# Gene Prediction in de novo Metagenomic Assemblyes

Ekaterina Sosa
Supervisor: Nikolay Vyahhi
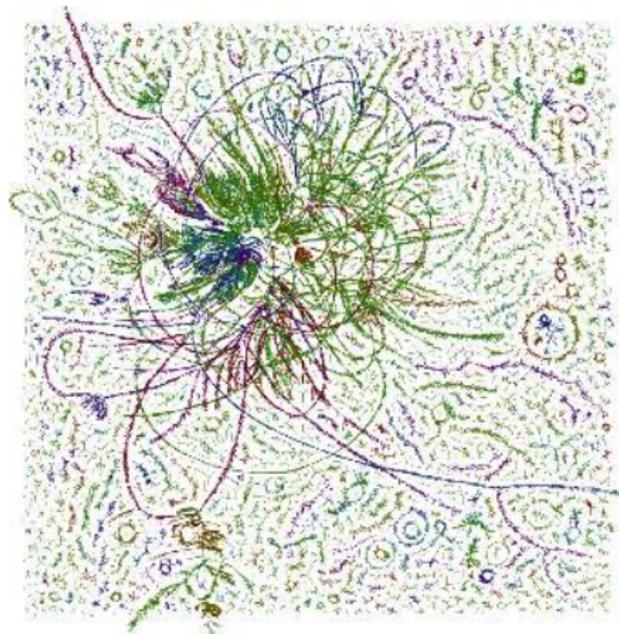
December, 2012

# Initial problem formulation

- to compare prokaryotic gene prediction tools on single-cell datasets
- to investigate metagenomics analysis methods
- to combine metagenomcs analysis and gene prediction

# Introduction: about metagenomics



**Meta-analysis** is the study of statistically combining separate analyses.
**Genomics** is the study of comprehensive analysis of organisms genetic material.
**Metagenomics** is the study of genomic material obtained directly from the environment, instead of from culture.

# Introduction: about prokaryotic gene prediction

A lot of gene finders are based on Hidden Markov Model. Where transition matrix is selected by GC-content.
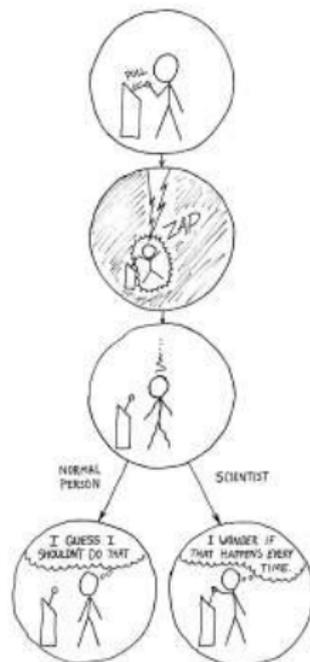
# Tools

- GeneMark
- GeneMark-S
- GeneMark.hmm
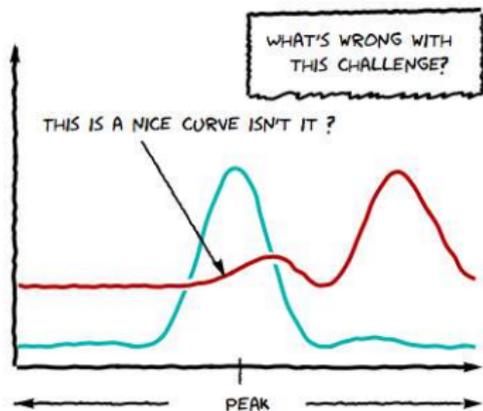- GLIMMER
- EasyGene

... with heuristic selection by

- certain GC-content for every contig
- total GC-content for assembly
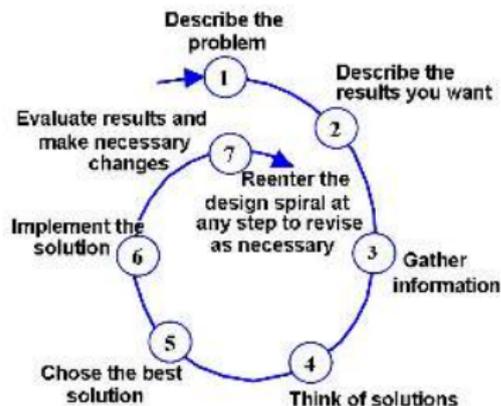
# Comparison

**GeFComp**
is a Python script for comparing
performance metrics of different gene
finding tools. In given a number of genome
assemblies in FASTA format, GeFComp
executes each tool on each of the genomes
and evaluates Type I (false positives) and
Type II (false negatives) errors.

# New problem formulation

- to compare prokaryotic gene prediction tools on single-cell datasets
- **to embed the best tool into QUAST**
- to investigate metagenomics analysis methods
- to combine metagenomcs analysis and gene prediction



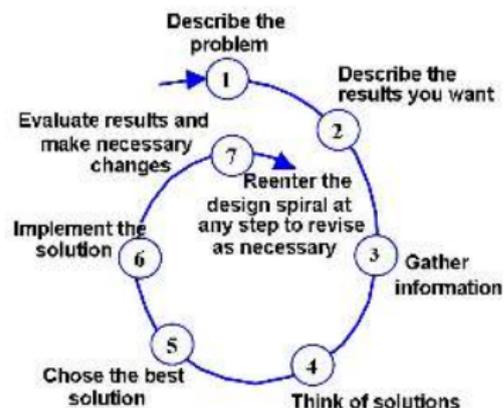**The Technological Method of Problem Solving**

And the best tool was embedded into QUAST
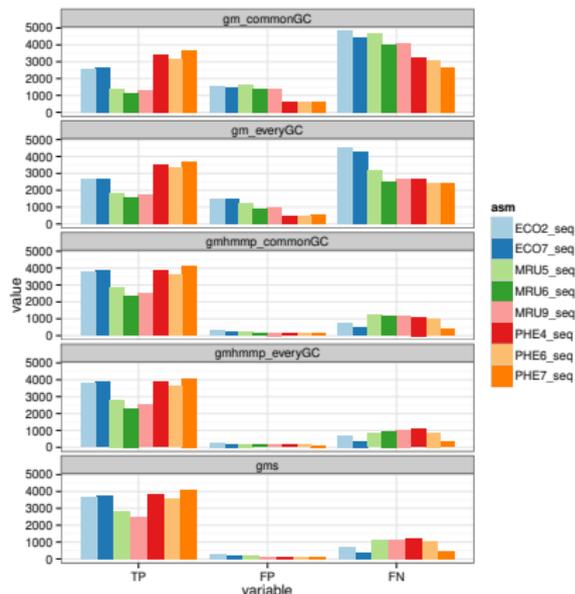
# New new problem formulation

- to compare prokaryotic gene prediction tools on single-cell datasets
- to embed the best tool into QUAST
- **to embed eukaryotic gene prediction into QUAST**
- to investigate metagenomics analysis methods
- to combine metagenomcs analysis and gene prediction

# Current results

**Hypothesis:** gene was predicted
**F**alse **P**ositives is Type I error
**F**alse **N**egatives is Type II error
**T**rue **P**ositives is correct outcome
**T**rue **N**egatives is correct outcome
(TN = ∞)



http://github.com/bioinf/GeneFinder

# Future directions



- to investigate metagenomics analysis methods
- to combine metagenomcs analysis and gene prediction
- to embed metagenomics gene prediction into QUAST
- to gene prediction tool