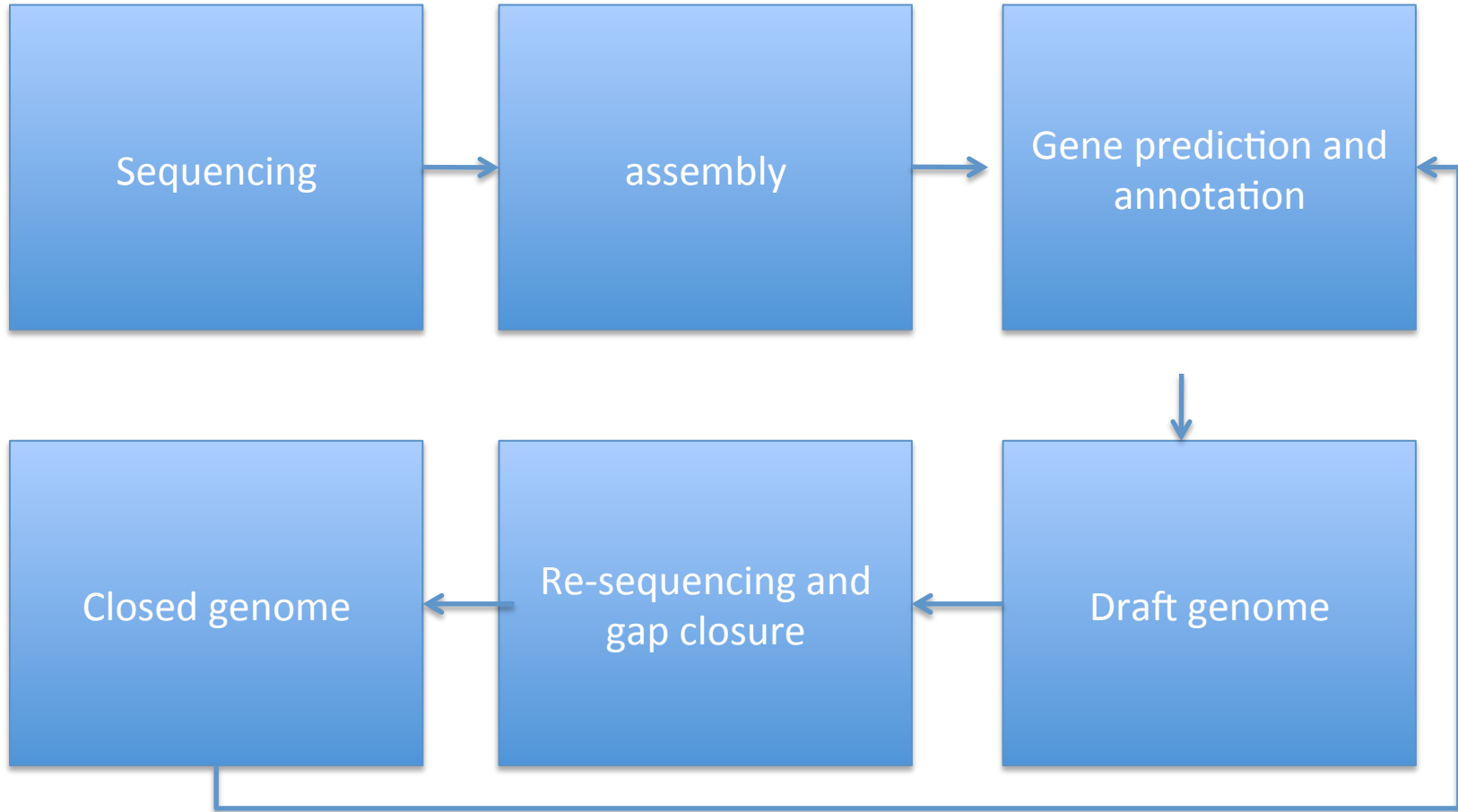# Lecture 7

Annotation

# The whole process

# Annotation

Finding genomic features using computation and experimental methods

1. Genes
2. Repeat sequences
3. Transcriptional regions
4. Other

# levels of annotation

- Gene prediction
- Function: Predictions of domains
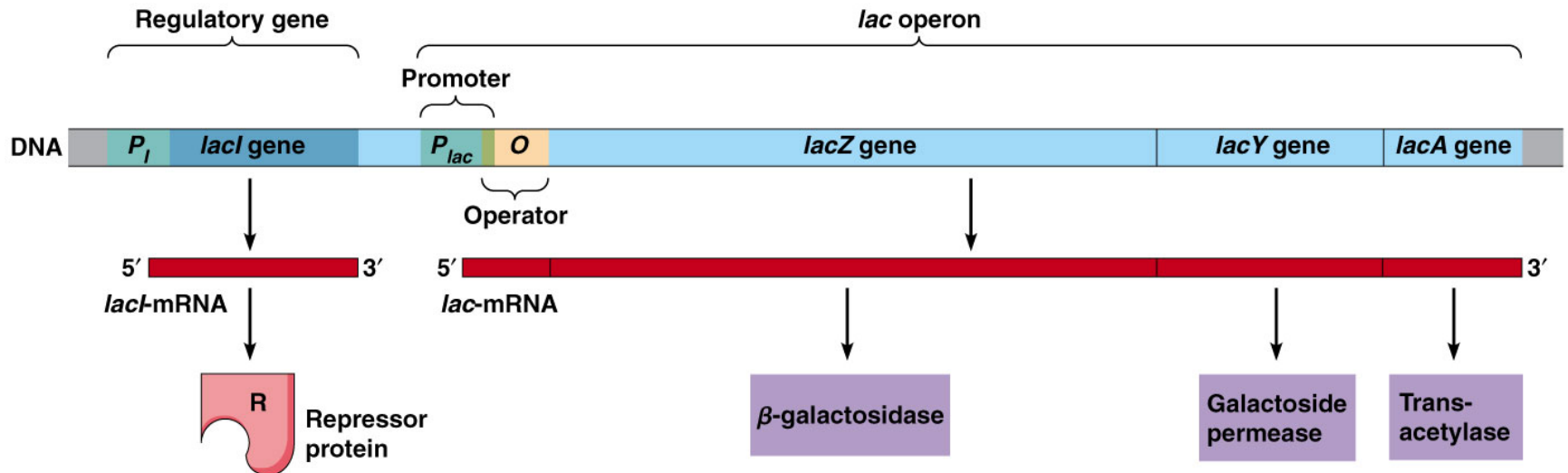- General Role and interaction: What pathways involved in?

# What are genes?

- DNA segments that produce functional products

- Manly proteins!
- RNA molecules!
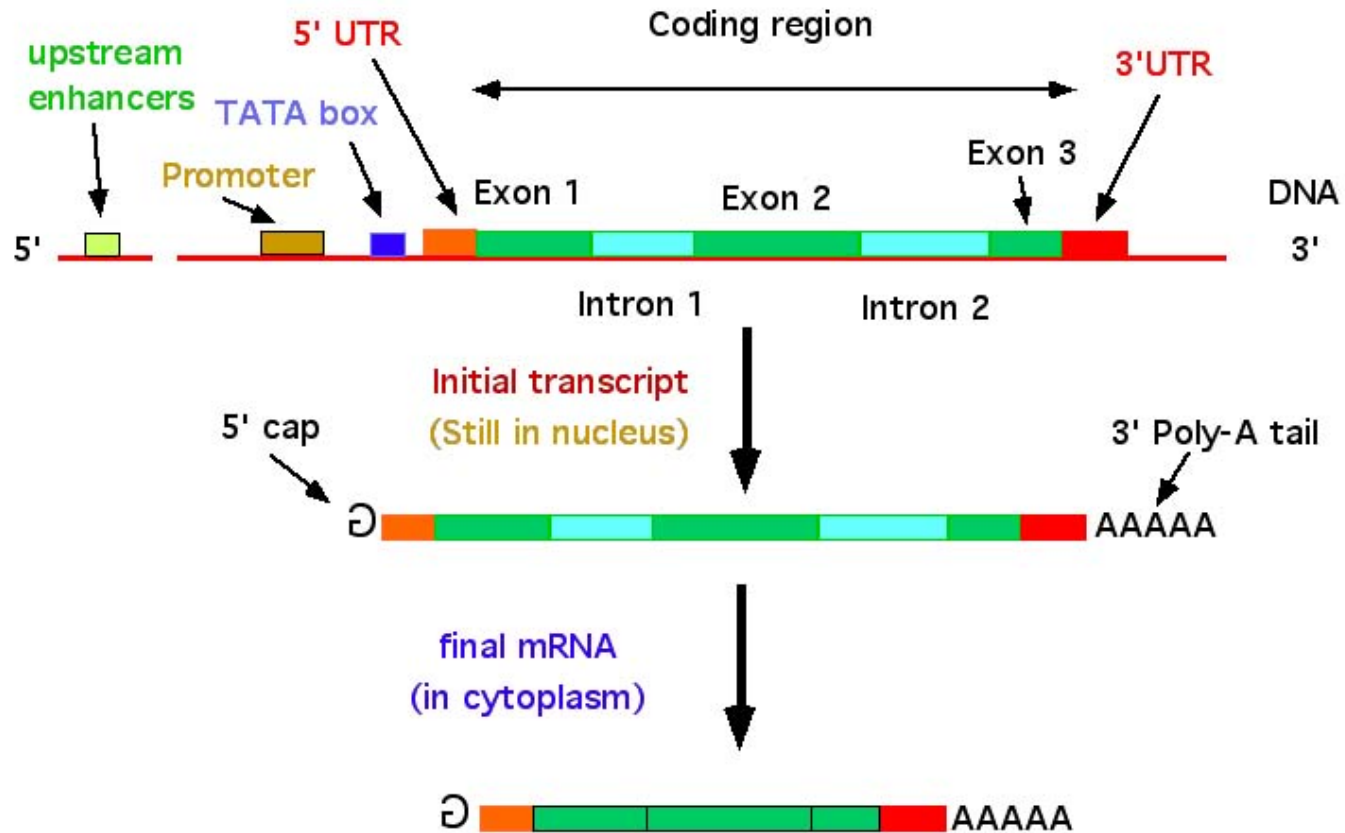- RNAi (interfering RNA)
- rRNA (ribosomal RNA)
- tRNA (trasnfer RNA)

# Prokaryotic vs Eukaryotic gene model

| Prokaryotic | Eukaryotic |
|---|---|
| Small genomes, high gene density | Large genomes, many regions without coding |
| Operons, many genes in one transcript | Operons, monocistronic. Some exceptions |
| Open reading frames. One OFR per genes. ORFs begin with start (ATG, GTG, TTG), end with stop codon (TAG, TAA, TGA) | Posttranscriptional modification |
| | |

# Prokaryotic gene model



© 2012 Pearson Education, Inc.

# Eukaryotic gene model

# Gene identification

1- Homology-based gene prediction

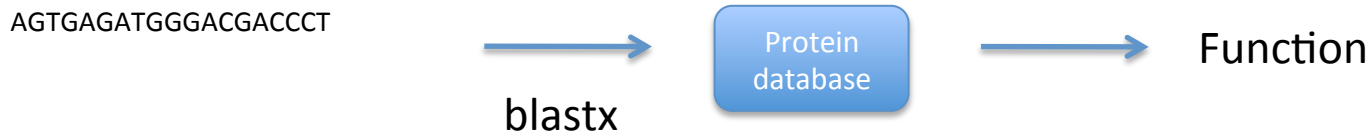-Similarity search (BLAST!)

-Genome browsers
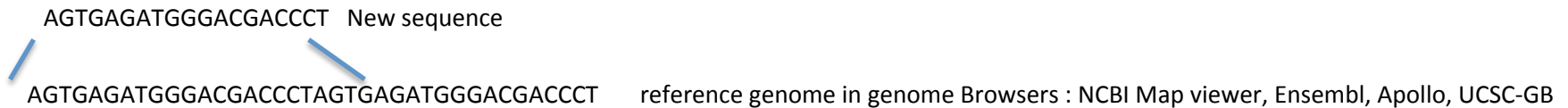
-RNA evidence (ETSs)


2- *Ab initio gene prediction*

Gene predictions programs

# Homology-based gene prediction

- ## Similarity search:

AGTGAGATGGGACGACCCT

blastx → Protein database → Function

- ## Genome browsers:

AGTGAGATGGGACGACCCT   New sequence

AGTGAGATGGGACGACCCTAGTGAGATGGGACGACCCT   reference genome in genome Browsers : NCBI Map viewer, Ensembl, Apollo, UCSC-GB

- ## RNA evidence: ETS

RNA Transcripts

AGTGAGATGGGACGACCCT

AGTGAGATGGGACGACCCT   match

# *Ab initio gene prediction*

- Prokaryotes

ORF detectors

ATGgtgtgg…………………………………..ttggggTGA

Start                                                          stop


- Eukaryotes

Position, extent and direction though promoter and polyA-signal

Promoter seq ………………..Splice-seq………….PolyAAA

# Tools

- ORF detectors

Gorf : http://www.ncbi.nih.gov/gorf/gorf.html

- Promoter predictors

ICG TATA-Box predictor

BDGP fruitfly.org/seq_tools/promoter.html

Virtual footprint (for bacteria) http://www.prodoric.de/vfp

Softberry (for bacteria)
http://linux1.softberry.com/berry.phtml?topic=bprom&group=programs&subgroup=gfindb

- PolyA signal predictors

CSHL: argon.cshl.org/tabaska/polyadq_form.html

- Splice site predictors

BDGP : http://www.fruitfly.org/seq_tools/splice.html

# Gene prediction programs

- Rule-based programs

Use specific set of rules to predict the genes(GeneFinder)

- Hidden Markov Model-based programs

Use probabilities of states and transitions to predict (genscan, genomeScan, glimmer)

# Annotation tools

- Blast searches (Blast)
- HHM models of specific genes or gene families (PFAM)
- Cellular location  (PBSORT, SignalP)
- Biochemical pathwayd/subsystem informtaion (KEGG)

# Gene Ontology  (GO)
# http://www.geneontology.org/

The GO consortium describes gene productys with a standard vocabulary.

Biological process

Cellular component

Molecular function

For example: cytochrome C is described with

Biological process: phosphorylation  GO:0016310

Cellular component: mitochondrial inner membrane GO:0005743

Molecular function: oxido-reductase activity GO:0016491

Each of these has a GO term

# COG and KOG
# Cluster of orthologs group

- http://www.ncbi.nlm.nih.gov/books/NBK21090/

- **Phylogenetic Classification of Proteins from Complete Genomes**

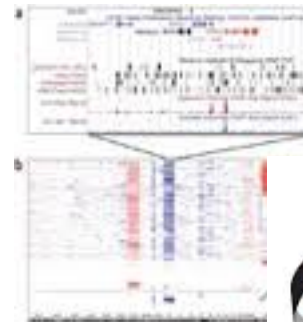# Tools for annotation

Prokka

# RAST server

- [http://rast.nmpdr.org/](http://rast.nmpdr.org/)
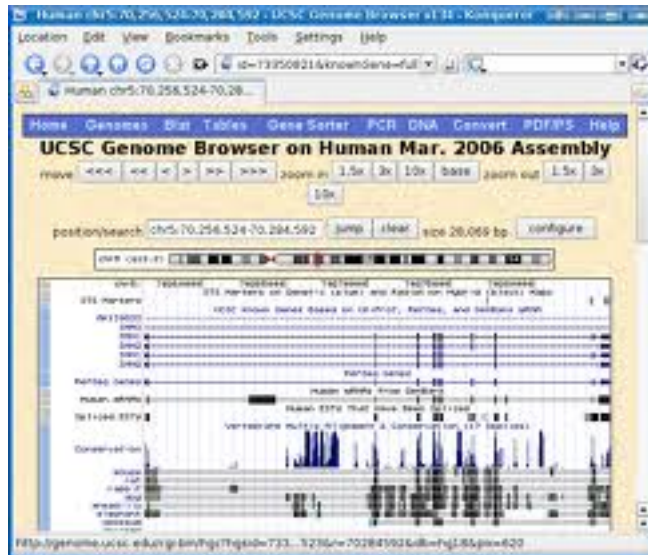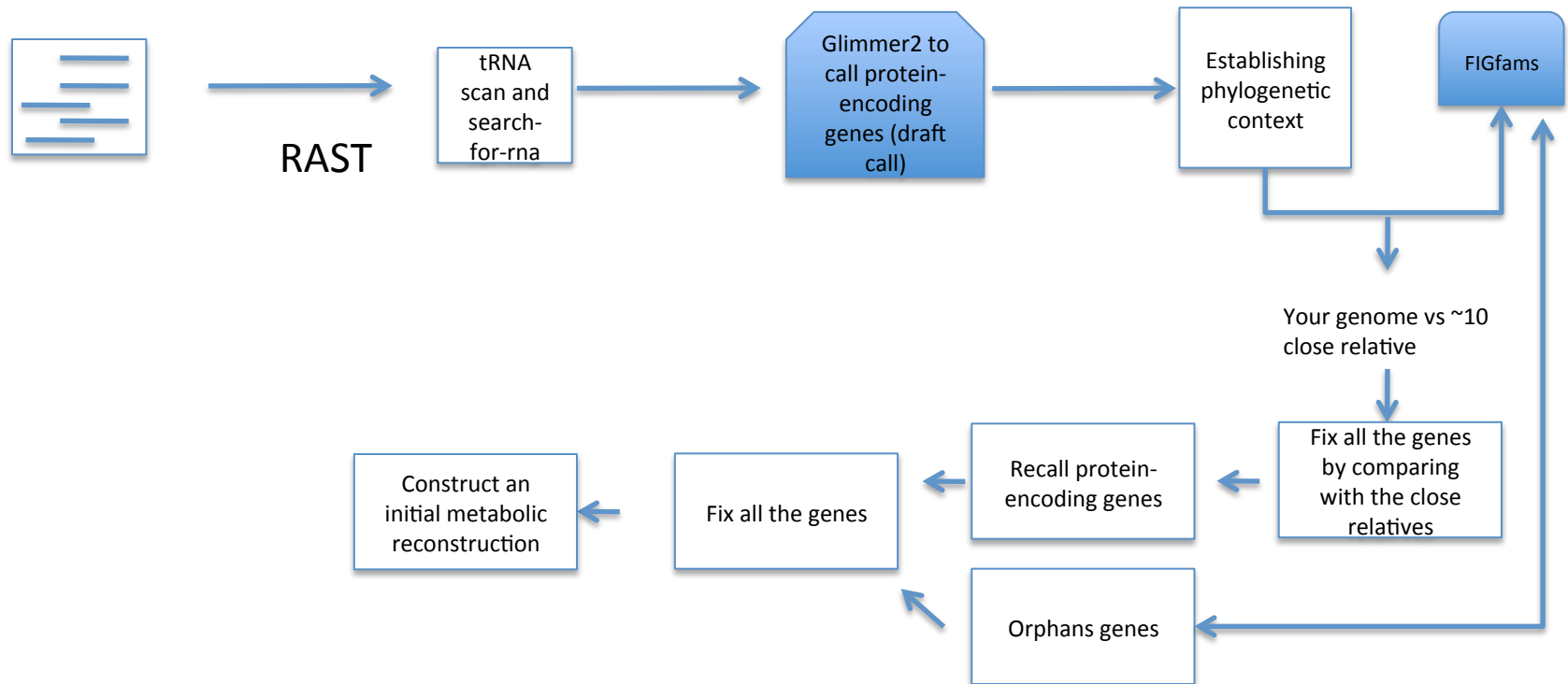- Rapid annotation using subsystem techology
- For prokaryotes.

# Genome browsers, visualization and curation tools

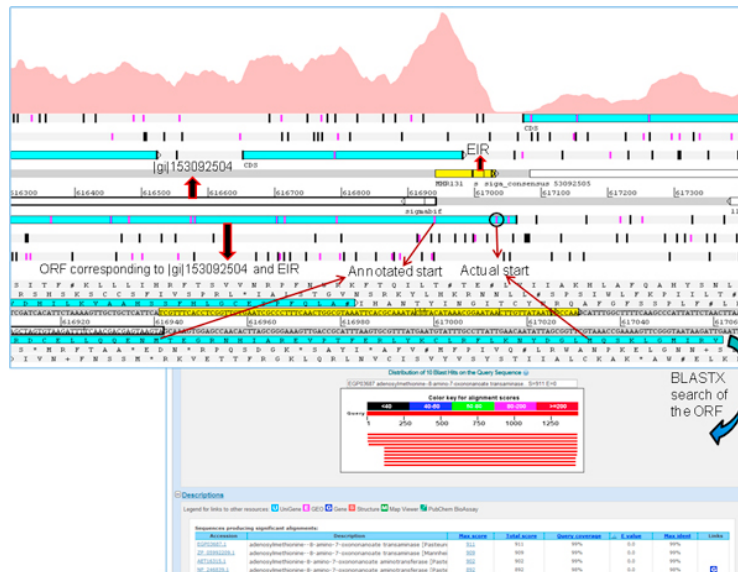# The basic steps in annotating a genome using RAST



RAST

tRNA scan and search-for-rna

Glimmer2 to call protein-encoding genes (draft call)

Establishing phylogenetic context

FIGfams

Your genome vs ~10 close relative

Fix all the genes by comparing with the close relatives

Recall protein-encoding genes

Fix all the genes

Construct an initial metabolic reconstruction

Orphans genes

# Artemis genome browser

- [https://www.sanger.ac.uk/resources/software/artemis/](https://www.sanger.ac.uk/resources/software/artemis/)