# Lecture 17

# Comparative genomics I:

# Genome annotation using evolutionary signatures

# Module V: Comparative genomics and evolution

- Today: Whole-genome comparative genomics
  - Evolutionary signatures for systematic genome annotation
- Next week: Phylogenetics and Phylogenomics
  - Distance-based and model-based phylogenetics approaches
  - Gene trees and species trees, reconciliation, coalescence
- Computational foundations:
  - Evolutionary rates and models of evolution
  - Dynamic programming on two-dimensional tree structures
  - Synteny-based alignment, genome assembly

# Key goal: Evolution preserves functional elements

Gal4

Gal10 — Gal1

GAL10

```
          Scer  TTATATTGAATTTTCAAAAATTCTTACTTTTTTTTTGGATGGACGCAAAGAAGTTTAATAATCATATTACATGGCATTACCACCATATACA
          Spar  CTATGTTGATCTTTTCAGAATTTTT-CACTATATTAAGATGGGTGCAAAGAAGTGTGATTATTATATTACATCGCTTTCCTATCATACACA
          Smik  GTATATTGAATTTTTCAGTTTTTTTTCACTATCTTCAAGGTTATGTAAAAAA-TGTCAAGATAATATTACATTTCGTTACTATCATACACA
          Sbay  TTTTTTTGATTTCTTTAGTTTTCTTTCTTTAACTTCAAAATTATAAAAGAAAGTGTAGTCACATCATGCTATCT-GTCACTATCACATATA
                 *    ** ****  * *   ** *   ** ** *   **         * *      **  ** **   *  * * ** ***  *
```

TBP

```
          Scer  TATCCATATCTAATCTTAC TTATA TGTTGT-GGAAAT-GTAAAGAGCCCCATTATCTTAGCCTAAAAAAACC--TTCTCTTTGGAACTTTCAGTAATACG
          Spar  TATCCATATCTAGTCTTAC TTATA TGTTGT-GAGAGT-GTTGATAACCCCAGTATCTTAACCCAAGAAAGCC--TT-TCTATGAAACTTGAACTG-TACG
          Smik  TACCGATGTCTAGTCTTAC TTATA TGTTAC-GGGAATTGTTGGTAATCCCAGTCTCCCAGATCAAAAAAGGT--CTTTCTATGGAGCTTTG-CTA-TATG
          Sbay  TAGATATTTCTGATCTTTC TTATA TATTATAGAGATGCCAATAAACGTGCTACCTCGAACAAAAGAAGGGGATTTTCTGTAGGGCTTTCCCTATTTTG
                **   ** *** **** ****** ** **   *  *      *     **     *         **  ** **       *** ****   *  *
```

GAL4          GAL4          GAL4

```
          Scer  CTTAACTGCTCATTGC-----TATATTGAAGTA CGG ATTAGAAGCCG CCG AG CGG GCGACAGCCCT CCGA CGG AAGACTCTCCT CCG TGCGTCCTCGTCT
          Spar  CTAAACTGCTCATTGC-----AATATTGAAGTA CGG ATCAGAAGCCG CCG AG CGG ACGACAGCCCT CCGA CGG AATATTCCCCT CCG TGCGTCGCCGTCT
          Smik  TTTAGCTGTTCAAG--------ATATTGAAATA CGG ATGAGAAGCCG CCG AA CGG ACGACAATTCC CCGA CGG AACATTCTCCT CCG CGCGGCGTCCTCT
          Sbay  TCTTATTGTCCATTACTTCGCAATGTTGAAATA CGG ATCAGAAGCTG CCG AC CGG ATGACAGTACT CCGG CGG AAAACTGTCCT CCG TGCGAAGTCGTCT
                       **  **          ** * ***** ** ***** ****** ** **   *** *    **** **   ***** *** *     ****** ***   * ***
```

GAL4

```
          Scer  TCACCGG-TCGCGTTCCTGAAACGCAGATGTGC CT CGC GCCGCACTGCT CCG AACAAT AAAGATTCTACAA-----TACTAGCTTTT--ATGGTTATGAA
          Spar  TCGTCGGGTTGTGTCCCTTAA-CATCGATGTAC CT CGC GCCGCCCTGCT CCG AACAAT AAGGATTCTACAAGAAA-TACTTGTTTTTTTATGGTTATGAC
          Smik  ACGTTGG-TCGCGTCCCTGAA-CATAGGTACGG CT CGC ACCACCGTGGT CCG AACTAT AATACTGGCATAAAGAGGTACTAATTTCT--ACGGTGATGCC
          Sbay  GTG-CGGATCACGTCCCTGAT-TACTGAAGCGT CT CGC CCCGCCATACC CCG AACAAT GCAAATGCAAGAACAAA-TGCCTGTAGTG--GCAGTTATGGT
                     ** *    ** *** *   *           * *****  * *     ****** *    *        *  *   *  *  *  *     ** ** ***
```

MIG1

```
          Scer  GAGGA-AAAATTGGCAGTAA----CCTGG CCCCACAAACCTT -CAAATTAACGAATCAAATTAACAACCATA-GGATGATAATGCGA------TTAG--T
          Spar  AGGAACAAAATAAGCAGCCC----ACTGA CCCCATATACCTT TCAAACTATTGAATCAAATTGGCCAGCATA-TGGTAATAGTACAG------TTAG--G
          Smik  CAACGCAAAATAAACAGTCC----CCCGG CCCCACATACCTT -CAAATCGATGCGTAAAACTGGCTAGCATA-GAATTTTGGTAGCAA-AATATTAG--G
          Sbay  GAACGTGAAATGACAATTCCTTGCCCCCT-CCCCAATATACTT TGTTCCGTGTACAGCACACTGGATAGAACAATGATGGGGTTGCGGTCAAGCCTACTCG
                      ****   *               * ********* ***         * * *   * ** *   *      *  *      **
```

MIG1                                                  TBP

```
          Scer  TTTTTAGCC TTATTTCTGGGG TAATTAATCAGCGAAGCG--ATGATTTTT-GATCTATTAACAGATA TATAA ATGGAAAAGCTGCATAACCAC-----TT
          Spar  GTTTT--TC TTATTCCTGAGA CAATTCATCCGCAAAAATAATGGTTTTT-GGTCTATTAGCAAACA TATAA ATGCAAAAGTTGCATAGCCAC-----TT
          Smik  TTCTCA--CC TTTCTCTGTGA TAATTCATCACCGAAATG--ATGGTTTA--GGACTATTAGCAAACA TATAA ATGCAAAAGTCGCAGAGATCA-----AT
          Sbay  TTTTCCGTT TTACTTCTGTAG TGGCTCAT--GCAGAAGTAATGGTGTTTTCTGTTCCTTTTGCAAACA TATAA ATATGAAAGTAAGATCGCCTCAATTGTA
                *  *      *    ***     * **   *  *     ** ** **   *      ******       **** *
```

```
          Scer  TAACTAATACTTTCAACATTTTCAGT--TTGTATTACTT-CTTATTCAAAT----GTCATAAAAGTATCAACA-AAAAATTGTTAATATACCTCTATACT
          Spar  TAAATAC-ATTTGCTCCTCCAAGATT--TTTAATTTCGT-TTTGTTTTATT----GTCATGGAAATATTAACA-ACAAGTAGTTAATATACATCTATACT
```

## We can 'read' evolution to reveal functional elements

Conservation island

Yeast (Kellis et al, Nature 2003), Mammals (Xie, Nature 2005), Fly (Stark et al, Nature 07)

# Comparative Genomics



Lecture 17 (Today):

Using evolution to study genomes

Evolution → Genomics

Using genomics to study evolution

Lectures 18-19 (Thursday):

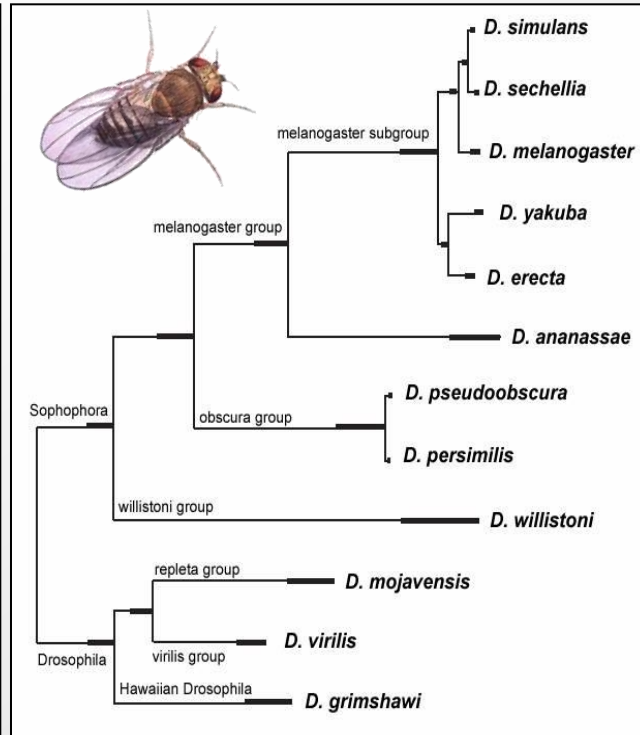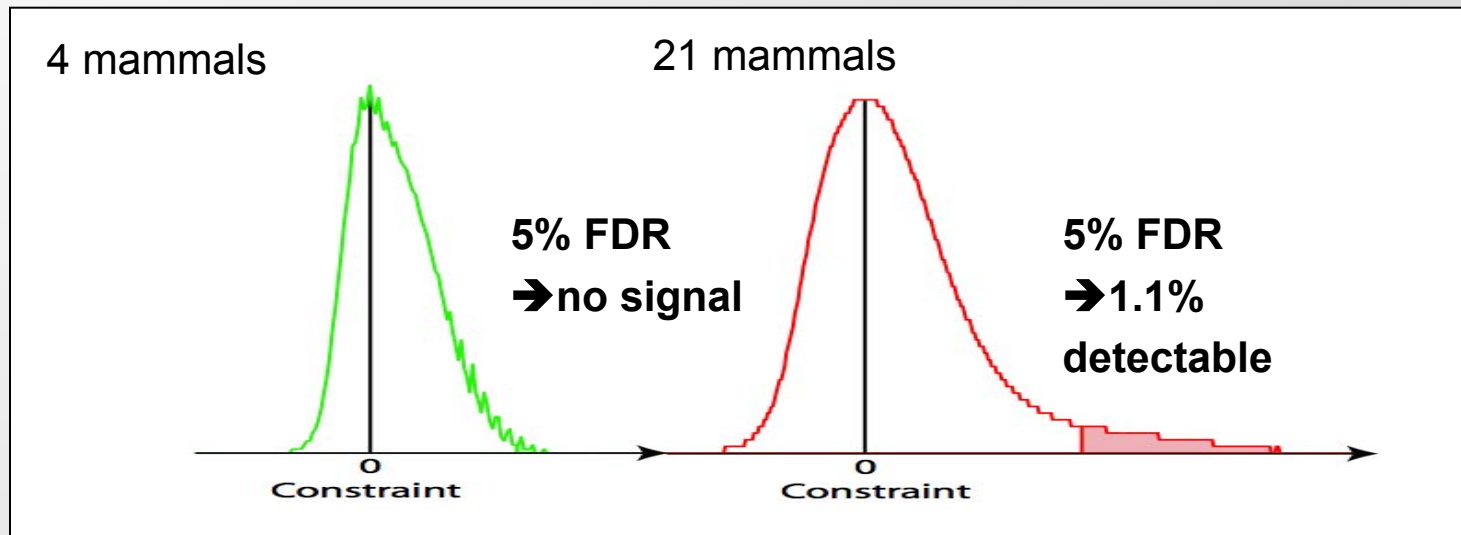# Comparative genomics I: Evolutionary signatures

- **Nucleotide conservation: evolutionary constraint**
  - Purifying selection, neutral branch length, discovery power
  - Detect constrained elements: nucleotides, windows, HMM
  - Estimate fraction constrained: signal vs. background
- **Evolutionary signatures: focus on pattern of change**
  - Different functions ⇔ Characteristic patterns of evolution
- **Signatures of protein-coding genes**
  - Reading-frame conservation, codon-substitution frequency
  - Likelihood ratio framework: Estimating $Q_C Q_N$, scoring
  - Revise genes, read-through, excess constraint regions
- **Signatures of microRNA genes**
  - Structural and evolutionary features of microRNAs
  - Combining features: decision trees, random forests
  - Sense/anti-sense miRNAs, mature/star arm cooperation
- **Measuring selection within the human lineage**

# Comparative genomics for genome annotation

## 29 mammals

## 12 flies

## 17 fungi

- **Compare related species to discover functional elmts**
- **Evolution process: random mutation, natural selection**
  - Non-functional regions: accumulate mutations, kept
  - Functional regions: accumulate mutations, decrease fitness
  - Evolutionary time: less fit organisms & their genes thin out

# Power of many closely related: total branch length

- **More branch length ➔ more events ➔ more power**
  - Goal: functional vs. non-functional based on # of mutations
  - Very close distances: no mutations in either region
  - Sufficient distance: ability to distinguish increases
  - Very far distances: functional regions no longer conserved
- **Many closely related species >> few distantly related**
  - For same total branch length: prefer many close species
  - Functional regions conserved for each pair of species
  - Non-functional regions accumulate noise **independently**
  - Analogy: recording a concert with multiple microphones

# Comparative genomics I: Evolutionary signatures

- **Nucleotide conservation: evolutionary constraint**
  - Purifying selection, neutral branch length, discovery power
  - Detect constrained elements: nucleotides, windows, HMM
  - Estimate fraction constrained: signal vs. background
- **Evolutionary signatures: focus on pattern of change**
  - Different functions ⇔ Characteristic patterns of evolution
- **Signatures of protein-coding genes**
  - Reading-frame conservation, codon-substitution frequency
  - Likelihood ratio framework: Estimating $Q_C Q_N$, scoring
  - Revise genes, read-through, excess constraint regions
- **Signatures of microRNA genes**
  - Structural and evolutionary features of microRNAs
  - Combining features: decision trees, random forests
  - Sense/anti-sense miRNAs, mature/star arm cooperation

# Genome-wide alignments reveal orthologous segments



Courtesy of Don Gilbert. Used with permission.

100 genes

- **Genome-wide alignments span entire genome**
- **Comparative identification of functional elements**

# Comparative genomics and evolutionary signatures

- **Comparative genomics can reveal functional elements**
  - For example: exons are deeply conserved to mouse, chicken, fish
  - Many other elements are also strongly conserved: exons / regulatory?
- **Develop methods for estimating the level of constraint**
  - Count the number of edit operations, number of substitutions and gaps
  - Estimate the number of mutations (including estimate of back-mutations)
  - Incorporate information about neighborhood: conservation 'windows'
  - Estimate the probability of a constrained 'hidden state': HMMs next week
  - Use phylogeny to estimate tree mutation rate, or 'rejected substitutions'
  - Allow different portions of the tree to have different rates: phylogenetics

# Detecting **rates** and **patterns** of selection ($\omega$/$\pi$)



**Neutral sequence**



**Decreased rate ω**



**Unusual patterns π**

## Estimating intensity of constraint ($\omega$):

• Probabilistic model of substitution rate

• Maximum Likelihood (ML) estimation of $\omega$

   - Report rate ω

   - Report log odds score that non-neutral

• Window-based vs. sitewise application

## Detect unusual substitution pattern ($\pi$):

• Probabilistic model of stationary distribution that is different from background.

• ML estimator ($\pi$) of this vector

   • Report PWM for each k-mer in genome.

   • Report log odds score that non-neutral

**Manuel Garber, Or Zuk, Xiaohui Xie**

# Measuring constraint at individual nucleotides

- **Reveal individual transcription factor binding sites**
- **Within motif instances reveal position-specific bias**
- **More species: motif consensus directly revealed**

# Detect SNPs that disrupt conserved regulatory motifs

- Functionally-associated SNPs enriched in states, constraint
- Prioritize candidates, increase resolution, disrupted motifs

13

# Comparative genomics I: Evolutionary signatures

- **Nucleotide conservation: evolutionary constraint**
  - Purifying selection, neutral branch length, discovery power
  - Detect constrained elements: nucleotides, windows, HMM
  - Estimate fraction constrained: signal vs. background
- **Evolutionary signatures: focus on pattern of change**
  - Different functions ⇔ Characteristic patterns of evolution
- **Signatures of protein-coding genes**
  - Reading-frame conservation, codon-substitution frequency
  - Likelihood ratio framework: Estimating $Q_C Q_N$, scoring
  - Revise genes, read-through, excess constraint regions
- **Signatures of microRNA genes**
  - Structural and evolutionary features of microRNAs
  - Combining features: decision trees, random forests
  - Sense/anti-sense miRNAs, mature/star arm cooperation

# Estimating portion of the genome under constraint

Constraint calculated over a **50mer**

| 4 mammals | 29 mammals |
|---|---|
| **5% FDR** **➜0.6%** **detectable** | **5% FDR** **➜1.8%** **detectable** |

Constraint calculated over a **12mer**

| 4 mammals | 21 mammals |
|---|---|
| **5% FDR** **➜no signal** | **5% FDR** **➜1.1%** **detectable** |

**Or Zuk, Manuel Garber**

# Estimating total fraction under constraint



PDF

Background

Cutoffs

Excess Constraint

True Positives

Signal (FG)

Conservation

False Positives

- Actual distribution of conservation scores (Signal) vs. expected distribution if no constraint (Background).

- At any cutoff: true positives (TP) and false predictions (FP)

- Can't **detect** all constrained elements since curves overlap

- But we can **estimate** the total amount of excess constraint by integrating over entire area between the two curves

16

# Detection of evolutionarily constrained elements



Most new elements in intronic/intergenic regions

Highest enrichment for coding transcripts

**Excess positive/purifying selection   Distribution of constraint**

# Coverage depth higher in functional regions



Legend:
- Exons, $\mu = 20.9$
- AR, $\mu = 11.4$
- HMRD elements, $\mu = 23.9$
- 29way elements, $\mu = 24.3$
- New elements, $\mu = 24.4$
- Masked genomic, $\mu = 17.1$

Y-axis: Frequency (0, 0.02, 0.04, 0.06, 0.08, 0.1, 0.12, 0.14, 0.16, 0.18)

X-axis: Number of aligned species (0, 5, 10, 15, 20, 25, 30)

Challenges of low-coverage genomes: varying aligment depth
Evidence of selection against deletions in functional regions

# Increase in power from HMRD to 29 mammals

| | π log-odds (12mers) | π log-odds (50mers) | ω (12mers) | ω (50mers) |
|---|---|---|---|---|
| 29 mammals | 7.1/1.5/4.6 | 6.8/1.8/4.1 | 5.7/ 1.1/3.8 | 5.7/1.8/3.0 |
| (HMRD) Human Mouse Rat Dog | 4.2/0.0/0.0 | 5.3/0.1/0.3 | 4.5/0.0/0.0 | 5.1/0.6/1.7 |

Estimated / kmers detectable at 5% FDR / base pairs detectable at 5% FDR

Small increase in estimate of genome percentage under constraint
Dramatic increase in power to detect small constrained elements

**Manuel Garber, Or Zuk** 19
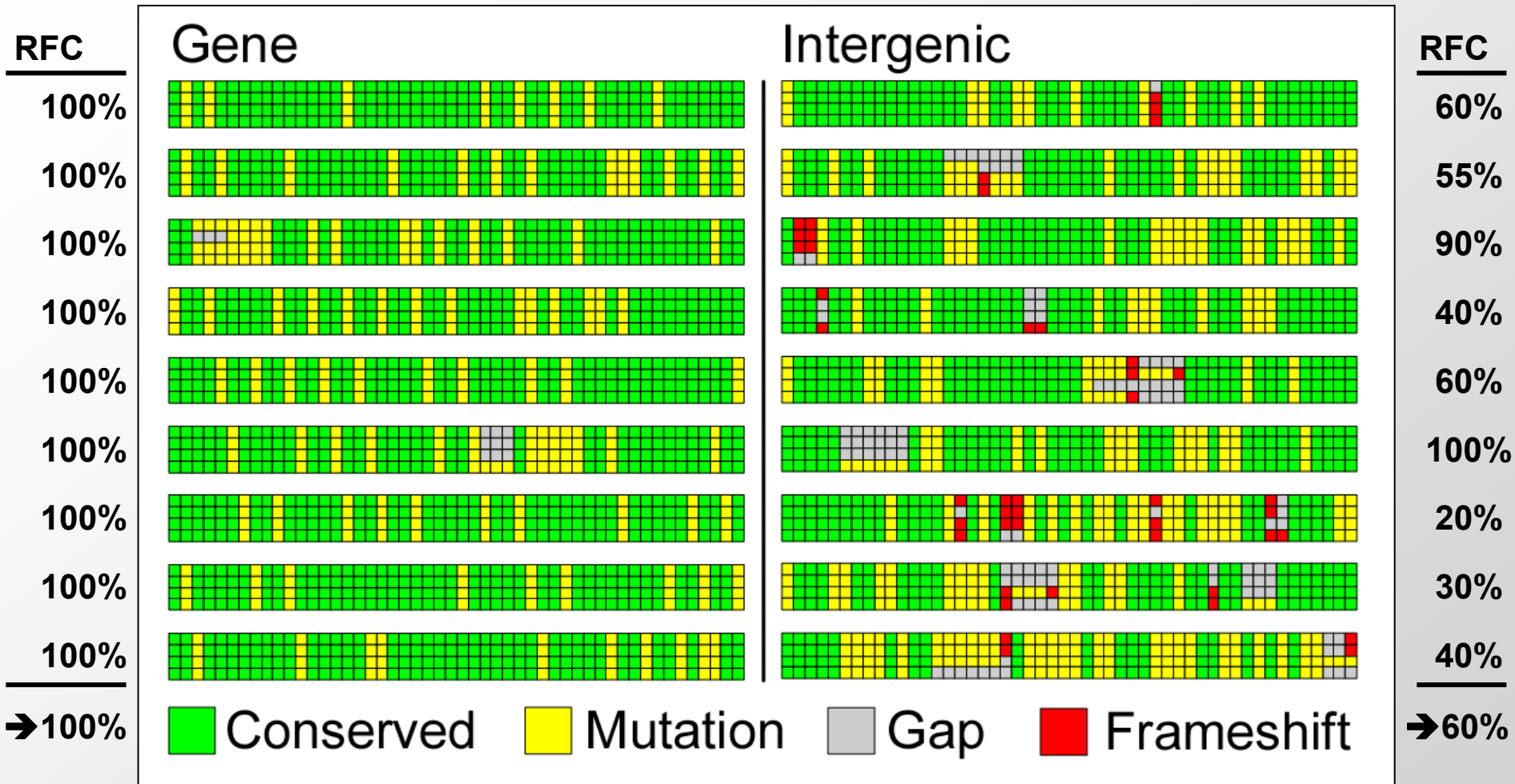
# Comparative genomics I: Evolutionary signatures

- **Nucleotide conservation: evolutionary constraint**
  - Purifying selection, neutral branch length, discovery power
  - Detect constrained elements: nucleotides, windows, HMM
  - Estimate fraction constrained: signal vs. background
- **Evolutionary signatures: focus on pattern of change**
  - Different functions ⇔ Characteristic patterns of evolution
- **Signatures of protein-coding genes**
  - Reading-frame conservation, codon-substitution frequency
  - Likelihood ratio framework: Estimating $Q_C Q_N$, scoring
  - Revise genes, read-through, excess constraint regions
- **Signatures of microRNA genes**
  - Structural and evolutionary features of microRNAs
  - Combining features: decision trees, random forests
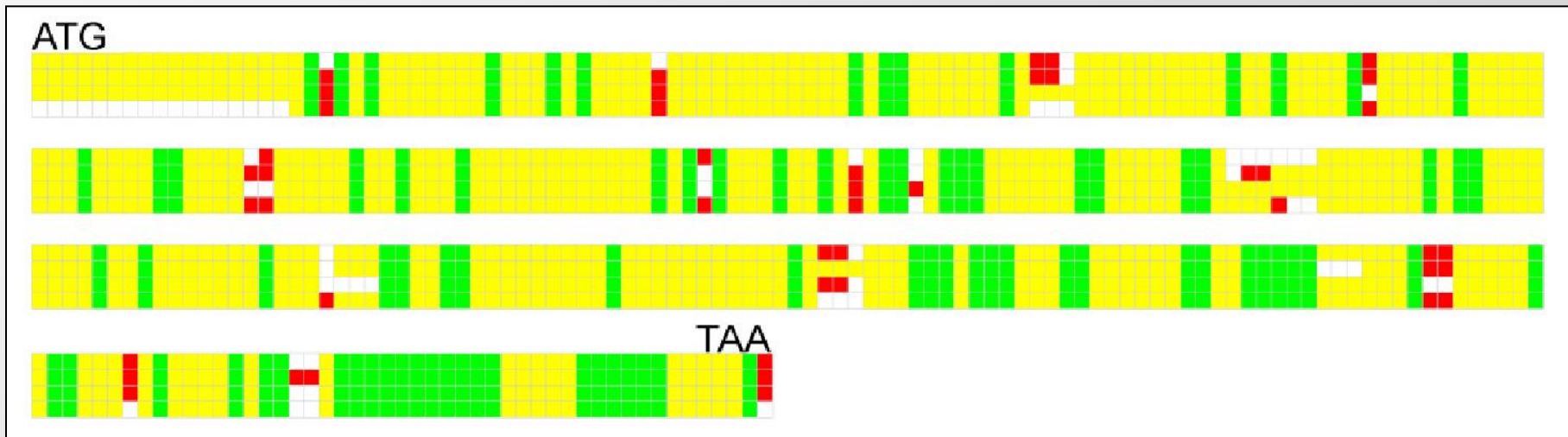  - Sense/anti-sense miRNAs, mature/star arm cooperation

# Comparative genomics and evolutionary signatures

- **Comparative genomics can reveal functional elements**
  - For example:  exons are deeply conserved to mouse, chicken, fish
  - Many other elements are also strongly conserved: exons / regulatory?

- **Can we also pinpoint specific functions of each region?  Yes!**
  - Patterns of change distinguish different types of functional elements
  - Specific function ⇔ Selective pressures ⇔ Patterns of mutation/inse/del

- **Develop evolutionary signatures characteristic of each function**

**Stark *et al*, Nature 2007**

# Evolutionary signatures for diverse functions



**Protein-coding genes**

- Codon Substitution Frequencies
- Reading Frame Conservation

**RNA structures**

- Compensatory changes
- Silent G-U substitutions

**microRNAs**

- Shape of conservation profile
- Structural features: loops, pairs
- Relationship with 3'UTR motifs

**Regulatory motifs**

- Mutations preserve consensus
- Increased Branch Length Score
- Genome-wide conservation

Source: Stark, Alexander et al. "Discovery of functional elements in 12 Drosophila genomes using evolutionary signatures." Nature 450, no. 7167 (2007): 219-232.

**Stark et al, Nature 2007**

# Implications for genome annotation / regulation



**Novel protein-coding genes**
**Revised gene annotations**
**Unusual gene structures**

**Novel structural families**
**Targeting, editing, stability**
**Riboswitches in mammals**

**Novel/expanded miR families**
**miR/miR* arm cooperation**
**Sense/anti-sense miR switches**

**Novel regulatory motifs**
**Regulatory motif instances**
**TF/miRNA regulatory networks**
**Single binding site resolution**

**Stark et al, Nature 2007**  23

# Comparative genomics I: Evolutionary signatures

- **Nucleotide conservation: evolutionary constraint**
  - Purifying selection, neutral branch length, discovery power
  - Detect constrained elements: nucleotides, windows, HMM
  - Estimate fraction constrained: signal vs. background
- **Evolutionary signatures: focus on pattern of change**
  - Different functions ⇔ Characteristic patterns of evolution
- **Signatures of protein-coding genes**
  - Reading-frame conservation, codon-substitution frequency
  - Likelihood ratio framework: Estimating $Q_C Q_N$, scoring
  - Revise genes, read-through, excess constraint regions
- **Signatures of microRNA genes**
  - Structural and evolutionary features of microRNAs
  - Combining features: decision trees, random forests
  - Sense/anti-sense miRNAs, mature/star arm cooperation

# Evolutionary signatures for protein-coding genes



**Frame-shifting indels**

Periodic mutations

Synonymous substs.

- **Same conservation levels, distinct patterns of divergence**
  - Gaps are multiples of three (preserve amino acid translation)
  - Mutations are largely 3-periodic (silent codon substitutions)
  - Specific triplets exchanged more frequently (conservative substs.)
  - Conservation boundaries are sharp (pinpoint individual splicing signals)

➔ **Evolutionary signatures of protein-coding selection**

# Evolutionary signatures of protein-coding genes

|      | the | fat | cat | sat |
|------|-----|-----|-----|-----|
| Δ1   | the | atc | ats | at  |
| Δ2   | the | tca | tsa | t   |
| Δ3   | the | cat | sat |     |



Second Letter

DNA insertions and deletions can either insert/remove AAs, or totally mangle the remainder of the protein (frameshift).

Some point mutations to the DNA sequence do not change its protein translation at all.

Natural selection tends to tolerate mutations with little/no effect on the protein.

# Protein-coding sequences tolerate distinctive types of change



protein-coding exon

conserved non-coding sequence

synonymous

conservative

non-conservative

frame-shifted

three stop codons

# Known genes stand out

# Comparative genomics I: Evolutionary signatures

- **Nucleotide conservation: evolutionary constraint**
  - Purifying selection, neutral branch length, discovery power
  - Detect constrained elements: nucleotides, windows, HMM
  - Estimate fraction constrained: signal vs. background
- **Evolutionary signatures: focus on pattern of change**
  - Different functions ⇔ Characteristic patterns of evolution
- **Signatures of protein-coding genes**
  - Reading-frame conservation, codon-substitution frequency
  - Likelihood ratio framework: Estimating $Q_C Q_N$, scoring
  - Revise genes, read-through, excess constraint regions
- **Signatures of microRNA genes**
  - Structural and evolutionary features of microRNAs
  - Combining features: decision trees, random forests
  - Sense/anti-sense miRNAs, mature/star arm cooperation

# Signature 1: Reading frame conservation

# Reading Frame Conservation Test

```
Scer      CTTCTAGATTTTCATCTT-GTCGATGTTCAAACAACGTGTTA-----TCAGAGAAACAGCTCTATGAGAAATCAGCTGATG
Scer_f1   123123123123123123-12312312312312312312312-----31231231231231231231231231231231231 23

Spar      TATTCATA-TCTCATCTTCATCAATGTTCAAACAGCGTGTTACAGACACAGAGAAACAGCTTC-TGAGAAGTCAGCCGGTG   RFC
Spar_f1   12312312-312312312312312312312312312312312312312312312312312312-31231231231231231 → 43%
Spar_f2   23123123-123123123123123123123123123123123123123123123123123123-12312312312312312 → 34%
Spar_f3   31231231-231231231231231231231231231231231231231231231231231231-23123123123123123 → 23%
```

←F1→  ←F2→  ←————— F1 —————→  ←——— F2 ———→  ←——— F3 ———→



Gene    Intergenic

100%  60%
100%  60%
100%  90%
100%  40%
100%  60%
100%  100%
100%  30%
100%  30%
100%  30%

🟩 Conserved  🟨 Mutation  ⬜ Gap  🟥 Frameshift

**100%**  **56%**

# Revisiting gene content with RFC test

|  | Accept | Reject |
|---|---|---|
| ~4000 named genes | 99.9% | 0.1% |
| ~300 intergenic regions | 1% | 99% |
| 2000 Hypothetical ORFs | 1500 | 500 |

High sensitivity and specificity

Example of a rejected ORF

# Comparative genomics I: Evolutionary signatures

- **Nucleotide conservation: evolutionary constraint**
  - Purifying selection, neutral branch length, discovery power
  - Detect constrained elements: nucleotides, windows, HMM
  - Estimate fraction constrained: signal vs. background
- **Evolutionary signatures: focus on pattern of change**
  - Different functions ⇔ Characteristic patterns of evolution
- **Signatures of protein-coding genes**
  - Reading-frame conservation, codon-substitution frequency
  - Likelihood ratio framework: Estimating $Q_C Q_N$, scoring
  - Revise genes, read-through, excess constraint regions
- **Signatures of microRNA genes**
  - Structural and evolutionary features of microRNAs
  - Combining features: decision trees, random forests
  - Sense/anti-sense miRNAs, mature/star arm cooperation

protein-coding exon

conserved non-coding sequence

# A method to distinguish these evolutionary signatures should:

- **Quantify the distinctiveness of all $64^2$ possible codon substitutions**

  - Synonymous: very frequent in protein-coding sequences

  - Nonsense: much more frequent in non-coding than coding regions

- **Model the phylogenetic relationship among the species**

  - Multiple apparent substitutions may be explained by one evolutionary event

- **Tolerate uncertainty in the input**

  - Unknown ancestral sequences

  - Alignment gaps, missing data

- **Report the [un]certainty of the result**

  - Quantify confidence that given alignment is protein-coding

  - Units: p-value, bits, decibans, etc.

# Codon evolution can be modeled as a Bayesian network



Conditional probability distribution (CPD) giving,
for all codons a & b, $\Pr(\text{dyak} = b | \text{Ancestor} = a)$

Each site (codon alignment column) is treated independently.

Given the topology and CPDs, we can simulate evolution of an ancestral sequence.

Additionally given extant (leaf) sequences, the ancestral sequences can be inferred.

For *L* leaves, CPDs total about $(2L - 2) \cdot 64^2$ parameters.

# The Bayes net is parameterized as a continuous-time Markov process



Rate matrix (**Q**)

Branch lengths <u>t</u>

Each CPD is determined by a rate matrix shared throughout the tree and a branch-specific 'time' (branch length):

$$\Pr(\text{child} = b | \text{parent} = a; t) = \left[ e^{\mathbf{Q}t} \right]_{a,b}$$

<u>Intuition</u>: The branch lengths specify how much 'time' passed between any two nodes. The rate matrix describes the relative frequencies of codon substitutions *per unit branch length*. Synonymous substitutions have high rates and nonsense substitutions have low rates.

We can obtain maximum likelihood estimates of $(2L - 2) + 64^2$ parameters using EM in training data.

36

# Example nucleotide (4x4) rate & substitution matrices

$$\mathbf{Q} = \begin{pmatrix} -4 & 2 & 1 & 1 \\ 2 & -4 & 1 & 1 \\ 1 & 1 & -4 & 2 \\ 1 & 1 & 2 & -4 \end{pmatrix} \begin{matrix} A \\ G \\ C \\ T \end{matrix}$$
$$\quad\quad A \quad G \quad C \quad T$$

$$\mathrm{Pr}(\text{child} = b | \text{parent} = a; t) = [e^{\mathbf{Q}t}]_{a,b}$$

$e^{\mathbf{Q}t} = \sum_{n=0}^{\infty} \dfrac{t^n}{n!} \mathbf{Q}^n$ is the solution to the system of differential equations describing the Markov process model of evolution.

`MatrixExp[Q * 0] // MatrixForm`

$$\begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

`MatrixExp[Q * 0.001] // MatrixForm // NumberForm[#, 4] &`

$$\begin{pmatrix} 0.996 & 0.001993 & 0.000998 & 0.000998 \\ 0.001993 & 0.996 & 0.000998 & 0.000998 \\ 0.000998 & 0.000998 & 0.996 & 0.001993 \\ 0.000998 & 0.000998 & 0.001993 & 0.996 \end{pmatrix}$$

`MatrixExp[Q * 0.01] // MatrixForm // NumberForm[#, 4] &`

$$\begin{pmatrix} 0.9611 & 0.01932 & 0.009803 & 0.009803 \\ 0.01932 & 0.9611 & 0.009803 & 0.009803 \\ 0.009803 & 0.009803 & 0.9611 & 0.01932 \\ 0.009803 & 0.009803 & 0.01932 & 0.9611 \end{pmatrix}$$

`MatrixExp[Q * 0.1] // MatrixForm // NumberForm[#, 4] &`

$$\begin{pmatrix} 0.692 & 0.1432 & 0.08242 & 0.08242 \\ 0.1432 & 0.692 & 0.08242 & 0.08242 \\ 0.08242 & 0.08242 & 0.692 & 0.1432 \\ 0.08242 & 0.08242 & 0.1432 & 0.692 \end{pmatrix}$$

`MatrixExp[Q * 1.0] // MatrixForm // NumberForm[#, 4] &`

$$\begin{pmatrix} 0.2558 & 0.2533 & 0.2454 & 0.2454 \\ 0.2533 & 0.2558 & 0.2454 & 0.2454 \\ 0.2454 & 0.2454 & 0.2558 & 0.2533 \\ 0.2454 & 0.2454 & 0.2533 & 0.2558 \end{pmatrix}$$

`MatrixExp[Q * 10.0] // MatrixForm // NumberForm[#, 4] &`

$$\begin{pmatrix} 0.25 & 0.25 & 0.25 & 0.25 \\ 0.25 & 0.25 & 0.25 & 0.25 \\ 0.25 & 0.25 & 0.25 & 0.25 \\ 0.25 & 0.25 & 0.25 & 0.25 \end{pmatrix}$$

Analogy: $y(t) = e^{qt}$

solves the differential equation

$$\dfrac{dy}{dt} = qy$$

`Plot[Exp[-t], {t, 0, 5}]`

Side note: Jukes-Cantor and Kimura models are set up so that the entries of $e^{Qt}$ have closed-form solutions.

# The hairy math: how do we estimate *Q*?

- Collect many alignments of <u>known</u> protein-coding sequences (training data)

- Consider the probability of the training data as <u>a function of *Q*</u>

$$\text{Likelihood}(\mathbf{Q}) = \text{Pr}(\text{Training Data}; \mathbf{Q}, \underline{t})$$

Still computed using Felsenstein algorithm

- Choose the *Q* that maximizes that probability:

$$\hat{\mathbf{Q}} = \underset{\mathbf{Q}}{\text{argmax}}\left(\text{Likelihood}(\mathbf{Q})\right)$$

Note: *Q* represents thousands of parameters

- Maximization strategies: expectation-maximization; gradient ascent; simulated annealing; spectral decomposition; others

- Branch lengths can also be optimized in the same way (simultaneously)

- Non-coding model estimated similarly, with random non-coding regions as training data.

Given this generative model
of codon evolution:



Rate matrix (**Q**)

Branch lengths t

We can compute the probability of any given alignment, marginalizing over all possible ancestral sequences, using Felsenstein's pruning algorithm.



protein-coding exon

$$\mathrm{Pr}(\mathrm{Leaves}; \mathbf{Q}, \underline{t}) = \frac{1}{10^{117}}$$

conserved non-coding sequence

$$\mathrm{Pr}(\mathrm{Leaves}; \mathbf{Q}, \underline{t}) = \frac{1}{10^{275}}$$

If I simulate alignments randomly according to the model, I'll get this
exact alignment once every $10^{117}$ samples

Now suppose we've estimated two rate matrices:



$Q_C$ estimated from known coding regions

$Q_N$ estimated from non-coding regions

These specify different rates of codon substitution, which in turn lead to different probabilities of any given alignment:



$$\Pr(\text{Leaves}; \mathbf{Q}_C, \underline{t}) = \frac{1}{10^{117}}$$

$$\Pr(\text{Leaves}; \mathbf{Q}_N, \underline{t}) = \frac{1}{10^{152}}$$

$$\Pr(\text{Leaves}; \mathbf{Q}_C, \underline{t}) = \frac{1}{10^{275}}$$

$$\Pr(\text{Leaves}; \mathbf{Q}_N, \underline{t}) = \frac{1}{10^{254}}$$

$$\frac{\text{Pr}(\text{Leaves}; \mathbf{Q}_C, \underline{t}) = \dfrac{1}{10^{117}}}{\text{Pr}(\text{Leaves}; \mathbf{Q}_N, \underline{t}) = \dfrac{1}{10^{152}}} = 10^{35}$$

This alignment is $10^{35}$ times <u>more probable</u> under the coding model than the non-coding model.



$$\frac{\text{Pr}(\text{Leaves}; \mathbf{Q}_C, \underline{t}) = \dfrac{1}{10^{275}}}{\text{Pr}(\text{Leaves}; \mathbf{Q}_N, \underline{t}) = \dfrac{1}{10^{254}}} = 10^{-21}$$
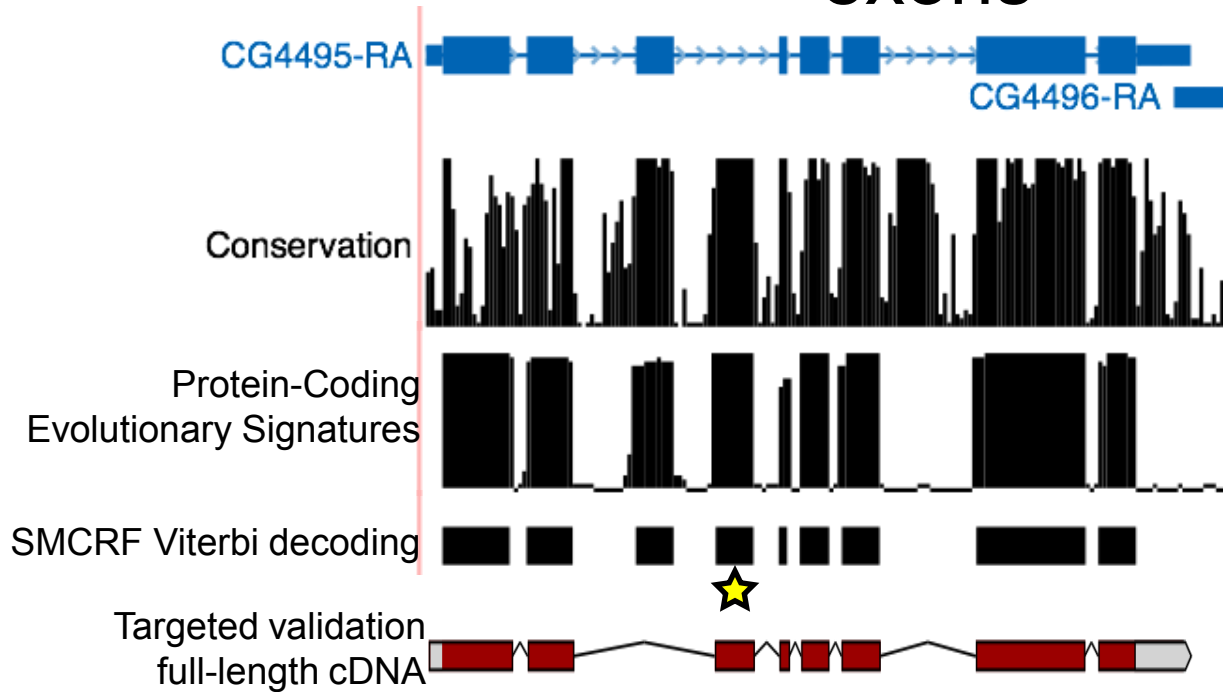
This alignment is $10^{21}$ times <u>less probable</u> under the coding model than the non-coding model.

This **likelihood ratio** $\dfrac{\text{Pr}(\text{Leaves}; \mathbf{Q}_C, \underline{t})}{\text{Pr}(\text{Leaves}; \mathbf{Q}_N, \underline{t})}$ is our measure of confidence that the alignment is protein-coding.
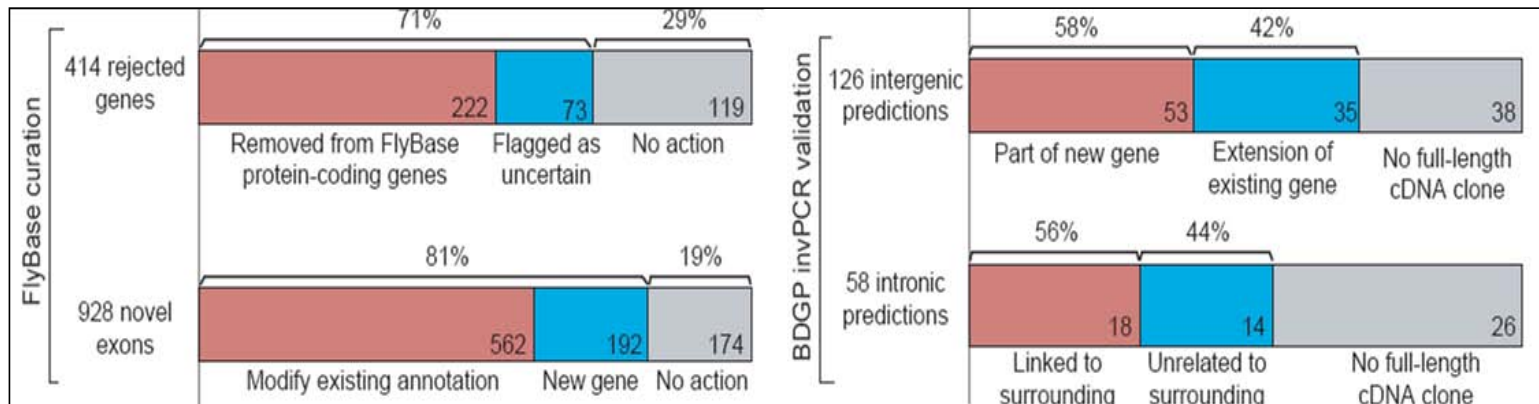
# Comparative genomics I: Evolutionary signatures

- **Nucleotide conservation: evolutionary constraint**
  - Purifying selection, neutral branch length, discovery power
  - Detect constrained elements: nucleotides, windows, HMM
  - Estimate fraction constrained: signal vs. background
- **Evolutionary signatures: focus on pattern of change**
  - Different functions ⇔ Characteristic patterns of evolution
- **Signatures of protein-coding genes**
  - Reading-frame conservation, codon-substitution frequency
  - Likelihood ratio framework: Estimating $Q_C Q_N$, scoring
  - Revise genes, read-through, excess constraint regions
- **Signatures of microRNA genes**
  - Structural and evolutionary features of microRNAs
  - Combining features: decision trees, random forests
  - Sense/anti-sense miRNAs, mature/star arm cooperation

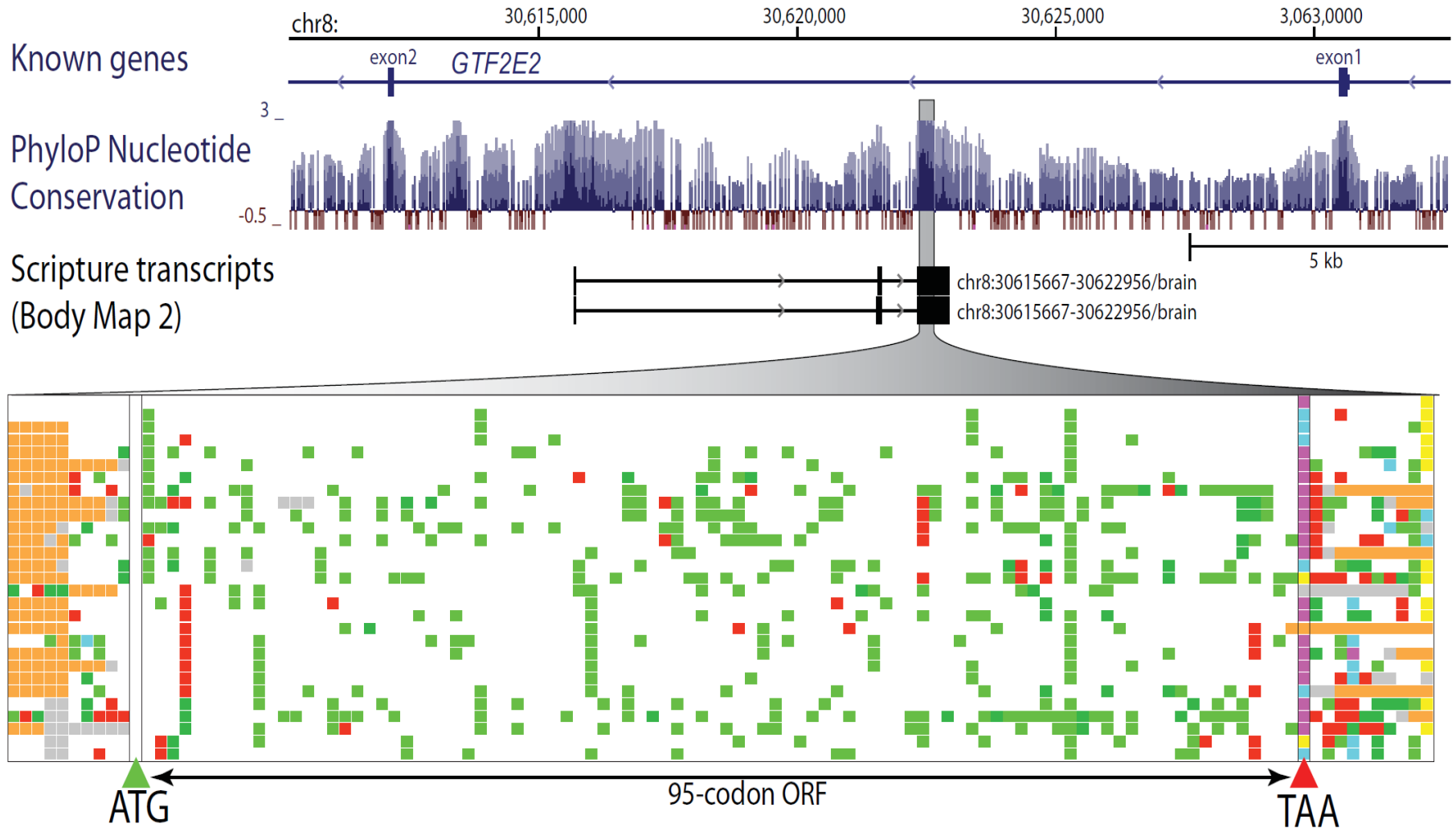# Evolutionary signatures can predict new genes and exons



*Evolutionary signatures built into a semi-Markov conditional random field to predict protein-coding exons*

Source: Stark, Alexander et al. "Discovery of functional elements in 12 Drosophila genomes using evolutionary signatures." Nature 450, no. 7167 (2007): 219-232.

# New protein-coding genes



Known genes

chr8: 30,615,000 · 30,620,000 · 30,625,000 · 3,063,0000

exon2 · GTF2E2 · exon1

PhyloP Nucleotide Conservation

3 —
-0.5 —

Scripture transcripts (Body Map 2)

chr8:30615667-30622956/brain
chr8:30615667-30622956/brain

5 kb

ATG

95-codon ORF

TAA

New genes supported by Illumina BodyAtlas transcripts
Submitted to GENCODE for validation / manual curation

44

# Translational read-through in flies and mammals

**One of four novel candidates in the human genome: OPRL1 neurotransmitter**



**Protein-coding conservation** | **Stop codon read through** | **Continued protein-coding conservation** | **2nd stop codon** | **No more conservation**

- ## New mechanism of post-transcriptional regulation?
  - Conserved in both mammals (4 candidates) and flies (350 candidates)
  - Strongly enriched for neurotransmitters, brain-expressed proteins, TF regulators
  - After correcting for gene length: TF enrichment remains

- ## Evidence suggestive of regulatory control
  - Read-through stop codon perfectly conserved in 93% of cases (24% at bkgrnd)
  - Upstream bases show increased conservation. Downstream is TGAC.
  - GCA triplet repeats
  - Increased RNA secondary structure

**Lin *et al*, Genome Research 2007**
**Jungreis *et al*, in preparation**
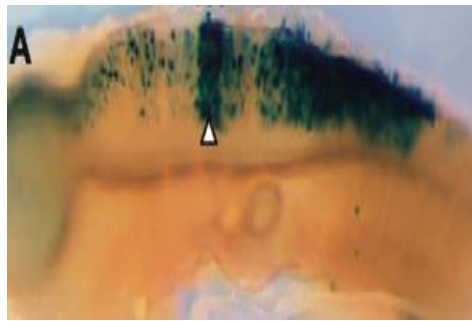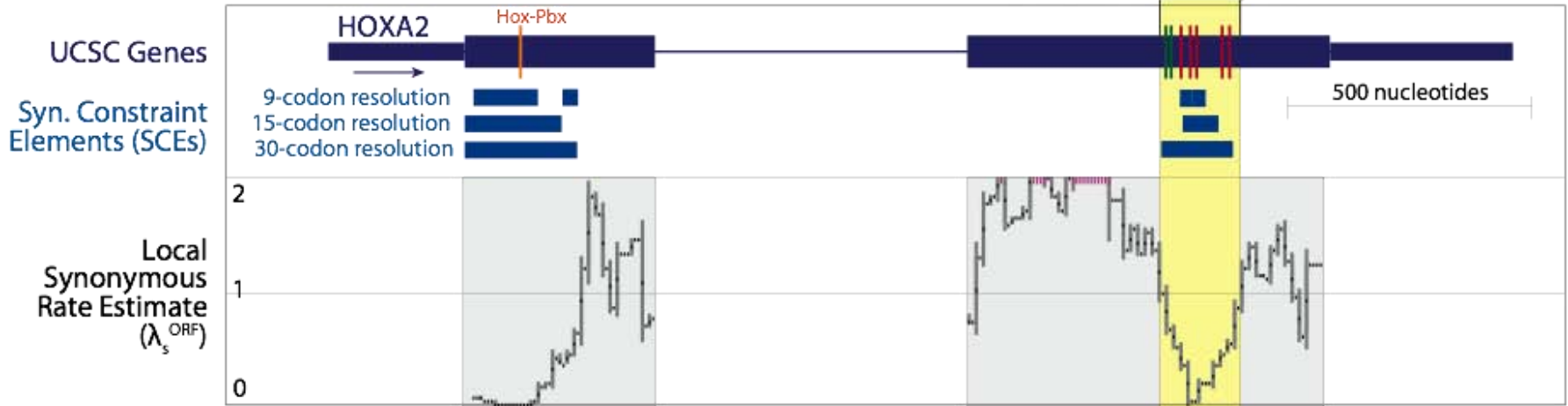
# Discover of translational readthrough genes

Discovery of 4 readthrough genes, abundant in many animal genomes

# Overlapping selection in protein-coding exons



rhombomere 4 expression
(Lampe *et al.*, NAR 2008)

rhombomere 2 expr.
(Tümpel PNAS 2008)

10,000 overlapping synonymous constrained elements
Roles in splicing, translation, regulation

# Codon-specific measures of positive selection

Gene-wide vs. punctate regions of exons positive selection

# Comparative genomics I: Evolutionary signatures

- **Nucleotide conservation: evolutionary constraint**
  - Purifying selection, neutral branch length, discovery power
  - Detect constrained elements: nucleotides, windows, HMM
  - Estimate fraction constrained: signal vs. background
- **Evolutionary signatures: focus on pattern of change**
  - Different functions ⇔ Characteristic patterns of evolution
- **Signatures of protein-coding genes**
  - Reading-frame conservation, codon-substitution frequency
  - Likelihood ratio framework: Estimating $Q_C Q_N$, scoring
  - Revise genes, read-through, excess constraint regions
- **Signatures of microRNA genes**
  - Structural and evolutionary features of microRNAs
  - Combining features: decision trees, random forests
  - Sense/anti-sense miRNAs, mature/star arm cooperation

# New RNA structures and families

| | No. of structures | No. of novel structures | No. of families | No. of novel families | EvoFold score | RNAz overlap enrichment (x) | DNAse hypersensitivity overlap (%) | Avg. correlation of tissue-specific expression within families |
|---|---|---|---|---|---|---|---|---|
| EvoFold all (no CDS) | 27,012 | 26,643 | n/a | n/a | 14 | 13.5 | 25 ($P \leq 5e{-}3$) | n/a |
| Unfiltered families | 3293 | 3081 | 1254 | 1215 | 18 | 17.3 | 25 ($P \leq 7e{-}3$) | 0.14 ($P \leq 1e{-}3$) |
| Filtered families | 725 | 526 | 220 | 181 | 18 | 29 | 32 ($P \leq 4e{-}3$) | 0.15 ($P \leq 1e{-}3$) |

# New structs fall in families, supported by evolut/energy

Ex: new struct in XIST long non-coding RNA
Known function in X-chromosome inactivation
Possible functional domain of XIST?

# RNA families: orthologous/paralogous conservation

Example of new structural 3'UTR family in MAT2A gene likely role in detecting S-adeosyl-methionic (SAM) level

# Computational challenge of miRNA discovery

**760,355
miRNA-like hairpins**

⟷

**60-100
true miRNAs**

A false positive rate of 0.5% ➔ 3800 spurious hairpins.
Need 99.99% specificity (>5,000-fold enrichment)

# Evolutionary signatures for microRNA genes



**(1) Conservation profile**

## miRNAs show characteristic conservation properties

# Distinguishing true miRNAs from random hairpins



**Evolutionary features**

(1) Correlation with conservation profile (−1 to 1)

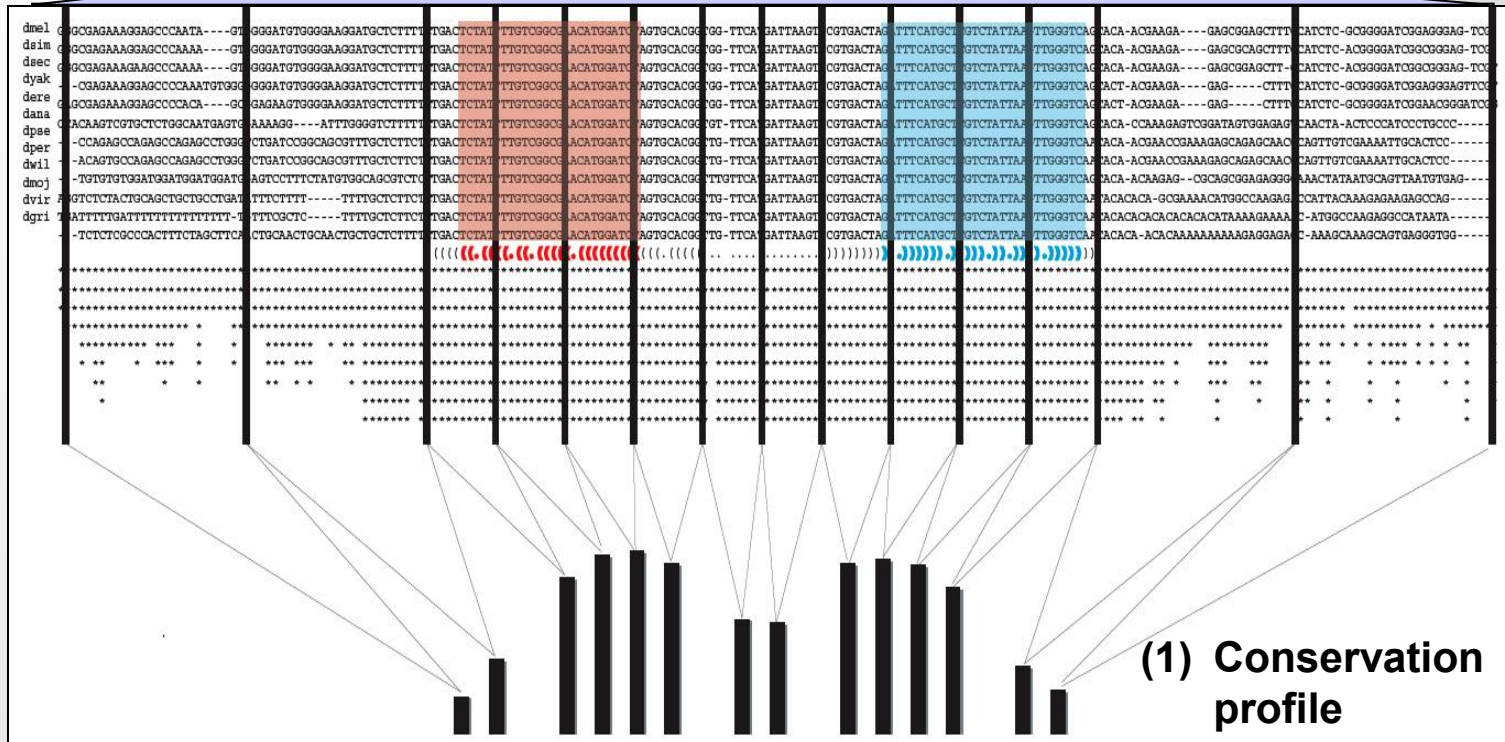(2) MFE of the consensus fold (80 to 10)

(3) Structure conservation index (0 to 2)

**Structural features**

(4) Hairpin stability (MFE z-score) (6 to 14)

(5) Number of asymmetric loops (0 to 14)

(6) Number of symmetric loops (0 to 16)

**Feature performance**

| | | Enrichment |
|---|---|---|
| Total | 51 / 142 | 4551 |
| Cons. Profile | 42 / 1,625 | 327 |
| Arm / Loop Cons. | 36 / 24,574 | 19 |
| Arm Cons. | 52 / 131,707 | 5 |
| Structure Cons. | 55 / 110,295 | 6 |
| Hairpin Energy | 43 / 13,885 | 39 |
| % Paired Bases | 48 / 262,543 | 2.3 |
| Arm length | 56 / 412,772 | 1.7 |
| Loop Length | 59 / 547,327 | 1.4 |
| Structure Length | 59 / 548,799 | 1.4 |
| Loop Symmetry | 35 / 147,352 | 3 |

■ Known (60)  ■ Random (760,355)

**Combination of features: > 4,500-fold enrichment**

# Comparative genomics I: Evolutionary signatures

- **Nucleotide conservation: evolutionary constraint**
  - Purifying selection, neutral branch length, discovery power
  - Detect constrained elements: nucleotides, windows, HMM
  - Estimate fraction constrained: signal vs. background
- **Evolutionary signatures: focus on pattern of change**
  - Different functions ⇔ Characteristic patterns of evolution
- **Signatures of protein-coding genes**
  - Reading-frame conservation, codon-substitution frequency
  - Likelihood ratio framework: Estimating $Q_C Q_N$, scoring
  - Revise genes, read-through, excess constraint regions
- **Signatures of microRNA genes**
  - Structural and evolutionary features of microRNAs
  - Combining features: decision trees, random forests
  - Sense/anti-sense miRNAs, mature/star arm cooperation

# miRNA detection using many decision trees

MFE<3?

yes     no

ProfileCorr<8     StrConsIndx>3

yes   no     yes   no

**NOT**    **miRNA**     **NOT**    #Loops>2

yes   no

Stability>4    **NOT**

yes   no

**NOT**    #Loops<5

yes   no

**miRNA**    **NOT**

- **For each tree:**
  – Randomly select:
    - Subset of features to base classification on
    - Subset of +/- training examples
    - Remainder of testing examples
  – Use to train a decision tree classifier:
    - Select a feature and cutoff at each level
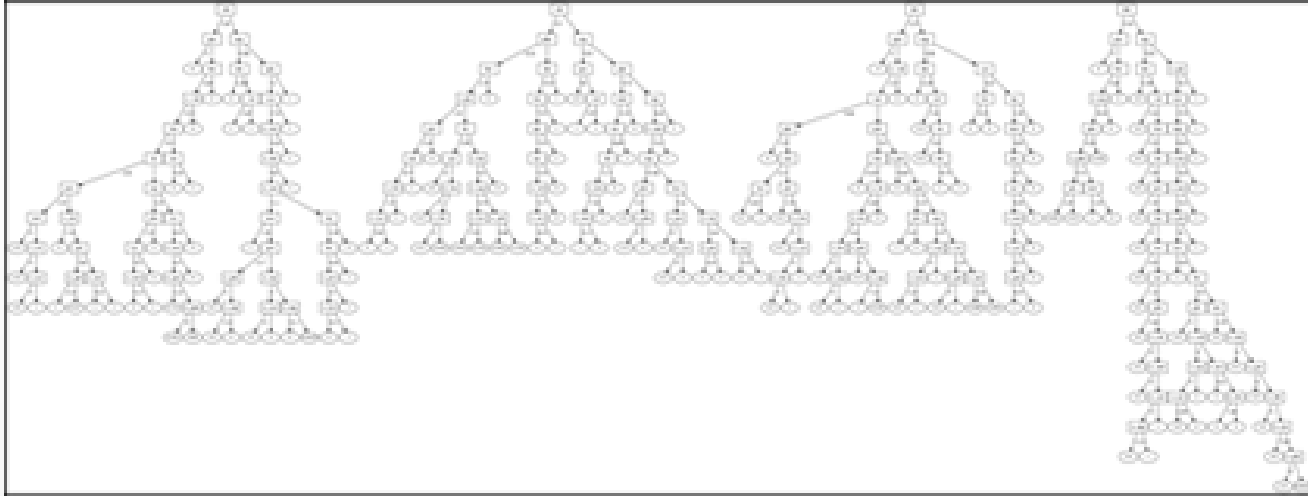    - Continue with feature/cutoff at next level
    - (…)
  – Evaluate performance on test set:
    - Push each element down the decision tree
    - Leaf label gives classification decision

- **To combine trees:**
  – Average prediction class across trees
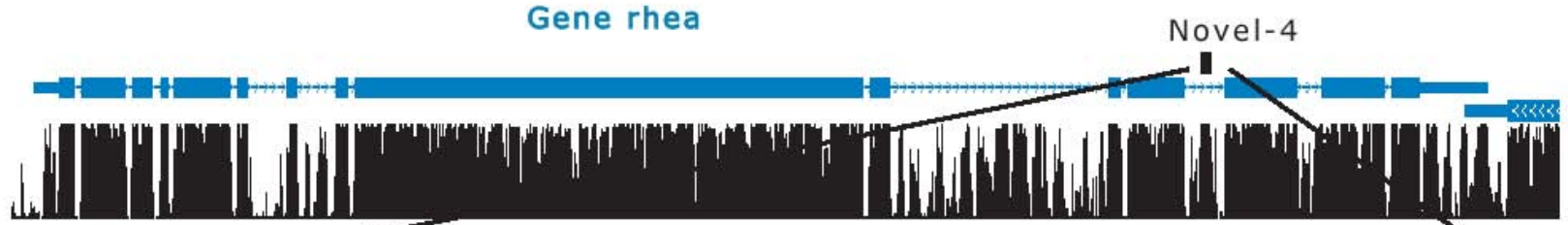  – Report class with maximum # of votes

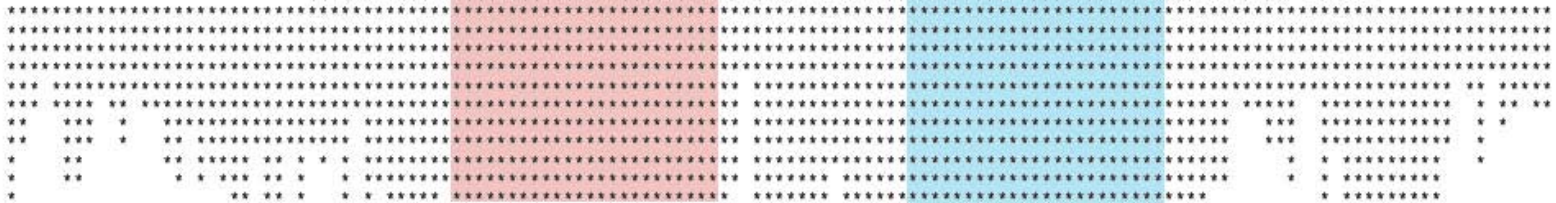# Random Forests: Combine many decision trees

- **Many decision trees:**
  - Each can select cutoffs and direction of cutoff
  - Each feature can be reused multiple times
  - Used serially (AND) and in parallel (OR)

- **Ensemble classifier**
  - Bagging: model averaging, combines predictions
  - Can take median of predictions

- **Advantages: Robustness, Feature importance**

**Ruby, Bartel, Lai**

# Evidence 2: Genomic properties typical of miRNAs



- **Novel miRNAs in introns of known genes**
- **Preference for + strand, transcription factors**



- **Genomic clustering with novel / known miRNAs**
- **Same family, common origin / same precursor**

# Two 'dubious' protein-coding genes are in fact miRNAs

## Two novel miRNAs overlap exons (5'UTR and coding!)

- **Both CG31044 and CG33311 were independently rejected as *dubious* based on their non-protein-coding conservation patterns (Lin *et al.*)**
- **Novel miRNA genes provide explanation for their transcripts, as their precursor miRNA**

# Comparative genomics I: Evolutionary signatures

- **Nucleotide conservation: evolutionary constraint**
  - Purifying selection, neutral branch length, discovery power
  - Detect constrained elements: nucleotides, windows, HMM
  - Estimate fraction constrained: signal vs. background
- **Evolutionary signatures: focus on pattern of change**
  - Different functions ⇔ Characteristic patterns of evolution
- **Signatures of protein-coding genes**
  - Reading-frame conservation, codon-substitution frequency
  - Likelihood ratio framework: Estimating $Q_C Q_N$, scoring
  - Revise genes, read-through, excess constraint regions
- **Signatures of microRNA genes**
  - Structural and evolutionary features of microRNAs
  - Combining features: decision trees, random forests
  - Sense/anti-sense miRNAs, mature/star arm cooperation

# Surprise 1: microRNA & microRNA* function

- **Both hairpin arms of a microRNA can be functional**
  - High scores, abundant processing, conserved targets
  - Hox miRNAs miR-10 and miR-iab-4 as master Hox regulators

**Stark *et al*, Genome Research 2007** 62

# Evidence of miR-iab-4 anti-sense (AS) function



© Cold Spring Harbor Laboratory Press. All rights reserved. This content is excluded from our Creative Commons license. For more information, see http://ocw.mit.edu/help/faq-fair-use/.

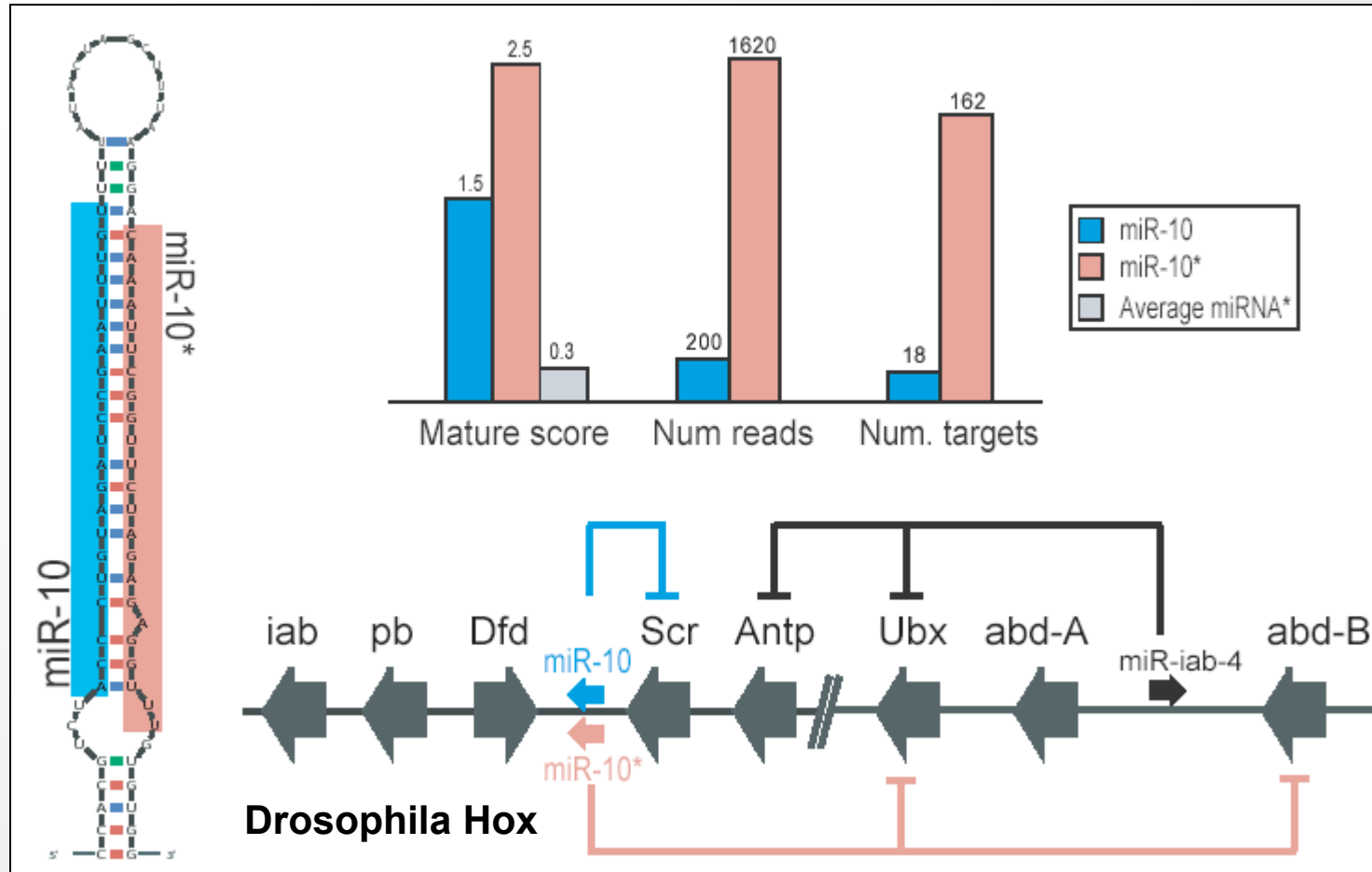Source: Stark, Alexander et al. "A single Hox locus in Drosophila produces functional microRNAs from opposite DNA strands." Genes & development 22, no. 1 (2008): 8-13.

- **A single miRNA locus transcribed from both strands**
- **The two transcripts show distinct expression domains (mutually exclusive)**
- **Both processed to mature miRNAs: mir-iab-4, miR-iab-4AS (anti-sense)**

# miR-iab-4AS leads to homeotic transformations



A →wing haltere

B B' Sensory bristles

Bx>mir-iab-4 anti-sense w1118

C haltere WT

w1118

D →wing sense

Bx>mir-iab-4 sense

E →wing w/bristles Antisense

Bx>mir-iab-4 anti-sense

**Note: C,D,E same magnification**

Source: Stark, Alexander et al. "A single Hox locus in Drosophila produces functional microRNAs from opposite DNA strands." Genes & development 22, no. 1 (2008): 8-13.

- **Mis-expression of mir-iab-4S & AS: alteres→wings homeotic transform.**

- **Stronger phenotype for AS miRNA**

- **Sense/anti-sense pairs as general building blocks for miRNA regulation**

- **10 sense/anti-sense miRNAs in mouse**

**Stark *et al*, Genes&Development 2008**

# Comparative genomics I: Evolutionary signatures

- **Nucleotide conservation: evolutionary constraint**
  - Purifying selection, neutral branch length, discovery power
  - Detect constrained elements: nucleotides, windows, HMM
  - Estimate fraction constrained: signal vs. background
- **Evolutionary signatures: focus on pattern of change**
  - Different functions ⇔ Characteristic patterns of evolution
- **Signatures of protein-coding genes**
  - Reading-frame conservation, codon-substitution frequency
  - Likelihood ratio framework: Estimating $Q_C Q_N$, scoring
  - Revise genes, read-through, excess constraint regions
- **Signatures of microRNA genes**
  - Structural and evolutionary features of microRNAs
  - Combining features: decision trees, random forests
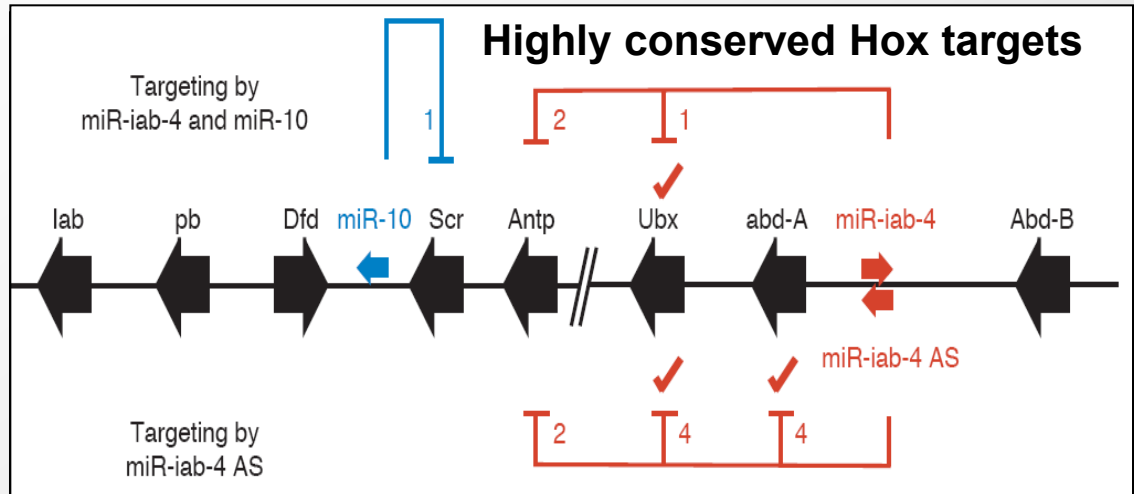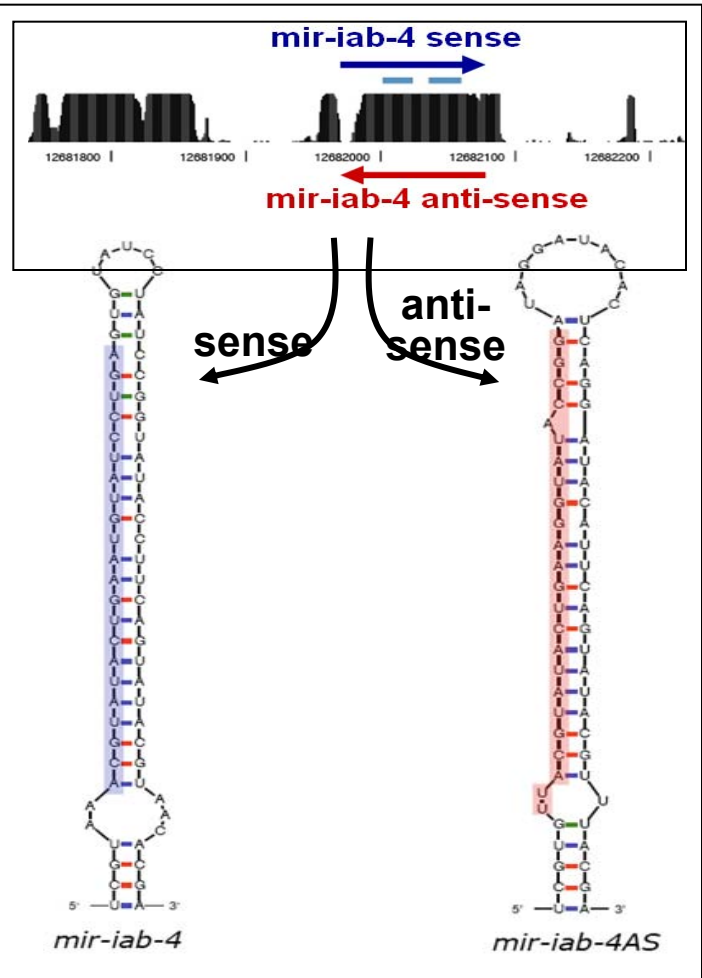  - Sense/anti-sense miRNAs, mature/star arm cooperation
- **Measuring selection within the human lineage**

# Mammalian constraint matches Human SNPs

Human SNPs match mammalian-wide twofold constraint

# Mammalian constraint matches Human SNPs

Genome-wide agreement of selection and polymorphisms

# Human constraint outside conserved regions

**Active regions**

Conserved

5.9

4.0

4.8

6.8

Average diversity (heterozygosity)

Aggregate over the genome

- **Non-conserved regions:**
  – ENCODE-active regions show reduced diversity
  ➔ Lineage-specific constraint in biochemically-active regions

- **Conserved regions:**
  – Non-ENCODE regions show increased diversity
  ➔ Loss of constraint in human when biochemically-inactive

# Strongest: motifs, short RNA, Dnase, ChIP, IncRNA



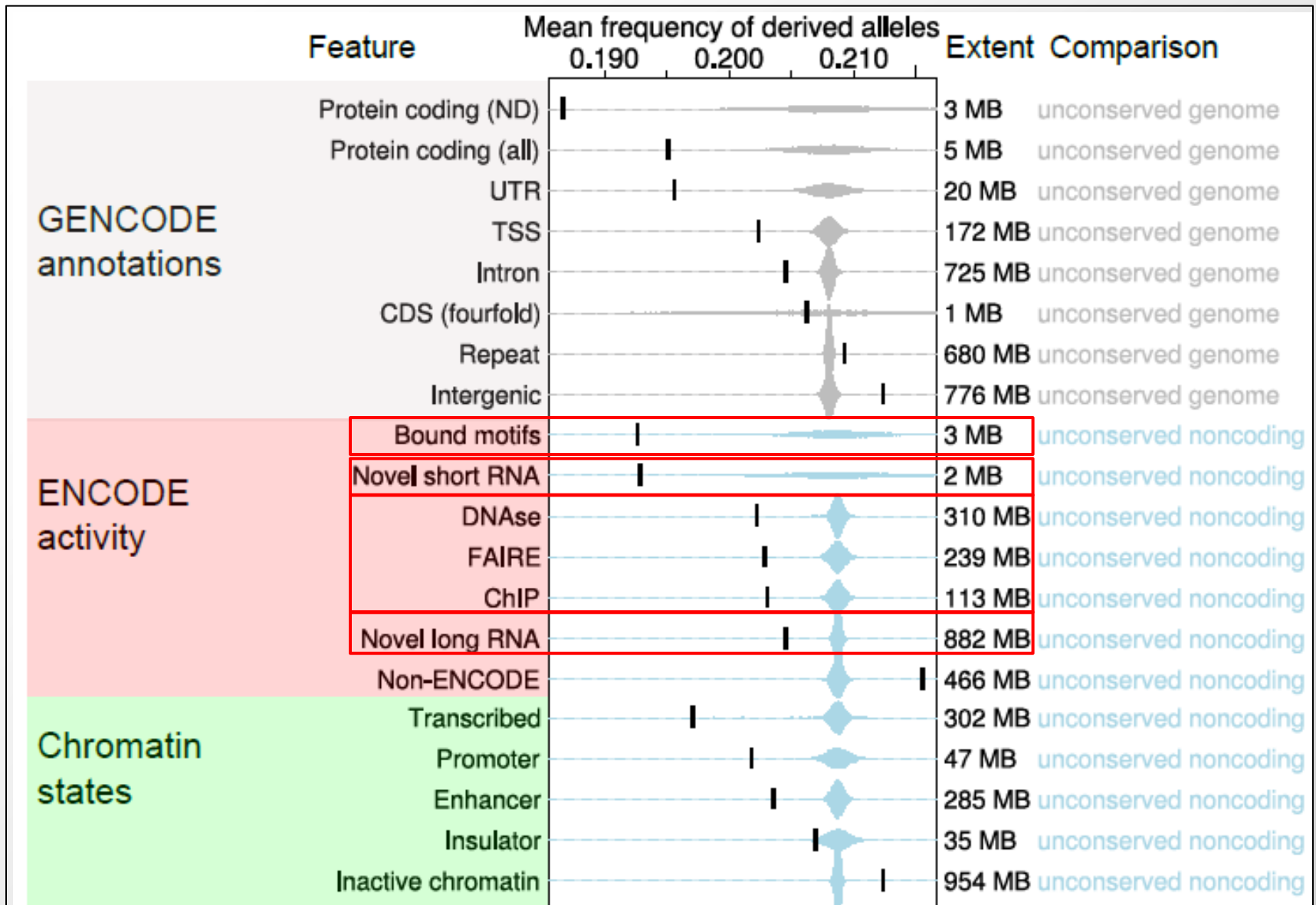| Feature | Mean frequency of derived alleles 0.190 0.200 0.210 | Extent | Comparison |
|---|---|---|---|
| **GENCODE annotations** | | | |
| Protein coding (ND) | | 3 MB | unconserved genome |
| Protein coding (all) | | 5 MB | unconserved genome |
| UTR | | 20 MB | unconserved genome |
| TSS | | 172 MB | unconserved genome |
| Intron | | 725 MB | unconserved genome |
| CDS (fourfold) | | 1 MB | unconserved genome |
| Repeat | | 680 MB | unconserved genome |
| Intergenic | | 776 MB | unconserved genome |
| **ENCODE activity** | | | |
| Bound motifs | | 3 MB | unconserved noncoding |
| Novel short RNA | | 2 MB | unconserved noncoding |
| DNAse | | 310 MB | unconserved noncoding |
| FAIRE | | 239 MB | unconserved noncoding |
| ChIP | | 113 MB | unconserved noncoding |
| Novel long RNA | | 882 MB | unconserved noncoding |
| Non-ENCODE | | 466 MB | unconserved noncoding |
| **Chromatin states** | | | |
| Transcribed | | 302 MB | unconserved noncoding |
| Promoter | | 47 MB | unconserved noncoding |
| Enhancer | | 285 MB | unconserved noncoding |
| Insulator | | 35 MB | unconserved noncoding |
| Inactive chromatin | | 954 MB | unconserved noncoding |

- **Significant derived allele depletion in active features**

# Bound motifs show increased human constraint

Position-specific reduction in bound motif heterozygosity
Aggregate across thousands of CTCF motif instances

# Most constrained human-specific enhancer functions



Transcription initiation from Pol2 promoter
Transcription coactivator activity
Transcription factor binding
Chromatin binding
Negative regulation of transcription, DNA-dependent
Transcription factor complex
Protein complex
Protein kinase activity
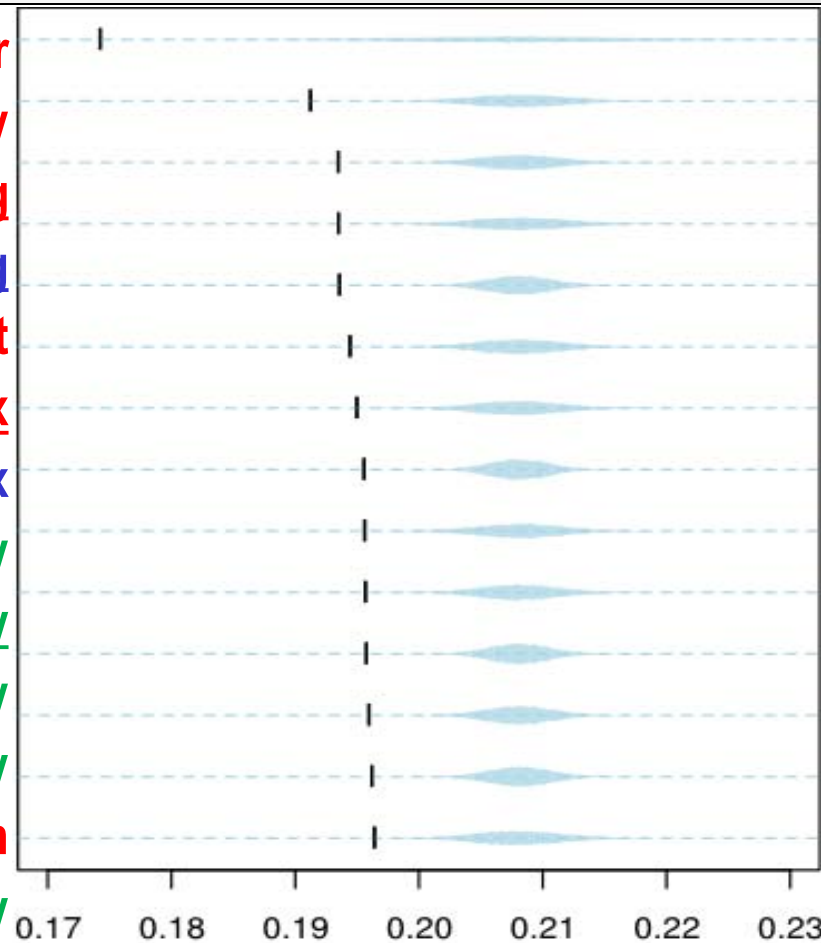Nerve growth factor receptor signaling pathway
Signal transducer activity
Protein serine/threonine kinase activity
Negative regulation of transcription from Pol2 prom
Protein tyrosine kinase activity
In utero embryonic development

0.17  0.18  0.19  0.20  0.21  0.22  0.23

**Regulatory genes: Transcription, Chromatin, Signaling.
Developmental enhancers: embryo, nerve growth**

# Comparative genomics I: Evolutionary signatures

- **Nucleotide conservation: evolutionary constraint**
  - Purifying selection, neutral branch length, discovery power
  - Detect constrained elements: nucleotides, windows, HMM
  - Estimate fraction constrained: signal vs. background
- **Evolutionary signatures: focus on pattern of change**
  - Different functions ⇔ Characteristic patterns of evolution
- **Signatures of protein-coding genes**
  - Reading-frame conservation, codon-substitution frequency
  - Likelihood ratio framework: Estimating $Q_C Q_N$, scoring
  - Revise genes, read-through, excess constraint regions
- **Signatures of microRNA genes**
  - Structural and evolutionary features of microRNAs
  - Combining features: decision trees, random forests
  - Sense/anti-sense miRNAs, mature/star arm cooperation
- **Measuring selection within the human lineage**

6.047 / 6.878 / HST.507 Computational Biology
Fall 2015