

The Data Integrator: a new way to combine data sources underlying the UCSC Genome Browser



Angie S. Hinrichs¹, Kate R. Rosenbloom¹, Matthew L. Speir¹, Donna Karolchik¹, Ann S. Zweig¹, David Haussler^{1,2}, Robert M. Kuhn¹, W. James Kent¹
¹University of California Santa Cruz Genomics Institute, Santa Cruz, CA; ²Howard Hughes Medical Institute, UCSC, Santa Cruz, CA

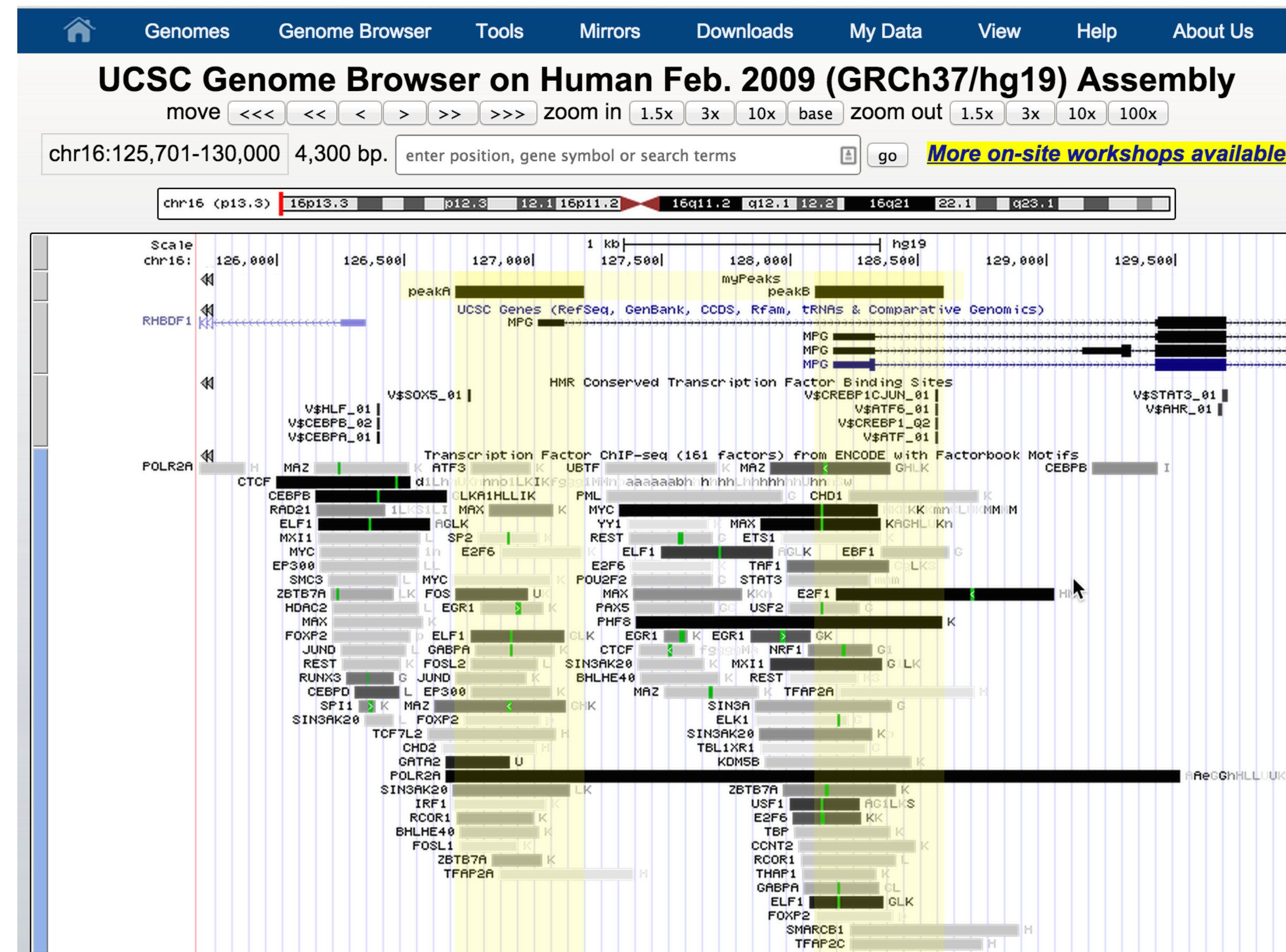
Introduction

Users of the UCSC Genome Browser have asked us variations of the same question many times over the years:

How can I get the $\left\{ \begin{matrix} \dots \text{score} \\ \text{gene name} \\ \text{TFBS} \end{matrix} \right\}$ for each of my $\left\{ \begin{matrix} \text{peaks} \\ \text{regions} \\ \text{probes} \end{matrix} \right\}$?

While the Table Browser tool can compute the intersections of two tracks by genomic position, it does so by reducing the second track to a set of genomic regions. It does not include any annotation columns of the second track in the output.

We have developed a new tool, the Data Integrator, that can combine data from up to five tracks based on overlap with items in the first track, and output all columns or a selected subset of columns. The output is tab-separated text that can be viewed in the web browser window or downloaded to a local file, optionally compressed by gzip. Like the Genome Browser and Table Browser, the Data Integrator can combine data from the browser database, user custom tracks and hub tracks.



Example: get a text file with the overlapping items shown here

Got variants?

The Data Integrator is a more flexible and open-ended tool complementary to the **Variant Annotation Integrator (VAI)**. If you are starting with a list of variant calls, and would like to get predicted functional effects on genes, then create a VCF or pgSnp-formatted custom track and use the VAI.

##fileformat=VCFv4.1

Data Integrator

```
track name=myPeaks
chr16 126700 127200 peakA
chr16 128100 128600 peakB
...
```

ct_myPeaks_9149.chrom	ct_myPeaks_9149.chromStart	ct_myPeaks_9149.chromEnd	ct_myPeaks_9149.name	knownGene.name
chr16	126700	127200	peakA	UC092CFA.4
chr16	126700	127200	peakA	VS50X5_01
chr16	126700	127200	peakA	VSAREB6_04
chr16	126700	127200	peakA	VS5P21_01
chr16	126700	127200	peakA	VS0LF1_01

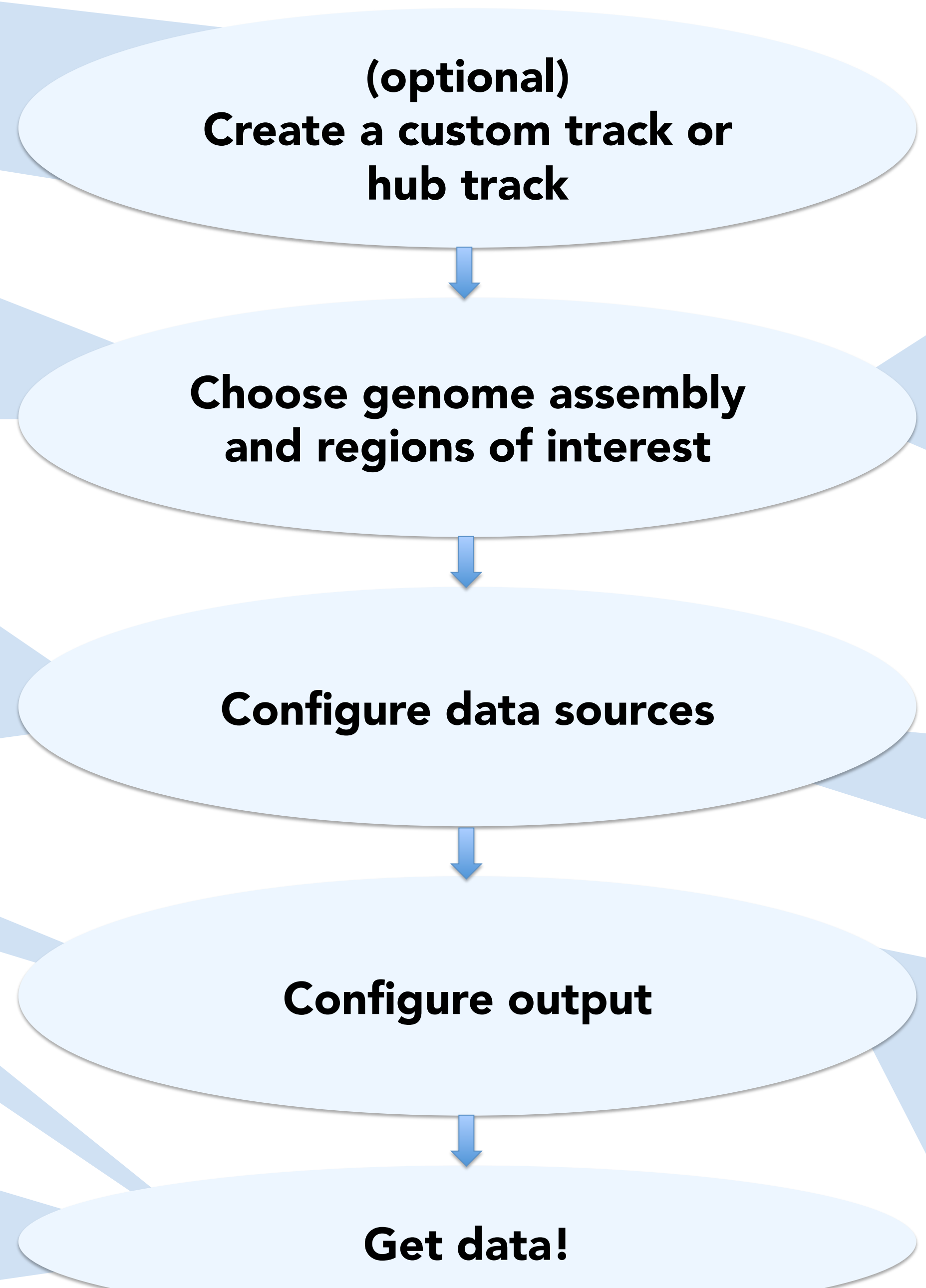
Future work

- Filter inputs by value
- Configure overlap rules
- Reorder output columns
- Add columns from related database tables

Reference:
 The UCSC Genome Browser database: 2015 update. Rosenbloom KR, Armstrong J, Barber GP, Casper J, Clawson H, Diekhans M, Dreszer TR, Fujita PA, Gurusudhan L, Haussler M, Harte RA, Heitner S, Hickey G, Hinrichs AS, Hubley R, Karolchik D, Learned K, Lee BT, Li CH, Miga KH, Nguyen N, Paten B, Raney BJ, Smit AF, Speir ML, Zweig AS, Haussler D, Kuhn RM, Kent WJ. *Nucleic Acids Res.* 2015 Jan;43(Database issue):D670-81.

Acknowledgements:
 This work was funded by National Human Genome Research Institute (5 U41 HG002371-15 to UCSC Center for Genomic Science). We would like to acknowledge the work of the UCSC Genome Bioinformatics technical staff (<http://genome.ucsc.edu/staff.html>), our many collaborators, and our users for their feedback and support.

How to build a query



Options

More information

- Search for answers in our mail list archives: <http://genome.ucsc.edu/contacts.html>
- Email a new question to our actively monitored list: genome@soe.ucsc.edu
- UCSC training (workshops, videos, tutorials): <http://genome.ucsc.edu/training/>
- Blog: <http://genome.ucsc.edu/blog/>



<http://genomewiki.ucsc.edu/index.php/BoG2015DataIntegratorPoster>