



Gene Prediction in Eukaryotes

Talk *Gene Prediction* on Feb. 22th, 2012

Introduction

Gene Finding Problem

Statistical Features of
Genes

Signal and Content Models

Training

Gene Finding as
Optimization Problem

HMMs/CRFs

Comparative Gene
Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

Mario Stanke

Institut für Mathematik und Informatik
Universität Greifswald



Introduction

Gene Finding Problem
Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

1 Introduction

Gene Finding Problem

Statistical Features of Genes

2 Signal and Content Models

3 Training

4 Gene Finding as Optimization Problem

5 HMMs/CRFs

6 Comparative Gene Finding

7 Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

8 Alternative Splicing

Gene Finding Problem

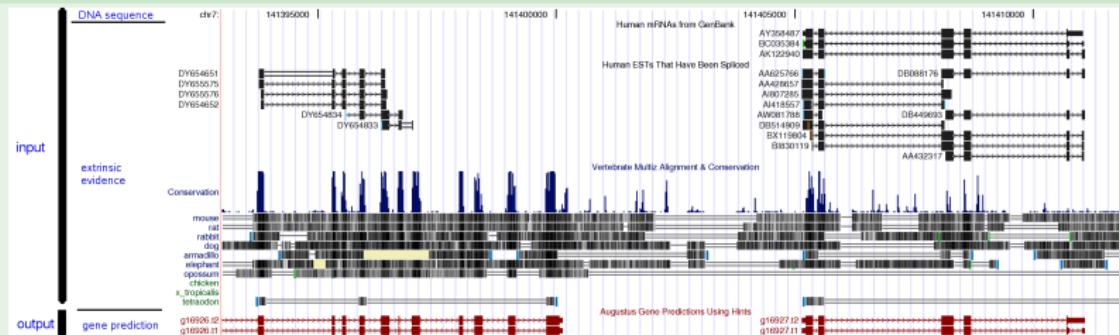
Input y

- target **DNA sequence**, e.g. a chromosome
 - optional: **extrinsic evidence**, e.g. from RNA-Seq

Output x

- start- and end positions of genes, coding regions, exons and introns
 - predicted amino acid sequences

Example (19000 bp of human chromosome 7)





Introduction

Gene Finding Problem

Statistical Features of
Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

Major Approaches to Gene Prediction

Approaches

approach	extrinsic evidence used
<i>ab initio</i>	-



Introduction

Gene Finding Problem

Statistical Features of
Genes

Signal and Content
Models

Training

Gene Finding as
Optimization Problem

HMMs/CRFs

Comparative Gene
Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

Major Approaches to Gene Prediction

Approaches

approach	extrinsic evidence used
<i>ab initio</i>	-
transcript based	transcript seqs (both native and alien)



Introduction

Gene Finding Problem

Statistical Features of
Genes

Signal and Content
Models

Training

Gene Finding as
Optimization Problem

HMMs/CRFs

Comparative Gene
Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

Major Approaches to Gene Prediction

Approaches

approach	extrinsic evidence used
<i>ab initio</i>	-
transcript based	transcript seqs (both native and alien)
protein homology	protein sequences



Introduction

Gene Finding Problem

Statistical Features of
Genes

Signal and Content
Models

Training

Gene Finding as
Optimization Problem

HMMs/CRFs

Comparative Gene
Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

Major Approaches to Gene Prediction

Approaches

approach	extrinsic evidence used
<i>ab initio</i>	-
transcript based	transcript seqs (both native and alien)
protein homology	protein sequences
comparative (<i>de novo</i>)	additional unannotated genomes



Major Approaches to Gene Prediction

Approaches

approach	extrinsic evidence used
<i>ab initio</i>	-
transcript based	transcript seqs (both native and alien)
protein homology	protein sequences
comparative (<i>de novo</i>)	additional unannotated genomes
proteogenomics	peptides from mass spectrometry

Introduction

Gene Finding Problem

Statistical Features of
Genes

Signal and Content
Models

Training

Gene Finding as
Optimization Problem

HMMs/CRFs

Comparative Gene
Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing



Major Approaches to Gene Prediction

Introduction

Gene Finding Problem

Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

Approaches

approach	extrinsic evidence used
<i>ab initio</i>	-
transcript based	transcript seqs (both native and alien)
protein homology	protein sequences
comparative (<i>de novo</i>)	additional unannotated genomes
proteogenomics	peptides from mass spectrometry
combiners	other gene predictions + transcript seqs + proteins + ?



Major Approaches to Gene Prediction

Introduction

Gene Finding Problem

Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

Approaches

approach	extrinsic evidence used
<i>ab initio</i>	-
transcript based	transcript seqs (both native and alien)
protein homology	protein sequences
comparative (<i>de novo</i>)	additional unannotated genomes
proteogenomics	peptides from mass spectrometry
combiners	other gene predictions + transcript seqs + proteins + ?

The annotation of most genomes requires a combination of approaches:

Use for every part of a gene all evidence available for that gene or region.



Introduction

Gene Finding Problem

Statistical Features of
Genes

Signal and Content
Models

Training

Gene Finding as
Optimization Problem

HMMs/CRFs

Comparative Gene
Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

1 Introduction

Gene Finding Problem

Statistical Features of Genes

2 Signal and Content Models

3 Training

4 Gene Finding as Optimization Problem

5 HMMs/CRFs

6 Comparative Gene Finding

7 Evidence Integration

Hints to AUGUSTUS

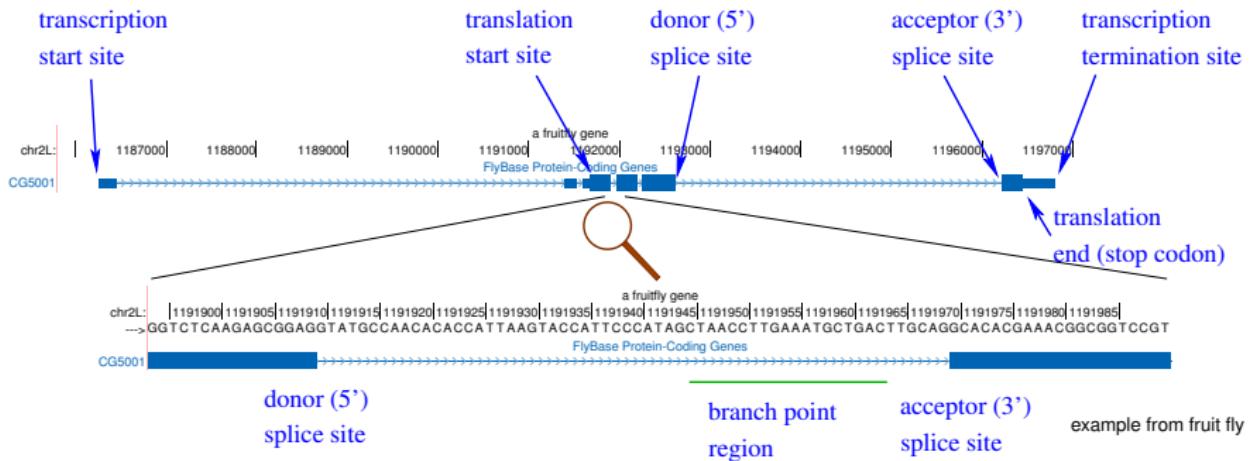
RNA-Seq

Proteogenomics

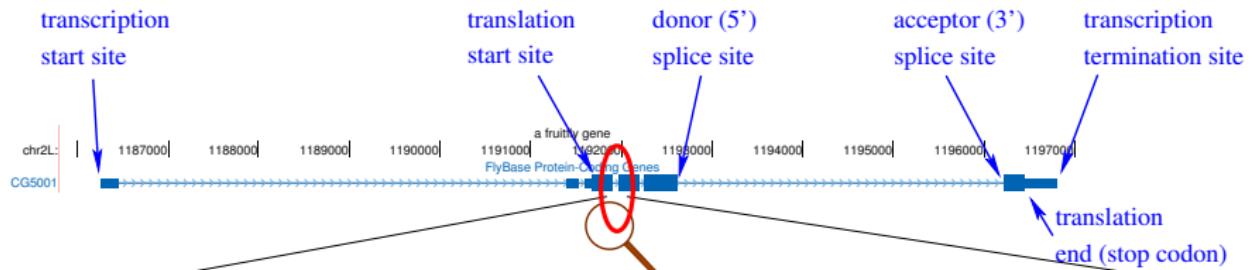
Protein Homology

8 Alternative Splicing

Signals



Signals



← exon

intron →



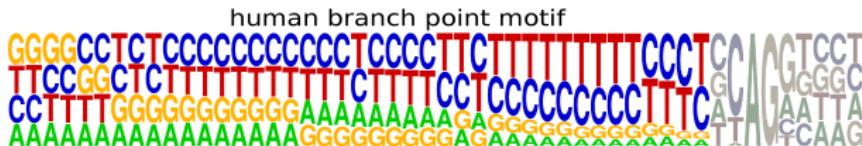
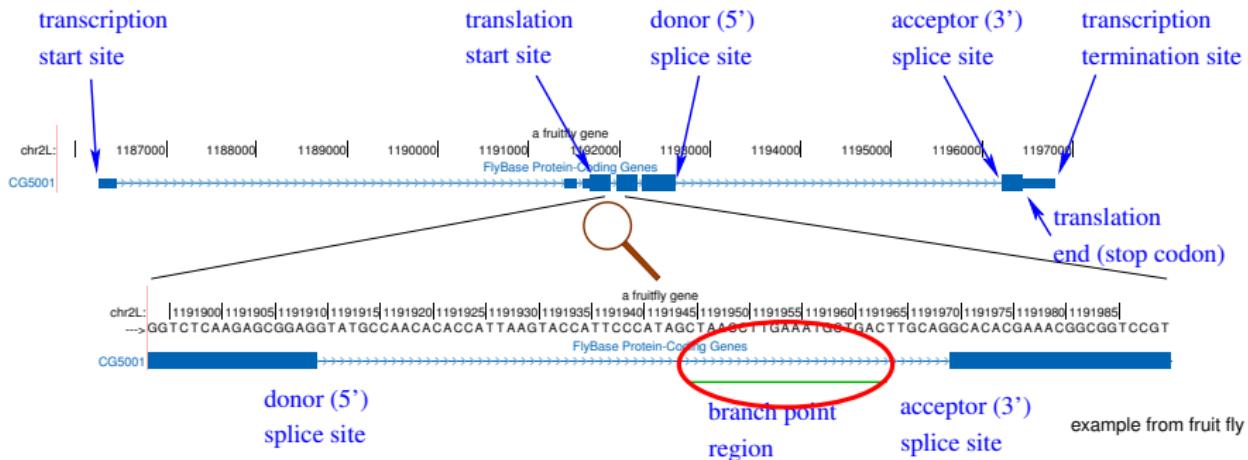
donor splice site (DSS) signal

Frequency of the nucleotides at positions relative to splice site.

acceptor splice site (ASS) signal

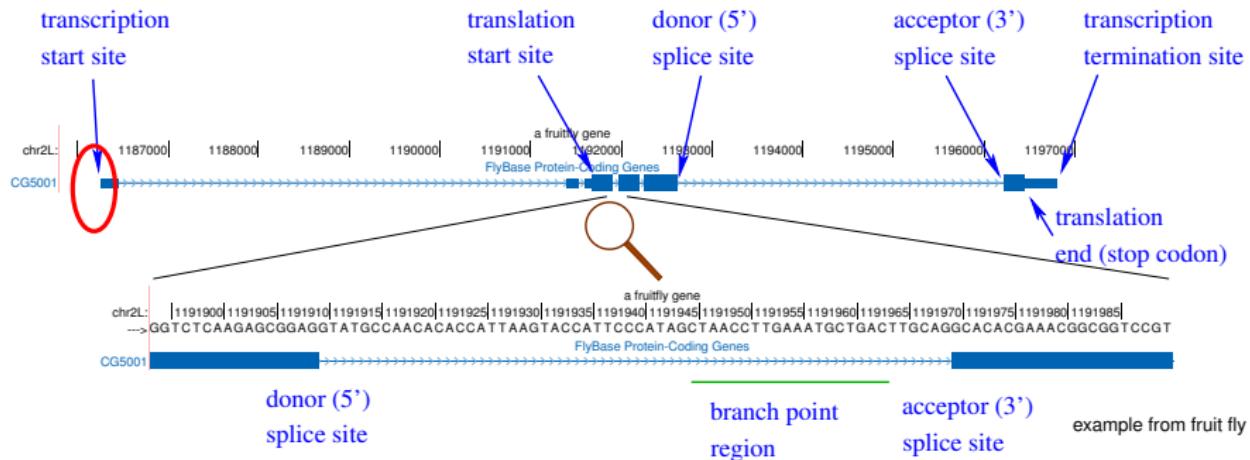
from green algae *Chlamydomonas*

Signals



Branch point: upstream of 3' splice site, a single **conserved adenine** at variable distance to 3' splice site (≈ -30), a splicing complex binds to it, pyrimidine (C,T) rich in human

Signals

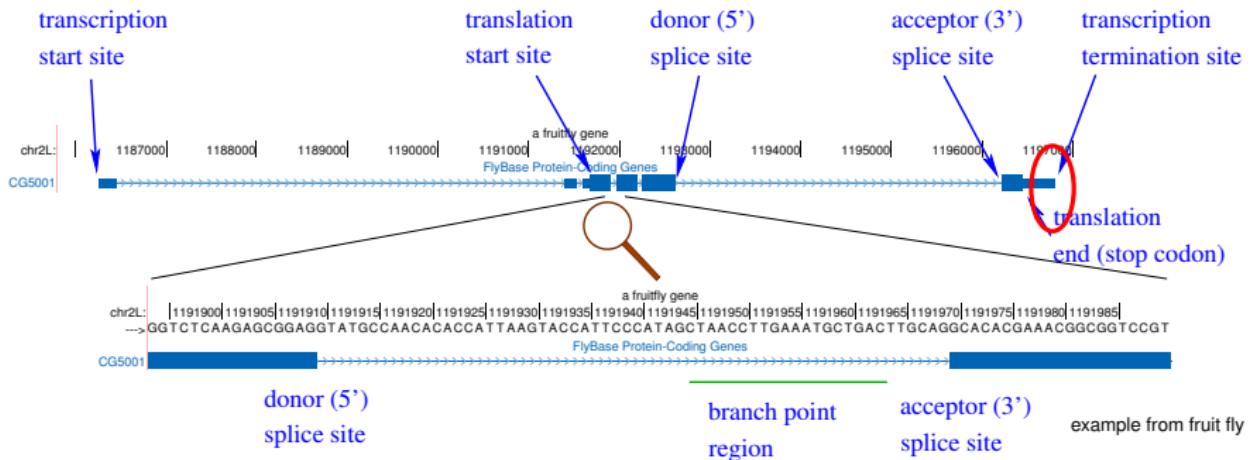


Transcription start site: Transcription from DNA to RNA by RNA polymerase starts here facilitated by **promoter** elements.

Promoter elements are diverse and their profiles tend to contain little info:

- diverse transcription factor binding sites at very variable positions
- sometimes TATA-box
- “CpG islands”

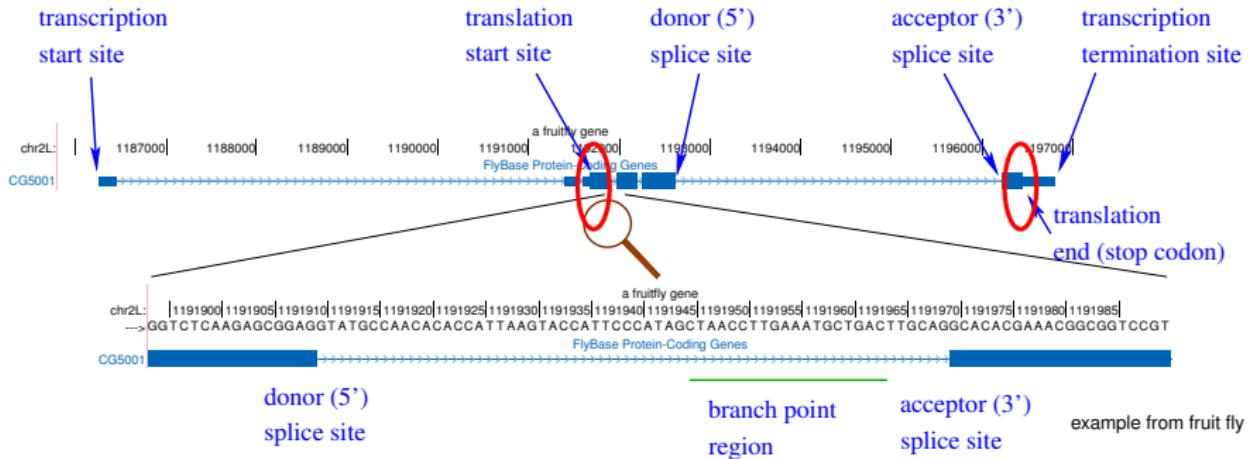
Signals



Transcription termination site (TTS):

- cleavage of the transcript.
- some non-templated A's are appended (polyadenylation).
- polyadenylation is triggered in many species in many genes by the hexamer **aataaa** roughly 15 bp upstream of the TTS.

Signals



Start and stop codon:

- start codon: **ATG**
- stop codons: **TAA, TAG, TGA**

In some species the genetic code is altered and a “stop codon” is actually coding for an amino acid.



Introduction

Gene Finding Problem

Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

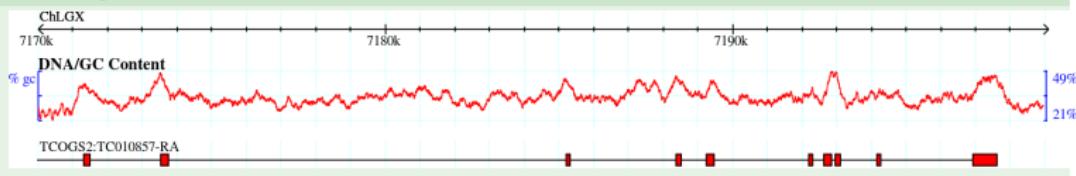
Alternative Splicing

Nucleotide Composition of Coding and Noncoding Regions

Sequence Content

Besides the signals, **position-unspecific** frequencies of **nucleotide patterns** can be used to guess biological classification (e.g. CDS, non-coding, CpG-island) of longer sequence intervals.

Example (GC content in red flour beetle)



Typically, higher order patterns are examined:
E.g. reading-frame dependent k -mer frequencies ($k = 5, 6$) for protein-coding regions.

Remark

Sequence content is usually only **indirect** evidence.



Introduction

Gene Finding Problem
Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

1 Introduction

Gene Finding Problem

Statistical Features of Genes

2 Signal and Content Models

3 Training

4 Gene Finding as Optimization Problem

5 HMMs/CRFs

6 Comparative Gene Finding

7 Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

8 Alternative Splicing



Introduction

Gene Finding Problem
Statistical Features of
Genes

Signal and Content Models

Training

Gene Finding as
Optimization Problem

HMMs/CRFs

Comparative Gene
Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

Signal and Content Models

Scores

- usually, we cannot decide with certainty whether something is an exon, splice site, etc
- ⇒ therefore assign a number (**score s** or probability) to exon candidates, splice site candidates
- aim:
 - the larger the score the more likely there is a true signal
 - the score is small for false signals



Introduction

Gene Finding Problem
Statistical Features of
Genes

Signal and Content Models

Training

Gene Finding as
Optimization Problem

HMMs/CRFs

Comparative Gene
Finding

Evidence Integration

Hints to AUGUSTUS
RNA-Seq
Proteogenomics
Protein Homology

Alternative Splicing

Signal and Content Models

Scores

- usually, we cannot decide with certainty whether something is an exon, splice site, etc
- ⇒ therefore assign a number (**score s** or probability) to exon candidates, splice site candidates
- aim:
 - the larger the score the more likely there is a true signal
 - the score is small for false signals

Signal and Content Models

- **Signal Model:** to score whether a biological signal (e.g. splice site) is present at a certain **position**
- **Content Model:** to score whether a sequence **region** belongs to a given class (e.g. coding)



Introduction

Gene Finding Problem
Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS
RNA-Seq
Proteogenomics
Protein Homology

Alternative Splicing

Example Signal Score

Example (DSS position weight matrix)

Have position specific frequency matrix for DSS

$$m(i, b) \quad (i = 1, 2, \dots, 7, b \in A, C, G, T),$$

$$m(i, A) + m(i, C) + m(i, G) + m(i, T) = 1$$



Have "background" distribution of nucleotides $q(b)$

$$q(A) + q(C) + q(G) + q(T) = 1.$$

Define log-odds score: $s = \log \prod_{i=1}^7 m(i, w_i) / q(w_i)$.

Decision rule

predict splice site : $\Leftrightarrow s \geq \text{threshold}$

guaranteed to have most true predictions given a limit on false positive rate under certain assumptions (Neyman-Pearson).



Introduction

Gene Finding Problem
Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

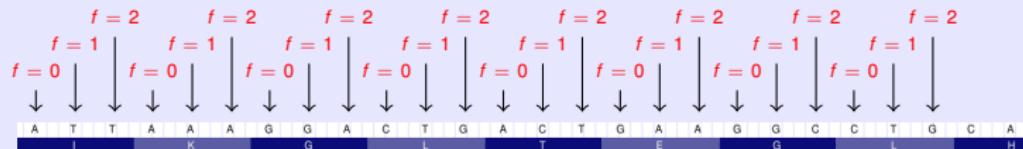
Evidence Integration

Hints to AUGUSTUS
RNA-Seq
Proteogenomics
Protein Homology

Alternative Splicing

Example Content Score

Base composition is frame-dependent



nucleotide frequencies in **human**:

	coding sequence			all f	noncoding sequence
	$f = 0$	$f = 1$	$f = 2$		
A	0.248	0.291	0.146	0.229	0.26
C	0.264	0.243	0.351	0.286	0.24
G	0.321	0.201	0.312	0.278	0.24
T	0.166	0.265	0.190	0.207	0.26



Introduction

Gene Finding Problem

Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

Example Content Score

Example (frame-dependent Markov chain of order k)

Let w be the DNA word of length n to be scored as CDS.
 Let $f \in \{0, 1, 2\}$ be the frame of the first position of w .

$$P(w) := p_f(w_1, \dots, w_k) \cdot \prod_{i=k+1}^n p_{f(i)}(w_i | w_{i-k}, \dots, w_{i-1})$$

- p_f is a start probability for the first k bases

Here: • $f(i) \in \{0, 1, 2\}$ such that $f(i) \equiv f - 1 + i \pmod{3}$
 is the frame of the i -th position of w

Define $s(w) = \log(P(w)/Q(w))$,
 where $Q(w)$ is the probability of w in a “background” model
 (e.g. non-coding).

Remark:

Division by background \Rightarrow good exon candidates get positive score



Introduction

Gene Finding Problem

Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

Example Content Score - Continued

Example

$$w = \text{ATTCTGC}$$

frame $f = 2$, i.e. with these codon breaks: A||TTC||TGC

$$k = 2$$

$$\begin{aligned} P(\text{ATTCTGC}) &= p_2(\text{AT})p_1(\text{T|AT})p_2(\text{C|TT}) \\ &\quad p_0(\text{T|TC})p_1(\text{G|CT})p_2(\text{C|TG}) \end{aligned}$$

Choice of order k

- probabilities $p_r(x | y_1, \dots, y_k)$ can be estimated on known coding sequences
- if $k \geq 2$ above content model can **reflect codon usage**
- usually: the higher k the bigger the difference between coding and noncoding



Introduction

Gene Finding Problem
Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

Example Content Score - Continued

Example

$$w = \text{ATTCTGC}$$

frame $f = 2$, i.e. with these codon breaks: A||TTC||TGC
 $k = 2$

$$\begin{aligned} P(\text{ATTCTGC}) &= p_2(\text{AT})p_1(\text{T|AT})p_2(\text{C|TT}) \\ &\quad p_0(\text{T|TC})p_1(\text{G|CT})p_2(\text{C|TG}) \end{aligned}$$

Choice of order k

- probabilities $p_r(x | y_1, \dots, y_k)$ can be estimated on known coding sequences
- if $k \geq 2$ above content model can **reflect codon usage**
- usually: the higher k the bigger the difference between coding and noncoding
- **but: content model \neq reality** (dead genes)
- typical: $k = 4$ or $k = 5$



Signal and Content Models of AUGUSTUS (without UTR)

Introduction

Gene Finding Problem
Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

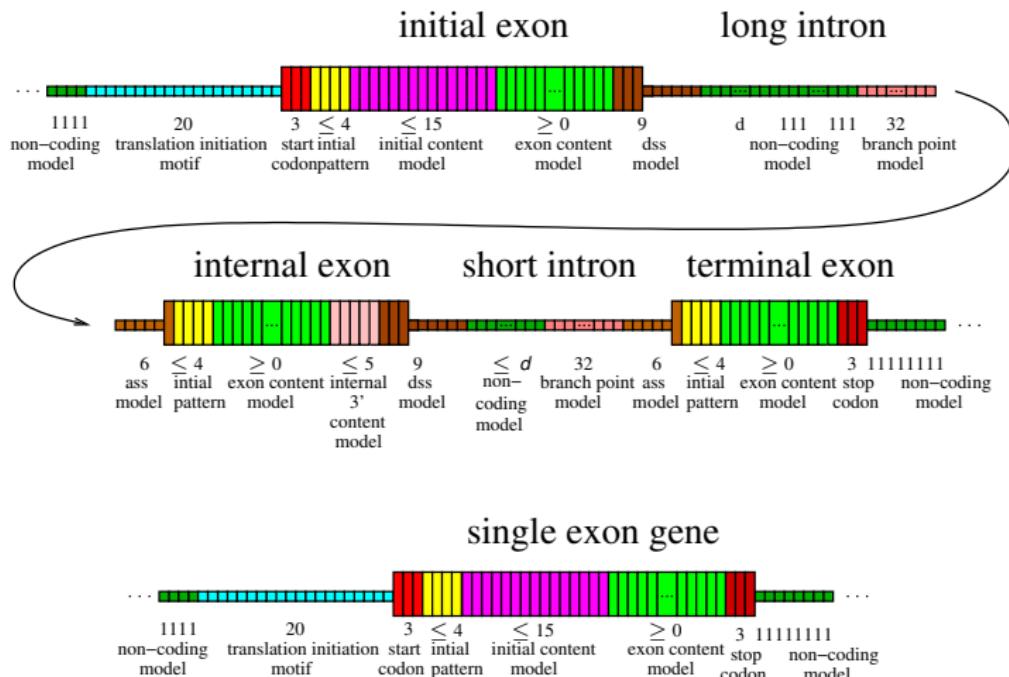
Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing





Introduction

Gene Finding Problem
Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

1 Introduction

Gene Finding Problem

Statistical Features of Genes

2 Signal and Content Models

3 Training

4 Gene Finding as Optimization Problem

5 HMMs/CRFs

6 Comparative Gene Finding

7 Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

8 Alternative Splicing



Introduction

Gene Finding Problem
Statistical Features of
Genes

Signal and Content Models

Training

Gene Finding as
Optimization Problem

HMMs/CRFs

Comparative Gene
Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

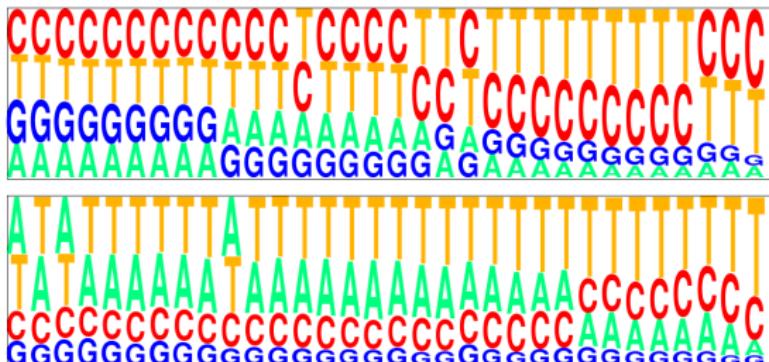
Alternative Splicing

How Species-Specific Must Gene Finding Models Be?

Differences:

- distribution at signals, e.g. branch point region

top: human / bottom: fly





Introduction

Gene Finding Problem

Statistical Features of
Genes

Signal and Content Models

Training

Gene Finding as
Optimization Problem

HMMs/CRFs

Comparative Gene
Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

How Species-Specific Must Gene Finding Models Be?

Differences:

- distribution at signals, e.g. branch point region
- GC content highly variable



Introduction

Gene Finding Problem

Statistical Features of
Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

How Species-Specific Must Gene Finding Models Be?

Differences:

- distribution at signals, e.g. branch point region
- GC content highly variable
- number and length distribution of introns

top: human / bottom: *C. elegans*





Introduction

Gene Finding Problem

Statistical Features of
Genes

Signal and Content Models

Training

Gene Finding as
Optimization Problem

HMMs/CRFs

Comparative Gene
Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

How Species-Specific Must Gene Finding Models Be?

Differences:

- distribution at signals, e.g. branch point region
- GC content highly variable
- number and length distribution of introns
- length distribution of UTRs



How Species-Specific Must Gene Finding Models Be?

Introduction

Gene Finding Problem

Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

Differences:

- distribution at signals, e.g. branch point region
- GC content highly variable
- number and length distribution of introns
- length distribution of UTRs
- gene density



Introduction

Gene Finding Problem

Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

Training: Estimate Species-Specific Parameters

“Training Set”

- input: set of annotated sequences

$$(x^{(k)}, y^{(k)})_{k=1, \dots, N},$$

such that the parse $x^{(k)}$ represents the gene structure of DNA sequence $y^{(k)}$.

- at least a few hundred genes (≥ 200 , more better up to ≈ 1000)
- used to estimate species-specific parameters
- self-training programs: use their own preliminary prediction as training set and iterate prediction and training (GeneMark-ES, Borodovsky et al.)
(does not work in all species)



Introduction

Gene Finding Problem

Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

Training: Estimate Species-Specific Parameters

“Training Set”

- input: set of annotated sequences

$$(x^{(k)}, y^{(k)})_{k=1, \dots, N},$$

such that the parse $x^{(k)}$ represents the gene structure of DNA sequence $y^{(k)}$.

- at least a few hundred genes (≥ 200 , more better up to ≈ 1000)
 - used to estimate species-specific parameters
 - self-training programs: use their own preliminary prediction as training set and iterate prediction and training (GeneMark-ES, Borodovsky et al.)
(does not work in all species)



Introduction

Gene Finding Problem

Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

Training: Estimate Species-Specific Parameters

“Training Set”

- input: set of annotated sequences

$$(x^{(k)}, y^{(k)})_{k=1, \dots, N},$$

such that the parse $x^{(k)}$ represents the gene structure of DNA sequence $y^{(k)}$.

- at least a few hundred genes (≥ 200 , more better up to ≈ 1000)
- used to estimate species-specific parameters
 - self-training programs: use their own preliminary prediction as training set and iterate prediction and training (GeneMark-ES, Borodovsky et al.)
(does not work in all species)



Introduction

Gene Finding Problem

Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

Training: Estimate Species-Specific Parameters

“Training Set”

- input: set of annotated sequences

$$(x^{(k)}, y^{(k)})_{k=1, \dots, N},$$

such that the parse $x^{(k)}$ represents the gene structure of DNA sequence $y^{(k)}$.

- at least a few hundred genes (≥ 200 , more better up to ≈ 1000)
- used to estimate species-specific parameters
- self-training programs: use their own preliminary prediction as training set and iterate prediction and training (GeneMark-ES, Borodovsky et al.)
(does not work in all species)



Introduction

Gene Finding Problem

Statistical Features of
Genes

Signal and Content
Models

Training

Gene Finding as
Optimization Problem

HMMs/CRFs

Comparative Gene
Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

Some Options to Construct a Training Set

Option 1: PASA

-



Introduction

Gene Finding Problem

Statistical Features of
Genes

Signal and Content
Models

Training

Gene Finding as
Optimization Problem

HMMs/CRFs

Comparative Gene
Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

Some Options to Construct a Training Set

Option 2: SCPIO

-



Introduction

Gene Finding Problem

Statistical Features of
Genes

Signal and Content
Models

Training

Gene Finding as
Optimization Problem

HMMs/CRFs

Comparative Gene
Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

Some Options to Construct a Training Set

Option 3: Predict and iterate

-



Introduction

Gene Finding Problem

Statistical Features of
Genes

Signal and Content
Models

Training

Gene Finding as
Optimization Problem

HMMs/CRFs

Comparative Gene
Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

Some Options to Construct a Training Set

Option 4: No retraining

-



Introduction

Gene Finding Problem
Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene
Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

1 Introduction

Gene Finding Problem

Statistical Features of Genes

2 Signal and Content Models

3 Training

4 Gene Finding as Optimization Problem

5 HMMs/CRFs

6 Comparative Gene Finding

7 Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

8 Alternative Splicing



Introduction

Gene Finding Problem
Statistical Features of
Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS
RNA-Seq
Proteogenomics
Protein Homology

Alternative Splicing

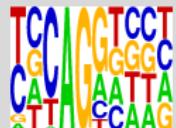
Problems and General Ansatz

Problems

- known signal models do not carry much information
- false positive signals because of low number of true positives
- sequence content can be misleading (pseudogenes, repeats)

Ansatz

- **combine** all individual weak info to boost discriminatory power
- **enforce standard** gene structure:
 - reading frame consistency between exons
 - minimal splice site consensus (GT/AG, maybe GC/AG)
 - no in-frame stop codons
 - minimal intron length (≈ 40 bp)



Introduction

Gene Finding Problem
Statistical Features of
Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

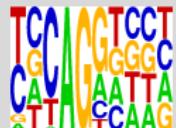
Problems and General Ansatz

Problems

- known signal models do not carry much information
- false positive signals because of low number of true positives
- sequence content can be misleading (pseudogenes, repeats)

Ansatz

- **combine** all individual weak info to boost discriminatory power
- **enforce standard** gene structure:
 - reading frame consistency between exons
 - minimal splice site consensus (GT/AG, maybe GC/AG)
 - no in-frame stop codons
 - minimal intron length (≈ 40 bp)



Introduction

Gene Finding Problem
Statistical Features of
Genes

Signal and Content
Models

Training

Gene Finding as
Optimization Problem

HMMs/CRFs

Comparative Gene
Finding

Evidence Integration

Hints to AUGUSTUS
RNA-Seq
Proteogenomics
Protein Homology

Alternative Splicing

Problems and General Ansatz

Problems

- known signal models do not carry much information
- false positive signals because of low number of true positives
- sequence content can be misleading (pseudogenes, repeats)

Ansatz

- **combine** all individual weak info to boost discriminatory power
- **enforce standard** gene structure:
 - reading frame consistency between exons
 - minimal splice site consensus (GT/AG, maybe GC/AG)
 - no in-frame stop codons
 - minimal intron length (≈ 40 bp)



Introduction

Gene Finding Problem
Statistical Features of
Genes

Signal and Content
Models

Training

Gene Finding as
Optimization Problem

HMMs/CRFs

Comparative Gene
Finding

Evidence Integration

Hints to AUGUSTUS
RNA-Seq
Proteogenomics
Protein Homology

Alternative Splicing

Problems and General Ansatz

Problems

- known signal models do not carry much information
- false positive signals because of low number of true positives
- sequence content can be misleading (pseudogenes, repeats)

Ansatz

- **combine** all individual weak info to boost discriminatory power
- **enforce standard** gene structure:
 - reading frame consistency between exons
 - minimal splice site consensus (GT/AG, maybe GC/AG)
 - no in-frame stop codons
 - minimal intron length (≈ 40 bp)



Introduction

Gene Finding Problem

Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

Problems and General Ansatz

Problems

- known signal models do not carry much information
- false positive signals because of low number of true positives
- sequence content can be misleading (pseudogenes, repeats)

Ansatz

- **combine** all individual weak info to boost discriminatory power
- **enforce standard** gene structure:
 - reading frame consistency between exons
 - minimal splice site consensus (GT/AG, maybe GC/AG)
 - no in-frame stop codons
 - minimal intron length (≈ 40 bp)

Gene Structure Formally

observed variables:

DNA
S

region of Fugu DNA

CTTCCTAACGGTCAAAGTCTGGAGGACAATTGAAACACTGACTGGGGTAAGAAAGCACCGACCCACTTTGGATATTGAAAGGTATAACAAGTTCTGGTCTTCAAGTGCTAACAGGCTCACAAAGTGCAGGGACCCCTCGGTGGGAGATGACCACCAAGGATGCTGTCCA

hints
H

regions with similarity to human protein

conservation of DNA with other fish



Gene Structure Formally

unobserved variables:

DNA
S

CTTCCTAACGGTCAAAGTCTGGAGGACAATTGAAACACTCTGACTGGGGTAAGAAAGCACCGACCTCCAAGGCACTTTGGATATTGAAAGTTATAACAAGTCTGGTCTTCCAAAGTGCTAACAGGCTCACAAAGTGCAGGGACCCCTCGGTGGGAGATGACCACCAAGGATGTCGTC

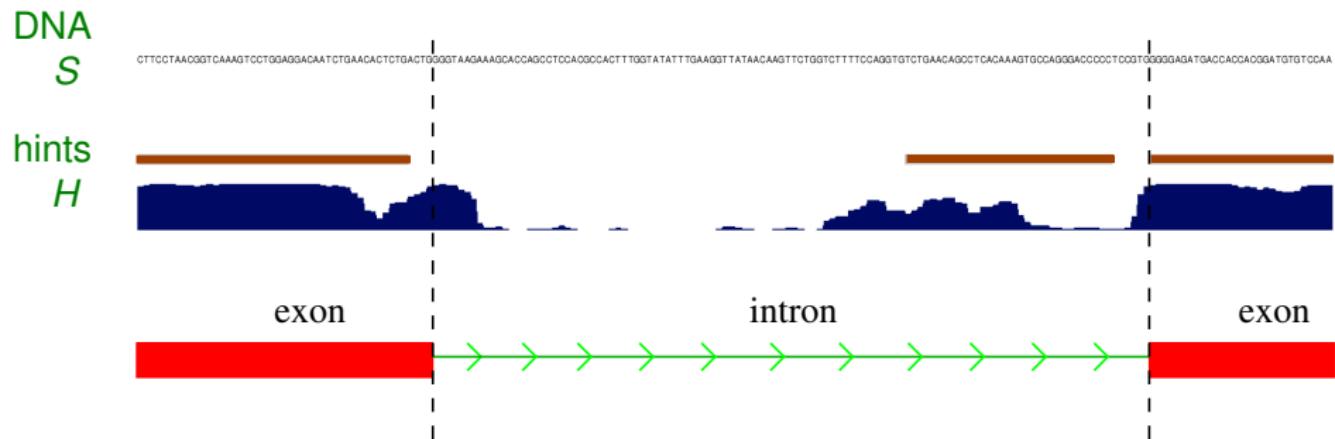
hints
H



exon

Gene Structure Formally

unobserved variables:



Gene Structure Formally

unobserved variables:

DNA
 S

CTTCCTAACGGTCAAAGTCTGGAGGACAATCTGAAACACTCTGACTGGGTAAAGAAAGCACCAGCCCAACGCCACTTGGTATATTGAAGGTATAACAAGTTCTGGCTTTCCAGGTGCTGAACAGGTCACAAAGTGCAGGGACCCCTCGTGAGGGAGATGACCAACGATGTGTCAA

hints
 H



exon

intron

exon



Gene Structure Formally

unobserved variables:

DNA
 S

CTTCCTTAAGGTCAAAGTCTGGAGGACAATCTGAAACACTCTGACTGGGTTAAGAAAGCACCGCCCAAGGCACTTGGTATATTGAAGGTATAACAGTTCTGGCTTTCCAGGTGTCCTGAACAGGCTCACAAAGTGCAAGGGACCCCTCGGTGGGGAGATGACCAACGGATGTGTCCAA

hints
 H



$x_1 = \text{Exon}$

$x_2 = \text{Intron}$

$x_3 = \text{Exon}$

b_1

b_2

gene structure $x = ((x_1, b_1), (x_2, b_2), \dots, (x_n, b_n))$ is a **labelled segmentation** of DNA

x_1, x_2, \dots biological labels, e.g. "exon", "intron", "intergenic region"

b_1, b_2, \dots segment boundaries



Introduction

Gene Finding Problem

Statistical Features of
Genes

Signal and Content
Models

Training

Gene Finding as
Optimization Problem

HMMs/CRFs

Comparative Gene
Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

Generic Approach

Input

sequence: DNA sequence to find genes in

evidence: optional additional evidence, e.g. conservation

Output

gene structure x : start end end positions of exons, genes

Gene finding as optimization problem

maximize

$s(\text{gene structure } x, \text{sequence } S, \text{evidence } H)$

subject to constraint

gene structure x makes biological sense



Introduction

Gene Finding Problem

Statistical Features of
Genes

Signal and Content
Models

Training

Gene Finding as
Optimization Problem

HMMs/CRFs

Comparative Gene
Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

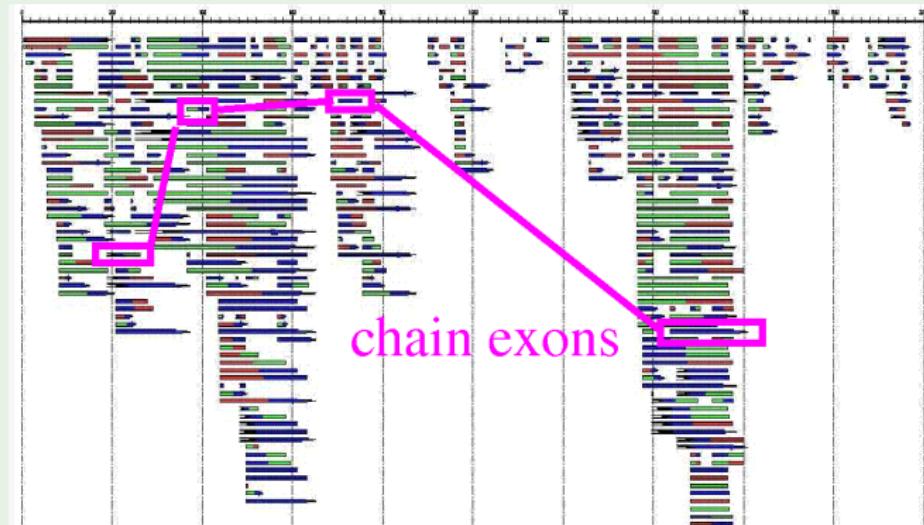
Alternative Splicing

A Hierarchical Approach Using Graphs

Exon Chaining

- find nonoverlapping compatible exon candidates with highest sum of scores
- e.g. done by GENEID (Parra et al., 2000)

Example (exon candidates in a DNA of length 2000)





Introduction

Gene Finding Problem
Statistical Features of
Genes

Signal and Content
Models

Training

Gene Finding as
Optimization Problem

HMMs/CRFs

Comparative Gene
Finding

Evidence Integration

Hints to AUGUSTUS
RNA-Seq
Proteogenomics
Protein Homology

Alternative Splicing

Explicit Probabilistic Models of DNA and Genes

Hidden Markov Model (HMM)

GENSCAN (Burge and Karlin, 1996),
GENIE (Kulp et al., 1996),
GENEMARK (Lukashin and Borodovsky, 1998),
FGENESH (Salamov and Solovyev, 2000),
AUGUSTUS (Stanke and Waack, 2003)

First publication on our gene finder AUGUSTUS

BIOINFORMATICS

Vol. 19 Suppl. 2 2003, pages ii215–ii225
DOI: 10.1093/bioinformatics/btg1080



Gene prediction with a hidden Markov model and a new intron submodel

Mario Stanke^{1,*} and Stephan Waack²

¹Institut für Mikrobiologie und Genetik, Abteilung Bioinformatik, Universität Göttingen, Goldschmidtstraße 1, Göttingen, 37077, Germany and ²Institut für Numerische und Angewandte Mathematik, Universität Göttingen, Lotzestraße 16-18, Göttingen, 37083, Germany

Some Genome Annotation Projects using AUGUSTUS

<i>Aedes aegypti</i>	yellow fever mosquito: dengue fever	<i>Science</i> , 2007
<i>Brugia malayi</i>	parasitic worm, causes elephantiasis	<i>Science</i> , 2007
<i>Tribolium castaneum</i>	red flour beetle, pest and model organism	<i>Nature</i> , 2008
<i>Schistosoma mansoni</i>	parasite causing bilharziosis	<i>Nature</i> , 2009
<i>Acyrtosiphon pisum</i>	pea aphid, agricultural pest	<i>PLOS Biology</i> , 2010
<i>Coprinus cinereus</i>	fungus	<i>PNAS</i> , 2010
<i>Nasonia vitripennis</i>	wasp	<i>Science</i> , 2010
<i>Amphimedon queenslandica</i>	sponge	<i>Nature</i> , 2010
<i>Culex pipiens</i>	common mosquito	<i>Science</i> , 2010
<i>Ricinus communis</i>	castor bean	<i>Nature Biotechnology</i> , 2010
<i>Chlamydomonas reinhardtii</i>	green algae	<i>Proteomics</i> , 2011
<i>Galdieria sulphuraria</i>	red algae	<i>Planta</i> , 2007
<i>Arabidopsis thaliana</i>	plant model organism	<i>PNAS</i> , 2008
<i>Leishmania tarentolae</i>	lizard parasite	<i>Nuc. Acids Res.</i> , 2011
<i>Verticillium longisporum</i>	fungus infecting oilseed rape	
<i>Tetrahymena thermophila</i>	model organisms in biomedical research	
<i>Nicotiana tabaccum</i>	tobacco	
<i>Toxoplasma gondii</i>	causes toxoplasmosis	
<i>Zea mays</i>	maize	
<i>Heliconius melpomene</i>	butterfly	
<i>Apis mellifera</i>	honey bee	

AUGUSTUS is open source.



Hidden Markov Model

Introduction

Gene Finding Problem

Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

HMM

$$P(x, y) = \prod_{i=1}^n P(x_i | x_{i-1}) \cdot P(y(b_{i-1}, b_i) | x_i)$$

is a **generative model**: everything that is used in prediction (y), must have a distribution.

E.g. in comparative gene finding: Must define probability of any observable alignment column.



Introduction

Gene Finding Problem

Statistical Features of
Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene
Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

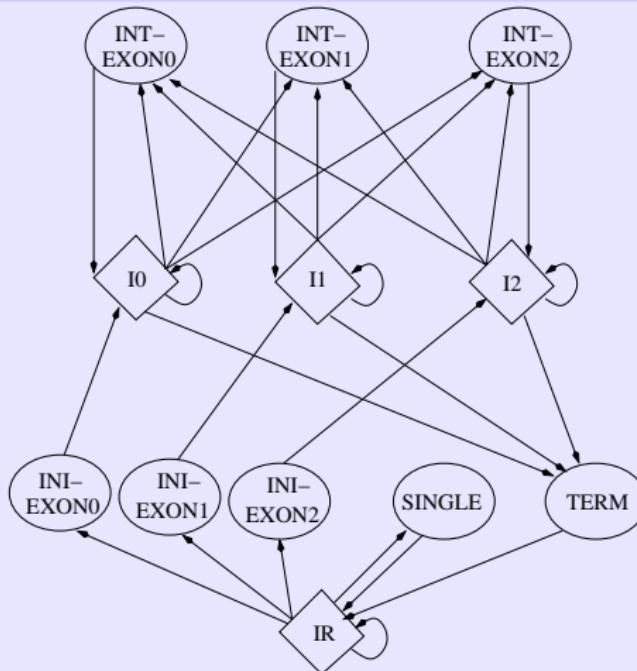
Proteogenomics

Protein Homology

Alternative Splicing

A Simple (Semi-Markov) HMM for Gene Finding: Model Topology

Model for (multiple) eukaryotic genes on forward strand



Transition Matrix of AUGUSTUS (without UTRs)

Mario Stanke



Introduction

Gene Finding Problem

Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

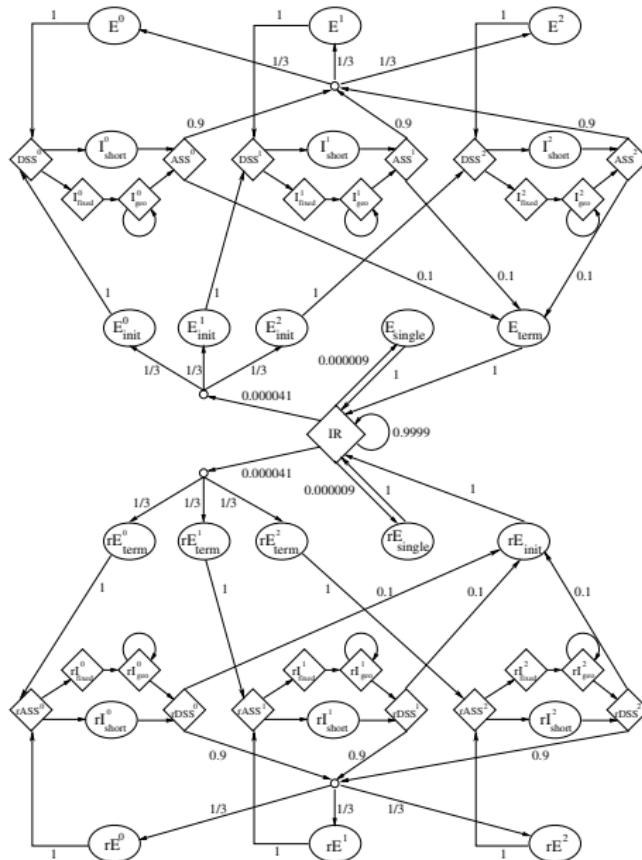
Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing





Discriminative Models of Genes given DNA

Introduction

Gene Finding Problem

Statistical Features of
Genes

Signal and Content
Models

Training

Gene Finding as
Optimization Problem

HMMs/CRFs

Comparative Gene
Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

Conditional Random Field (CRF)

CRAIG (Bernal et al., 2007),
CONTRAST (Gross et al., 2007),
AUGUSTUS (unpublished)



Introduction

Gene Finding Problem
Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

1 Introduction

Gene Finding Problem

Statistical Features of Genes

2 Signal and Content Models

3 Training

4 Gene Finding as Optimization Problem

5 HMMs/CRFs

6 Comparative Gene Finding

7 Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

8 Alternative Splicing



Introduction

Gene Finding Problem
Statistical Features of
Genes

Signal and Content
Models

Training

Gene Finding as
Optimization Problem

HMMs/CRFs

Comparative Gene
Finding

Evidence Integration

Hints to AUGUSTUS
RNA-Seq
Proteogenomics
Protein Homology
Alternative Splicing

Conditional Random Field for Gene Prediction

Definition (CRF)

A (Semi-)Conditional-Random-Field is a conditional distribution of gene structures x given the observation y of the following form

$$P(x | y) = \frac{1}{Z(y)} \exp \sum_{i=1}^n f(x_{i-1}, x_i, b_{i-1}, b_i, y).$$

CRF

- gene structure is scored locally:
score may only locally depend on one labelled segment
- not scorable, e.g.: gene should have a total exon length of 1000bp



Introduction

Gene Finding Problem

Statistical Features of
Genes

Signal and Content
Models

Training

Gene Finding as
Optimization Problem

HMMs/CRFs

Comparative Gene
Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

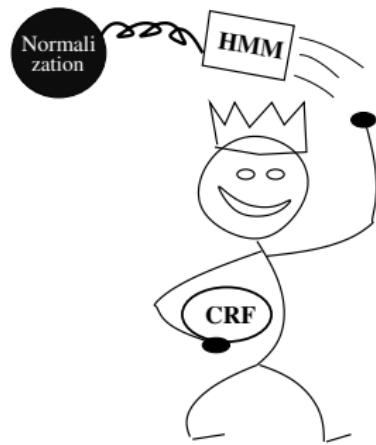
Protein Homology

Alternative Splicing

Conditional Random Field versus Hidden Markov Model

CRF

- for every HMM there is a CRF with same $P(x|y)$
- only *globally* normalization necessary
- can easier be **discriminatively** trained

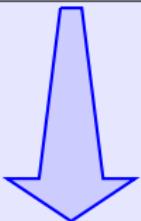


Motivation for Model Change

Signal models: disadvantage of HMM parameter estimation

Call splice site pattern **strong** when relative frequency in real genes large.

relative frequency in range	strength	counts of correct DSS	HMM: counts of predicted DSS
$[0, 7.72 \cdot 10^{-6})$		0	108
$[7.72 \cdot 10^{-6}, 0.000139)$		57	83
$[0.000139, 0.000275)$		175	183
$[0.000275, 0.00071)$		532	465
$[0.00071, 1]$		3053	2186
sum		3817	3025



⇒ Predicted gene structures have **fewer strong** splice sites than correct gene structures and **more weak** splice sites.

One should be able to improve predictions by putting more weight on the splice site model than classical HMMs do.

[Introduction](#)[Gene Finding Problem](#)[Statistical Features of Genes](#)[Signal and Content Models](#)[Training](#)[Gene Finding as Optimization Problem](#)[HMMs/CRFs](#)[Comparative Gene Finding](#)[Evidence Integration](#)[Hints to AUGUSTUS](#)[RNA-Seq](#)[Proteogenomics](#)[Protein Homology](#)[Alternative Splicing](#)

Conditional Random Field

Online large-margin training

Iteratively update parameters θ (human: 8272 parameters) using a training set of gene structures g_{train} (200-2000 genes) so that

$$P_\theta(g_{\text{train}} \mid \text{DNA}) > P_\theta(g \mid \text{DNA}) + \underbrace{L(g_{\text{train}}, g)}_{\geq 0} \text{ punishes important errors}$$

for all **false gene structures** $g \neq g_{\text{train}}$.

Conditional Random Field versus GHMM training

Accuracy improvement due to Semi-CRF training

ab initio accuracy	human		<i>Drosophila</i>		<i>Chlamydomonas</i>	
	GHMM	Semi-CRF	GHMM	Semi-CRF	GHMM	Semi-CRF
base sn	.754	.817	.907	.960	.889	.986
base sp	.659	.782	.870	.886	.994	.977
exon sn	.601	.731	.807	.811	.875	.958
exon sp	.567	.712	.743	.771	.918	.952

(sn = sensitivity = true positives / actual positives, sp = specificity = true positives / predicted positives)



Introduction

Gene Finding Problem
Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

1 Introduction

Gene Finding Problem

Statistical Features of Genes

2 Signal and Content Models

3 Training

4 Gene Finding as Optimization Problem

5 HMMs/CRFs

6 Comparative Gene Finding

7 Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

8 Alternative Splicing



Introduction

Gene Finding Problem
Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

Conservation of Gene Structure and Sequence

Observation

Protein sequences and **rough structure** of genes are often **conserved** between species that are tens of millions of years separated.

Example (Human-Mouse: 75 million years)

- 95% of orthologous gene pairs have same number of exons in human and mouse
- coding sequence to $\approx 85\%$ identical

Human: GCAC TTTCTTAAGGAAAGTAATGCCAGTGAA GTGTGGGGAGCATTAAGGACTGACTGAAGGCCCTGCATGGATCCATGTTCA TGAGTTGGAGATAATACAGCAAGTGGGTG
 Mouse: GTATTTCTTAAAGGAAAGTAATGCCAGTGAA GTGTGGGGAGCATTAAGGACTGACTGAAGGCCCTGCATGGATCCATGTTCA TGAGTTGGAGATAATACAGCAAGTGGGTG
 SOD1
 Gap:

- noncoding sequence to $\approx 35\%$ identical

Human: AGTGTGGAAACAAGATTACCATCTCCCTTTGAGGACACAGGCTAGAGCAAGTTAACAGCCTTGCTGGAGGTTCACTGGCTAGAAAGTGGTCAGCCTGGATTGGACACAGATTTCC
 Mouse: AGTGTAGGAA...GTT...TGGAGAGAGGCT...AGAGCTAGCST...CTCCAGGCCAC...CCTGAGGAGTGGCTAAGATGAAACATAGGTTTCT
 SOD1
 Gap:



Introduction

Gene Finding Problem
Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

Conservation of Gene Structure and Sequence

Observation

Protein sequences and **rough structure** of genes are often **conserved** between species that are tens of millions of years separated.

Example (Human-Mouse: 75 million years)

- 95% of orthologous gene pairs have same number of exons in human and mouse
- coding sequence to $\approx 85\%$ identical

→GCAC TTTCTTAAGGAAAGTAAATGCCAGTGAA GGTTGGGGAGCATTAAGGACTGACTGAAGGCCCTGATGGATTCCATGTTCA TGAGTTGGAGATAATACAGCAAGTGGGTG
SOD1 Human: GCAC TTTCTTAAGGAAAGTAAATGCCAGTGAA GGTTGGGGAGCATTAAGGACTGACTGAAGGCCCTGATGGATTCCATGTTCA TGAGTTGGAGATAATACAGCAAGTGGGTG
Mouse: GTATTTTTCAGGCAAGCGGTGAAACAGTTGTGTGAGGAAATTACAGGTTAACGAGCCAGATGGATTCCATGTTCA TGAGTTGGAGATAATACAGCAAGTGGGTG
→ESNGPVKWGSISIKGLTEGLHGFHVHEFODNTA
→SOD1 Human: GAGCTGGAAACAAGTACCATCTCCCTTTGAGGACACAGGCTAGAGCAAGCTTGAGGAGTTCACTGGCTAGAAAGTGTCA GCGCTGGATTGGACACAGATTTTCC
Mouse: AGTGTGAGGAA...GTG...TGTGAGGAGGGCT...AGAGCTAGCTTGCTGGAGGTTCACTGGCTAGAAAGTGTCA GCGCTGGATTGGACACAGATTTTCC
→SOD1 Human: AGTGTGAGGAAACAAGTACCATCTCCCTTTGAGGACACAGGCTAGAGCAAGCTTGCTGGAGGTTCACTGGCTAGAAAGTGTCA GCGCTGGATTGGACACAGATTTTCC
Mouse: AGTGTGAGGAA...GTG...TGTGAGGAGGGCT...AGAGCTAGCTTGCTGGAGGTTCACTGGCTAGAAAGTGTCA GCGCTGGATTGGACACAGATTTTCC

- noncoding sequence to $\approx 35\%$ identical

chr21: 33035800| 33035900| 33035910| 33035920| 33035930| 33035940| 33035950| 33035960| 33035970| 33035980| 33035990|
→AGTGTGAGGAAACAAGTACCATCTCCCTTTGAGGACACAGGCTAGAGCAAGCTTGCTGGAGGTTCACTGGCTAGAAAGTGTCA GCGCTGGATTGGACACAGATTTTCC
SOD1 Human: AGTGTGAGGAAACAAGTACCATCTCCCTTTGAGGACACAGGCTAGAGCAAGCTTGCTGGAGGTTCACTGGCTAGAAAGTGTCA GCGCTGGATTGGACACAGATTTTCC
Mouse: AGTGTGAGGAA...GTG...TGTGAGGAGGGCT...AGAGCTAGCTTGCTGGAGGTTCACTGGCTAGAAAGTGTCA GCGCTGGATTGGACACAGATTTTCC
→SOD1 Human: AGTGTGAGGAAACAAGTACCATCTCCCTTTGAGGACACAGGCTAGAGCAAGCTTGCTGGAGGTTCACTGGCTAGAAAGTGTCA GCGCTGGATTGGACACAGATTTTCC
Mouse: AGTGTGAGGAA...GTG...TGTGAGGAGGGCT...AGAGCTAGCTTGCTGGAGGTTCACTGGCTAGAAAGTGTCA GCGCTGGATTGGACACAGATTTTCC



Introduction

Gene Finding Problem
Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

Conservation of Gene Structure and Sequence

Observation

Protein sequences and **rough structure** of genes are often **conserved** between species that are tens of millions of years separated.

Example (Human-Mouse: 75 million years)

- 95% of orthologous gene pairs have same number of exons in human and mouse
- coding sequence to $\approx 85\%$ identical

Human: GCAC TTTCTTAAGGAAAGTAATGCCAGTGAA GGTTGGGGAGCATTAAGGACTGACTGAAGGCCCTGCATGGATTCATGTTCA TGAGTTGGAGATAATACAGCAAGTGGGTG
Mouse: GTATTTCTTAAGGAAAGTAATGCCAGTGAA GGTTGGGGAGCATTAAGGACTGACTGAAGGCCCTGCATGGATTCATGTTCA TGAGTTGGAGATAATACAGCAAGTGGGTG

- noncoding sequence to $\approx 35\%$ identical

Human: AGTGTGGAAACAAGATTACCATCTCCCTTTGAGGACACAGGCTAGAGCAGTTAACGAGCTTGCTGGAGGTTCACTGGCTAGAAAGTGTCA GCGCTGGATTGGACACAGATTTCC
Mouse: AGTGTGGAAACAAGATTACCATCTCCCTTTGAGGACACAGGCTAGAGCAGTTAACGAGCTTGCTGGAGGTTCACTGGCTAGAAAGTGTCA GCGCTGGATTGGACACAGATTTCC



Introduction

Gene Finding Problem
Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

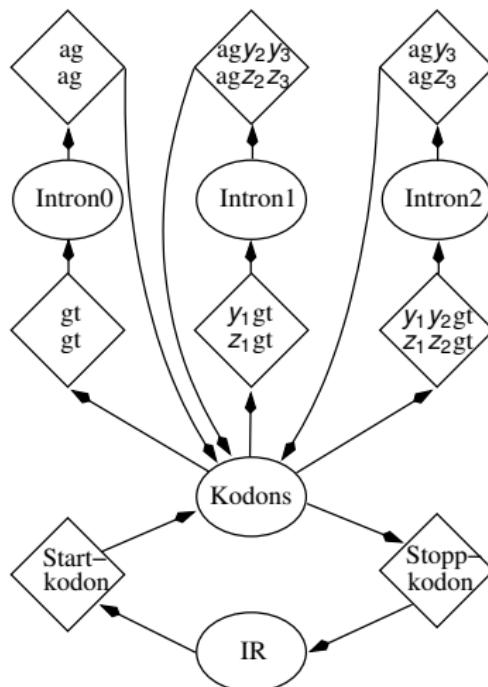
Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS
RNA-Seq
Proteogenomics
Protein Homology

Alternative Splicing

Pair HMMs for Eukaryotic Gene Finding



- DOUBLESCAN
(Meyer and Durbin, 2002)
- does alignment of two genomics sequences and gene prediction in both at same time
- running time quadratic in sequence length
- does not generalize efficiently to >2 species
- not popular



Introduction

Gene Finding Problem

Statistical Features of
GenesSignal and Content
Models

Training

Gene Finding as
Optimization Problem

HMMs/CRFs

Comparative Gene
Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

CONTRAST (Gross et al., 2007)

CONTRAST

- still best *de novo* predictor for human (?)
- discriminatory model (SVM+CRF)
- uses 11 other vertebrate genomes as informants
- roughly for 50% of human genes at least one splice form is **correctly predicted**

human	ACAGGTGAGGAGGCG
macaque
mouse
rat	ACAGGTGAGAAAG..
rabbit
dog	ACAGGTGAGGAGTCG
cow	ACAGGTGAGCAGTCG
armadillo	ACAGGTGAGGAG_CA
elephant
tenrec
opossum	CCAGGGAAG.....
chicken	CCAGGTGA.....
EST	SSSSIIIIIIIIII



Introduction

Gene Finding Problem
Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

1 Introduction

Gene Finding Problem

Statistical Features of Genes

2 Signal and Content Models

3 Training

4 Gene Finding as Optimization Problem

5 HMMs/CRFs

6 Comparative Gene Finding

7 Evidence Integration

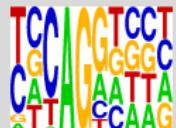
Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

8 Alternative Splicing



Introduction

Gene Finding Problem

Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

Using Extrinsic Evidence to Improve the Accuracy of AUGUSTUS

Extrinsic evidence

additional information to aid gene prediction, not contained in the sequence

AUGUSTUS currently can use

- mRNA/EST/454 alignments, RNA-Seq
- genomic conservation
- annotation of related species using a syntenic alignment
- protein alignments, peptide mass spectrometry

Hint

piece of additional local information about gene structure, in general uncertain



Types of Hints

Introduction

Gene Finding Problem
Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq
Proteogenomics
Protein Homology

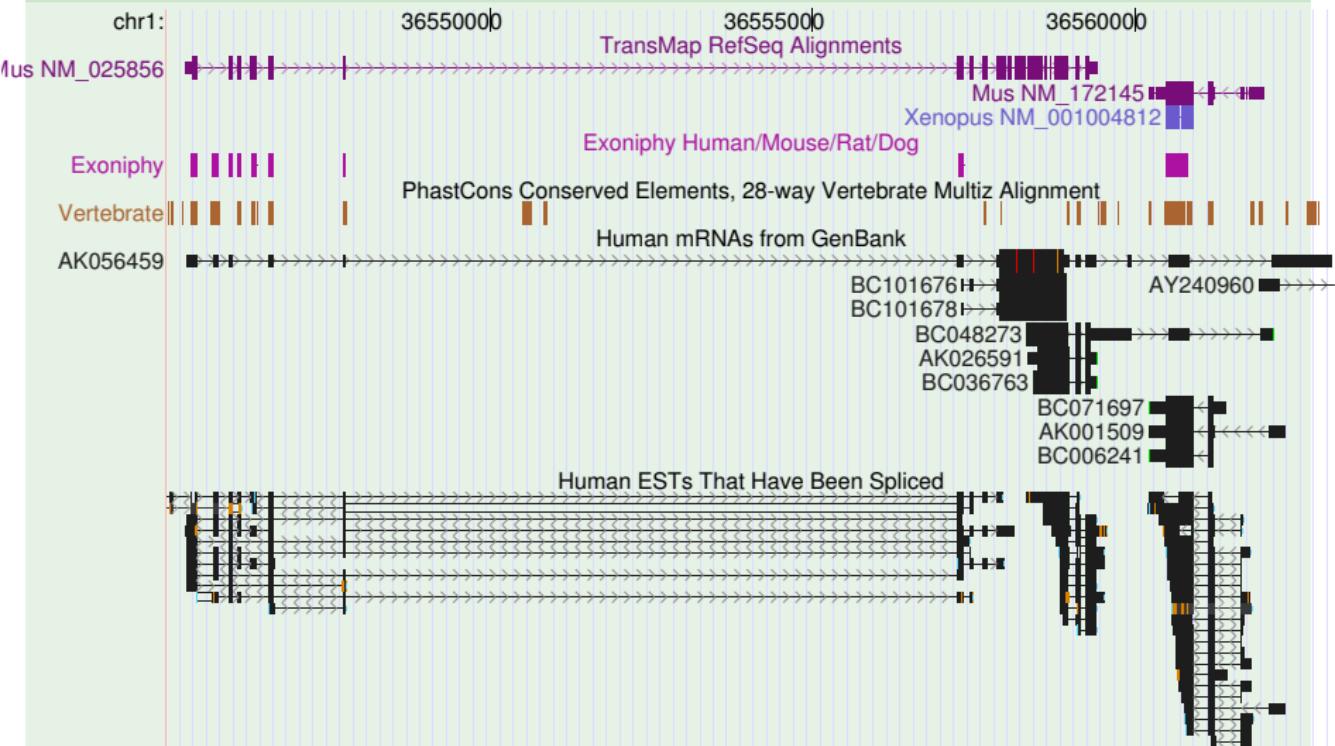
Alternative Splicing

AUGUSTUS understands these hints...

type	meaning	source of information
start	start codon	protein, TransMap
stop	stop codon	protein, TransMap
tss	transcription start site	TransMap, CpG islands, promoter predictor
tts	transcription termination site	TransMap, 3' ESTs
ass	acceptor splice site	cDNA, proteins, TransMap, conservation
dss	donor splice site	cDNA, proteins, TransMap, conservation
exonpart	part of biological exon	cDNA
exon	biological exon	cDNA
intronpart	part of intron	TransMap
intron	intron	cDNA, proteins, TransMap
CDSpart	part of coding part of exon	protein, TransMap, conservation
CDS	coding part of exon	protein, conservation
UTRpart	part of non-coding part of exon	TransMap
UTR	non-coding part of exon	
irpart	part of intergenic region	mRNA, TransMap, cDNA
nonexonpart	intron or intergenic region	retro genes, repeats

Sources of Hints - Example

Example



Integration of Hints into the Model

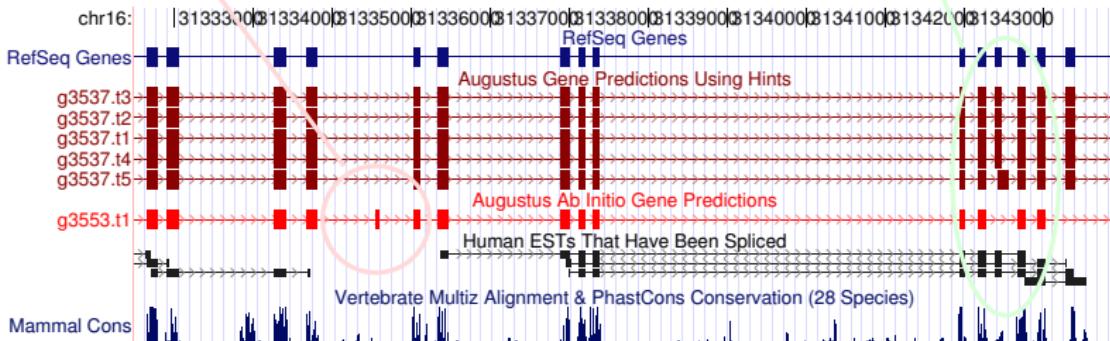
Hints modify the joint probability

$$P(\text{gene structure}, \text{DNA, hints})$$

by a factor that depends on the compatibility of the gene structure with the hints.

Hints have two effects

- ① gene structure candidates that **obey** a hint get a relative **bonus**.
- ② gene structure candidates with **unsupported** exons, introns or signals get a relative **malus** (penalty)



Integration of Hints into the Model

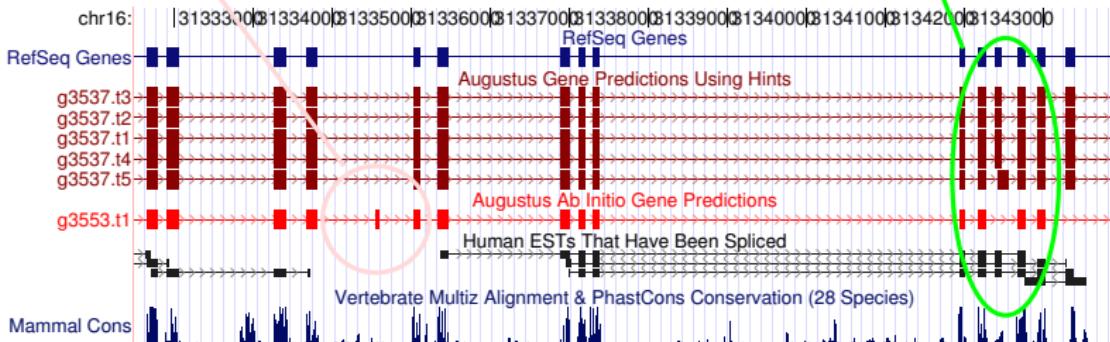
Hints modify the joint probability

$$P(\text{gene structure}, \text{DNA, hints})$$

by a factor that depends on the compatibility of the gene structure with the hints.

Hints have two effects

- ① gene structure candidates that **obey** a hint get a relative **bonus**.
- ② gene structure candidates with **unsupported** exons, introns or signals get a relative **malus** (penalty)



Integration of Hints into the Model

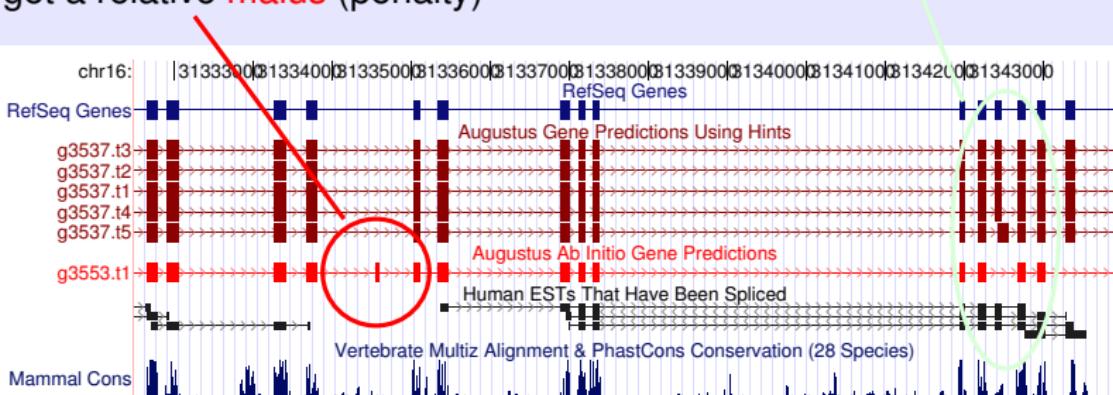
Hints modify the joint probability

$$P(\text{gene structure}, \text{DNA, hints})$$

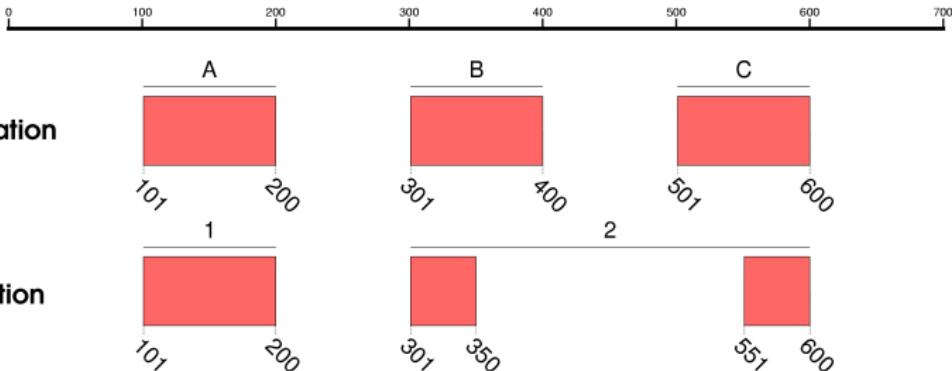
by a factor that depends on the compatibility of the gene structure with the hints.

Hints have two effects

- ① gene structure candidates that **obey** a hint get a relative **bonus**.
- ② gene structure candidates with **unsupported** exons, introns or signals get a relative **malus** (penalty)



Prediction Evaluation Measures



$$\text{Sensitivity sn} = \frac{\text{true positives}}{\text{annotated positives}}$$

$$\text{Specificity sp} = \frac{\text{true positives}}{\text{predicted positives}}$$

Gene sn	33.3%
Gene sp	50%
Transcript sn	33.3%
Transcript sp	50%
Exon sn	33.3%
Exon sp	33.3%
Base sn	66.7%
Base sp	100%

Accuracy of Genome Annotation

Two independent published gene prediction assessments since 2000:

- EGASP on **human** ENCODE regions, 2005:

genes correctly predicted:

ab initio	24% (AUGUSTUS)
using any available info	73% (JIGSAW)

(Guigo et al., *Genome Biology*, 2006)

- nGASP on ***C. elegans***, 2007:

genes correctly predicted:

ab initio	61% (AUGUSTUS, mGENE, FGENESH)
using ESTs & proteins	80% (AUGUSTUS, FGENESH, mGENE)
best-performing combiner	80% (JIGSAW)
for comparison:	65% (GLEAN)

(Coghlan et al., *BMC Bioinformatics*, 2008)

Note: Most species have less available evidence than these two!



Introduction

Gene Finding Problem
Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

1 Introduction

Gene Finding Problem

Statistical Features of Genes

2 Signal and Content Models

3 Training

4 Gene Finding as Optimization Problem

5 HMMs/CRFs

6 Comparative Gene Finding

7 Evidence Integration

Hints to AUGUSTUS

RNA-Seq

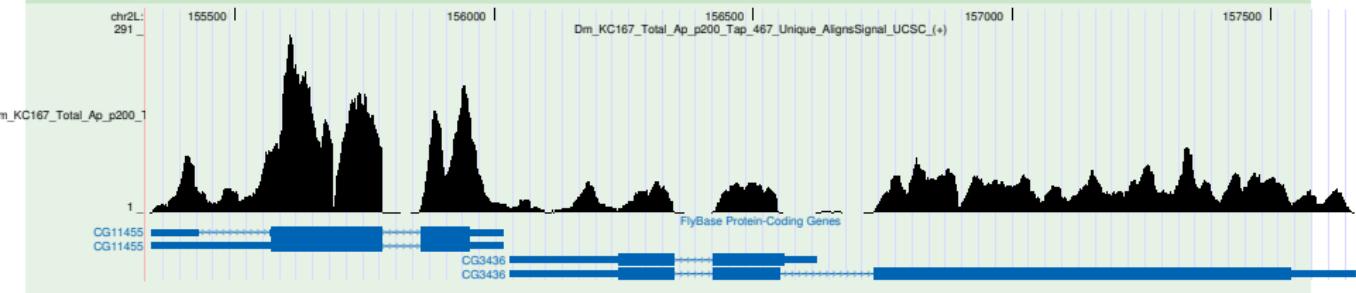
Proteogenomics

Protein Homology

8 Alternative Splicing

Next Generation Sequencing of Transcriptomes (RNA-Seq)

Example (coverage per position in genome)



Chances and Challenges

- genes expressed in sample have usually good coverage
- intuitively, should give boost to gene prediction accuracy
- but how?



Introduction

Gene Finding Problem
Statistical Features of
Genes

Signal and Content Models

Training

Gene Finding as
Optimization Problem

HMMs/CRFs

Comparative Gene
Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics
Protein Homology

Alternative Splicing

RGASP: RNA-Seq Genome Annotation Assessment Project

RGASP

- assess the current progress of automatic gene building using RNAseq
- part of ENCODE project,
<http://www.gencodegenes.org/rgasp/>
- 17 participating groups submitting 2 years ago,
all on same data:
 - on human (41Gb), fly (13Gb), worm (21Gb)
 - Illumina, SOLiD, Helicos
 - paired, unpaired reads
 - unstranded, stranded reads
- accuracy results on site <https://compgen.bio.ub.es> (Josep Abril)



Introduction

Gene Finding Problem

Statistical Features of
Genes

Signal and Content
Models

Training

Gene Finding as
Optimization Problem

HMMs/CRFs

Comparative Gene
Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

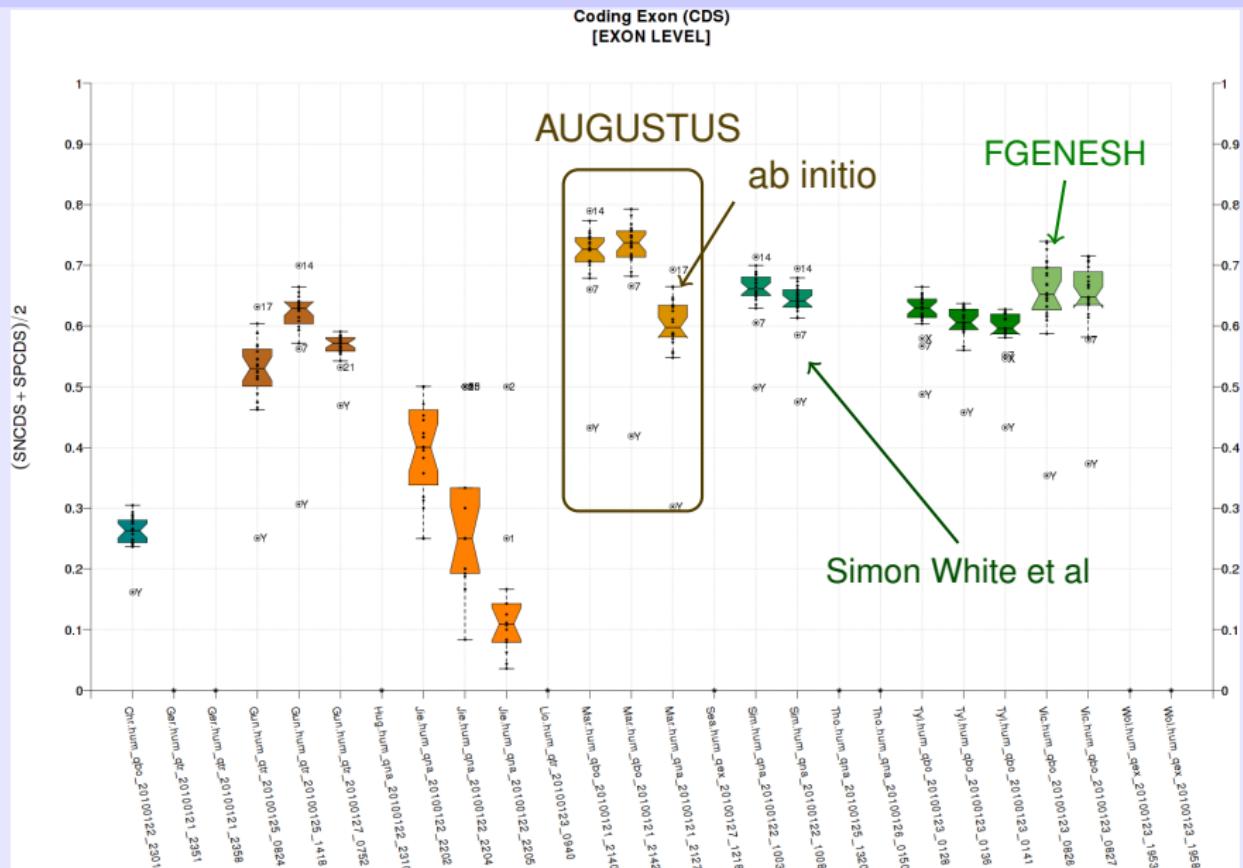
RGASP: RNA-Seq Genome Annotation Assessment Project

Three major approaches to transcript identification with RNA-Seq

- evidence integration into gene finder
(e.g. AUGUSTUS, FGENESH, mGENE, GENEID)
 - ① align reads to genome first
 - ② integrate evidence from **coverage** and **spliced alignments** into gene finder
- purely alignment-based (e.g. Cufflinks)
 - ① align reads to genome first
 - ② construct transcripts from spliced alignments (no gene finding)
- de novo assembly of reads
 - ① assemble transcriptome reads into transcript contigs
 - ② use contigs for gene finding or just align them

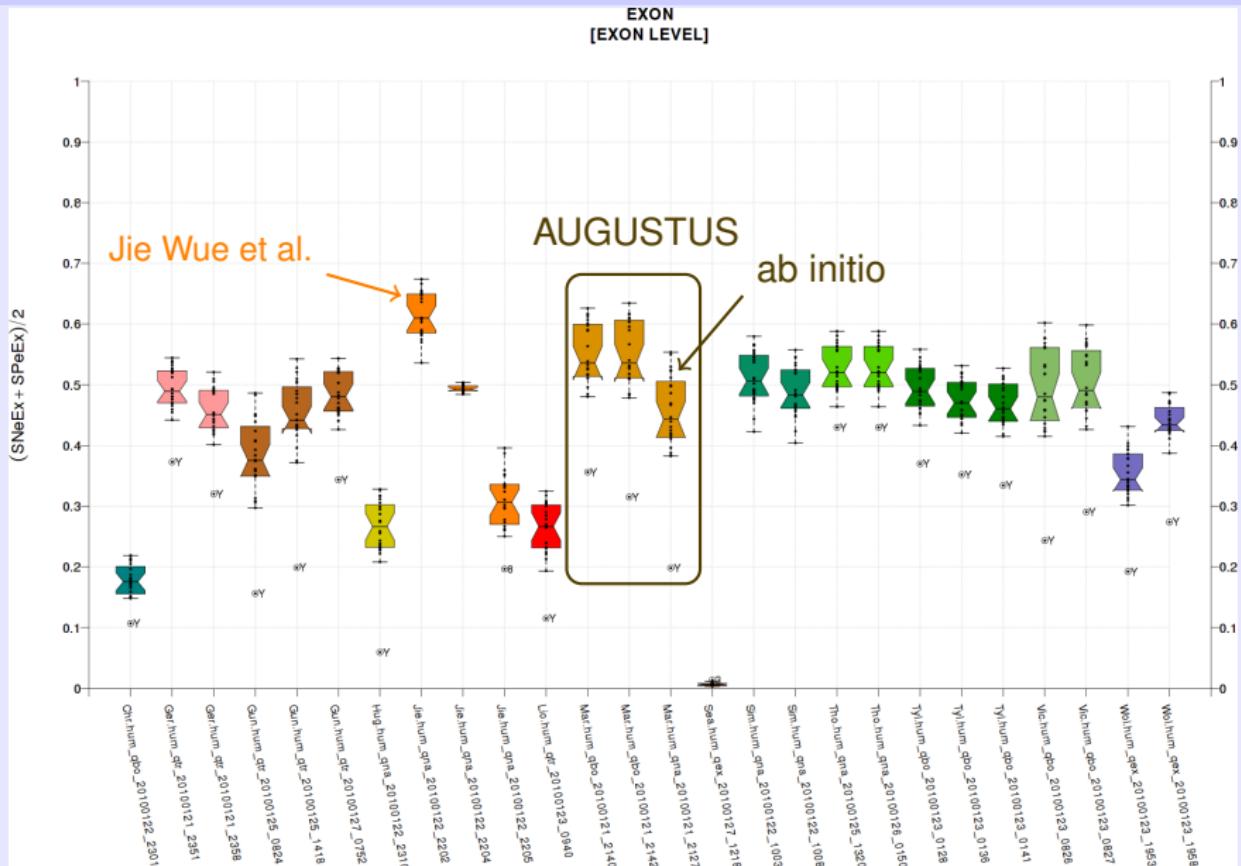
RGASP: results by submissions

human, coding exon level



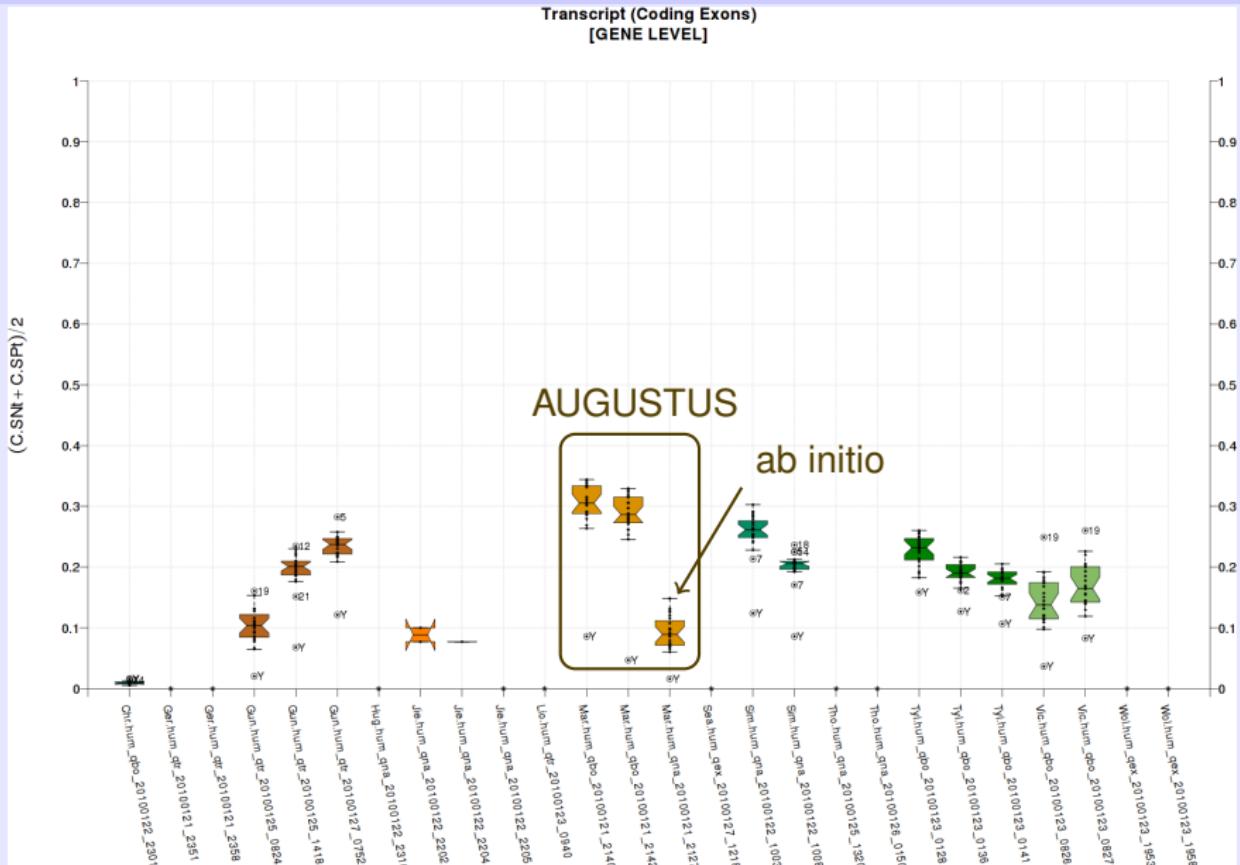
RGASP: results by submissions

human, exon level (no coding sequence predicted)



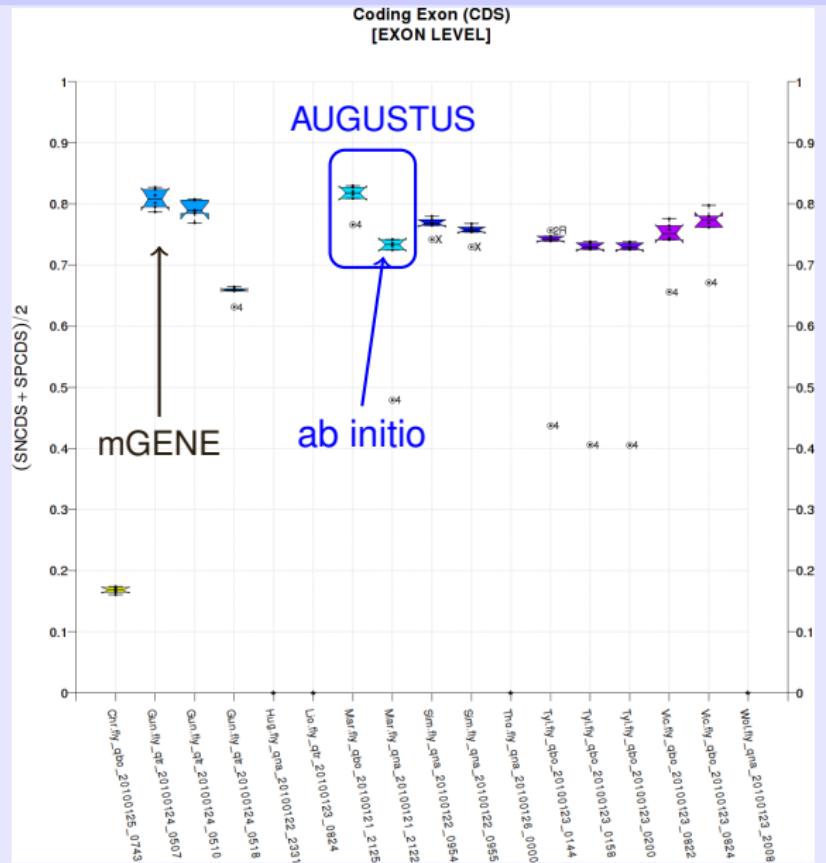
RGASP: results by submissions

human, transcript level, coding sequence



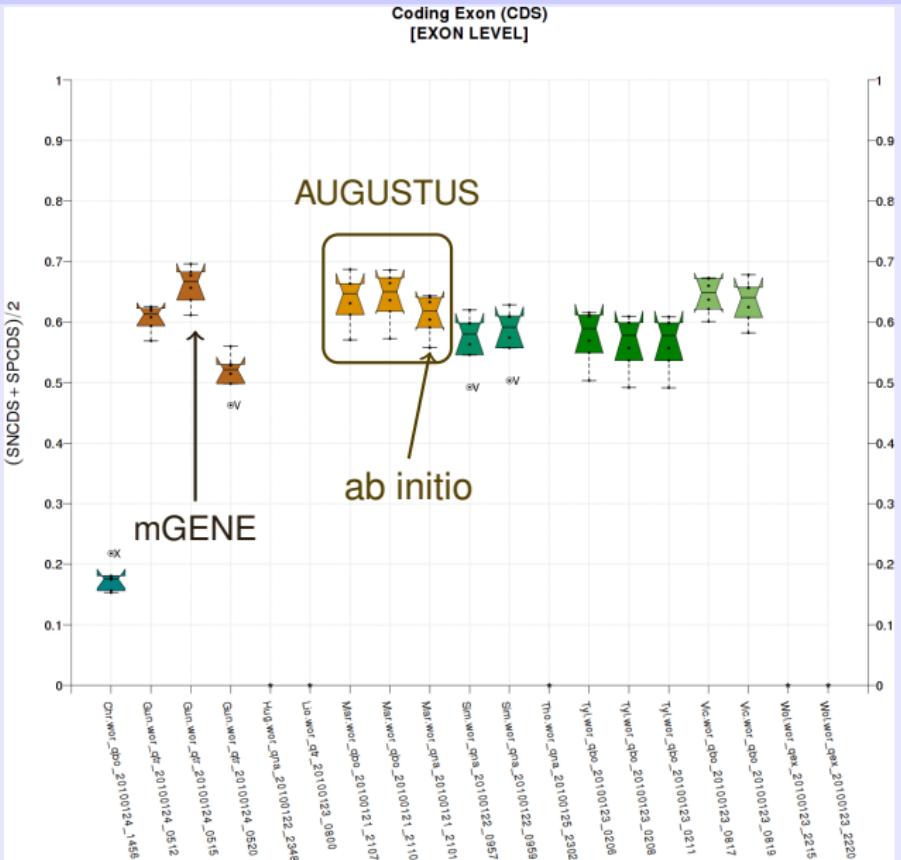
RGASP: results by submissions

Drosophila, coding exon level



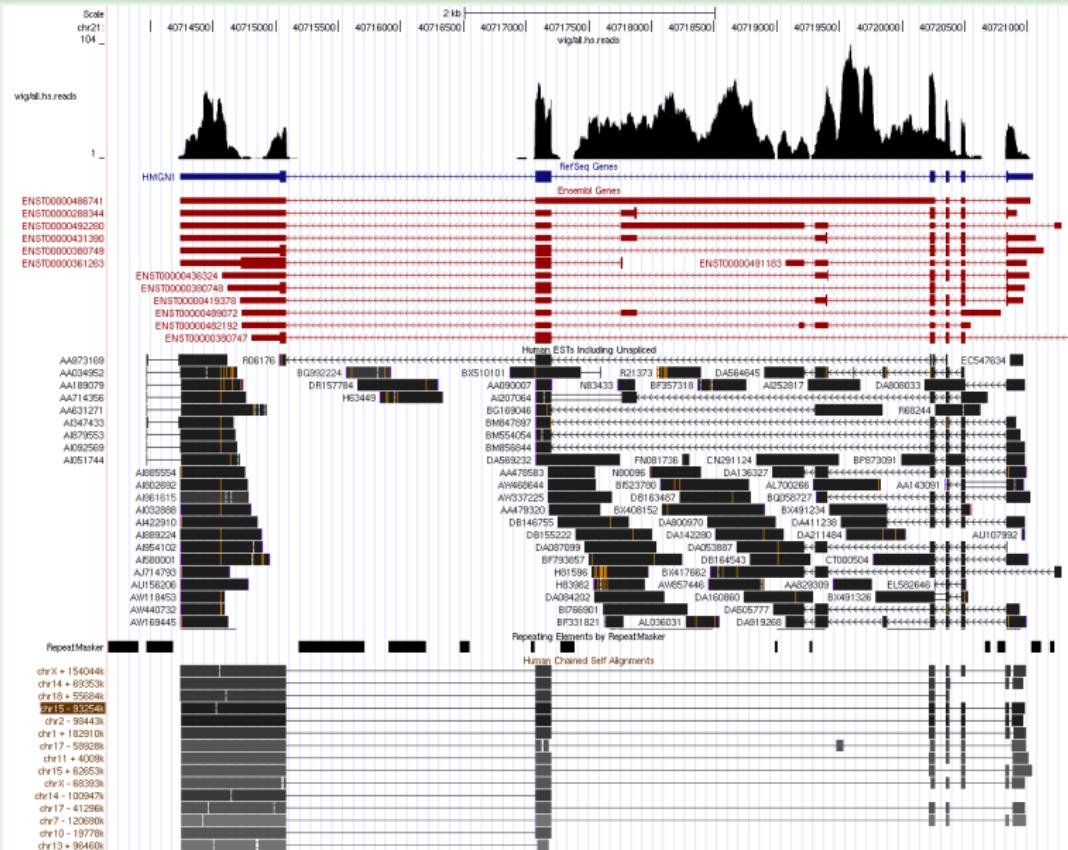
RGASP: results by submissions

C. elegans, coding exon level



Why was the accuracy not better?

Problems 1 and 2: intronic transcription, self-similarity of genome





Introduction

Gene Finding Problem

Statistical Features of
Genes

Signal and Content Models

Training

Gene Finding as
Optimization Problem

HMMs/CRFs

Comparative Gene
Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

Why was the accuracy not better?

Other problems with RNA-Seq

Problem 3 repeats

Problem 4 alignment errors

Problem 5 mRNA does not determine protein sequence
(see Natalie's talk)



Introduction

Gene Finding Problem

Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

Conclusions from RGASP

- ① even with RNA-Seq the best methods **still** make **a lot of errors** or miss genes
- ② **improving** an underlying **ab initio** model also improves RNA-Seq-based performance



Introduction

Gene Finding Problem

Statistical Features of
Genes

Signal and Content
Models

Training

Gene Finding as
Optimization Problem

HMMs/CRFs

Comparative Gene
Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

RGASP: preliminary results

My Conclusions

- RNA-Seq did not help as much as expected:
 - substantial intronic transcription
 - mapping ambiguity due to self-similarity of genomes
- still a lot of potential for improvement



Introduction

Gene Finding Problem
Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

1 Introduction

Gene Finding Problem

Statistical Features of Genes

2 Signal and Content Models

3 Training

4 Gene Finding as Optimization Problem

5 HMMs/CRFs

6 Comparative Gene Finding

7 Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

8 Alternative Splicing



Introduction

Gene Finding Problem

Statistical Features of
Genes

Signal and Content
Models

Training

Gene Finding as
Optimization Problem

HMMs/CRFs

Comparative Gene
Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

Proteomics and Genomics: Proteogenomics

Peptide from Tandem Mass Spectrometry (MS/MS) for Gene Prediction

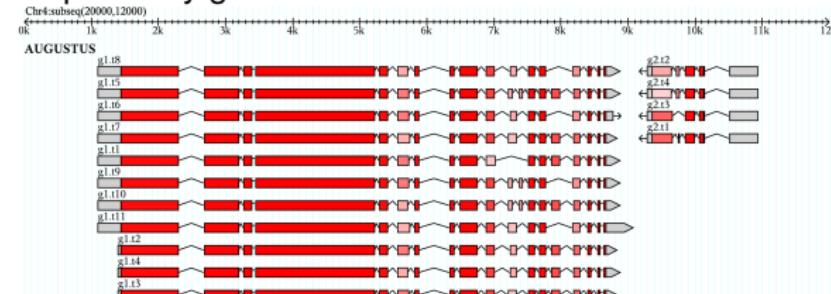
- unknown proteins of interest cleaved into smaller peptides
- masses of peptides measured with a mass spectrometer
- masses compared *in silico* with peptides from a candidate protein database
- result: set of peptides that is actually translated



Candidate Protein Database

should be very **inclusive**

- ① six-frame translation of genome
- ② known genes
- ③ predicted genes:
 - sample many gene structures with AUGUSTUS



[Introduction](#)

[Gene Finding Problem](#)

[Statistical Features of Genes](#)

[Signal and Content Models](#)

[Training](#)

[Gene Finding as Optimization Problem](#)

[HMMs/CRFs](#)

[Comparative Gene Finding](#)

[Evidence Integration](#)

[Hints to AUGUSTUS](#)

[RNA-Seq](#)

[Proteogenomics](#)

[Protein Homology](#)

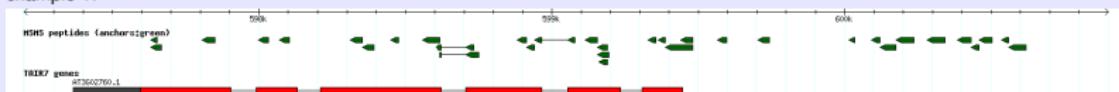
[Alternative Splicing](#)

Map Peptides to Genome

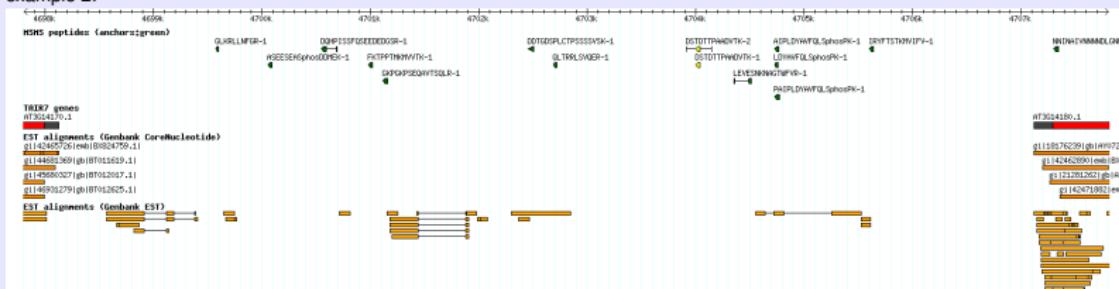
Peptides can

- match protein from set of “known genes”: **Partially confirm gene structure, existence of gene**
- not match any known gene: What then?
Extend known gene? Join two known genes? New splice form?
Replace splice form of “known” gene?

example 1:

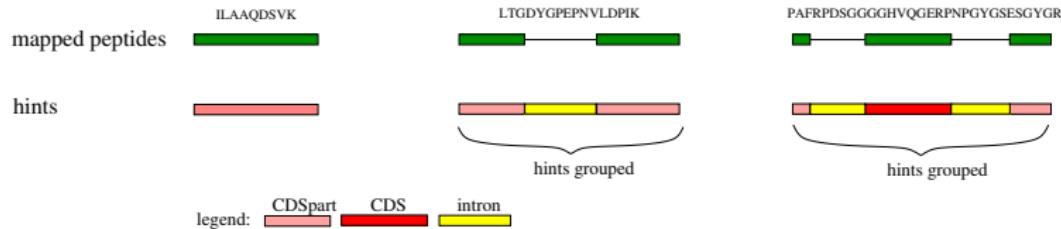


example 2:





Peptides to “Hints”



(average peptide length: 15.2)

CDSpart hint: region is likely protein-coding

CDS hint: region is exact protein-coding exon

intron hint: region is likely an intron with exact boundaries

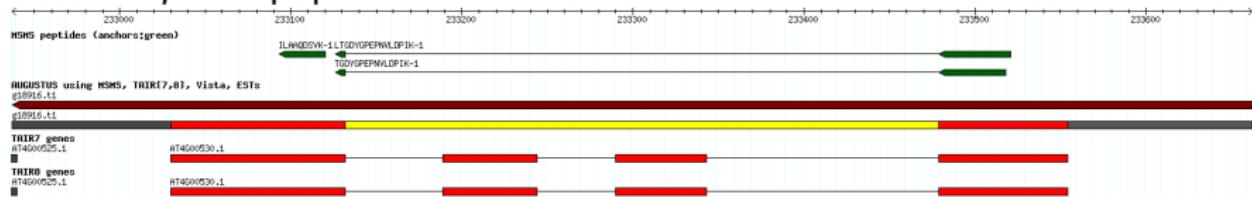
Hints can include strand and frame information.

Examples for Gene Corrections

Who would have doubted this gene?



precious spliced peptides





Introduction

Gene Finding Problem

Statistical Features of
Genes

Signal and Content
Models

Training

Gene Finding as
Optimization Problem

HMMs/CRFs

Comparative Gene
Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

Conclusions on Proteogenomic Evidence

MS/MS

- peptides not only good for “verification” of given gene structures but also valuable, additional **orthogonal** source of evidence for predicting gene structures
- new translation-level info (reading frame)
- not enough coverage for identification of alternative splicing based on MS/MS alone



Introduction

Gene Finding Problem
Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

1 Introduction

Gene Finding Problem

Statistical Features of Genes

2 Signal and Content Models

3 Training

4 Gene Finding as Optimization Problem

5 HMMs/CRFs

6 Comparative Gene Finding

7 Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

8 Alternative Splicing



Introduction

Gene Finding Problem

Statistical Features of
Genes

Signal and Content
Models

Training

Gene Finding as
Optimization Problem

HMMs/CRFs

Comparative Gene
Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

Incorporating Protein Homology: "The Pipeline Way"

Use known proteins from other species

- ① align proteins to genome
(e.g. with BLASTX, TBLASTN, Exonerate)
- ② use matches from step 1 as evidence for coding regions in
gene finding
(many gene finders)

Disadvantages

- step 1 misses info: no weak homology or many false positives
- gene finders as subject to split gene errors
- protein family MSAs hold more info than individual proteins



Introduction

Gene Finding Problem

Statistical Features of
Genes

Signal and Content
Models

Training

Gene Finding as
Optimization Problem

HMMs/CRFs

Comparative Gene
Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

Incorporating Protein Homology: "The Pipeline Way"

Use known proteins from other species

- ① align proteins to genome
(e.g. with BLASTX, TBLASTN, Exonerate)
- ② use matches from step 1 as evidence for coding regions in
gene finding
(many gene finders)

Disadvantages

- step 1 misses info: no weak homology or many false positives
- gene finders as subject to split gene errors
- protein family MSAs hold more info than individual proteins



Introduction

Gene Finding Problem

Statistical Features of
Genes

Signal and Content
Models

Training

Gene Finding as
Optimization Problem

HMMs/CRFs

Comparative Gene
Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

Protein Homology in One Step

Genewise (Birney 2004)

- marry (multiply) two HMMs:
 - simple (species-unspecific) HMM for gene prediction with
 - profile HMM for protein family (like HMMer)
- one big (product) HMM that does both:
 - gene finding and
 - protein to protein-family alignment + membership decision

Disadvantages

- slow
- poor performance for weak homology
- cannot cope with missing homology
- no integration with other evidence (RNA-Seq, ESTs)



Introduction

Gene Finding Problem

Statistical Features of
Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

Protein Homology in One Step

Genewise (Birney 2004)

- marry (multiply) two HMMs:
 - simple (species-unspecific) HMM for gene prediction with
 - profile HMM for protein family (like HMMer)
- one big (product) HMM that does both:
 - gene finding and
 - protein to protein-family alignment + membership decision

Disadvantages

- slow
- poor performance for weak homology
- cannot cope with missing homology
- no integration with other evidence (RNA-Seq, ESTs)



Introduction

Gene Finding Problem

Statistical Features of
Genes

Signal and Content
Models

Training

Gene Finding as
Optimization Problem

HMMs/CRFs

Comparative Gene
Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

Protein Profile Extension of AUGUSTUS

Bioinformatics Advance Access published January 6, 2011

A novel hybrid gene prediction method employing protein multiple sequence alignments

Oliver Keller^{1,*}, Martin Kollmar², Mario Stanke^{3,*} and Stephan Waack¹

AUGUSTUS-PPX: Protein Profile eXtension

- find members of given protein family (MSA) in DNA
- seamless extension of AUGUSTUS:
 - same program, same species parameters, etc.
 - UTR prediction possible
 - additional hints integration possible, e.g. RNA-Seq
 - defaults to normal AUGUSTUS for non-member genes and gene parts without homology



Introduction

Gene Finding Problem

Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

Block

Ungapped and highly conserved section of a MSA

Example

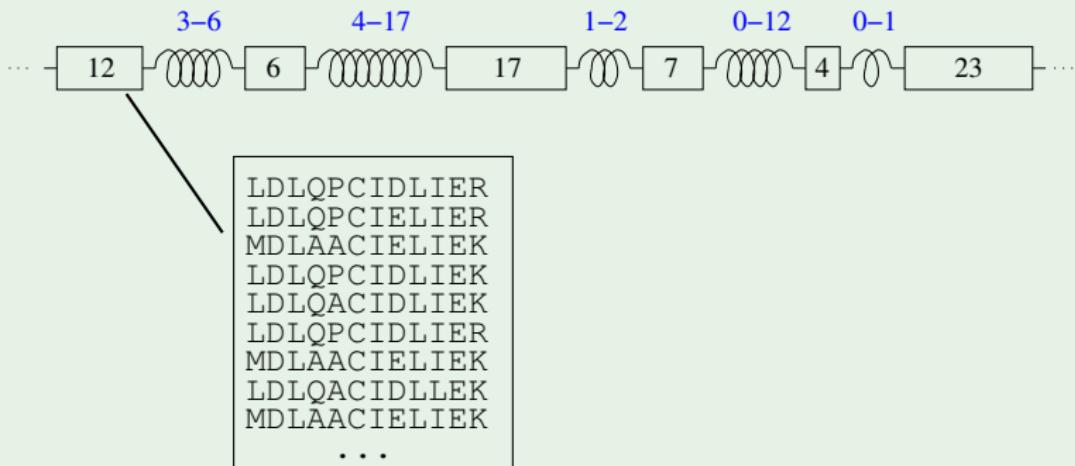
HsDHC1 f1	[CNPPTDP--GRKPLSHRFLRHVP--VYVVDYPGPASLTQIYGTFNRAMLRLIPSLRT---
HsDHC2 f1	MSAGGRL--GRKLTTTRFTSIVR--LCSIDYPEREQLQTIIYGAYLEPVLHKNLKNHSIWG
HsDHC3A f1	MIHPGG--GRNDIPQRLKRQFT--VFNCTLPESNASIDKIFGIIIGCGYFDPCRSFKPQIC
HsDHC3B f1	MIHPGG--GRNDIPQRLKRQFS--IFNCTLPESEASVDKIFGVIGVGHYCTQRGKFSEEVRI
HsDHC4A f1	MNPMV---GSHTINPRLQRHFT--VFAFNFPSDLALNTIYGQIFSFHFO--QQAFAPSI
HsDHC4B f1	MNPTS---GSETIDSRLQRHFC--VFAVSFPQEALETTIINTLTQHLA--FRSVSMAI
HsDHC4C f1	MNPTA---GSETINPRLQRHFS--VFVLSFFPGADALSSIYIISIILTQHLK--LGNFPAWL
HsDHC5 f1	MGKAGG--GRNEVDPRFISLFS--VFNVPPFSEESLHLIYSSILKGHTSFHESIVA--
HsDHC6 f1	MGPPGG--GRTVISPRLRSRFN--IINMTFFPKSIIIRIFGTMINQKLQDFEEEVKP--
HsDHC7A f1	MGPPGG--GRNPVTPRYYMRHFN--IITINEFSDKSMYTIFSRLTWHLIECYKFPDEFI
HsDHC7B f1	MGPPGG--GRNDITGRFTRHLN--IISINAFEDDILTKIFSSIVDWHFGKG--FDVMFL
HsDHC7C f1	MGPPGG--GRNPVTPRRCIRHFN--ICSINSFSDETMVRIFSSIVAFYLRTHE--FPPEYF
HsDHC8 f1	MGPPGG--GRNTVTIPRLMRHFN--YLSFAEMDEVSKKRIFSTILGNWLDGLGEKSYRE
HsDHC9 f1	CAPPGG--GRNPVTPRPFIRHFS--MLCLPMPSEHSLKQIFQAILNGFLS--DFPPAVKQ
HsDHC10 f1	CVPVVN----DISPRLLKHF--MLVLPHPSQDILCTIFQAHLGIVFSINNFTPEVQK
HsDHC11 f1	VTVPGYC---ERPLCPRLFRET--VIALESMTQATLIERHVPTIQAWLERFPSVERERA

AUGUSTUS-PPX

Block Profile

- frequency table for each block column
- distances between blocks constrained by intervals $[d_{\min}, d_{\max}]$
- no insertions or deletions
- order of blocks preserved

Example (Myosins)





Introduction

Gene Finding Problem

Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

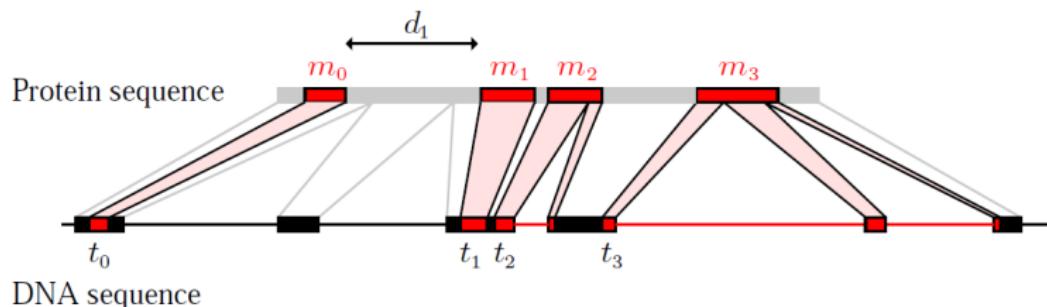
Protein Homology

Alternative Splicing

AUGUSTUS-PPX

Problem:

Find gene structure in genome under **hypothesis**, that predicted **protein sequence fits into family**.



Scoring

$$\begin{aligned} & \text{score}(\text{gene structure}, \text{DNA}, \text{block profile}) \\ = & \text{score}(\text{gene structure}, \text{DNA}) \\ & \cdot \underbrace{\text{score}(\text{block profile} \mid \text{gene structure}, \text{DNA})}_{\rightarrow \text{protein sequence cand.}} \end{aligned}$$



Introduction

Gene Finding Problem
Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

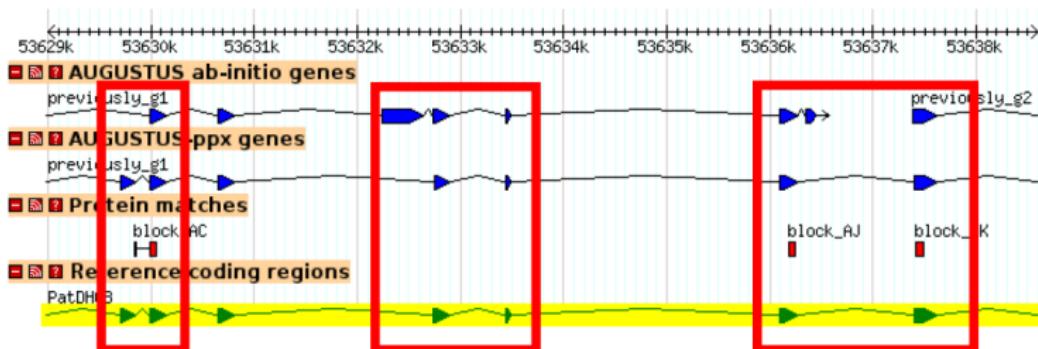
Proteogenomics

Protein Homology

Alternative Splicing

AUGUSTUS-PPX: How it helps

Example



Corrections w.r.t. ab initio

- add exons containing blocks
- remove false positive exons (constraint on inter-block lengths)
- join split genes



Introduction

Gene Finding Problem
Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

AUGUSTUS-PPX: Results

Test on: Dynein Heavy Chain (DHC) family

- > 4000 amino acids, up to 100 exons, several 100 000 bp genomic range
- 42 Blocks with 1214 sites

Accuracy Results

species	AUGUSTUS-PPX		AUG. ab-initio	Genewise	
	full	ex-ortho		full	ex-ortho
highly accurate genes (%)					
human	62.5	62.5	31.3	93.8	6.3
mouse	72.2	61.1	22.2	55.6	0.0
chicken	58.3	50.0	8.3	16.7	0.0
frog	44.4	38.9	0.0	11.1	0.0
zebrafish	57.1	35.7	57.1	14.3	0.0
snail	44.4	44.4	5.6	11.1	0.0
total	56.3	49.0	11.5	34.4	1.0



Introduction

Gene Finding Problem
Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS
RNA-Seq
Proteogenomics
Protein Homology

Alternative Splicing

1 Introduction

Gene Finding Problem
Statistical Features of Genes

2 Signal and Content Models

3 Training

4 Gene Finding as Optimization Problem

5 HMMs/CRFs

6 Comparative Gene Finding

7 Evidence Integration

Hints to AUGUSTUS
RNA-Seq
Proteogenomics
Protein Homology

8 Alternative Splicing

Evidence Suggests Alternative Splicing





Introduction

Gene Finding Problem

Statistical Features of Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS

RNA-Seq

Proteogenomics

Protein Homology

Alternative Splicing

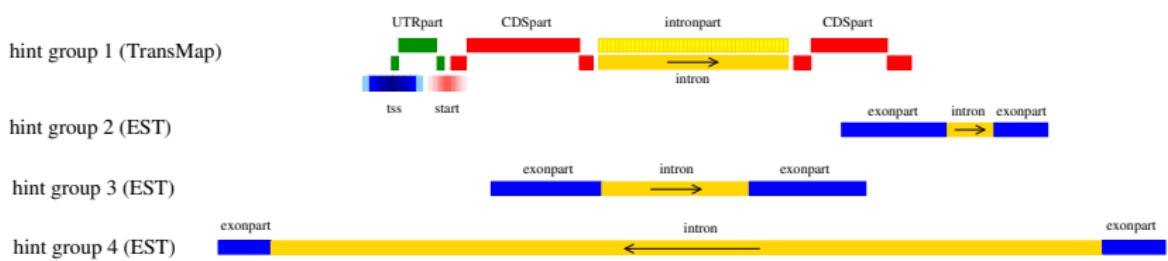
AUGUSTUS: Alternative Transcripts Based on Evidence

Hints can be **grouped** to tell AUGUSTUS they are coming from the **same transcript**.

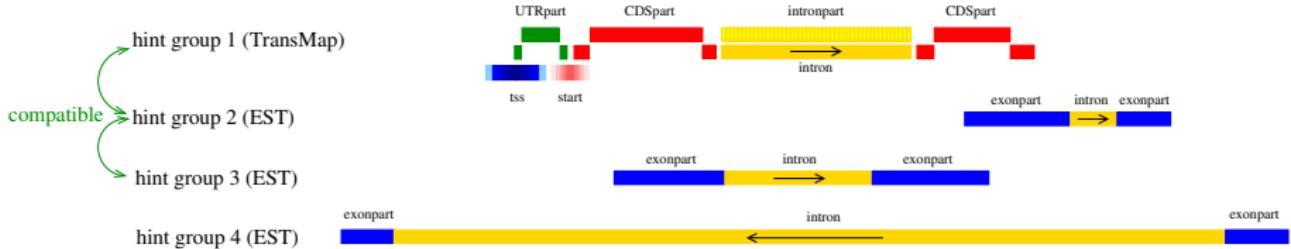
For example: All hints from one cDNA alignment belong to one group. Based on evidence predict **alternative splice forms**.

- ① Split input sequence between genes based on hint groups and preliminary gene prediction. For each piece, do:
- ② Determine **compatibility** between hint groups.
- ③ Filter out outlier hint groups.
- ④ For each hint group g , do:
 - ① Deactivate all hint groups h that are incompatible with g .
 - ② Run Viterbi algorithm to predict single-transcript genes.
- ⑤ Join transcripts to genes

Alternative Transcripts Based on Evidence



Alternative Transcripts Based on Evidence



Alternative Transcripts Based on Evidence

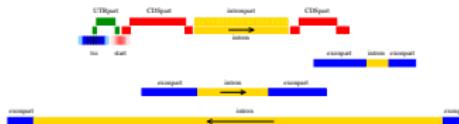
prediction
run 1



prediction
run 2



prediction
run 3

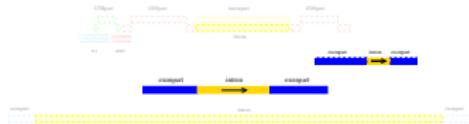


Alternative Transcripts Based on Evidence

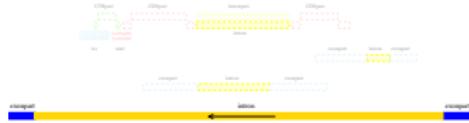
prediction
run 1



prediction
run 2

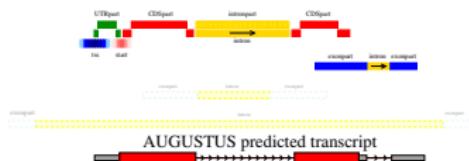


prediction
run 3

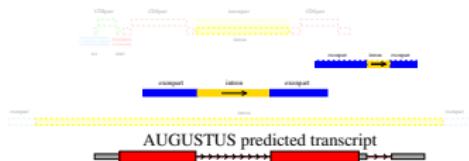


Alternative Transcripts Based on Evidence

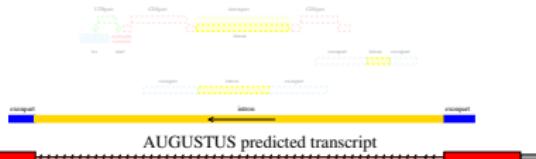
prediction
run 1



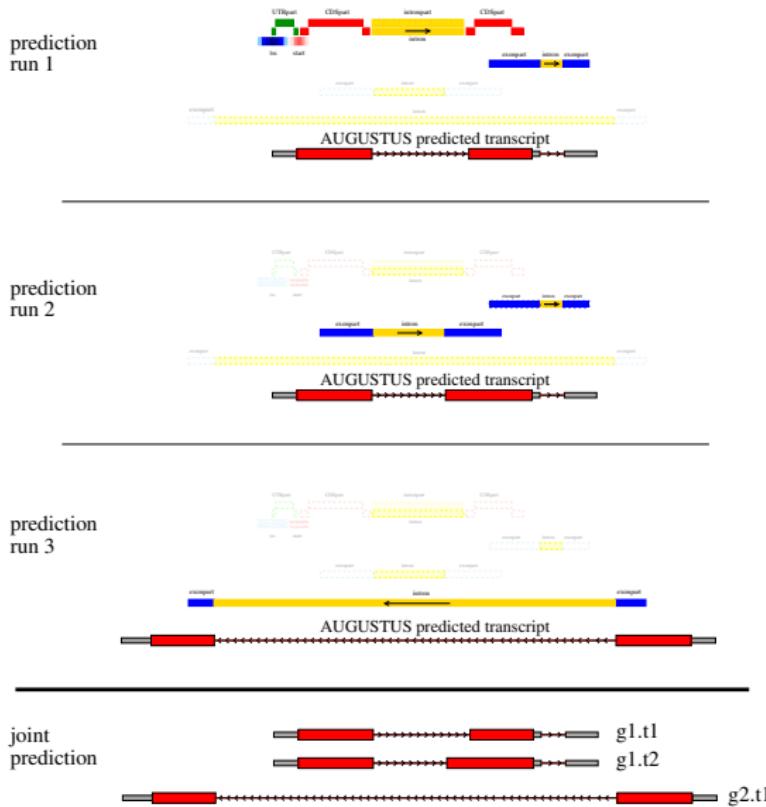
prediction
run 2



prediction
run 3



Alternative Transcripts Based on Evidence



What else?

- new training web server:

<http://bioinf.uni-greifswald.de>

- upload genome & transcriptome (currently, ESTs & 454 only)
- get annotation & trained AUGUSTUS parameter set

- use other evidence

- e.g. peptides from MS/MS (Natalie's talk)



- AUGUSTUS-PPX extension (predicts protein family members in genome)

HCEVGLIIGTGNAKYMEELRQIDFGWGT.....DDGRMCINTEWDLGDDGSLIEDRKFEDFRRRG. SLNPGKQRFERMVSGRYMEDVVRLLVVK
EPWICISIFGIGTGNTGCYMEEEINITKLPEQLDRDKLKEKTHMININVEWSF. DNEHLKHLPITKDYVVIDOKLSTNPFGHFLFEKRVSGFLGEVLRNILVD
EAKMGLFSITGNCNGAYDOWDNIPKLEGKVP. .DDIKSSSPMAINECYGAF. DNEHLVIIPLTKYDIODIEE. SPPPGQOAFEMISGYYLGEVLRLALLD
ETKMGSVIFGIGVNGAYDOWDNIPKLEGKVP. .DDIPSSPMAINECYGAF. DNEHLVIIPLTKYDIODIEE. SPPPGQOAFEMISGYYLGEVLRLALLD
TCEIVLIGVGNSACYMEEEINVEWEG. .NORQMCINHEWGAFGNDGCGSDOIRIDFDKVVDEE. SLNSGNQRFENMISGYYLGEIVRNILID
GTEIVLIGVGNSACYMEEEINVEWEG. .DFFPDEMNNIICEWCDF. DNQHVLVPLTKYDVIADEE. SPPPGQOAFEMISGYYLGEIVRNILID
ETKMGSVIFGIGVNGAYDOWDSIEKLEGKLA. .DDIPSNSPMAINECYGAF. DNEHLVIIPLTKYDVIADEE. SPPPGQOAFEMISGYYLGEELLRLVLE
PCEVGLVVDITGNAKYMEEARHVAVLDE.DRGRVCVSVENGSFSDDGGLGPVLTFTDHTLDOE. SLNPGKQRFERMISGGLYLGEIVRNILAH
RCEIIGLIVGIGTGNAKYMEELRNVAVGPG.DSGRMCINHEWGAFGDDGSLAMLSFRPDASVOA. SINPGKQRFERMISGGNYLGEIVRHILLH
PCYIIGLIGVGNSACYMEEPWKKYYKAG.KINIEFGNF.DK. .DLPTSPIDLVMDVY. SANRSRQLFEMISGAYLGEIVRRIIFVN
HCEVGLIIGTGNSACYMEEMRNVELVEG.EEGRMCVSMEWGAFGNDGCLDDFRTEFDVADEL. SLNPGKQRFERMISGGMLGEIVRNILID



Acknowledgements

Introduction

Gene Finding Problem
Statistical Features of
Genes

Signal and Content Models

Training

Gene Finding as Optimization Problem

HMMs/CRFs

Comparative Gene Finding

Evidence Integration

Hints to AUGUSTUS
RNA-Seq
Proteogenomics
Protein Homology

Alternative Splicing

Thanks to

Stephan Waack

|| Göttingen

Oliver Keller

Burkhard Morgenstern

David Haussler

|| UC Santa Cruz

Mark Diekhans

Michael Hippler

|| Münster

Michael Specht

Vineet Bafna

|| UC San Diego

Natalie Castellana

Sam Payne

Katharina Hoff

|| Greifswald