

Let's Talk About

---

# BEING GREEDY



# **Algorithmen in der Bioinformatik**

## **2. Genom-Umordnung - Gierige Algorithmen -**

**Prof. Dr. Gunnar Klau**

**Nach Jones & Pevzner: An Introduction to Bioinformatics Algorithms, Ch.**

# Hinweise zum Urheberrecht



Bitte verzichten Sie in dieser Lehrveranstaltung  
auf Bild- und Tonaufzeichnungen.

Weitere Informationen finden Sie im Studierendenportal

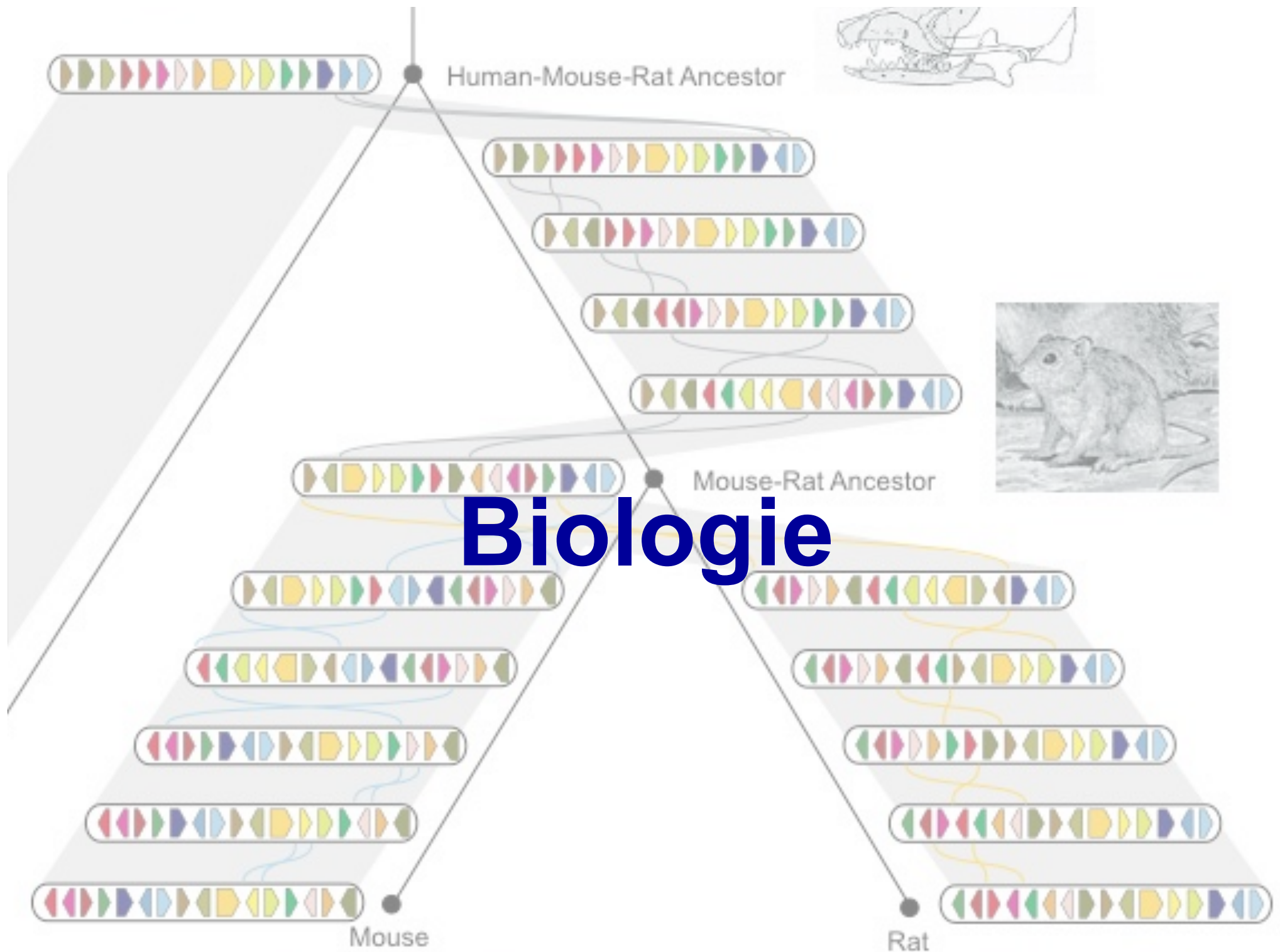
# Hinweise zum Urheberrecht



Das Veröffentlichen oder Teilen von Bild- und Tonaufzeichnungen dieser Lehrveranstaltung ist nicht gestattet.

Weitere Informationen finden Sie im Studierendenportal





# Weißkohl

vs.

# Kohlrabi



*Brassica oleracea var. capitata f. alba*

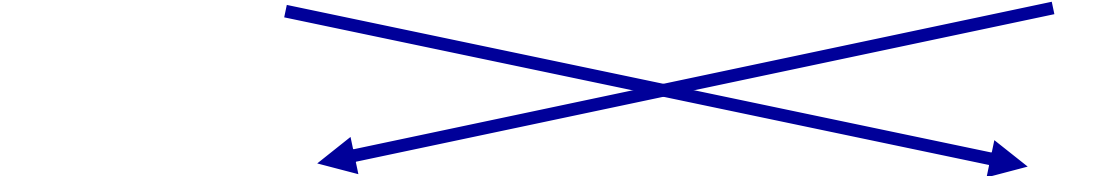


*Brassica oleracea gongylodes*

- Gemeinsamer Vorfahre
- Unterschiedliches Aussehen und Geschmack
- Mitochondriale DNA ist zu 99% gleich
- Unterschied: Gen-Anordnung!

# mtDNA: Gen-Ordnung

**Vorher**

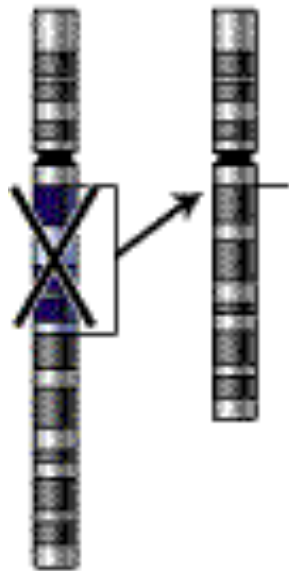


**Nachher**

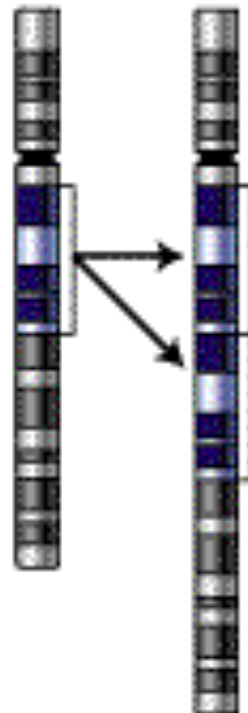


# Genom-Umordnung: Mutationen

Deletion



Duplication



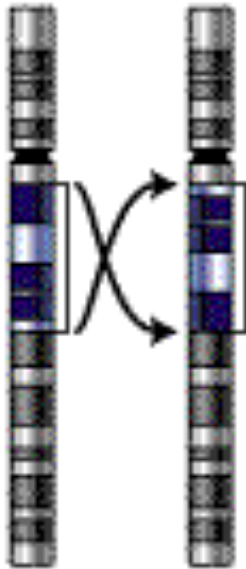
Insertion



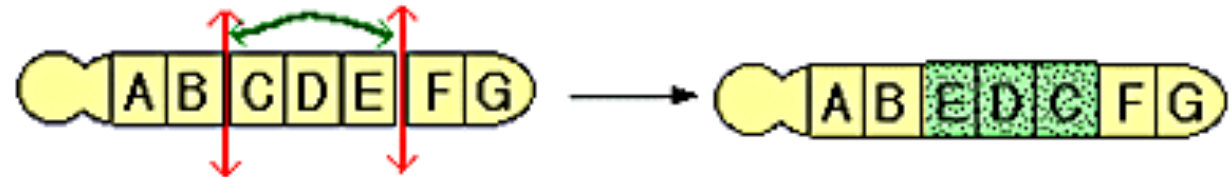


# Genom-Umordnung: Mutationen

Inversion



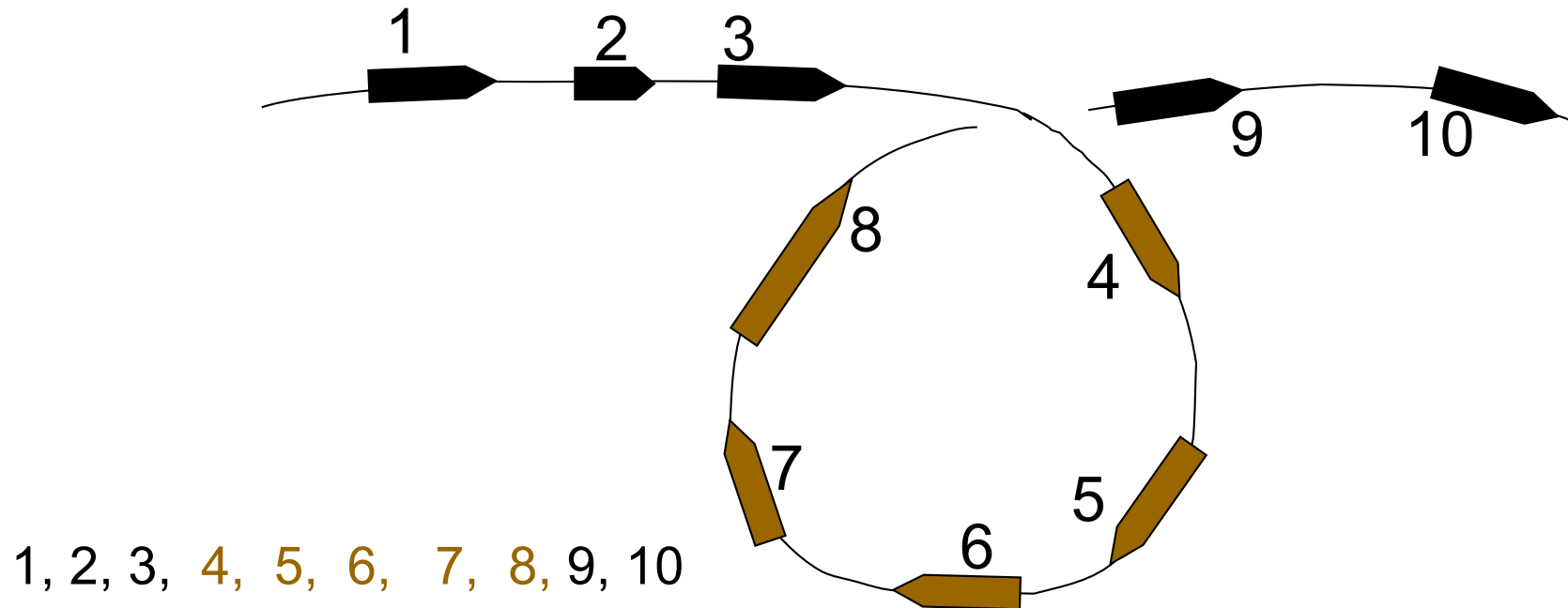
Inversion



1 2 3 4 5 6 7 → 1 2 -5 -4 -3 6 7

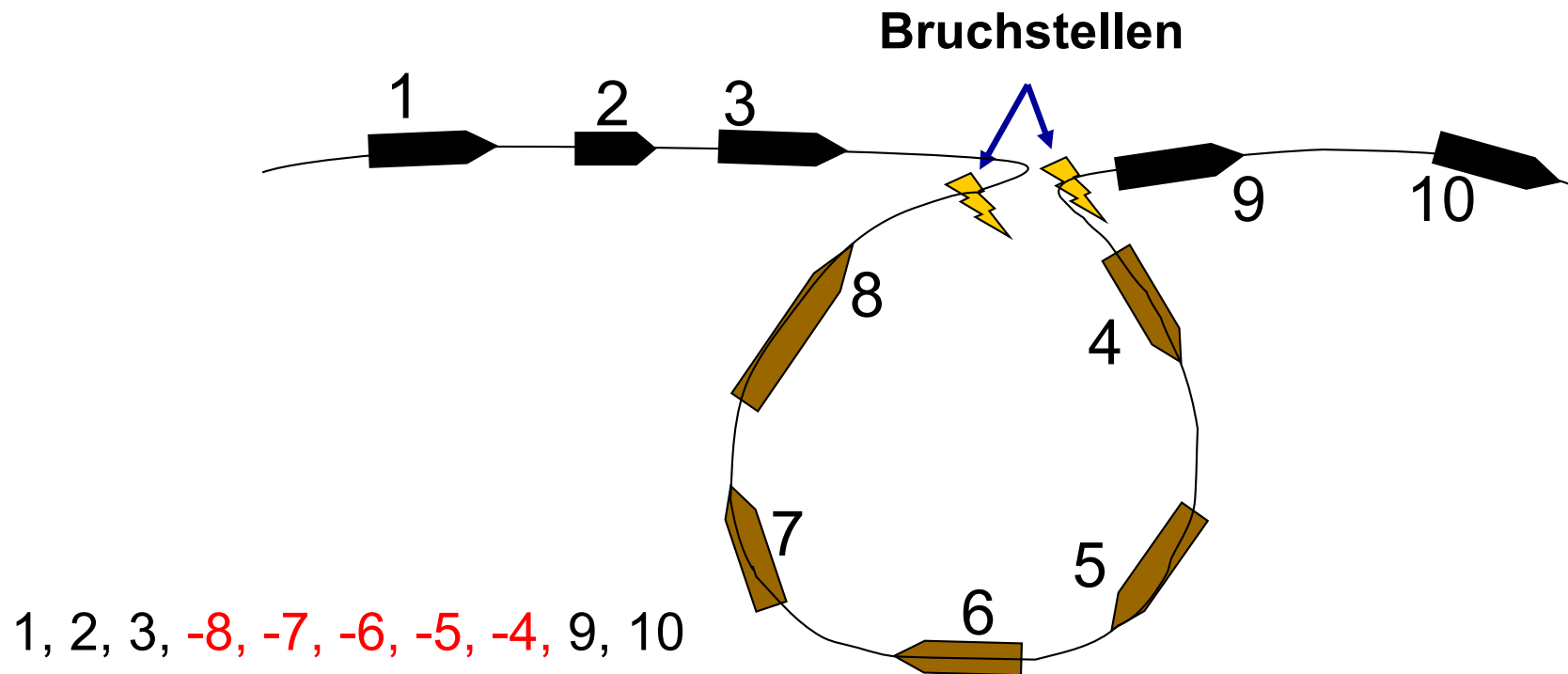
**Inversionen (Umkehrungen) sind die häufigste Ordnungs-Mutation!**

# Umkehrungen sind falsch aufgelöste Knoten



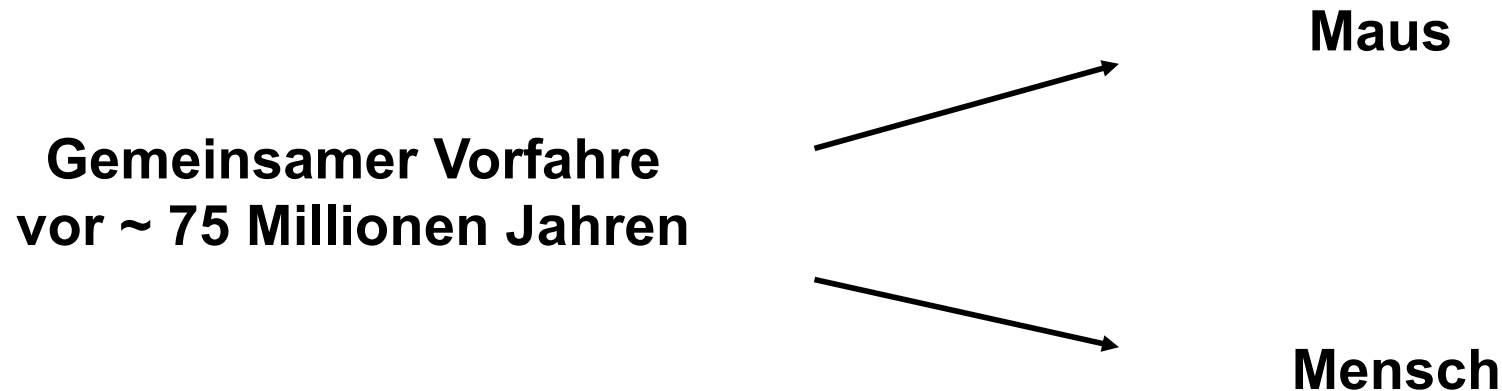
[Blöcke sind Gene]

# Umkehrungen sind falsch aufgelöste Knoten



[Blöcke sind Gene]

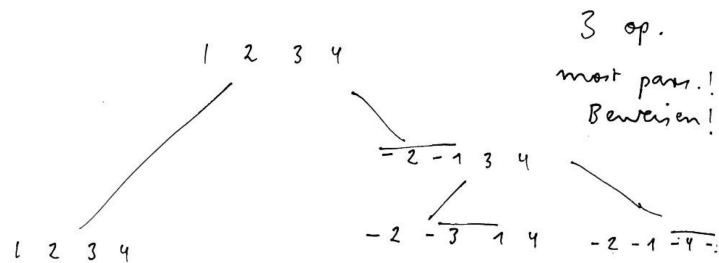
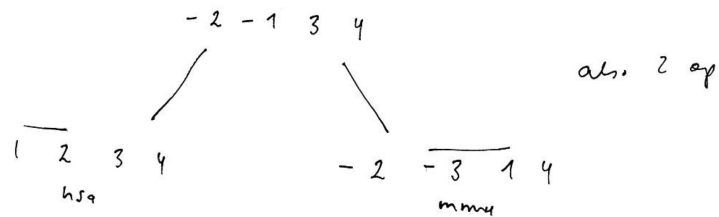
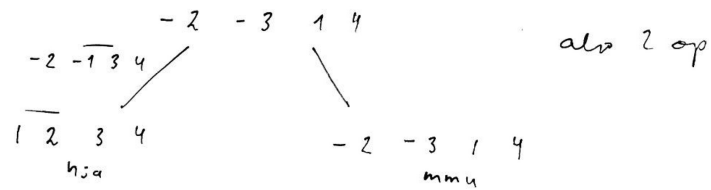
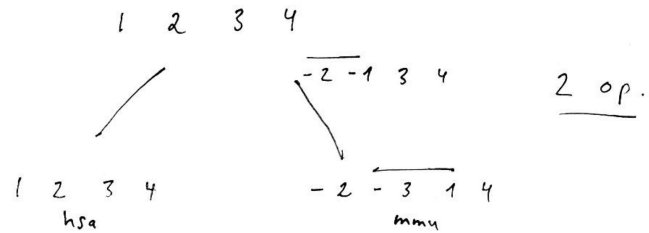
# Mensch-Maus: X-Chromosom



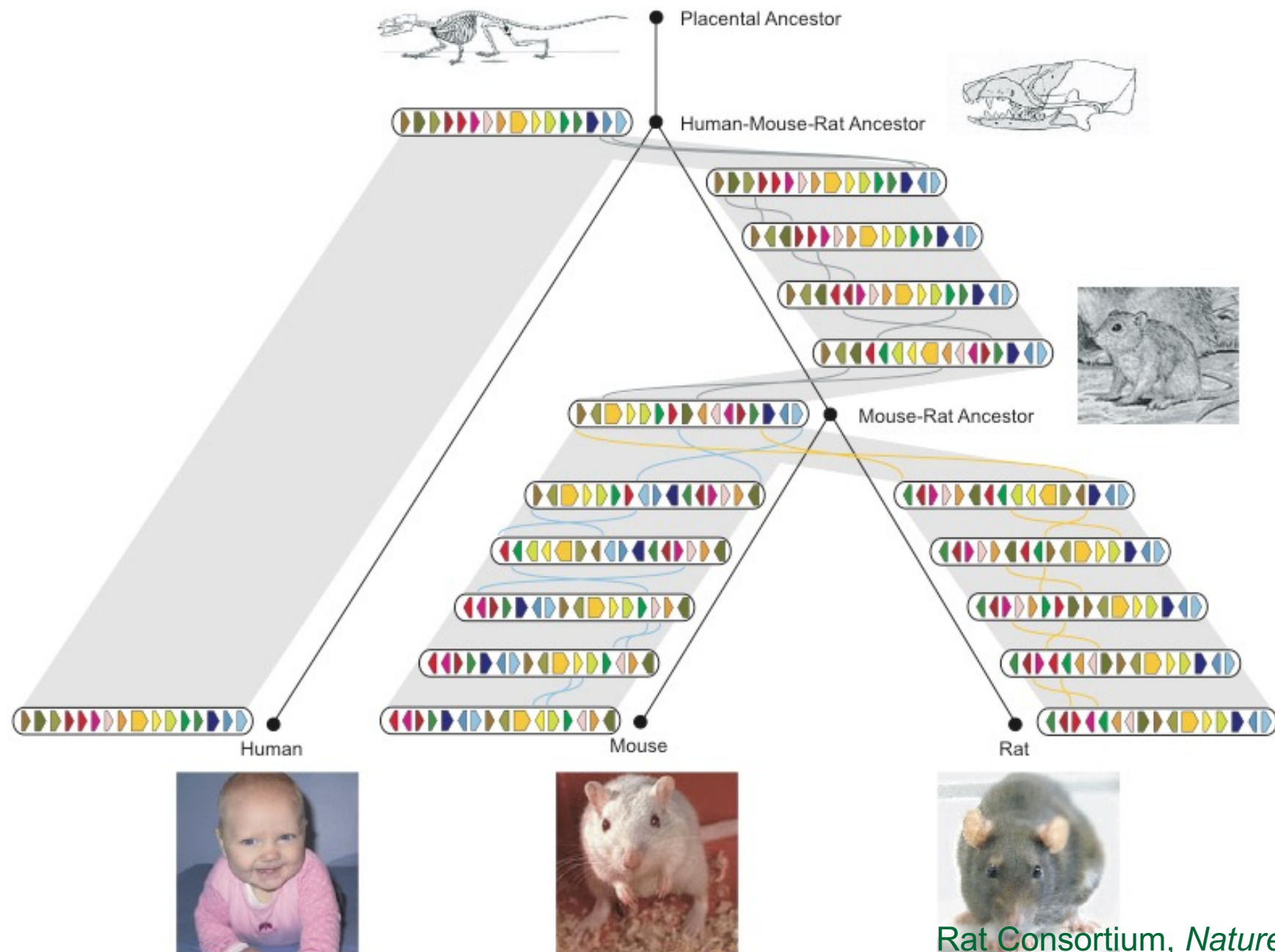
## Fragen:

- Wie identifizieren wir Blöcke, die ähnlich (unverändert) geblieben sind?  
→ ☐ schnelles Alignment (siehe später im Kurs)
- Wie sah das Genom des Vorfahren aus?  
→ ☐ nur durch Vergleich mit weiteren Genomen zu beantworten
- Welche Schritte haben zu den unterschiedlichen Anordnungen geführt?



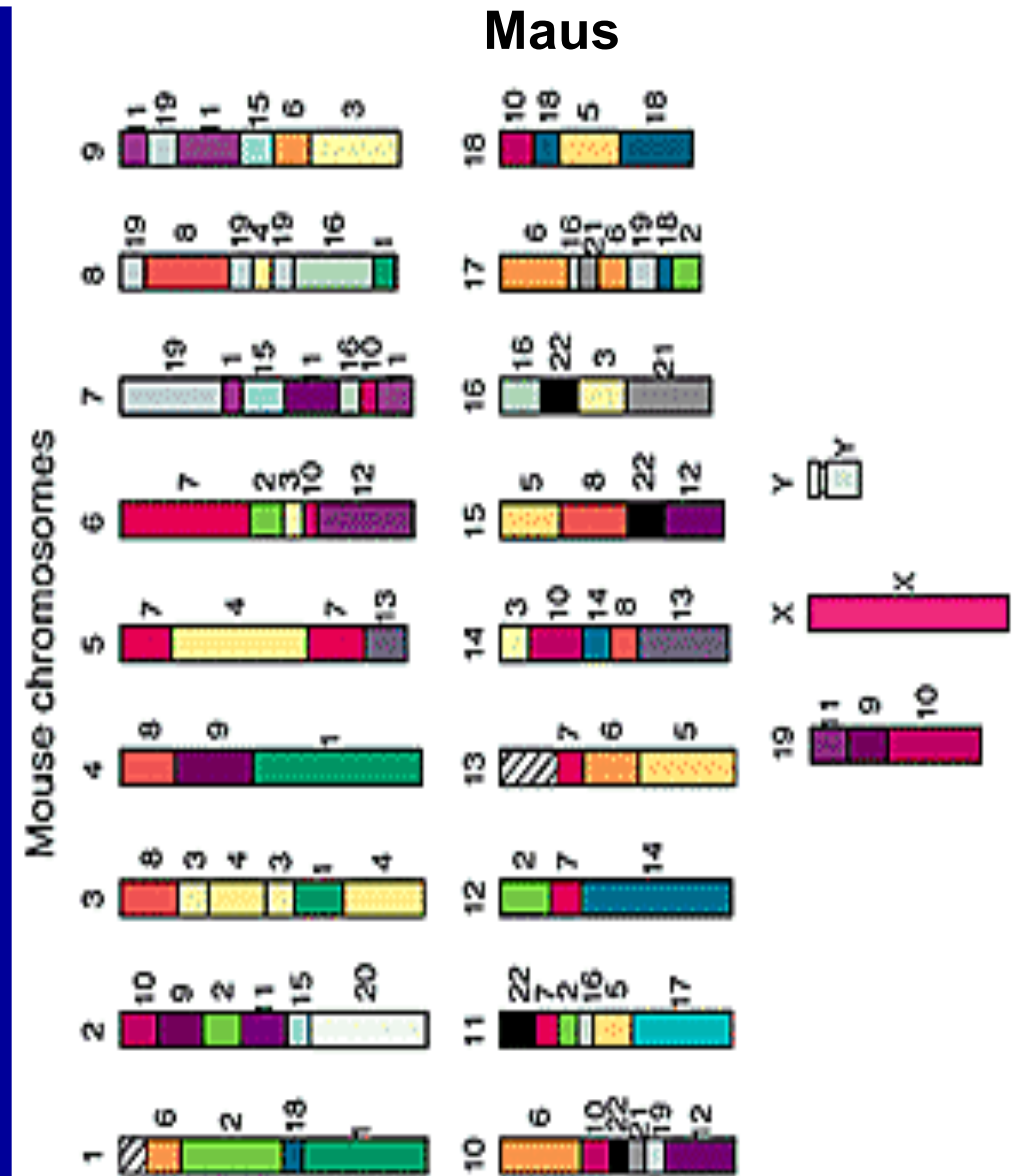
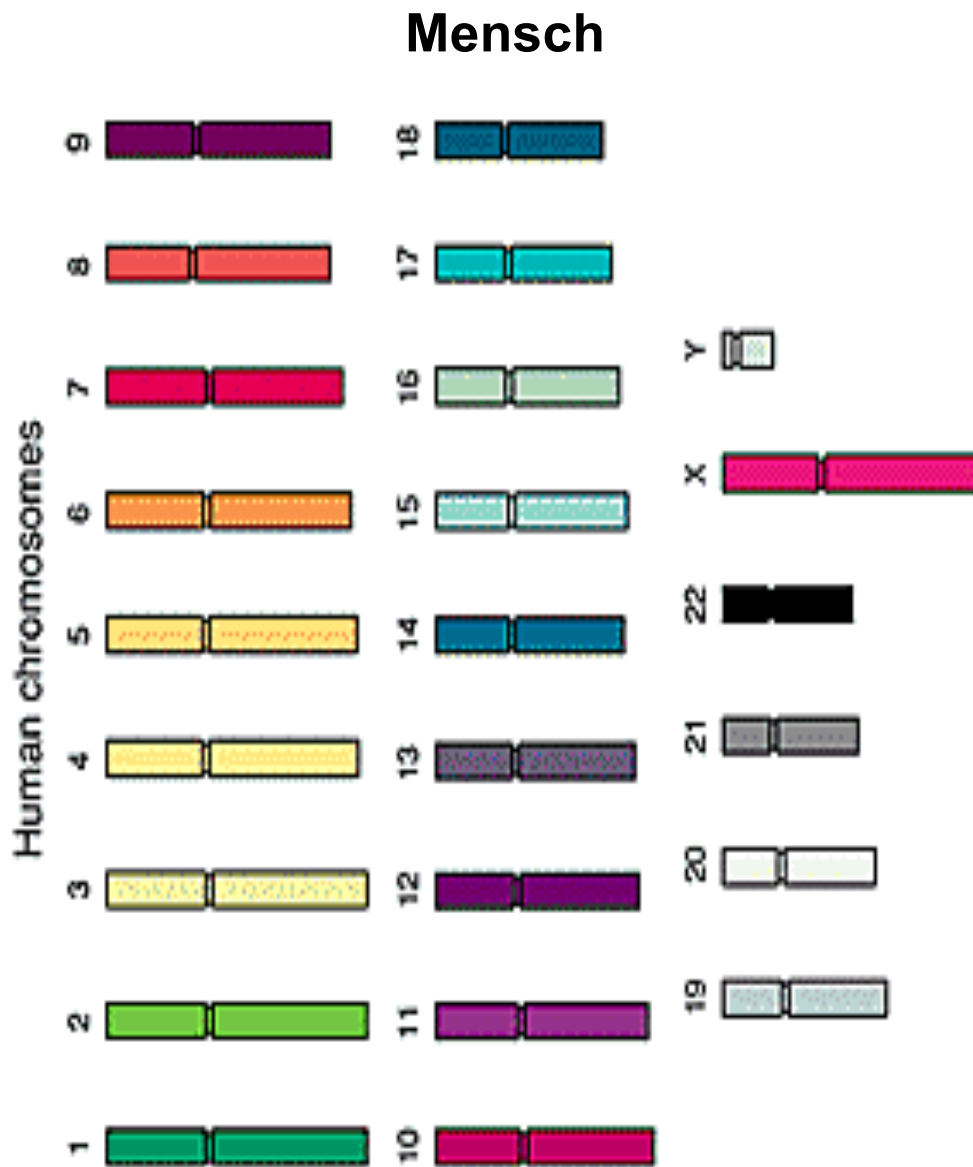


# Die Geschichte des X-Chromsoms



# Mensch-Maus: Genom-Vergleich

~245 Umordnungs-Mutationen:

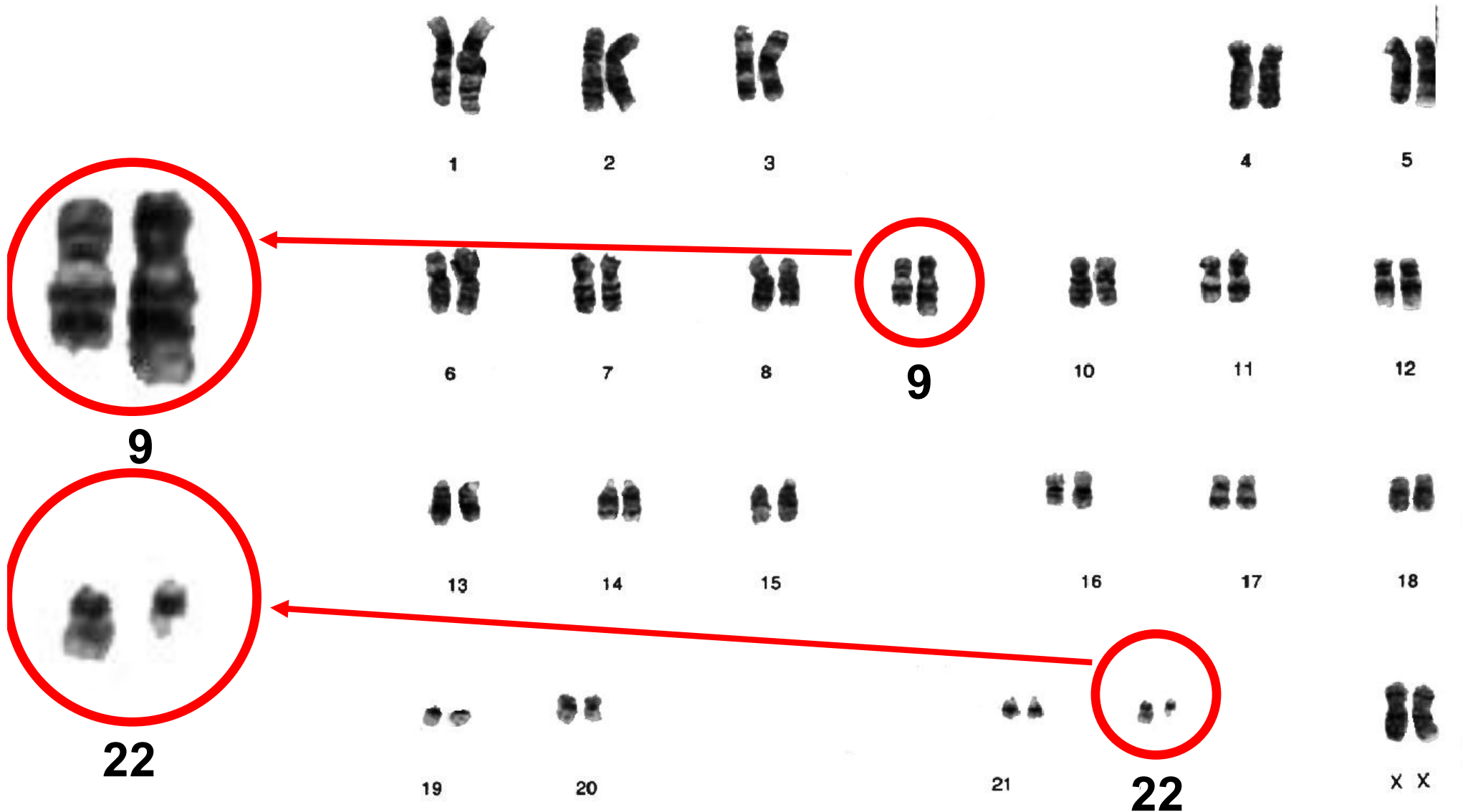


**Wen interessiert das?**

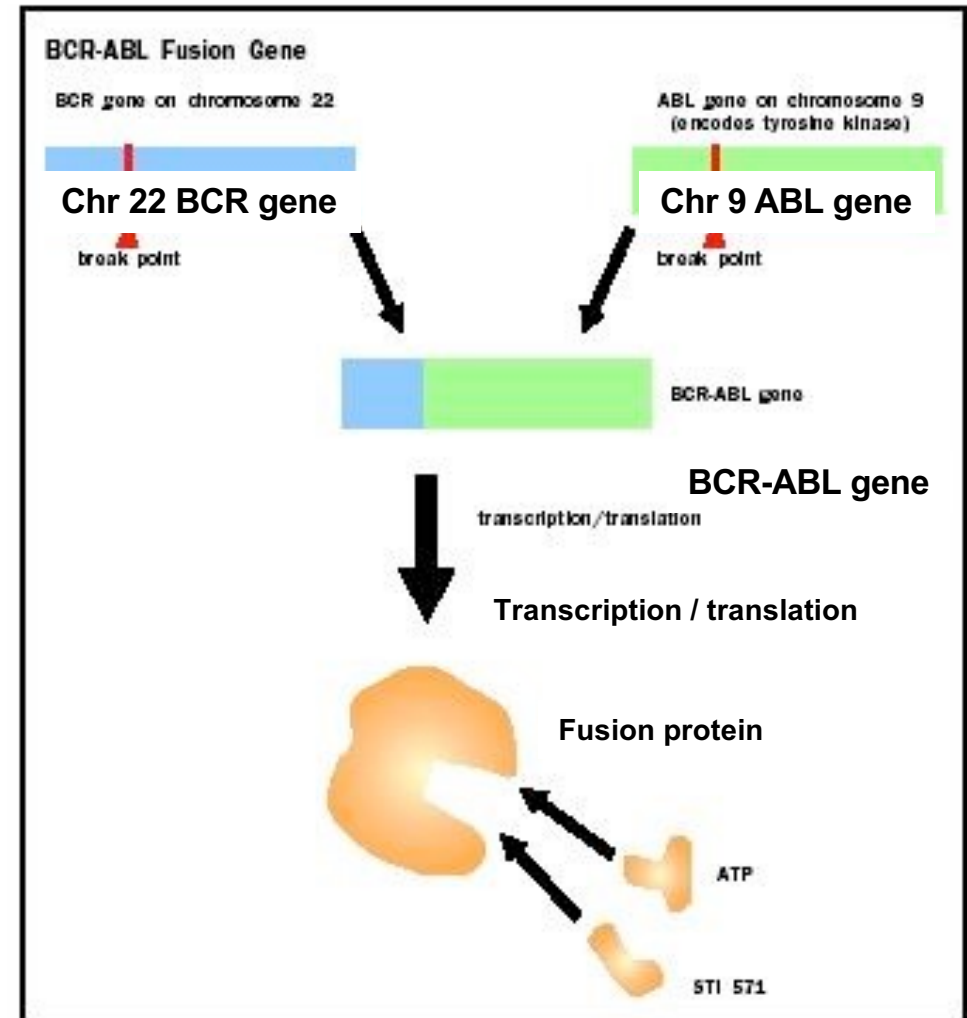
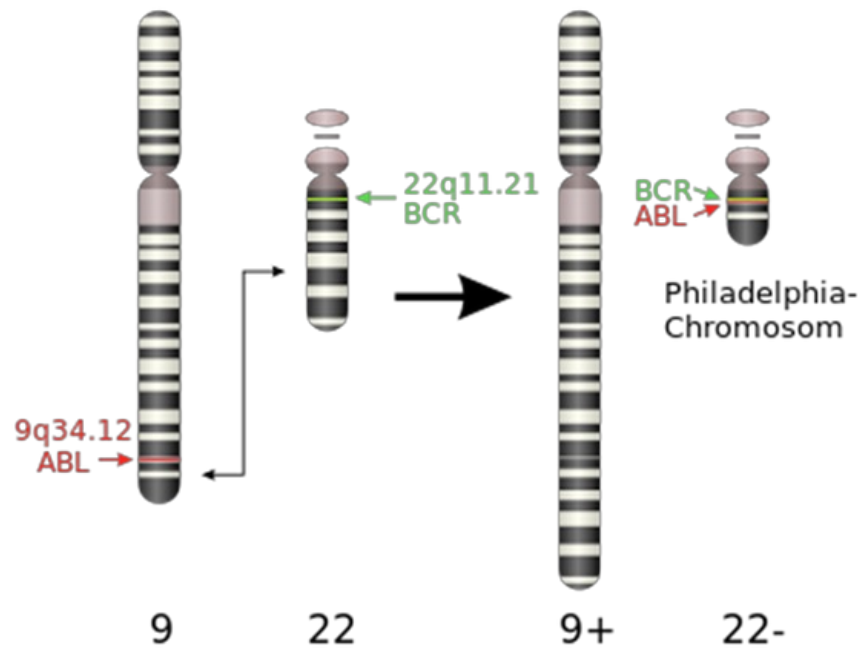


# Philadelphia-Chromosom

- Translokation zwischen Chr. 22 → 9
- Haben 90% aller Patienten mit Chronic Myelogenous Leukemia (CML)

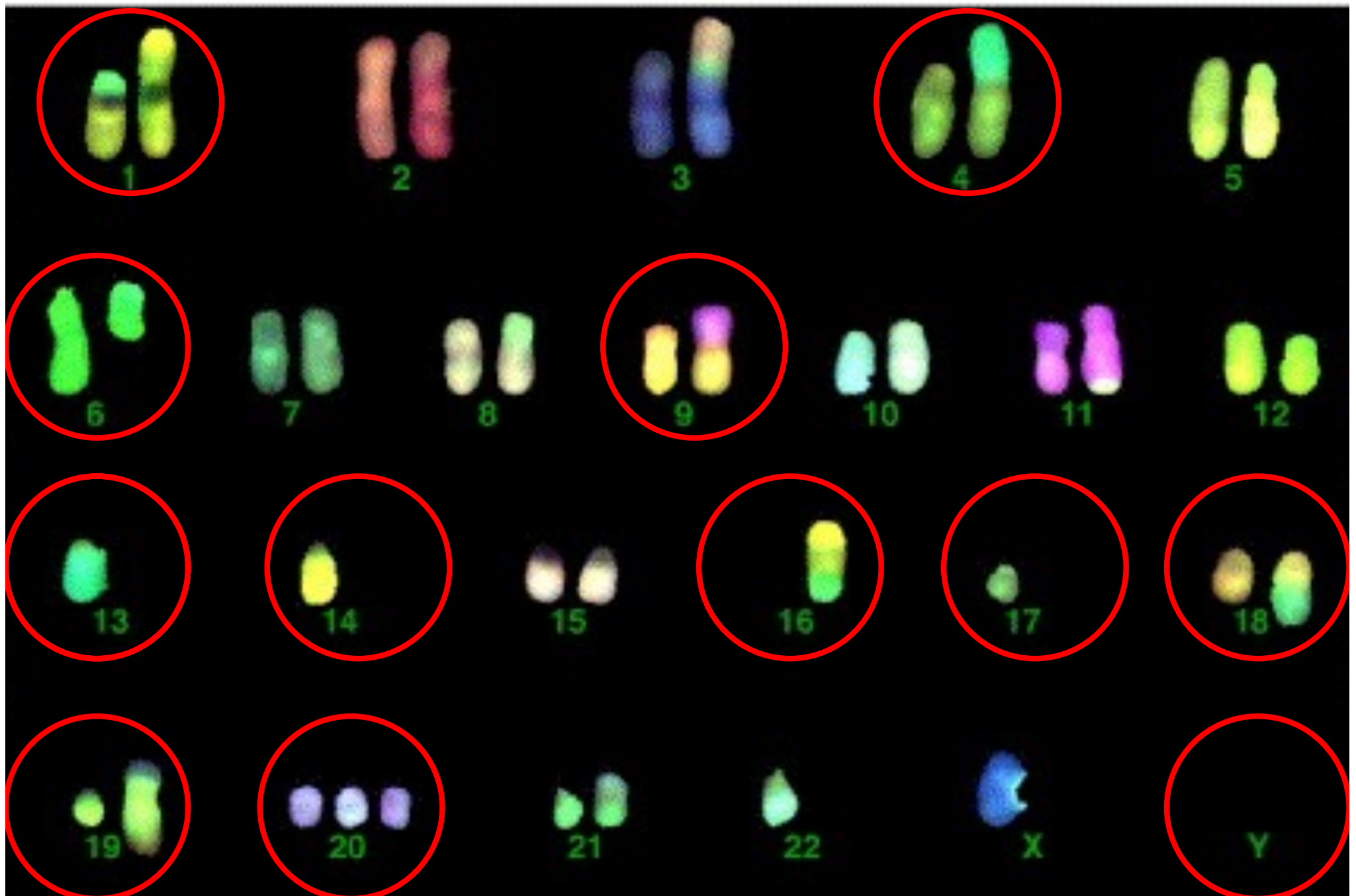


# Fusion gene auf dem Philadelphia-Chromosom



The fusion gene activates signal transduction pathways such as RAS/MAPK, PI-3 kinase, c-CBL and CRKL pathways, JAK-STAT, and the Src pathway.

# Darmkrebs

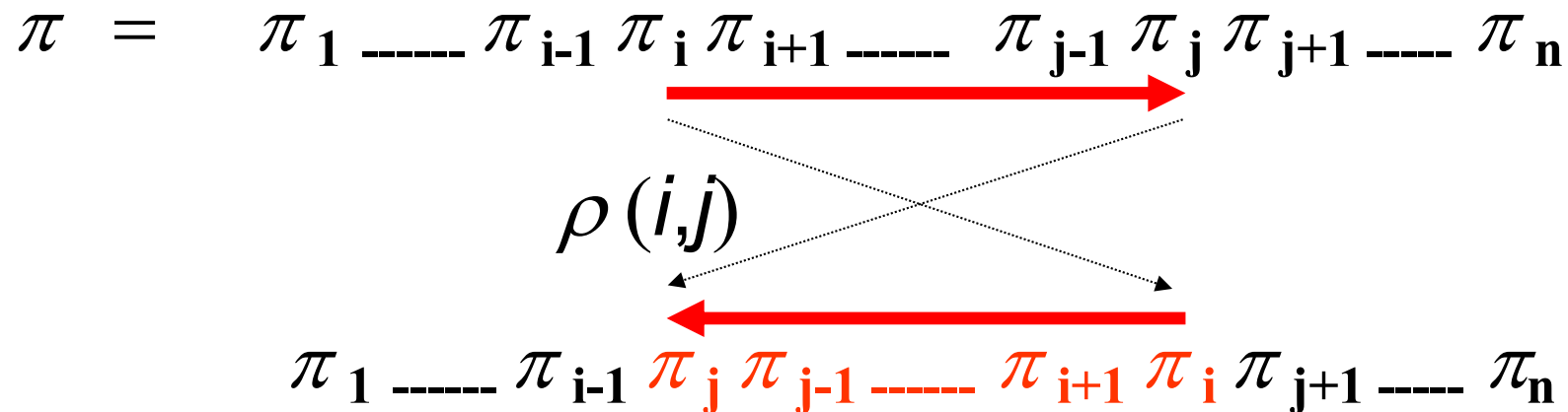






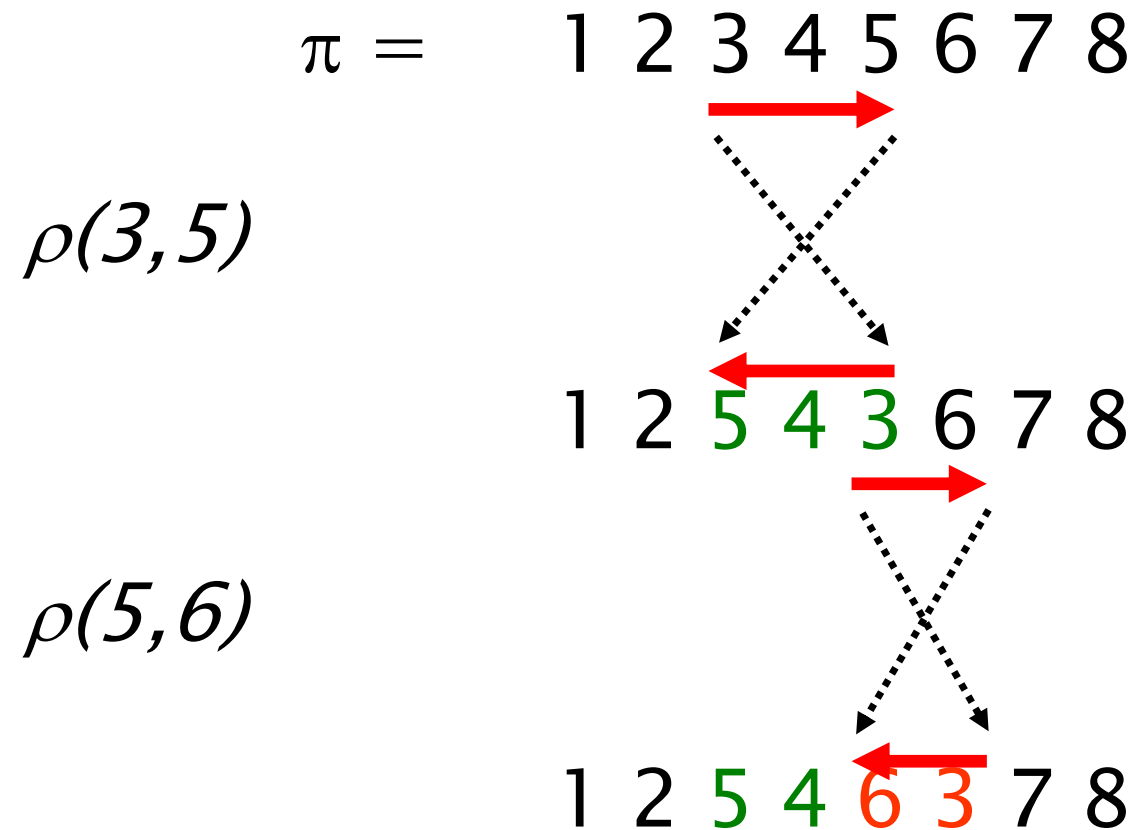
# Umkehrungen und Genordnung

Die Genordnung wird durch eine Permutation  $\pi$  beschrieben:



Die Umkehrung  $\rho(i, j)$  verdreht ein 'Bruchstück' von  $i$  bis  $j$  in  $\pi$ .

# Umkehrungen: Beispiel



# Umkehrungsabstand: Das Problem

Ziel: Für zwei vorgegebene Permutationen: finde die kürzeste Abfolge von Umkehrungen, die eine Permutation in die andere überführt.

Eingabe: Permutationen  $\pi$  und  $\sigma$

Ausgabe: Eine Abfolge von Umkehrungen  $\rho_1, \dots, \rho_t$ , die  $\pi$  in  $\sigma$  überführt, so dass  $t$  minimal ist (für alle möglichen Abfolgen  $\rho_1, \dots, \rho_v$ )

Definition:

$d(\pi, \sigma)$  := Umkehrungsabstand  
:= kleinstmöglicher Wert für  $t$  für gegebene  $\pi$  und  $\sigma$

# Sortieren durch Umkehrungen: Das Problem

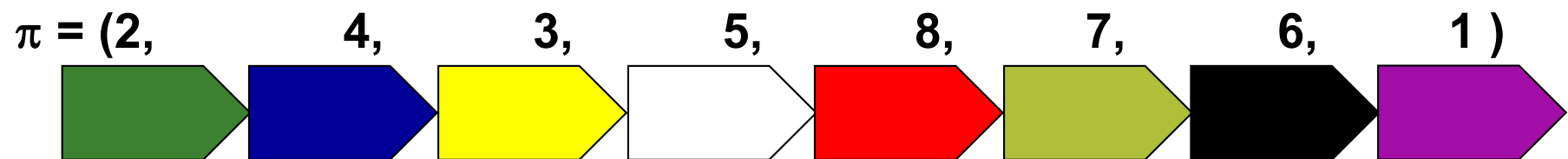
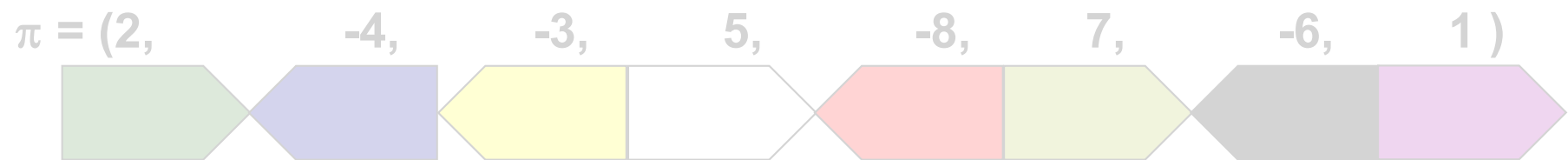
Ziel: Für **eine** vorgegebene Permutationen: finde die kürzeste Abfolge von Umkehrungen, die sie in die Identitäts-Permutation  $(1, \dots, n)$  überführt.

Eingabe: Permutation  $\pi$

Ausgabe: Eine Abfolge von Umkehrungen  $\rho_1, \dots, \rho_t$ , die  $\pi$  in die Einheits-Permutation überführt, so dass  $t$  minimal ist (für alle möglichen Abfolgen  $\rho_1, \dots, \rho_v$ )



# Vorzeichen?



# Das Pfannkuchen-Problem

- Der Koch ist nachlässig:  
seine Pfannkuchen haben alle unterschiedliche Größen.
- Der Kellner will sie ordnen:  
der Größte unten, der Kleinste oben.
- Dazu nimmt er jeweils eine Gruppe von Pfannkuchen hoch und dreht sie um (*Flip*)



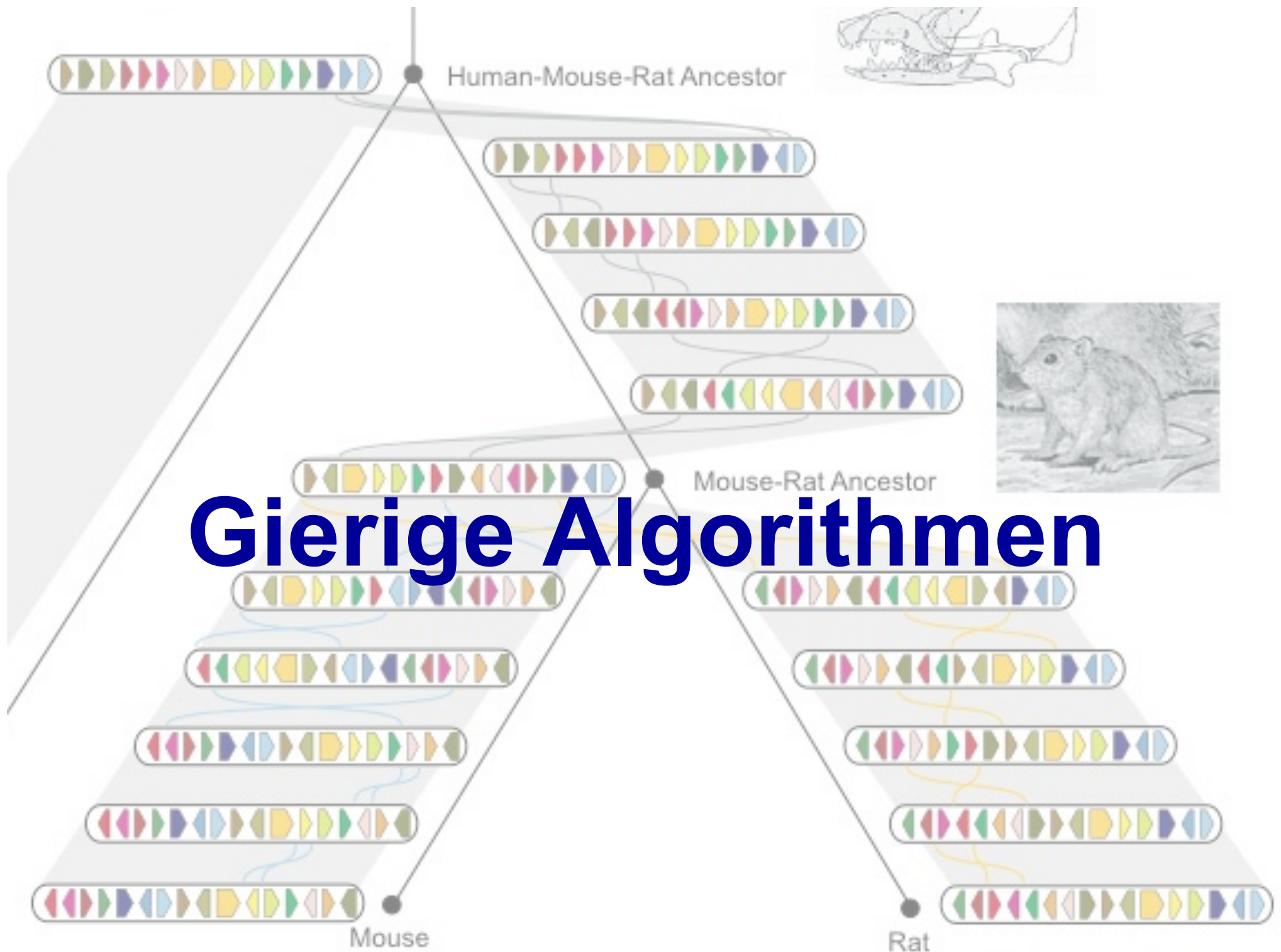
Gegeben ein beliebiger Stapel, löse das mit der minimalen Anzahl von *Flips*.

# Das Pfannkuchen-Problem

Ziel: Für einen vorgegebenen Stapel von  $n$  Pfannkuchen: Was ist die minimale Anzahl von Flips, um den Stapel zu ordnen?

Eingabe: Permutation  $\pi$

Ausgabe: Eine Abfolge von Präfix-Umkehrungen  $\rho_1, \dots, \rho_t$ , die  $\pi$  in die Identitäts-Permutation überführen, so dass  $t$  minimal ist (für alle möglichen  $\rho_1, \dots, \rho_v$ ).



# Pfannkuchen-Problem: Gieriger Algorithmus

Idee: Sortiere den jeweils größten unsortierten Pfannkuchen ein: zuerst nach oben (1. Flip), dann an seinen Platz (2. Flip)

⇒ max.  $2(n - 1)$  Flips insgesamt  
(der letzte Pfannkuchen stimmt automatisch)



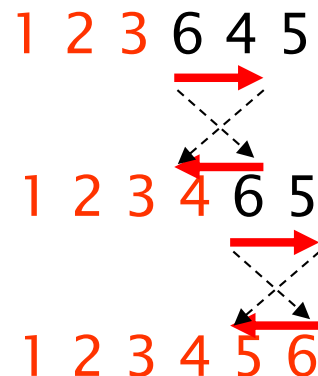
*Christos Papadimitriou &  
William Gates:  
höchstens  $\frac{5}{3}(n + 1)$  Flips*

# Sortieren durch Umkehrungen: Gieriger Ansatz

In  $\pi = 1\ 2\ 3\ 6\ 4\ 5$  sind die ersten drei Elemente sortiert  
→ lassen wir so

Die Länge des bereits sortierten Präfixes nennen wir  $prefix(\pi)$  (=3)

Idee: erhöhe  $prefix(\pi)$  bei jedem Schritt



⇒ Anzahl von Umkehrungen ist max.  $(n - 1)$

# Gieriger Algorithmus: Pseudocode

SimpleReversalSort( $\pi$ )

```
1  for  $i \leftarrow 1$  to  $n - 1$ 
2       $j \leftarrow$  Position des  $i$ -ten Elements in  $\pi$  (i.e.,  $\pi_j = i$ )
3      if  $j \neq i$ 
4           $\pi \leftarrow \pi * \rho(i, j)$ 
5          output  $\pi$ 
6      if ( $\pi$  ist Identiät) return
```

**Liefert SimpleReversalSort die optimale Lösung?**

$\pi = (6, 1, 2, 3, 4, 5)$



# Approximationsalgorithmen (Minimierung)

Sei  $\text{OPT}(\pi) \neq 0$  der optimale Lösungswert bei Eingabe  $\pi$ , zum Beispiel die optimale Anzahl Reversals beim Sortieren durch Umkehren oder die minimale Distanz (TotalDistance) beim Median-String-Problem.

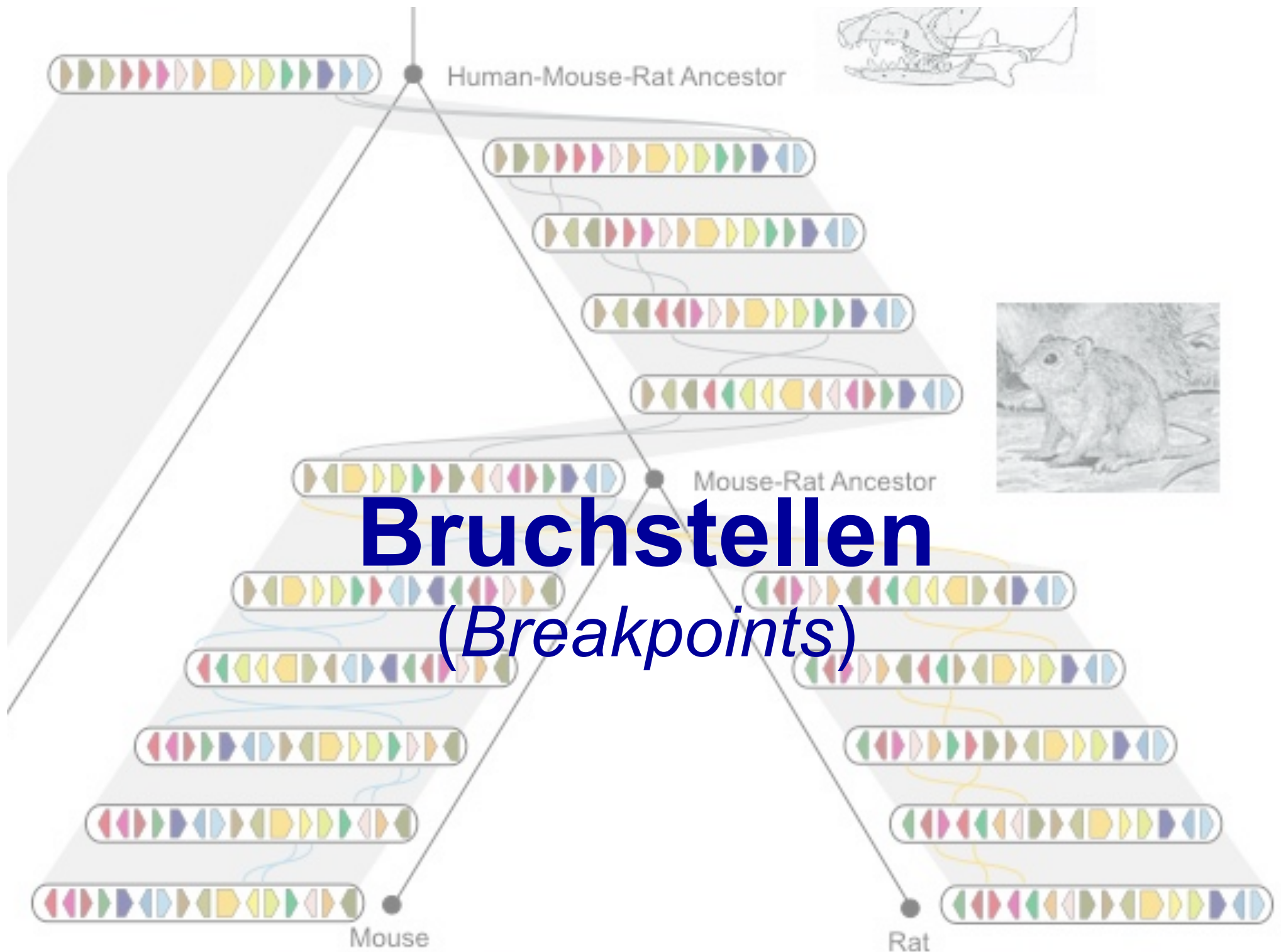
Dann ist der **Approximationsfaktor** eines Algorithmus A **auf Eingabe  $\pi$** :

$$A(\pi) / \text{OPT}(\pi)$$

**Approximationsfaktor** eines Algorithmus A:  
**schlechtester** Approximationsfaktor über **alle** Eingaben. Bei Minimierungsproblemen:

$$\alpha = \max_{\pi} A(\pi) / \text{OPT}(\pi)$$

Bemerkung: bei Maximierungsproblemen gilt  $\alpha = \min_{\pi} A(\pi) / \text{OPT}(\pi)$



# Nachbarschaften & Bruchstellen

Permutation  $\pi = \pi_1\pi_2\pi_3\ldots\pi_{n-1}\pi_n$

## Definitionen:

- Zwei Elemente  $\pi_i$  und  $\pi_{i+1}$  sind **benachbart**  
 $:\Leftrightarrow \pi_{i+1} = \pi_i \pm 1$
- Zwischen zwei Elementen  $\pi_i$  und  $\pi_{i+1}$  befindet sich eine **Bruchstelle**  
 $:\Leftrightarrow \pi_i$  und  $\pi_{i+1}$  sind nicht benachbart
- $b(\pi) := \#$  Bruchstellen in  $\pi$

Beispiel:  $\pi = 1 \ 9 \ 3 \ 4 \ 7 \ 8 \ 2 \ 6 \ 5$

# Erweiterte Permutationen

Damit die Bruchstellen vollständige Auskunft über die Sortierung geben, fügen wir zwei Elemente  $\pi_0 = 0$  und  $\pi_{n+1} = n+1$  an.

Beispiel:

$$\begin{array}{ccccccccccc} \pi = & 1 & | & 9 & | & 3 & 4 & | & 7 & 8 & | & 2 & | & 6 & 5 \\ & & & & & & & & \downarrow & & & & & & \\ \pi = & 0 & | & 1 & | & 9 & | & 3 & 4 & | & 7 & 8 & | & 2 & | & 6 & 5 & | & 10 \end{array}$$

# Umkehrungsabstand & Bruchstellen

Jede Umkehrung entfernt maximal zwei Bruchstellen

Beispiel:

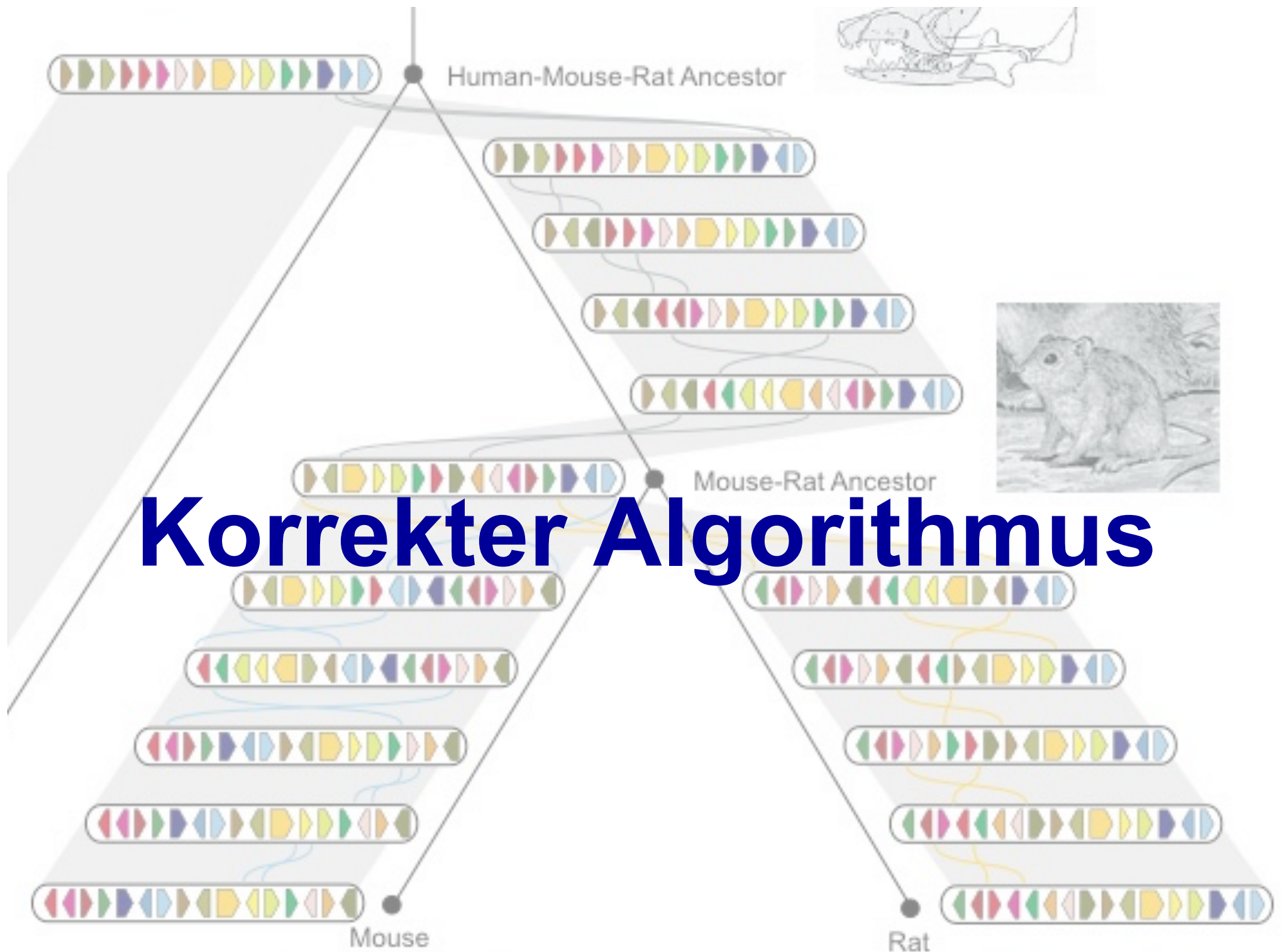
$$\pi = 0 \ 2 \ 3 \ 1 \ 4 \ 6 \ 5 \ 7$$

$\Rightarrow$  Umkehrungsabstand  $\geq \# \text{ Bruchstellen} / 2$

# Besser Sortieren durch Bruchstellen?

BreakPointReversalSort?( $\pi$ )

```
1  while  $b(\pi) > 0$ 
2      Finde die Umkehrung  $\rho(i,j)$ , die  $b(\pi \cdot \rho(i,j))$  minimiert
3       $\pi \leftarrow \pi \cdot \rho(i, j)$ 
4      output  $\pi$ 
5  return
```





# Strips

## Definitionen:

- Streifen: ein Intervall zwischen zwei aufeinanderfolgenden Bruchstellen
- Abnehmender Streifen: Elemente in abnehmender Folge  
(z. B. 6 5 oder 4 3 2)
- Zunehmender Streifen: Elemente in zunehmender Folge  
(z. B. 5 6 oder 2 3 4)
- 1-elementige Streifen können zu- oder abnehmend sein.  
Wir definieren sie als abnehmend (mit Ausnahme von 0 und  $n + 1$ )

## Beispiel:



# Reduzierung der Bruchstellen

## Theorem:

Wenn die Permutation  $\pi$  wenigstens einen abnehmenden Streifen hat, dann existiert eine Umkehrung  $\rho$ , die die Anzahl der Breakpoints vermindert (d.h.:  $\exists i \ b(\pi \cdot \rho) < b(\pi)$ ).

## Beweis:

Wähle denjenigen abnehmenden Streifen in  $\pi$  mit dem kleinsten Element  $k$ . Dann ist  $k$  das letzte Element in diesem Streifen. Dann muss  $k-1$  das letzte Element in einem zunehmenden Streifen sein. Umkehrung des Intervalls zwischen  $k-1$  und  $k$  einschließlich  $k$  (wenn  $k-1$  links von  $k$ ; sonst zwischen  $k$  und  $k-1$  einschließlich  $k-1$ ) bringt die beiden zusammen und beseitigt damit mind. eine Bruchstelle. QED

## Was, wenn alle Streifen zunehmend sind?

Dann drehen wir einfach einen um!

# Korrektes Sortieren durch Umkehrungen!

CorrectBreakpointReversalSort( $\pi$ )

```
1  while  $b(\pi) > 0$ 
2      if  $\pi$  hat einen abnehmenden Streifen
3          Wähle Umkehrung  $\rho$ , die  $b(\pi \cdot \rho)$  minimiert
4      else
5          Wähle Umkehrung  $\rho$  für zunehmenden Streifen in  $\pi$ 
6           $\pi \leftarrow \pi \cdot \rho$ 
7      output  $\pi$ 
8  return
```

# Performance-Garantie

$$= \frac{\max(\text{Schritte in Approximation})}{\min(\text{Schritte in optimaler Lösung})}$$

Approximation durch CorrectBreakpointReversalSort():

Mindestens jeder 2. Schritt eliminiert mindestens eine Bruchstelle

$\Rightarrow$  max.  $2b(\pi)$  Schritte

Bestmögliche Lösung:

Jeder Schritt eliminiert höchstens 2 Bruchstellen

$\Rightarrow$  min.  $b(\pi)/2$  Schritte

➤ Performance-Garantie:  $A = \frac{2b(\pi)}{\frac{b(\pi)}{2}} = 4$

# Gierige Motif-Suche



## Take home messages

- Verständnis von Genomumordnungen ist wichtig
- Eine Version (nicht die relevanteste)
  - Sortieren von Permutationen ohne Vorzeichen
- Approximationsalgorithmen
- Gierige Algorithmen
- Sortieren über Bruchstellen ergibt Faktor 4, viel besser als SimpleReversalSort

# Das Motivsucheproblem

Ziel: Finde genau ein  $\ell$ -mer für jede DNA-Sequenz einer vorgegebenen Menge, so dass der Konsensusscore über die  $\ell$ -mere maximal ist.

Eingabe: Eine  $t \times n$  Matrix **DNA** sowie  $\ell$  (die Länge des gesuchten Motivs)

Ausgabe: Ein Vektor mit  $t$  Startpositionen  $\mathbf{s} = (s_1, s_2, \dots, s_t)$ , so dass  $\text{Score}(\mathbf{s}, t, \mathbf{DNA})$  maximal ist.

Brachialmethode:  $O(\ell t n^t)$ , etwas eleganter  $O(\ell t n^4)$ : [letzte Vorlesung]

➤ zu langsam für realistische Probleme



# Gierige Motivsuche: CONSENSUS

**Idee:** Suche die beste Lösung für Sequenz  $i$  unter der Annahme, dass die beste Lösung für  $1 \dots i-1$  schon gefunden ist!

Starte mit 2 der  $t$  Sequenzen

- Finde die 2 ähnlichsten  $l$ -mere darin, d.h.:  
finde  $\mathbf{s}=(s_1, s_2)$  mit maximalem  $\text{Score}(\mathbf{s}, 2, \text{DNA})$
- Iterativ für Sequenz  $i = 3, \dots, t$ :  
Finde das dazu ähnlichste  $l$ -mer in Sequenz  $i$ , d.h.:  
finde  $s_i$  mit maximalem  $\text{Score}(\mathbf{s}, i, \text{DNA})$  für  $\mathbf{s}=(s_1, \dots, s_{i-1}, s_i, \dots)$

Score berechnet  
für 2 Sequenzen

**CONSENSUS:** Permutiere Reihenfolge **1000** mal  
und merke bestes Motiv.

Score berechnet  
für  $i$  Sequenzen

**Resultat:**

- Nicht exakt (evtl. gibt es bessere Motive)
- Schlechter Approximationsfaktor (Übung)
- Sehr schnell (die meiste Zeit für die 2 ersten  $l$ -mere)
- Ergebnis häufig sehr gut

# GreedyMotifSearch()

GreedyMotifSearch( $DNA, t, n, l$ )

1. **bestMotif**  $\leftarrow \text{?}(1, 1, \dots, 1)$

2. **s**  $\leftarrow \text{?}(1, 1, \dots, 1)$

3. **for**  $s_1 \leftarrow \text{?} 1$  **to**  $n - l + 1$

4. **for**  $s_2 \leftarrow \text{?} 1$  **to**  $n - l + 1$

5. **if**  $\text{Score}(\mathbf{s}, 2, DNA) > \text{Score}(\mathbf{bestMotif}, 2, DNA)$

6.  $bestMotif_1 \leftarrow \text{?} s_1$

7.  $bestMotif_2 \leftarrow \text{?} s_2$

8.  $s_1 \leftarrow \text{?} bestMotif_1$

9.  $s_2 \leftarrow \text{?} bestMotif_2$

10. **for**  $i \leftarrow \text{?} 3$  **to**  $t$

11. **for**  $s_i \leftarrow \text{?} 2$  **to**  $n - l + \text{?}$

12. **if**  $\text{Score}(\mathbf{s}, i, DNA) > \text{Score}(\mathbf{bestMotif}, i, DNA)$

13.  $\text{?} bestMotif_i \leftarrow \text{?} s_i$

14.  $s_i \leftarrow \text{?} bestMotif_i$

15. **return** **bestMotif**

Bisher bester  
Startvektor

Aktueller  
Startvektor

Score berechnet  
für  $t'=2$

Score berechnet  
für  $t'=i$

# Programme zur Motivsuche

- **CONSENSUS**

*Hertz, Stormo (1989-1999)*

- **GibbsDNA**

*Lawrence et al (1993)*

- **MITRA**

*Eskin, Pevzner (2002)*

- **Pattern Branching**

*Price, Pevzner (2003)*

- **Grisotto**

*Carvallo, Olivera (2011)*

- **MEME**

*Bailey et al. (1995-2011)*

- ...

## Take home messages

- Verständnis von Genomumordnungen ist wichtig
- Eine Version (nicht die relevanteste)
  - Sortieren von Permutationen ohne Vorzeichen
- Approximationsalgorithmen
- Gierige Algorithmen
- Sortieren über Bruchstellen ergibt Faktor 4, viel besser als SimpleReversalSort
- Motivsuche geht viel schneller (aber nicht mehr exakt) mit gierigen Algorithmen.