



### Eukaryotic gene predictions have high error rates

- Gene finders generally do a **poor job (<60%)** predicting genes in eukaryotes
- More **variations** in the gene models
  - Alternative splicing (multiple isoforms)
  - Non-canonical splice sites (e.g., *toy*)
  - Non-canonical start codon (e.g., *Fmr1*)
  - Stop codon read through (e.g., *gish*)
  - Nested genes (e.g., *ko*)
  - Trans-splicing (e.g., *mod(mdg4)*)
  - Pseudogenes (e.g., *swaPsi*)

7

### Types of eukaryotic gene predictors

- Ab initio**
  - GENSCAN, geneid, SNAP, GlimmerHMM
- Evidence-based (extrinsic)
  - Augustus, genBlastG, GeMoMa, Exonerate, GenomeScan
- Comparative genomics
  - Twinscan/N-SCAN, SGP2
- Transcriptome-based (RNA-Seq)
  - Cufflinks, StringTie, Trinity, CodingQuarry
- Combine *ab initio* and evidence-based approaches
  - Gnomon, MAKER, EASEL, EVM, JIGSAW, IPred, GLEAN

8

### Ab initio gene prediction

- Ab initio* = from the beginning
- Predict genes using only the genomic DNA sequence
  - Search for **signals** of protein coding regions
  - Based on a probabilistic model
    - Hidden Markov Models (HMM)
    - Support Vector Machines (SVM)
- GENSCAN
  - Burge C. and Karlin S. *Prediction of complete gene structures in human genomic DNA*, JMB. (1997) **268**, 78-94

9

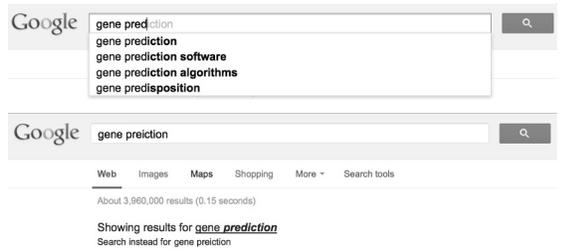
### Hidden Markov Models (HMM)



- A type of supervised machine learning algorithm
  - Uses **Bayesian statistics**
  - Makes classifications based on characteristics of **training data**
- Many types of applications
  - Speech and gesture recognition
  - Bioinformatics
    - Gene predictions
    - Sequence alignments
    - ChIP-seq analysis
    - Protein folding

10

### Supervised machine learning

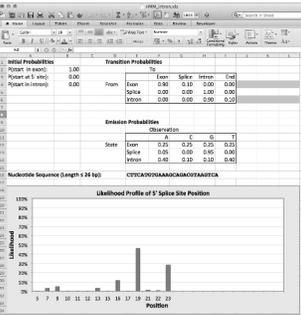


Use previous search results to predict search terms and correct spelling errors

Norvig P. How to write a spelling corrector. <https://www.norvig.com/spell-correct.html>

11

### GEP curriculum on HMM



- Use an HMM to predict a splice donor site
  - Use Excel to experiment with different emission and transition probabilities
- See the **Curriculum** section of the GEP website
  - Also available on CourseSource

Weisstein AE et al. *A Hands-on Introduction to Hidden Markov Models*. CourseSource. (2016). <https://doi.org/10.24918/cs.2016.8>

12

### Ways to create **training sets** to estimate transition and emission parameters

- ☉ Manually curated genes for the target species
- ☉ Bootstrap with *ab initio* gene predictions
  - ☉ GeneMark-ES, GENSCAN
- ☉ Sequence similarity to orthologs in informant species
  - ☉ BUSCO, BRAKER2
- ☉ Whole genome conservation profiles
  - ☉ Augustus-cgp, N-SCAN, SGP2
- ☉ RNA-Seq (splice junctions, assembled transcripts)
  - ☉ BRAKER1

13

### BRAKER2 training protocols

**Training with genome assembly only**

**Training with proteins and RNA-Seq alignments**

Hoff KJ. BRAKER 2 User Guide. <https://github.com/Gaius-Augustus/BRAKER>

14

### GENSCAN HMM Model

- ☉ GENSCAN considers:
  - ☉ Promoter, splice sites and polyadenylation signals
  - ☉ Hexamer frequencies and base compositions
  - ☉ Probability of coding and non-coding DNA
  - ☉ Distributions of gene, exon and intron lengths

Burge C. and Karlin S. Prediction of complete gene structures in human genomic DNA, JMB. (1997) 268, 78-94

15

### Use multiple HMMs to describe different parts of a gene

Stanke M. and Waack S. Gene prediction with a hidden Markov model and a new intron submodel. Bioinformatics. (2003) 19, Suppl 2:ii215-25.

16

### Evidence-based gene predictions

- ☉ Use sequence alignments to improve predictions
  - ☉ EST, cDNA or protein from closely-related species

**Exon sensitivity:**  
Percent of real exons identified

**Exon specificity:**  
Percent of predicted exons that are correct

Yeh RF, et al. Computational Inference of Homologous Gene Structures in the Human Genome, Genome Res. (2001) 11, 803-816

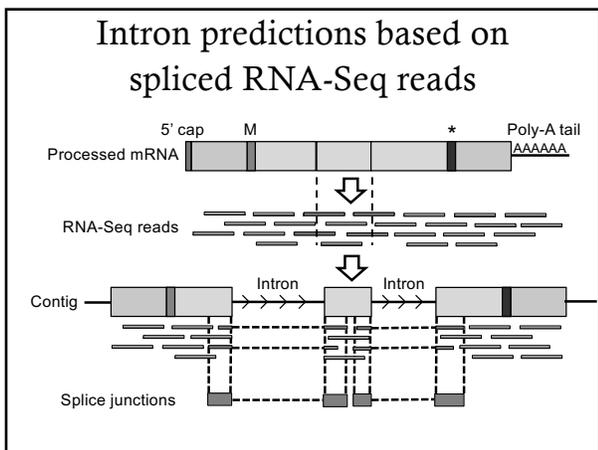
17

### Predictions using comparative genomics

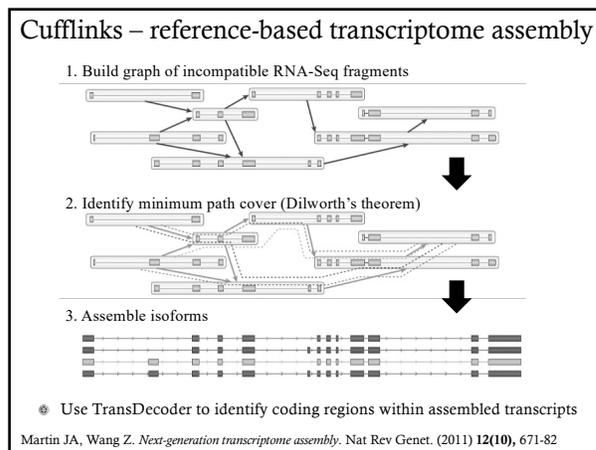
- ☉ Use whole genome alignments from one or more **informant species**
- ☉ CONTRAST predicts 50% of genes correctly
- ☉ Requires **high quality whole genome alignments and training data**

Flicek P. Gene prediction: compare and CONTRAST. Genome Biology (2007) 8, 233

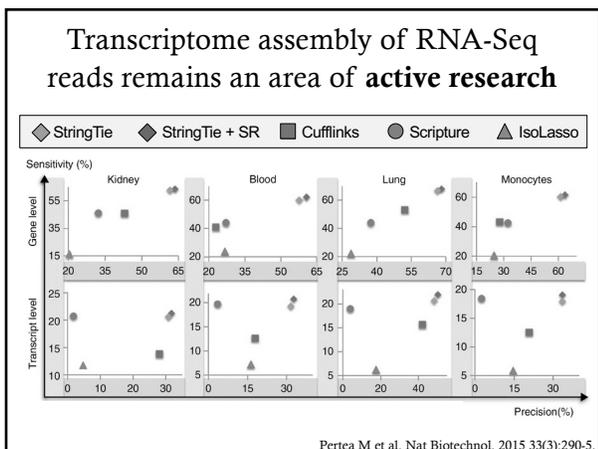
18



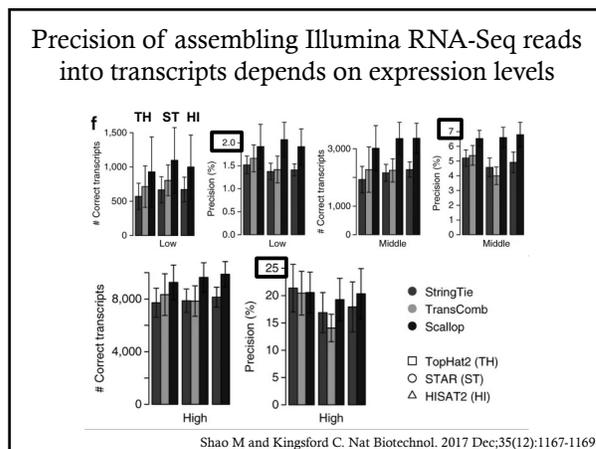
19



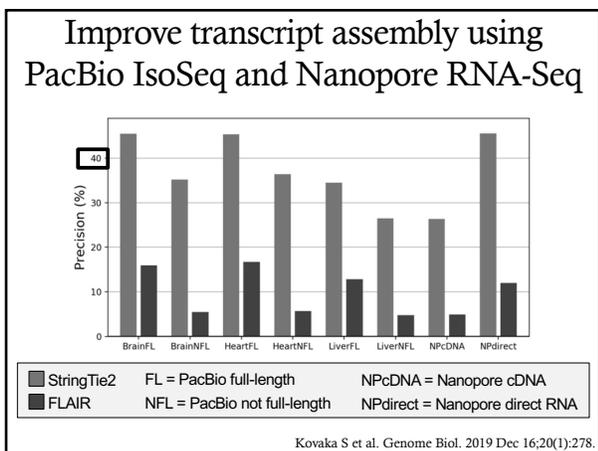
20



21



22



23

### Generate consensus gene models

- Gene predictors have different strengths and weaknesses
- Create **consensus gene models** by combining results from multiple gene finders and sequence alignments
  - EvidenceModeler (EVM)
    - Haas BJ et al. *Automated eukaryotic gene structure annotation using EvidenceModeler and the Program to Assemble Spliced Alignments*. Genome Biology (2008) 9(1), R7
  - TSEBRA
    - Gabriel L et al. *TSEBRA: transcript selector for BRAKER*. BMC Bioinformatics (2021) 22(1), 566

24

### Automated annotation pipelines

NCBI Gnomon gene prediction pipeline

- Integrate **biological evidence** into the predicted gene models
- Examples:
  - NCBI Gnomon
  - Ensembl
  - UCSC Gene Build
- EGASP results for the Ensembl pipeline:
  - 71.6% gene sensitivity
  - 67.3% gene specificity

[https://www.ncbi.nlm.nih.gov/genome/annotation\\_euk/gnomon/](https://www.ncbi.nlm.nih.gov/genome/annotation_euk/gnomon/)

25

### Eukaryotic genomes annotated by NCBI

- RefSeq annotations available for more than **1,100 species**

[https://www.ncbi.nlm.nih.gov/genome/annotation\\_euk/#graphs](https://www.ncbi.nlm.nih.gov/genome/annotation_euk/#graphs)

26

### Drosophila RefSeq gene predictions

- Based on **RNA-Seq** data from either the same or closely-related species
- Predictions include **untranslated regions** and **multiple isoforms**
- Gnomon gene predictions are available through the NCBI RefSeq database:
  - [https://www.ncbi.nlm.nih.gov/genome/annotation\\_euk/all/](https://www.ncbi.nlm.nih.gov/genome/annotation_euk/all/)

27

### EASEL — incorporate feature filtering and functional annotations into the workflow

Efficient, Accurate, Scalable Eukaryotic models

<https://github.com/PlantGenomicsLab/easel>

28

### Most genomes can benefit from manual gene annotations

- Average sensitivity (Sn) and sensitivity (Sp) of gene predictions for eight eukaryotic genomes

Gabriell L et al. bioRxiv 2023.06.10.544449; doi:10.1101/2023.06.10.544449

29

### Common problems with gene finders

- Split single gene into multiple predictions
- Fused with neighboring genes
- Missing exons
- Over predict exons or genes
- Missing isoforms

30

### Non-canonical splice donors and acceptors

- Many gene predictors strongly prefer models with canonical splice donor (**GT**) and acceptor (**AG**) sites
- Check **Gene Record Finder** or FlyBase for genes that use non-canonical splice sites in *D. melanogaster*

Introns with Non-canonical Splice Sites			
Transcript Name	FlyBase ID	Splice Donor	Splice Acceptor
Cadps-RD	intron_Cadps:22_Cadps:23	GC	AG

Frequency of non-canonical splice sites in FlyBase Release 6.55 (Number of unique introns: 72,061)

Donor site	Count	Acceptor site	Count
GC	603	AC	34
AT	30	TG	28
GA	15	AT	16

31

### Annotate unusual features in gene models using *D. melanogaster* as a reference

- Examine the “**Comments on Gene Model**” and the “**Sequence Ontology**” sections of the FlyBase Gene Report

Non-canonical start codon:

Comments on Gene Model: Gene model reviewed during 6.02  
Unconventional translation start (CUG) postulated; FBr0213401.

Sequence Ontology: Class of Gene: gene\_with\_start\_codon\_CUG

Stop codon read through:

Comments on Gene Model: Gene model reviewed during 5.44  
Stop-codon suppression (UGA) postulated; FBr0218884.

Sequence Ontology: Class of Gene: gene\_with\_edited\_transcript (Rodriguez et al., 2012)  
gene\_with\_stop\_codon\_read\_through (Jungras et al., 2011)

32

### Nested genes in *Drosophila*

33

### Trans-spliced gene in *Drosophila*

A special type of RNA processing where exons from two primary transcripts are ligated together

34

### Gene prediction results for the GEP annotation projects

- Gene prediction results are available through the GEP UCSC Genome Browser mirror
  - Under the **Genes and Gene Prediction Tracks** section
- Access the predicted peptide sequence:
  - Click on the feature, and then click on the **Predicted Protein** link

35

### Summary

- Gene predictors can quickly identify potentially interesting features within a genomic sequence
- The predictions are hypotheses that must be confirmed experimentally
- Eukaryotic gene predictors generally can accurately identify internal exons
  - Much lower sensitivity and specificity when predicting complete gene models

36