

# Populationsgenetik 3: Kopplungsungleichgewicht (LD)

Peter N. Robinson

Institut für medizinische Genetik  
Charité Universitätsmedizin Berlin

1. Dezember 2014

# Outline

- 1 LD
- 2 LD-Koeffiziente
- 3 Normalisierende Selektion
- 4 Die molekulare Uhr

# Hardy-Weinberg-Gesetz

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

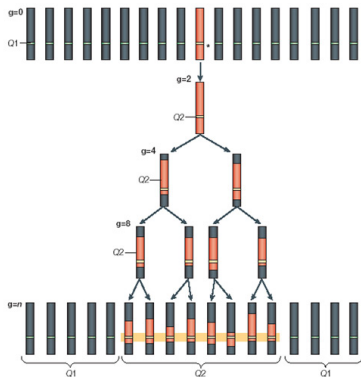
$$p^2 + 2pq + q^2$$

# Kopplungsungleichgewicht

- Dieses Mal ...
- Eigenschaften von Gruppen von Genorten
- Haplotypen
- Kopplungsgleichgewicht
- Kopplungsungleichgewicht ([englisch: Linkage Disequilibrium, LD](#))

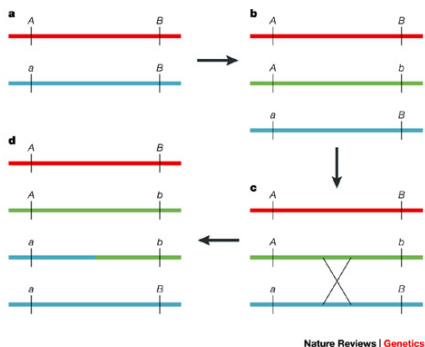
# Eine Geschichte zweier Mutationen

- Heute existierende Haplotypen sind durch weit zurückliegende Mutationsereignisse entstanden



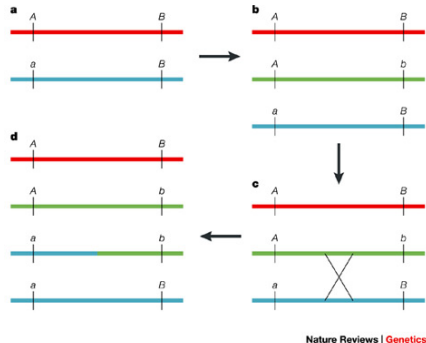
- Das Allel **Q2** entsteht in der Generation 0 durch Mutation von Allele **Q1**.
- Zunächst besteht ein komplettes LD zwischen Q2 und den Allelen aller anderen Loci im ersten Gameten, welcher Q2 trägt
- Dieses LD verringert sich allmählich im Laufe der Generationen auf Grund von **Rekombination**
- LD bleibt für eng benachbarte Loci bestehen
- In Generation  $n$  kann noch ein minimaler Haplotyp um Q2 bei mehreren Individuen identifiziert werden (gelber Balken)

# Eine Geschichte zweier Mutationen (2)



- a Anfangs besteht ein polymorpher Locus mit zwei Allelen, **A** und **a**.

# Eine Geschichte zweier Mutationen (3)

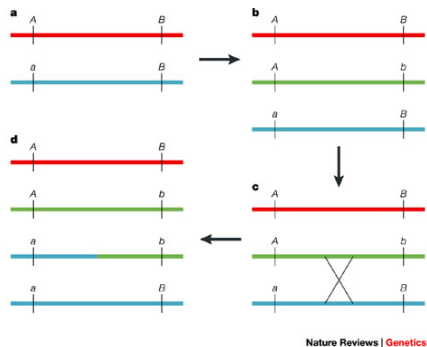


- b** Tritt eine Mutation an einem benachbarten Locus auf (**B** → **b**), betrifft die Mutation ein einzelnes Chromosom mit entweder dem **A**- oder dem **a**-Allel.

Daher werden früh in der Geschichte der Mutation nur drei der möglichen vier Haplotypen in der Population

beobachtet (In diesen Beispiel: **AB**, **aB** und **Ab** aber nicht **ab**).

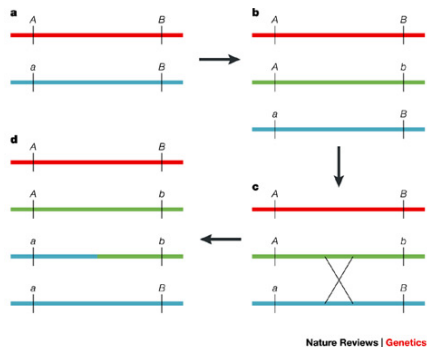
# Eine Geschichte zweier Mutationen (4)



- c Die anfänglich absolute Assoziation zwischen den Allelen an den beiden Loci wird im Laufe der Generationen durch Rekombination zwischen den Loci verringert.



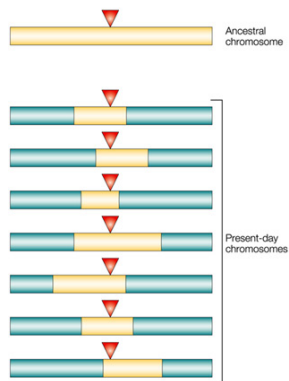
# Eine Geschichte zweier Mutationen (5)



- d Das Ergebnis ist ein vierter Haplotyp und eine Abnahme des LD durch eine Zunahme der Frequenz des rekombinanten Haplotypes **ab** in der Population.

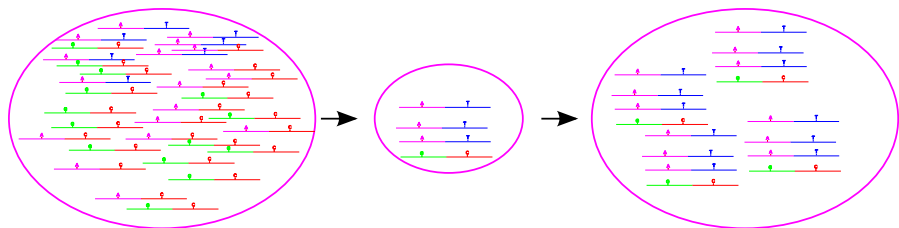
# Kopplungsungleichgewicht

- 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



Nature Reviews | **Genetics**

# 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

# LD und Kartierung (Vorschau)

- 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 zwischen 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

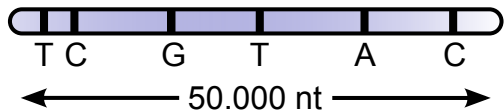
- ▶ "haploider Genotyp"
- ▶ Eine Rekombination wird nur selten zwei Loci trennen, die nahe beieinander auf einem Chromosom liegen
- ▶ Deshalb werden Gruppen von Allelen, die auf demselben Chromosomenabschnitt liegen eher zusammen (als durch Rekombination getrennt) als **Block** übertragen werden

## 2 Genotyp

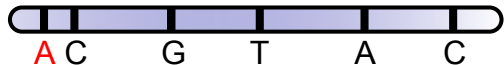
- ▶ die (diploide) genetische Ausstattung eines Individuums an einem oder mehreren Loci. Beide Exemplare eines Allels werden berücksichtigt, z.B. der Genotyp an einem Locus mit Allelen  $A$  und  $a$  kann  $AA$ ,  $Aa$  oder  $aa$  sein.

# Haplotyp

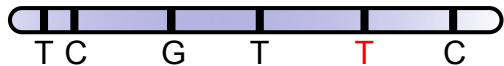
Ursprünglicher  
Haplotyp



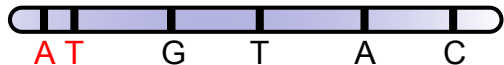
Haplotyp 1



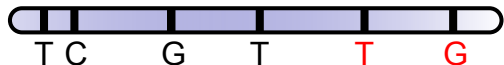
Haplotyp 2



Haplotyp 3

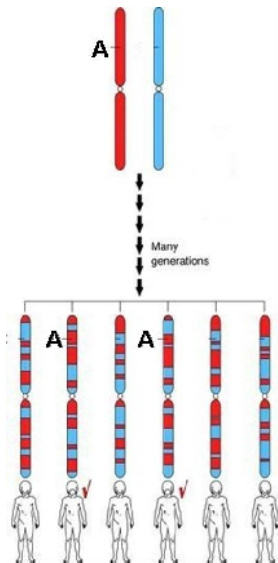


Haplotyp 4



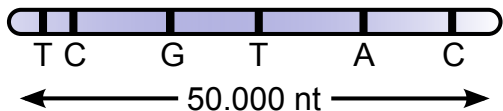
# Haplotyp & Krankheitsassoziation

- Variante A tritt als Mutation auf einem bestimmten Haplotyp auf
- Rekombination über viele Generationen erzeugt unterschiedliche Chromosomen
- Falls die Variante A im ursprünglichen Chromosom das Risiko für eine bestimmte Erkrankung erhöht, haben die Personen, welche dieses Allel geerbt haben, demnach ein erhöhtes Risiko

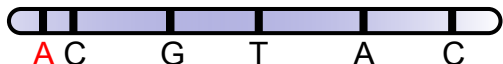


# Haplotyp

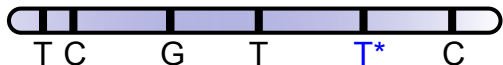
Ursprünglicher  
Haplotyp



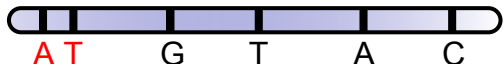
Haplotyp 1



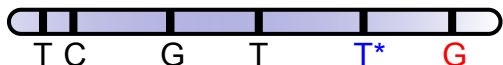
Haplotyp 2



Haplotyp 3



Haplotyp 4



- Mit Erkrankung assoziiert: Allel **T\***
- **T\*** ist an andere SNPs auf diesem Haplotyp (z.B. das T an Position 1) physikalisch gekoppelt



# Kopplungsungleichgewicht

- 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
  - ▶  $D$ ,  $D'$ ,  $r^2$

# Outline

1 LD

2 LD-Koeffiziente

3 Normalisierende Selektion

4 Die molekulare Uhr

- 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<sup>†</sup>:  $p_{AB}$ ,  $p_{Ab}$ ,  $p_{aB}$  und  $p_{ab}$ .

---

<sup>†</sup>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<sup>†</sup>, dann ist die Wahrscheinlichkeit, dass ein Gamet das Allel  $a$  aufweist, unabhängig von der Wahrscheinlichkeit, dass er das Allel  $b$  aufweist

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

---

<sup>†</sup>zum Beispiel weil die Genloci auf unterschiedlichen Chromosomen gelegen sind

# Kopplungsungleichgewicht

- 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:

$$p_{ab} = p_a p_b + D \tag{1}$$

# 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<sup>†</sup>

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

---

<sup>†</sup>s. Skript

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

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

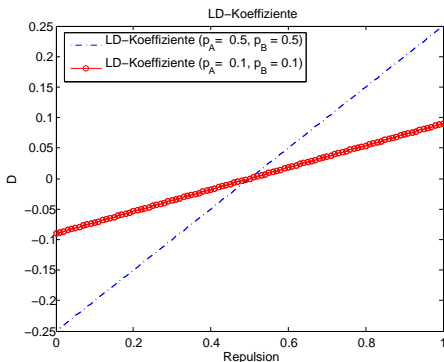


Diese Formel leitet sich von der Definition (1) ab:

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

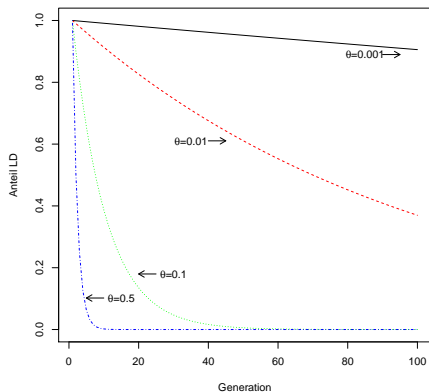
# D und Allelfrequenzen

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 Kopplung von  $AB$  bzw.  $ab$ , 1,0 = komplette Repulsion (Überschuss an  $Ab$  und  $aB$ ).

# Wie lange dauert es, bis ein Kopplungsungleichgewicht verschwunden ist?



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

# R

```
generations <- 1:100  ## 1,2,...,100
theta <- c(0.001,0.01,0.1,0.5)
n <- 100 ## Reihen
m <- 4   ## Spalten
D <- matrix(data=rep(0,n*m),nrow=n,ncol=m)
D[1,] = rep(1,4)  ## erst Reihe mit "1" initialisieren

for (k in 2:n) {
  D[k,] = D[k-1,] * (1-theta) ## elementweise Multiplikation
}
```

- `rep(x, y)` Element "x" y-mal wiederholen
- `matrix(data=rep(0, n*m), nrow=n, ncol=m)`  
initialisiere  $n \times m$  Matrix, setze alle Elemente auf 0.
- `D[1, ]` die erste Reihe von D

```
plot (generations ,D[ ,1] ,type='l' ,lty=1 ,col='black' ,  
      ylim=c(0,1) ,  
      ylab="Anteil LD" ,  
      xlab="Generation")  
lines (generations ,D[ ,2] ,type='l' ,lty=2 ,col='red')  
lines (generations ,D[ ,3] ,type='l' ,lty=3 ,col='green')  
lines (generations ,D[ ,4] ,type='l' ,lty=4 ,col='blue')
```

- S. ?arrows und ?text für Hilfe zur Platzierung von Text und Pfeilen in Plots



# Outline

1 LD

2 LD-Koeffiziente

**3 Normalisierende Selektion**

4 Die molekulare Uhr

# 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 phänotypisches Merkmals, die dem Mittelwert des Merkmals in der Bevölkerung nahe sind. Dieses Phänomen ist zunächst Hermon Bumpus 1898 aufgefallen<sup>†</sup>

---

<sup>†</sup>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 (Länge in cm)	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

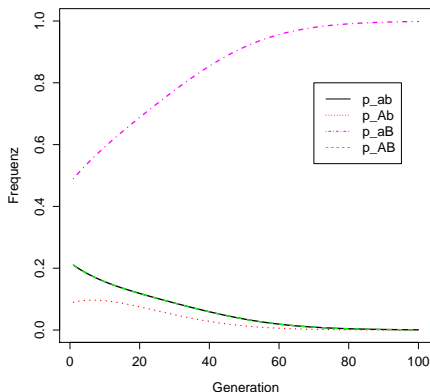


Abbildung: Normalisierende Selektion

- Im folgenden wird der R-Code erklärt, womit die Simulation durchgeführt wurde

# R-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: 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
```

# R-Code

Der nächste Code-Abschnitt definiert die Vektoren  $d$  als  $100 \times 1$ -Vektor und  $freq$  als  $100 \times 4$ -Matrix. Diese Variablen werden für Generationen  $1 \dots 100$  die Werte für die LD-Koeffiziente  $D$  und die Frequenzen der vier Genotypen festhalten.

```
d <- vector(mode="numeric", ngenerations)
freq <- matrix(data=rep(0, ngenerations*4),
               nrow=ngenerations, ncol=4)
```

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  $freq$ .

```
D <- p_AB*p_ab - p_Ab*p_aB
d[1] <- D
freq[1,] <- c(p_ab,p_aB,p_Ab,p_AB)
```

# R-Code

- Ab jetzt simulieren wir eine normalisierende Selektion mit der Funktion `gtypesel`

```
## Ergebnis fuer die uebrigen Generationen berechnen
P <- vector(mode="numeric",4)
for (i in 2:ngenerations) {
  ## Calculate and store genotype frequencies
  #c(p_ab,p_aB,p_Ab,p_AB)
  P <- gtypesel(p_ab,p_aB,p_Ab,p_AB);
  #print(P)
  freq[i,] <- P # c(p_ab,p_aB,p_Ab,p_AB)
  p_ab <- P[1]
  p_aB <- P[2]
  p_Ab <- P[3]
  p_AB <- P[4]
  ## Calculate and store LD
  d[i] <- p_AB*p_ab - p_Ab*p_aB
}
```



# R-Code

- gtypesel

```
gtypesel <- function(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  
  ...  
}
```

# R-Code

- gtypesel

```
# 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.)
```

# R-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
p_ab_Ab <- p_ab_Ab / total
... (usw.)
```

# Rekombination & Gameten

Einige, aber nicht alle Rekombinationen führen zu neuen Haplotypen<sup>1</sup>:

Genotyp  $\xrightarrow{\theta}$  Gameten

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

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

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

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

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

---

<sup>1</sup>Bemerke, 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)

## R-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.

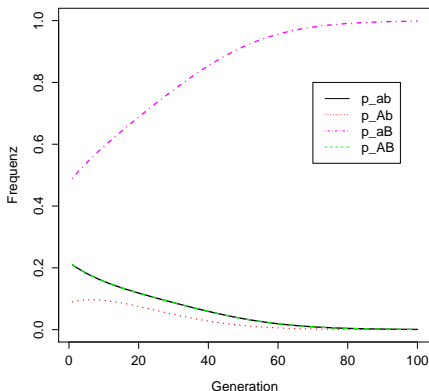
# matlab/octave-Code

```
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 +
... (usw.)
```

# R-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  $b$ .



# Outline

1 LD

2 LD-Koeffiziente

3 Normalisierende Selektion

4 Die molekulare Uhr



# 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 Aminosäurepositionen 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

 $\mu^0$ 

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

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  $\frac{N_e\mu^0}{N_e} = \mu^0$

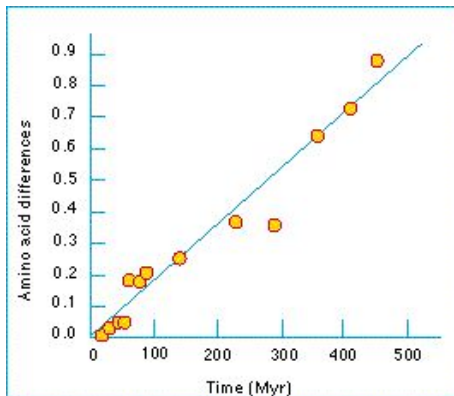
# Die neutrale Theorie der molekularen Evolution

- 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  $\mu^0$  die molekulare Evolutionsgeschwindigkeit unabhängig von der Populationsgröße
- Dies ist ein wichtiges Ergebnis (Vorhersage):

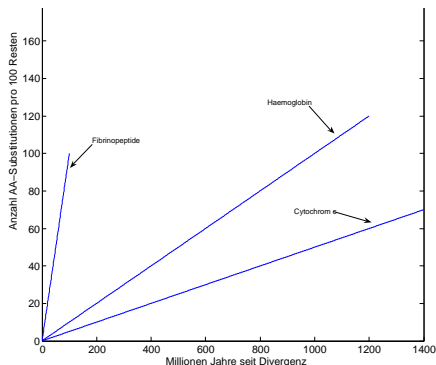
Die molekulare Evolutionsrate in einer Spezies ist dieselbe wie die neutrale Mutationsrate in Individuen<sup>a</sup>

<sup>a</sup>Merke dass die Mutationsrate und auch die Evolutionsrate sich für unterschiedliche Proteine unterscheiden.

# Die evolutionäre Zeit mit der molekularen Uhr messen



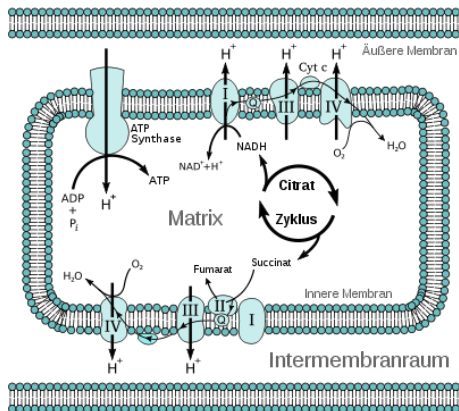
# Die evolutionäre Zeit mit der molekularen Uhr messen



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

# Cytochrome c

	20	30	40	50
Human	TVEKGGKHK	TGPNLHGLFGRK	TGQA	GFSY
Chicken	TVEKGGKHK	TGPNLHGLFGRK	TGQA	GFSY
Cow	TVEKGGKHK	TGPNLHGLFGRK	TGQA	GFSY
Dog	TVEKGGKHK	TGPNLHGLFGRK	TGQA	GFSY
Mouse	TVEKGGKHK	TGPNLHGLFGRK	TGQA	GFSY
Rat	TVEKGGKHK	TGPNLHGLFGRK	TGQA	GFSY
Mosquito	TVEGGKHK	TGPNLHGLFGRK	TGQA	GFSY
Fruitfly	TVEGGKHK	TGPNLHGLFGRK	TGQA	GFSY
Worm	VDS	TAUK	TGPNLHGLFGRK	TGQA
Honeybee	TE	ESGGKHK	TGPNLHGLFGRK	TGQA
consensus	*****	!!!!!!	*****	!!!!!!



- Wesentlicher Bestandteil der Elektronentransportkette
- Extrem hoch konserviert



# Welche sind die wichtigen Positionen in einem multiplen Alignment?

```

Q5E940 BOVIN -----M*PREDRATWKSNSYFLKIIQLDDVVKCFVIGADNVGKKMOQIIRMSLRGK-AVYLMGKNTMMRKAIRGHLENN--PALE 76
RLAO_HUMAN -----M*PREDRATWKSNSYFLKIIQLDDVVKCFVIGADNVGKKMOQIIRMSLRGK-AVYLMGKNTMMRKAIRGHLENN--PALE 76
RLAO_MOUSE -----M*PREDRATWKSNSYFLKIIQLDDVVKCFVIGADNVGKKMOQIIRMSLRGK-AVYLMGKNTMMRKAIRGHLENN--PALE 76
RLAO_RAT -----M*PREDRATWKSNSYFLKIIQLDDVVKCFVIGADNVGKKMOQIIRMSLRGK-AVYLMGKNTMMRKAIRGHLENN--PALE 76
RLAO_CHICK -----M*PREDRATWKSNSYFMKIIQLDDVVKCFVIGADNVGKKMOQIIRMSLRGK-AVYLMGKNTMMRKAIRGHLENN--PALE 76
RLAO_RANSY -----M*PREDRATWKSNSYFLKIIQLDDVVKCFVIGADNVGKKMOQIIRMSLRGK-AVYLMGKNTMMRKAIRGHLENN--SALE 76
Q7ZUG3 BRARE -----M*PREDRATWKSNSYFLKIIQLDDVVKCFVIGADNVGKKMOTIRLSLRGK-AVYLMGKNTMMRKAIRGHLENN--PALE 76
RLAO ICTPU -----M*PREDRATWKSNSYFLKIIQLDDVVKCFVIGADNVGKKMOTIRLSLRGK-AVYLMGKNTMMRKAIRGHLENN--PALE 76
RLAO_DROME -----M*VRENKRAWKQAFYIKVYVLEDFEYKCFVIGADNVGKKMOTIRLSLRGK-AVYLMGKNTMMRKAIRGHLENN--POLB 76
RLAO_DICDI -----MS*AAE-SKRKRLFEKAKYKLFYTKDKMIVAEADYVGSOLQKIRKKSIRGI-GAVLMGKNTMIRKRYIDLADSK--PELD 75
Q54LP0 DICDI -----MS*AAE-SKRKNVFEKAKYKLFYTKDKMIVAEADYVGSOLQKIRKKSIRGI-GAVLMGKNTMIRKRYIDLADSK--PELD 75
RLAO_PLAFB -----M*KLSEKQKQMYTEKLSLITQYYSKLIIVHVDVGNQMSAVKKSIRGK-AITLMGKNTIRITAKKLNLA--PQTE 76
RLAO_SULAC -----M*IGLAVITTKKIKAKWKYDEVAELTKLTKHTIIANIEGFPADKHEIRKSLRGK-ADIVKKNLNFNIALKNAG--YDKE 79
RLAO_SULTO -----M*IRLMAVITQERKIAKWKIEVKLEKLRVHTIIANIEGFPADKHDIRKMRGM-AEIKVKNLFLGIAAKNAG--LDVS 80
RLAO_SULSO -----M*KRILALALKQRKVASWKEVKEKLETLIKNSMTLIGNLEGFPADKHEIRKSLRGK-AIKVYKNTLFLKIAAKNAG--LDIE 80
RLAO_AERPE MS*VVSIVGQMYKREKFEPEKTLMLRELEFLSKIRVYVFDLTGTPPEFVYQVYKSLWKK-VPMVAKKRILLAKKAGGLE--LDDN 86
RLAO_PYRAE M*MAATGCKRYYVYRTQYFARKVKVISEATLQKQPYVYVFDLQGLSQRILBEYVYRIRRYGVKTIKPLFKIAPTKVYGG--LPAP 85
RLAO_METAC -----M*AEERHTEHLPQWKDELEIKELIQSKVYEMVLEGLATKMKCIRRLQDVAVLQVSRNTLEKALMGGLS--ETIP 78
RLAO_METMA -----M*AEERHTEHLPQWKDELEIKELIQSKVYEMVLEGLATKMKCIRRLQDVAVLQVSRNTLEKALMGGLS--ETIP 78
RLAO_ARCFU -----M*AAVRS-----DPEYVYRAVEIKRMISKYPVAIVSRNVPAEQKIRBEFIRGK-AEIKVYKNTLLEKALDGLS--DGLV 75
RLAO_METKA MA*VKAKEQPSSEYEPKVAEKWREKELKLMDEYENGVLDLEIPAPLOEIRAKLERDGLIIRSRNTLMRLIAEKEDDER--PEL 88
RLAO_METH -----M*HAVVAEKWKEVEQLHDLIKVEYVGIANLADAPARLOKMRQTLRDS-ALIRMKKGLISLALIEKAGREL--FNVD 74
RLAO_METL M*ITAESEHKIAPWKEEYENKLELLKNQOIVALVOMMEVPARLOEIRDKIR-GLMKMARTLEKAIKVAEYETGPEFA 82
RLAO_METVA M*IDAKSEHKIAPWKEEYENKLELLKSANVIALIDMMEVPARLOEIRDKIR-DQMLKMRRTLEKIKVAEYETGPEFA 82
RLAO_METJA METKVAHVAPWKEEYVTKGLIKSPVVAIVDMMVPAPELOEIRDKIR-DKVKLIRSRNTLMRLIAEKEDDER--PEL 81
RLAO_PYRAB -----M*HAVVAEKWKEVEELANLIKSPYVIALVDFSSHPAYPSQMRRLIRENGGLLQVSRNTLEKAIKVAEYETGPELE 77
RLAO_PTHO -----M*HAVVAEKWKEVEELAKLIKSPYVIALVDFSSHPAYPSQMRRLIRENGGLLQVSRNTLEKAIKVAEYETGPELE 77
RLAO_PYHFO -----M*HAVVAEKWKEVEELANLIKSPYVIALVDFSSHPAYPSQMRRLIRENGGLLQVSRNTLEKAIKVAEYETGPELE 77
RLAO_PYKO -----M*HAVVAEKWKEVEELANLIKSPYVIALVDFSSHPAYPSQMRRLIRENGGLLQVSRNTLEKAIKVAEYETGPELE 76
RLAO_HALMA MS*SEBEKTEIDPKWQVEADYVMISESQVYNIAGIDRLODMRRDLHGT-ALRYSRNTLEKALDDYD--DGL 79
RLAO_HALVO MS*SEVEQTEVFPQWKRVEVDLYDIESYESQVYVAGIDRLODMRRDLHGT-ALRYSRNTLEKALDDYD--DGL 79
RLAO_HALSA MS*EEQTTEVFPQWKRVEVDLYDIESYESQVYVNTGIDRLODMRRDLHGT-ALRYSRNTLEKALDDYD--DGL 79
RLAO_THEAC -----M*KEVSQKKELVNETHIRIKASRSVAIVDLAGIR-ROTODIRGKNRKG-INKVTKKLLFKALDNLGD--EKLS 72
RLAO_THEVO -----M*RIIDPKKIEIVSEADITKSKVAIVDTRCVRRQMDIRAKNRDK-VKIVYKLLFKALDIND--EKLT 72
RLAO_PICTO M*EPQKIDFVKNLENEINSRKVAIVSIRGLRNNFEPKIRMSIRDK-ARIKVRGALLRLAIEEYK--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)

# Neutralist-Selektionist-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

# Proteine vs. DNA

- 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 heute an ca. jedem zweiten Nukleotid
- Weniger als 1% der  $\sim 22.000$  proteinkodierende Gene in der Maus haben kein homologes Gen beim Menschen. Viele Proteinsequenzen zeigen eine Übereinstimmung von über 90%



# Bei der Betrachtung eines multiplen Alignments...

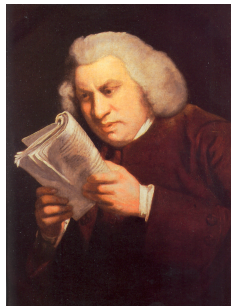
```
08CUN7_OCEIH/23-262
08ELC4_OCEIH/23-266
09K8LU_BACHD/22-288
05WEW6_BACSK/22-279
08ER73_OCEIH/22-286
CTAG_BACSU/21-280
065K08_BACLD/22-281
02B4U5_9BACI/1-263
05L109_GEOKA/22-284
02ASN4_9BACI/33-295
041CW2_9BACI/26-285
01H1C4_METFK/43-289
001RA2_50LUE/42-282
0131J4_BURXL/64-313
A064F4_9BURK/59-310
00ML06_9RHIZ/37-286
04IX09_AZOVI/37-282
A0HY77_PSEME/39-261
05P902_AZOSE/147-370
00LFY3_HERAU/14-252
A0IGY9_9CHLR/14-253
067MK8_5YTH/42-281
03J094_BURP1/20-265
06XH36_RHOER/28-280
07K40_9ACT0/355-608
073XL2_MYCPA/390-645
A1TCF0_MYCVP/380-635
01B5W0_LGSCS/373-630
06ABY8_LEIXX/366-618
```

```
CLLYIICYFFVL.....PVTE..SGSMKKA VLTILGVTILFIAIGSPINIIAR.LTF0GHMIOILLTVVSA P L L V A
VVVIAI IYASSI.....IFLTDV..KVYHROPILFFLSLIFYVINGSP L ATISH.LSFSLHMIQMSILYFIVP L L L I T
LAHVGFVYFMIA.TKWRERFTNSE..PVPIRKKZYFVLGLIAYLWGSPPYVAGH.LMITFHMAQMVFAFYIAPV L L L L
LVGVAFIYSH.....FFRRTN..AGPKRKLFFFLGLAALYLWGSPLVYTGH.VMNTLHMLQMVFAFYIAPV L L L L
YILLGTAYLIT..BPLREFQND..KPTAKQOSHFYIALLLWYFKGAPD L L L S H . I T L T A H M I Q M A I L L V L V P L L M I K
LIGITALYFYI.....RRMSSKPN..RITGKEHWCFLSAHLFLYAEQSPV D L L G H . I M F S A H M I Q M A V L L V L V P P L L I A
VYFITALYFLH.....KRLGSEGE..RASRKEIGLFLTAHILLYASKQSPV D L L G H . I M F S A H M I Q M A V L L V L V P P L L I I L
MLAAILAYFLLT.VKYRQRFSGST..PLAKOAWLFTAAILIYVAVKQSPV D L L M S H . I T F Y A H M I Q M S I L Y L V I P P L L I T
LAAVALLYIGIT..GPHRQRFGLGD..AVSEKQAKYFLTGIALLYICKQSPV D L L M S H . L T F T A H M I Q M A V L L V L V P P C F I L
HISILISYFLIT..GPYRTRFENAT..KVSKKQIFPYTGTGIVLLVYVKGQPD L I G H . I I F S A H M F E M A H V Y I A V P P L L L I
LIGIYVLYAVLT.....EKIRRPDEA..ETLQKQFSHLAALFVYIGFQSPV D L V L A H . I T F S A H M L Q M V F Y H V A P P L L M I
YGLALVLYLA.....GLYRHTRKIGRPTSDAPHRAKAFLLAHLTVVALLFSPV D L L G N . A Y F S A H M H V Q H E L L M V I A A P L F V M
L L L T A V L Y F . . . . . R G A S R Q R . . . . . G V S L K O T F F F W A G S I L C L A L L S P L H P L G E . A L F S A H M A Q H E L L M V A A P L L V L
M L A S T L A Y A V G Y V R L R L R G S P R S R . . . . . A T R A M H A S A F A A G M A A L V F A L C S P L D L S L A . A L F S A H M V Q H E T H M L I A A P L L V L
M A M S A A A Y A A G Y W R L R A R A S P R S R S A R V R A V H L I A F V S G W V A L A L A L F S P L D L S G . A L F S A H M V Q H E S H M L I A A P L L V A
P L A L T A L A Y A I G H R R L W S A S A R Q Q . . . . . T I H L O R A V C F A A G W L F A A A L V S P L D R L A T . Q L F T A H M I E H E I L M V I A A P L F V L
L G A L L A A A W I G Y E G G C R H R P A A R . . . . . R R A L L H G G L L L A A L S L F G P L D E A E . S S A A H M A Q H M L L M A V A P L L A L
V L L G S A W L Y I T . . . . . G C R K V R P H G R . . . . . A L W M H L A H V I T V F A V F G P I D O W A E . N T S L H H V Q H L M F N I V I A P L W L
L I A A L L G A A C G L Y G L G A R R V P P G R . . . . . Q A T W F C A A M A I G A L A V F G P L D R W A E . T S T L H H V Q H L I V V V A P L G A L
L I A A T V G Y L W A V . G P A R K L G G P A . . . . . A F P V G A V A F L S G L L A L G L S I M S P T G I W A D R Y L F T H H M H V Q H M I L T M F C A P H I L L
L A L I A G Y L C V T G P L R R F F P G S A . . . . . P P T P V Q V R L F Y V G W V L F I A L A S P I D S L A S . Y L L T H M L Q H L L A M V A P P F L L L
T Y L L N A A Y L L L . N L W R A F N W G P . . . . . P V P W R O V L F C L G L W T V Y L S E G T P I H I M S E L Y F S V H M V Q H T L T W M V P P L I L L
Y L V A G V L F A . . . . . R G A R K A . . . . . K V S A S R R V A F W F G L V A L Y Y A L H T R L D Y F F E . H E F F H R A Q H L L H H L G P F F I A L
L P I G V L L A A W Y C G Y V R Y T A S G R . . . . . V M P T R T A S F L G C I L L I V T G L A V E Y B Y . E L F S Z W F Q H L T S M A I P P L L V L
Y A A V A G G Y V V . . . . . A A R R E R A S G H . R W P I R R T V S W V G C A V V V S T S S G L K A Y G N . A F S V H H V A H M L S H L V I P L L V G
F G T A A I V L A G L Y V A G V V R L R R R G D . . . . . R W P P G R T S S M L L G C L V L L F V T S S G V G R Y M P . A M F S H H H V H M C L S H L V P I L L A L
F G T A A I V F A L Y L A G V R L R R R G D . . . . . A M P I G R V W A M L G C L V L L A T S S G I R Y M P . A M F S H H H I A H M L S H L A P I L L V L
L G S A A I I A L V Y L A G W M L R R R G D . . . . . A M P A G R T V A M L G C A L L I T T S S G L R Y M P . A M F S V H H V A H M L S H L V I P L L V L
L A C A F A L F F Y L . . . . . A G V W R L R K R G D . . . . . R W P V H R I L T F G Z V L L F F V T S G G V N Y E K . Y I F L V H M S A H M V L T M A V P L L V P
```

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

- Kontakt:  
peter.robinson@charite.de
- Vorlesungsskript Kapitel 4,  
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).