

### Abstract

The Harvard Chan Bioinformatics Core (HBC) provides best practice bioinformatics support and training. Researchers at Harvard Medical School (HMS) continue to make significant additions to biology using state of the art methods and novel experimental designs. The accompanying rapid rate of technological development requires complex analyses of large high throughput data sets and presents a challenge: expertise may not be readily available within experimental labs, making it difficult to determine and implement best practices to ensure accurate and reproducible results.

The HBC provides a single point of contact for researchers at HMS interested in using bioinformatics in their research. The HBC directly supports researchers data analysis (and training), with expertise in **study design, analysis and interpretation of next generation sequencing technologies** such as **RNA-seq, Single Cell RNA-seq, Variant sequencing, Bisulfite sequencing, ChIP-seq and ATAC-seq.** Grant support and support for **functional analyses by gene set enrichment or network mapping** are also available. Applying both established and developing methodologies in genomics, bioinformatics and biostatistics, the HBC handles projects of all sizes, from small expression studies to studies involving thousands of whole genomes, helping with the management, integration, and contextual analysis of high-throughput biological data. The HBC follows best practices and uses documented tools wherever possible, but can also adapt or develop new solutions if required. The HBC also provides bioinformatics training for HMS at multiple levels, from introductory workshops for technologies like RNA-seq and ChIP-seq to in depth training in programming at the command line. For more details about training see our training poster.

Through its work with HMS and the wider Harvard community, the HBC aims to provide a central institutional source of bioinformatics knowledge to help HMS researchers attain their research

### Expertise

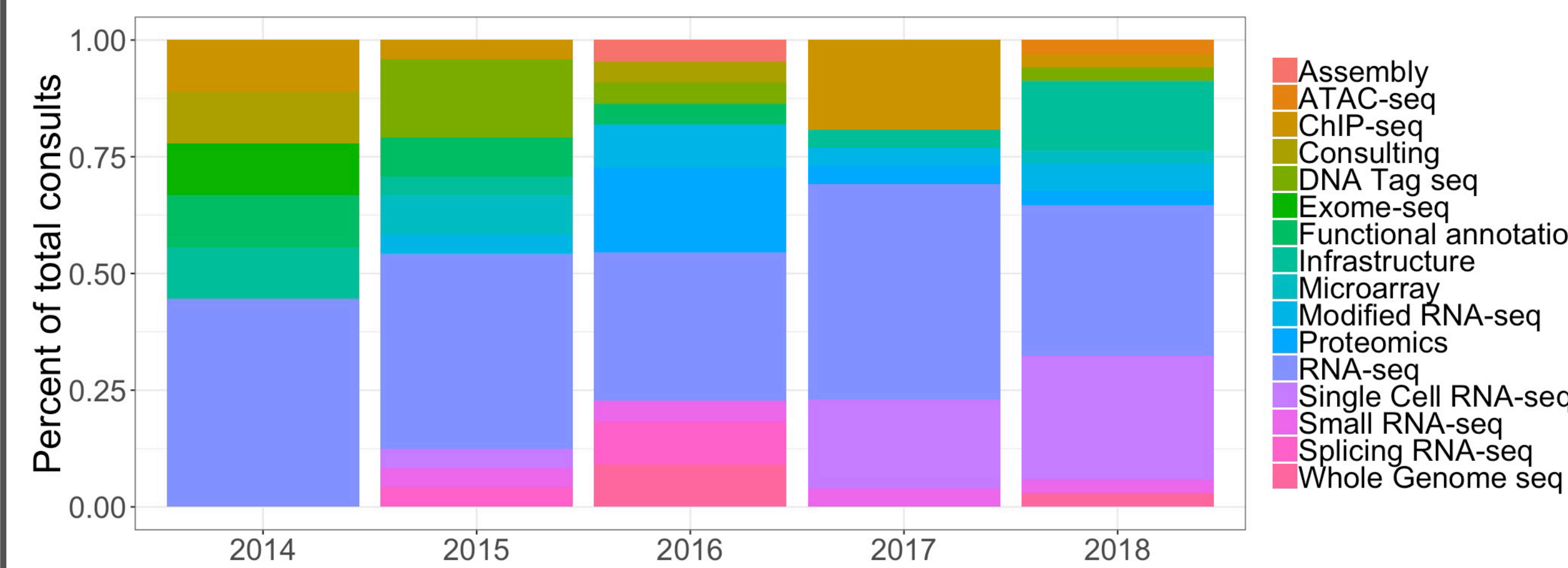
Appointment	Role	Skill
Professor	Analysis	Bisulfite-seq
Research Associate	Infrastructure	ChIP/ATAC-seq
Research Data Analyst	Leadership	Data Integration
Research Scientist	Training	Data Management
Senior Research Scientist		Functional Annotation
		RNA-seq
		Single Cell RNA-seq
		Small RNA-seq
		Splicing RNA-seq
		Whole Genome seq
		Variant-seq

### Projects

the HBC has worked with over 55 labs at HMS on more than 105 projects



- analysis of gene expression by RNA-seq remains the most common bioinformatics need
- we continue to see large increases in demand for Single Cell RNA-seq, including Nuc-seq
  - 33 InDrops consultations (13 experimental design, 20 analyses)
  - 10 10X Chromium consultations

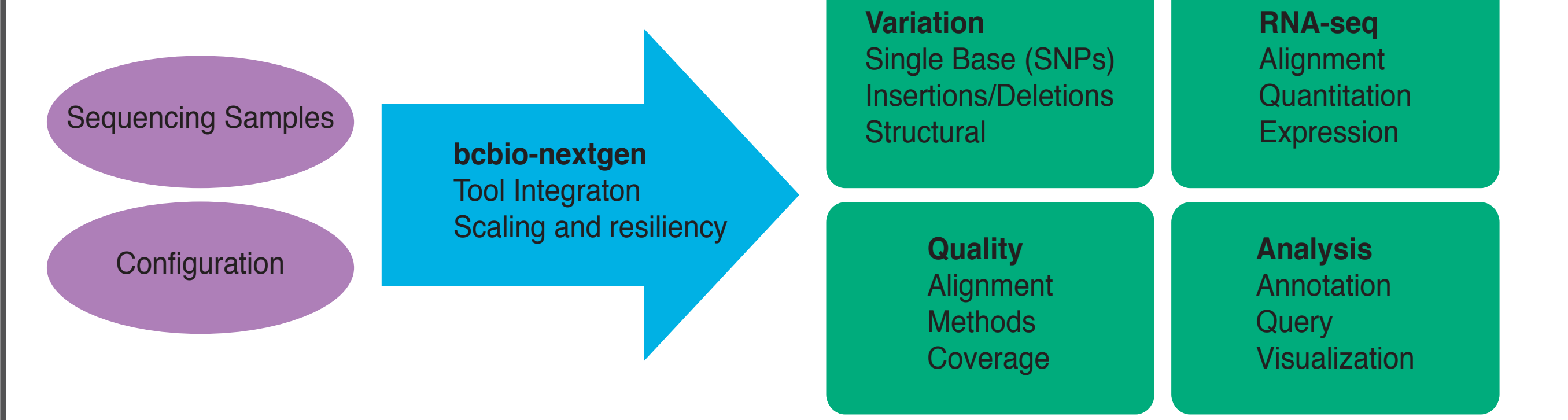


### Infrastructure

**Pipelines**

<https://bcbio-nextgen.readthedocs.io>

- Approach**
- Reproducible, Scalable, Automated, Documented, Self-contained, Interoperable
  - Open Source and Community Driven
- Functions**
- Variant Calling (exome, whole genome, structural, CNVs, cancer)
  - RNA-seq (bulk, single cell, small RNAs)
  - ChIP-seq



### R Packages

**DEGreport**

bioconductor version: Release (3.6)

Creation of a HTML report of differential expression analysis of count data. It integrates some of the code mentioned in DESeq2 and edgeR vignettes, and report a ranked list of genes according to the fold changes mean and variability for each selected gene.

Author: Lorena Pantano (aut, cre), John Hutchinson (ctb), Victor Barrera (ctb), Mary Piper (ctb), Kenneth Dalry (ctb), Thameez Malai Perumal (ctb), Rory Kirchner (ctb), Michael Steinbaugh (ctb)

Maintainer: Lorena Pantano <lorena.pantano at gmail.com>

Citation (from within R, enter `citation("DEGreport")`):

Pantano L. (2017). DEGreport: Report of DEG analysis. R package version 1.14.0.

---

**bcbioRNASeq: R package for bcbio RNA-seq analysis [version 1; referees: 1 approved with reservations]**

Michael J. Steinbaugh <sup>1</sup>, Lorena Pantano <sup>1\*</sup>, Rory D. Kirchner <sup>1</sup>, Victor Barrera <sup>1</sup>, Brad A. Chapman <sup>1</sup>, Mary E. Piper <sup>1</sup>, Meeta Mistry <sup>1</sup>, Radhika S. Khetani <sup>1</sup>, Kayleigh D. Rutherford <sup>1</sup>, Oliver Hofmann <sup>2</sup>, John N. Hutchinson <sup>2</sup>, Shannan Ho Sui <sup>1</sup>

<sup>1</sup>Harvard T.H. Chan School of Public Health, Boston, MA, 02115, USA  
<sup>2</sup>University of Melbourne, Centre for Cancer Research, Melbourne, VIC, 3000, Australia

\* Equal contributors

### Working With Us

- contact us early for help with experimental design
- schedule an initial meeting, free of charge
- we will create a timeline with deliverables
- we will provide a quote covering personnel, data storage and compute costs
- subsidized rates are available for HMS researchers on the quad**
- progress is regularly documented on a secure project site
- all data sets, results and documentation are shared

### Programs

```
bmtools,2.4.0
bcbio-nextgen,0.9.8a0-0183767
bcbio-variation,0.2.6
bcftools,1.3
bedtools,2.24.0
biobambam,2.0.42
bioconductor-bubbletree,2.1.5
bowtie2,2.2.8
bwa,0.7.13
chanjo,
cnvkit,0.7.11
cufflinks,2.2.1
cutadapt,1.9.1
fastqc,0.11.5
featurecounts,1.4.4
freebayes,1.0.2
gatk,3.2-2-gec30cee
gatk-framework,3.5.21
gemin,0.18.3
grabix,0.1.6
hisat2,2.0.3beta
htseq,0.6.1p1
lumpy-sv,0.2.12
manta,0.25.6
metasv,0.4.0
mutect,1.1.5
novaalign,3.04.04
novosort,03.00.02
oncofuse,1.1.0
phyloWgs,20150714
picard,1.141
platypus-variant,0.8.1
qualimap,2.1.3
rna-star,2.4.1d
rtg-tools,3.6
sailfish,0.9.0
salmon,0.6.0
sambamba,0.6.1
samblaster,0.1.22
samtools,1.3.1
scalpel,0.5.1
snpeff,4.2
vardict,2016.02.19
vardict-java,1.4.5
variant-effect-predictor,83
varscan,2.4.1
vcflib,1.0.0_rc0
vt,2015.11.10
wham,1.7.0.162
```

### Rmarkdown report with code

3.1 Over-representation analysis

- for differentially expressed genes and top fold change genes

3.1.1 gprofiler

- gprofiler will look for overrepresentation of a group of genes among multiple functional gene groups derived from databases including the Gene Ontologies, KEGG pathways, Reactome and others
- did a first pass with just the DE genes as defined by log2 fold change (1.5) and adjusted pvalue cutoff (0.2)
- pvalues for the the gprofiler results are all adjusted for multiple testing

```
> top.results.df.annot <- subset(results.df.annot, pad) <- qval.cutoff & abs(log2(foldchange) > 1.5 & cutoff)
> # number for life
> top.results.df.annot <- top.results.df.annot[order(absolute(top.results.df.annot$log2(foldchange)
+ decreasing = TRUE)), ]
>
> # run gprofiler with ordered query and background set of genes
> gprofiler.results <- gprofiler(query = as.vector(top.results.df.annot$term_id), organism
+ "musmusculi", ordered_query = TRUE, exclude_lea = F, correction_method = "gpc")
> #str:table(gprofiler.results, rownames = FALSE, caption = "gprofiler results for DE gene
+ ")
```

gprofiler results for DE genes

querynumber	significant	pvalue	termname	querysize	overlapsize	recall	precision	term.id	domain	subgraph.number	term.name	relative
1	TRUE	0.032	NSC	75	37	0.493	0.055	TS040868.1	df	1	factor C70, mod NSC/C70/C96, match class: 1	

- this first pass shows very few enrichments of any categories for the differentially expressed genes
- for the next pass I used the top 200 genes as determined by sorting by log2 fold change

```
> top.results.df.annot <- results.df.annot
> # number for genes with adjusted pvalue
> top.results.df.annot <- top.results.df.annot[order(absolute(top.results.df.annot$log2(foldchange)
+ decreasing = TRUE), 1:[1000]), ]
>
> # number for life
> top.results.df.annot <- top.results.df.annot[order(absolute(top.results.df.annot$log2(foldchange)
+ decreasing = TRUE), 1:[1000]), ]
>
> # run gprofiler with ordered query and background set of genes
> gprofiler.results <- gprofiler(query = as.vector(top.results.df.annot$term_id), organism
+ "musmusculi", ordered_query = TRUE, exclude_lea = F, correction_method = "gpc")
> #str:table(gprofiler.results, rownames = FALSE, caption = "gprofiler results for top 200
+ genes")
```

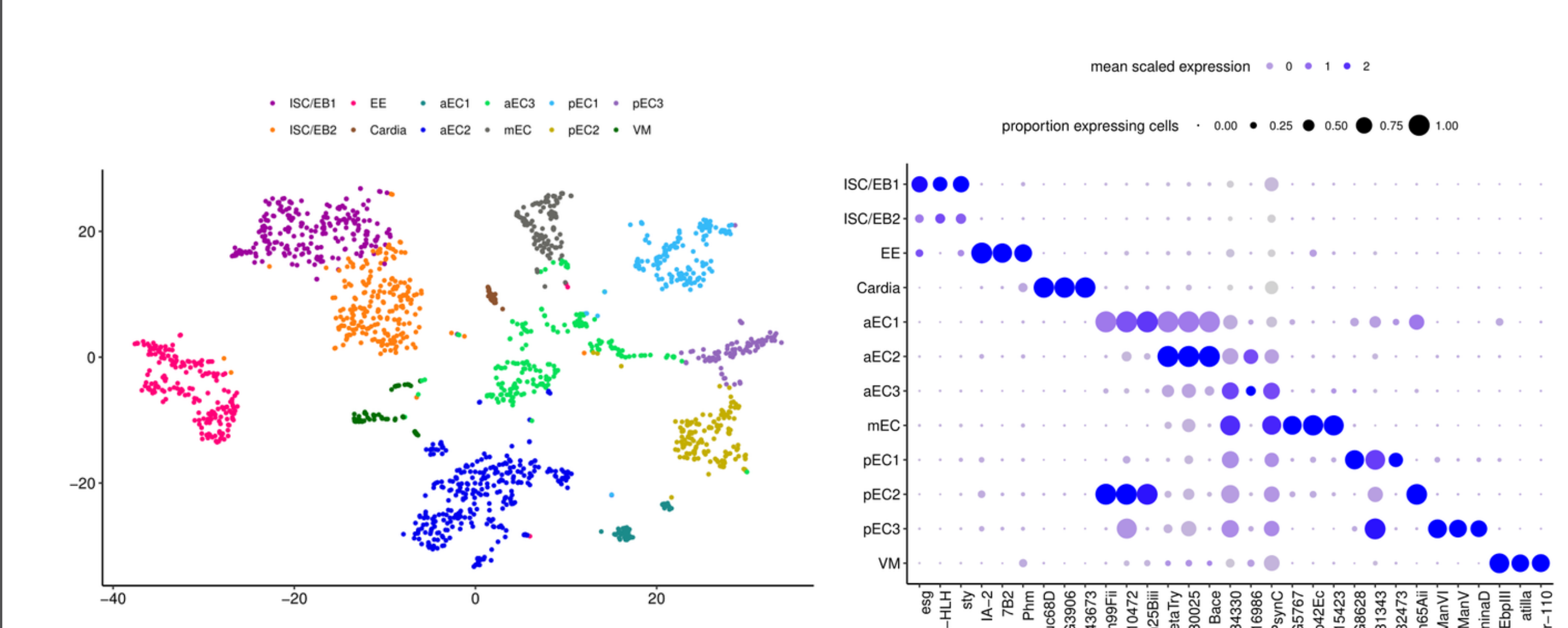
gprofiler results for top 200 genes

querynumber	significant	pvalue	termname	querysize	overlapsize	recall	precision	term.id	domain	subgraph.number	term.name	rela
1	TRUE	0.00279	5	169	3	0.018	0.000	GO:0041110	BP	15	intermediate filament bundle assembly	
1	TRUE	0.00529	46	155	5	0.032	0.109	GO:0019750	BP	14	apoptosis	
1	TRUE	0.01660	1	30	1	0.033	1.000	CORUM:1997	cor	2	Adiponectin homotrimer complex	
1	TRUE	0.03350	1	34	1	0.029	1.000	CORUM:549	cor	2	Sp1/3-DNA/12	

- we provide methods, publication quality figures and help with GEO submissions

### Example Results

Ruei-Jiun Hung, Yanhui Hu, **Rory Kirchner**, Fangge Li, Chiwei Xu, Aram Comjean, Sudhir Gopal Tattikota, Wei Roc Song, **Shannan Ho Sui**, Norbert Perrimon. (2018). A cell atlas of the adult *Drosophila* midgut. *bioRxiv* (doi: <https://doi.org/10.1101/410423>)



T-SNE visualization of 1753 adult fly intestinal gut cells: ISC/EB, intestinal stem cells/enteroblast; EE, enteroendocrine; aEC, anterior enterocyte; mEC, middle enterocyte; pEC, posterior enterocyte.

Expression levels and percentage of cells expressing the top 3 markers in each cluster as a dot plot

### Contact us:

We are located at the Harvard School of Public Health, SPH2, 2nd floor, Room 215.  
 Projects: [bioinformatics@hsph.harvard.edu](mailto:bioinformatics@hsph.harvard.edu)  
 Training: [hbctraining@hsph.harvard.edu](mailto:hbctraining@hsph.harvard.edu) (see our other poster for details!)  
 Website: [bioinformatics.hms.harvard.edu](http://bioinformatics.hms.harvard.edu)

