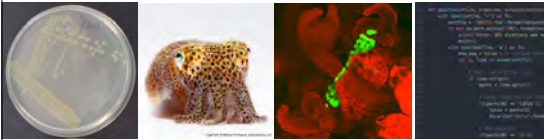
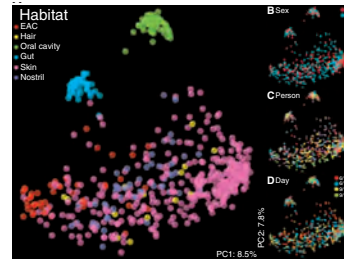


## Fancy genetics and simple scripts: Manipulating DNA data and becoming more proficient with Python



Mark Mandel  
@markjmandel  
Northwestern University Feinberg School of Medicine  
Department of Microbiology-Immunology

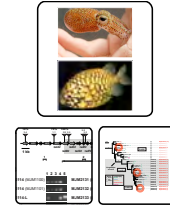
## Host specialization in humans



Costello 2009 Science

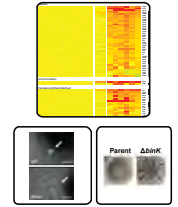
## Lab Interests

### Colonization Specificity



Comparative Genomics  
on Natural Isolates

### Symbiotic Development



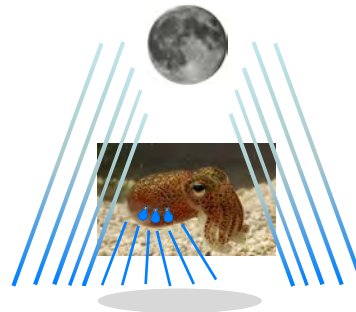
Functional Genomics  
on a Model Strain

## Study system



*Euprymna scolopes*  
Hawaiian bobtail squid

*Vibrio fischeri*



## *Euprymna scolopes*: Hawaiian bobtail squid



Adult



Hatchlings

light organ / ink sac

Adult image: Chris Frazee

## Stages of colonization

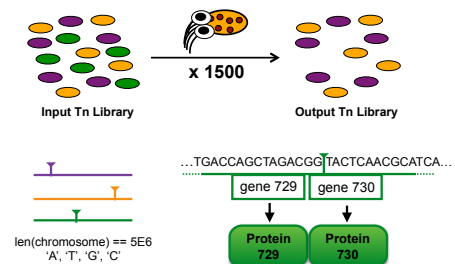
### What genes (and proteins) are required for bacterial colonization?



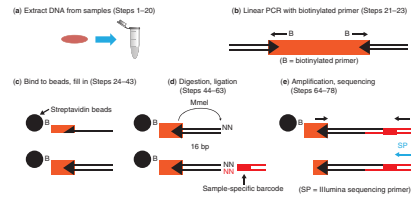
## Objectives

- Establish better laboratory practices for reproducibility of data analyses and for harnessing of large data sets.
- Expand the analyses that we can perform... Stop leaving (so much) unanalyzed data on the table...
- Example: Insertion Sequencing (INSeq)

## Screen for novel factors

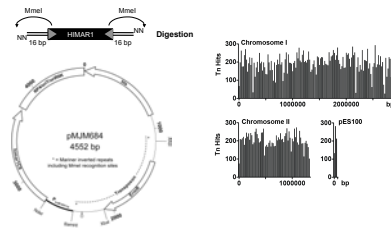


## Insertion Sequencing (INSeq)



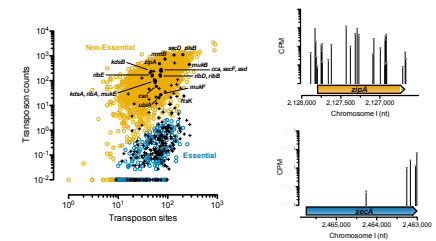
Goodman et al. 2011 Nature Protocols

## Tools for INSeq in *V. fischeri*



mBio 2014, PNAS 2014

## Essential genes in *V. fischeri*

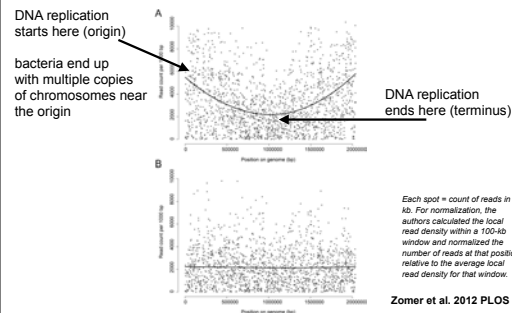


PNAS 2014

## What problems do I hope to address?

- Extend the current pipeline:
  - New transposon
  - Transposon with new functionality i.e., L != R
- Provide more robust data normalization
- Automate file format conversions
- Automate analyses
- Compare datasets robustly
- Examine down to individual genes

## LOESS Normalization



Zomer et al. 2012 PLOS One

## The pyinseq pipeline

- input files:
  - genome file = chromosomal reference sequences
  - reads file = experimental DNA sequences (190M reads)
  - sample file = DNA barcode key

## Genome file (.gbk)

```
LOCUS CP000020 2997536 bp DNA circular BCT 02-APR-2008
DEFINITION Vibrio fischeri RB114 chromosome I, complete sequence.
ACCESSION CP000020
[41 lines]
...
gene 19150..20790
/locus_tag="VF_0017"
CDS 19150..20790
/locus_tag="VF_0017"
[50443 lines]
...
BASE COUNT 882812 a 565193 c 563483 g 886048 t
ORIGIN
1 aagatcactt aatatatata agatctttaa aagagatctt tattatagtc tattatag
61 atotgtatgc tcgtgtgata agtgataaat gatcaatagg atcatatact ttatagtgat
121 coaaagtgtt tatcttttct tgatcttoga toggacagct tggagcaaaa agagttagtt
181 atocaaaggg gggggggggg ttgagctctt ttocegggat aactataact tpatctctgg
241 atctttctat agttatocaa atagtagtta toactotatt aataactttt atagatogga
[48294 lines]
```

## genome.fna

```
>CP000020
aagatcactt aatatatata aagatctt ttttaagagatctt ttattagatctt tattatag
atcgtgatct cctgtgataagtgataaatgat caataggatcatatact ttatagtgat
coaaagtgtt tatcttcttcttgatctcgtacggacagct tggagcaaaa agagttagtt
atccacaagg gggggggggcgttagatctt tcaatggataactaaacttgatcactgg
atctttctat agttatccacatagtaggtatcatctatttaaaacttttatagatcgga
caaacacttt tattaacaaatgtgtgttttagccaaactctcgtggtcttcagggat
actattttgagt tacatctatctcctaataagaggtacttttcgatagcgcaatgctcttt
taaggtattttgaagtgtttctcgtcgcgagaaagtgatcataacttgaatcaccgat
agcaacaacag caaattcaacgctcgataatggctgagtcacgcttctcaactgttgaat
aaatggtttaatgtgtcagggtatcaccacgacagcgtgagttgatcaacagatgaacca
taagctatca atatcaatatcatctaaattgggtgtattgaatatcggtggaataatc
catttctcttaataaactcagcaagatggtcaccacactactcagacaccactagagtgct
tctgtataatagatcacttttctcatgaattatctataaaaataaaaaatgggcc
tacaatggccattattaatcttatttaatttggtttatttaccacacagaatgaagt
aaatacgtcccaagtaaatcatcagaggttaattcaccggtgatttcaacttaattggt
```

# genome.ftt

| LOCUS   | CFP000020 |                         |            |            |         |         |     |           |  |
|---|-----------|-------------------------|------------|------------|---------|---------|-----|-----------|--|
| Location  | Strand    | Length                  | PID        | Gene       | Synonym | Code    | COD | Product   |  |
| 313..747  | -         | 435                     | AAM84499.1 | miOC       | VF_0001 | -       | -   | FMR-      |  |
| Binding protein MIOc                            |           |                         |            |            |         |         |     |           |  |
| E17..2184                                       |           | 1368                    | AAM84497.1 | TrmE       | VF_0002 | -       | -   | tRNA      |  |
| modification GTPase TrmE                        |           |                         |            |            |         |         |     |           |  |
| 2281..3906                                      |           | 1626                    | AAM84498.1 | yigC       | VF_0003 | -       | -   |           |  |
| cytoplasmic insertase into membrane protein, de |           |                         |            |            |         |         |     |           |  |
| 4133..4486                                      |           | 354                     | AAM84499.1 | snaB       | VF_0004 | -       | -   |           |  |
| ribonuclease P protein component                |           |                         |            |            |         |         |     |           |  |
| 4502..4636                                      |           | 135                     | AAM84500.1 | rpsH       | VF_0005 | -       | -   | SOS       |  |
| ribosomal protein L38                           |           |                         |            |            |         |         |     |           |  |
| 4822..5559                                      |           | 738                     | AAM84501.1 | yecC       | VF_0006 | -       | -   | cystine   |  |
| transport ATP-binding protein                   |           |                         |            |            |         |         |     |           |  |
| 5556..6227                                      | -         | 672                     | AAM84502.1 | yecS       | VF_0007 | -       | -   | cystine   |  |
| transport system peremease protein              |           |                         |            |            |         |         |     |           |  |
| 6352..7104                                      |           | 752                     | AAM84503.1 | fljY       | VF_0008 | -       | -   | cysteine- |  |
| binding protein                                 |           |                         |            |            |         |         |     |           |  |
| 7401..8810                                      |           | 1410                    | AAM84504.1 | dnaA       | VF_0009 | -       | -   |           |  |
| chromosomal replication initiator protein DnaA  |           |                         |            |            |         |         |     |           |  |
| 8843..9943                                      |           | 1101                    | AAM84505.1 | dnaN       | VF_0010 | -       | -   | DNA       |  |
| polymerase III, beta subunit                    |           |                         |            |            |         |         |     |           |  |
| 9957..110136                                    |           | +<br>gap repair protein | 1080       | AAM84506.1 | recF    | VF_0011 | -   | -         |  |

# Reads file (.fastq)

```
@DGL9ZZQ1:720:C6YD0ACXX:2:1101:1246:2185 1:N:0:
GAGACACATTAACCTCAATTACAGGTTGGATGATAAGTCCCCGGTCTTCG
+
CCCCFFFDHHHHHCGHHHIJDHIJJFHJJJJJJJJJDHIJJHIAIGIJJJ
@DGL9ZZQ1:720:C6YD0ACXX:2:1101:1246:2185 1:N:0:
CTTTCGACACCGCAGCAGGTAACAGGTTGGATGATAAGTCCCCGGTCTTCG
+
CCCCFFFDHHHHHCGHHHIJDHIJJFHJJJJJJJJJDHIJJHIAIGIJJJ
```

4 lines x 109,000,000 reads (760,000,000 lines in file)

```
@identifier
barcodechromosomestransposon
+
quality scores
```

| Sample file (.txt) |      |
|--------------------|------|
| E001_01            | GAAG |
| E001_02            | CTTT |

# The Results...

```
$ python2.7 pyinseq.py -i reads.fastq -s samples.txt
-g genome.gb -e exp01

$ python2.7 pyinseq.py -i _exampleData/
example01.fastq -s _exampleData/example01.txt -g
_exampleData/ES114v2.gb -e example01
```

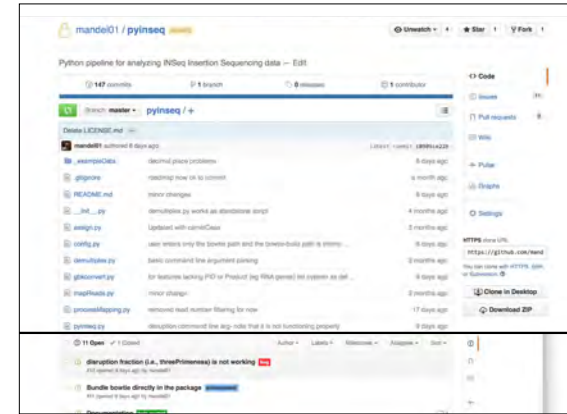
| Contig   | Start | End   | Strand | Length | PID        | Gene  | Synonym  | Code | CDG | Product                    | example01.CTTT | example01.GAAG |
|----------|-------|-------|--------|--------|------------|-------|----------|------|-----|----------------------------|----------------|----------------|
| CP000001 | 13554 | 13471 | +      | 2418   | AAW84607.1 | gfp9  | VF_0012  | -    | -   | DNA gyrase, subunit        | 100000         | 100000         |
| CP000001 | 89622 | 62081 | -      | 680    | AAW84616.1 | gfp8  | VF_0019  | -    | -   | thiamine synthase          | 100000         | 50000          |
| CP000001 | 56145 | 56748 | -      | 1108   | AAW84762.1 | gmetK | VF_00482 | -    | -   | phenylglyoxylate isomerase | 200000         | 300000         |
| CP000002 | 3775  | 4208  | -      | 498    | AAW85245.1 | -     | VF_00007 | -    | -   | DNA repair protein         | 300000         | 300000         |

Real results would have ~ 4000 rows; one for each gene in the genome.

# What have I learned in the past 6 months?

- Python :)
- Test data set (and data sets < 50,000 reads)
- generator to read in large data set
- argparse
- basic logging into the terminal
- combining results into an informative table
  - Previous Perl script would output 3 files/  
sample x 16 samples, each of which would be run with  
a separate shell command, then output separate lists.
- Git/Github

```
with open(infile, 'r') as fi:  
    for line in fi:  
        ...
```



# Next Steps

- scaling up to work for larger files (don't pass through data so many times!)
- string formatting so it works in Python3
- modularize... then optimize and expand.

```
graph LR; 1[1. Prepare Reads] --> 2[2. Basic INSeq mapping]; 2 --> 3[3. Analysis];
```

- 1**  
**Prepare Reads:**  
Demultiplex by barcode; separate into folders
- 2**  
**Basic INSeq mapping:**  
Map to genome, normalization(s), map to gene, output tables.
- 3**  
**Analysis:**  
Statistical comparisons, plotting, quality checking  
*Connect to other DBs*

# Some of the many things I don't know (but will soon!)

- buffered reading & writing
- numpy
- pandas and dataframes
- plotting with matplotlib
- unit testing
- classes (oh, the shame...)
- And much, much more

# Acknowledgments



John Brooks  
Mattias Gyllborg  
David Cronin



Andrew Goodman (Yale)  
Ben Johnson (MSU)



