

WS2018/2019

„Genomforschung und Sequenzanalyse

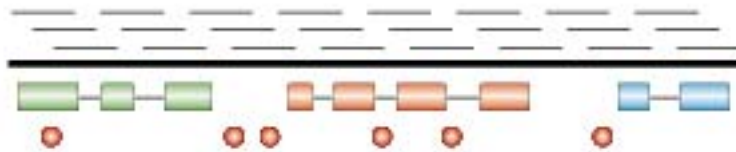
- Einführung in Methoden der Bioinformatik- “

Thomas Hankeln



Genvorhersage & Genom- Annotation

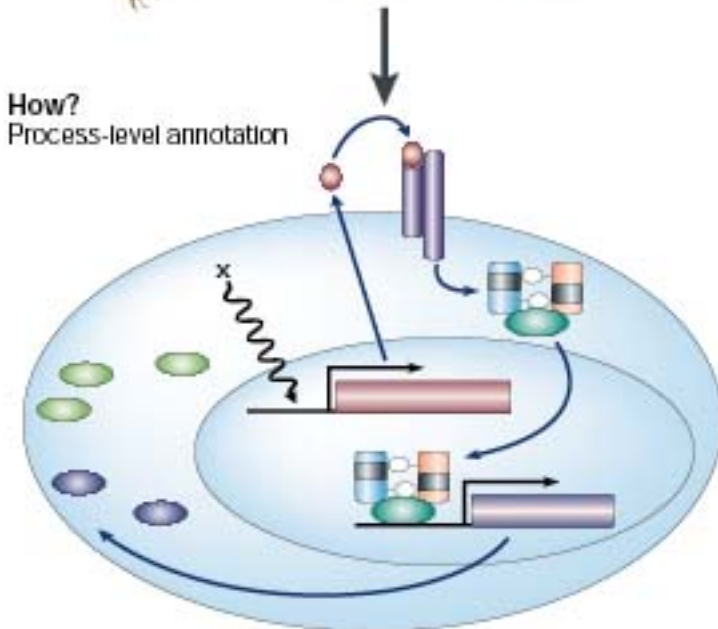
Where?
Nucleotide-level annotation



What?
Protein-level annotation



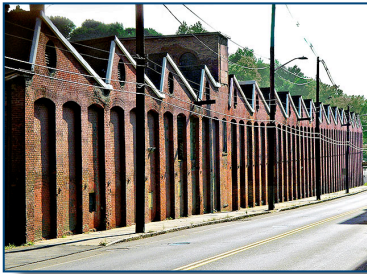
How?
Process-level annotation



Ebenen der Annotation

- Genstruktur (Exons/Introns, UTR 's, Promoter, andere regulator. Elemente)
- Orthologe and Paraloge
- Gen-Varianten (SNPs etc.)
- Transkripte (Gewebe, Zelltypen?)
- Protein-Domänen und –motive,
- Kristallstruktur, Faltung
- RNA/Protein: Funktion, intrazell. Lokalisation, Aktivitätsmuster
- interagierende Proteine
- Stoffwechselwege, metabol. Produkte
- Krankheitsrelevanz
-

Annotationsstrategien

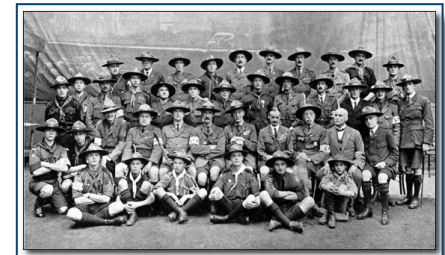


- „Factory“ > automatisch



- „museum“ > durch Kuratoren

- „party“ / jamboree



- „TPA“ = 3rd party annotation

„GeneOntology“: Ein gemeinsames Vokabular

Bsp: Cytochrome C

„Molecular function“
7393

electron transporter

„Biological process“
10370

oxidative
phosphorylation

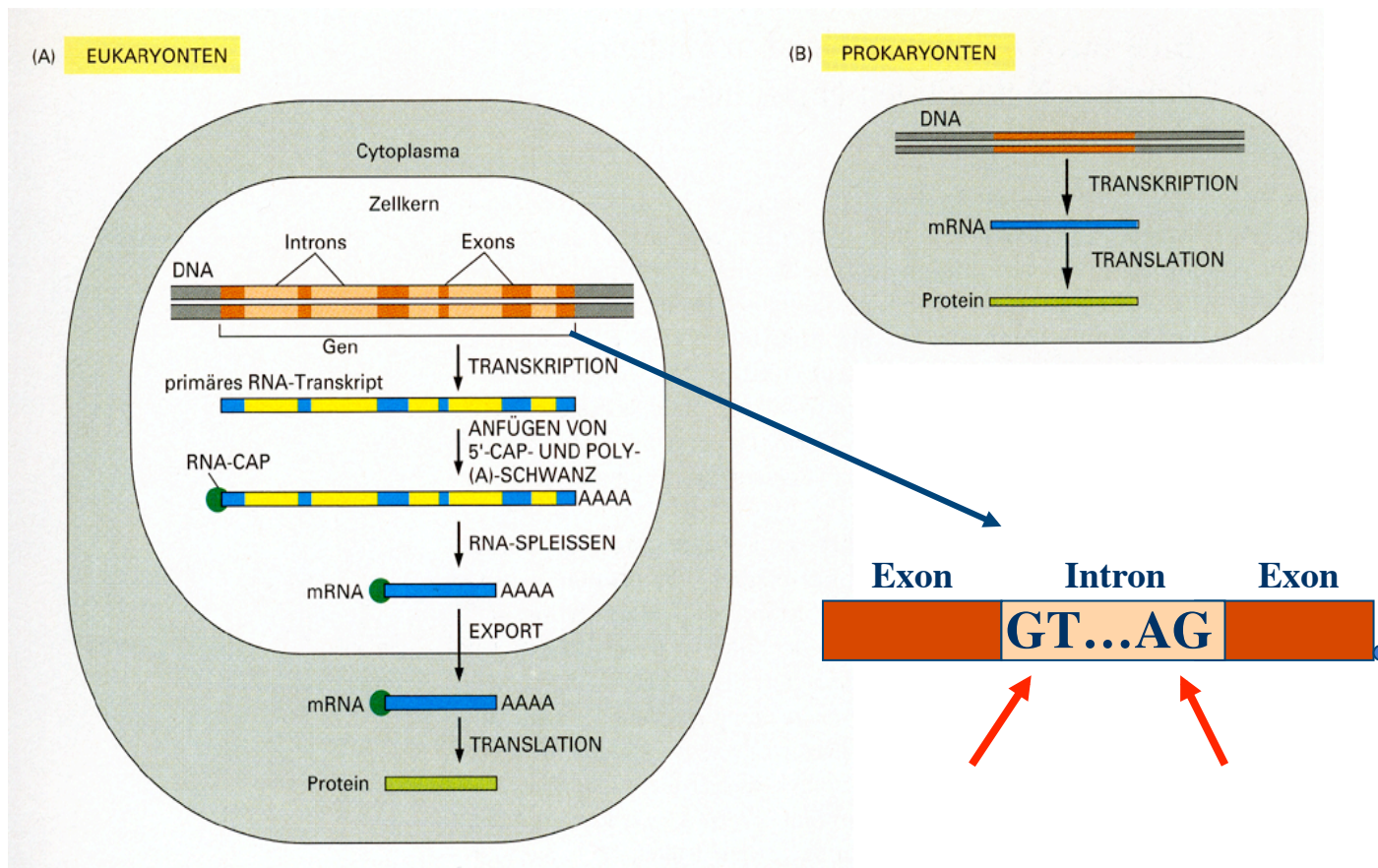
„Cellular component“
1681

mitochondrial
inner membrane

Wo steckt denn nun die genetische Information?

1	ccgaacgctt	atagagagct	atagagtga	agctgagaag	aaccaaaccg	gagcataaac
61	atgaacagcg	atgaggtgca	actgatcaag	aagacctggg	aatccccgt	ggcaacacca
121	acagattctg	gagcggcgat	actgacgcag	tttttcaacc	gctttccgtc	caacttggag
181	aagttcccct	tccgcgatgt	tcctttggag	gagctaagt	tgagttgtac	cttacacata
241	ggtcttcaat	taactcaaga	ttaacttgat	ctgttttctt	tcagggaaat	gctcgcttcc
301	gagcacatgc	cggcagaatc	ataaggggtct	ttgacgagtc	catccaggtc	ctgggccagg
361	atggcgatct	ggagaagctg	gacgagatct	ggaccaaact	tgccgttagt	cacattccgc
421	ggaccgtttc	caaggagtct	tacaacgtaa	ggtgaacact	gcagtcgagc	tctcgacttt
481	gagatacctg	ttggtcagat	agtggaagtt	gaaagctata	tgacatttaa	aaattcaatt
541	gcatttaaaa	catcatttta	tttttttttag	caactgaaag	gagttatcct	ggatgtgctg
601	acagctgcct	gcagtctgga	cgagagtcaa	gcggccacgt	ggccaagct	ggtggaccat
661	gtctacgcaa	tcattcttcaa	ggcgatcgac	gacgacggca	acgccaagta	gatgaggcag
721	ctggaggtgg	agatgcaacc	gaatccgcgg	a		

Bei Eukaryoten-Genomen ist Generkennung besonders schwierig



Die („vereinfachte“) Aufgabe:

- gegeben sind uncharakterisierte Genom-DNA-Sequenzen

- **FINDE...**

Protein-kodierende Regionen

Exon/Intron-Grenzen

mögliche genregulatorische Abschnitte

Mache daraus ein Modell für die Struktur des Gens!

Warum „vereinfacht“?

- nicht alle Gene werden in Proteine übersetzt!
(RNA-Gene!!)
- auch nicht alle Genregionen proteinkodierender Gene werden in Proteine übersetzt
(5' und 3'-untranslatierte Exons)
- Gene werden alternativ gespleißt.
Die ALT-mRNAs können unterschiedliche Proteine kodieren.
- RNA editing, translational readthrough, ribosome frameshift

Die drei Wege zum Gen

1. Datenbanksuchen/Alignments

extrinsisch

„es gibt bereits passende Sequenzen in den Datenbanken“

2. „*ab initio*“-Genvorhersage

intrinsisch

„Signale“ und „Inhalt“ in der DNA zeigen: hier ist ein Gen“

3. Vergleich von verwandten Genomen

(„comparative genomics“)

extrinsisch

„hier ist eine evolutionär konservierte Region“

Der direkte Beweis: Gene finden durch DB-Suchen

Ein vermutliches Gen liegt in einem Genom-DNA-Abschnitt vor, wenn...

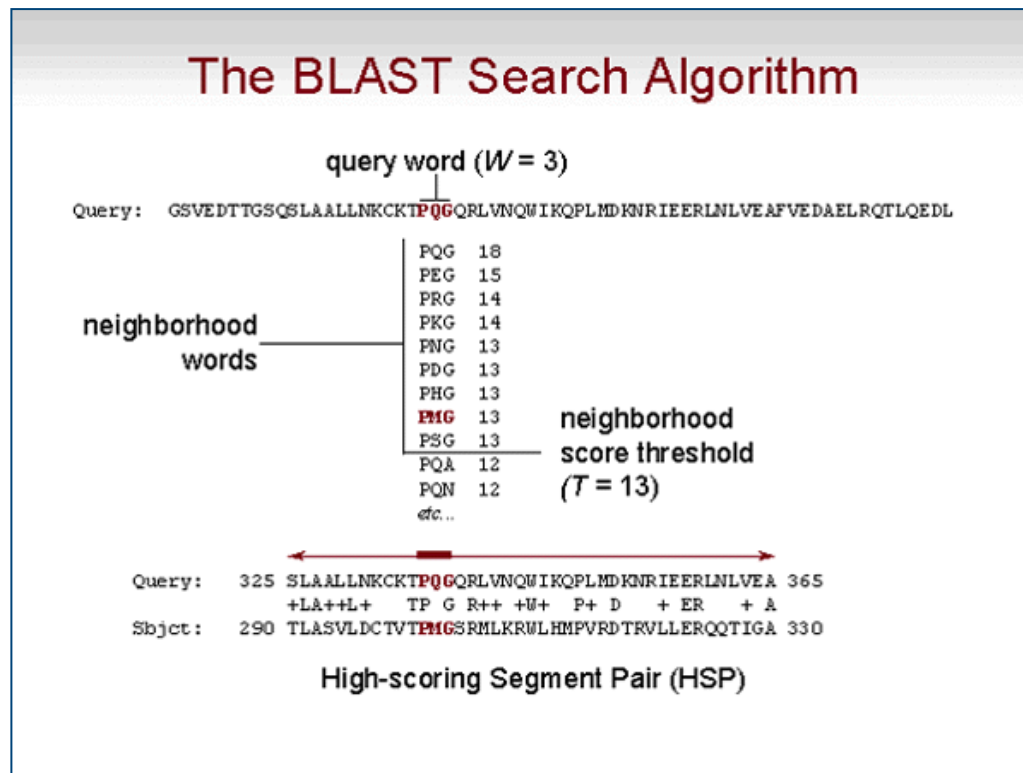
- eine **komplette cDNA** des Gens bekannt ist
- passende partielle cDNA-Sequenzen (z.B. **ESTs**) existieren
- auf Proteinebene (also nach Translation der Genom-Sequenz!) eine Ähnlichkeit zu einem bekannten **Protein** gefunden wird

Vorzugsweise wird zuerst nach Datenbankeinträgen desselben oder naher verwandter Organismen gesucht (auf DNA-Ebene), dann auf Proteinebene nach Ähnlichkeiten in entfernten Organismen (oder entfernt verwandten Proteinen).

Wir erinnern uns....

Der BLAST-Algorithmus

- liefert lokale alignments



- zunächst wird nach kurzen lokal passenden Abschnitten („words“) gesucht,

dann versucht BLAST2.0, die Bereiche neben den „matching words“ unter Einbeziehung von Lücken zu optimieren

(word size $W = 11$ bei DNA)

BLAST :

Entdecke die Möglichkeiten!

blast**n**

DNA-Sequenz ÷ DNA-DB

> nur nahe Verwandtschaft; beide Stränge verglichen

blast**p**

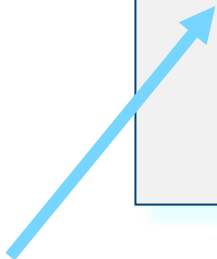
As-Sequenz ÷ Protein-DB

> entfernte Verwandtschaft

blast**x**

DNA-Seq > in 6 Leserahmen translatiert
÷ Protein-DB

> findet mögliche Proteine in einer nicht-
charakterisierten DNA-Sequenz



BLAST :

Entdecke die Möglichkeiten!

tblastn

As-Seq ÷ DNA-DB (6-frame translatiert!)

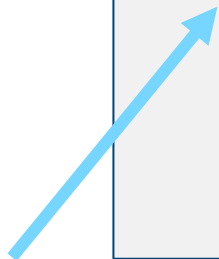
> findet nicht-annotierte Genregionen in DNA-DB-Sequenzen

tblastx

6-frame-Translation einer DNA-Seq ÷
6-frame-Translation einer DNA-DB

> Analyse z. B. von ESTs auf Proteinebene zur Detektion entfernter Verwandtschaften

(kann nicht mit nr-DB benutzt werden (zu aufwändig))



Bei der Suche nach neuen Genen verwenden wir...

BLAST^N > BLAST^X > TBLAST^X

↙
cDNAs, ESTs
aus der gleichen
Spezies
(nr, refseq, dbEST
GSS etc.)

↓ Translation
Proteine der
gleichen oder
anderer Arten
(nr, sp-trembl)

↘ Translation
Uncharakterisierte
DNA-Sequenzen
aus anderen Arten
(dbEST, GSS)

Kleiner Exkurs:

Charakterisierung eines gefundenen Gens/Proteins durch BLAST

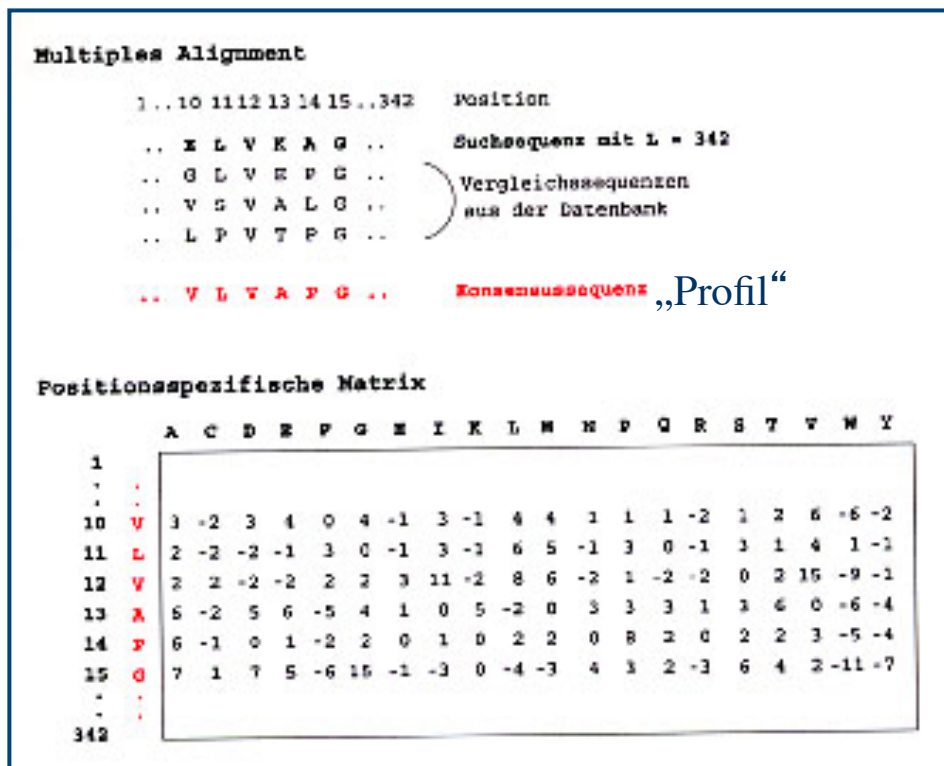
- meist zuerst über **BLASTP**
- BLASTP findet jedoch keine sehr entfernt verwandten Übereinstimmungen mehr.

Hierfür existieren weitere BLAST-Formen...

PSI-BLAST

Position-specific iterated BLAST

- speziell für die Suche sehr entfernt verwandter Proteine



1. Erste Suche = einfacher BLAST

2. Matches untereinander schreiben,
> Konsensussequenz errechnen
(„Profil“)

3. „Positions-spezifische“ Substitutions-
Matrix errechnen

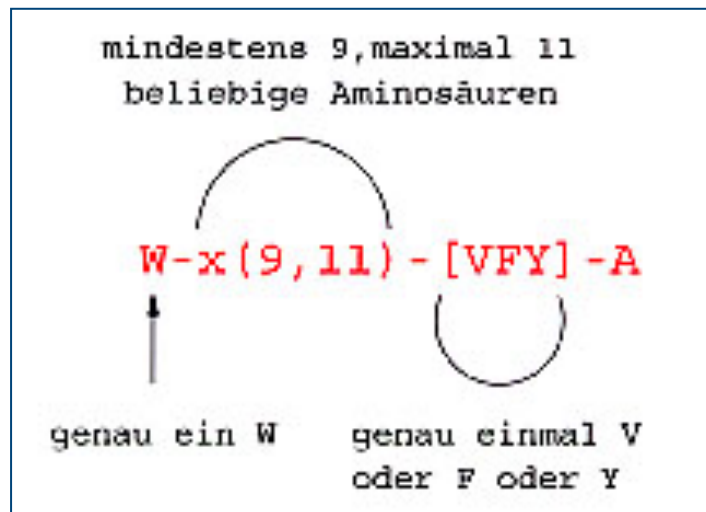
4. BLAST mit dem „Profil“ und der
PSSM mehrfach wiederholen

➤ Gezieltere Suche nach entfernt
verwandten Proteinen durch
Positionsinformation

PHI-BLAST

Pattern-hit initiated BLAST

- sucht Sequenzmuster („Signatur“), das typisch für Proteindomäne ist
- Suche über „qualitatives“ Sequenzmotif (PSI-Blast über quantitatives Motiv)



- Muster zusammen mit Suchsequenz gegen DB laufen lassen
- Treffer = Proteine mit Ähnlichkeit zur Suchsequenz und das Motiv enthaltend

Anwendung: „Suche z.B. alle Proteine mit Zinkfinger-Motiven in der Datenbank“

Zurück zum Problem des „gene findings“ in Genom-DNA per Datenbanksuche...

Am einfachsten ist es natürlich, wenn wir eine passende komplette cDNA (oder zumindest einige ESTs) zu unserem Gen in den Datenbanken finden!

Ultraschnelles Alignment (Genom-DNA/cDNA) über BLAT

„BLAST-like alignment tool“

The screenshot shows the BLAT Search Genome web interface in a Netscape browser window. The address bar shows the URL <http://genome.ucsc.edu/cgi-bin/hgBlat?command=start>. The page has a navigation bar with links: Home, Genome Browser, Blat Search, FAQ, and User Guide. The main heading is "BLAT Search Genome". Below this, there are dropdown menus for "Genome:" (set to Human), "Assembly:" (set to Human June 2002), "Query type:" (set to BLAT's guess), "Sort output:" (set to query_score), and "Output type:" (set to hyperlink). A text box for pasting a query sequence is present, with a note: "Please paste in a query sequence to see where it is located in the genome. Multiple sequences can be searched at once if starting with > and the sequence name." Below the text box is a large empty area for results. At the bottom, there is a section for uploading a file: "Rather than pasting a sequence, you can choose to upload a text file containing the sequence. Upload sequence:" followed by a "Browse..." button and a "Submit File" button.

- DNA-BLAT findet 40 Bp (> 95% id) bzw. perfekte matches von >33 Bp

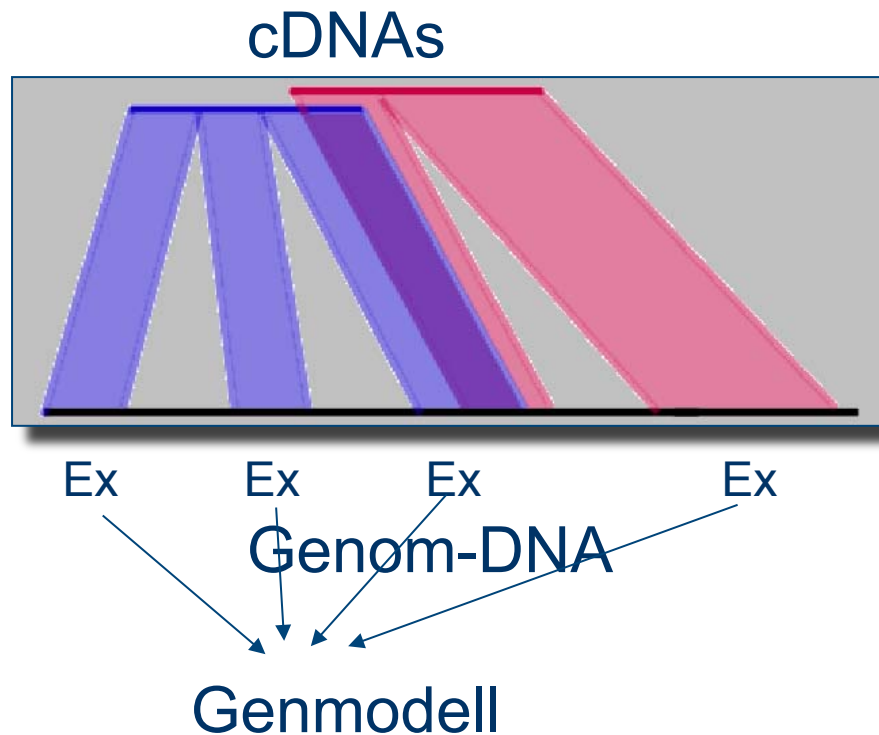
- Protein-BLAT findet 20 aa (< 80% id)

- Index (DNA) enthält alle nicht-überlappenden 11-mere des Genoms (1 Gb RAM)!!!

- Index wird gebraucht um passende Regionen im Genom schnell zu identifizieren, die dann für genaueren Vergleich „hochgeladen“ werden

<http://genome.ucsc.edu/cgi-bin/hgBlat?command=start>

„Spliced alignment“ von Genomsequenz und cDNA



1. BLAST (mRNA vs. Genom)
> Identifizierung des „Genomfensters“
2. Konsistenz der alignen Abschnitte?
(selber DNA-Strang, lineare Abfolge)
3. „low stringency“-BLAST um komplette mRNA zu alignen
4. Anpassen der Lage von Spleiß-Signalen an den Enden der Alignments

• <http://bioweb.pasteur.fr/seqanal/interfaces/est2genome.html>
<https://www.ncbi.nlm.nih.gov/sutils/splign/splign.cgi>

Die drei Wege zum Gen

1. Datenbanksuchen/Alignments

„es gibt bereits passende Sequenzen in den Datenbanken“

2. „*ab initio*“-Genvorhersage

„Signale“ in der DNA zeigen: hier ist ein Gen“

3. Vergleich von verwandten Genomen

(„comparative genomics“)

„hier ist eine evolutionär konservierte Region“

- **Genvorhersage** macht Modelle für Gene
- **Gen-Annotation** benutzt diese Modelle, fügt experimentelle Evidenzen hinzu und bewertet das Ergebnis

Edgar Allen Poe und die DNA-Linguistik

Zum Schatz von Captain Kidd... („The Gold-Bug“)

5 3 †††3 0 5)) 6 * ; 4 8 2 6) 4 †.) 4 † : 8 0 6 * ; 4 8 † 8 ¶ 6 0)) 8
 5 ; 1 † (; : † * 8 † 8 3 (8 8) 5 * † ; 4 6 (8 8 * 9 6 * ? ; 8) * † (; 4 8 5
) ; 5 * † 2 : * † (; 4 9 5 6 * 2 (5 * - - 4) 8 ¶ 8 * ; 4 0 6 9 2 8 5) ;) 6 † 8
) 4 †† ; 1 († 9 ; 4 8 0 8 1 ; 8 : 8 † 1 ; 4 8 † 8 5 ; 4) 4 8 5 † 5 2 8 8 0 6
 * 8 1 († 9 ; 4 8 ; (8 8 ; 4 († ? 3 4 ; 4 8) 4 † ; 1 6 1 ; : 1 8 8 ; † ? ;

- häufigstes engl. Wort? ;48 the

5 3 †††3 0 5)) 6 * T H E 2 6) H †.) H † : E 0 6 * T H E † E ¶ 6 0)) E
 5 T 1 † (T : † * E † E 3 (E E) 5 * † T 4 6 (E E * 9 6 * ? T E) * † (T H E 5
) T 5 * † 2 : * † (T H 9 5 6 * 2 (5 * - - H) E ¶ E * T H 0 6 9 2 E 5) T) 6 † E
) H †† T 1 († 9 T H E 0 E 1 T E : E † 1 T H E † E 5 T H) H E 5 † 5 2 E E 0 6
 * E 1 († 9 T H E T (E E T H († ? 3 H T H E) H † T 1 6 1 T : 1 E E T † ? T

Welche „Signale“ von Genen kennen wir?

- **Repeats** = meist keine Gene > also filtern/maskieren
-

- **Startkodons, Stopkodons** > ORFS („open reading frames“)
- **Spleiß-Donor/Akzeptor-Stellen** („GT-intron-AG“)
- Promoter: Bindemotive für Transkriptionsfaktoren („**Boxen**“)
Startpunkt der Transkription (+1, cap site)
CpG-Inseln
- **Polyadenylierungssignal** (AATAAA) am Ende des Transkripts

Voraussetzung für viele weitere Analysen:

Masked Sequence:

„REPEATMASKER“

Probleme:

- Probleme:**
- Repeats sind artspezifisch
 - Repeats mutieren rasch zur Unkenntlichkeit und werden nicht maskiert

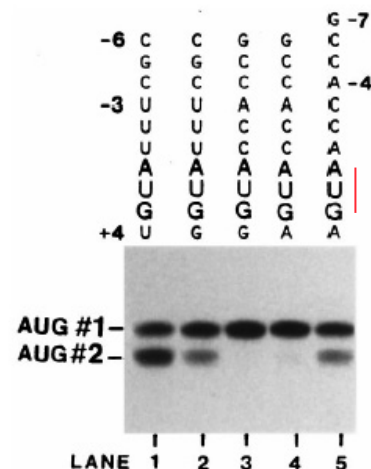
Manche „Signale“ sind recht eindeutig definiert...

Bsp. Startkodon

Kozak's rule“:

Meist gilt das erste „in frame“-ATG vom 5`-Ende der mRNA her gesehen als Startkodon

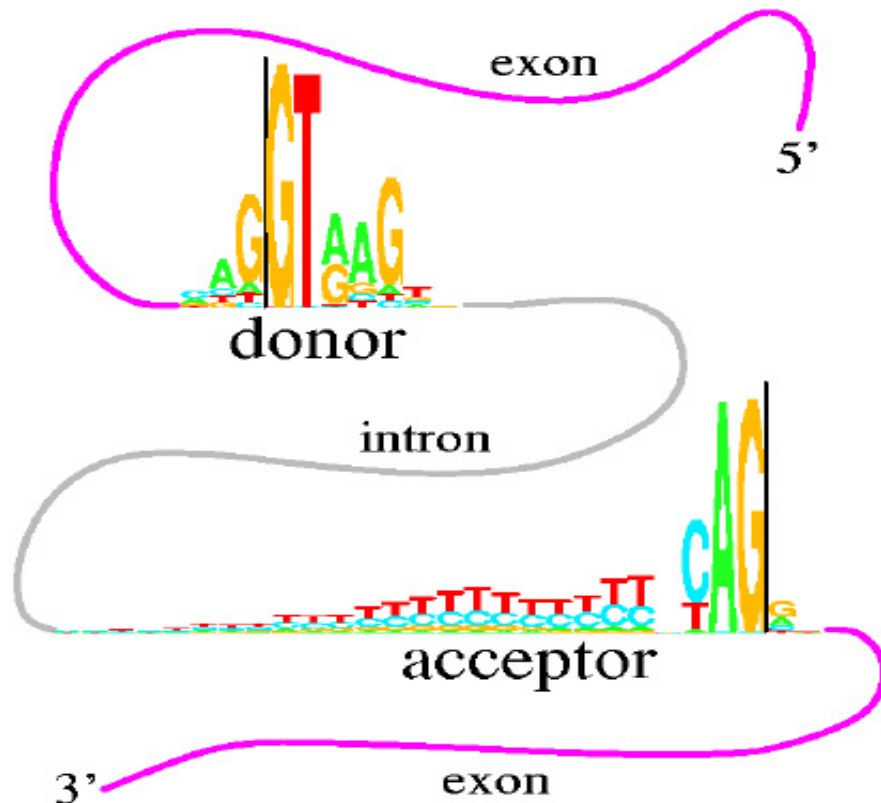
Aber es gibt Ausnahmen:



Gene 234 (1999)
187-208

Viele „Signale“ sind jedoch unscharf definiert...

Bsp. Spleißstellen



- was unterscheidet eine effiziente von einer kryptischen Spleiß-Stelle??

„Position-specific weight matrices“

erkennen Signale besser als Konsensus-Sequenzen

Problem von
Konsensus-
Sequenzen:

MELON
MANGO
HONEY
SWEET
COOKY

MONEY

Bsp. „branch site“ von Introns

• PSWM:

	-3	-2	-1	0	+1
A	1	0	39	99	4
C	76	8	15	1	45
G	2	0	42	0	6
T	21	91	4	0	38

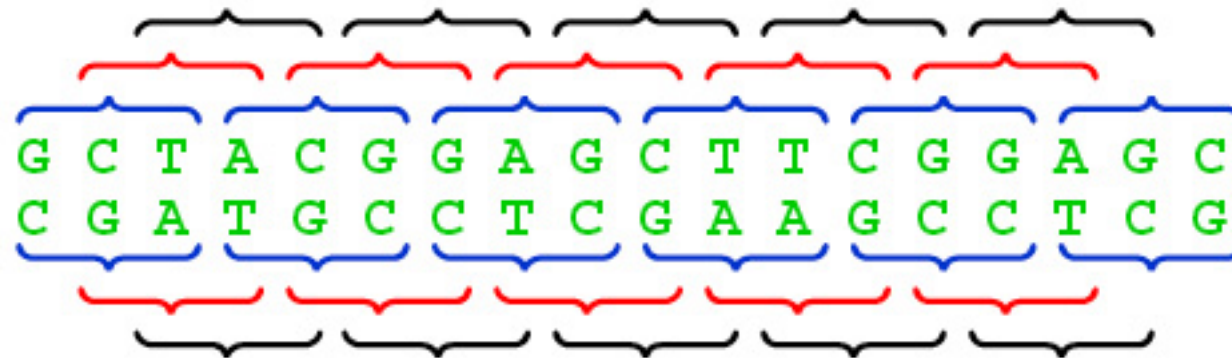
Proteinkodierende Gene haben einen „besonderen Inhalt“

- sie lassen sich als einen „offenen Leserahmen“ (ORF) lesen

Frame 3

Frame 2

Frame 1



Frame 4

Frame 6

Frame 5

Suche nach ORFs



| Start
| Stop

Offene Leserahmen (ORFs)

Der NCBI-ORFfinder

NCBI Resources How To Sign in to NCBI

ORFfinder PubMed Search

Open Reading Frame Finder

ORF finder searches for open reading frames (ORFs) in the DNA sequence you enter. The program returns the range of each ORF, along with its protein translation. Use ORF finder to search newly sequenced DNA for potential protein encoding segments, verify predicted protein using newly developed SMART BLAST or regular BLASTP.

This web version of the ORF finder is limited to the subrange of the query sequence up to 50 kb long. Stand-alone version, which doesn't have query sequence length limitation, is available for [Linux x64](#).

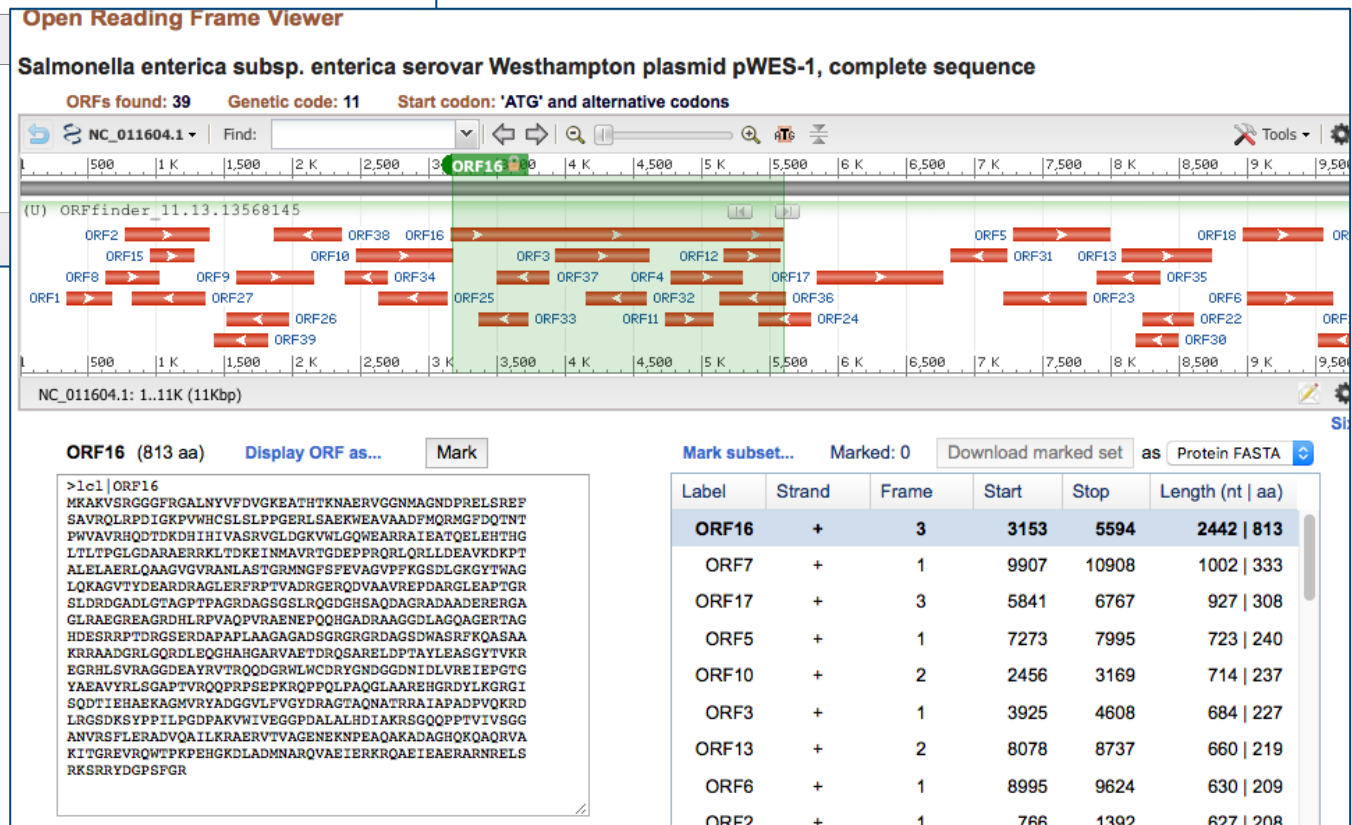
Examples (click to set values, then click Submit button) :

- NC_011604 Salmonella enterica plasmid pWES-1; genetic code: 11; 'ATG' and alternative initiation codons; minimal ORF length: 300 nt
- NM_000059; genetic code: 1; start codon: 'ATG' only; minimal ORF length: 150 nt

Enter Query Sequence

Enter accession number, gi, or nucleotide sequence in FASTA format:

From: To:



Proteinkodierende Gene haben einen „besonderen Inhalt“

- ihre Nukleotid-Komposition unterscheidet sich von nicht-kodierender DNA:
 - > bestimmte Aminosäuren/Codons sind häufiger als andere (z. B. ist Leu häufiger als Trp)
 - > es gibt unterschiedliche Kodonanzahlen für verschiedene As (Leu:6 Trp:1)
 - > für eine bestimmte As werden bestimmte Kodons häufiger gebraucht als andere („codon usage“)

Proteinkodierende Gene haben einen „besonderen Inhalt“

- Di-Nukleotid-Häufigkeiten: z.B. „AT“ in proteinkodierenden Sequenzen unterrepräsentiert

```
cctcacctctgagaaaacctctttgccaccaataccatgaagctctgcgtgactgtcctgtctctc
ctcgtgctagtagctgccttctgctctctagcactctcagcaccaagtaagtctacttttgcagct
gctatttcgagtcaggtgtaggcagagtccttttttctagtcattggctggcaaacagtgggatct
ggggatgggacaaaaggcagctaggaagattgccatgtagtctgctgctaaatgtagagtctagta
gatattcagtaacattcaagttcctattttcttaagaattagcaaccagcagaggaaaacgatggg
ctggaagtcagactggtgaattggctctgcctttaattattgttcaagcaagcccctgtccctct
ctgtgccttggtttccccatctgtcatatgaaggagtgcgatgtgttctgagactgaatccagtt
ccaatcttctagatttctttctcgttcttctctgaagatccactattcagaataagactcctgctc
atgttaggtgggaatggatacaagggaccatatgtggggttctggtagctccacagggatgctcaa
tgaagatgcaaaattagaagtcaaaataaacagctcccattgggcagtgttgatctcaccctggcct
ttcctttcagtgggctcagaccctcccaccgcctgctgcttttcttacaccgcgaggaagcttcct
cgcaactttgtggttagattactatgagaccagcagcctctgctcccagccagctgtggtgtgagta
tcaacccttggtgctgcctgggaggcaaggggtgagggctggatttttaaagggggcctgttttgggg
agggggtgatgagcgtggtgggagggcagctctcagggctgaagccttccctgacagcagtgaggtca
caggtcatgaactcacttttcaagtgtgaaggcggctgagtggcagccgagacagaagggggttc
ctggggaggaagttattcagaggacaggggaagcaggggaaggcagacaggtcccattgagatatgga
ccaattccttaaaccatgctagaaaaaacatgtgaaaagtactaccaggtggcaggggaatgggg
caatctattcactgattgcaatgccactggttcctaattctgggcaacccttggggcccacagc
taaattccagtgagtggaagttacagggagtgcttccagtgctgctcgaggaaggatcccattcca
ccagagctgccccacatggaccatggtcaggcagaggaagatgcctaccacaggcaagggaataag
ccagatgacctcaaaggtcccattgggatttctaattctgtctgctccttgttctacagattccaaacc
aaaagaggcaagcaagtctgcgctgacccagtgagtcctgggtccaggagtacgtgtatgacctg
gaactgaactgagctgctcagagacaggaagtcttc
```

Sequenzbeispiel:
M. Stanke, Greifswald

Codon usage

- Kodons werden unterschiedlich häufig verwendet
- dies kann sogar je nach Spezies unterschiedlich sein
(> artspezifische Trainingsdatensätze erforderlich!)

Bsp. E. coli

AA	codon	/1000

Gly	GGG	1.89
Gly	GGA	0.44
Gly	GGU	52.99
Gly	GGC	34.55
Glu	GAG	15.68
Glu	GAA	57.20
Asp	GAU	21.63
Asp	GAC	43.26

Eine codon-usage-Methode zur Erkennung einer proteinkodierenden Sequenz

- Staden and Mc Lachlan 1982:

Bestimme die Wahrscheinlichkeit P , dass eine Sequenz S „so aussieht“, **weil** sie ein Protein kodiert

- Die Sequenz in einem Sequenzfenster fällt immer in eine von 7 Kategorien (frame 1...frame 6, non-cod)
- Benutze Bayes-Satz, um die Wahrscheinlichkeit jeder Kategorie zu bestimmen:

$$\Pr(\text{coding}_i \mid S) = \frac{\Pr(S \mid \text{coding}_i) \Pr(\text{coding}_i)}{\Pr(S)}$$

Wahrscheinlichkeit, dass Sequenz ein Protein in frame i kodiert

Eine codon-usage-Methode zur Erkennung einer proteinkodierenden Sequenz

- Kalkuliere:
$$\Pr(S \mid \text{coding}_i) \approx \prod_{j=1}^n \Pr(S_i(j) \mid \text{coding}_i)$$

Wahrscheinlichkeit des Kodons j in frame i, gegeben die Sequenz ist kodierend



Diagram showing a DNA sequence with codons grouped by blue brackets:

```

G C T A C G G A G C T T C G G A G C
C G A T G C C T C G A A G C C T C G
  
```

$\Pr(S \mid \text{coding}_0) \approx$
 $\Pr(\text{GCT} \mid \text{coding}_0) \times \Pr(\text{ACG} \mid \text{coding}_0) \dots$

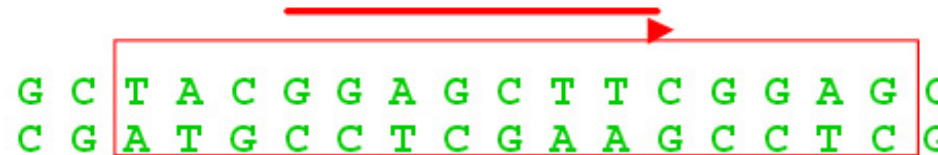


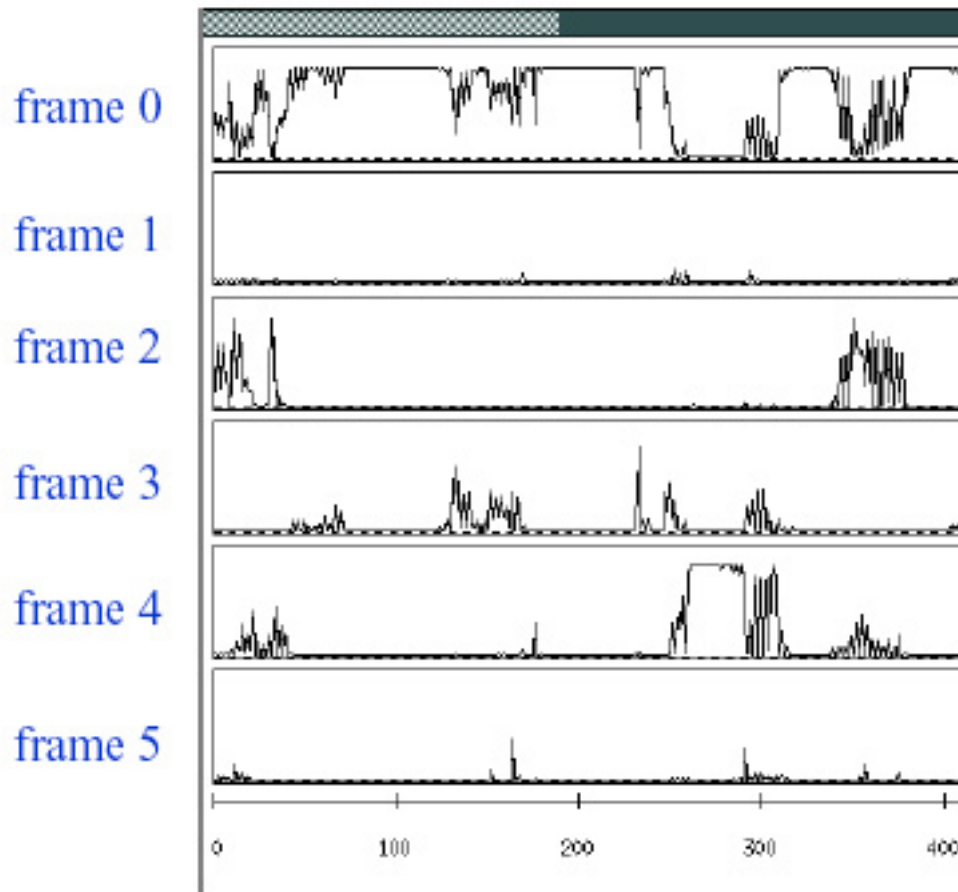
Diagram showing a DNA sequence with a sliding window highlighted in red:

```

G C T A C G G A G C T T C G G A G C
C G A T G C C T C G A A G C C T C G
  
```

Sliding
window

Eine codon-usage-Methode zur Erkennung einer proteinkodierenden Sequenz



- aber: recht grobe Analyse der Kodierungskapazität

- Exakte Begrenzung der ORFs nicht festgelegt

Bisherige Suchen sind nicht ausreichend, um komplizierte Genmodelle exakt vorherzusagen!

Moderne *integrierte* Genvorhersage-Programme verbinden Suche nach Signalen und Inhalt mit statistischem Vorab-Wissen über Gene...

...Hidden Markov Models (HMM)

Markov WER??

- Andrei Andreyevich Markov (1856-1922)
- Markov-Kette:



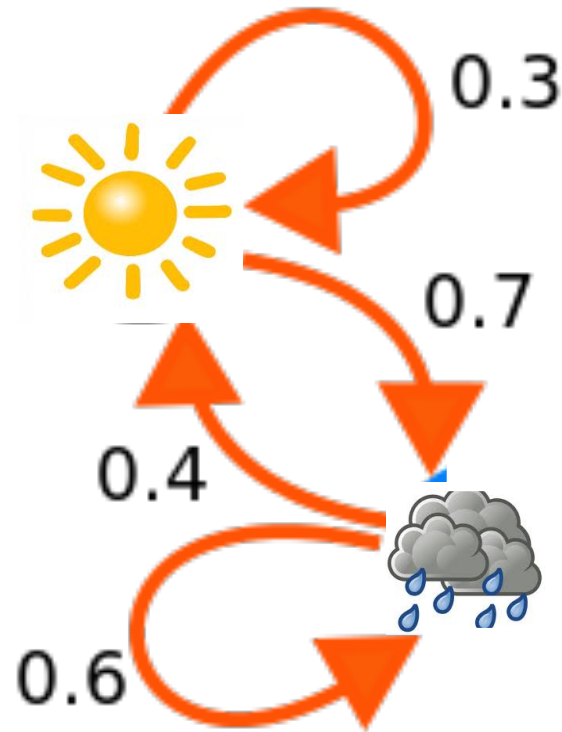
Eine **Markovkette** ist ein stochastischer Prozess, der nacheinander eine Reihe von Zuständen mit einer gewissen Wahrscheinlichkeit durchläuft. Dabei hängt die Wahrscheinlichkeit für den jeweils nächsten Zustand nur vom aktuellen Zustand ab:

$$P(t_{i+1} | t_i, t_{i-1}, \dots, t_j) = P(t_{i+1} | t_i)$$



Pfeile geben
Übergangswahrscheinlichkeiten an

Markov-Kette „Wettervorhersage“



Hidden Markov Models

- verwende **statistische Informationen**, um Abfolgen (z. B. Sequenzen) zu klassifizieren

- Analogie:

„Automatische Erkennung der Sprache eines Textes“

In einem typischen deutschen Text macht der Buchstabe ,e‘ ca. 16,55% aller Buchstaben aus, in einem schwedischen nur ca. 9.77%.

⇒ zähle die e‘s im Text, um zu berechnen mit welcher Wahrscheinlichkeit es sich um einen deutschen Text handelt

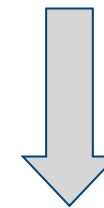
Hidden Markov Models

Was ist denn da „hidden“??

- wir *sehen* nur die „e ‘s“
- dahinter *versteckt* sich die Information:

„dies ist ein deutscher Text“

„emission“



„state“

„The Occasionally Dishonest Casino“

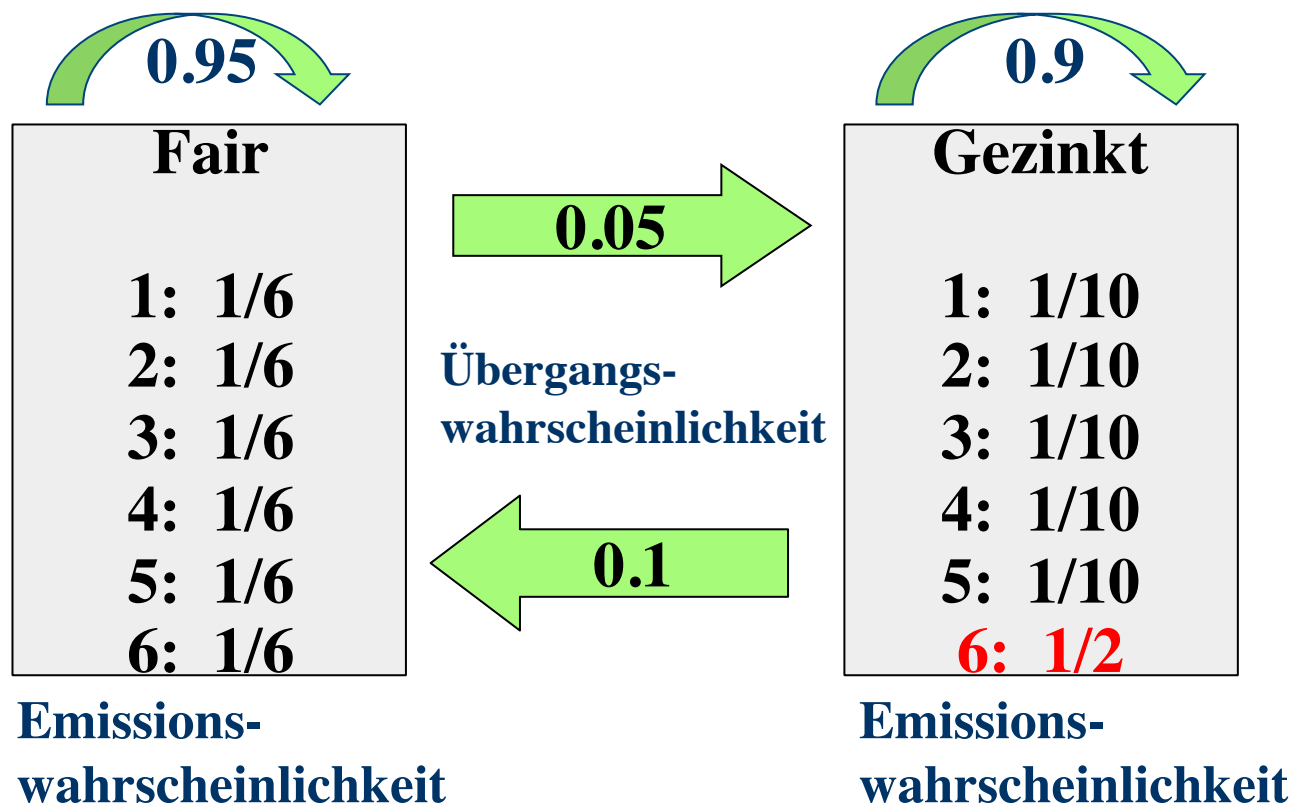
- Das Casino besitzt gezinkte und faire Würfel
- Der gezinkte Würfel hat 50% Chance eine 6 zu bekommen und gleiche Chancen (10%) eine der anderen Zahlen zu bekommen.
- Das Casino wechselt von Zeit zu Zeit die Würfel, um die Spieler zu betrügen.

Zahlen: 315116246446644245311321363165662656666651166453132651245

Wann ist der Würfel gewechselt worden?

The Occasionally Dishonest Casino

Annahme: wir kennen die Wahrscheinlichkeit des Würfel-Wechsels und die Wahrscheinlichkeiten der von diesen Würfeln gewürfelten Zahlen



The Occasionally Dishonest Casino

- Der Viterbi-Algorithmus (Viterbi 1967) berechnet die Würfelsituation auf Basis der gewürfelten Zahlen.
- Das HMM ist die Grundlage des Algorithmus.

Zahlen: 315116246446644245311321363165662656666651166453132651245
Würfel: FFFFFFFFFFFFFFFFFFFFFFFFFFGGGGGGGGGGGGGGGGGFFFFFFFFFFF
Viterbi: FFFFFFFFFFFFFFFFFFFFFFFFFFGGGGGGGGGGGGGGGGGFFFFFFFFFFF

Hidden Markov Models

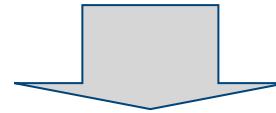
- Anwendungsgebiete in der Bioinformatik:
 - > **Vorhersage der Genstruktur (Exons/Introns)**
 - > **Vorhersage von Promoterbereichen**
 - > Erstellung von Modellen für Proteinfamilien
zum Suchen nach entfernt verwandten Proteinen
in DB („profile HMMs“)

Von der Textsuche zum HMM

1 ACA---ATG
2 TCAACTATC
3 ACAC--AGC
4 AGA---ATC
5 ACCG--ATC

Bsp.: Fünf Sequenzen, die
ein Signal definieren

→ Einfache Textsuche würde erfolgen nach:
(AT) (GC) (AC) (ACGT)* A(TG) (GC)



Dies kann bei Suche nicht unterscheiden zwischen...

...einer plausiblen Sequenz (z.B. der Konsensus-S.)

ACAC--ATC

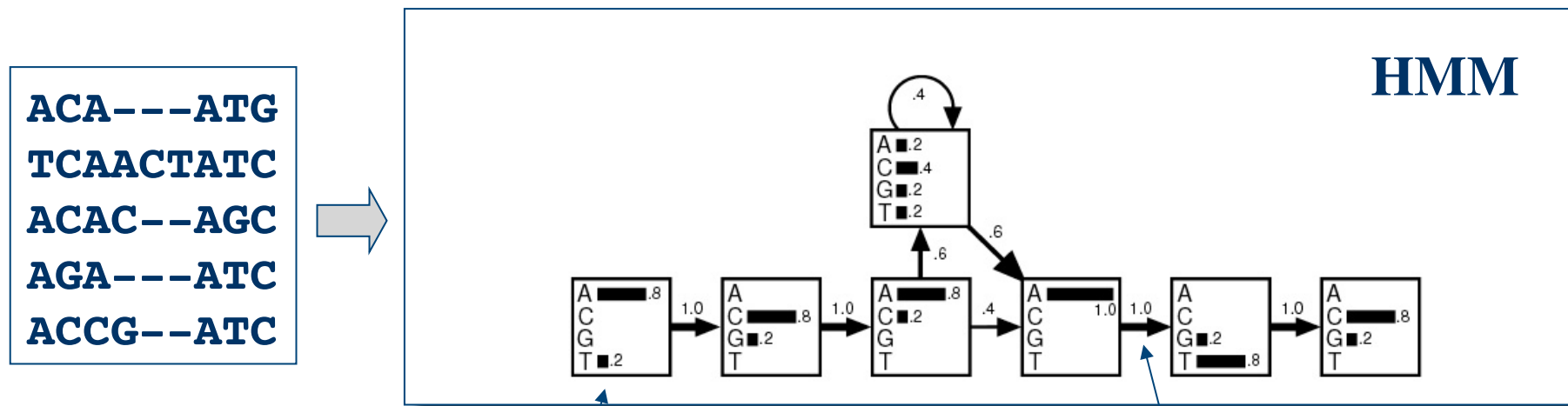
...und einer höchst unwahrscheinlichen Sequenz

TGCT--AGG

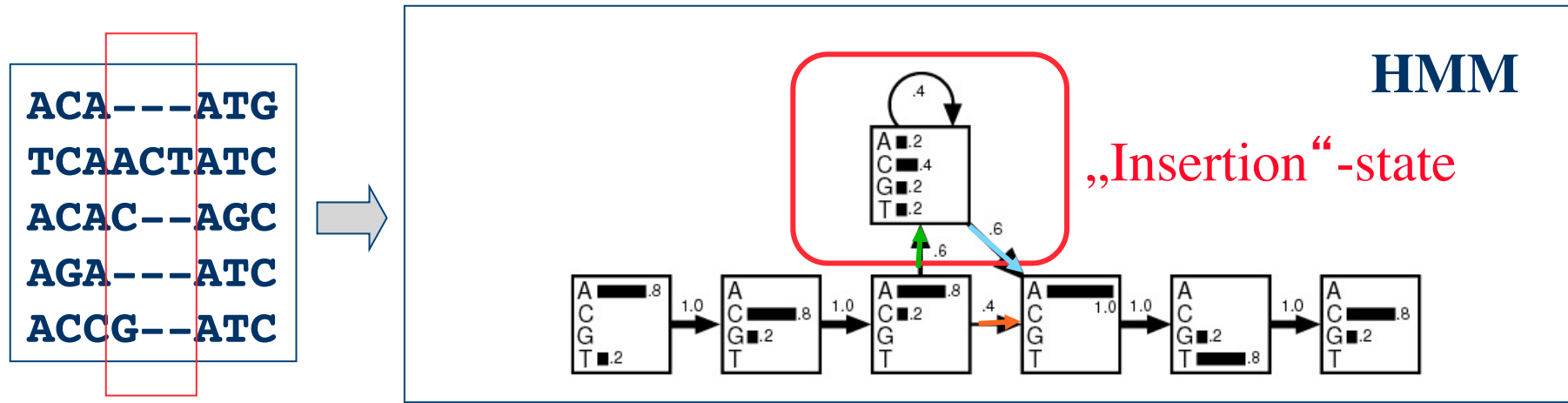
HMM

Also besser:

Bewerten, ob Sequenzabfolge „gut“ in das Alignment passt...

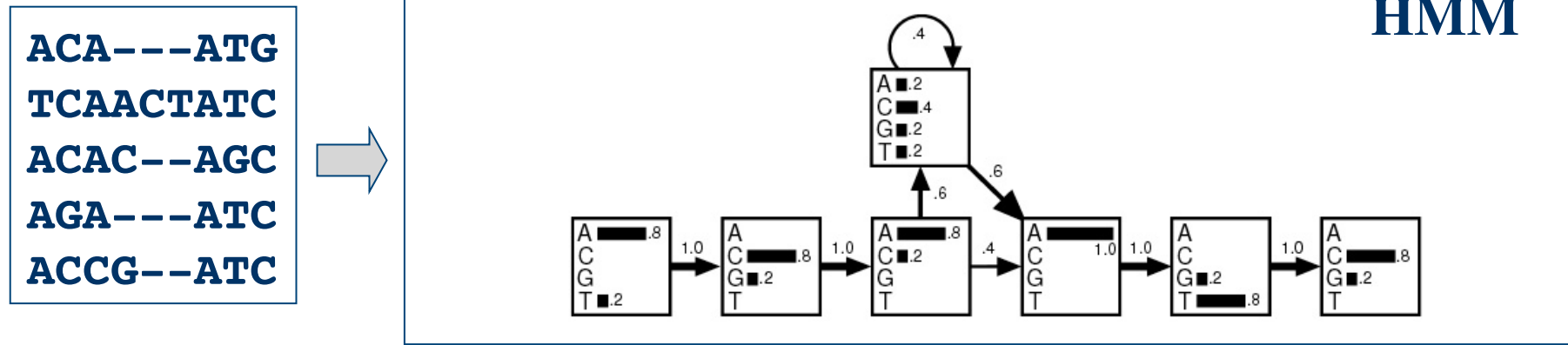


HMM



- nach Pos. 3 haben 3/5 Sequenzen eine Insertion (=0.6)
- diese 3/5 (=0.6) Sequenzen gehen von Insertionstate zu Pos. 4
- 2/5 (=0,4) Sequenzen gehen von Pos. 3 zu Pos.4
- nach der ersten Insertion (Sequ 2,3,5) kehrt Sequ 2 zweimal (=0.4) in den Insertion-State zurück

HMM



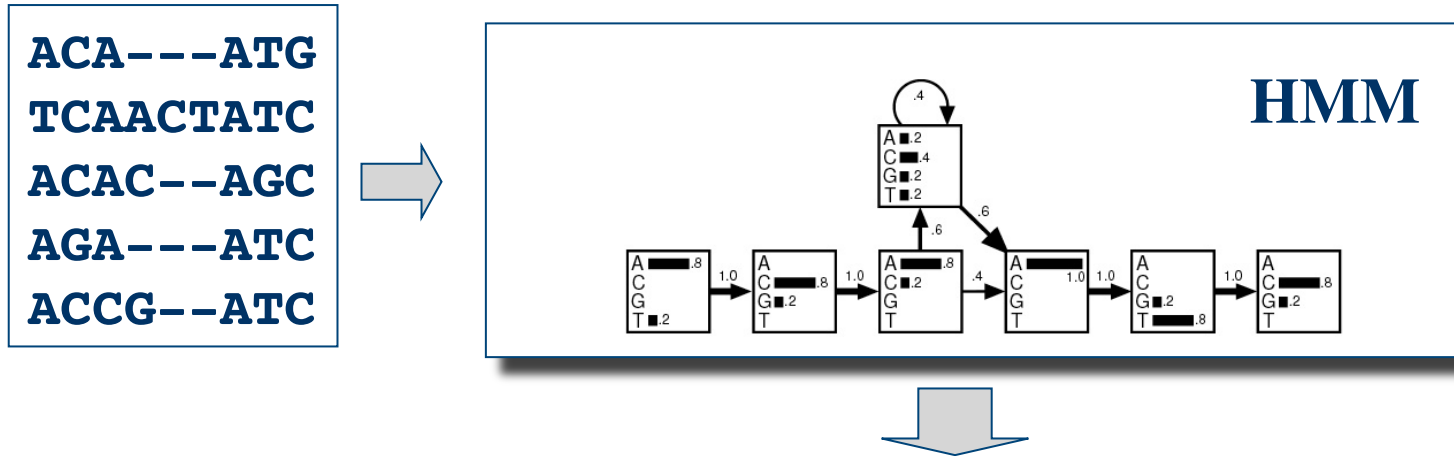
Jetzt bewerten wir damit als Beispiel die Konsensussequenz:

ACAC--ATC ➔

$$\begin{aligned}
 P(\text{ACACATC}) &= 0.8 \times 1 \times 0.8 \times 1 \times 0.8 \times 0.6 \times \\
 &\quad 0.4 \times 0.6 \times 1 \times 1 \times 0.8 \times 1 \times 0.8 \\
 &\simeq 4.7 \times 10^{-2}.
 \end{aligned}$$

Die „unplausible“ Sequenz TGCT--AGG ergibt: $0,0023 \times 10^{-2}$

HMM



Bewertung aller Sequenzen des Alignments:

	Sequence	Probability $\times 100$	Log odds
Consensus	A C A C - - A T C	4.7	6.7
Original sequences	A C A - - - A T G	3.3	4.9
	T C A A C T A T C	0.0075	3.0
	A C A C - - A G C	1.2	5.3
	A G A - - - A T C	3.3	4.9
	A C C G - - A T C	0.59	4.6
Exceptional	T G C T - - A G G	0.0023	-0.97

HMM

	Sequence	Probability $\times 100$	Log odds
Consensus	A C A C - - A T C	4.7	6.7
Original sequences	A C A - - - A T G	3.3	4.9
	T C A A C T A T C	0.0075	3.0
	A C A C - - A G C	1.2	5.3
	A G A - - - A T C	3.3	4.9
	A C C G - - A T C	0.59	4.6
Exceptional	T G C T - - A G G	0.0023	-0.97

Da der P-Wert stark von der Länge der Sequenz abhängt, erfolgt eine Normalisierung der P-Werte auf die Sequenzlänge (und bessere Skalierung durch Umrechnung in „log odds“-Werte)

$$\text{log-odds for sequence } S = \log \frac{P(S)}{0.25^L} = \log P(S) - L \log 0.25.$$

Wahrscheinlichkeit einer DNA-Sequenz der Länge L

Diskriminierung passender Sequenzen gegenüber schlechten ist nach Normalisierung und Umrechnung sehr viel besser: vergleiche **Sequenz 2** im Alignment (s.o.)

HMM

Die Umrechnung von P-Werten auf Log odds-Werte erfolgt primär aus mathematischen Gründen :

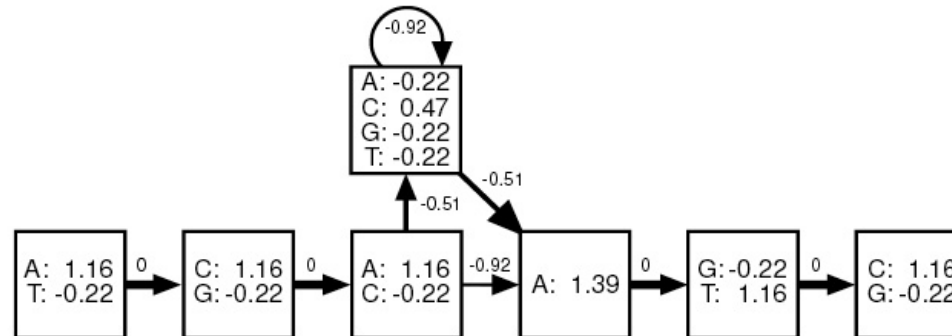


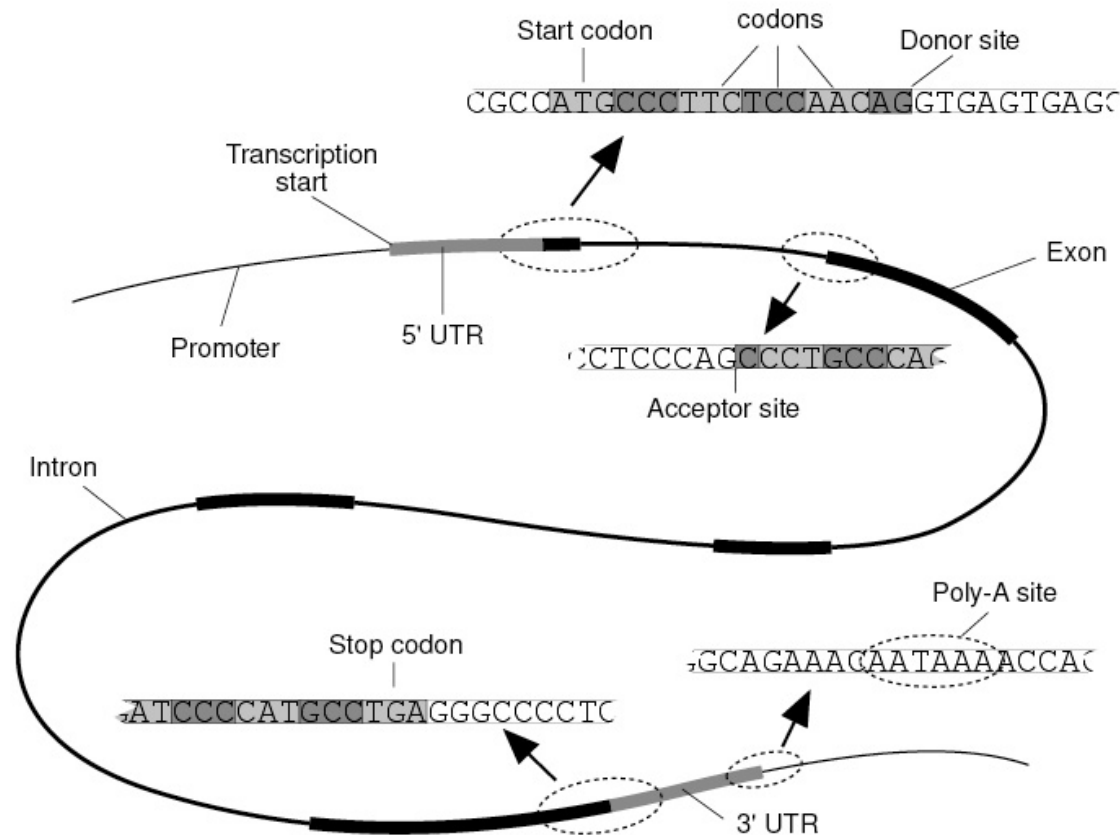
Figure 4.2: The probabilities of the model in Fig. 4.1 have been turned into log-odds by taking the logarithm of each nucleotide probability and subtracting $\log(0.25)$. The transition probabilities have been converted to simple logs.

Nun ist eine
einfache
Summation möglich:

$$\begin{aligned} \log\text{-odds}(\text{ACACATC}) &= 1.16 + 0 + 1.16 + 0 + 1.16 - 0.51 + \\ &\quad 0.47 - 0.51 + 1.39 + 0 + 1.16 + 0 + 1.16 \\ &= 6.64. \end{aligned}$$

HMMs zur Genvorhersage

Finde die „Grammatik“ des Genoms...



HMMs zur Genvorhersage

Beispiel: ein HMM-Sensor für die 3 'Spleißstelle

Intron															Exon			
C	T	C	C	C	T	G	T	G	T	C	C	A	C	A	G	G	C	T
T	A	T	T	G	T	T	T	T	C	T	T	A	C	A	G	G	G	C
G	T	T	C	C	T	T	T	G	T	T	T	C	T	A	G	C	A	C
T	G	C	C	T	C	T	C	T	T	T	T	C	A	A	G	C	A	C
T	T	C	C	T	A	T	A	T	G	T	T	G	A	C	A	G	G	T
T	T	C	C	T	G	T	T	C	C	G	A	T	G	C	A	G	G	C
T	T	G	G	G	T	T	T	C	C	T	T	T	G	C	A	G	A	C
C	A	C	T	T	T	G	C	T	C	C	C	A	C	A	G	C	G	T
C	C	C	A	T	G	T	G	A	C	C	C	T	G	C	A	G	C	T
T	A	T	T	T	A	T	T	T	A	A	C	C	C	A	G	G	A	G
A	T	G	T	G	C	A	T	C	C	C	C	C	C	A	G	A	A	T
T	T	T	T	C	C	T	T	T	T	C	T	A	C	A	G	A	C	C
T	T	C	C	A	T	G	T	C	C	T	G	A	C	A	G	C	C	C
T	T	C	C	A	T	G	T	C	C	T	G	A	C	A	G	G	T	G
A	C	G	A	C	A	T	T	T	T	C	C	A	C	A	G	G	A	G
G	T	G	C	C	T	C	T	C	C	C	T	C	C	A	G	A	T	T

Bisheriges HMM war zu einfach, da in der DNA starke Präferenzen für Di-Nukleotide bekannt sind!

d.h., ein Nt in der Folge hängt stark von seinem Nachbarn ab



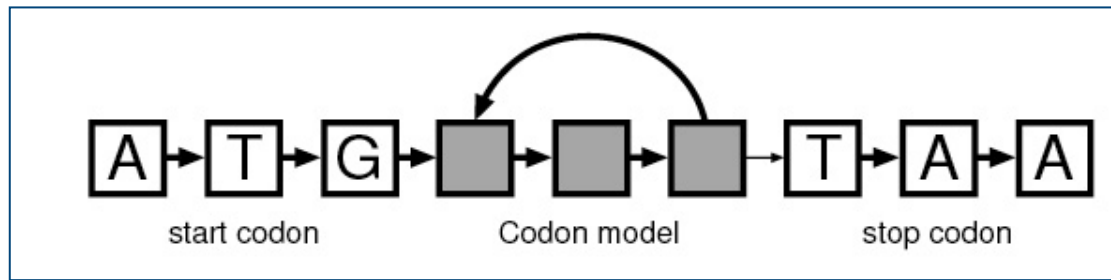
Wir müssen konditionale Wahrscheinlichkeiten für das HMM kalkulieren:

$$P(ACTGTC...) = p_1(A) \times p_2(C|A) \times p_3(T|C) \times p_4(G|T) \times p_5(T|G) \times p_6(C|T) \times \dots$$

Benutzung konditionaler Wahrscheinlichkeiten > HMM „erster Ordnung“

HMMs zur Genvorhersage

Beispiel: ein HMM-Modell für ein ungespleißtes
proteinkodierendes Gen



HMM zweiter Ordnung:
Kodonpos. 3 ist
abhängig von Pos. 1 und 2

- zähle die Anzahl jedes der 64 Kodons in einem Datensatz von bekannt-kodierenden Genregionen aus und berechne deren Wahrscheinlichkeiten

z.B. für Kodons CAA, CAC, CAG und CAT

$$p(A|CA) = c(CAA)/[c(CAA) + c(CAC) + c(CAG) + c(CAT)]$$

$$p(C|CA) = c(CAC)/[c(CAA) + c(CAC) + c(CAG) + c(CAT)]$$

$$p(G|CA) = c(CAG)/[c(CAA) + c(CAC) + c(CAG) + c(CAT)]$$

$$p(T|CA) = c(CAT)/[c(CAA) + c(CAC) + c(CAG) + c(CAT)]$$

C =
Codonanzahl

HMMs zur Genvorhersage

Kombination der einzelnen Sensoren:

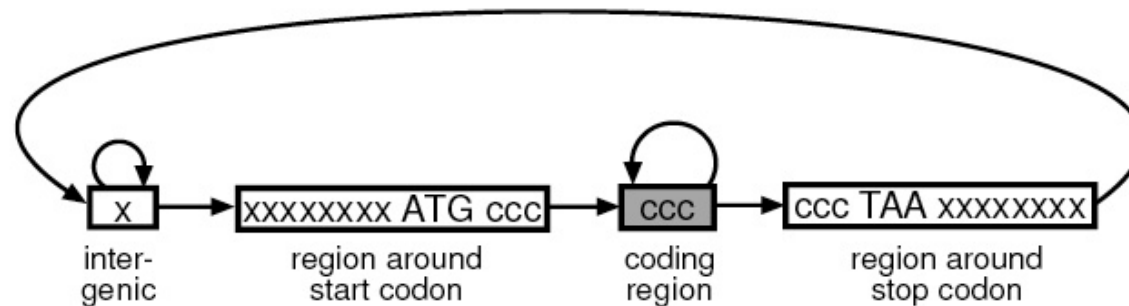
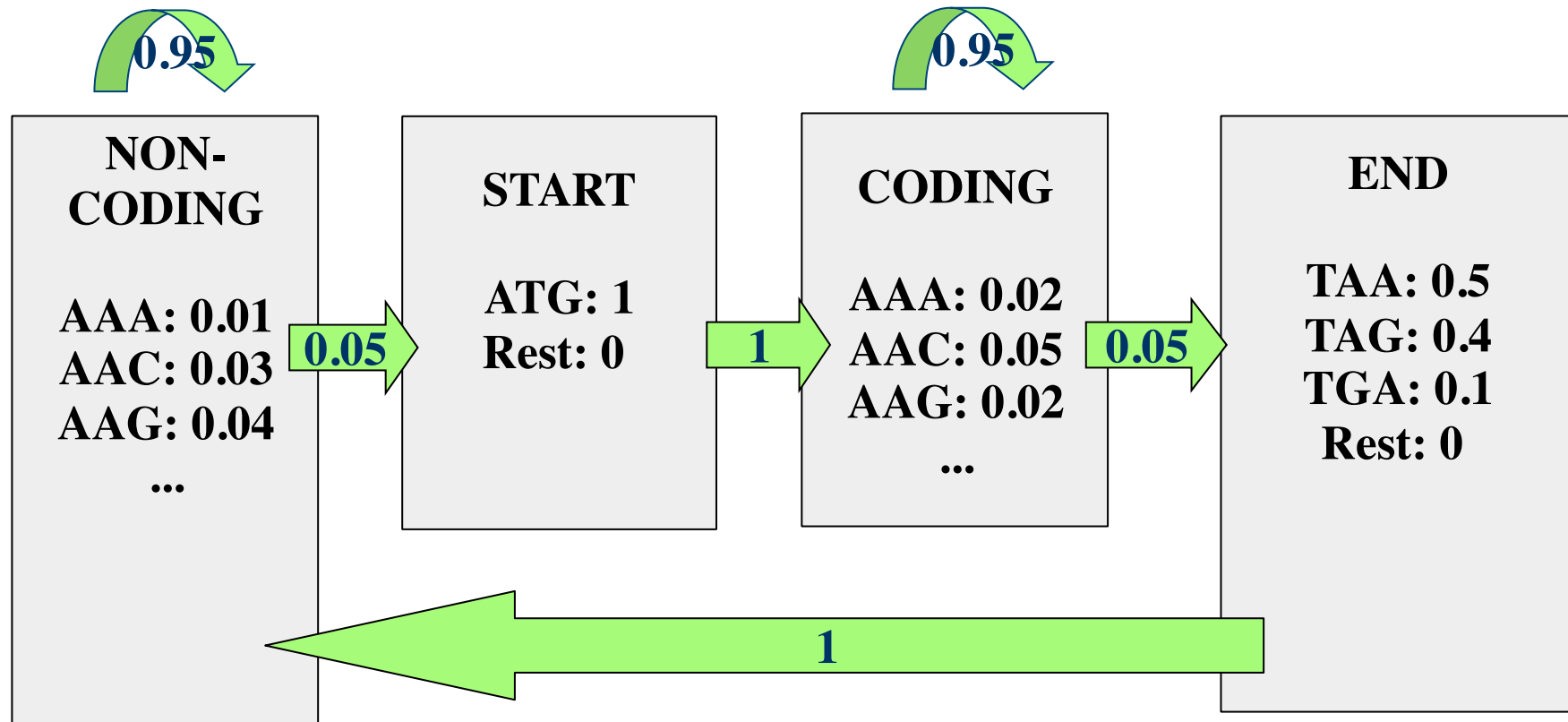


Figure 4.11: A hidden Markov model for unspliced genes. In this drawing an 'x' means a state for non-coding DNA, and a 'c' a state for coding DNA. Only one of the three possible stop codons are shown in the model of the region around the stop codon.

HMM für prokaryotisches Gen

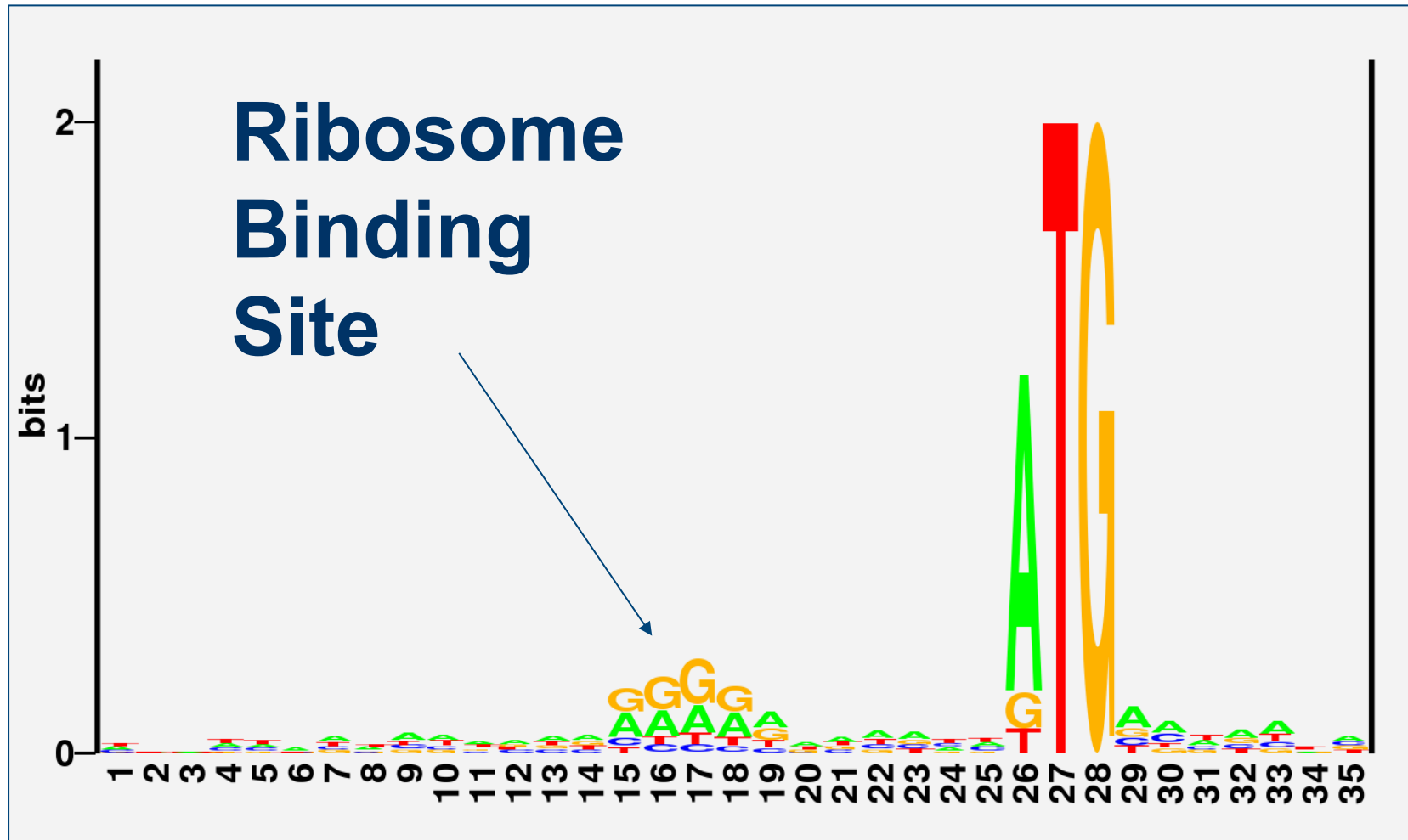


HMM für prokaryotisches Gen

Ein noch relativ übersichtliches Problem....

- **Kleine Genome (0.5-10 Mbp)**
- **Hohe Gendichte (>90%)**
- **1 Gen = 1 ORF**
- **Gene hintereinander in Operons**
- **Ribosome Binding Sites (RBS, Shine-Dalgarno) and TK-Terminatoren.**

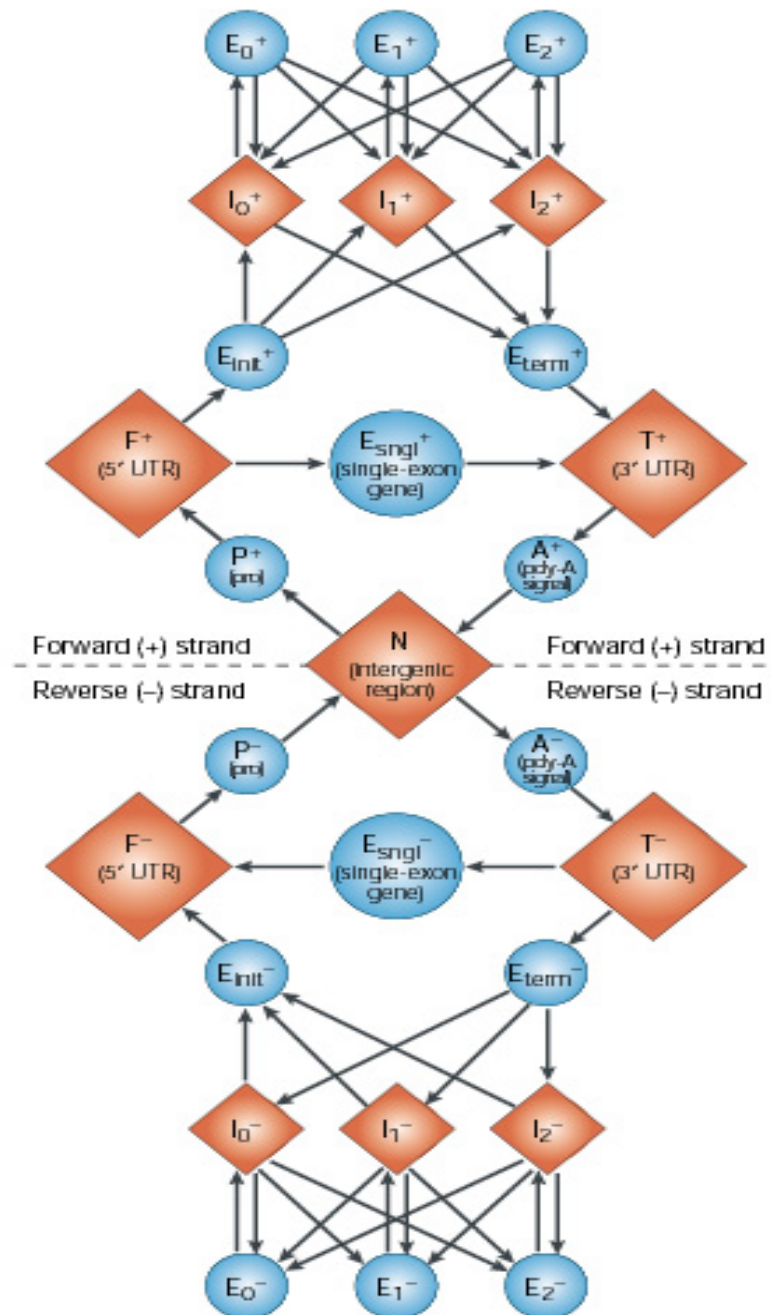
Statistisches Vorwissen für HMM



HMM für prokaryotisches Gen: „GeneMark“

- Betrachtet sliding window (96 nt) in DNA-Sequenz
- Berechnet Kodierungspotential für alle 6 Leserahmen (**Inhalt**)
- Lokalisiert RBS (**Signal**), um Startkodon genauer zu identifizieren.

Genvorhersage in Eukaryoten



GENSCAN
(Karlin and Burge 1997)

Genvorhersage mit GENSCAN

- eine Beispiel

genome.mit.edu/GENSCAN.html

- max 200 kb mit mehreren Genen möglich

Predicted genes/exons:

Gn.Ex	Type	S	.Begin	...End	.Len	Fr	Ph	I/Ac	Do/T	CodRg	P....	Tscr..
1.01	Init	+	1579	1663	85	0	1	114	54	131	0.741	13.24
1.02	Intr	+	2540	2635	96	0	0	1	100	134	0.698	6.38
1.03	Intr	+	3455	3588	134	0	2	101	81	136	0.999	15.07
1.04	Intr	+	4820	5042	223	1	1	85	56	432	0.998	37.93
1.05	Intr	+	5153	5350	198	0	0	73	81	371	0.999	34.74
1.06	Intr	+	5688	5889	202	1	1	53	69	378	0.979	31.27
1.07	Intr	+	6318	6426	109	0	1	62	80	20	0.843	-0.61
1.08	Intr	+	6576	6634	59	2	2	105	77	51	0.888	3.87
1.09	Term	+	6723	6792	70	0	1	63	54	98	0.785	1.61
1.10	PlyA	+	6853	6858	6							1.05

Predicted peptide sequence(s):

>gi|GENSCAN_predicted_peptide_1|391_aa

MAMQKIFAREILDSRGNPTVEVDLHTAKGRFRAAVPSGASTGIYEALERDGDGKGRYLKG
AKFGANAILGVSLAVCKAGAAEKGVPLYRHIADLAGNPDILPVPFNVINGGSHAGNKL
AMQEFMILPVGASSFKAMRIGAEVYHHLKGVIAKAKYGKDATNVGDEGGFAPNILENNEA
LELLKTAIQAAGYPDKVVIGMDVAASEFYRNGKYDLDFKSPDDPARHITGEKLGELYKSF
IKNYPVVSIEDPFDQDDWATWTSFLSGVNIQIVGDDLTVTNPKRIAQAVEKKACNCLLLK
VNQIGSVTESIQACKLAQSNWGVMSVSHRSGETEDTFIADLVVGLCTGQIKTGAPCRSER
LAKYNQLMRIEEALGDKAIFAGRKFRNPKAK

Column Description

Gn.Ex	Description
Type	Init = Initial exon Intr = Internal exon Term = Terminal exon Sngl = Single-exon gene Prom = Promoter PlyA = poly-A signal
S	DNA strand (+ = input strand; - = opposite strand)
Begin	beginning of exon or signal (numbered on input strand)
End	end point of exon or signal (numbered on input strand)
Len	length of exon or signal (bp)
Fr	reading frame (a codon ending at x is in frame f = x modulo 3)
Ph	net phase of exon (exon length modulo 3)
I/Ac	initiation signal or acceptor splice site score (x 10)
Do/T	donor splice site or termination signal score (x 10)
CodRg	coding region score (x 10)
P	probability of exon (sum over all parses containing exon)
Tscr	exon score (depends on length, B/Ac, Do/T and CodRg scores)

Predicted Exons	Annotated Exons
1579 1663	1579 1663
2540 2635	2540 2635
	2796 2854
	3016 3085
3455 3588	3455 3588
4820 5042	4820 5042
5153 5350	5153 5350
5688 5889	5688 5889
6318 6426	6318 6426
6576 6634	6576 6634
6723 6792	6723 6792

Bewertung der Genvorhersage

Sensitivität

wieviele vorhandenen Exons
werden richtig erkannt?
(„false negatives?“)

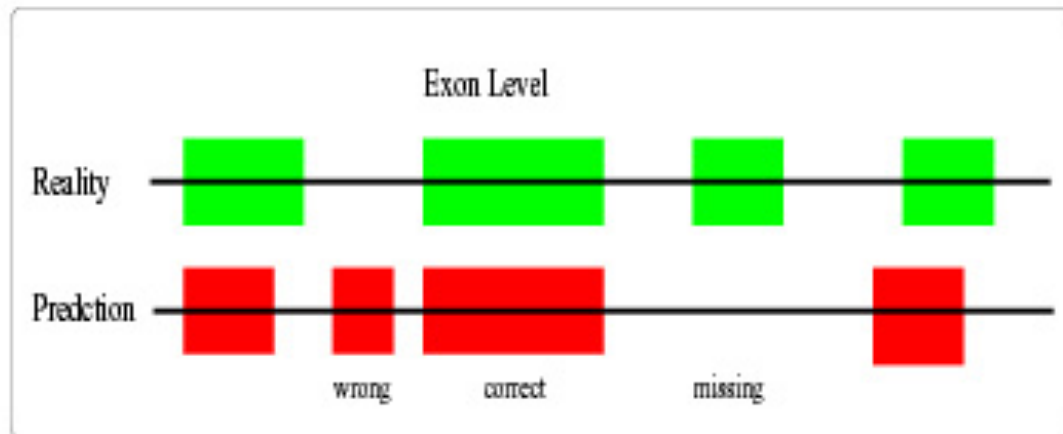
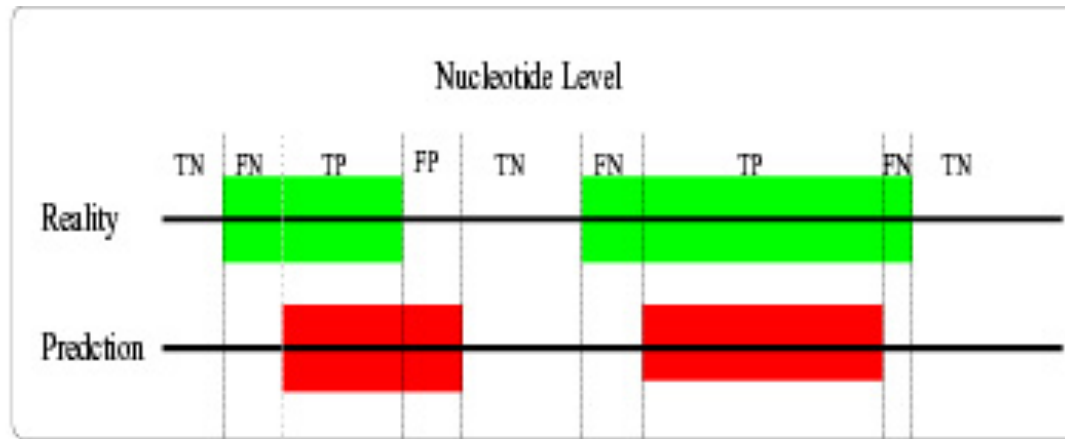
Spezifität

wie gut wird gegenüber falschen
Exons diskriminiert?
(„false positives?“)

Die meisten Genvorhersageprogramme erreichen mittlerweile eine **Sensitivität von ca. 70 %**.

Ein großes Problem bleibt auch die **oft geringe Spezifität!**

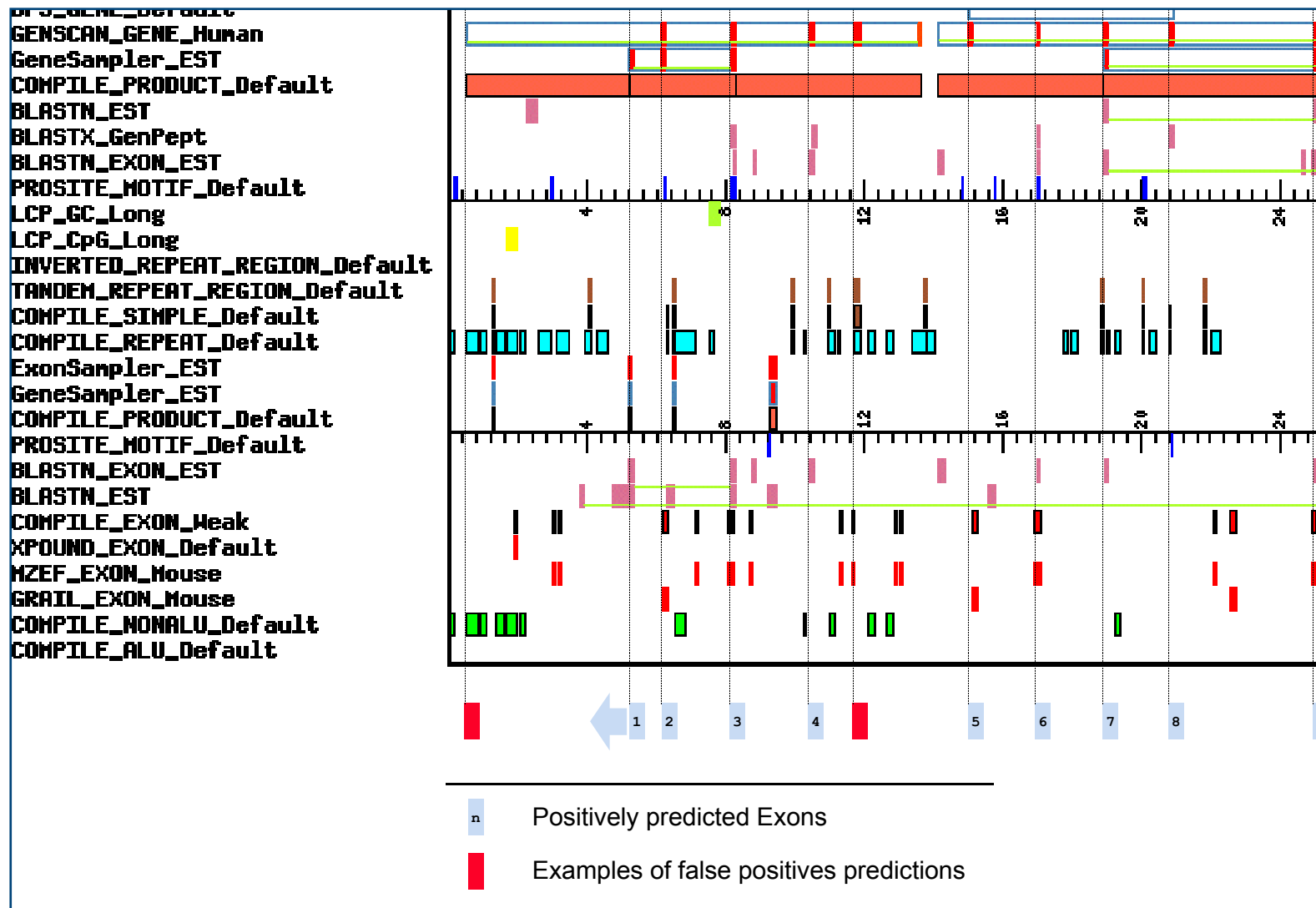
Bewertung der Genvorhersage



TN true negatives
FN false negatives
TP true positives
FP false positives

Eine mögliche Lösung:

mehrere Programme verwenden!



Genvorhersage- Programme für eukaryotische Gene

GRAIL	nur HMM: berechnet sieben Parameter kodierender Regionen
GRAIL-EX	HMM plus BLASTN
GenelD	HMM: erst Signale berechnet, dann Kodierungspotenzial dazwischen gesucht
Genscan	HMM mit „higher order parameters“: Gendichte, typische Exon-Grösse und Anzahl, GC-Verteilung
GenomeScan	= Genscan plus BLASTX-Match zu bekannten Proteinen (Nachteil: geringere scores für Exons ohne Proteinmatch)
FGENESH	HMM , linear discriminant analysis
Genewise	HMM und parallel Proteinhomologievergleich (gut, aufwändig)
Twinscan	Genscan plus Genomvergleich
SGP-2	GenelD plus Genomvergleich

Literatur & Links

http://www.wikiwand.com/en/List_of_gene_prediction_software

Literatur:

Guigo et al., Genome Biol. 2006; 7:S2.1-31

Brent & Guigo, Curr Opin Struct Biol. 2004; 14:264-72

Yandell & Ence, Nature Reviews Genetics 2012; 13: 329

Forschungsfelder in der *ab initio*-Genvorhersage

- sehr kleine Exons
- ALT-splicing
- überlappende Gene
- Promoter-Erkennung
- 5 ' und 3 ' UTRs
- RNA-kodierende Gene
- „confidence values“ angeben

Die drei Wege zum Gen

1. Datenbanksuchen/Alignments

„es gibt bereits passende Sequenzen in den Datenbanken“

2. „*ab initio*“-Genvorhersage

„Signale“ in der DNA zeigen: hier ist ein Gen“

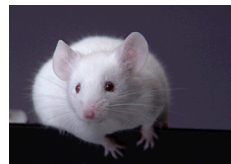
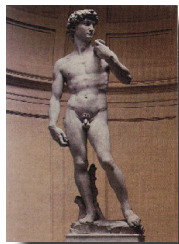
3. Vergleich von verwandten Genomen

(„comparative genomics“)

„hier ist eine evolutionär konservierte Region“

Vergleichende Genomanalyse:

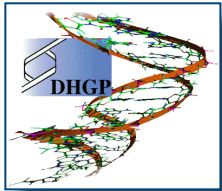
die Lösung des Genvorhersageproblems



80-100 Mio. Jahre



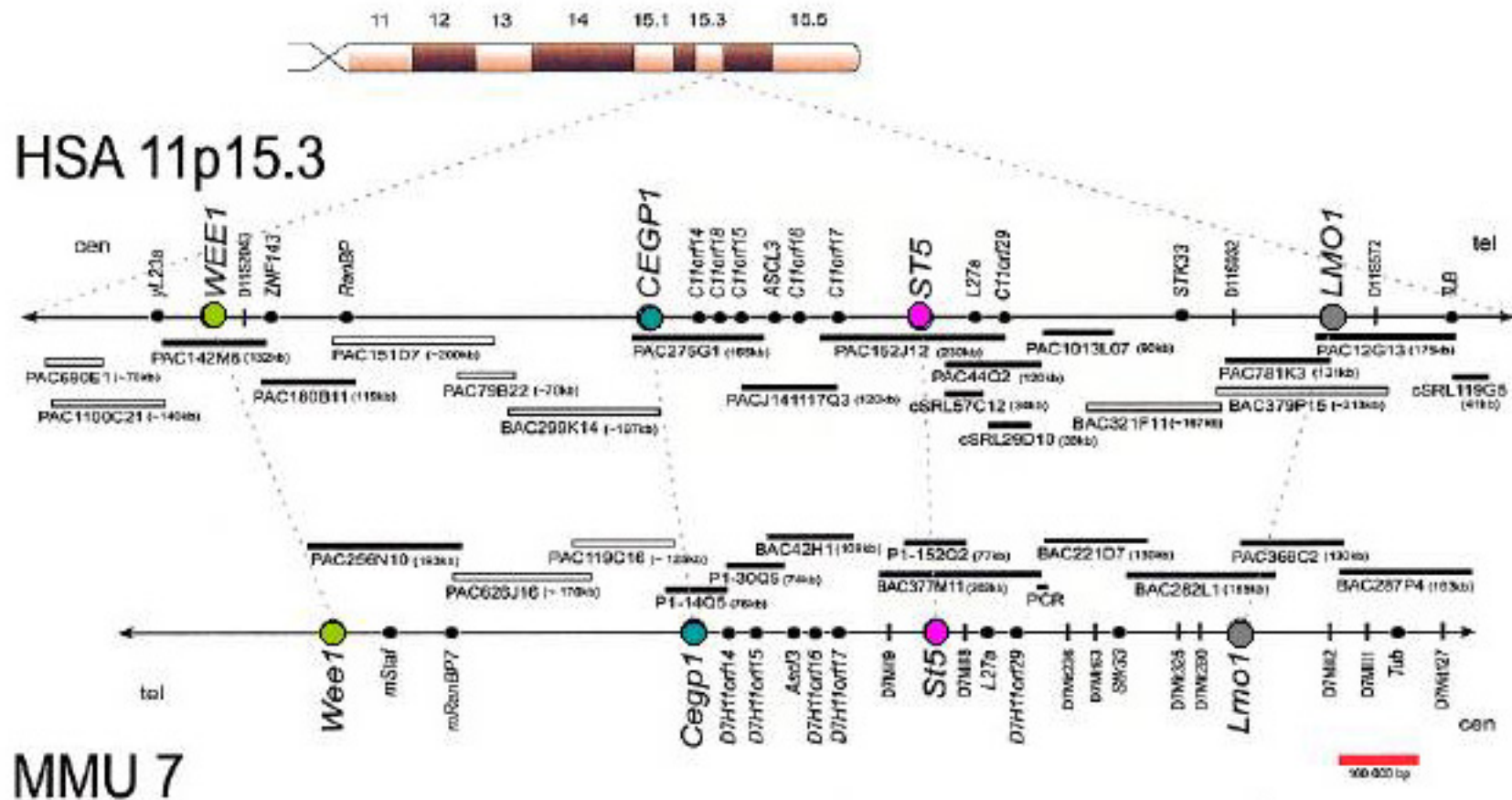
„nur funktionell wichtige DNA-Sequenzen (Exons, regulatorische Bereiche) bleiben während der Evolution konserviert“



Unser Ziel damals (1995):

1 Mio Bp in Mensch und Maus

in hoher Qualität ($\ll 1$ Fehler / 10 000) und annotiert



Fehlerhaftigkeit der *ab initio*-Genvorhersage

CEGP1

EXON	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	fp	%
Human	-	-	+	+	+	+	+	-	+	+	+	-	+	+	+	+	+	+	-	+	+	+	3	78
Mouse	+	+	+	+	+	-	+	+	+	+	+	+	+	+	+	-	+	+	+	+	+	+	4	91

C11ORF1

EXON	1	2	3	4	5	6	7	fp	%
Human	+	+	+	+	+	+	-	0	86
Mouse	+	+	+	+	+	+	+	0	100

C11ORF2

EXON	1	2	3	4	5	fp	%
Human	+	+	+	+	+	0	100
Mouse	+	+	+	+	+	0	100

OK!

ASCL1

EXON	2	fp	%
Human	+	1	100
Mouse	+	0	100

C11ORF3

EXON	2	3	4	5	fp	%
Human	-	+	-	-	2	25
Mouse	-	+	-	-	2	25

C11ORF4

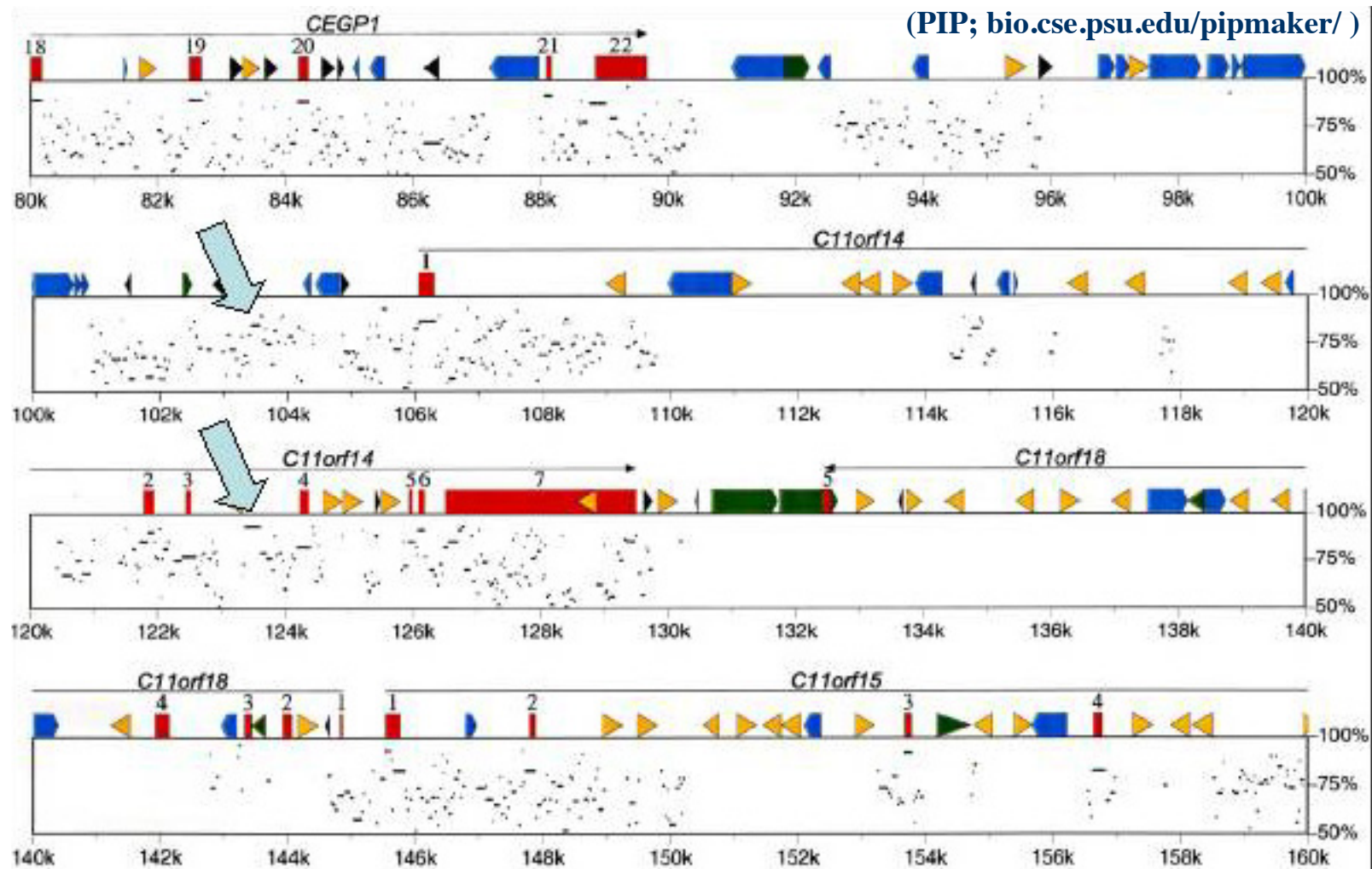
EXON	2	3	4	5	6	fp	%
Human	-	+	-	-	+	0	40
Mouse	+	+	+	-	+	0	80

+ korrekt vorhergesagt

-- falsch negativ

fp falsch-positive Vorhersagen

Mensch-Maus-Sequenzvergleich mit „Percent-Identity-Plot“



Mensch-Sequenz als ‚Vorlage‘ dargestellt. Skala rechts zeigt Sequenzidentität in der Maus.

Interspezies-Vergleich verbessert Genvorhersage

Mensch	TCTTCACTCATCCCT CAG AGTGGGTTT CAG CTCCCCTCTCCCTTTGACCA	54138
Maus	CCTCTACTCACATGG CAG ACGTGGCAT CAG CTCCCCCTCTCTTTGATAA	53437
<hr/>		
Mensch	TTAC CAG ATGTGGATGAGTGCTCTTTGGATAGGACCTGTGACCACAGCTGC	54188
Maus	TTAC CAG ATGTGGATGAATGTTTCGTTGGAAAGGACCTGTGACCACAGTTGC	53487
<hr/>		
Mensch	<u>ATCAACCACCCTGGCACATTTGCTTGTGCTTGCAACCGAGGGTACACCCT</u>	54238
Maus	<u>ATCAACCACCCTGGCACATTTATTTGTGCCTGCAACCCAGGGTACACTCT</u>	53537
<hr/>		
Mensch	<u>GTATGGCTTCACCCACTGTGGAGGTGAGCAGACTGCCCCTGTCAAGTTGG</u>	54288
Maus	<u>CTATAGCTTCACCCACTGTGGAGGTGAGTGGGTGGCCCCTCTATGGTGAG</u>	53587
<hr/>		

Aber: die exakte Vorhersage der korrekten Spleißstellen ist so nicht möglich

Kombination von Genom-Vergleich und *ab initio*-Genvorhersage verbessert das Ergebnis enorm...

Bsp: Annotation von Chromosom HSA 14

Table 3 Specificity and sensitivity of the non-expressed resources in the annotation of human chromosome 14						
	Exofish	BLAT-mouse	GENSCAN exons	FGENESH exons	Exofish + BLAT-mouse	Exofish + BLAT-mouse + GENSCAN exons + FGENESH exons
Annotation features	4,396	18,583	10,192	7,901	3,862	2,609
Specificity*	96%	31%	52%	61%	96%	99%
Sensitivity†	42%	70%	68%	63%	74%	83%

*The fraction of the annotation features of the various resources or a combination of resources that were included in the exons of the final gene annotation is indicated.

†The total number of annotated exons for all gene categories is 7,305 exons. The fraction of exons identified with the various resources or with a combination of resources is indicated.

GENSCAN, FGENESH
EXOFISH
BLAT-mouse

ab-initio Vorhersage
Fugu-Genomvergleich
Maus-Genomvergleich

Es wurde hier keine Info über cDNAs/ESTs/Proteine verwendet!!

Vergleichende Genomanalyse:

unerwartete Erkenntnisse

Ultraconserved Elements in the Human Genome

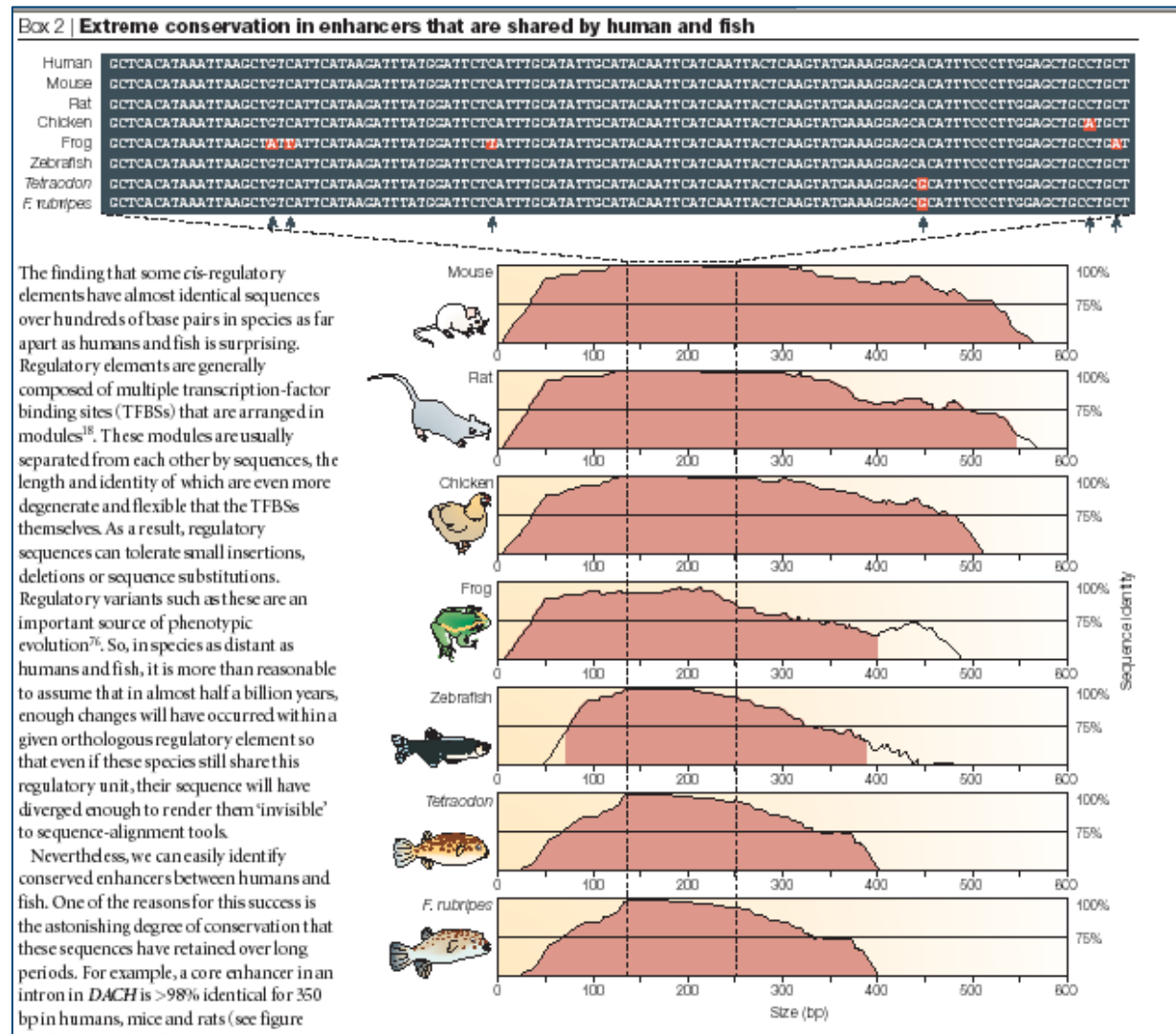
Gill Bejerano,^{1*} Michael Pheasant,³ Igor Makunin,³
Stuart Stephen,³ W. James Kent,¹ John S. Mattick,³
David Haussler^{2*}

There are 481 segments longer than 200 base pairs (bp) that are absolutely conserved (100% identity with no insertions or deletions) between orthologous regions of the human, rat, and mouse genomes. Nearly all of these segments are also conserved in the chicken and dog genomes, with an average of 95 and 99% identity, respectively. Many are also significantly conserved in fish. These ultraconserved elements of the human genome are most often located either overlapping exons in genes involved in RNA processing or in introns or nearby genes involved in the regulation of transcription and development. Along with more than 5000 sequences of over 100 bp that are absolutely conserved among the three sequenced mammals, these represent a class of genetic elements whose functions and evolutionary origins are yet to be determined, but which are more highly conserved between these species than are proteins and appear to be essential for the ontogeny of mammals and other vertebrates.

Vergleichende Genomanalyse: Ultraconserved Elements

Beispiel eines UCE mit möglicher Enhancer-Funktion

Warum ist UCE besser konserviert als die kodierende DNA?



a look at what is essentially an ordinary earthquake," Ellsworth says.

The mine is a challenging work environment. Moab Khotsong claims to be home to the world's deepest mine shaft: a 1000-meter descent at some 15 meters per second in a shaking cage, with only the arcs of head lamps piercing the darkness. At the drilling level, 95 floors down, it's a 10-minute ride in a rail carriage to the drill site. It can be eerily quiet, and the air carries an acrid scent of burnt rock and ammonia from recent dynamiting. "It's a smell that'd be associated with Hades," says Tullis Onstott, a geomicrobiologist at Princeton University who has also joined the project.

Onstott's goal is to learn whether earthquakes can favor the microbes that live in deep rock. In an experiment he's been trying to do for a decade, ever since it failed at a San Andreas drill site, he will install an automated sampling system, triggered by seismicity in the borehole at the fault. Past work has shown that an earthquake can release a pulse of hydrogen gas, which might have been trapped in the rock or generated by chemical reactions caused by the fracturing. Onstott hopes, as more aftershocks strike, to discover a cascade of microbial populations feeding off this chemical energy. If so, similar faults could become targets for a search for life on Mars, which might influence the selection of a landing site for NASA's Mars 2020 rover mission, Onstott says.

Mine operators might also get what they are seeking: a clearer picture of how human activity can trigger quakes. In the United States, the underground disposal of wastewater from oil and gas drilling is known to trigger quakes by boosting pressure in the rock pores. Pore pressure is unlikely to be a factor in South Africa, says Zélie Roches, a geophysicist at the University of Oklahoma in Norman. But he notes that the Oriskany quake, like some of the largest induced quakes in Oklahoma, occurred deep in basement rock along a previously unmapped fault. (Roches is planning to drill into a fault in Oklahoma, but he'll have to go 4.5 kilometers down.)

Plenty can still go wrong for DEIS. In May, operators shuttered another mine that the project planned to use, citing an increased risk of falling rocks. But Ogasawara says they have plenty of candidate earthquake faults for their campaign, including small, shallow ones accessible from nearby mines, which may allow them to target the very heart of an earthquake: the hypocenter, the exact place where a fault first starts to give way. "That's what we want to see," Ogasawara says. "We're expecting to show what a hypocenter looks like." ■

GENOMICS

Mysterious unchanging DNA finds a purpose in life

Ultraconserved elements have developmental roles

By Elizabeth Pennisi
in Cold Spring Harbor, New York

"First, thank you. Many of us have been waiting for this." That's how Eric Landry, director of the Broad Institute in Cambridge, Massachusetts, greeted genome scientist Diane Dickel earlier this month after she finished a talk on mysterious genetic sequences known as ultraconserved elements (UCEs). Speaking at the Biology of Genomes meeting here, Dickel had presented what may be the first compelling evidence of what these much-discussed stretches of DNA do. Most biologists are convinced they are important, as they have barely changed over tens of millions of years. But clues to their function had been scarce until Dickel's findings, which suggest they may play key roles in animal development.

Computational scientists first identified these anomalous stretches of DNA 13 years ago, by comparing the genomes of diverse mammals (*Science*, 28 May 2004, p. 1321). The University of California (UC), Santa Cruz, researchers found 481 genome segments, 100 to more than 700 bases long, that were 100% identical in humans, rats, and mice. These same segments were at least 99% identical in chickens and dogs, and many were also strongly conserved in fish. Many UCEs were also located outside genes, in regions of chromosomes many researchers considered nonfunctional. Why would this "junk" be conserved? "That was quite baffling," recalls Kelly Frazer, a genome scientist at UC San Diego.

Apparently off limits to evolutionary tinkering, UCEs presumably had some crucial function. Change a UCE, and an animal should die. Or so researchers thought. Three years later, a team from Lawrence Berkeley National Laboratory (LBNL) in Berkeley, California, created four genetically engineered mouse strains, each with a different UCE knocked out. The results were startling: The mice seemed to be just fine.

A decade later, Dickel, part of the same

LBNL lab group that couldn't originally find a purpose for UCEs, tried again, this time with the gene-editing technique called CRISPR. She wondered whether some UCEs are redundant, which would allow intact elements to fill in for deleted ones, so she used CRISPR to delete more than one UCE at a time.

She focused on four UCEs near the *ARX* gene, which codes for a transcription factor that regulates genes for proper brain and sexual development. One of those had been deleted by the original LBNL group, and all four elements are thought to be enhancers, stretches of regulatory DNA that drive gene activity, two in the upper forebrain of a developing mouse and two in the lower forebrain.

"The bottom line: These things are important in biology."

Diane Dickel, Lawrence Berkeley National Laboratory

At first, knocking out these elements, whether singly or in pairs, seemed to leave the mice unaffected. But when Dickel and a developmental neuroscientist took a close look at the brains of the mice, they detected abnormal numbers of neurons or a shrunken hippocampus—effects visible in both the double and, to a lesser degree, the single knockouts. "This was very satisfying in that it resolved a long-standing observation that really bothered us," Dickel says.

Although the brain differences may be too small to be significant for well-fed, coddled lab mice, they might be catastrophic for mice in the wild, where they struggle to find food and avoid predation, Dickel notes. And any evidence of a role for UCEs is welcome. As Harvard University geneticist Chao-ting Wu puts it: "It's really lovely that there is some function."

Dickel now plans to use CRISPR to surgically knock out individual bases in the ultraconserved *ARX* regions. Other researchers have suggested that UCEs in enhancers vary so little because they need to bind to multiple proteins to regulate gene activity—and the loss of single base would disable them. Dickel isn't convinced yet that they're right, but she's happy to have finally provided some purpose for UCEs. "The bottom line: These things are important in biology." ■

Downloaded from <http://www.sciencemag.org/> on November 13, 2016