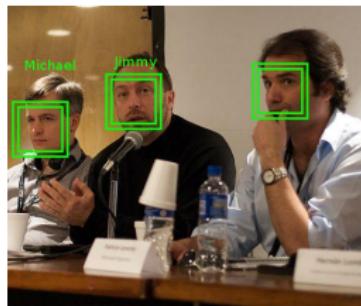
```
[-g GSIZEx  
[-s TSIZE] [--bw BW] [-m MFOLD MFOLD] [--fix-bimodal]  
[--nomodel] [--extsize EXTSIZE | --shiftsize SHIFTSIZE]  
[-q QVALUE | -p PVALUE | -F FOLDENRICHMENT] [--to-large]  
[--down-sample] [--seed SEED] [--nolambda]  
[--slocal SMALLLOCAL] [--llocal LARGELOCAL]  
[--shift-control] [--half-ext] [--broad]  
[--broad-cutoff BROADCUTOFF] [--call-summits]
```

10 parameters for Histone modifications in cancer (HMCan), Ashoor et al, 2013.

```
minLength 145  
medLength 150  
maxLength 155  
smallBinLength 50  
largeBinLength 100000  
pvalueThreshold 0.01  
mergeDistance 200  
iterationThreshold 5  
finalThreshold 0  
maxIter 20
```

Previous work in computer vision: look and add labels to...

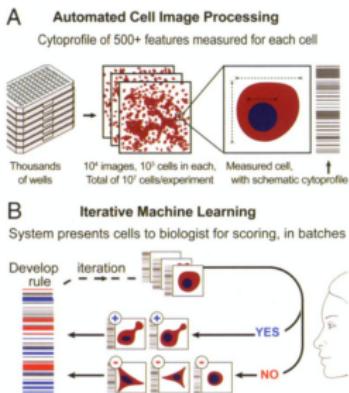
Photos



Labels: names

CVPR 2013
246 papers

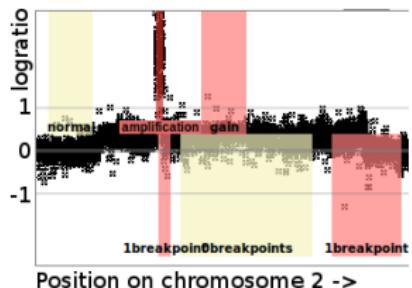
Cell images



phenotypes

CellProfiler
873 citations

Copy number profiles



alterations

SegAnnDB
Hocking et al, 2014.

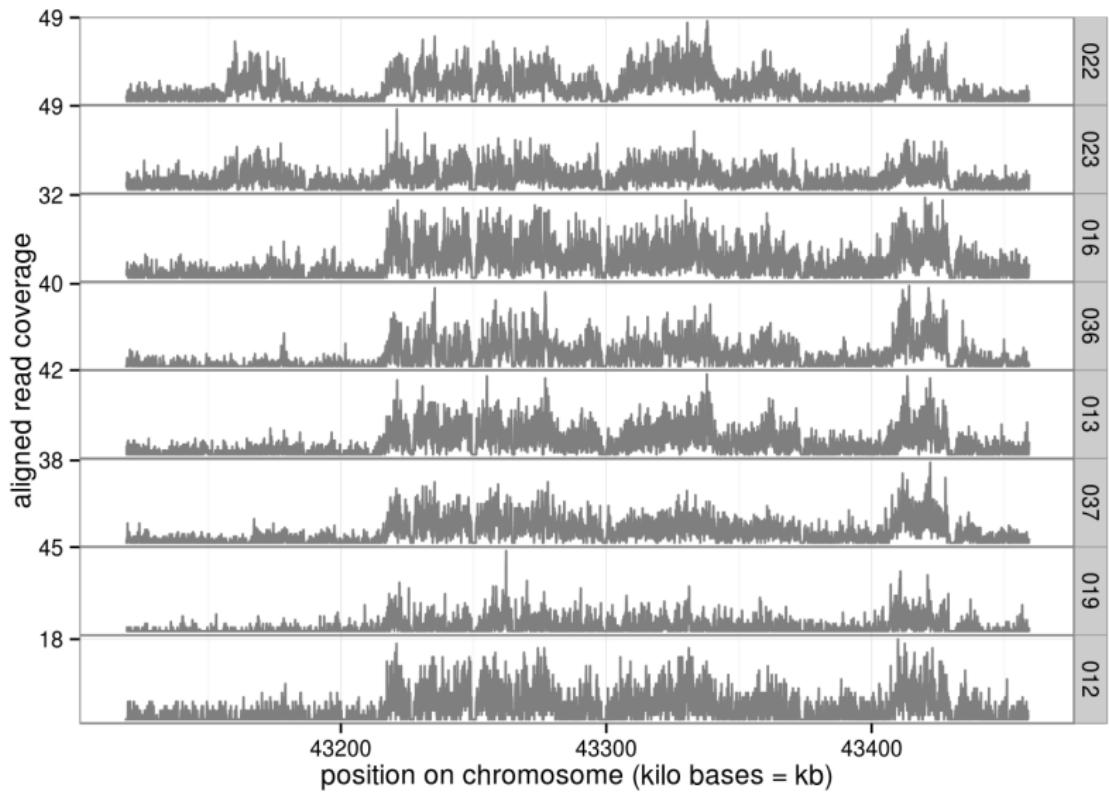
Sources: http://en.wikipedia.org/wiki/Face_detection
Jones et al PNAS 2009. Scoring diverse cellular morphologies in image-based screens with iterative feedback and machine learning.

ChIP-seq data and previous work on peak detection

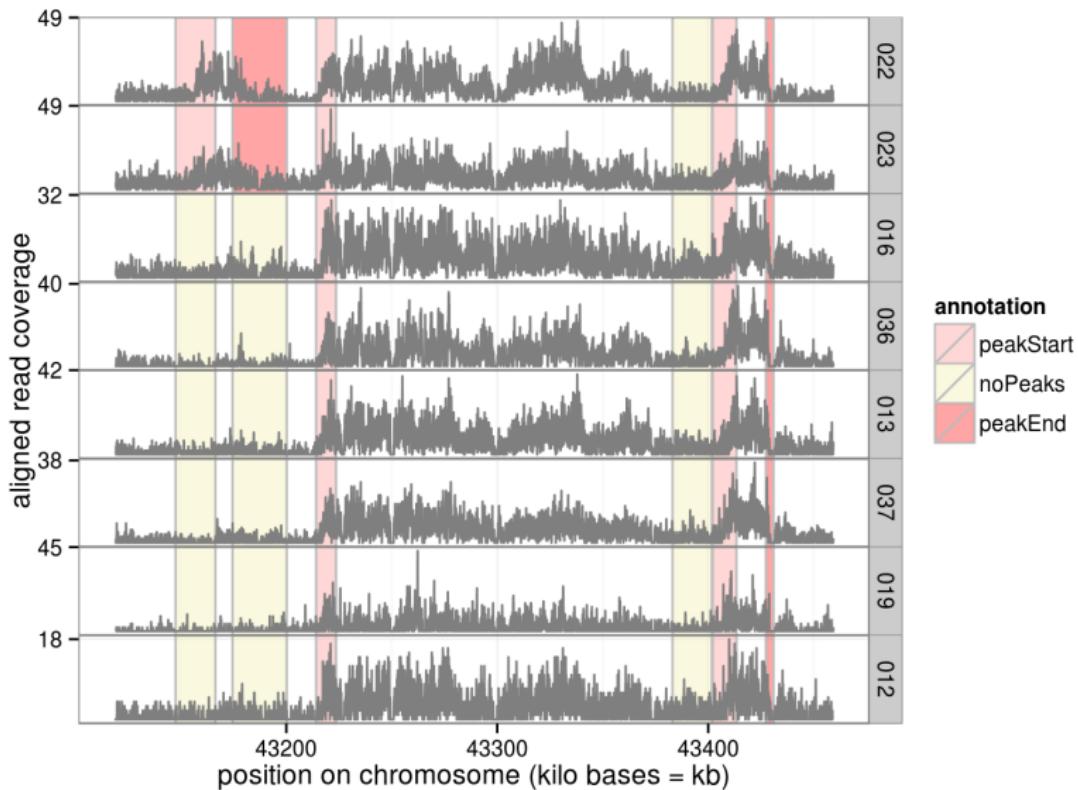
The PeakSegJoint model

Conclusions

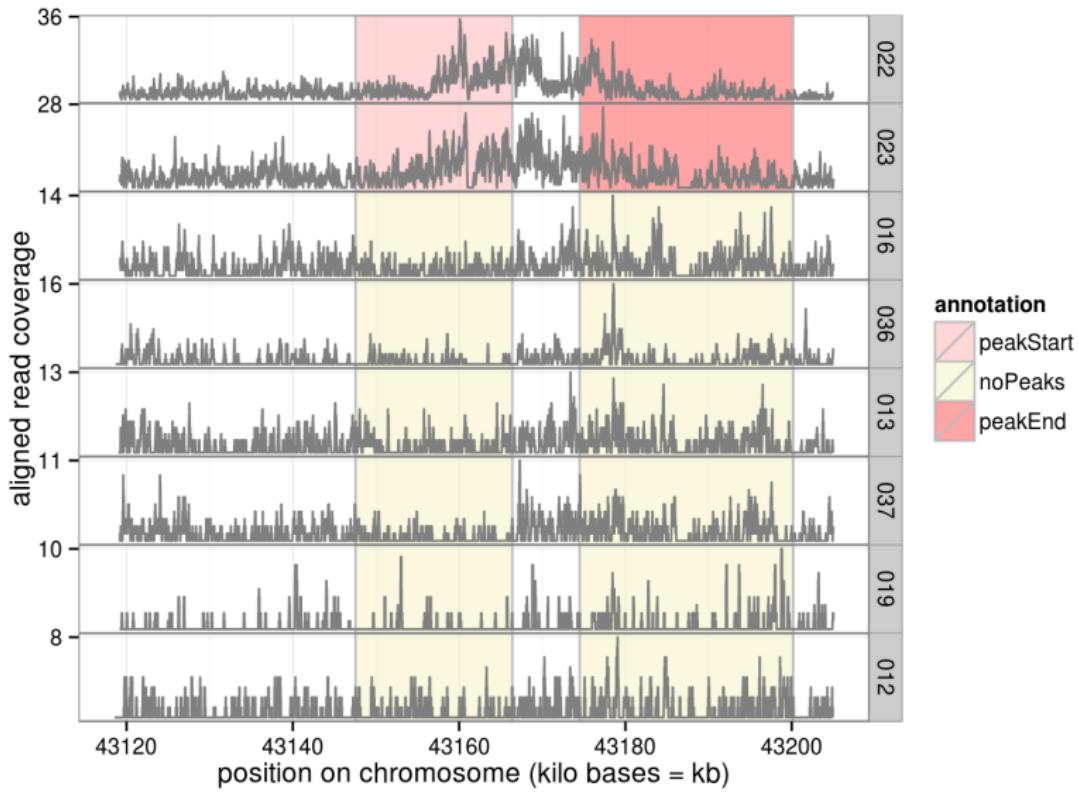
Peaks visually obvious in H3K36me3 data



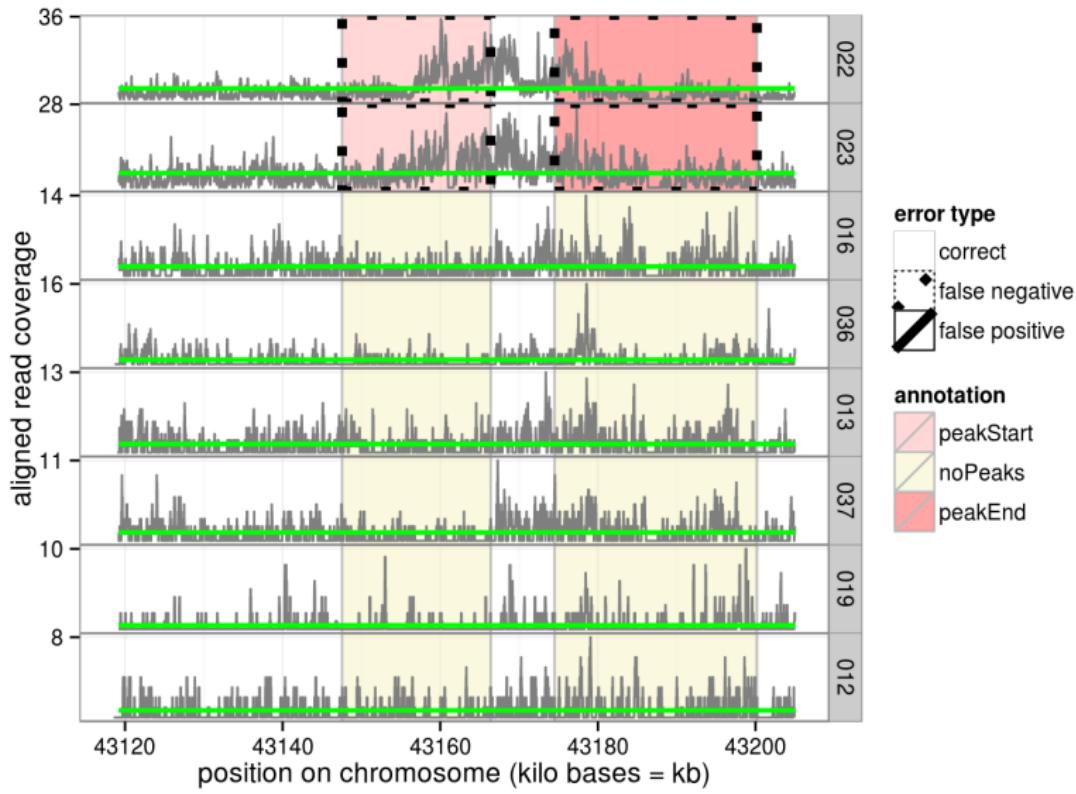
H3K36me3 data and visually determined labels



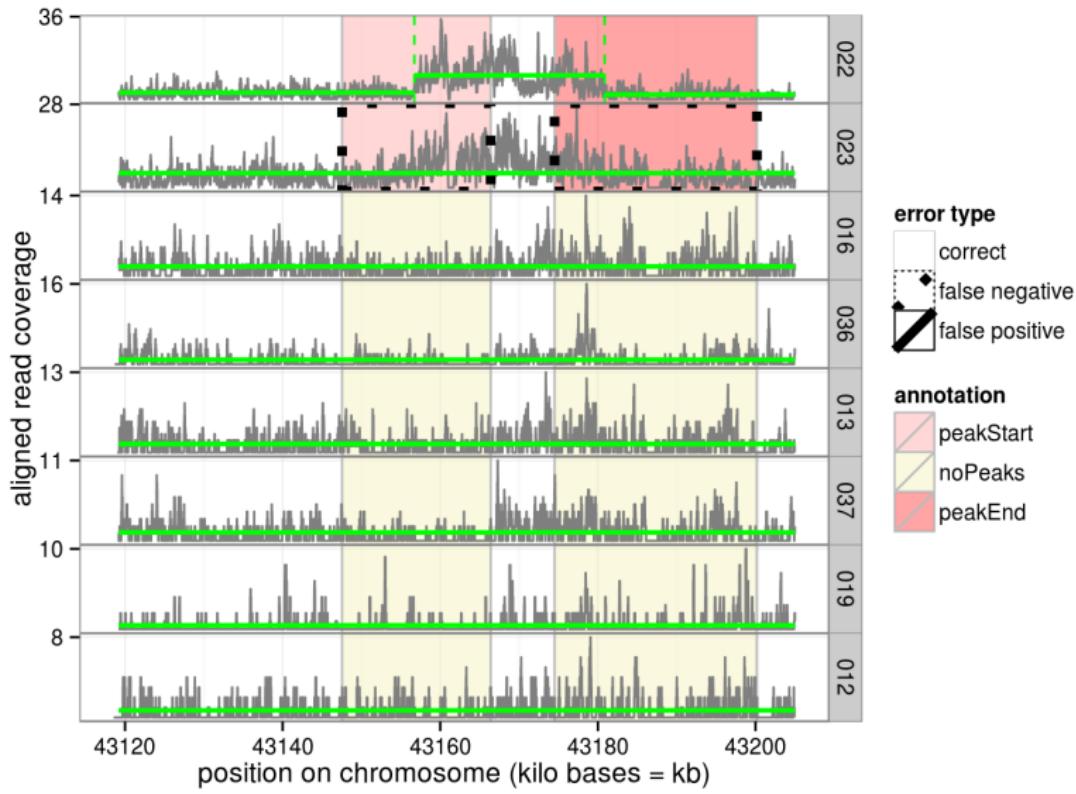
H3K36me3 data and labels (zoom to one peak)



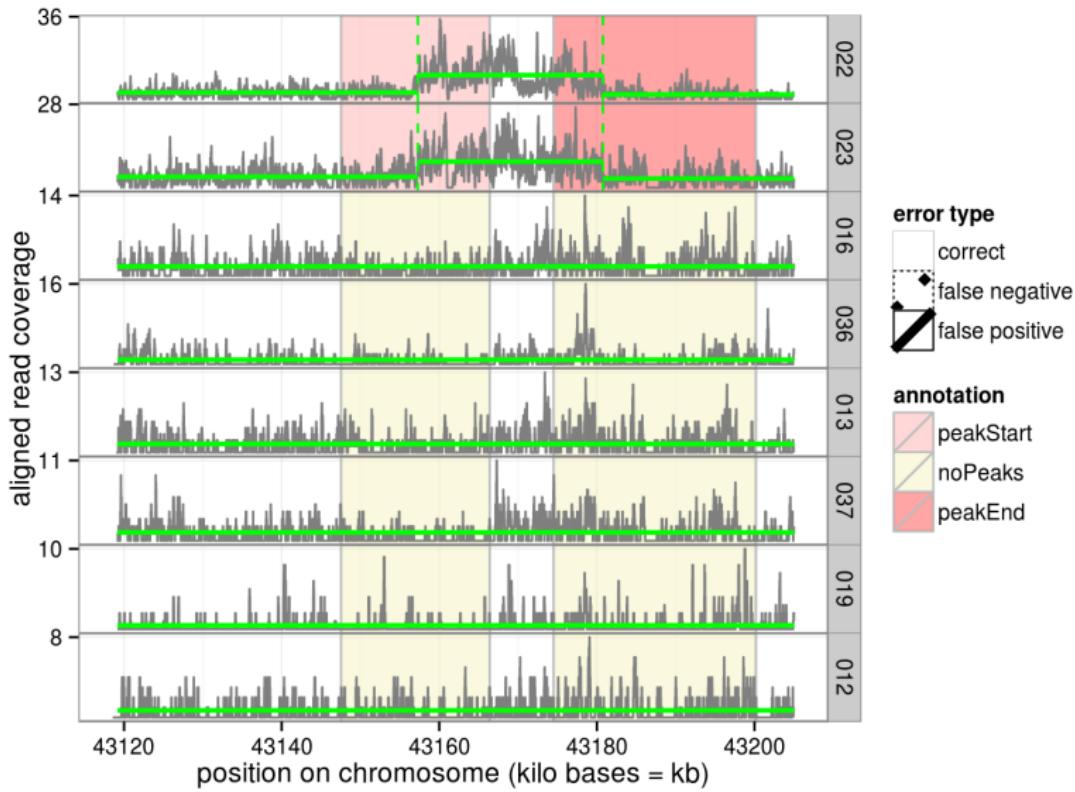
PeakSegJoint model with 0 peaks



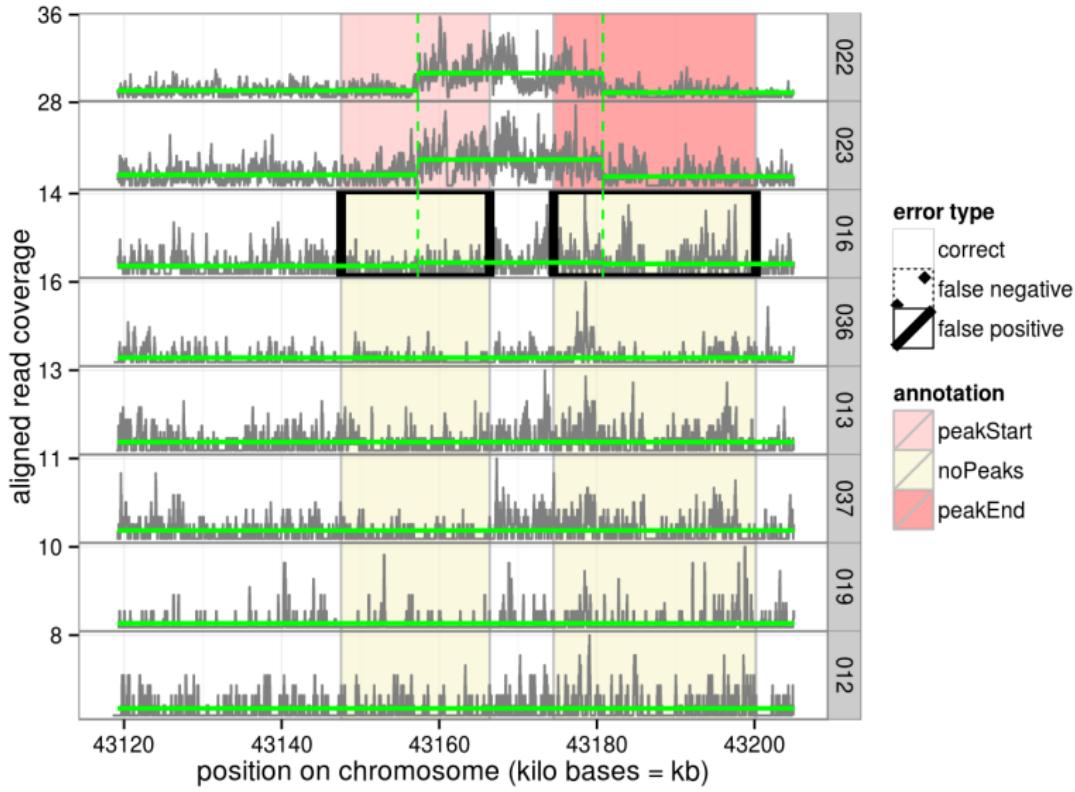
PeakSegJoint model with 1 peak



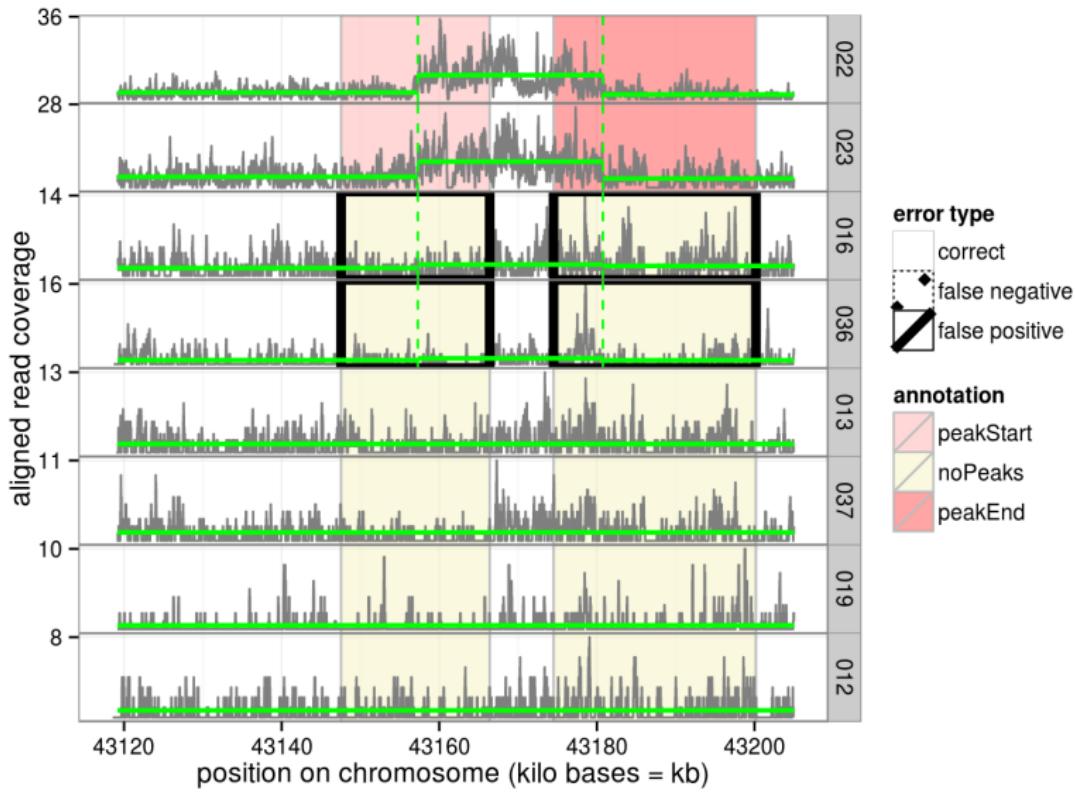
PeakSegJoint model with 2 peaks



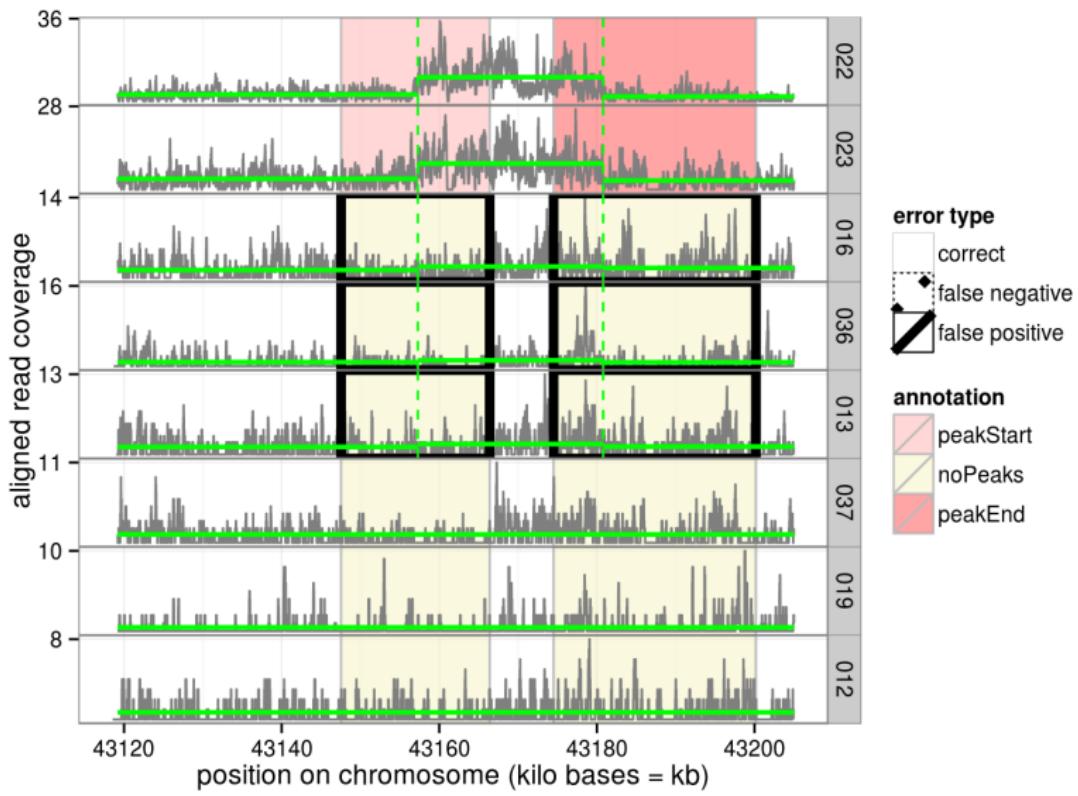
PeakSegJoint model with 3 peaks



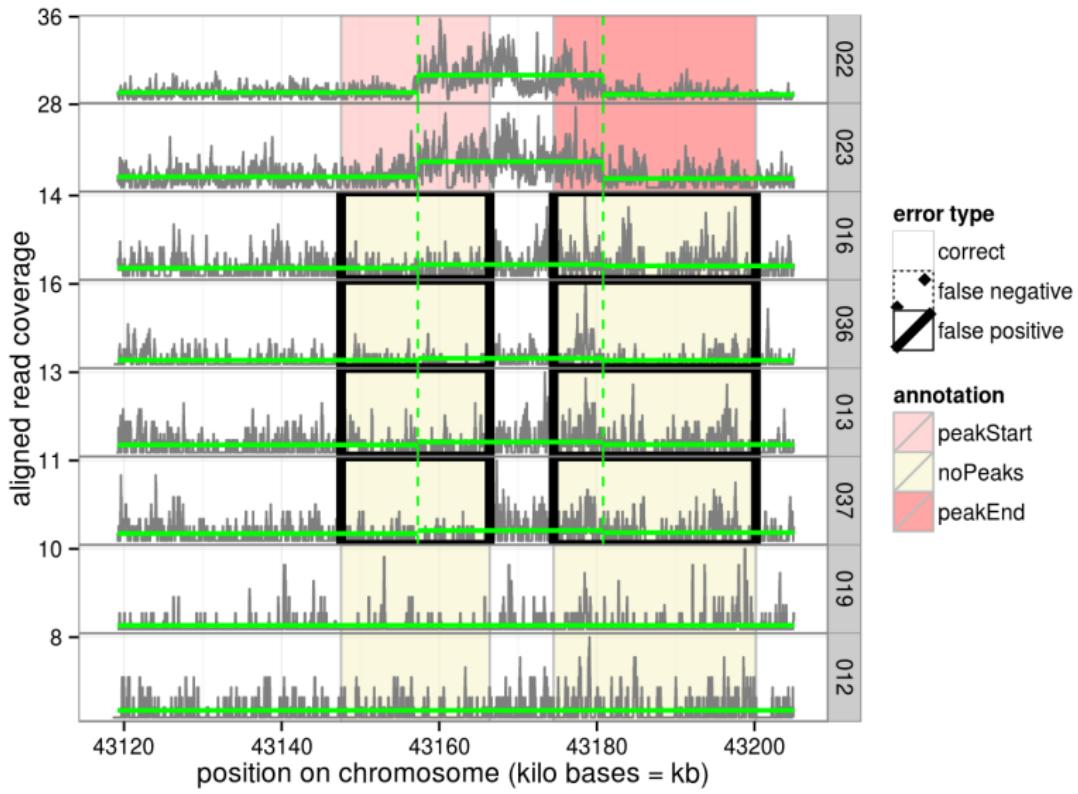
PeakSegJoint model with 4 peaks



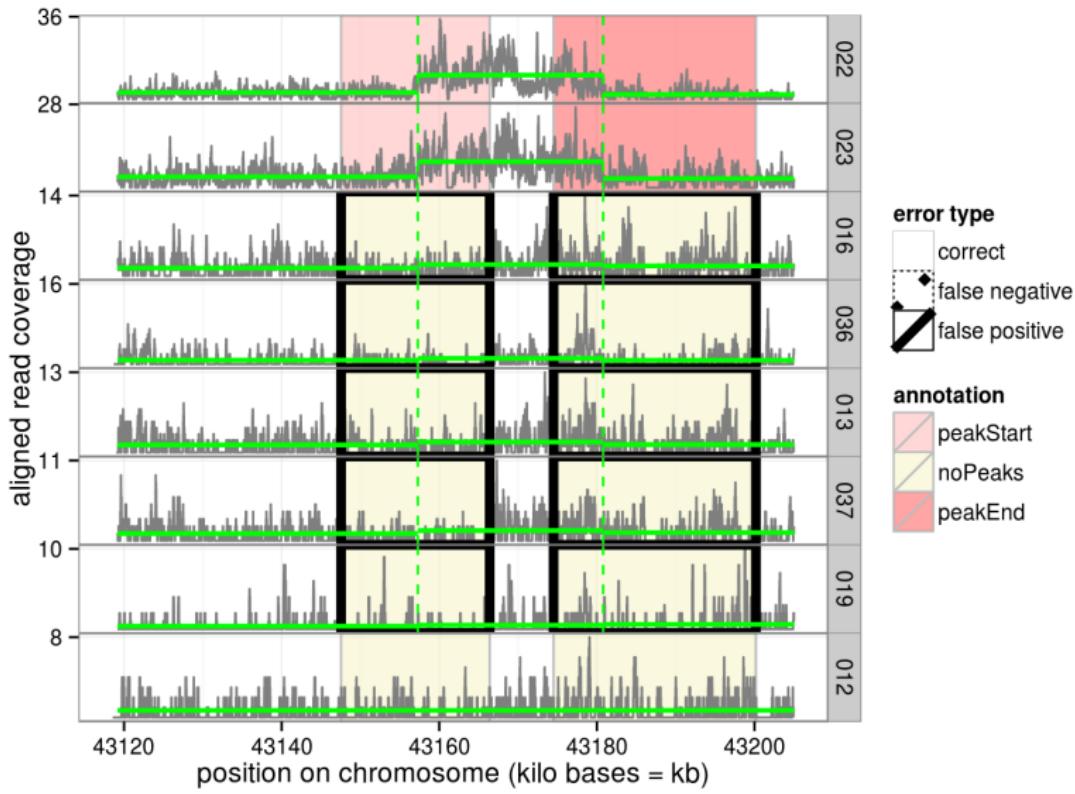
PeakSegJoint model with 5 peaks



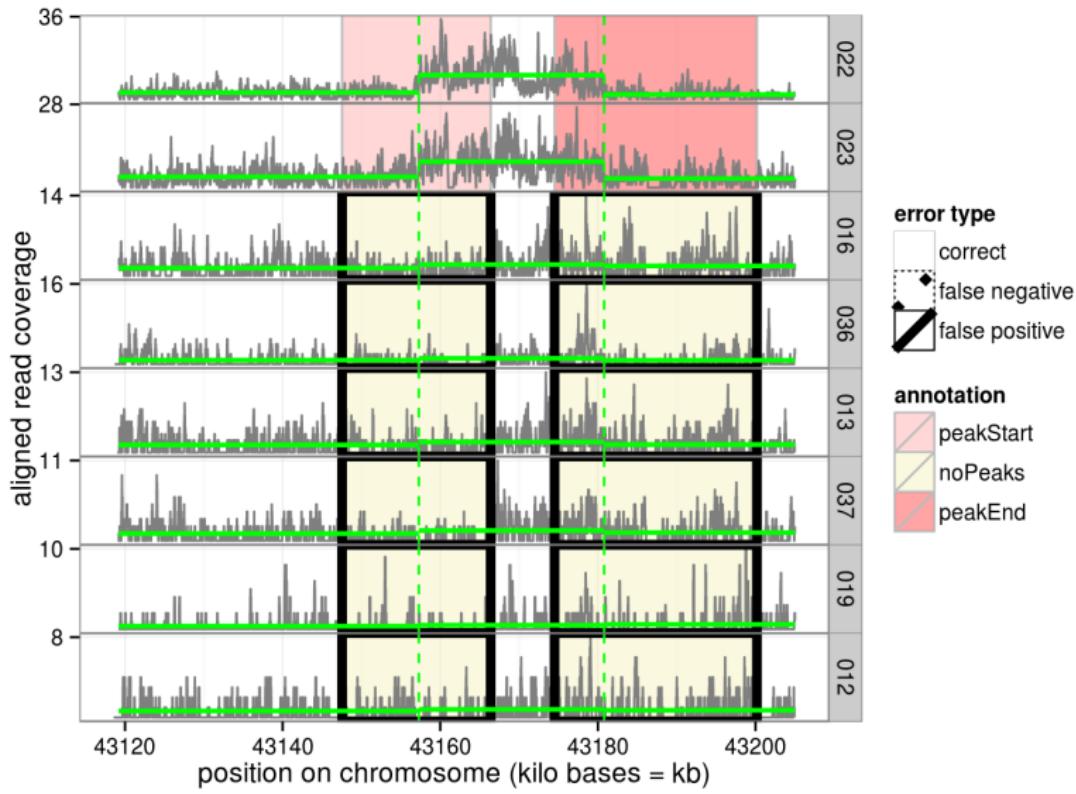
PeakSegJoint model with 6 peaks



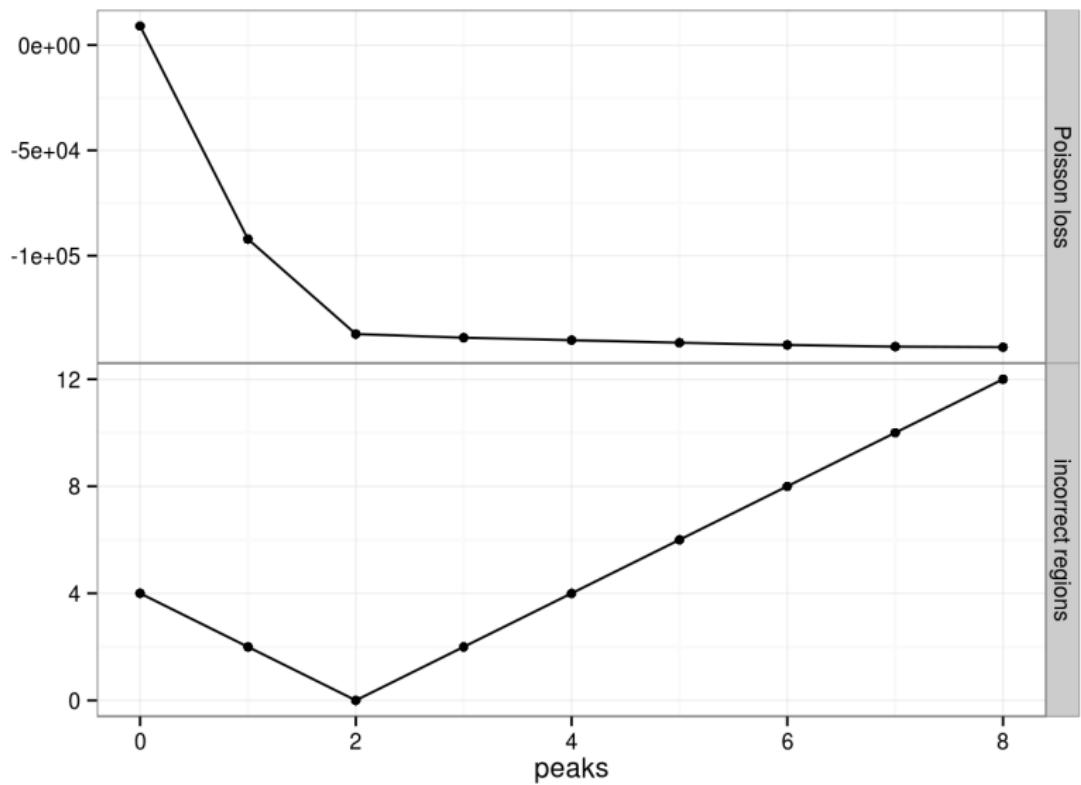
PeakSegJoint model with 7 peaks



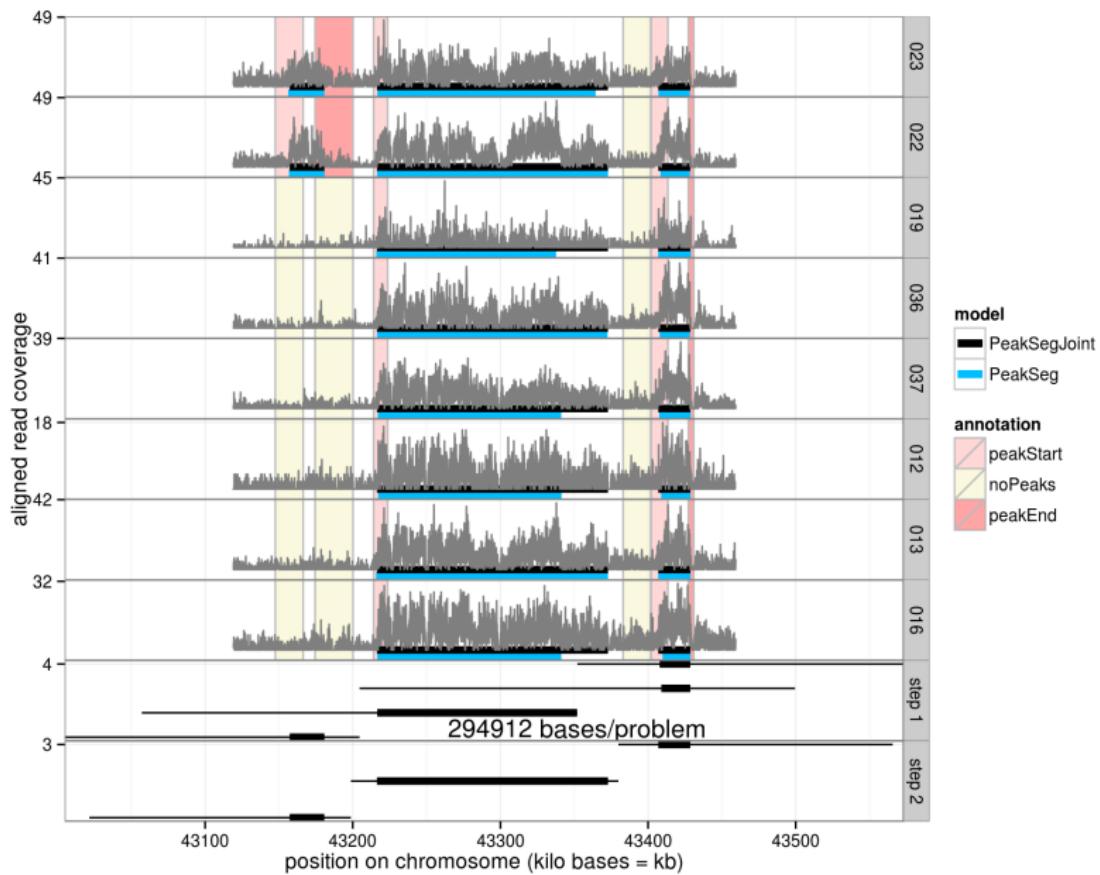
PeakSegJoint model with 8 peaks



Select model with minimal number of incorrect regions



H3K36me3 data, PeakSeg and Joint model



Timings on example H3K36me3 data

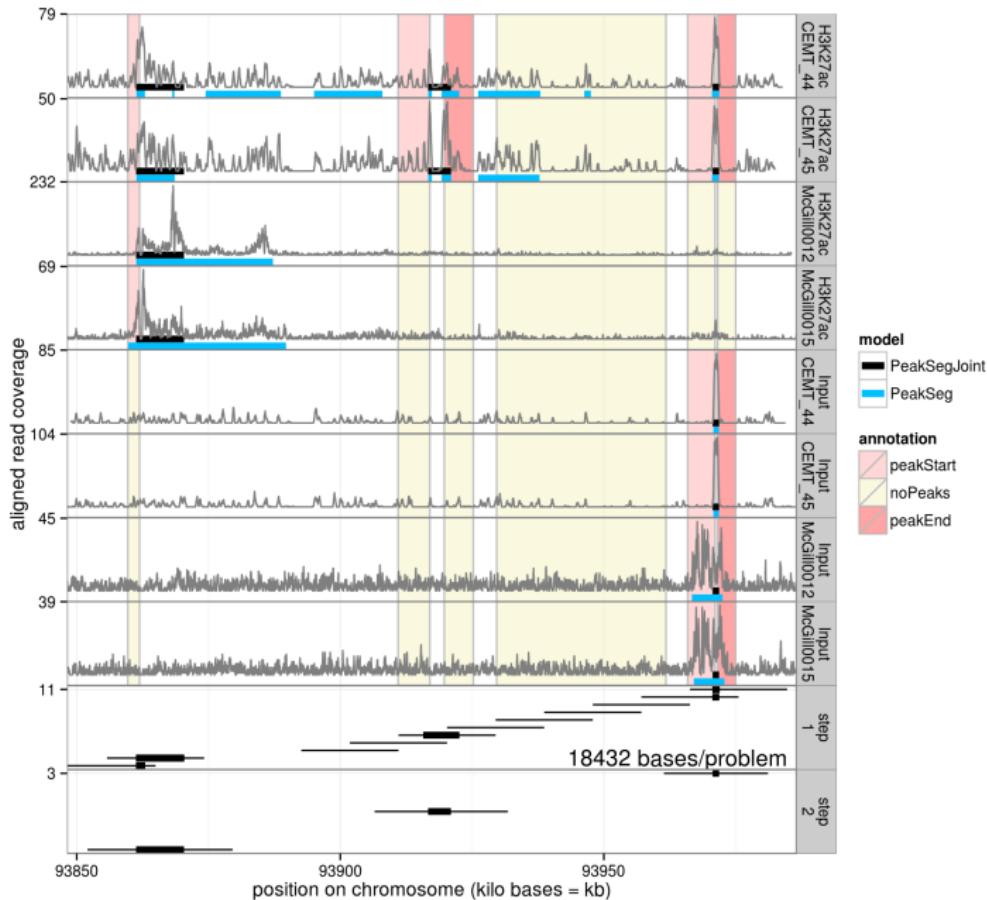
Find best 0,...,9 peaks in each of 8 samples (80 PeakSeg models):

| seconds | sample.id |
|---------|------------|
| 0.75 | McGill0023 |
| 0.77 | McGill0022 |
| 0.75 | McGill0019 |
| 0.79 | McGill0036 |
| 0.77 | McGill0037 |
| 0.77 | McGill0012 |
| 0.77 | McGill0013 |
| 0.76 | McGill0016 |
| 6.14 | total |

Find best common peak in 0,...,8 samples in each of 4 genomic regions (36 PeakSegJoint models):

| | data | seconds |
|-------------------------|--------|---------|
| chr21:42909696-43204608 | 23595 | 0.04 |
| chr21:43057152-43352064 | 162129 | 0.11 |
| chr21:43204608-43499520 | 206437 | 0.12 |
| chr21:43352064-43646976 | 67903 | 0.06 |
| total | | 0.33 |

H3K27ac and Input data, PeakSeg and Joint model



Timings on example H3K27ac data

Find best
0,...,9 peaks
in each of 8 samples
(80 PeakSeg models)

| seconds | sample.id |
|---------|--------------------|
| 0.99 | H3K27ac CEMT_44 |
| 0.96 | H3K27ac CEMT_45 |
| 1.00 | H3K27ac McGill0012 |
| 1.00 | H3K27ac McGill0015 |
| 0.99 | Input CEMT_44 |
| 1.00 | Input CEMT_45 |
| 1.01 | Input McGill0012 |
| 1.00 | Input McGill0015 |
| 7.94 | total |

Find best common peak
in 0,...,8 samples
in each of 11 genomic regions
(99 PeakSegJoint models)

| | data | seconds |
|-------------------------|-------|---------|
| chr11:93846528-93864960 | 7510 | 0.03 |
| chr11:93855744-93874176 | 11675 | 0.03 |
| chr11:93892608-93911040 | 5619 | 0.03 |
| chr11:93901824-93920256 | 6236 | 0.03 |
| chr11:93911040-93929472 | 5559 | 0.03 |
| chr11:93920256-93938688 | 5149 | 0.04 |
| chr11:93929472-93947904 | 4359 | 0.01 |
| chr11:93938688-93957120 | 3071 | 0.03 |
| chr11:93947904-93966336 | 3030 | 0.02 |
| chr11:93957120-93975552 | 7184 | 0.04 |
| chr11:93966336-93984768 | 7446 | 0.04 |
| total | | 0.32 |

ChIP-seq data and previous work on peak detection

The PeakSegJoint model

Conclusions

Thanks for your attention!

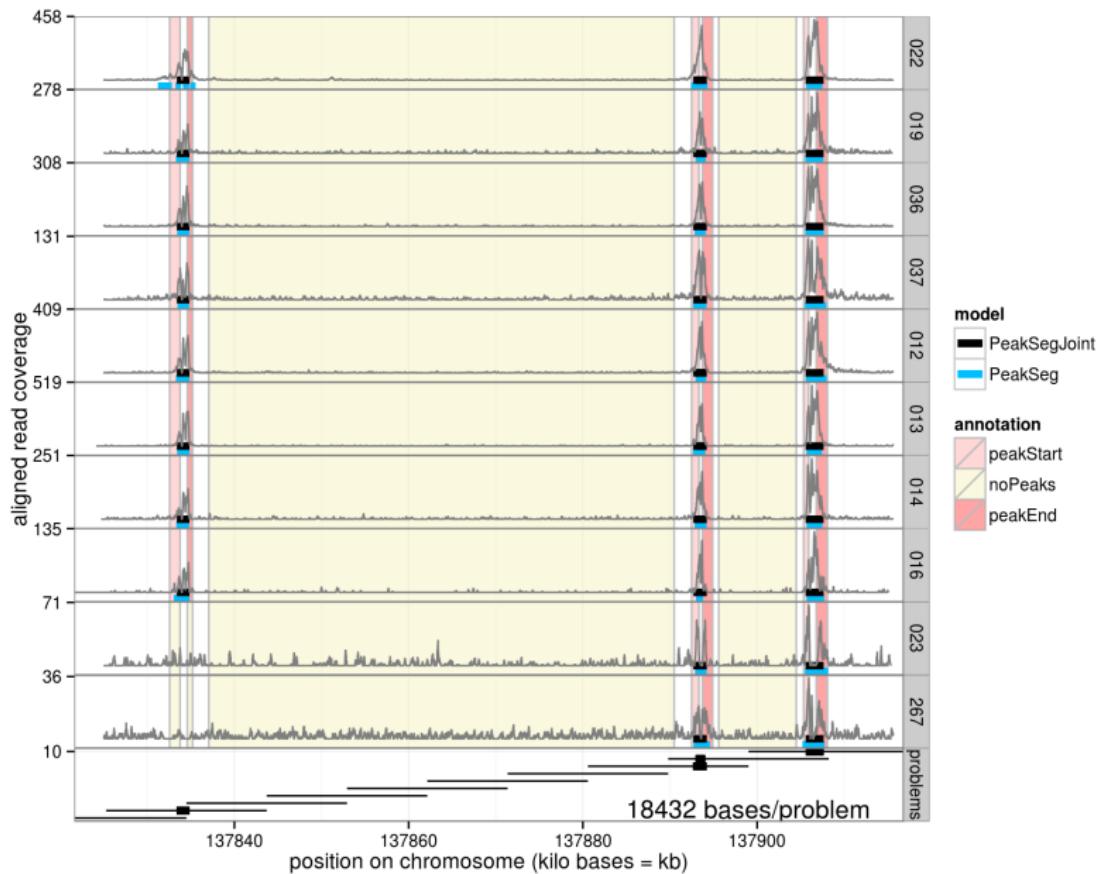
Write me at toby.hocking@mail.mcgill.ca to collaborate!

Source code for slides, figures, paper online!

<https://github.com/tdhock/PeakSegJoint-paper>

Supplementary slides appear after this one.

H3K4me3 data, PeakSeg and Joint model



Timings on example H3K4me3 data

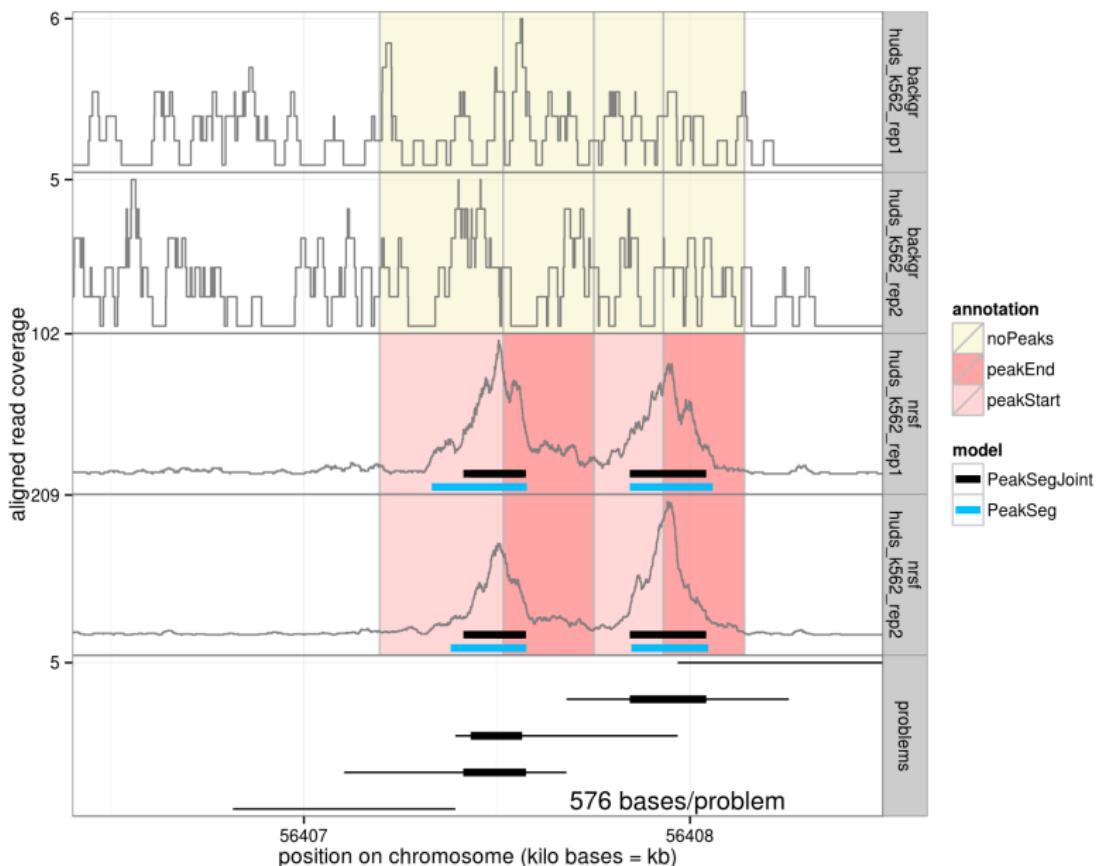
Find best
0,...,9 peaks
in each of 10 samples
(100 PeakSeg models)

| seconds | sample.id |
|---------|------------|
| 0.72 | McGill0022 |
| 0.71 | McGill0019 |
| 0.72 | McGill0036 |
| 0.72 | McGill0037 |
| 0.74 | McGill0012 |
| 0.76 | McGill0013 |
| 0.72 | McGill0014 |
| 0.72 | McGill0016 |
| 0.73 | McGill0023 |
| 0.75 | McGill0267 |
| 7.30 | total |

Find best common peak
in 0,...,10 samples
in each of 10 genomic regions
(110 PeakSegJoint models)

| data | seconds |
|-------|---------|
| 7603 | 0.01 |
| 12420 | 0.05 |
| 7023 | 0.01 |
| 3915 | 0.04 |
| 3597 | 0.03 |
| 3588 | 0.03 |
| 4255 | 0.05 |
| 13317 | 0.05 |
| 26436 | 0.05 |
| 19644 | 0.05 |
| total | 0.36 |

NRSF transcription factor data, PeakSeg and Joint model



Timings on example transcription factor data

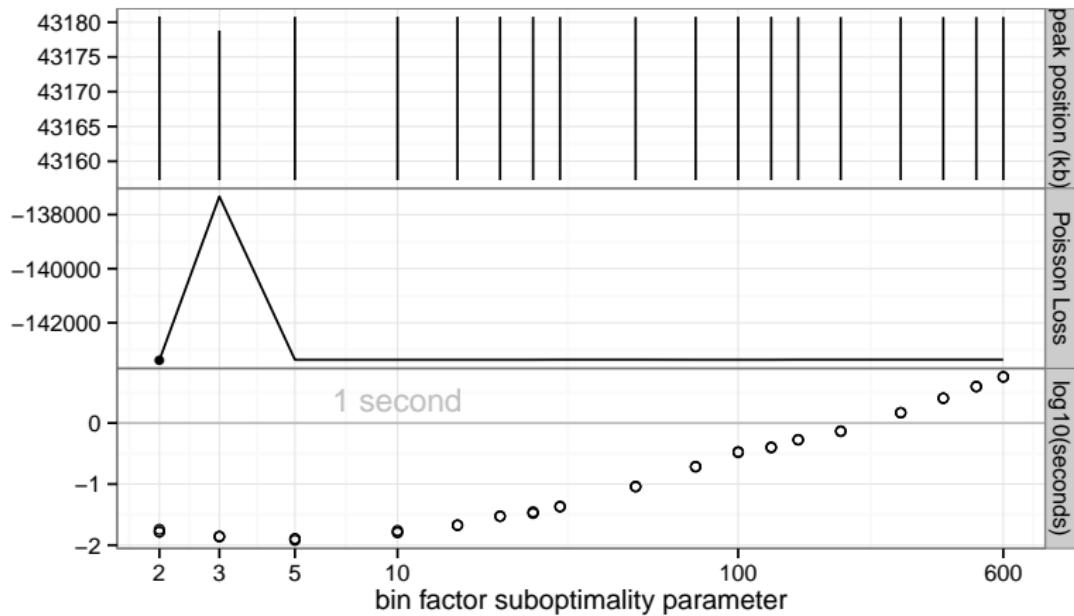
Find best
0,...,9 peaks
in each of 4 samples
(40 PeakSeg models)

| seconds | sample.id |
|---------|-----------------------|
| 0.26 | backgr huds_k562_rep1 |
| 0.24 | backgr huds_k562_rep2 |
| 0.30 | nrsf huds_k562_rep1 |
| 0.31 | nrsf huds_k562_rep2 |
| 1.10 | total |

Find best common peak
in 0,...,4 samples
in each of 5 genomic regions
(25 PeakSegJoint models)

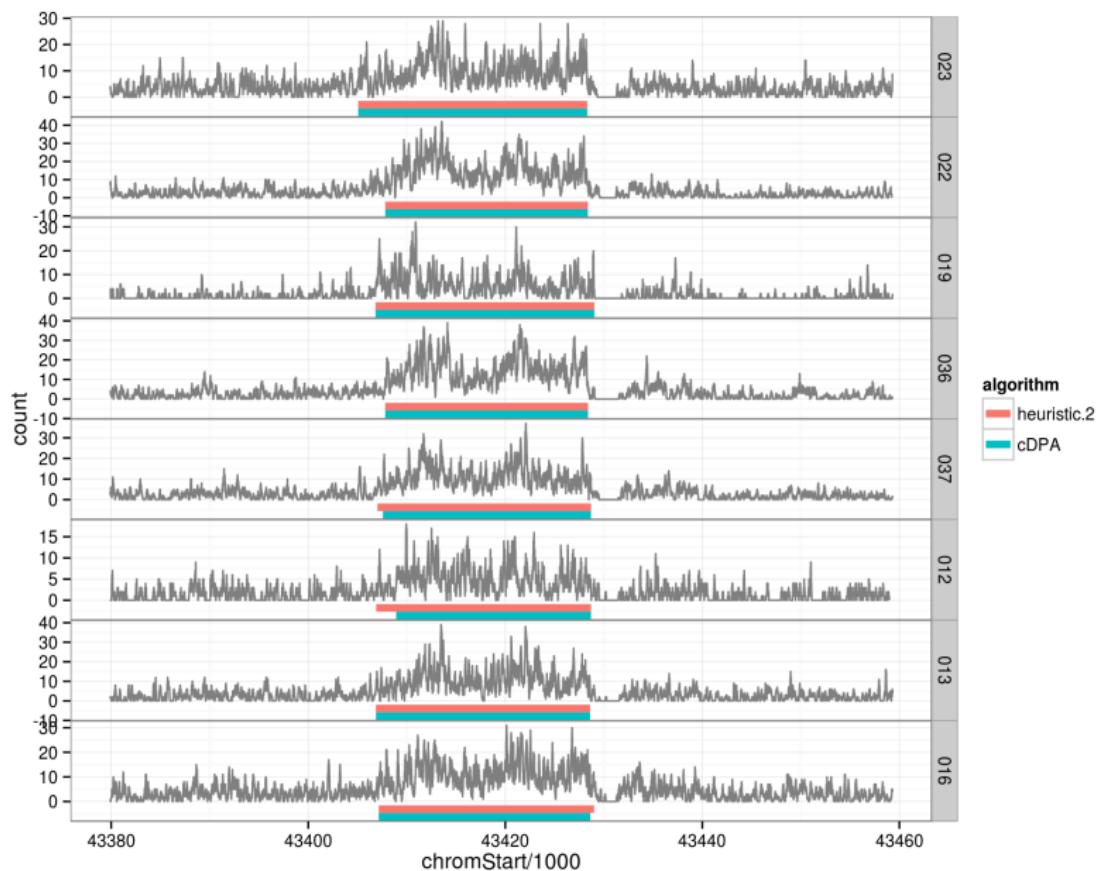
| | data | seconds |
|-------------------------|------|---------|
| chr21:56406816-56407392 | 345 | 0.01 |
| chr21:56407104-56407680 | 761 | 0.02 |
| chr21:56407392-56407968 | 975 | 0.01 |
| chr21:56407680-56408256 | 709 | 0.02 |
| chr21:56407968-56408544 | 298 | 0.01 |
| total | | 0.07 |

Bin factor parameter controls optimality and speed

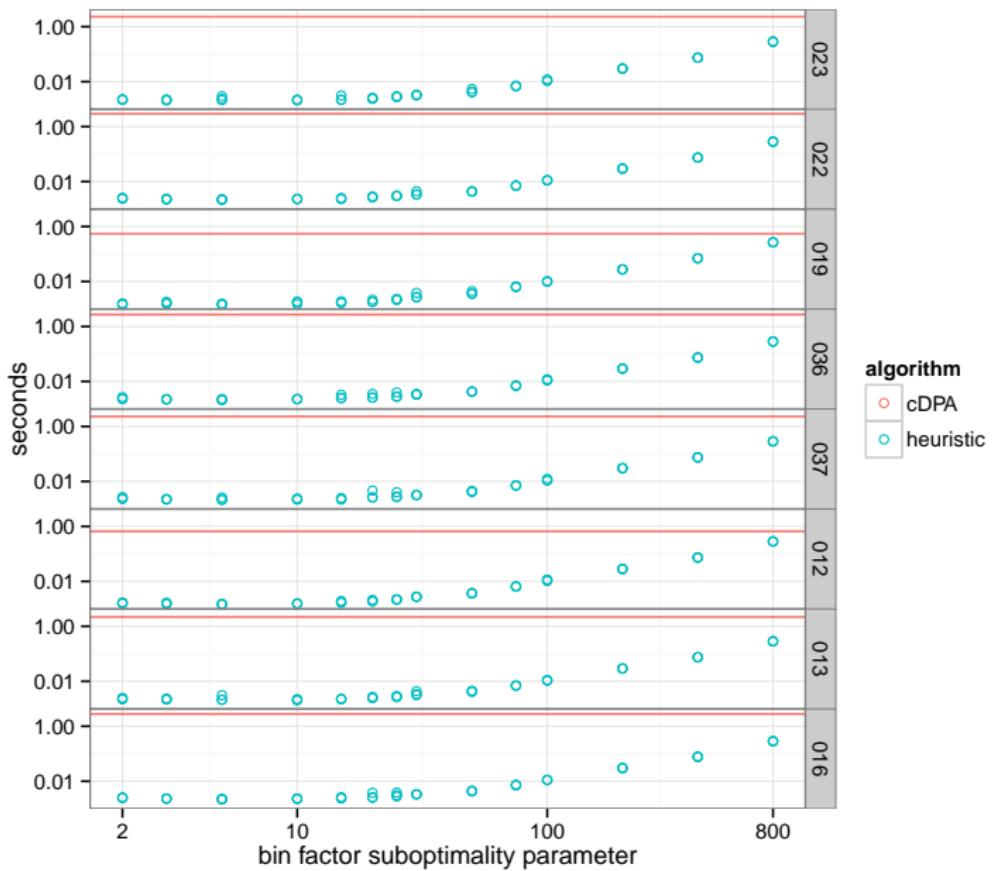


H3K36me3 example data set, PeakSegJoint model with 2 peaks.

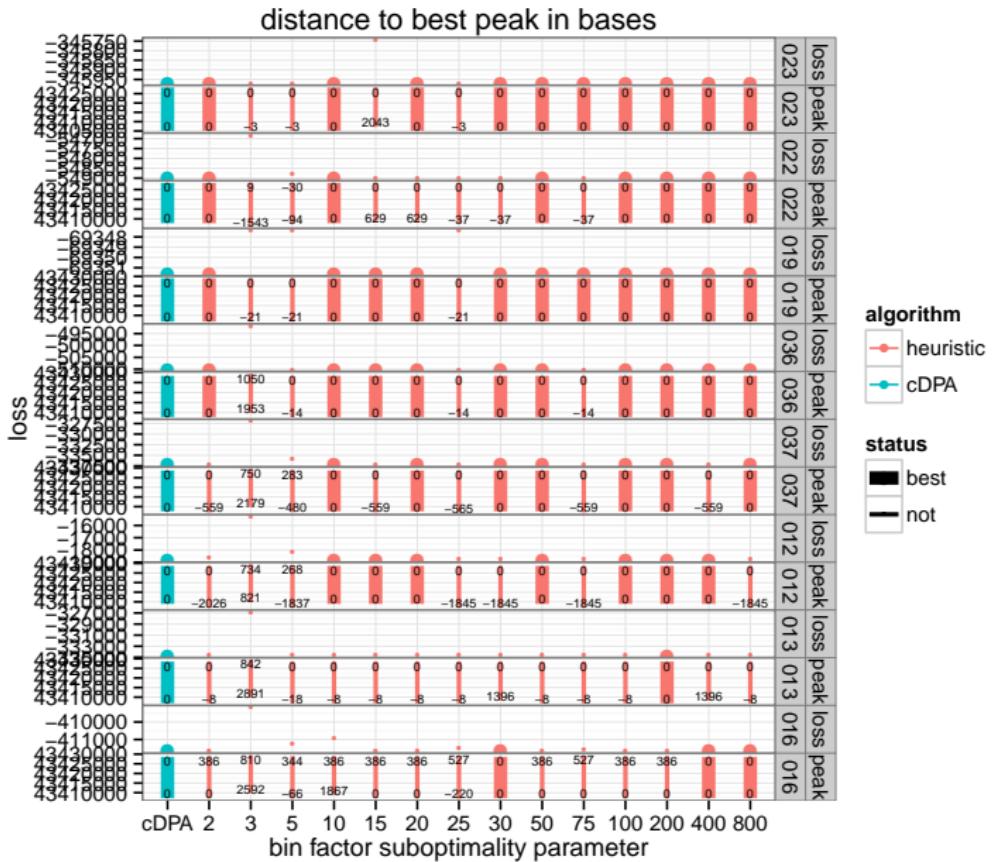
H3K36me3 data, cDPA and heuristic algorithms



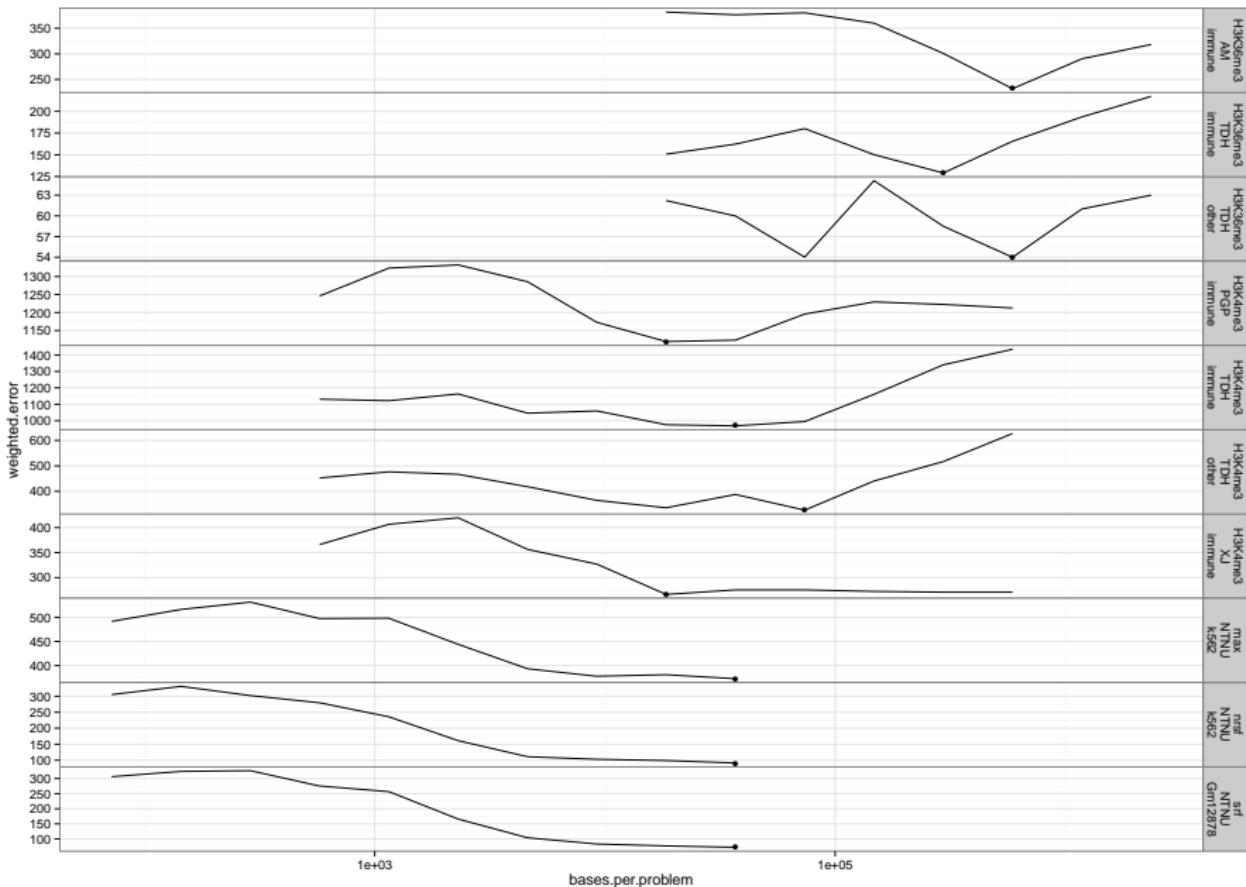
Heuristic is much faster than cDPA



Heuristic often as good as cDPA



Weighted train error not good for model selection



Select L1-regularized model with minimal validation error

