

Spring 2024  
**BIOS 477/877**  
*Bioinformatics and Molecular Evolution*  
**Lecture 30**

BIOS477/877 L30 - 1

1

**TODAY'S TOPICS**

- Gene prediction
- Ab initio & Combination methods
- Genome annotation pipelines
- **Course Evaluation on Canvas**  
Submit by 11:59 PM (Saturday)

BIOS477/877 L30 - 2

2

**Gene prediction methods**

- Various pattern recognition methods are used
  - Decision trees [e.g., MORGAN, GlimmerM]
  - Discriminant function analysis (LDA, QDA) [e.g., MZEF, FGENES, HEXON]
  - Neural networks (NN) [e.g., GRAIL II]
  - **Hidden Markov Models (HMMs)** [e.g., GENSCAN, HMMGene, GeneMark.hmm, Glimmer, Augustus]
  - **Combiner or ensemble methods** [e.g., JIGSAW, GLEAN, GeneComber, EvidenceModeler, TSEBRA]
  - **Combinations with similarity or evidence** [e.g., NSCAN, CONTRAST, Augustus, GeMoMa, Maker2, Braker3]

BIOS477/877 L30 - 3

3

**A benchmark study of ab initio gene prediction methods in diverse eukaryotic organisms**

Nicolaas Scalzitti, Anne Jeannin-Girardon, Pierre Collet, Olivier Poch & Julie D. Thompson  
 BMC Genomics 21, Article number: 293 (2020) | Cite this article

**[Benchmark data]**  
 889 confirmed protein-coding genes, including 7,968 exons  
 Collected from diverse organisms (protists ~ human)

| Method                  | Augustus    | GENSCAN | GeneID | SNAP | GlimmerHMM |
|-------------------------|-------------|---------|--------|------|------------|
| Augustus (w/o evidence) | 263         | 47      | 38     | 144  | 51         |
| GENSCAN                 | 47          | 174     | 33     | 161  | 25         |
| GeneID                  | 38          | 33      | 49     | 13   | 5          |
| SNAP                    | 144         | 161     | 13     | 52   | 18         |
| GlimmerHMM              | 51          | 25      | 5      | 18   | 76         |
| <b>Total</b>            | <b>5461</b> |         |        |      |            |

**Correctly predicted exons**

**Sensitivity=TP/(TP+FN); Specificity=TP/(TP+FP)**  
 UDT: undetermined (including 'N')

BIOS477/877 L30 - 4

4

**Gene prediction methods**

- Various pattern recognition methods are used
  - Decision trees [e.g., MORGAN, GlimmerM]
  - Discriminant function analysis (LDA, QDA) [e.g., MZEF, FGENES, HEXON]
  - Neural networks (NN) [e.g., GRAIL II]
  - Hidden Markov Models (HMMs) [e.g., GENSCAN, HMMGene, GeneMark.hmm, Glimmer, Augustus]
  - **Combiner or ensemble methods** [e.g., JIGSAW, GLEAN, GeneComber, EvidenceModeler, TSEBRA]
  - **Combinations with similarity or evidence** [e.g., NSCAN, CONTRAST, Augustus, GeMoMa, Maker2, Braker3]

BIOS477/877 L30 - 5

5

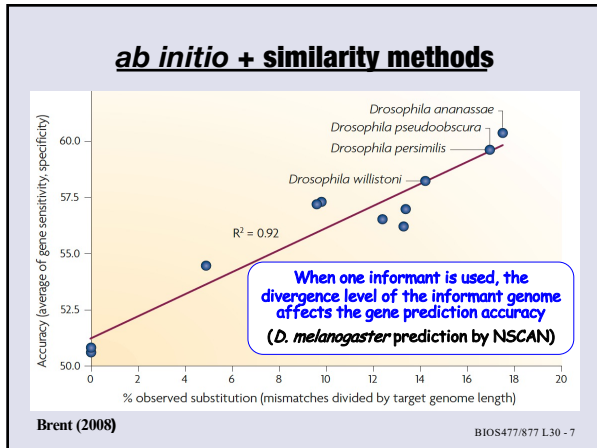
**ab initio prediction + evidence**

Box 2 | Gene prediction versus gene annotation

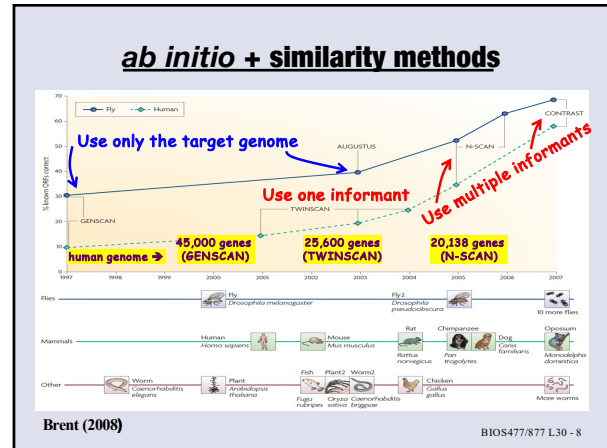
**Various evidence can be used to improve gene prediction results: e.g., 5' and 3' untranslated regions (UTRs), alternatively spliced variants (isoforms)**

Yandell & Ence (2012)  
 BIOS477/877 L30 - 6

6



7



8

### Simultaneous gene finding in multiple genomes

Stefanie König\*, Lars W. Romoth, Lizzy Gerischer and Mario Stanke\* *Bioinformatics* 2016 32: 3388-95

Sensitivity (Sn) = TP/(TP+FN), Specificity (Sp) = TN/(TN+FP) at Gene, Exon, or Nucleotide level

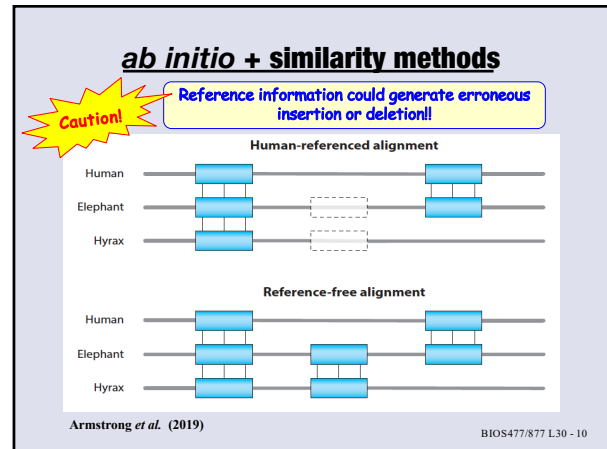
| No RNAseq              | k-way | Gene Sn | Gene Sp | Exon Sn | Exon Sp | Nuc. Sn | Nuc. Sp |
|------------------------|-------|---------|---------|---------|---------|---------|---------|
| <b>D. melanogaster</b> |       |         |         |         |         |         |         |
| AUGUSTUS               | -     | 56.48   | 60.33   | 72.88   | 82.52   | 92.32   | 96.90   |
| AUGUSTUS <sub>cp</sub> | 12    | 66.20   | 65.79   | 77.69   | 85.79   | 93.57   | 96.67   |
| N-SCAN <sub>p</sub>    | 2     | 47.91   | 54.44   | 68.89   | 77.20   | 94.01   | 91.65   |
| CONTRAST               | 15    | 70.49   | 71.99   | 78.07   | 90.76   | 91.84   | 98.30   |

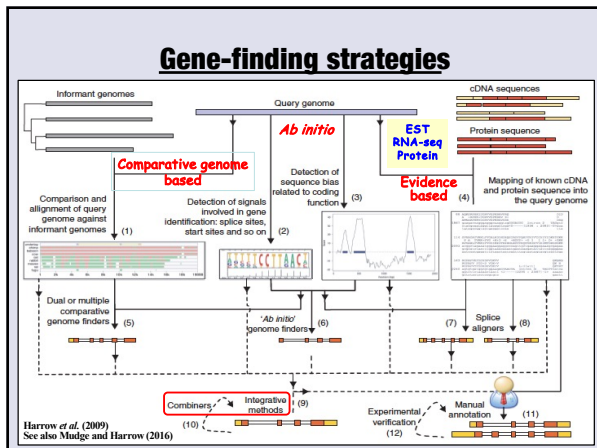
| With RNAseq            | D. mel | Gene Sn | Gene Sp | Exon Sn | Exon Sp | Nuc. Sn | Nuc. Sp |
|------------------------|--------|---------|---------|---------|---------|---------|---------|
| AUGUSTUS               | -      | 68.01   | 73.26   | 76.46   | 90.20   | 93.39   | 97.57   |
| AUGUSTUS <sub>cp</sub> | D. mel | 71.88   | 72.20   | 79.27   | 89.09   | 96.46   | 97.49   |
|                        | D. sim | 67.10   | 67.43   | 77.64   | 86.67   | 95.96   | 97.71   |
|                        | 2 Dros | 74.33   | 73.13   | 80.22   | 89.74   | 96.77   | 97.57   |
|                        | 4 Dros | 74.46   | 73.18   | 80.31   | 89.79   | 96.79   | 97.56   |

CGP: Comparative Gene Prediction BIOS477/877 L30 - 9

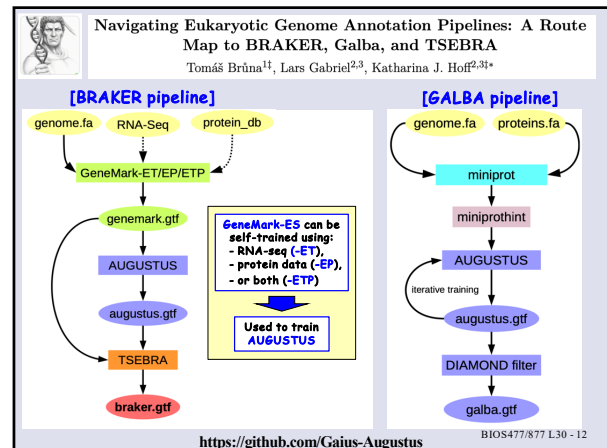
9



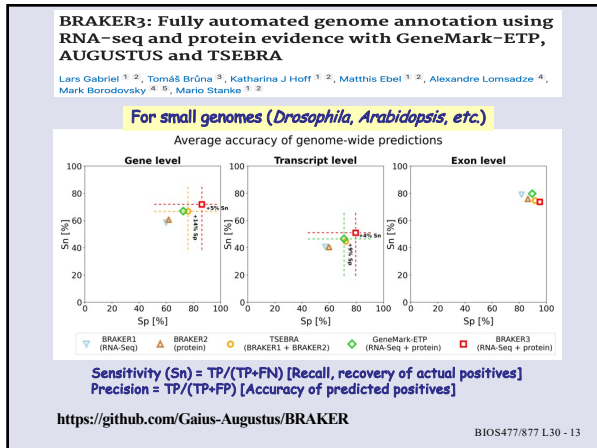
10



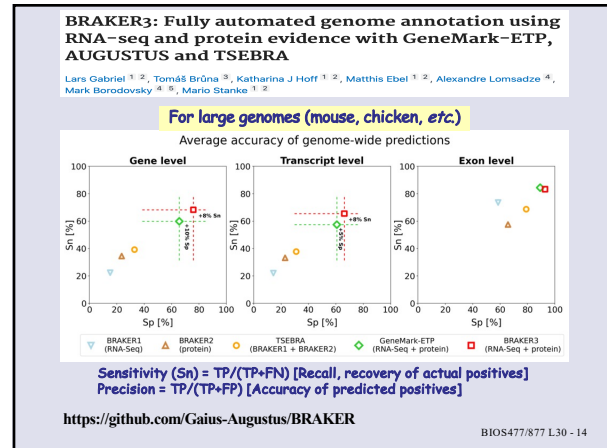
11



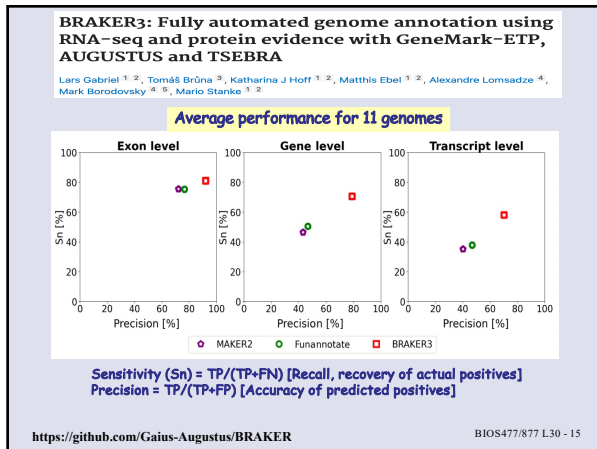
12



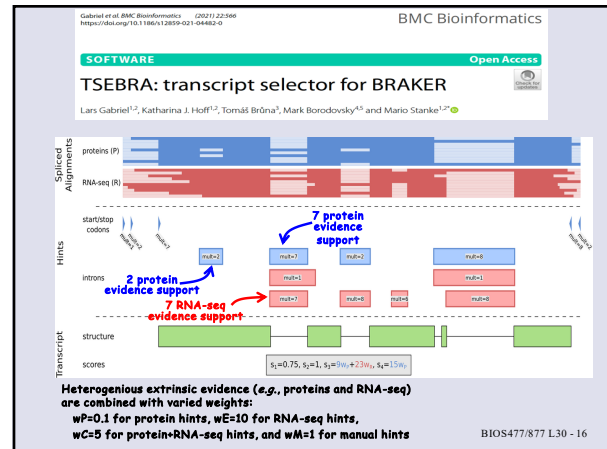
13



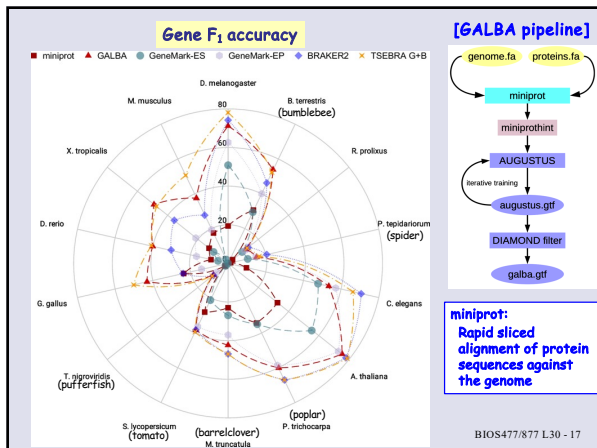
14



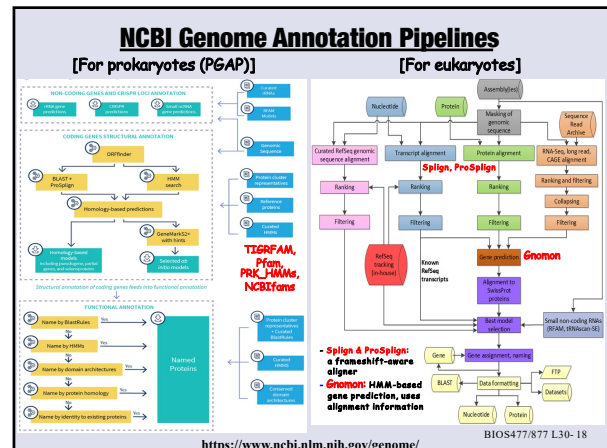
15



16



17



18

## Prokaryotic genome annotation tools

- **PATRIC** (Pathosystems Resource Integration Center) <https://www.bv-brc.org>  
Assembly, RAST-based annotation, etc.
- **RAST** (Rapid Annotation using Subsystem Technology) <https://www.me-rast.org/>  
Gene calling (Glimmer3, GeneMarkS, Prodigal), protein annotation (FIGfams), etc.
- **Prokka** (Rapid prokaryotic genome annotation) <https://github.com/seemann/prokka>  
Gene calling (Prodigal), protein annotation using HMMER, etc.
- **IMG/M** (Integrated Microbial Genomes & Microbiomes) <https://img.jgi.doe.gov/m/>  
Assembly, protein-coding gene calling (Prodigal, GeneMark), protein annotation (COGs, Pfam, TIGRFAM, SUPERFAMILY, SMART, CATH-FunFam), etc.
- **KBase** (The DOE Systems Biology Knowledgebase) <http://kbase.us/>  
Assembly, annotation with RAST, Prokka, many other sequence analyses.

Also in Kimbrel et al. (2022) *Methods in Molecular Biology*, vol 2349. BIOS477/877 L30-19

19

## No one tool to rule them all: prokaryotic gene prediction tool annotations are highly dependent on the organism of study

Nicholas J Dilonaco, Wayne Aubrey, Kim Kenobi, Amanda Clare, Christopher J Creevey Author Notes  
*Bioinformatics*, Volume 38, Issue 5, 1 March 2022, Pages 1198–1207.

Tested against six genomes: *B. subtilis*, *C. crescentus*, *E. coli*, *M. genitalium*, *P. fluorescens*, and *S. aureus*

**Overall rankings**

**Rank**

**Metric**

M1: % genes detected (TP/factual P)

M2: % predicted CDS (TP/predicted P)

M3: % difference of # predicted CDSs

M4: % difference of median CDS length

M5: % perfect matches

M6: Median start difference

M7: Median stop difference

M8: % difference overlapping CDS

M9: % difference of short CDS (< 100 bp)

M10: Precision

M11: Recall

M12: False discovery rate

Strict model based

BIOS477/877 L30-20

20

## Controversies in modern evolutionary biology: the imperative for error detection and quality control

Francisco Prosdocimi<sup>1,2\*</sup>, Benjamin Linard<sup>1\*</sup>, Pierre Pontarotti<sup>3</sup>, Olivier Poch<sup>4</sup> and Julie D Thompson<sup>1\*</sup>

On average 41% of sequences were erroneous

[Error types] N/C-deletion: N/C-term deletion; N/C-extension: N/C-term extension; Segment: suspicious seq. segment; Deletion: internal deletion; Insertion: internal insertion

Percentage of predicted sequence errors in 19,778 protein families. Blue: the percentage of sequences with at least one error. Red: the percentage of total errors observed.

BIOS477/877 L30 - 21

21

## Understanding the causes of errors in eukaryotic protein-coding gene prediction: a case study of primate proteomes

Corentin Meyer, Nicolas Scabitti, Anne Jeannin-Grardou, Pierre Collet, Olivier Poch and Julie D. Thompson\*

Frequency of gene prediction errors per protein sequence, in each of the Uniprot primate proteomes

BIOS477/877 L30 - 22

22

## Remember...

- Check your data (sequences)!
- Check your alignments!
- Check your phylogenies!
- Check your outputs!!!

BIOS477/877 L30-23

23

## Assignment #12: due May 13

## Final Graduate Report: due May 14

## Course Evaluation on Canvas Complete by 11:59 PM, May 11 (Sat)

BIOS477/877 L30 - 24

24