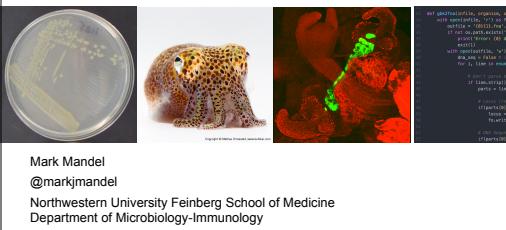
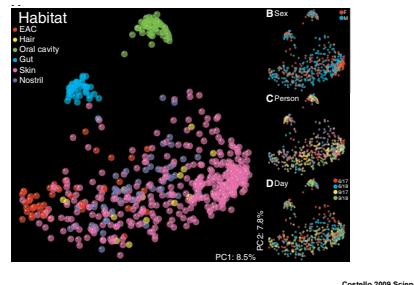


Fancy genetics and simple scripts: Manipulating DNA data and becoming more proficient with Python



Host specialization in humans



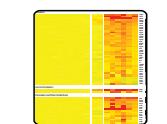
Lab Interests

Colonization Specificity



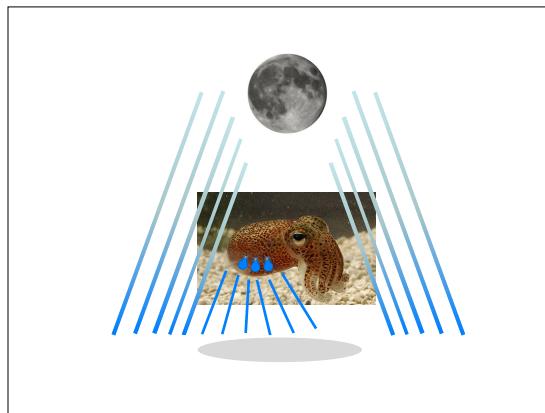
Comparative Genomics
on Natural Isolates

Symbiotic Development

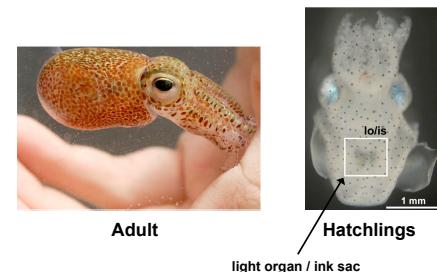


Functional Genomics
on a Model Strain

Study system



Euprymna scolopes: Hawaiian bobtail squid



Stages of colonization

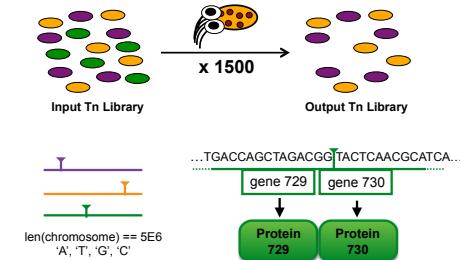
What genes (and proteins) are required for bacterial colonization?



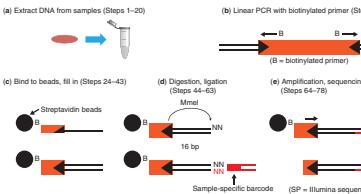
Objectives

- 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)

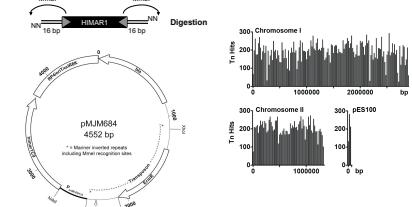


Goodman et al. 2011 Nature Protocols

What problems do I hope to address?

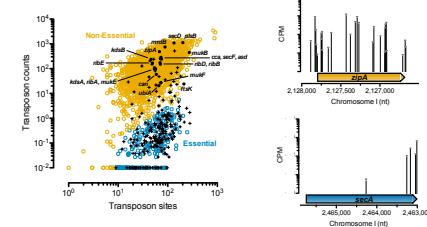
- Extend the current pipeline:
 - New transposon
 - Transposon with new functionality i.e., $L \neq R$
 - Provide more robust data normalization
 - Automate file format conversions
 - Automate analyses
 - Compare datasets robustly
 - Examine down to individual genes

Tools for INSeq in *V. fischeri*



—Blz. 2214, PNAC 2214

Essential genes in *V. fischeri*



DNAG 2014

Genome file (.gbk)

```

LOCUS   CP000020          2897536 pp   DNA      circular BCT 02-APR-2008
DEFINITION Vibrio fischeri ES114 chromosome I, complete sequence.
ACCESSION CP000020
[41 lines]
...
  gene           19150..20790
    /gene="malS"
    /locus_tag="VF_00017"
CDS     19150..20790
    /gene="malS"
    /locus_tag="VF_00017"
[50443 lines]
...
BASE COUNT  882812 a 565193 c 563483 g 886048 t
ORIGIN
  1 aatataatataatataatataatccgtttt aacggatctt ttatagatc tttatatag
  61 attcgtatcgtatcggtatcgtatcaatcgtatccatgtt atccatatact ttataggatg
  121 ccgtatccatgtt tgcattttttt tgatccatgttccgtatccatgtt ccgtatccatgtt
  181 atccatccaaaggccggggaggccgtt tgatccatgtt ccgtatccatgtt aactatcaatc
  241 atccatccatgttccggggaggccgtt tgatccatgtt ccgtatccatgtt aactatcaatc
  301 tgcatttttttccgtatccatgtt ccgtatccatgtt aactatcaatc tttataggatg
  361 atccatccatgttccggggaggccgtt tgatccatgtt ccgtatccatgtt aactatcaatc
  421 atccatccatgttccggggaggccgtt tgatccatgtt ccgtatccatgtt aactatcaatc
[48294 lines]

```

```
gbkconvert.py

32
33 def gbk2fna(infile, organism, outputdirectory=''):
34     with open(infile, 'r') as fi:
35         outfile = '{0}({1}).fna'.format(outputdirectory, organism)
36         if not os.path.exists('{0}'.format(outfile)):
37             printError('{0} directory was not created.'.format(outputdirectory))
38         else:
39             exit(1)
40     with open(outfile, 'w') as fo:
41         dna_seq = False
42         for i, line in enumerate(fi):
43             # Don't parse blank lines
44             if line.strip():
45                 parts = line.split()
46
47                 # Locus (replicon) as header
48                 if(parts[0] == 'LOCUS'):
49                     locus = parts[1]
50                     fo.write('>{0}\n'.format(locus))
51
52                     # DNA Sequence
53                     if(parts[0] == '/'):
54                         dna_seq = False
55                     if dna_seq:
56                         sequence = ''.join(n for n in line.strip() if n.isalpha())
57                         fo.write('{0}\n'.format(sequence))
58                     if(parts[0] == 'ORIGIN'):
59                         dna_seq = True
60
61     fo.close()
62
63     print('File {0} has been created in {1}'.format(outfile, outputdirectory))

abacmeyer 118-89
```

genome.fna

| genome.fasta | | | | | | | | | |
|--------------|----------|--------|------------|------|---------|---------|------|-----|--|
| LOCUS | CP000020 | Strand | Length | PID | Gene | Synonym | Code | COG | Product |
| 313..747 | - | 435 | AAN84496.1 | mioC | VF_0001 | - | - | - | FMN- |
| 817..2184 | - | 1368 | AAN84497.1 | tRNA | VF_0002 | - | - | - | tRNA |
| 2281..3906 | - | 1626 | AAN84498.1 | yidC | VF_0003 | - | - | - | cysteine |
| 4133..4486 | - | 354 | AAN84499.1 | rnpA | VF_0004 | - | - | - | ribonucleic P protein component |
| 4502..4636 | - | 135 | AAN84500.1 | rpmB | VF_0005 | - | - | - | 50S ribosomal protein L34 |
| 4822..5585 | - | 763 | AAN84501.1 | yeoC | VF_0006 | - | - | - | cysteine |
| 5556..6227 | - | 672 | AAN84502.1 | yeoB | VF_0007 | - | - | - | transport system permease protein |
| 6352..7104 | - | 753 | AAN84503.1 | fliY | VF_0008 | - | - | - | cysteine |
| 7401..8884 | - | 1480 | AAN84504.1 | dnaA | VF_0009 | - | - | - | chromosomal replication initiator protein DnaA |
| 8843..9943 | + | 1101 | AAN84505.1 | dnaN | VF_0010 | - | - | - | DNA polymerase III, beta subunit |
| 9957..11036 | + | 1080 | AAN84506.1 | recF | VF_0011 | - | - | - | gap repair protein |

Reads file (.fastq)

```
@DGL9ZZQ1:720:C6YD0ACXX:2:1101:1246:2185 1:N:0:
GAAGCACATTAAACCTCACTTACACAGGTGGATGATAAGTCCCCGGTCTCG
+
CCCCFFFDHHHHHCGHHHIJDHIIJJFHHJIIJJJIIJDHIIJJHIA吉JJ
@DGL9ZZQ1:720:C6YD0ACXX:2:1101:1246:2185 1:N:0:
CTTCGACACCGAACACCGTAACACGGTGGATGATAAGTCCCCGGTCTCG
+
CCCCFFFDHHHHHCGHHHIJDHIIJJFHHJIIJJJIIJDHIIJJHIA吉JJ

4 lines x 190,000,000 reads (760,000,000 lines in file)

@Identifier
barcodechromosometransposon
+
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

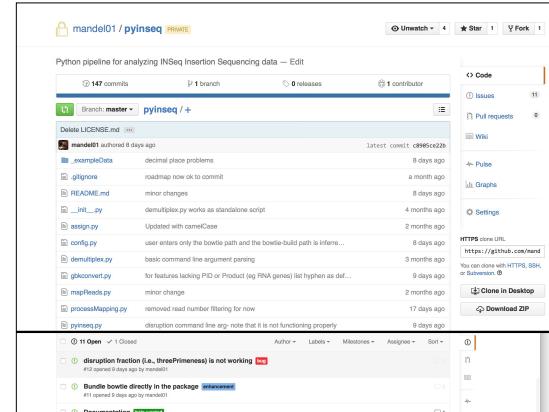
$ python2.7 pyinseq.py -i _exampleData/example01.fastq -s _exampleData/example01.txt -g _exampleData/ES114v2.gb -e example01

Config Start End Strand Length PID Gene Synonym Code COG Product example01:CTT example01:GAAG
CP000020 11024 13471 + 2418 AAWB8249.1 VF_0001 - DNA repair protel 100000 100000
CP000020 600 600 + 600 AAWB8249.1 VF_0002 - DNA repair protel 100000 100000
CP000021 56045 56724 + 1504 AAWB8250.1 VF_0003 - amyo-N-acyl 250000 300000
CP000022 3771 4268 - 498 AAWB8249.1 VF_0007 - DNA repair protel 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
 - with open(infile, 'r') as f:
for line in f:
 - ...
- 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



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.

| | | |
|--|---|---|
| 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

