

Populationsgenetik 3: Kopplungsungleichgewicht (LD)

Peter N. Robinson

Institut für medizinische Genetik
Charité Universitätsmedizin Berlin

28. Juni 2008

Hardy-Weinberg-Gesetz

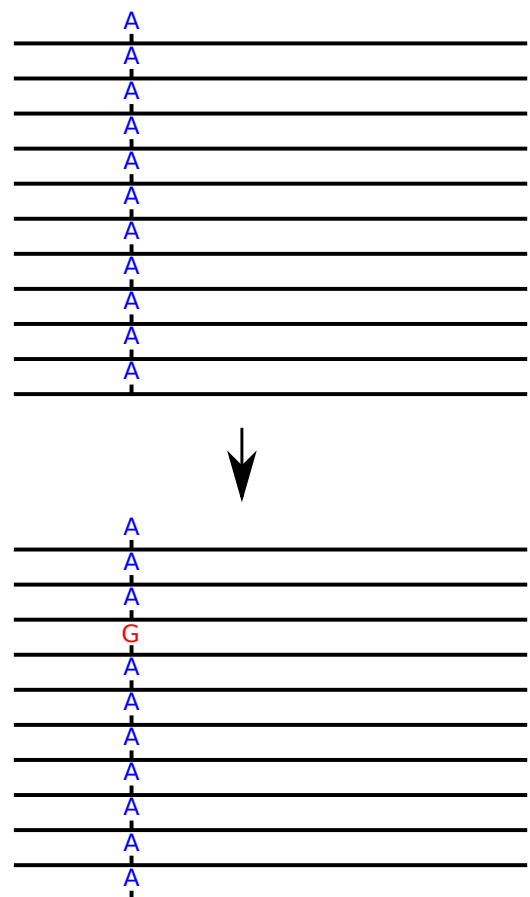
- Letztes Mal ...
- Eigenschaften eines *einzelnen* Genlocus:
Allelfrequenzen, Genotypfrequenzen
- Hardy-Weinberg-Gesetz: Beziehung
zwischen Allel- und Genotypfrequenzen

$$p^2 + 2pq + q^2$$

- Dieses Mal ...
- Eigenschaften von Gruppen von Genorten
- Haplotypen
- Kopplungsgleichgewicht
- Kopplungsungleichgewicht

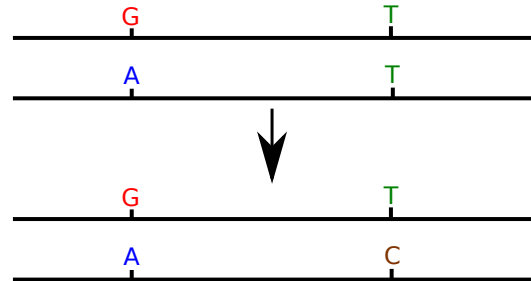
Eine Geschichte zweier Mutationen

- Heute existierende Allele
entstanden durch weit
zurückliegende
Mutationsereignisse



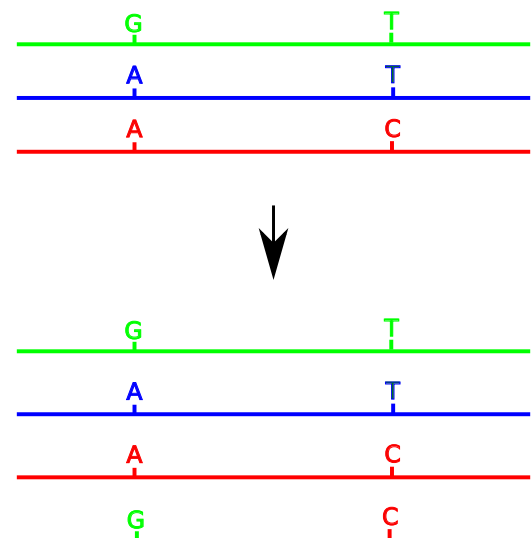
Eine Geschichte zweier Mutationen (2)

- Nach einem ersten Mutationsereignis entstand eine weitere Mutation auf demselben Chromosom ...



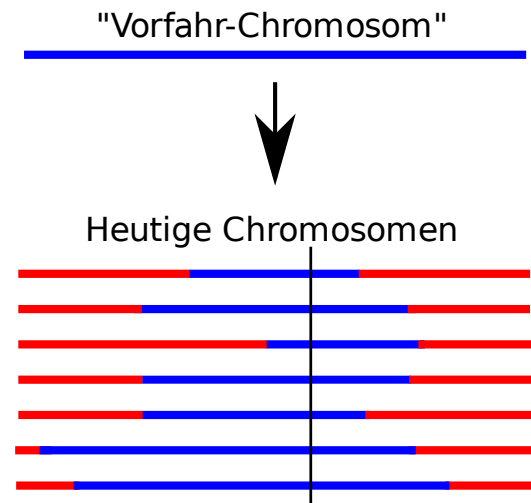
Eine Geschichte zweier Mutationen (2)

- Rekombination führt zu neuen Kombinationen der Allele

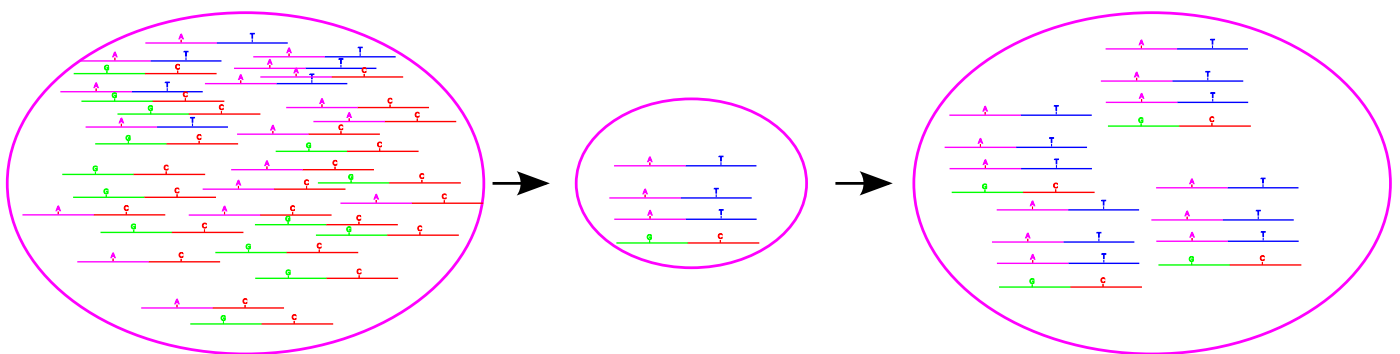


Kopplungsungleichgewicht

- Linkage Disequilibrium: **LD**
- Chromosomen sind Mosaik
 - ▶ Rekombination
 - ▶ Mutation
 - ▶ Genetische Drift
 - ▶ Natürliche Auslese (Selektion)
- Kombinationen von Allelen in nah beieinander liegenden Loci: weit zurückliegende ("ancestral") Haplotypen



Gründereffekt und LD



- Eine in der Stammpopulation bestehende, große genetische Variabilität reduziert sich bei der Gründung einer Kolonie durch wenig Individuen
- Häufige Ursache von LD in menschlichen Populationen

- In einer späteren Vorlesung werden wir die Bedeutung des LD für die Entdeckung von genetischen Varianten, die mit einer erhöhten Anfälligkeit für häufige Krankheiten wie Diabetes, Herzinfarkt, Schlaganfall, Allergie, ..., eingehen
- Dieses Mal wollen wir die mathematischen Hintergründe und die biologischen Grundlagen erklären
- *Wichtig*: Unterscheide Linkage und Linkage Disequilibrium (LD):
 - ▶ Linkage: gemeinsame Vererbung zweier Loci in Familien
 - ▶ LD: Beziehung zwischen zwei Allelen an 2 Loci in einer Population

Genotyp vs. Haplotyp

1 Haplotyp

- ▶ "haploider Genotyp"
- ▶ Eine Rekombination wird nur selten zwei Loci trennen, die nahe beieinander auf einem Chromosom liegen
- ▶ Deshalb werden Gruppen von Allelen, die auf demselben Chromosomenabschnitt liegen eher zusammen (als durch Rekombination getrennt) als **Block** übertragen werden

2 Genotyp

- ▶ die (diploide) genetische Ausstattung eines Individuums an einem oder mehreren Loci. Beide Exemplare eines Allels werden berücksichtigt, z.B. der Genotyp an einem Locus mit Allelen A und a kann AA , Aa oder aa sein.

- Ein Kopplungsungleichgewicht (**LD**) besteht zwischen zwei Genloci, die auf einem Chromosom eng beieinander liegen, und deshalb zusammenvererbt werden.
- Zwei eng beieinander liegende Loci werden dann nicht zusammen vererbt, wenn zwischen ihnen eine Rekombination erfolgt.
- Begriffe
 - ▶ Haplotypfrequenz
 - ▶ D , D' , r^2

LD: D

- A und a: zwei Allele von einem Locus (Allelfrequenz p_A und p_a)
- B und b zwei Allele eines anderen Locus. (Allelfrequenz p_B und p_b)
- Häufigkeiten von Kombinationen dieser Allele innerhalb einer Population von Gameten[†]: p_{AB} , p_{Ab} , p_{aB} und p_{ab} .

[†] zur Erinnerung sind Gameten Keimzellen, d.h. haploide Zellen, die im Gegensatz zu diploiden Zellen jeweils nur ein Exemplar jedes Locus haben.

- Die entsprechenden Allelfrequenzen ergeben sich aus der Summe der Genotypfrequenzen:

$$p_a = p_{ab} + p_{aB}$$

$$p_A = p_{Ab} + p_{AB} = 1 - p_a$$

und

$$p_b = p_{ab} + p_{Ab}$$

$$p_B = p_{AB} + p_{aB} = 1 - p_b$$

Kopplungsgleichgewicht

- Sind die beiden Genloci untereinander im Kopplungsgleichgewicht[†], dann ist die Wahrscheinlichkeit, dass eine Gamete das Allel a aufweist, unabhängig von der Wahrscheinlichkeit, dass sie das Allel b aufweist

$$p_{ab} = p_a \times p_b$$

[†]zum Beispiel weil die Genloci auf unterschiedlichen Chromosomen gelegen sind

- Sind die Loci nicht im Kopplungsgleichgewicht, dann gilt

$$p_{ab} \neq p_a \times p_b$$

.

- Wir führen die Variable D ein, um die Abweichung vom Kopplungsgleichgewicht zu beschreiben:

$$p_{ab} = p_a p_b + D \quad (1)$$

LD

Hieraus folgt

$$\begin{aligned} p_{aB} &= p_a - p_{ab} \\ &= p_a - p_a p_b - D \\ &= p_a (1 - p_b) - D \\ &= p_a p_B - D \end{aligned}$$

Eine analoge Berechnung zeigt:

$$p_{Ab} = p_A p_b - D$$

und[†]

$$p_{AB} = p_A p_B + D$$

[†]s. Skript

Die am häufigsten verwendete Definition der LD-Koeffiziente D ist jedoch

$$D = p_{AB}p_{ab} - p_{Ab}p_{aB} \quad (2)$$

Diese Formel leitet sich von der Definition (1) ab:

$$\begin{aligned}
 D &= p_{ab} - p_a p_b \\
 &= p_{ab} - (p_{aB} + p_{ab})(p_{Ab} + p_{ab}) \\
 &= p_{ab} - p_{aB}p_{Ab} - p_{ab}p_{ab} - p_{aB}p_{ab} - p_{Ab}p_{ab} \\
 &= p_{ab}(1 - p_{ab} - p_{aB} - p_{Ab}) - p_{aB}p_{Ab} \\
 &= p_{ab}(p_{AB}) - p_{aB}p_{Ab} \\
 &= p_{AB}p_{ab} - p_{Ab}p_{aB}
 \end{aligned}$$

D und Allelfrequenzen

Die Abbildung zeigt den Einfluss von unterschiedlichen Allelfrequenzen auf die Spannweite von D .

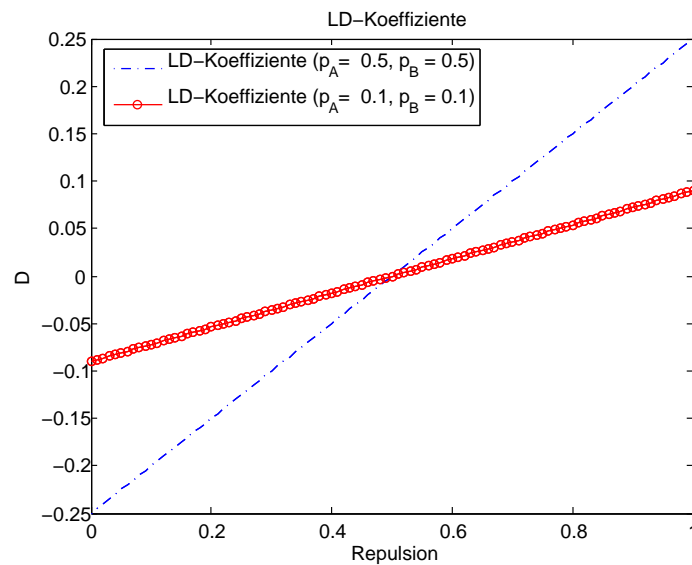
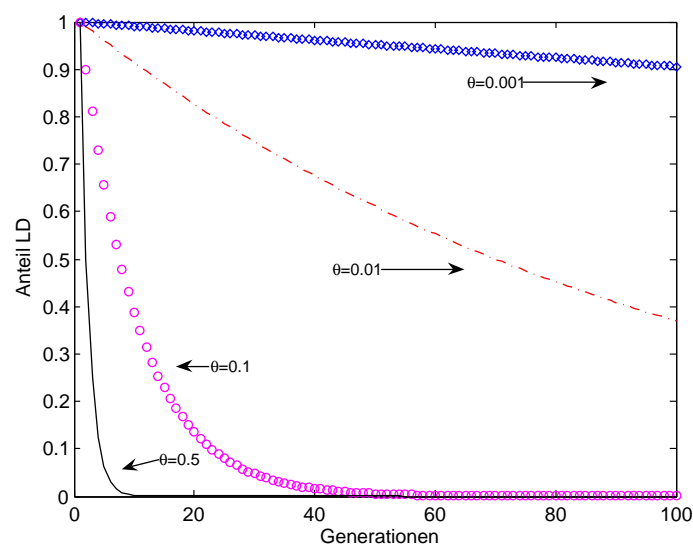


Abbildung: Abhängigkeit der LD-Koeffiziente D von den Allelfrequenzen und vom Grad an Repulsion. 0 = komplette Coupling von AB und ab, 1,0 = komplette Repulsion (Überschuss an Ab und aB).

Wie lange dauert es, bis ein Kopplungsungleichgewicht verschwunden ist?



$$D^i = (1 - \theta)^i D^0$$

```

generationen=[1:100];

Theta=[0.001 0.01 0.1 0.5];
D=zeros(100,4); % 100 Reihen <--> 100 Generationen
                % 4Spalten <--> vier Werte von Theta

D(1,:) = ones(1,4); %Initialisierung

for k=2:100
    D(k,:) = D(k-1,:) .* (1-Theta);
end

```

```

plot(generationen,D(:,1),'bd');
hold on;
plot(generationen,D(:,2),'r-');
plot(generationen,D(:,3),'mo');
plot(generationen,D(:,4),'k-');
xlabel('Generationen');
ylabel('Anteil LD');

```

Beispiel: Normalisierende Selektion

- Als etwas ausführlicheres Beispiel des Einflusses des LD wollen wir die normalisierende Selektion untersuchen.
- Die natürliche Auslese wirkt häufig gegen Individuen an den Extremen des phänotypischen Spektrums und begünstigt Ausprägungen eines phänotypischen Merkmals, die dem Mittelwert des Merkmals in der Bevölkerung nahe sind. Dieses Phänomen ist zunächst Hermon Bumpus 1898 aufgefallen[†]

[†]Bumpus, Hermon C. 1898. Eleventh lecture. The elimination of the unfit as illustrated by the introduced sparrow, *Passer domesticus*. (A fourth contribution to the study of variation.) Biol. Lectures: Woods Hole Marine Biological Laboratory, 209-225.

Beispiel: Normalisierende Selektion



- Nach einem schweren Wintersturm wurden 136 Hausschwalben untersucht, wovon die Hälfte überlebte
- Unter den überlebenden fand sich ein Überschuss an Vögeln mit durchschnittlichen Maßen hinsichtlich Flügellänge, während Vögel mit kurzen oder langen Flügeln öfter als erwartet gestorben waren.

Normalisierende Selektion: Eine Simulation

Phänotyp	8	9	10	11	12
Genotypen	$\frac{ab}{ab}$	$\frac{aB}{ab}$ $\frac{Ab}{ab}$	$\frac{aB}{aB}$ $\frac{aB}{aB}$ $\frac{Ab}{Ab}$ $\frac{Ab}{AB}$ $\frac{ab}{ab}$	$\frac{aB}{AB}$ $\frac{AB}{AB}$ $\frac{Ab}{Ab}$	$\frac{AB}{AB}$
Fitness	0,8	0,9	1,0	0,9	0,8

- a bzw. b: +2 cm
- A bzw. B: +3 cm

Normalisierende Selektion: Eine Simulation

- Anfangs Kopplungsgleichgewicht mit $p(a) = 0,55$, $p(A) = 0,45$, $p(b) = 0,6$ und $p(B) = 0,4$
- Ohne LD gilt $p(ab) = p(a)p(b)$ usw.
- Die diploiden Genotypfrequenzen für erwachsenen Schwalben können aus den haploiden Genotypfrequenzen der Gameten wie folgt berechnet werden[†]:

$$p\left(\frac{ab}{ab}\right) = p(ab)p(ab) \qquad p\left(\frac{aB}{ab}\right) = 2p(aB)p(ab)$$

[†] Der Faktor 2 kommt daher, dass der Gamet aB von der Mutter und der Gamet ab vom Vater kommen kann oder auch umgekehrt.

Normalisierende Selektion: Eine Simulation

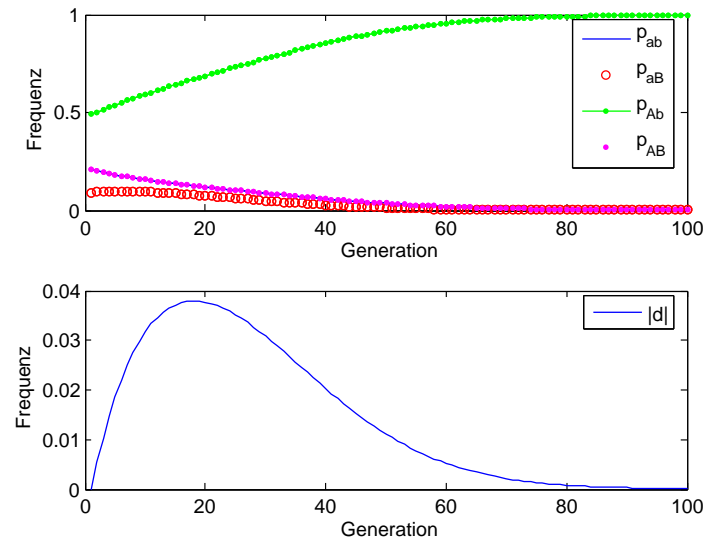


Abbildung: Normalisierende Selektion

- Im folgenden wird der matlab-Code¹ erklärt, womit die Simulation durchgeführt wurde.

¹kann von <http://compbio.charite.de> heruntergeladen werden.

matlab/octave-Code

Unter der Annahme eines Kopplungsgleichgewichts gilt $p(ab) = p(a)p(b)$ usw.

```
p_a=0.3;  
p_b=0.7;  
ngenerations=100;
```

```
p_A=1-p_a;  
p_B=1-p_b;
```

%Am Anfang: Kopplungsgleichgewicht $\rightarrow p(ab)=p(a)p(b)$ usw.

```
p_ab=p_a*p_b;  
p_aB=p_a*p_B;  
p_Ab=p_A*p_b;  
p_AB=p_A*p_B;
```

Der nächste Code-Abschnitt definiert die Vektoren d als 100×1 -Vektor und `genotype_freq` als 100×4 -Matrix. Diese Variablen werden für Generationen $1 \dots 100$ die Werte für die LD-Koeffiziente D und die Frequenzen der vier Genotypen festhalten. Die matlab-Funktion `zeros(M,N)` alloziert Speicher für eine $M \times N$ -Matrix.

```
d= zeros(ngenerations ,1);  
genotype_freq = zeros(ngenerations ,4);
```

matlab/octave-Code

Im folgenden berechnen wir für die erste Generation D :

$$D = p(AB)p(ab) - p(Ab)p(aB)$$

und speichern das Ergebnis im ersten Feld von d . Wir speichern die Genotypfrequenzen der ersten Generation in der ersten Reihe von `genotype_freq`.

```
D=p_AB*p_ab - p_Ab*p_aB;  
d(1)=D;  
genotype_freq(1,:)=[p_ab,p_aB,p_Ab,p_AB];
```

- Ab jetzt simulieren wir eine normalisierende Selektion mit der Funktion `gtype_select`
- Rekombinationsfrequenz $\theta = 0,1$

```
for i=2:ngenerations
    % Calculate and store genotype frequencies
    [p_ab,p_aB,p_Ab,p_AB] = gtype_select(p_ab,p_aB,p_Ab,p_AB);
    genotype_freq(i,:)=[p_ab,p_aB,p_Ab,p_AB];
    %Calculate and store LD
    d(i)=p_AB*p_ab - p_Ab*p_aB;
end
```

matlab/octave-Code

- `gtype_select`

```
function [p_ab,p_aB,p_Ab,p_AB] = ...
    gtype_select(p_ab,p_aB,p_Ab,p_AB)

%Rekombinationsfrequenz 0.1
theta=0.1;

%% Selektion auf Grund des Phaenotyps
%% 8–9–10–11–12 cm Fluegellaenge
fitness_8 = 0.8;
fitness_9 = 0.9;
fitness_10 = 1.0;
fitness_11 = 0.9;
fitness_12 = 0.8;
```


matlab/octave-Code

- gtype_select

```
%% phenotype = 8 cm
p_ab_ab=p_ab^2 * fitness_8 ;

%% phenotype = 9 cm
p_ab_aB = 2*p_ab*p_aB * fitness_9 ;
p_ab_Ab = 2*p_ab*p_Ab * fitness_9 ;

%% phenotype = 10 cm
p_Ab_aB = 2* p_Ab * p_aB * fitness_10 ;
p_AB_ab = 2*p_AB*p_ab * fitness_10 ;
p_Ab_Ab = p_Ab^2 * fitness_10 ;
p_aB_aB = p_aB^2 * fitness_10 ;

... (usw.)
```

matlab/octave-Code

Die Summe der einzelnen Häufigkeiten muss 1 ergeben, weshalb wir renormalisieren müssen:

```
%Renormalize
total = p_ab_ab + p_ab_aB + p_ab_Ab + p_Ab_aB + p_AB_ab + ...
      + p_Ab_Ab + p_aB_aB + p_Ab_AB + p_aB_AB + p_AB_AB ;
p_ab_ab = p_ab_ab / total ;
p_ab_aB = p_ab_aB / total ;
... (usw.)
```

Rekombination & Gameten

Einige, aber nicht alle Rekombinationen führen zu neuen Haplotypen²:

Genotyp $\xrightarrow{\theta}$ Gameten

$$\boxed{AB/AB} \xrightarrow{\theta} \boxed{AB} \& \boxed{AB}$$

$$\boxed{ab/Ab} \xrightarrow{\theta} \boxed{Ab} \& \boxed{ab}$$

$$\boxed{aB/ab} \xrightarrow{\theta} \boxed{ab} \& \boxed{aB}$$

$$\boxed{aB/Ab} \xrightarrow{\theta} \boxed{ab} \& \boxed{AB}$$

$$\boxed{ab/ab} \xrightarrow{\theta} \boxed{ab} \& \boxed{ab}$$

²Bemerke, dass wir der Einfachheit halber die Rekombination so modellieren, dass bei einer Rekombination alle Chromatiden rekombinieren und nicht nur zwei der vier Chromatiden (vgl. Abb. 2.11 von Strachan und Read)

matlab/octave-Code

Wir können nun die Frequenz des Haplotyps ab unter den Gameten berechnen als

$$\begin{aligned} p(ab) &= p\left(\frac{ab}{ab}\right) && \bullet \text{ Jeder Gamet: } ab \\ &+ 0.5 \times p\left(\frac{ab}{aB}\right) && \bullet \text{ Jeder 2. Gamet: } ab \\ &+ 0.5 \times p\left(\frac{ab}{Ab}\right) && \bullet \text{ Jeder 2. Gamet: } ab \\ &+ (1 - \theta) \times 0.5 \times p\left(\frac{AB}{ab}\right) && \bullet \text{ Jeder 2. nicht rek. Gamet: } ab \\ &+ \theta \times 0.5 \times p\left(\frac{Ab}{aB}\right) && \bullet \text{ Jeder 2. rek. Gamet: } ab \end{aligned}$$

Die Berechnungen für die übrigen drei Gametengenotypen erfolgen analog.

```
p_ab = p_ab_ab ...
      + 0.5 * p_ab_aB ...
      + 0.5 * p_ab_Ab ...
      + (1-theta) * 0.5 * p_AB_ab ...
      + theta * 0.5 * p_Ab_aB;

p_aB = 0.5 * p_ab_aB ...
      + 0.5 * (1-theta)* p_Ab_aB ...
      + p_aB_aB ...
      + 0.5*p_aB_AB...
      + 0.5*theta*p_AB_ab;
... (usw.)
```

matlab/octave-Code

- Nach 90 Generationen hat fast jedes Individuum den Genotyp Ab/Ab und somit den günstigsten Phänotyp (Flügelänge 10 cm).
- D steigt anfangs und sinkt mit zunehmender Fixation der Allele A und b .

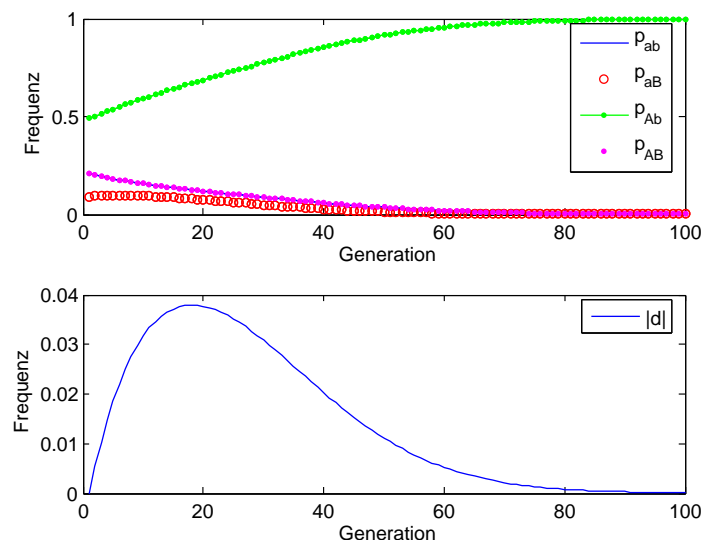


Abbildung: Normalisierende Selektion

Die neutrale Theorie der molekularen Evolution

- Die Neutrale Theorie der molekularen Evolution bzw. die verwandte Idee einer molekularen Uhr wurden in den 1960er–1980er Jahren von Motoo Kimura eingeführt
- Die Evolutionsrate der Aminosäuresequenzen bestimmter Proteine weist über lange evolutionäre Zeiträume eine konstante Rate auf, was sich als Folge der Genetischen Drift erklären lässt.



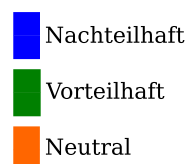
Die neutrale Theorie der molekularen Evolution

- Frühe Darwinistische Theorien gingen davon aus, dass alle Sequenzveränderungen einen Einfluss auf die Fitness haben und somit vorteilhaft oder nachteilhaft sind
- Nach der neutralen Theorie der molekularen Evolution sind die meisten Aminosäurenpositionen neutral, Veränderungen haben keinen wesentlich Einfluss auf die Fitness

Selektion-Theorie



Neutrale Theorie



Die neutrale Theorie der molekularen Evolution

- Vorteilhafte Mutationen: Relativ selten
- Nachteilhafte Mutation dagegen werden durch die natürliche Auslese schnell vom Genpool entfernt
- Ein relativ großer Anteil der denkbaren Veränderungen der Aminosäuresequenz eines Proteins hat keinen wesentlichen Effekt auf die Funktion des Proteins
- Die Anhäufung (Akkumulation) dieser Mutation hängt demnach von der Mutationsrate ab

Die neutrale Theorie der molekularen Evolution

μ^0

Die Mutationsrate μ ergibt sich aus die Rate für nachteilhafte (μ^-), vorteilhafte (μ^+) und neutrale (μ^0) Mutationen. Da vorteilhafte Mutationen selten sind und nachteilhafte durch Selektion schnell aus der Population entfernt werden, fokussieren wir uns auf μ^0 .

Sei N_e die effektive Populationsgröße. Für eine Population einer haploiden Spezies beträgt die Anzahl Mutationen pro Generation $N_e\mu^0$.

Es kann gezeigt werden, dass die Wahrscheinlichkeit, dass eine neutrale Mutation durch genetische Drift fixiert wird, $1/N_e$ beträgt. Daher beträgt die Anzahl von neutralen Mutationen, die Pro Generation fixiert werden $\frac{N_e\mu^0}{N_e} = \mu^0$

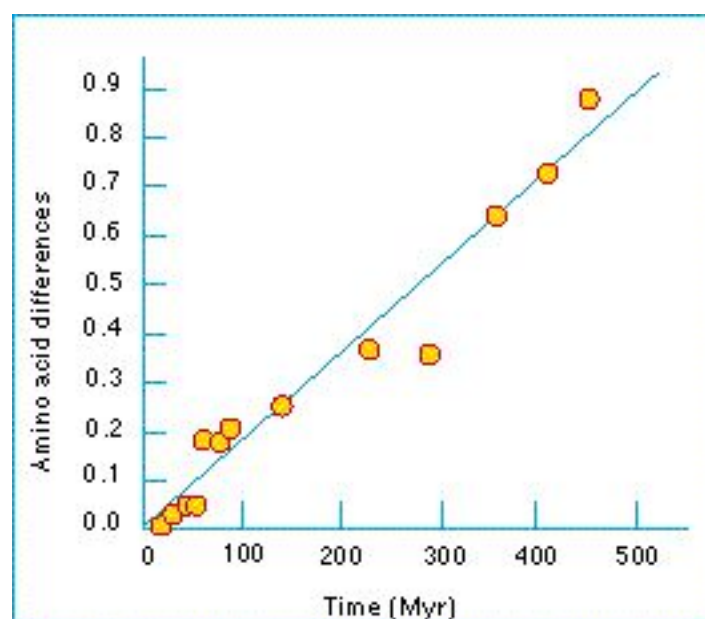
Die neutrale Theorie der molekularen Evolution

- Intuitiv: Obwohl in einer größeren Population mehr Mutationen auftreten, die Wahrscheinlichkeit dass eine spezifische Mutation in der Population fixiert wird sinkt proportional zur Populationsgröße
- Nach dem neutralen Modell bestimmt daher die Mutationsrate μ^0 die molekulare Evolutionsgeschwindigkeit unabhängig von der Populationsgröße
- Dies ist ein wichtiges Ergebnis (Vorhersage):

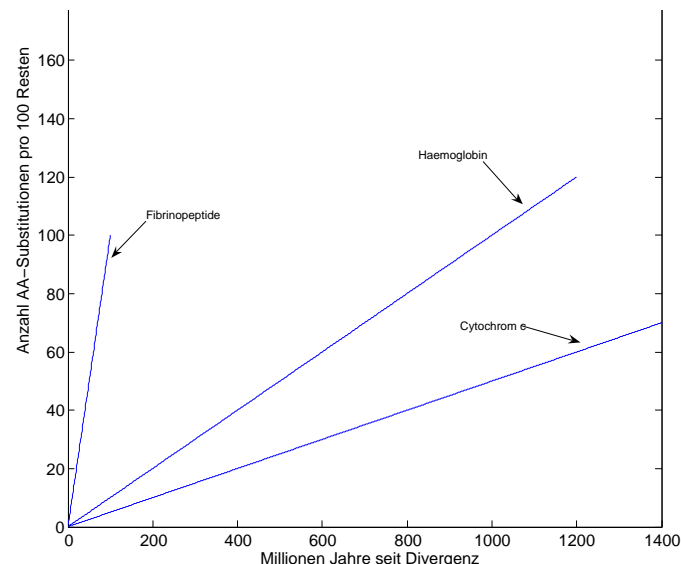
Die molekulare Evolutionsrate in einer Spezies ist dieselbe wie die neutrale Mutationsrate in Individuen^a

^aMerke dass die Mutationsrate und auch die Evolutionsrate sich für unterschiedliche Proteine unterscheiden.

Die evolutionäre Zeit mit der molekularen Uhr messen



Die evolutionäre Zeit mit der molekularen Uhr messen



- Die evolutionäre Rate ist unterschiedlich für unterschiedliche Proteine (unterschiedlicher Anteil an neutralen Resten, andere Faktoren)

Welche sind die wichtigen Positionen in einem multiplen Alignment?

Q5E940 BOVIN	-----MPREDRATWKSNYFLKTIQLDDYPKCFIVGADNVGSKMQQIRMSLRGK-AVVLMGKNTMMRKAIRGHLENN--PALE	76
RLA0 HUMAN	-----MPREDRATWKSNYFLKTIQLDDYPKCFIVGADNVGSKMQQIRMSLRGK-AVVLMGKNTMMRKAIRGHLENN--PALE	76
RLA0 MOUSE	-----MPREDRATWKSNYFLKTIQLDDYPKCFIVGADNVGSKMQQIRMSLRGK-AVVLMGKNTMMRKAIRGHLENN--PALE	76
RLA0 RAT	-----MPREDRATWKSNYFLKTIQLDDYPKCFIVGADNVGSKMQQIRMSLRGK-AVVLMGKNTMMRKAIRGHLENN--PALE	76
RLA0 CHICK	-----MPREDRATWKSNYFLKTIQLDDYPKCFIVGADNVGSKMQQIRMSLRGK-AVVLMGKNTMMRKAIRGHLENN--PALE	76
RLA0 RANSY	-----MPREDRATWKSNYFLKTIQLDDYPKCFIVGADNVGSKMQQIRMSLRGK-AVVLMGKNTMMRKAIRGHLENN--SALE	76
Q7ZUG3 BRARE	-----MPREDRATWKSNYFLKTIQLDDYPKCFIVGADNVGSKMQQIRMSLRGK-AVVLMGKNTMMRKAIRGHLENN--PALE	76
RLA0 ICTPU	-----MPREDRATWKSNYFLKTIQLDDYPKCFIVGADNVGSKMQQIRMSLRGK-AVVLMGKNTMMRKAIRGHLENN--PALE	76
RLA0 DROME	-----MVRENKAANKAQYFIKVVLFDEFPKCFIVGADNVGSKMQQIRMSLRGK-AVVLMGKNTMMRKAIRGHLENN--PALE	76
RLA0 DICDI	-----MSGAG-SKRKKLFIEKATKLFITDKMIVAEADFGSSLOKIRKRSIRGI-GAVLMGKNTMIRKVIDRLADSK--PELD	75
Q54LP0 DICDI	-----MSGAG-SKRKNVFIEKATKLFITDKMIVAEADFGSSLOKIRKRSIRGI-GAVLMGKNTMIRKVIDRLADSK--PELD	75
RLA0 PLAF8	-----MAKLSKQKKQMYIEKLSSLIQQSKILIVHVDNVGSKMASVVRKSLRGK-ATILMGKNTIRITALKKNLQAV--DQIE	76
RLA0 SULAC	-----MIGLAVTTTKKIAKWKVDEVAELTEKLTHTKTIITIANIEGFPADKLHEIRKKLRGK-ADIKVTKNLNFNIALKNAG--VDIK	79
RLA0 SULTO	-----MRIMAVITQERKIAKWKIEVKELEKIREYHTIITIANIEGFPADKLHDIRKKMRGM-AEIKVTKNLFLGIAAKNAG--LDVS	80
RLA0 SULSO	-----MKRLALALKQKRVASWKELEVKELEIKNSNTILIGNIEGFPADKLHEIRKKLRGK-ATIKVTKNLFLKIAAKNAG--LDIE	80
RLA0 AERPE	MSVSVLVGQMYKREKIDPEWKTLMLELEELFSKHVVVLFADLTCTPFVVRVRKRLWKK-VDMVAKKRITLIRAKKAAGLE--LDDN	86
RLA0 PYRAE	MMLAICGKRRYVTRQYDARKVKIVSEATELLQKYDYYVLFDFLHGLSRILHEYRYRLRRY-GVIKIIPKLFKIAFTKVYGG--IPAE	85
RLA0 METAC	MAEERHHEHIDPQWKDEIENIKELIQSHKVFQMVRIEGILATKIQIRDLKDV-AVLKVSNTLTERALNQLG--ETIP	78
RLA0 METMA	MAEERHHEHIDPQWKDEIENIKELIQSHKVFQMVRIEGILATKIQIRDLKDV-AVLKVSNTLTERALNQLG--ETIP	78
RLA0 ARCFU	MAAVRGS--PPEYKVRAVEETKRMISSKPVVAIVSFERNVPAGQMKIRREFRGK-AEIKVVKNTLLEKALDALG--GDIL	75
RLA0 METKA	MAVKAKGQPPSGYEPKVAEWKRREVKELKLMDEYENVGLVLEGTPAPLOEIRAKLRERD-TIRMSRNTLMIRALEKEDER--PELE	88
RLA0 METTH	MAHVAEWKKKEVEELANLKSYPVIALVDVSSMPAYPLSQMRRLIRENGGLLRVSNTLLELAIKKAAGELGKPELE--ENVY	74
RLA0 METTL	MITAESEHKIAPWKIEEYVKLLELLKNGQIVVALVDMMEVPAROLOEIRDKIR-GTMLKMSRNTLLELAIKKAAGELGKPELE--ENVY	82
RLA0 METVA	MIDAKSEHKIAPWKIEEYVKLLELLKNSANVIALIDMEVPAVLOEIRDKIR-DQMLKMSRNTLLELAIKKAAGELGKPELE--ENVY	82
RLA0 METJA	METKVKAHVAPWKIEEYVKLLELLKNSANVIALIDMEVPAVLOEIRDKIR-DKVKLMSRNTLLELAIKKAAGELGKPELE--ENVY	81
RLA0 PYRAB	MAHVAEWKKKEVEELANLKSYPVIALVDVSSMPAYPLSQMRRLIRENGGLLRVSNTLLELAIKKAAGELGKPELE--ENVY	77
RLA0 PYRFO	MAHVAEWKKKEVEELANLKSYPVIALVDVSSMPAYPLSQMRRLIRENGGLLRVSNTLLELAIKKAAGELGKPELE--ENVY	77
RLA0 PYRKO	MAHVAEWKKKEVEELANLKSYPVIALVDVSSMPAYPLSQMRRLIRENGGLLRVSNTLLELAIKKAAGELGKPELE--ENVY	76
RLA0 HALMA	MSAESEKRTETIDPKKQEVDAIVEMIESVGVVNLGTPSRRLQDMRDLHGT-AELKVSNTLLELAIKKAAGELGKPELE--ENVY	79
RLA0 HALVO	MSAESEVQRTETIDPKKQEVDAIVEMIESVGVVNLGTPSRRLQDMRDLHGT-AELKVSNTLLELAIKKAAGELGKPELE--ENVY	79
RLA0 HALSA	MSAESEVQRTETIDPKKQEVDAIVEMIESVGVVNLGTPSRRLQDMRDLHGT-AELKVSNTLLELAIKKAAGELGKPELE--ENVY	79
RLA0 THEAC	MSAESEVQRTETIDPKKQEVDAIVEMIESVGVVNLGTPSRRLQDMRDLHGT-AELKVSNTLLELAIKKAAGELGKPELE--ENVY	79
RLA0 THEVO	MRKINIKKKEIYSELAADITTSKAVAIYDIKVRIRMODIRAKNRDK-VKIKVVKKLLLELAIKKAAGELGKPELE--ENVY	72
RLA0 PICTO	MTEPAQNKIDFVKNLNEHINSRKVAIVSISGLRNNFQKIRNSIRDK-ARIKVSARLLRLALENIGK--NNIV	72
ruler	1.....10.....20.....30.....40.....50.....60.....70.....80.....90	

Welche sind die wichtigen Positionen in einem multiplen Alignment?

”Neodarwinistische Antwort”

Die Positionen, welche Veränderung aufweisen, sind besonders interessant, weil sie uns zeigen, wo die positive Selektion gewirkt hat

Synthetische Theorie der Evolution (1950–1960

G.G: Simpson, 1964)

Der Konsens ist, dass neutrale Gene oder Allele sehr selten sein müssen, falls sie überhaupt existieren. Für einen Evolutionsbiologen erscheint es hochunwahrscheinlich, dass Proteine, die ja durch Gene bestimmt werden, funktionslose Teile haben, dass schlummernde Gene über viele Generationen hinweg existieren oder dass Moleküle sich auf eine regelmäßige aber nicht adaptive Art und Weise verändern sollen. Die natürliche Auslese ist der Komponist des genetischen Codes, und die DNA, RNA und Protein seine Boten

- Nach dieser Ansicht: Unterschiede in Alignments → Folge der positiven Selektion

Welche sind die wichtigen Positionen in einem multiplen Alignment?

Antwort nach der Theorie der neutralen molekularen Evolution

Die Positionen, welche (über ausreichend lange evolutionäre Zeiträume) unverändert geblieben sind, stellen die interessantesten Positionen dar, weil sie uns zeigen, wo die negative Selektion gewirkt hat (d.h., Mutationen in diesen Positionen sind nachteilhaft)

Neutralisten-Selektionisten-Debatte

- Lange Zeit wurde angezweifelt, dass es überhaupt neutrale Mutationen geben kann
- Das erscheint heute klar. Auch Darwin hat die neutrale Theorie der molekularen Evolution vorweggenommen:

Charles Darwin, *On the Origin of Species by Means of Natural Selection*, 6th ed., 1872

Variations neither useful nor injurious would not be affected by natural selection, and would be left either a fluctuating element, as perhaps we see in certain polymorphic species, or would ultimately become fixed...

Synonyme und nichtsynonyme Substitutionen

Synonyme Substitutionen

Nukleotidsubstitutionen in einem Kodon, welche die kodierte Aminosäure nicht verändern. Zum Beispiel CTT=Leucin. CTT→CTA=Leucin, CTT→CTC=Leucin und CTT→CTG=Leucin

Nichtsynonyme Substitutionen

Nukleotidsubstitutionen in einem Kodon, welche die kodierte Aminosäure verändern. Zum Beispiel CTT=Leucin. CTT→ATT=Isoleucin, CTT→GTT=Valin und CTT→TTT=Phenylalanin

Proteine vs. DNA

- Proteine sind die moleküle, die biologische Funktionen erfüllen
- Annahme: Die natürliche Auslese wirkt daher auf Proteinebene und viel weniger auf DNA-Ebene
- Schlussfolgerung: Die Mutationsrate für **synonyme** Substitutionen gibt (ungefähr) die neutrale Mutationsrate an
- Die Mutationsrate für **nichtsynonyme** Substitutionen variiert dagegen je nach Typ und Stärke der natürlichen Auslese

- Sei bei einem Alignment zweier DNA-Sequenzen d_S die Anzahl der synonymen Substitutionen und d_N die Anzahl der nichtsynonymen
- Dann ist $d_N > d_S$ ein Hinweis auf **positive** Selektion
- $d_N < d_S$ ein Hinweis auf **negative** Selektion
- Zahlreiche Methoden sind entwickelt worden, um z.B. solche Methoden auf multiple Alignments anzuwenden, um Hinweise auf Neutralität in bestimmten Codons/Abschnitten zu suchen.

Mäuse und Menschen

- Der letzte gemeinsame Vorfahr von Mäusen und Menschen lebte vor ca. 75 Million Jahren
- Die Genomsequenzen dieser Organismen unterscheiden sich an ca. jedem zweiten Nukleotid
- Weniger als 1% der 25.000 proteinkodierende Gene in der Maus haben kein homologes Gen beim Menschen. Viele Proteinsequenzen zeigen eine Übereinstimmung von über 90%

Urheberrechtlich geschütztes Bild entfernt

X-Chromosom: Syntenieblöcke. Pevzner et al. (2003) *Genome Res.***13**: 37–45.

Bei der Betrachtung eines multiplen Alignments...

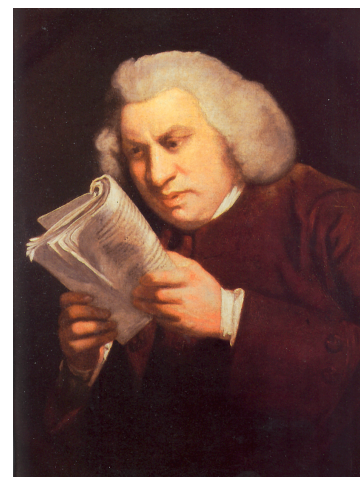
```
Q8CUN7_OCEIH/23-262
Q8ELC4_OCEIH/23-266
Q9K8U7_BACHD/22-288
Q5WEW6_BACSK/22-279
Q8ER73_OCEIH/22-286
CTAG_BACSU/21-280
Q65K08_BACLD/22-281
Q2B4U5_9BACI/1-263
Q5L109_GEOKA/22-284
Q2ASN4_9BACI/33-295
Q41CW2_9BACI/26-285
Q1H1C4_METFK/43-289
Q01RA2_SOLUE/42-282
Q13I14_BURXL/64-313
A0G4F4_9BURK/59-310
Q0ML06_9RHIZ/37-286
Q4IX09_AZ0VI/37-262
A0HY77_PSEME/39-261
Q5P902_AZ0SE/147-370
Q0LFY3_HERAU/14-252
A0I6Y9_9CHLR/14-253
Q67HK8_SYMTH/42-281
Q3J094_BURP1/20-265
Q6XN36_RHOER/28-280
Q70K40_9ACT0/355-608
Q73XL2_MYCPA/390-645
A1TCF0_MYCVP/380-635
Q1B5W0_MYCSS/375-630
Q6ABY8_LEIXX/366-618
```

```
CLLYIICYFFVL.....PVTE..SGSMKKAVLTLGVTLFIAIGSPLNIAR.LTFQGHMIOILLTVVSAPLLVA
VVVIAIIYASSI.....IFLTDV..KVYHROPILFFLSLSLFYIMGSPLATISH.LSFSLHMIOISILYIVPPLLIT
LANVGFVYFHIA.TKWRERFTNSE..PVPTRKKIYFVLGLIALYIGWGSFFVYAGH.LMITFHMAQMVFAFIAVPLFLL
LVGVAFIYSW.....FFRRTN..AGPKRKPLFFFLGLAALYLWGGSPLYVTGH.VHMTLHMLQMVFAFIAVPLLLL
VILLGTAYYLIT.GPLRREFGDND..KPTAKQSMFYIALLLLYFVKGAPIDLISH.ITLTAHMIQMAIYLLVFPMLMIK
LLGITALLYFY.....RRMSSKPN..RIITGKEMVCFLSAMFLYAAEGSPVDLLGH.IMFSAHMQMAVLYLVVPPLLIA
VVFITALYFFL.....KRLGSEGE..RASRKEIGLFLTAMILLYASKGSPVDLLGH.IMFSAHMQMAVLYLVVPPLLIT
MLAILAAYFLLT.VKYRQRFSGST..PLSAKQAVLFTAAILLYVAVKGSPPVDLHSH.ITFYAHMIQMSILYIVPPLLIT
LAAVALLYGIT.GPWRQRFGLGD..AVSEKQKAYFLTGIALLYICKGSPDLMLGH.LTFTAHMVQMAVLYLVPOCFIL
MISILISYFLII.GPYRTREFENAT..KVSKKQIFYFTGIVLLYVVKGGPIDLIGH.IIFSAAHMFHMAVMIYVPPLLL
LTGIYVLYAVLT.....EKIRRPDEA..ETTLGOKFSMLAALFVYYIGFGSPLDVLAH.ITFSAHMLQMVFVYVMAVPLLM
VGLALVLYLA.....GLYRMTKRIGRPTSDAPMRKAFLLAWLTVVALFSPVDLTGN.AYFSMHMVQHELLMIVAAPLFVM
LLLTAVLYF.....RGASRRR..GVSLKOTFFFWAGWSILCLALLSPLHPLGE.ALFSAHMQHELLMVAAPLLVL
MLASTLAYAVGYVRLRLRGSPRSR..ATRAWHASAAGMAALVFALCSPDLSLA.ALFSAHMVQHEHMLIAAPLLVL
MAMSAAYAGGYVRLRARASPRSRARVRAVHLIAFVSGHVALALALFSPDLTSG.ALFSAHMVQHEHMLIAAPLLVA
PLALTALAYAIGHRLWSASARQQ..TIHLQRAVCFAGWFLAALVSPDLRLAT.QLFTAHHIEHEILMVIAAPLFVL
LGAALLAAWIGYEGGCRHRAAR..R...RRALLHGGLLAALSLFGLDEAAE.SSAAAHMAQHMLMLAVAPLLAL
VLLGSALLYII..GCRKVPRHGR..E....ALWMLHAMVITVFAVFGPIDDAE.TSTSLHMHVQHMLMIVIAPLWAL
LIAALLGAACGLYGLGARRVPPGR..V....QATWFCAMAIGALAVFGPLDRAE.NSTALHMHVQHMLIVVAPLAL
LIAATVGYLWAV.GPARKRLGGPA..AFPYKAVAFLSGLLALGLSIMPIGIADRYLFTMHMVQHMLTMFCAPMLLI
LALIAAGVALCVTGPLRRFPFSA..PPTPVQVRLFYGVVVLFIALASPIDSLAS.YLLTMHMLQHLMLAMVAPFLLL
TVLLNAAVLLV..NLWRRAFNWGP..PVPVWRQVLFCLGLWTVYLSGPTPIHISELYLFSVHHVQHTLLTHVMPILL
VLVAGVLF.....RGARKA..KVSASRRVAFWFLVALYVALHTRLDYFFE.HEFFMHRAQHLVHLHGLPFPIAL
LPLIGVLLAAWYCWGVYRVTSAGR..VMPWTRTASFLFGCILLILVTGLAVEGYGY.ELFSIWMFQHLTSLMAIPPLLVL
VAAVVGAVGYV..AARRERASGH..RWPTRRTVSWVVGCAVVVVTSSGLKAYGN.ALFSVHMAEHTALTHIAPVLLVG
FGTAAIVLAGLYVAGVRLRRRGD..RWPPGRGSSWLLGCLVLLFVTSSGVGRYMP.AMFSMHMVQHMLMLAPILLAL
FGTAAIVFALVYLAGVRLRRRGD..AWPIGRVVAWLLGCLVLLATSSGVGRYMP.AMFSMHMVQHMLMLAPILLVL
LGSAAIILALVYLAGVRLRRRGD..AWPAGRTVAWLLGCATLLITSSGLGRYMP.AMFSVHMAHMLMLAPILLVL
LACAFALFFYL..AGVWRLKRGRD..RWPVHRTILWTFGLVLLFFVTSSGVNVYEK.YIFLVHMSAHMVLTHAVPLLLV
```

Wichtigste Voraussetzung: Die Divergenz ist ausreichend hoch, so dass man funktionell bedeutsame Elemente durch ihren hohen Grad an Konservierung erkennen kann

The End of the Lecture as We Know It

- Diese Vorlesungsdiass stehen unter der GNU-Lizenz für freie Dokumentation^b
- Kontakt: peter.robinson@charite.de
- Vorlesungsskript Kapitel 3, Strachan & Read Kapitel 15.4
- Bromham L, Penny D (2003) The modern molecular clock. Nature Reviews Genet 4:216–224.



Lectures were once useful; but now, when all can read, and books are so numerous, lectures are unnecessary. If your attention fails, and you miss a part of a lecture, it is lost; you cannot go back as you do upon a book... People have nowadays got a strange opinion that everything should be taught by lectures. Now, I cannot see that lectures can do as much good as reading the books from which the lectures are taken. I know nothing that can be best taught by lectures, except where experiments are to be shown. You may teach chymistry by lectures. You might teach making shoes by lectures!

Samuel Johnson, quoted in Boswell's Life of Johnson (1791).

^b <http://www.gnu.org/licenses/fdl.txt>