# Assignment 8: Metagenomics
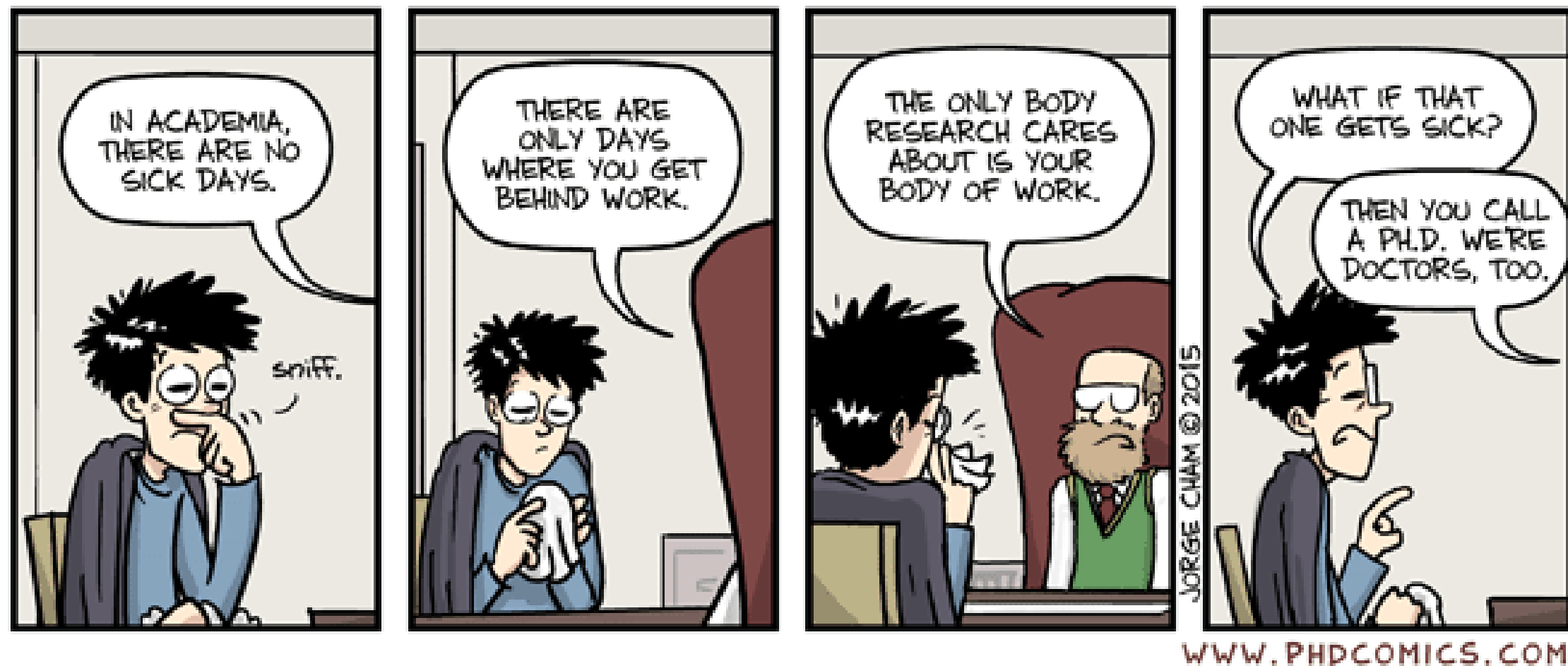


*Bio5488*
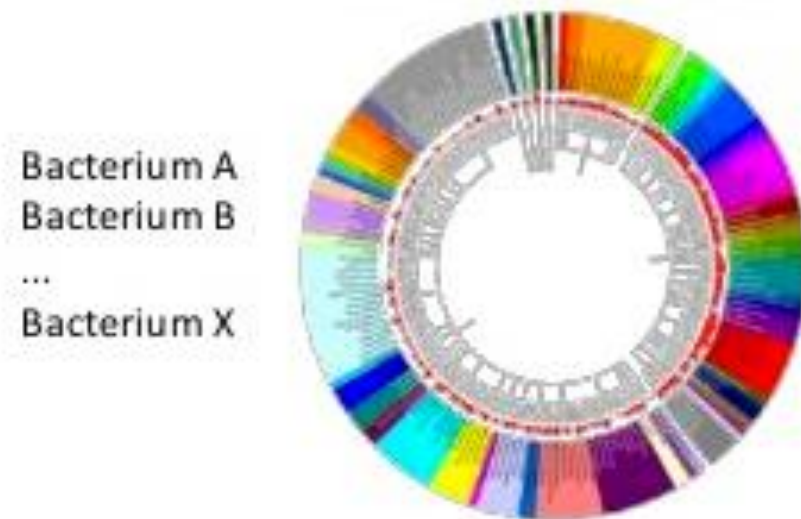*Spring '16*

# Extracting genomes from metagenomes
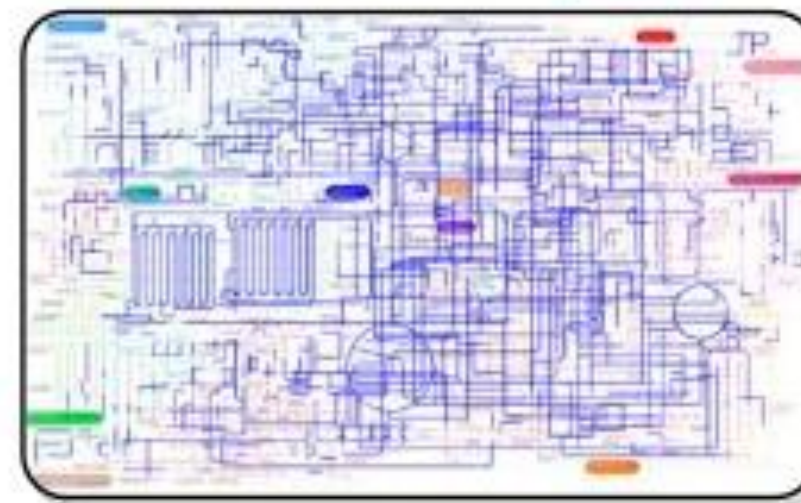
# Metagenomics process

# Six open reading frames (ORFs)

5'                                                                    3'

atgcccaagctgaatagcgtagagggggttttcatcatttgaggacgatgtataa
---------+---------+---------+---------+---------+----
tacgggttcgacttatcgcatctcccccaaaagtagtaaactcctgctacatatt

3'                                                                    5'

This DNA fragment can be read in **six reading frames:**

**A)** Three in the forward direction (5'→3')

| 1 | Atg | ccc | aag | ctg | aat | agc | gta | gag | ggg | ttt | tca | tca | ttt | gag | gac | gat | gta | taa |
|---|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
|   | M   | P   | K   | L   | N   | S   | V   | E   | G   | F   | S   | S   | F   | E   | D   | D   | V   | *   |

| 2 | a Tgc | cca | agc | tga | ata | gcg | tag | agg | ggt | ttt | cat | cat | ttg | agg | acg | atg | tat | aa |
|---|-------|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|----|
|   | C     | P   | S   | *   | I   | A   | *   | R   | G   | F   | H   | H   | L   | R   | T   | M   | Y   |    |

| 3 | at Gcc | caa | gct | gaa | tag | cgt | aga | ggg | gtt | ttc | atc | att | tga | gga | cga | tgt | ata | a |
|---|--------|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|---|
|   | A      | Q   | A   | E   | *   | R   | R   | G   | V   | F   | I   | I   | *   | G   | R   | C   | I   |   |

**B)** Three in the reverse direction (5'→3')

| 4 | tac | ggg | ttc | gac | tta | tcg | cat | ctc | ccc | aaa | agt | agt | aaa | ctc | ctg | cta | cat | atT |
|---|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
|   | H   | G   | L   | Q   | I   | A   | Y   | L   | P   | K   | *   | *   | K   | L   | V   | I   | Y   | L   |

| 5 | ta | cgg | gtt | cga | ctt | atc | gca | tct | ccc | caa | aag | tag | taa | act | cct | gct | aca | taT | t |
|---|----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|---|
|   |    | G   | L   | S   | F   | L   | T   | S   | P   | N   | E   | D   | N   | S   | S   | S   | T   | Y   |   |

| 6 | t | acg | ggt | tcg | act | tat | cgc | atc | tcc | cca | aaa | gta | gta | aac | tcc | tgc | tac | atA | tt |
|---|---|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|----|
|   |   | A   | W   | A   | S   | Y   | R   | L   | P   | T   | K   | M   | M   | Q   | P   | R   | H   | I   |    |

The stop codons are indicated with "*"

Adapted from: http://en.bioinformatica-na-escola.org/media/images/research3/orf.png

# Search using blastp

```
CD-Length = 395 residues,  99.7% aligned
Score =  538 bits (1388), Expect = 3e-154
```

```
Query:  2    ADIQLSKYHVSKDIGFLLEPLQDVLPDYFAPWNRLAKSLPDLVASHKFRDAVKEMPLLDS  61
Sbjct:  1    SLPILEKYHISEDVGFLLPPLQRLLPDKYMPWEEIAKDLPSLIESGKLREVVEKLPVLDL  60

Query:  62   SKLAGYRQKRLAHLQLVLITSGYLWQEGEGGAVQRLPECVAKPLWNVSNDLGLKPVLTYG  121
Sbjct:  61   DELGDHREQRLAHLILGFITMAYVWASGTGDVRKVLPECIAVPLCELSHKLGLPPILTYA  120

Query:  122  DVCLTNCRVKG-------GDIEVMYNLPGGAGTEWFLKVCGLVELTLGKGAQSVQNVLDG  174
Sbjct:  121  DCVLANWKVKDPNGPLTYENIDVLFSFPGGDCEKWFFLVSLLVEIAASAAIKAIPTVLRA  180

Query:  175  AKANDKAKMTSGLTELTTTIGNMQAALAKMNDNLTPDHFYNVLRPFLGGFGGPASPISGG  234
Sbjct:  181  IRSQDKANLIKGLEDLAATIEKASKALMRMEDKVEPNVFYFVLRPFLSGWKGMSSMLSPG  240

Query:  235  LIYEGVSDAPVTMIGGSAAQSSAMQLLDNLLGVTHSPDKQ---AFLDEISNYMIPAHKQL  291
Sbjct:  241  LVYEGVWDQPKIFSGGSAAQSSLFQTLDVLLGIKHTAGKAHSANFLDEMRKYMPPAHRNF  300

Query:  292  LADLTKMPRKVPQIVAEAKDANLSKAYSGCVAALTQYRTYHIQVVTKYIVTASK------  345
Sbjct:  301  LYHLESVPNIVREVVRSASNAALTEAYNRCVSALVSFRSYHIQIVTKYIILPSNSKPKPN  360

Query:  346  --SDSPKSLAYKDTGKSDLIPFLKEVRDDTEKMQ  377
Sbjct:  361  VLSEIPSNLEAKGTGGTDLMPFLKQVRDTTEKTL  394
```
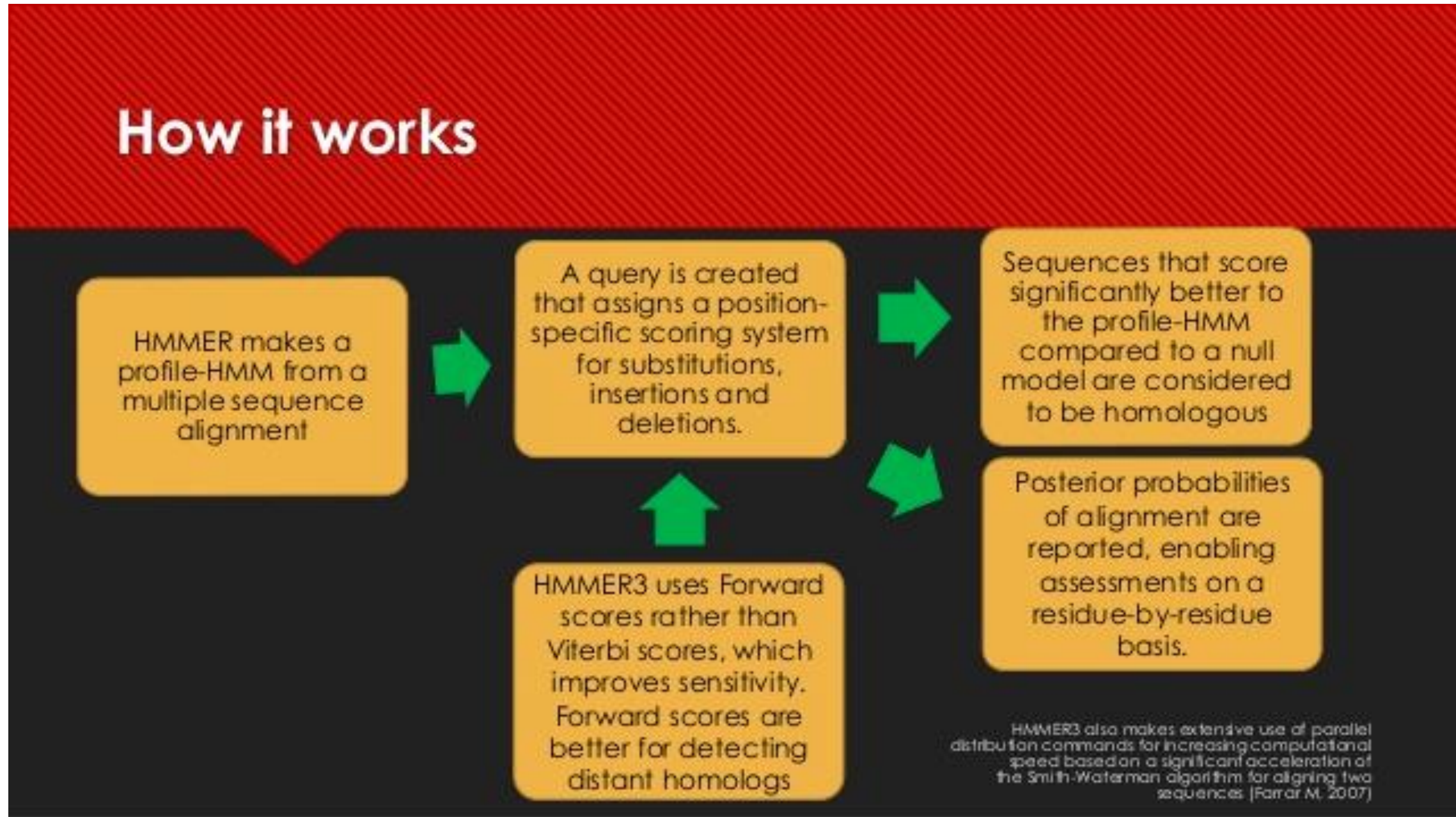
# Search using superior Profile HMM models via HMMer

# What is HMMer?

- HMMer is a open source program suite used to implement profile HMM for biological sequence analysis.

- Used to make the Pfam database of protein families

- Introduction to HMMer: http://pt.slideshare.net/anaxfotopoulos/introduction-to-hmmer-a-hidden-markov

# …for those interested

# Please turn in:

- A README.txt with the answers to the questions and the commands you used to answer the questions.
- A **commented** call_orfs.py to identify ORFs in a fasta of contigs
- A **commented** compare_orf_callers.py to compare MetaGeneMark output to your call_orfs.py output.
- A **commented** count_ar_genes_from_blast.py that filters your blast output.
- A **commented** count_ar_genes_from_resfams.py that counts the number of genes in a hmmscan output file.
- Optional: a **commented** gff_to_nt_aa.py that takes a MetaGeneMark GFF file and outputs the nucleotide and amino acid sequences from the predicted ORFs.
- All files created from the above scripts or commands: all_orfs.fna, all_proteins.faa, mgm_predictions.gff, mgm_orfs.fna, mgm_orfs.faa, mgm_predictions.gff, blast_to_card.txt, and resfams_annotations.txt
- **Due in 3 weeks (3/23/15) at 10 AM**

# Questions

# Good luck on your (metagenomic) journey! :)