

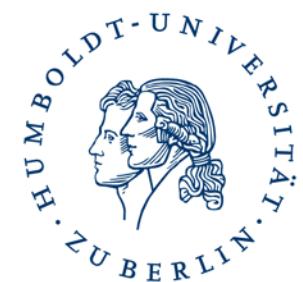
Algorithmische Bioinformatik

Wintersemester 2015 / 2016

Master: 10 SP Modul



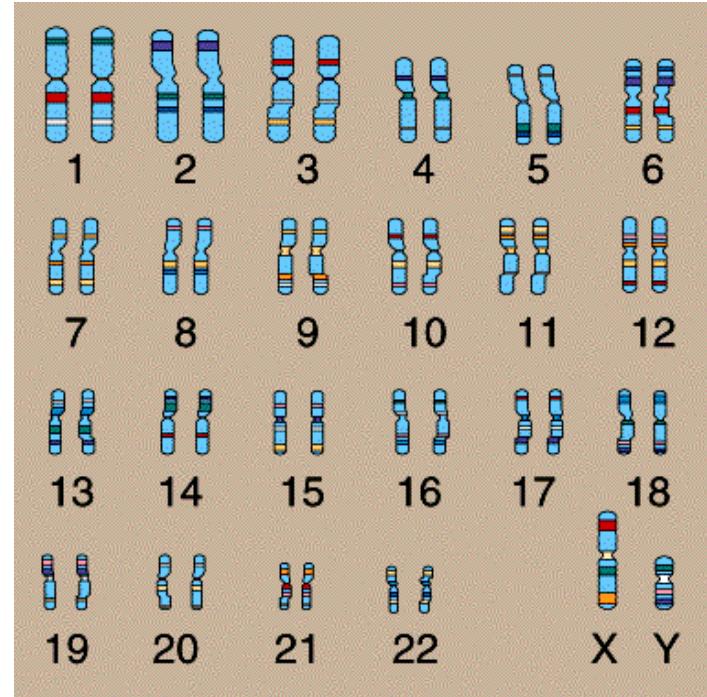
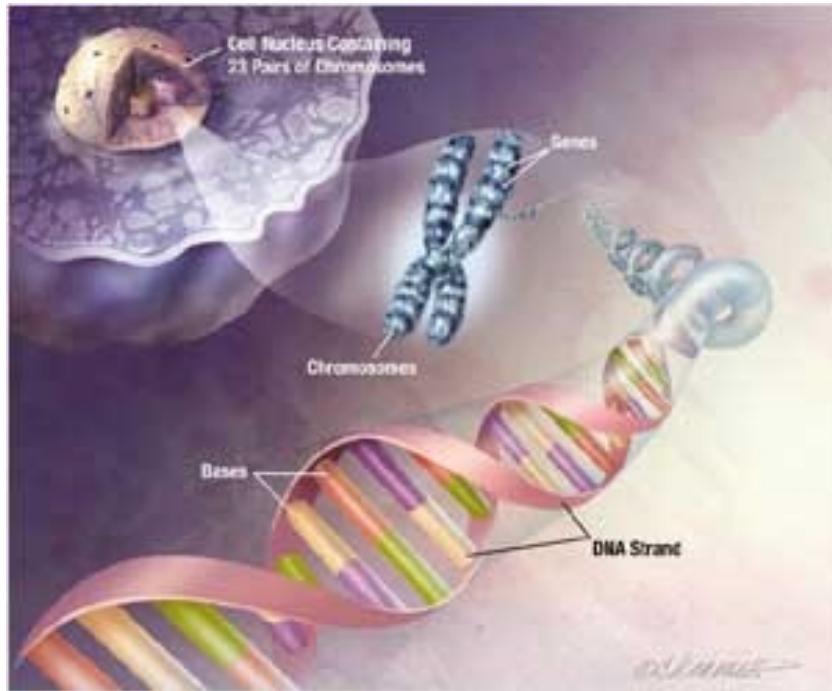
Ulf Leser
Wissensmanagement in der
Bioinformatik



Ziele für heute

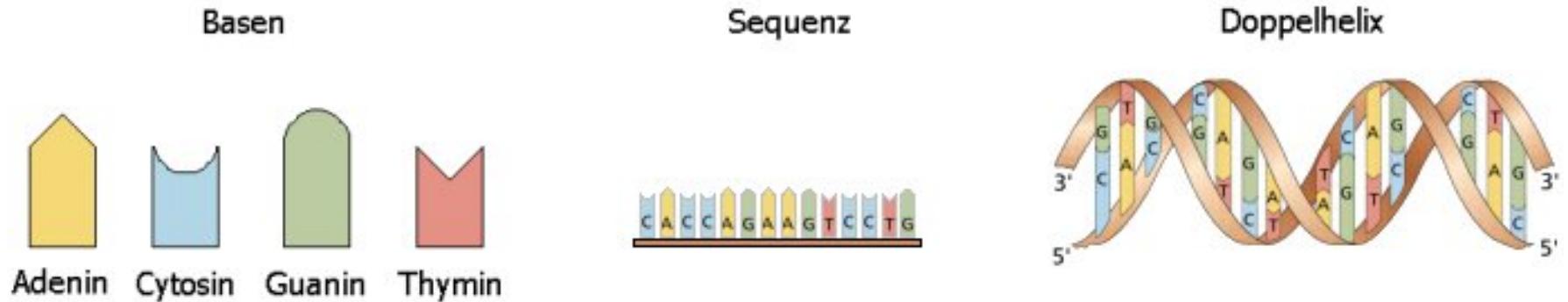
- Lust auf das Thema machen
- Gefühl für Rasanz der Entwicklung geben
- Überblick über die Vorlesung

Genome \approx String



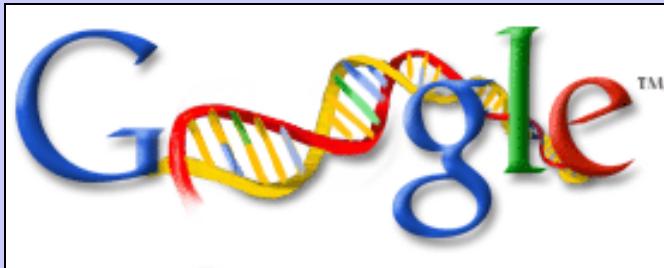
- Human genome: app. 3.000.000.000 letters $\in \{A,C,G,T\}$

DesoxyriboNucleicAcid



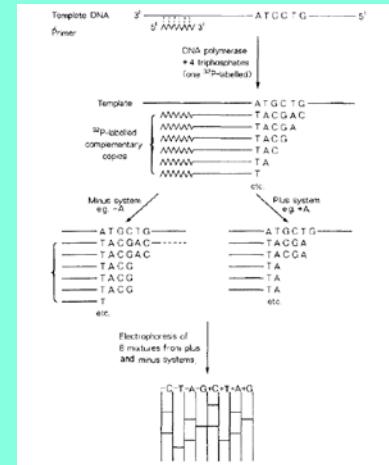
- DNA: Desoxyribonukleinsäure
- Träger der vererbten Information – Genom
- **Alles Leben** verwendet DNA (RNA) aus den selben 4 (5) Molekülen

Fast Development



1953
Double helix structure of DNA,
Watson/Crick

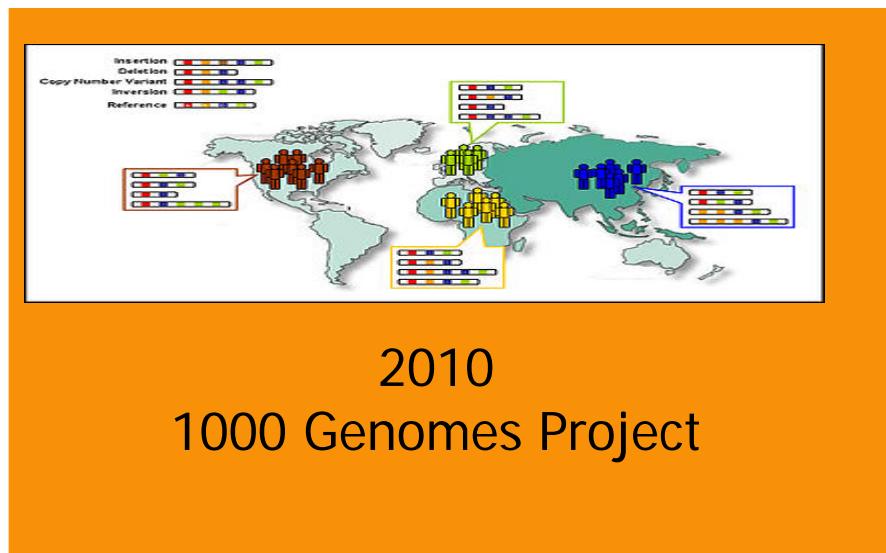
1975
High-throughput sequencing
Sanger/ Coulson



2003
First human genome sequenced
Took ~14 years, ~3 billion USD

Fast Development 2

The screenshot shows the nature news homepage. At the top, there is a red header with the 'nature news' logo. Below it, a navigation bar includes links for 'nature news home', 'news archive', 'specials', 'opinion', 'features', 'news blog', and 'events blog'. A 'comments on this story' link is also present. The main headline reads 'James Watson's genome sequenced at high speed'. Below the headline, it says 'Published online 16 April 2008 | 452, 788 (2008) | doi:10.1038/452788b News'. The date '2008' is displayed prominently below the headline. The text 'Genome of J. Watson finished 4 Months, 1.5 Million USD' is also present.



1000GP releases more data in first 6 months than EMBL collected in the 25 years before

Large Scale Sequencing Projects



50.000 samples: To obtain a comprehensive description of genomic, transcriptomic and epigenomic changes in 50 different tumor types and/or subtypes which are of clinical and societal importance across the globe.



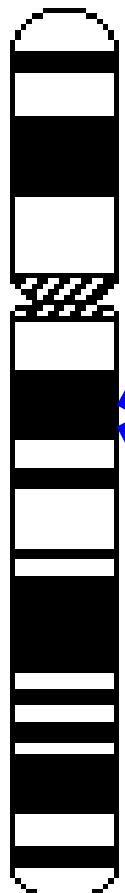
Genomics England ... is creating a lasting legacy for patients, the NHS and the UK economy through the sequencing of 100,000 genomes: [the 100,000 Genomes Project](#).



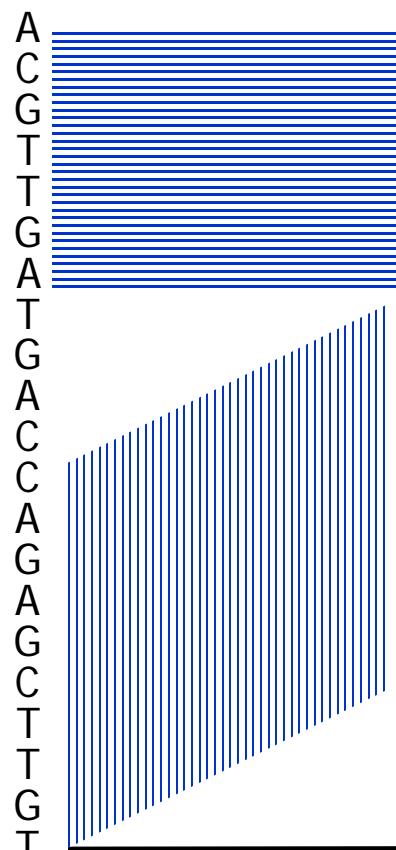
The Veterans Affairs (VA) Office of Research and Development is launching the [Million Veteran Program \(MVP\)](#) The goal of MVP is to better understand how genes affect health and illness in order to improve health care.

Was ist ein Gen?

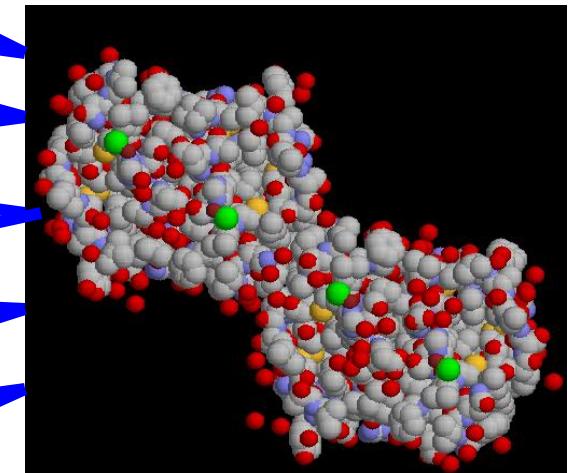
Chromosom DNA



RNA



Protein

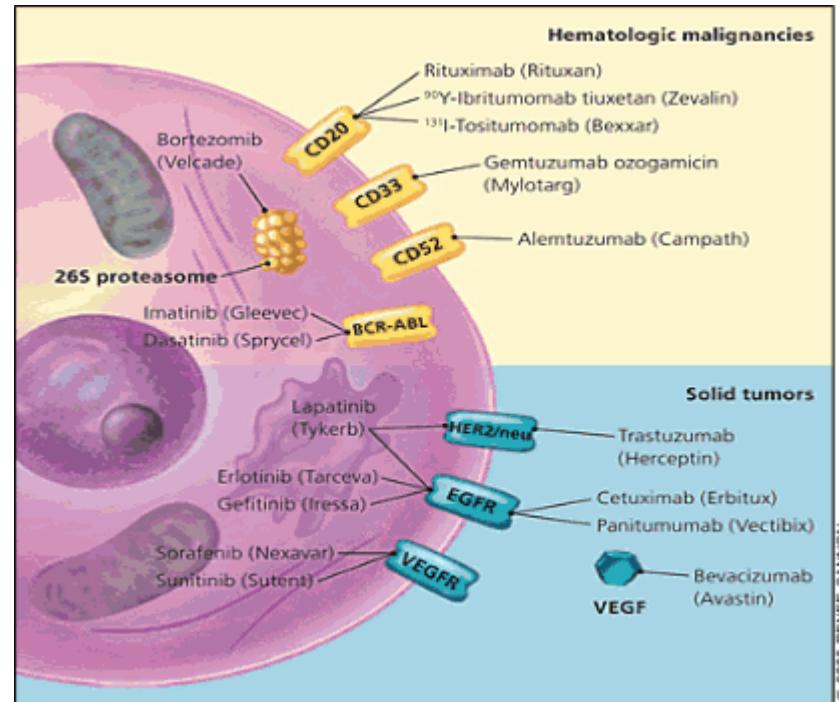


Genomics in a Nutshell

- ~2% are coding – genes being translated into proteins
 - Whole Genome Sequencing – WGS
 - Whole Exome Sequencing - WES
- ~20.000 genes, forming maybe 500K different proteins
 - ~3000 genes are conserved since ever (yeast)
 - We share ~95% of our genes with mice, rats, dogs, ...
 - ~25% of our genes have a still unknown function
- It's not only genes: miRNA, enhancer, binding sites, chromatin structure, epigenomics, ...

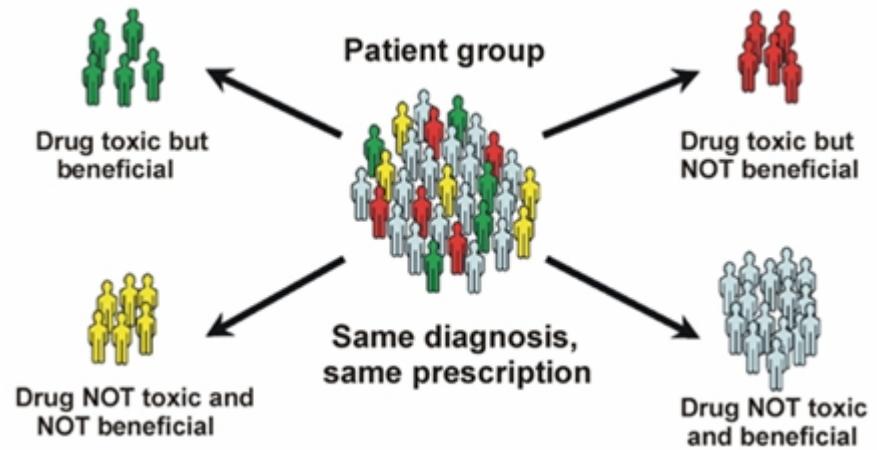
Genomics for Medicine

- Cancer, immunology, genetic diseases, infections
- Cancer
 - Cells proliferating uncontrolled, leaving their tissue
 - What goes wrong? Cell division, DNA repair, surface adhesion, cellular signaling
 - ~200 core cancer genes
- Targeted therapy: Drug attacking a mutated gene
- “Cancer is becoming a chronic diseases”



Precision Medicine, Personalized Medicine

- Tailor treatment to the **individual patient's genome**
- “Genome” – **mutation profile**
 - We know 10s of Millions of human mutations
 - Mutation – deviation from the norm?
 - Mutation – genomic subsequence rarely seen
- Requires **many genomes**
 - What is rare?
 - Often enough to obtain a **statistically robust association**
 - Most effects involve many mutations / genes
 - Combinatorial explosion



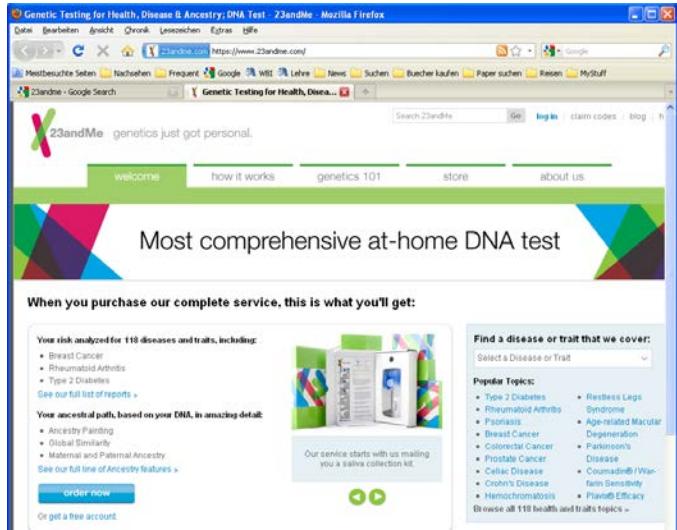
Wofür kann man sie benutzen?

- Kommerzielle Gentests im WWW
- 2 Firmen
- 32 Tests
- Preise: 100–1400€

Tabelle
Gentests, die im Internet in Deutschland bestellbar sind (Stand: Juli 2002)

Indikation*	Anbieter**	Untersuchungsgegenstand	Preis (inkl. MwSt.)
Alkoholverträglichkeit	2	keine Angaben (k. A.)	207,79 €
Alzheimer	2	k. A.	134,06 €
Alzheimer ¹⁾	1	E4-Allel des Apolipoprotein-E-Gens auf Chromosom 10	650,00 €
Angelman-Syndrom ²⁾	1	Deletion auf dem Chromosom 15	850,00 €
Anti-Aging-Risikoprofil	2	k. A.	653,61 €
Asthma/Chronische Bronchitis/Herzinfarkt/Schlaganfall	2	k. A.	512,81 €
Azoturie	1	31 Mutationen einschließlich einer 5T-Variante auf dem CFTR-Gen auf dem Chromosom 7	850,00 €
Bluthochdruck	2	k. A.	127,40 € 439,24 €
Diabetes ³⁾	2	k. A.	127,40 € 194,39 €
Dickdarmkrebs ³⁾	1	MLH1- und MSH2-Mutationen	1600,00 €
Entgiftungsfähigkeit	2	k. A.	811,10 €
Faktor V Leiden-Mutation	1	Gerinnungsfaktor-V auf dem langen Arm von Chromosom 1	400,00 €
Familiäre Hypercholesterinämie	1	Mutationen im Low-Density-Lipoprotein-Rezeptor-Gen und im Exon 26 Apolipoprotein-B-Gen	850,00 €
Familiäre Hyperlipoproteinämie Typ III	1	E2-Allel des Apolipoprotein-E-Gens auf Chromosom 19	500,00 €
Familiärer Brustkrebs ³⁾	1	BCRA1- und BCRA2-Mutationen	1400,00 €
Fettgen/Adipositas	2	k. A.	241,35 € 576,44 €
Fettstoffwechsel/Cholesterin	2	k. A.	395,48 €
Fragiles X-Syndrom ⁴⁾	1	FMR1-(fragile X mental retardation-)Gen des X-Chromosoms (Region Xq27.3)	950,00 €
Hämochromatose	2	k. A.	207,84 €
Hämochromatose	1	Austausch der DNS-Basen Guanin zu Adenin an der Position 845 und von Cytosin zu Guanin an der Position 187 des HFE-Gens auf dem Chromosom 6	500,00 €
Hyperhomocysteinämie	1	k. A.	550,00 €
Mukoviszidose (Cystische Fibrose)	1	Mutation eines Gens auf Chromosom 7	850,00 €
Muskeldystrophie	1	Deletionen (Verlust von DNA-Teilsequenzen) im Dystrophin-Gen auf dem X-Chromosom	850,00 €
Osteoporose	2	k. A.	103,89 € 191,01 €
Osteoporose	1	Mutation (Basenaustausch von Guanin zu Thymin) im Intron 1 des Kollagen Typ I Alpha 1-Gens	650,00 €
Ovarialkarzinom ³⁾	1	BCRA1- und BCRA2-Mutationen	850,00 €
Persönliches Ernährungsprofil	2	k. A.	841,32 €
Prader-Willi-Syndrom	1	Deletion oder Translokation auf dem langen Arm des Chromosoms 15 (15q11)	850,00 €
Prothrombin-Mutation	1	Austausch der DNS-Basen Guanin zu Adenin an der Position 20210 des Prothrombingens auf dem Chromosom 11	550,00 €
Risiko Alkohol- und Drogenabhängigkeit	2	k. A.	274,86 €
Thrombose	2	k. A.	134,06 € 281,52 €

State of the “Art”



- 6/2010: „Gentest-Firma vertauscht DNA-Ergebnisse ihrer Kunden“ (Nature Blog)
- 7/2010: US general accounting office compared 15 (4) companies: totally **contradicting results**

Modul Algorithmische Bioinformatik

- Vorlesung 4 SWS
 - Übung 2 SWS
-
- Sprechstunde: Nach Vereinbarung
Ulf Leser
Raum: IV.105
Tel: (030) 2093 – 3902
eMail: leser(..) informatik . hu-berlin . de

Termine und Prüfung

- Vorlesung
 - Dienstag, 13-15 Uhr
 - Donnerstag, 13-15 Uhr
- Übung
 - Donnerstag, 15-17 Uhr
- Erste Übung für alle: Donnerstag, 22.10.2015, 15 Uhr
- Voraussetzung für Prüfung
 - Bestehen aller Übungsaufgaben
 - Verständnis der Algorithmen

Gäste

- Die vier nächsten Doppelstunden
 - Raik Otto
 - Zellen, Chromosomen, Gene, Transcription und Translation, differentielles Splicen, ...
- Im Verlaufe des Semesters
 - ?

Literatur

- Primär
 - Dan Gusfield: „Algorithms on Strings, Trees, and Sequences“, Cambridge University Press, 1997 (ca. 60 Euro)
- Weitere
 - Ohlebusch: "Bioinformatics Algorithms", Verlag Enno Ohlebusch.
 - David Mount: „Bioinformatics. Sequence and Genome Analysis“, Cold Spring Harbour Press, 2001 (ca. 70 Euro)
 - Gibson & Muse: „A primer of genome sciences“, Sinauer Associates, 2001 (ca. 50 Euro)
- Sowie Originalliteratur

Webseite

Halbkurs Algorithmische Bioinformatik -- Mozilla Firefox

Datei Bearbeiten Ansicht Chronik Lesezeichen Extras Hilfe

Halbkurs Algorithmische Bioinformatik - zope.informatik.hu-berlin.de/forschung/gebiete/wbi/teaching/archive/ws1112/hk_algobio/ Meistbesuchte Seiten Nachsehen Frequent WBI Lehr Google News Projekte Buecher kaufen Paper suchen Reisen MyStuff Lymphomexplorer

WS 11/12

Halbkurs Algorithmische Bioinformatik

Professor Ulf Leser

Der Halbkurs "Algorithmische Bioinformatik" behandelt Algorithmen zur Lösung grundlegender Fragestellungen moderner Molekularbiologie. Nach einer Einführung in die Grundlagen der Molekularbiologie (Gene und Genome, Expression, Proteine, Regulation und Transkription) werden die folgenden algorithmischen Probleme behandelt: Exaktes Stringmatching, Stringmatching mit mehreren Patten, approximatives Matching, Indexstrukturen für Sequenzdatenbanken, Editabstand und Alignment, Multiples Alignment, Phylogenetische Bäume. Die Algorithmen werden jeweils anhand der zugrundeliegenden biologischen Fragestellung erklärt, wie z.B. Patternsuche in DNA- und Proteinsequenzen, Assembly von Teilsequenzen, Homologiesuche in Sequenzdatenbanken, und Berechnung evolutionärer Stammbäume. Die Vorlesung wird durch eine Übung begleitet.

Erste Vorlesung ist am Mittwoch, den 19.10.2011.

Voraussetzungen

Voraussetzung für den Besuch sind gute Kenntnisse in Algorithmen (z.B. Modul Algorithmen & Datenstrukturen). Kenntnisse in der Molekularbiologie werden nicht vorausgesetzt, sondern vermittelt.

Prüfungen

Prüfungen sind mündlich. Die Vorlesung ist anrechenbar für:

- Diplomstudiengang Informatik, Halbkurs praktischen Informatik, 8SP
- Master Informatik, 8SP
- Masterstudiengang Biophysik, Vertiefungsrichtung Bioinformatik, 8SP

Voraussetzung für die Zulassung zur Prüfung ist das Bestehen der Übung.

Literatur zur Vorlesung

Dan Gusfield: "Algorithms on Strings, Trees, and Sequences", Cambridge University Press. Die Vorlesung folgt in grossen Teilen diesem Buch. Zusätzliche Literatur wird in den jeweiligen Stunden angegeben.

Themen und Termine im Einzelnen

(Folien sind hier jeweils vor der Vorlesung als PDF verfügbar. Änderungen möglich).

- 19.10.11: Einleitung und Überblick

Suchen: publication Abwärts Aufwärts Hervorheben Groß-/Kleinschreibung

WBI Humboldt Universität Informatik

Institut für Informatik

English

Wissensmanagement in der Bioinformatik

Kontakt

Mitarbeiter

Veranstaltungen

Lehre

Archiv

WS 11/12

Modul Data Warehousing und Data Mining

Übung Data Warehousing and Data Mining

Halbkurs Algorithmische Bioinformatik

Übung Algorithmische Bioinformatik

Seminar Data and Text Mining

Forschungsseminar WBI - DBIS

SS 11

WS 10/11

SS 10

WS 09/10

SS 09

WS 08/09

SS 08

WS 07/08

SS 07

Ihre Bewertung

	Freundlich	Fragen	Sprache	Präsentation	Beispiele	Konzeption	Überblick	Viel neues	Kritische Auseinan	Nützlich	Lernziele	Materialien	Tempo	Schwerigkeit	Arbeitsaufwand	Dozent	Vorlesung	Abweichung vom C	Abweichung pro F _i	Alter	Geschlecht	Gefehlt	Teilnehmerzahl	Warum kommen?	Studiengang	Fachsemester	Korrektur 6gr	Korrektur 3gr			
1																															
4	5	5	5	5	4	5	4	5	5	4	5	5	3	4	4	2	2	20	1,18	23	M	1	3	3,4	M	7	0	0			
5	6	6	6	5	5	6	6	6	6	6	6	6	3	4	4	1	1	5	0,31	20	M	0	3	2,3	BA	3	1	0			
6	6	6	6	6	6	6	6	6	6	6	6	6	3	4	4	1	1	2	0,12	27	M	3	2	2,3	DI	13	0	0			
7	6	5	6	5	5	5	6	5	5	6	5	6	4	4	4	1	2	11	0,65	23	M	2	2	2,3	M	9	0	0			
8	Durchschnitt	5,83	5,67	5,83	5,50	5,33	5,67	5,67	5,67	4,80	5,40	5,60	5,33	3,17	3,67	3,67	1,17	1,33									25,7	1,8	2,7	6,2	
9	Wunschzahl	6,00	6,00	6,00	6,00	6,00	6,00	6,00	6,00	6,00	6,00	6,00	3,00	3,00	3,00	1,00	1,00											3,0	0,33		
10	Abweichung	0,17	0,33	0,17	0,50	0,67	0,33	0,33	0,33	1,20	0,60	0,40	0,67	-0,17	-0,67	-0,67	-0,17	-0,33													
11		0,17	0,33	0,17	0,50	0,67	0,33	0,33	0,33	1,20	0,60	0,40	0,67	0,17	0,67	0,67	0,17	0,33	7,70												
12																															
13	Besonders gut																														
14																															
15																															
16																															
17																															
18																															
19																															
20																															

Was wir ändern wollen

- Compressive Genomics
- Ev. Burrows-Wheeler Transform, Read Mapping
- Folien zu Editabstand kürzen, aus 3 macht 2
- Folien Neighbor Joining ausführlicher
- Übung: Aufgabe 5 wird ersetzt

Fragen

- Diplominformatiker?
 - Biophysiker?
 - Bachelor?
 - Semester?
-
- Spezielle Erwartungen?

Inhalte der Vorlesung

- Einführung
- Ein konkretes Beispiel
- Themen der Vorlesung

Beispiel 1: H5N1

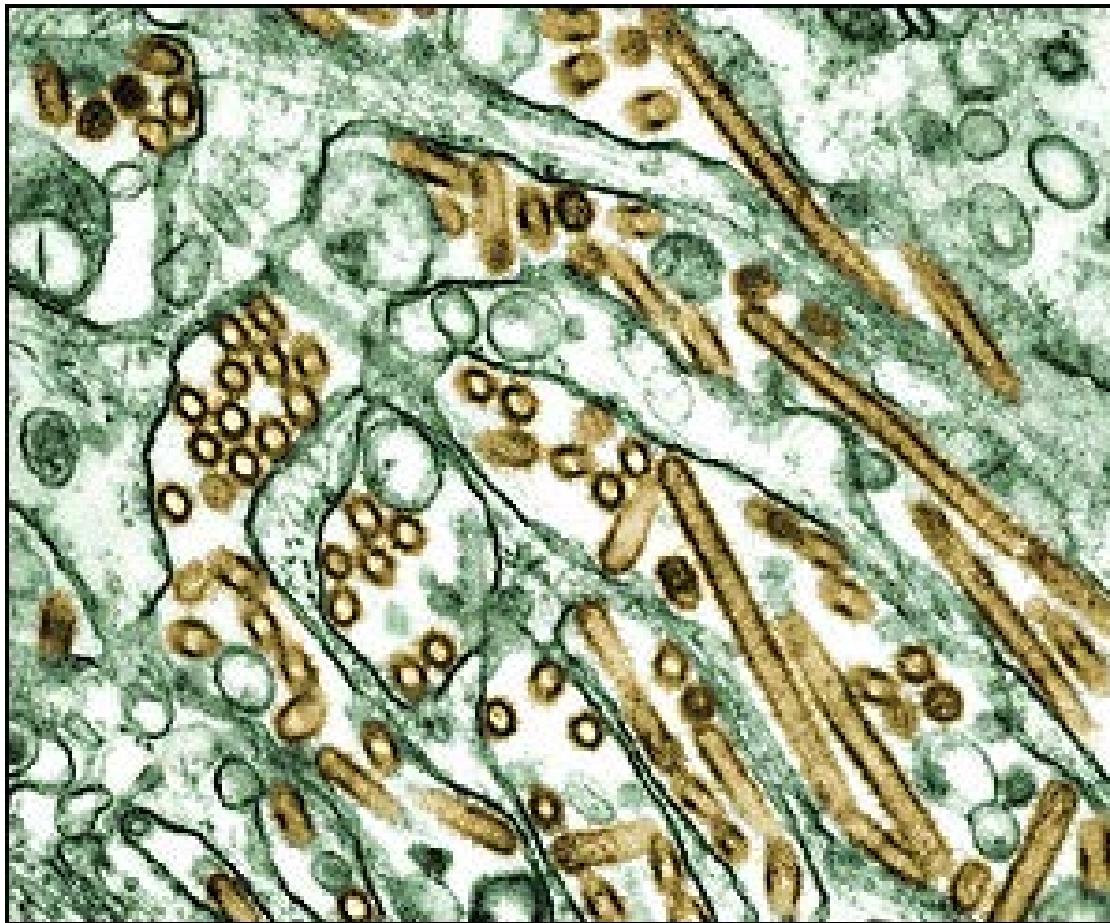
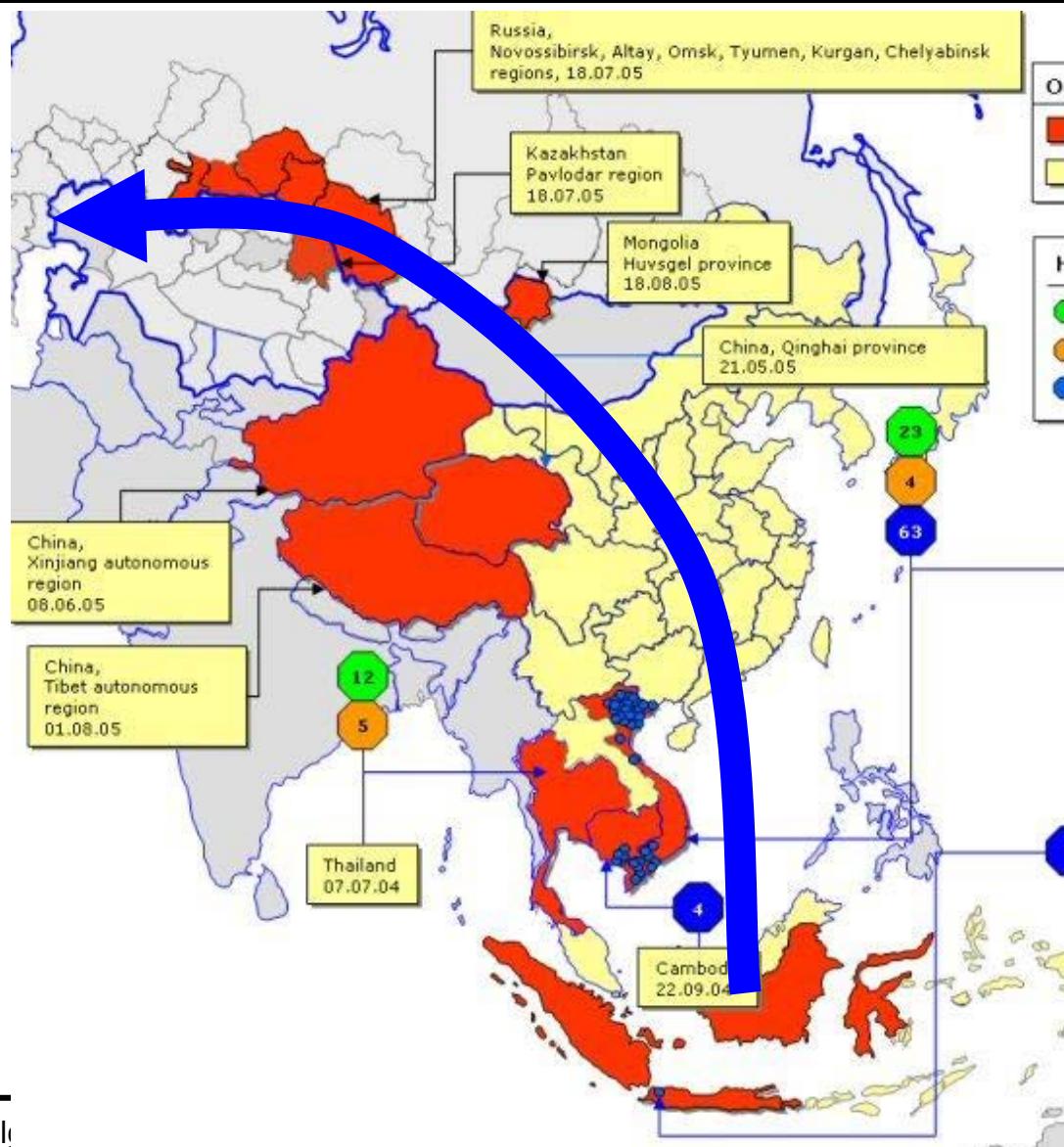
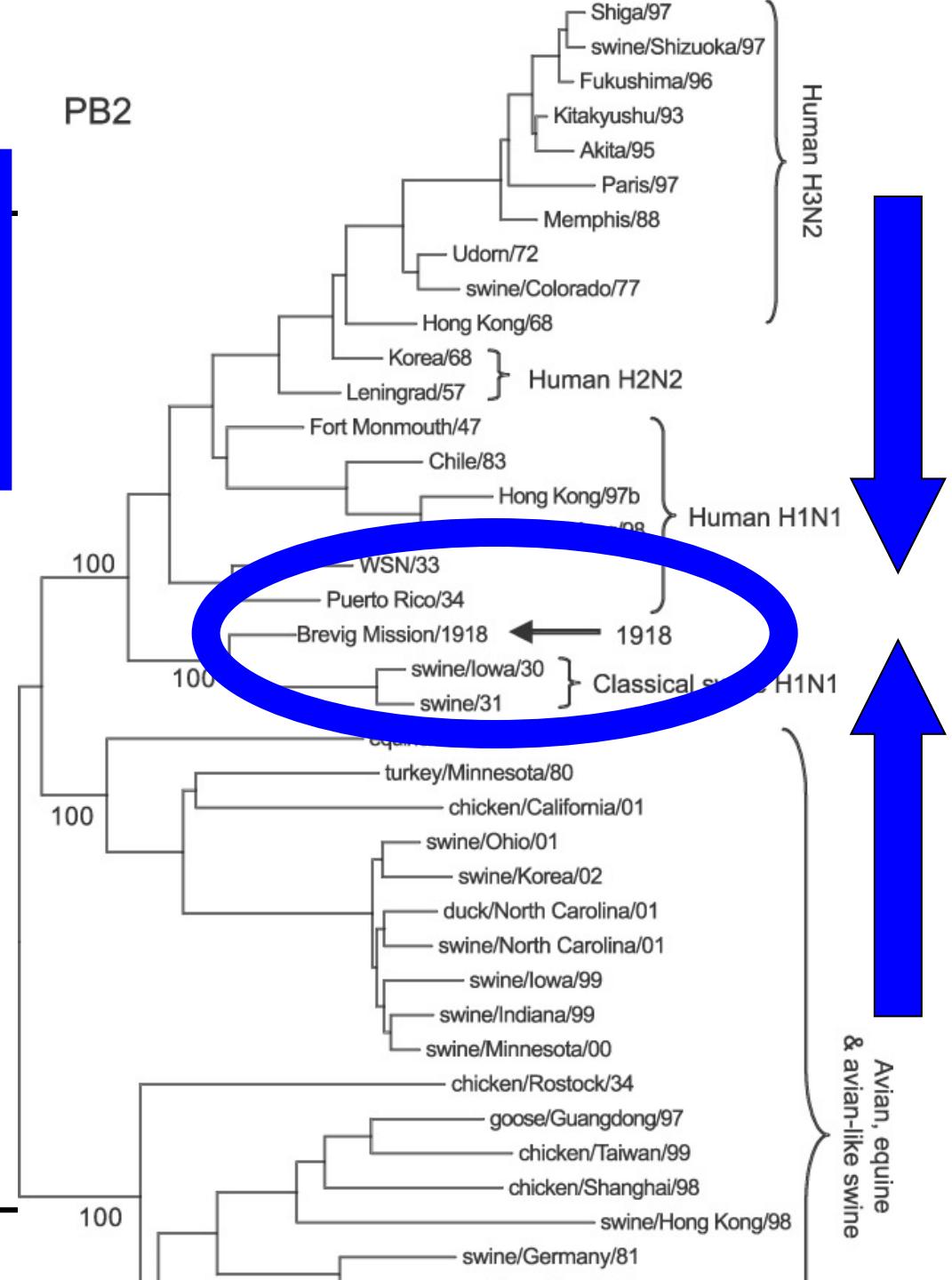


Foto: Centers for Disease Control

Migration

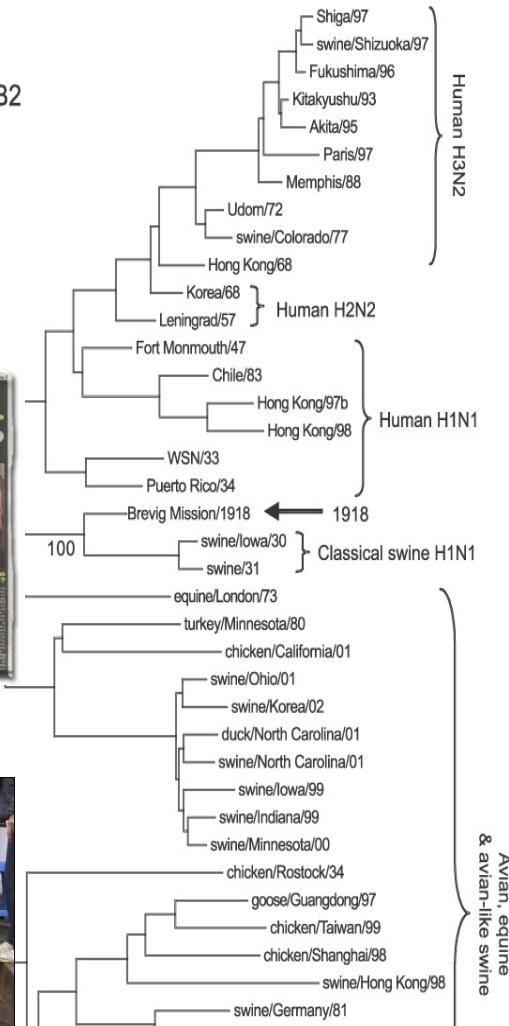
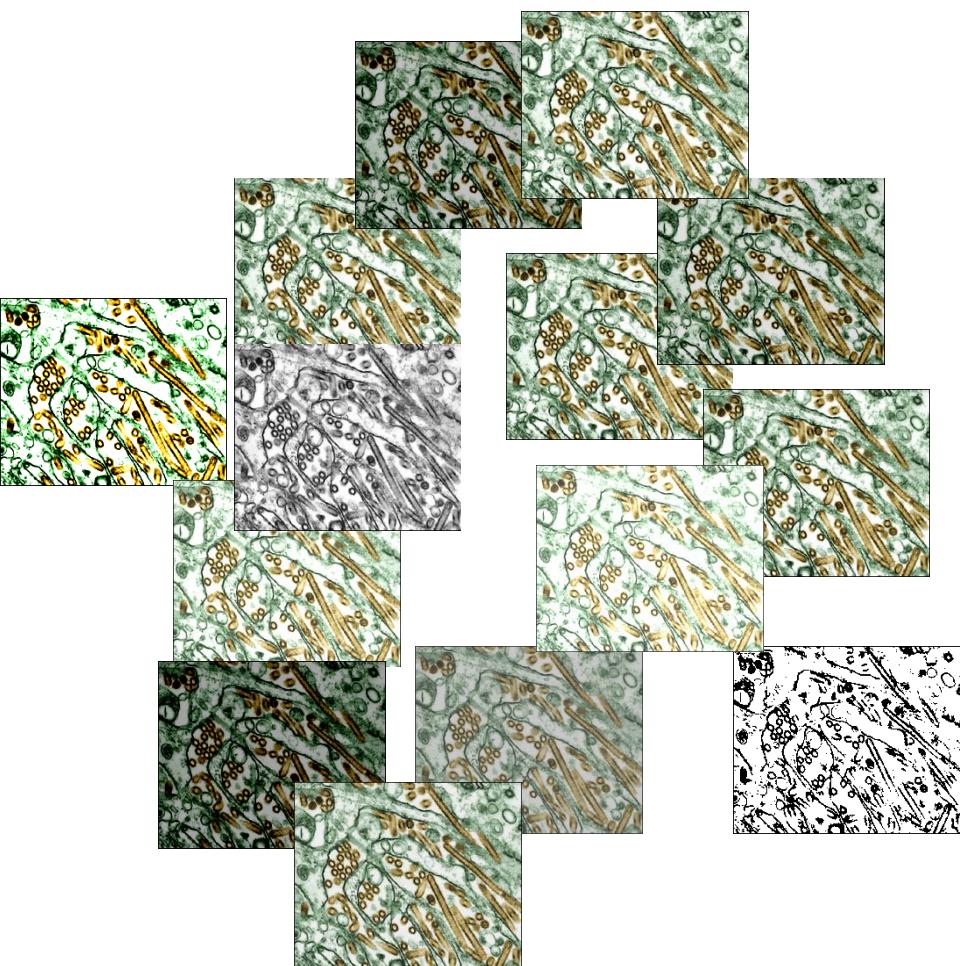


Vogelgrippe beim Menschen?

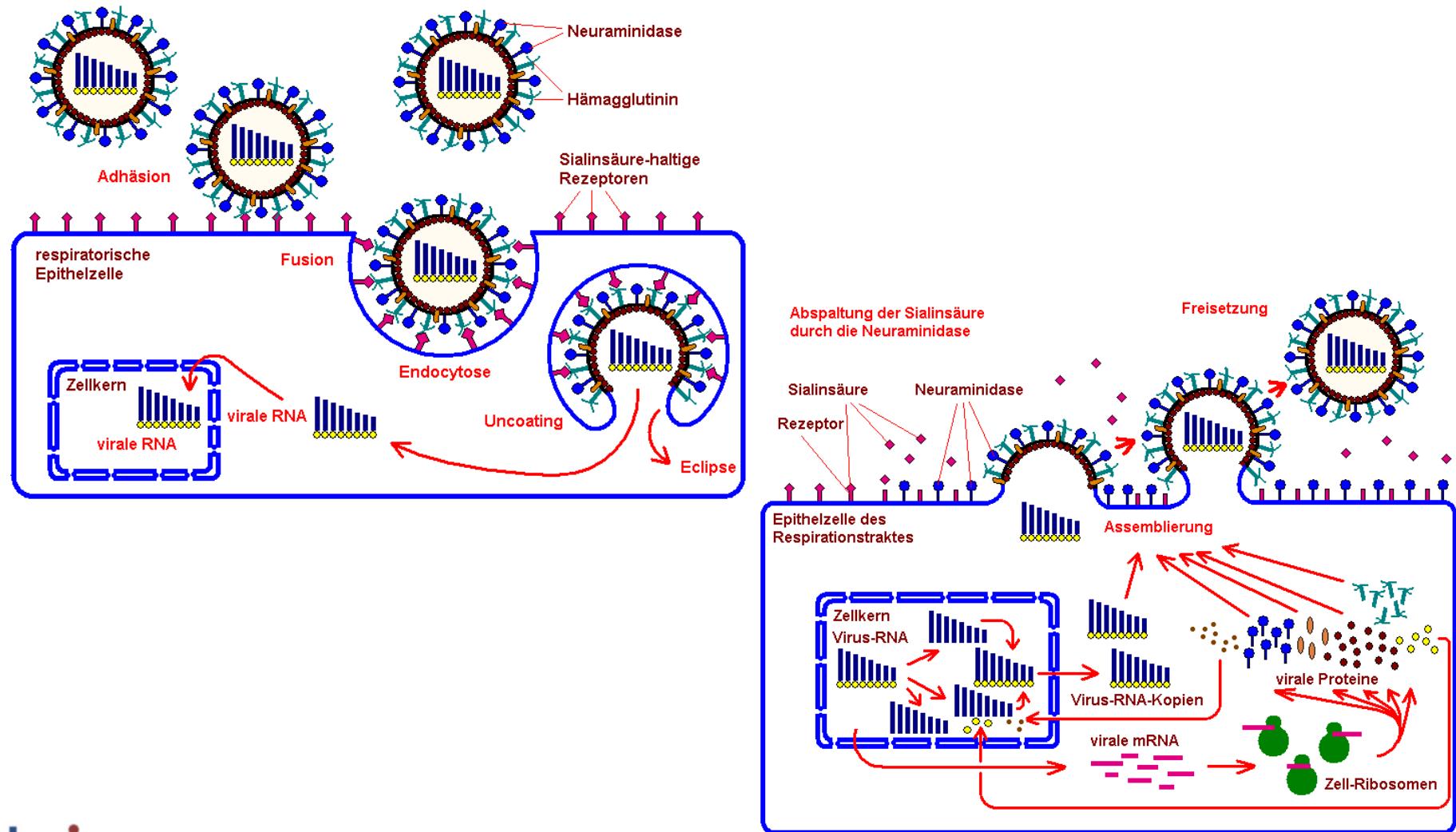


Nature. 2005 Oct 6;437(7060):889-93.

Wo kommt der Stammbaum her?

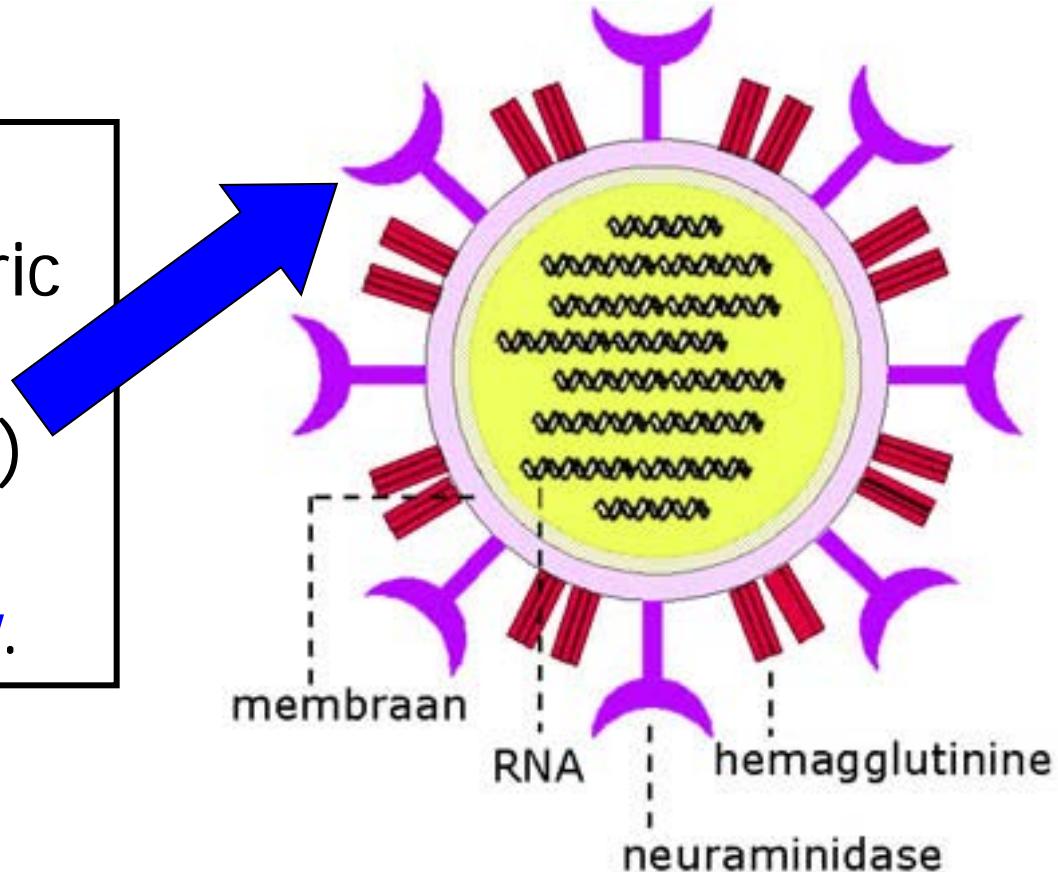


Viren



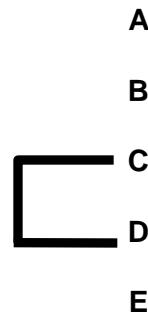
Grundlage für Stammbaumberechnung

The influenza A
viral heterotrimeric
polymerase
complex (... , PB2)
... having a **role**
in host specificity.

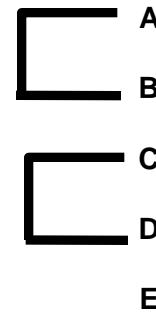


Konstruktion des Guide Trees

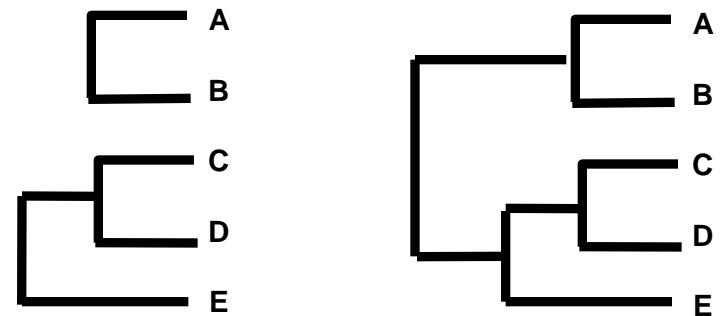
	A	B	C	D	E
A		17	59	59	77
B			37	61	53
C				13	41
D					21



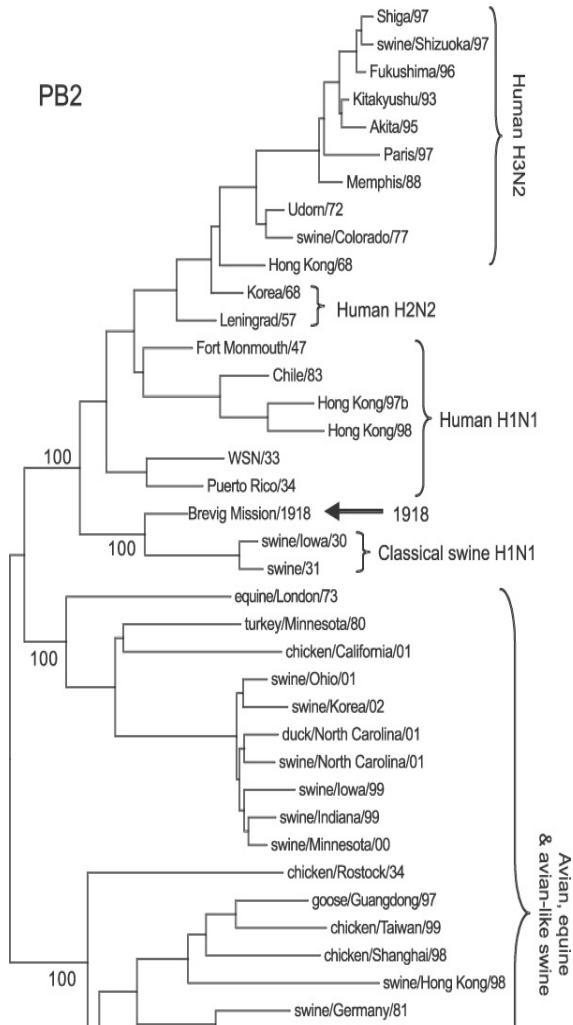
	A	B	E	CD
A		17	77	59
B			53	49
E				31



	E	CD	AB
E			65
CD			54



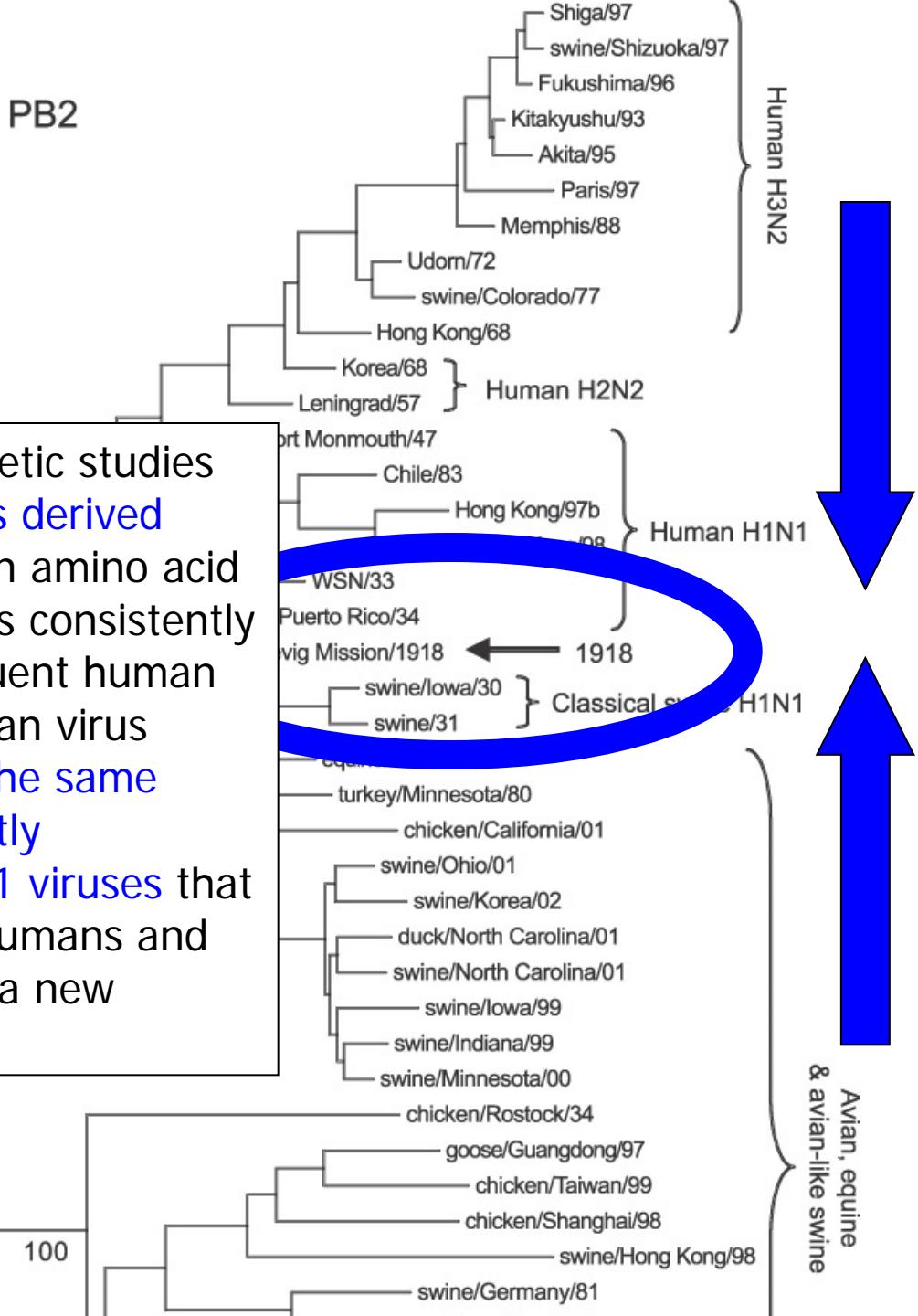
Was bedeutet der Stammbaum?



- Grundidee ist Evolution
 - Ur-Virus und Speziation
 - Richtungslose Mutationen
 - Selektion durch Umwelt (Host!)
- Benachbarte Stämme haben ähnliche Sequenzen
- Ähnliche Sequenzen bedingen ähnliche Funktion der Proteine und wahrscheinlich auch ähnliche Wirkung
- Um benachbarte Knoten zu „erreichen“, sind **nur noch wenige Mutationen** notwendig

Vogelgrippe beim Menschen?

These data support prior phylogenetic studies suggesting that the **1918 virus was derived from an avian source**. A total of ten amino acid changes in the polymerase proteins consistently differentiate the 1918 and subsequent human influenza virus sequences from avian virus sequences. Notably, a number of the same changes have been found in recently circulating, highly pathogenic H5N1 viruses that have caused illness and death in humans and are feared to be the precursors of a new influenza pandemic.

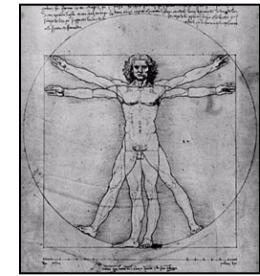
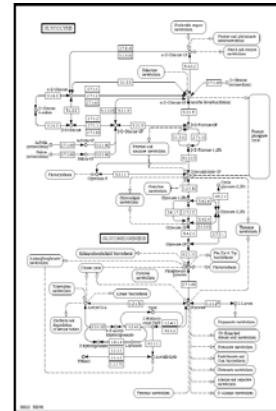
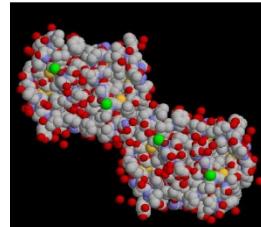
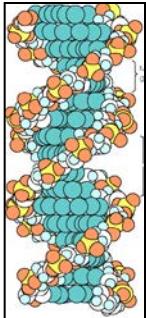


Nature. 2005 Oct 6;437(7060):889-93.

Inhalte der Vorlesung

- Einführung
- Ein konkretes Beispiel
- Themen der Vorlesung

Diese Vorlesung



Sequenzierung
Erkennung von
Genen
Verwandtschaft
zw. Spezies
Regulation &
Expression
RNA Gene –
„Dark matter“

Dreidimensionale
Faltung
Strukturvergleich
und -ähnlichkeit
Interaktion
Sekundärstruktur
Proteinidenti-
fikation

Netzwerkanalyse
Geschwindigkeit
von
Reaktionsketten
Stoffumsatz
Kompartamente
Muster und
Redundanz

Korrelation
Phänotyp/Genotyp
Arzneimittel-
empfänglichkeit
Studienstatistik
und -verwaltung

Stringalgorithmen

- Gegeben ein Template T und ein Pattern P. Finde alle Vorkommen von P in T in möglichst kurzer Zeit
 - Exaktes Matching
 - Z-Box
 - Boyer-Moore
 - Knuth-Morris-Prath
 - Varianten
 - Suche nach mehreren P: Aho-Corasick, Keyword Trees
 - Suche mit Wildcards
 - Suche mit regulären Ausdrücken (= endlichen Automaten)
- Vorverarbeitung von P
Schnellster in Praxis
Elegante Analyse; Erweiterbar

Indexstrukturen

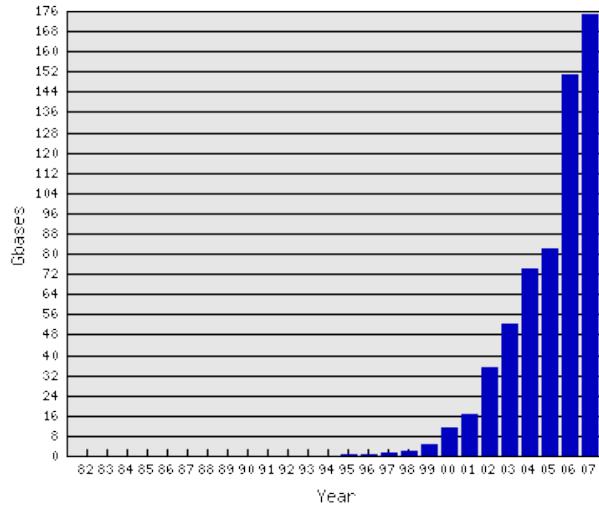
- Gegeben ein festes T und dauernd wechselnde P. Finde eine Datenstruktur für T, die die Suche nach allen P in möglichst kurzer Zeit gestattet
- Grundlegende Datenstruktur: **Suffixbäume**
- Suche und Konstruktion
 - Ukkonen's linearer Algorithmus
- Verschiedene Anwendungen
 - Längster Substring
 - Längster Repeat
- **Suffixarrays**

Approximatives Stringmatching

- Gegeben ein Template T und ein Pattern P. Finde alle Vorkommen von Substrings „ähnlich“ zu P in T in möglichst kurzer Zeit
- Was heißt überhaupt ähnlich?
 - Ähnlichkeitsmaße, Edit-Abstand, Alignierung
- Naiver Algorithmus benötigt exponentielle Laufzeit
 - Verbesserung durch dynamische Programmierung
 - Erreicht quadratische Laufzeit
- Viele Varianten: Globale, lokal, end-free, ...

Heuristiken

- Quadratische Laufzeit ist zu teuer
 - Genomanalyse benötigt Suche auf allen bekannten Sequenzen
 - Celera Sequenzierung:
All-against-all Vergleich von 28.000.000 Teilsequenzen
- Also: Heuristiken, z.B. BLAST
 - Suche nach „Seeds“ mit exakten Matches
 - Verlängerung und Zusammenfügung der Seeds zu Matches
 - Findet nicht alle Hits, aber die meisten „interessanten“
 - „BLAST“ ist fast **Synonym für Bioinformatik** geworden

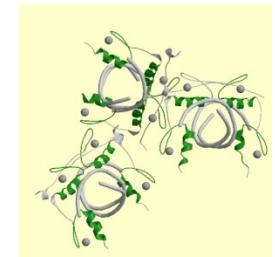


Multiples Alignment

- Gegeben eine Menge von Strings. Ein Multiple Sequence Alignment (MSA) ist eine Anordnung der Strings mit Spaces untereinander

YVCK..	.LCN.	.FAFKTKGNL	TKHMKSK	.AH
YRCPR.	ENCD..	.RTYTTKFNLKSHI	LTLT..	.FH
FRCGY.	KCGG..	.RLYTTAAHHLKVHERA	..H	
YRCE..	KCG..	.KMYKTERCLKVHNLV	..H	
FSCS..	QCD..	.ESFVQRSELELHRQL	..H	
FPCE..	QCD..	.EKFKTEKQLERHVKLT	..H	
FQCN..	QCG..	.ASFTQKGNLLRRHIKL	..H	
FKCH..	LCY..	.RCFQQQTNLDRHLKK	..H	
FRCK..	RCR..	.TRFRQQSELKKHHMKT	..H	
FECN..	VCG..	.SAFRLQLYLSEHQKT	..H	
MSCKV..	CD..	.RVFYRLDNLRSHLKQ	..H	
FSCQ..	HCH..	.RAFADRSNLRRAHLQT	..H	
FRCG..	YCG..	.RAFTVKDYLNKHLTT	..H	
HVCWV..	PGCH..	.RAFSRSNDNLNAHYTK	..TH	
ITCAH..	CD..	.WSFDNVMKLVRHRCV	..H	

Quelle: Pfam, Zinc finger domain



- Hauptziel von MSAs: Finde das „Gemeinsame“ der Sequenzen
 - Funktionen werden oft von sehr kurzen Sequenzstücken bestimmt
 - Welcher Teil eines Proteins bestimmt die Funktion?
 - Wie kann man Proteine in Familien anordnen?
- „Gute MSAs“ sind nicht klar definiert
 - Konkretes Maß zur Güte hängt von der Anwendung ab

Genomanalyse - Genvorhersage

- Welche Elemente eines Genoms sind interessant?
- Gene und deren Struktur
 - Promoterregionen, Start Site, Exons, Introns, ...
 - Verschiedene Regionen haben verschiedene Eigenschaften
 - Modellierung als Features
- Aufgabe: Finde die wahrscheinlichste Modulanordnung, gegeben ein Modell

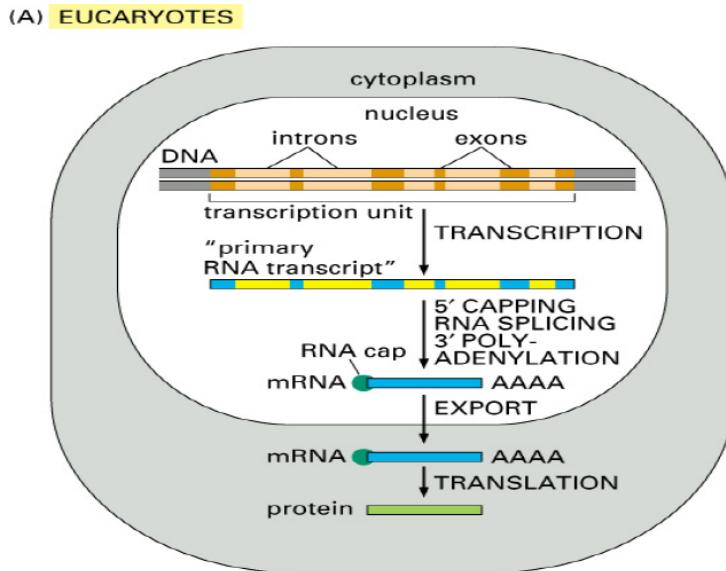
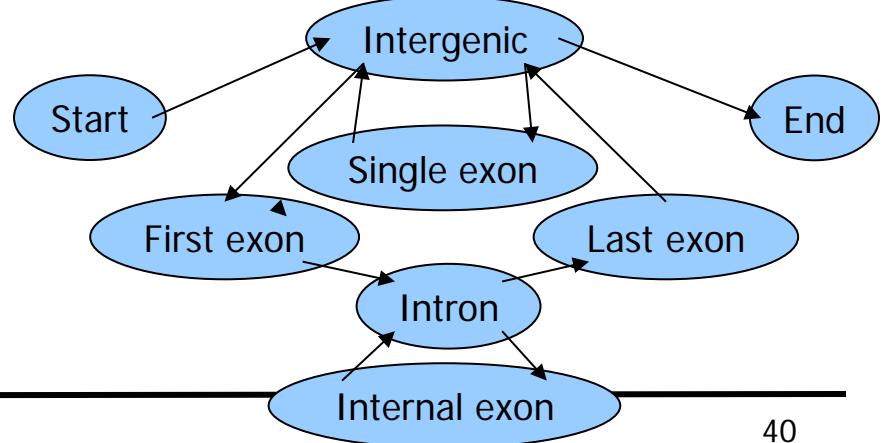
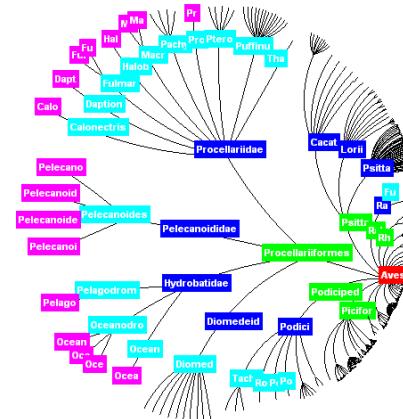


Figure 6-21 part 1 of 2. Molecular Biology of the Cell, 4th Edition.



Phylogenetische Bäume

- Grundannahme
 - Spezies entstehen durch Evolution
 - Also gibt es gemeinsame Vorfahren; Spezies stehen in Vater-Kind Beziehungen
- Phylogenie = „Baum der Evolution“
 - Auch: Berechnung des Evolutionsbaums
 - Beantwortung des Taxonomie-problems auf molekularer Basis



Selbsttest

- Wie viele Gene gibt es ungefähr in einem menschlichen Genom? Wie viele Basen? Wie viele Proteine?
- Was ist ein phylogenetischer Baum?
- Wie heißen die vier Basen eines Genoms?
- Wie viel Speicherplatz braucht man für 1000 Genome? Wie kann man sparen?