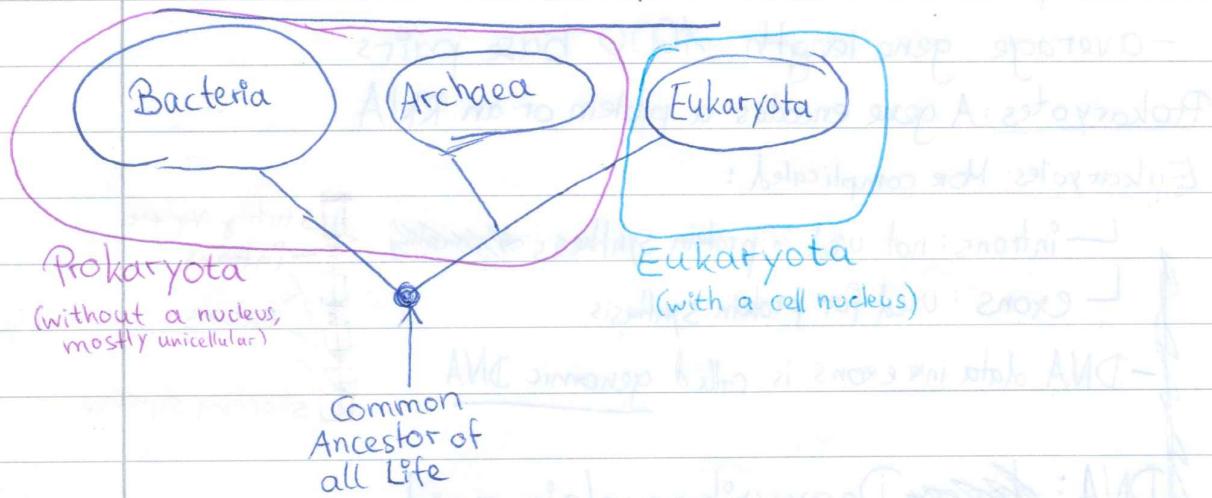


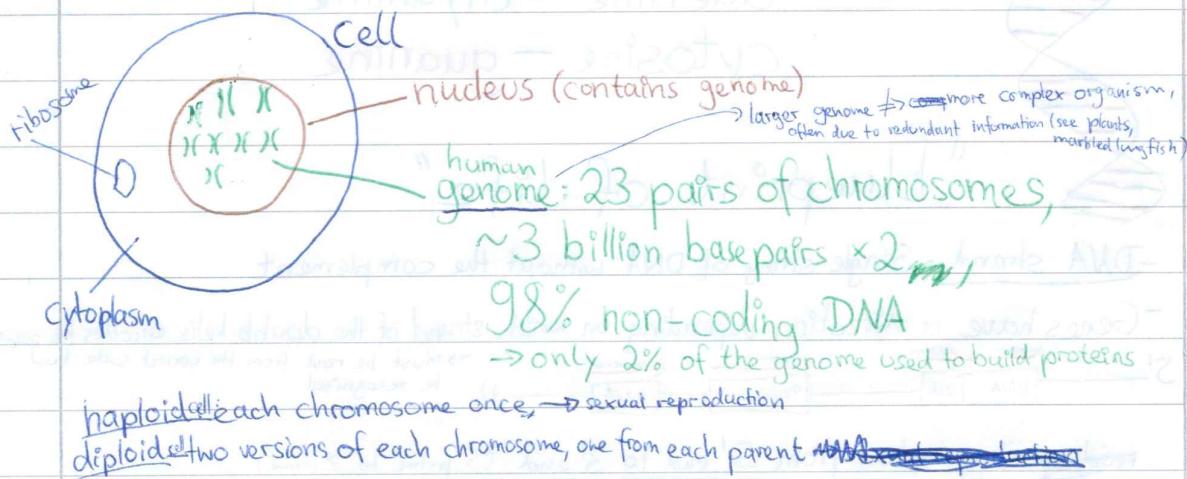
# Biological Background Knowledge

## 3 domains of life



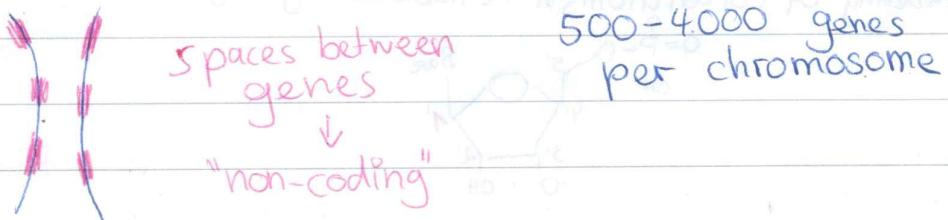
Phylogenetic tree based on the 16S rRNA gene

## Eukaryota, human example



## Chromosomes (long DNA molecules)

- Prokaryotes: one chromosome
- Eukaryotes: many chromosomes, organized in pairs (paternal/maternal)
- ~~paired~~ homologous: paired chromosomes



Genes :- coding parts of DNA, each gene encodes for either RNA or a protein

- average gene length: 1000 base pairs

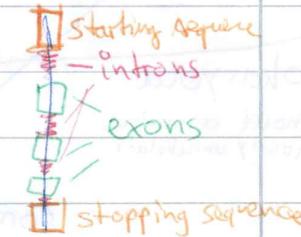
Prokaryotes: A gene encodes a protein or an RNA

Eukaryotes: More complicated:

└ introns: not used in protein synthesis

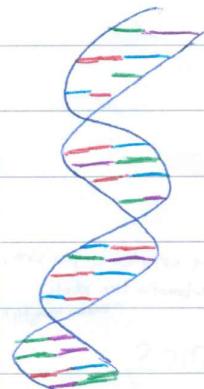
└ exons: used for protein synthesis

└ DNA data in exons is called genomic DNA



DNA: Deoxyribonucleic acid

- double helix structure, organized in chromosomes



4 nucleic acids/base pairs/nucleotides:

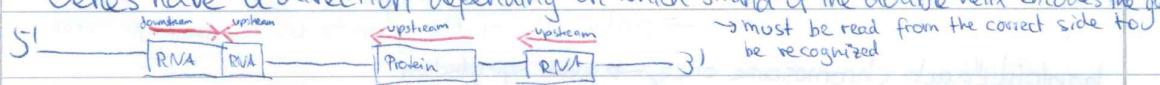
adenine - thymine

cytosine - guanine

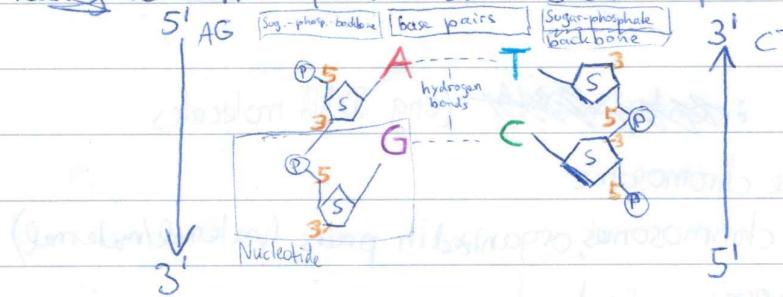
"blueprint of life"

- DNA strand: Single string of DNA without the complement

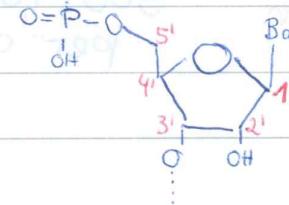
- Genes have a direction depending on which strand of the double helix encodes the gene



Reading convention: from 5' end to 3' end ("5 prime to 3 prime")



Labelling of carbon atoms in the nucleotide sugar-ring:



# Biological Background Knowledge II

## RNA: Ribonucleic acid

- Single-stranded sequence
- Thymine is replaced by Uracil

- RNA is a copy of a coding DNA strand (a gene)

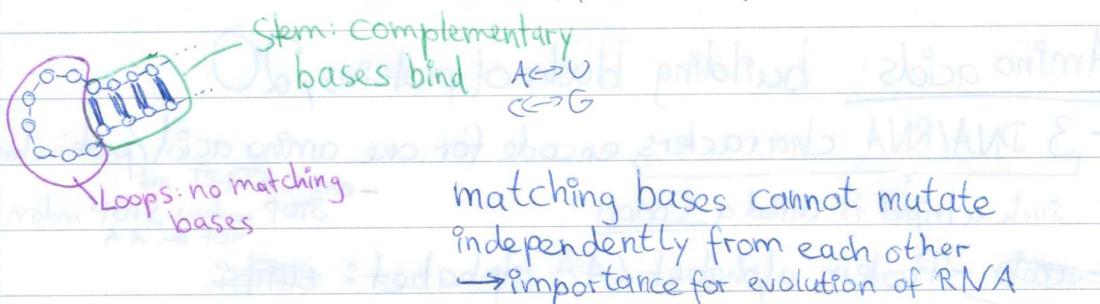
- involved in transcription to construct either:

↳ Coding RNA: ~~also~~ involved in ~~transcription~~ translation (protein synthesis)

↳ non-coding RNA: some other direct function in the cell

- has a secondary structure (and even a tertiary):

↳ influence on its function



## - types of RNA

~3% Coding RNA:

mRNA

(messenger RNA)

transports RNA

data to the ribosome  
for protein synthesis

non-coding RNA

tRNA

(ribosomal RNA)

carries out  
the translation  
in the ribosome  
via catalysis

tRNA

(transfer RNA)

brings in  
the amino  
acids

{  
One of a few  
genes that are  
common to all species

↳ most well-known:

16S rRNA gene

↳ can be used to  
infer evolutionary  
relationships amongst  
all species, useful for  
classification of bacteria

## - Transcriptome: Set of all RNA molecules in a cell

↳ reflects the activity in a cell → temporal and spatial component

- Specialized cells → different genes are active

- point of time changes transcriptome

- Meta-Genome: blind sequencing of all genetic material of a bacterial community → many species,

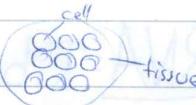
figure out microbial diversity. Can be done at whole-genome level → metagenomics

↳ single gene level → metagenetics, e.g. 16S RNA for

Bacteria

## Proteins: do lots of stuff in the body

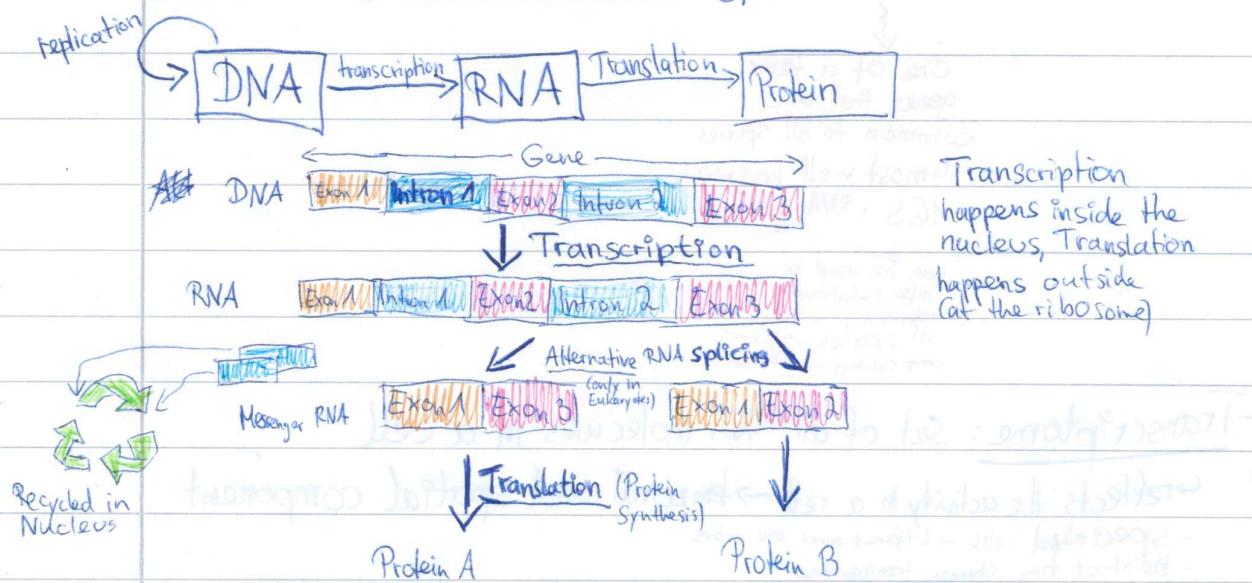
- └ Structural proteins: tissue building blocks
- └ Enzymatic proteins: catalysts of biochemical reactions in the body
  - examples: oxygen transport, immune defense, provide & store energy, ...
  - human: ~ 20.000 proteins
  - Sequence/string of amino acid characters, protein letters are called residues
  - have primary, secondary, tertiary and quaternary structure



## Amino acids: building blocks of proteins, 20

- 3 DNA/RNA characters encode for one amino acid/protein character
- + Such a triplet is called a Codon
  - extra START and STOP codons, STOP codon is not an AA
- ~~Protein alphabet / AA alphabet:~~
- ↳ DNA → Protein: not ambiguous, but redundant
- ↳ Protein → DNA: ambiguous, several DNA triplets encode for the same AA
- Synonymous substitutions/mutations ( $\text{GCC} \rightarrow \text{GCT} \equiv \text{Alanine} \rightarrow \text{Alanine}$ )
- non-synonymous substitutions/mutations ( $\text{GGT} \rightarrow \text{GTT} \equiv \text{Glycine} \rightarrow \text{Valine}$ )
- 3rd codon position is less vulnerable to mutations
  - ↳ mostly results in the same AA after mutation

## Central Dogma of Molecular Biology



# Biological Background Knowledge III

DNA Sequencing: reading the nucleotide bases in a DNA molecule

- Extended DNA ~~sequence~~ alphabet

Since sequencing methods are not ~~exact~~ exact

→ character set to denote "A", "A or C", "A or C or G",

"A or C or G or T", ... and a gap-character

(needed for alignment)

- Sanger Sequencing (first generation sequencing, ~1977)

↳ Chain Termination Method

- Informal game analogy:

You think of a ① You think of a 200 digit number, others must guess it

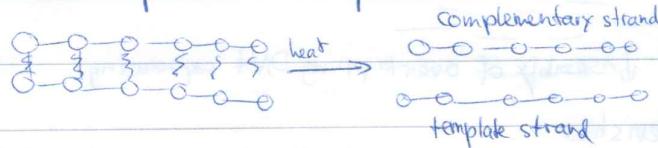
② You ~~read~~ read the number in your head

③ Someone randomly interrupts

④ You answer current position and digit.

⑤ Start again with step 2.

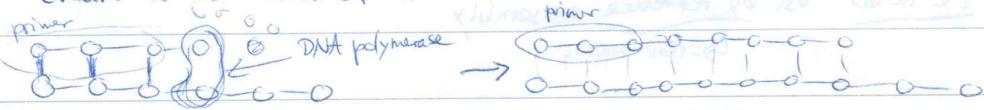
- ① heat up DNA to separate its strands



② lower temperature to allow the primer sequence to bind to its complementary sequence in the template DNA



③ raise temperature s.t. the enzyme polymerase 3 combines to the DNA and creates a new strand of DNA



④ dNTPs are added by the enzyme until a ddNTP is added which stops synthesis

↳ creates DNA strands of different length

colored & chemically altered version of A, C, G, T

⑤ apply heat again to separate

⑥ let the pieces run a race on ~~separate~~ a gel with 4 lanes, (gel electrophoresis)

Smaller pieces are faster than longer ones, letters can be seen by color

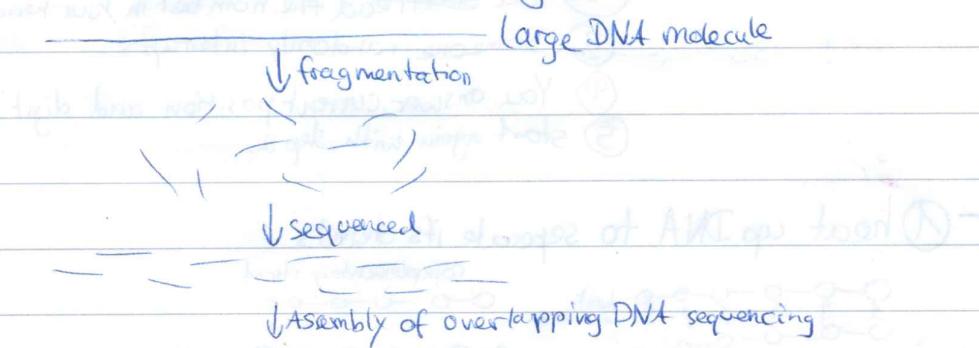
### III Sequencing Methods

Sanger Sequencing	Next Generation Sequencing (pyrosequencing)
<del>fast, selective</del>	- cheap
- expensive	- fast, many sequences
- slow, few sequences	- short sequences
- long sequences	- less accurate
- high accuracy	

pyrosequencing: Die Dinger geben direkt Licht ab was gemessen werden kann → kein Abbruch

#### Shotgun Sequencing: reading the genome

- ① Break up genome randomly into fragments
- ② Read fragments
- ③ Assemble fragments into genome



#### Important characteristics

- Coverage (how many fragments/reads cover one nucleotide on the genome)
- Fragment length
- Paired-end (ATG-----GTC) vs. single-end reads
  - ↳ additional information makes assembly process easier
  - we know the length of this gaps, same for all fragments
- De novo vs. by reference assembly
  - ↳ read mapping

De novo assembly: a library mapping all yet unknown DNA → assembling separate blocks → it's not guaranteed that blocks fit together → assembly of single read blocks → what does it mean?

# Pairwise Sequence Alignment

- Sequence Alignment: compare strings to infer their similarity
- Online sequence alignment: no preprocessing for index building
- Edit distance ( $\geq 0$ ,  $= 0$  iff same string, symmetric,  $\Delta$ -inequality)
  - ↳ Damerau-Levenshtein:  $\text{sub}(a, b) = \text{sub}(b, a) = \text{del}(a) = \text{ins}(b) = 1$ ,  $a \neq b$
  - ↳ Hamming distance: ~~del(a) = ins(a) = 0~~, count only substitutions  
minimum cost of a set of edit operations to transform a into b

- alignment: for example

(A C G T P A)  
 (A T G C T A) no ~~gap~~ (-) allowed  
 here:  $x = \text{ACGA}$ ,  
 $y = \text{ATGCTA}$   
 → cost 3

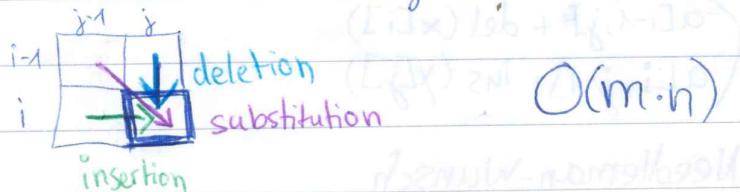
## Edit distance via dynamic programming

$$a[0, 0] = 0$$

$$a[i, 0] = a[i-1, 0] + \text{del}(x[i])$$

$$a[0, j] = a[0, j-1] + \text{ins}(x[j])$$

$$a[i, j] = \min \begin{cases} a[i-1, j-1] + \text{sub}(x[i], y[j]) \\ a[i-1, j] + \text{del}(x[i]) \\ a[i, j-1] + \text{ins}(x[i]) \end{cases}$$



## Needleman-Wunsch: Global alignment

- maximize similarity,  $\text{sub}(a, a) > 0$ ,  $\text{ins}(a), \text{del}(a), \text{sub}(a, b) < 0$

indel

## Substitution Matrices: BLOSUM

(BLOCKS of AMINO ACID Substitution Matrix)

→ for sequence alignment of protein sequences

→ based on observed (local) alignments

→ got data from conserved regions of protein families w/o gaps

→ different matrices for closely/distantly related sequences

$\text{Sub}(a, b) = \log \left( \frac{\text{P}_{\text{biological sense}}}{\text{P}_{\text{by chance}}} \right) \in \text{log-odds score}$

→ positive score: more likely

## Hamming distance

Input: String  $x$ , String  $y$ , integer  $k$

Wanted: Alignments with hamming distance  $\leq k$

$$a[i,j] = \begin{cases} k+1, & j=0, i>0 \\ 0, & i=0 \\ a[i-1,j-1], & x[i]=y[i] \\ a[i-1,j-1]+1, & x[i]\neq y[i] \end{cases}$$



traceback from all values in the last row that are  $\leq k$  to a zero.

## Smith-Waterman: Local Alignment

- Needleman-Wunsch with free taxi drives

$$a[i,j] = \max \begin{cases} 0, & i=0 \vee j=0 \\ a[i-1,j-1] + \text{sub}(x[i], y[j]), \\ a[i-1,j-1] + \text{del}(x[i]), \\ a[i, j-1] + \text{ins}(y[j]) \end{cases}$$

free taxi drive



differences to Needleman-Wunsch

- other initialization (free drive to your starting position)
- best value stands anywhere in the table (free drive to your ending position)

## Global vs Local Alignment

- find the highest scoring alignment between two strings

- good for similar strings with ~~small~~ Similar size

- find the highest scoring alignment between two substrings

↳ global alignment in a subrectangle of the DP-matrix  
- idea: free taxi rides



- better for different strings with small similar regions (i.e. genes in different species)  
↳ expected to have small regions of similarity

# Searching Similar Sequences

- Assumption: Similar sequences are evolutionary and/or functionally related

- Naive approach: Pair-wise alignment of query sequence

with every sequence in the database via Smith-Waterman, report best matches

→  $\Theta(m \cdot n \cdot d)$  → too slow, since database size  $d$  often large

- BLAST: Basic Local Alignment Search Tool

↳ heuristic to find similar sequences, widely-used

↳ good for large, fast evolving database; small number of queries (compared to DB size); interested in distantly related sequences

example:  
→ found new gene in mouse,  
perform BLAST search on human genome  
to check if humans have a similar gene

↳ different flavors: (database, query) ∈ {nucleotide sequence, protein sequence,  
nucleotides translated into protein sequences}

↳ Idea: reduce search space by ignoring dissimilar regions

① Seeding: find common subwords between query and database sequences → seeds

- collect all substrings of length  $k$  of the query sequence
- for each of these substrings, build neighborhood with "similar" subwords of length  $k$

↳ defined via substitution matrix, BLOSUM for proteins,  
+2/-3 (or +5/-4) for DNA match/mismatch

add only subwords with a similarity threshold  $T$ .

→ neighbors are only used to find more positions to look at

~~Scanning~~ Search

② Scanning: Scan the database for exact matches of ~~sets~~ seeds:

③ Extension: Starting from seeds, extend alignment in both directions → high-scoring segment pairs (HSP)

↳ try to extend the alignment to the left and to the right from the seed, stop if current total score < current MaxScore -  $X$ , trim back the alignment to curr MaxScore

6/36

### ③ Evaluation: Assess the statistical significance of each HSP

↳ are the sequences biologically related or just by chance similar?

↳ Smith-Waterman local alignment scores between two

random sequences follow Gumbel extreme value distribution,  
 $p(S \geq x) = 1 - \exp(-e^{-x(x-p)})$ , depending on substitution matrix,  
 gap penalties, sequence length, nucleotide frequencies

↳ BLAST reports e-value (expectation value):

expected number of times that an unrelated database sequence would obtain a score  $S$

higher than  $x$  by chance

$$E \approx 1 - e^{-p(S > x) \cdot d}$$

### Genome assembly: ~~reconstructing~~

taking a large number of short DNA sequences and putting them back together, reconstruct genome

Genome (multiple copies)      A long DNA molecule

ACTTGCTA...

↓ Shotgun sequencing

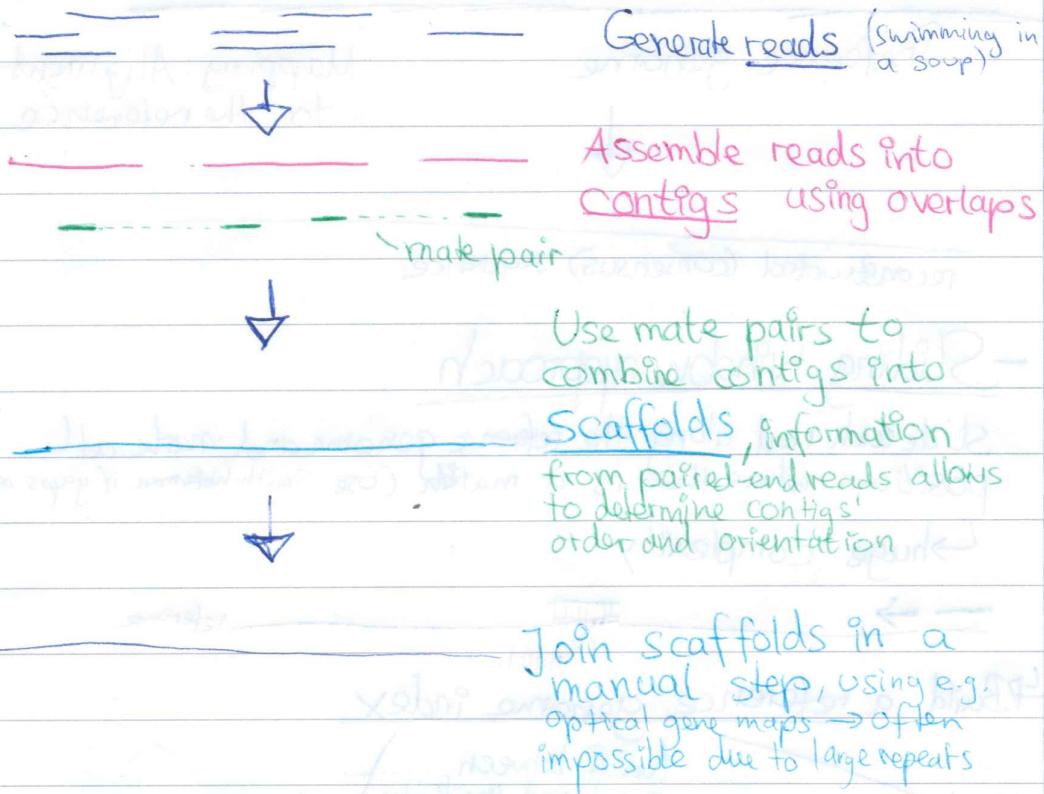
Short reads

Assembler

Reconstructed (consensus) sequence

# De novo assembly

- exploit read overlaps to build a (novel) genome sequence from scratch



- length of contigs and scaffolds is an important metric of assembly quality
- 2 approaches for de novo assembly/reconstructing sequences:

## Overlap graphs

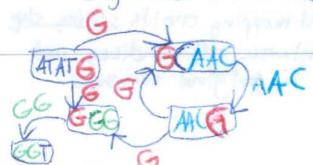
nodes: each read

edges: overlaps between reads, found by pair-wise alignments

wanted: hamiltonian path  
(visit each node exactly once)

- no efficient algorithm for finding hamiltonian path known, quadratic time complexity of pair-wise alignments

→ good if small number of reads with significant overlap, e.g. Sanger sequencing



## de Bruijn graphs

- each read is decomposed into k-mers

- each unique k-mer is represented by an edge in the graph

- nodes are (k-1)-mers: prefix + suffix of the k-mer associated with the edge connecting them

- wanted: eulerian path (visit each edge exactly once)

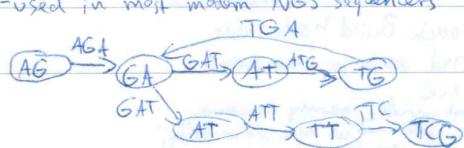
- efficient algorithm for eulerian path exists

- compact representation of repeats (duplicate k-mers within single node, longer repeats form a single series of adjacent nodes)

- Graph building time is linear in number of reads

- due sequencing errors require tip: bubble removal

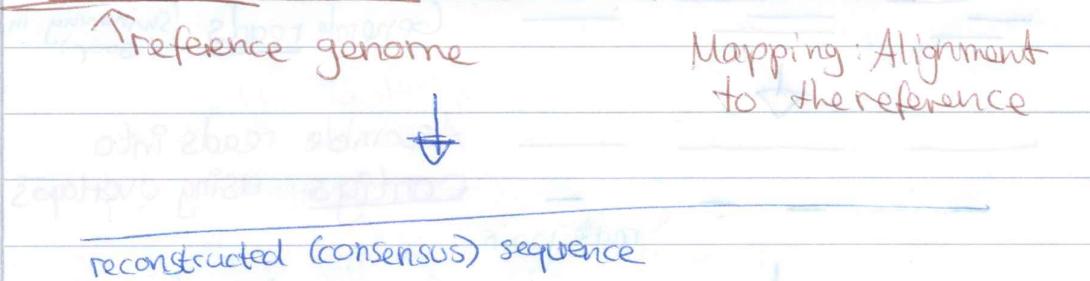
- used in most modern NGS sequencers



# By-reference assembly

- map reads to the known reference sequence - usually genome of the same species or a close relative / close relatives
- faster and less memory intensive than de novo

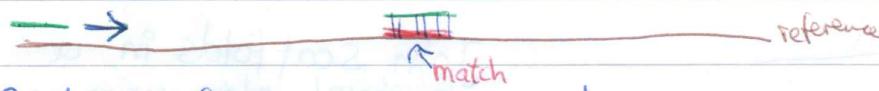
short reads



## Sliding Window approach

Slide each read along the reference genome and mark all positions where there is a match (use Smith-Waterman if gaps are allowed)

↳ huge complexity



## Build a reference genome index

### Hashing-based

use hash table to store positions of all  $k$ -mers in a genome,  $k \ll$  read size  
→ linked list of positions for each  $k$ -mer

- ① for each read, select one "proxy"  $k$ -mer (leftmost or middle part since better read quality)
- ② find all positions of this  $k$ -mer in the genome → seeds
- ③ for each seed, try to extend the alignment (map the rest of the read), with mismatches/gaps like in Smith-Waterman

### Variants & Optimizations:

#### Inexact matches with spaced seeds

↳ binary mask defines positions where mismatches are allowed  
↳ allows overlapping seeds with mismatches

#### Multiple $k$ -mers per read

↳ require at least  $n$  seed matches for a mapping location to be considered

#### Inverted approach: Build hash table from reads and search for $k$ -mers present in the reference

↳ smaller and variable memory requirements, depending on input read set size and variability; may be lower if very few reads are unique

tradeoff between  
speed and sensitivity  
→ heuristics  
using the index

### Burrows-Wheeler transform (BWT) - based,

- better in both time and memory usage; used for data compression

- ① Write down all cyclic rotations of the source string  $S$  (has end-of-line marker)
- ② Sort the rows lexicographically
- ③ Store the last column → BWT( $S$ )

- rows of the matrix form a sorted list of suffixes → allows efficient search for substring occurrences

- BWT is reversible → everything can be reconstructed from the last column

↳ trick: in order to obtain first column, sort last one, in order to obtain second column, sort last one + first one, ...

↳ rotate & sort

↳ any rows having  $P$  as a prefix are consecutive

- use FM-index to make matrix traversal efficient

↳ last-to-first mapping avoids sorting step

↳ mapping between suffix indices and positions in original sequence

# Multiple Sequence Alignment

13/36

Why?: compare more than two sequences at a time

↳ input for phylogenetic reconstruction

↳ discover important (conserved) parts of a protein family

↳ group of evolutionary related genes/proteins  
in different species with similar functions/structure

## ~~Orthology vs Homology~~

~~Orthologous sequences~~: Sequences in different species that have evolved from the same ancestral gene

↳ share a common evolutionary history

~~Homologous sequences~~: Sequences in different species that do not share a common evolutionary history, evolved independently

↳ convergent evolution (like bat wings, bird wings)

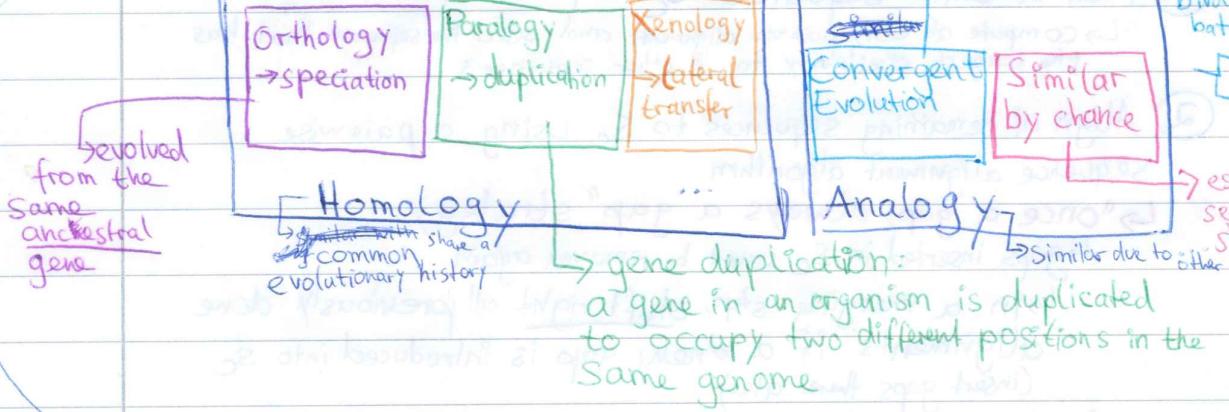
↳ similar by chance (if sequences are short)



Mice

Birds

↳ especially if sequences are short



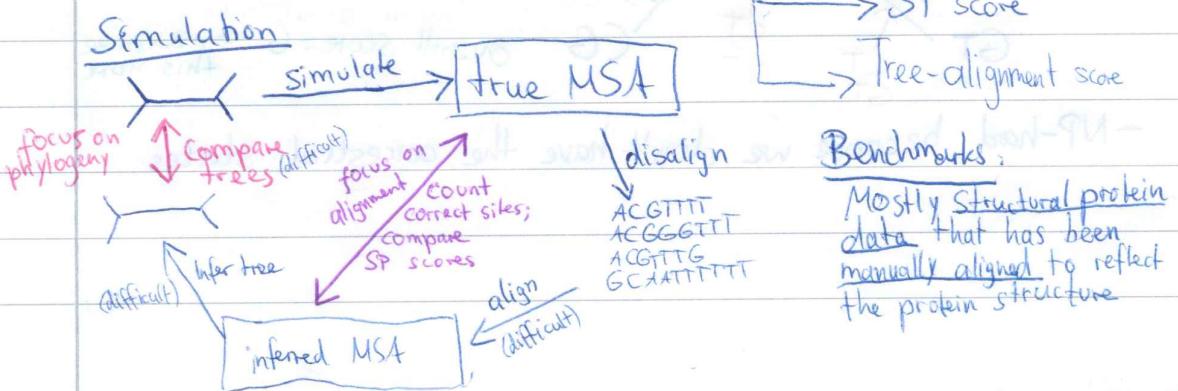
## Similarity

Given: n orthologous sequences of different length

Wanted: optimal alignment (allowed to insert gaps, but no columns/sites that consist entirely of gaps; corresponding homologous characters are aligned to each other)

- ~~Alignment Quality~~: unclear definition different criteria

↳ SP (sum of pairs) measure, real data benchmarks, evolutionary measures, simulations



→ SP score

→ Tree-alignment score

Benchmarks:

Mostly structural protein data that has been manually aligned to reflect the protein structure

# The SP-measure

- Sum of pairs score: sum of all pairwise scores between characters of a site; favor matches, penalize mismatches and gaps

example with edit distance for scoring:

s1: A A G A A - A A

s2: A T - A A T G B

s3: C T G - G - G C

$$p(-, -) = 0 \quad \text{no penalty}$$

for ~~two~~ gaps

$$2 + 2 + 2 + 2 + 2 + 2 + 3 = \underline{\underline{17}}$$

→ same as sum of pair-wise edit distances.

→ depends on choice of scoring function, not granted that it is biologically plausible

- NP-complete problem to compute SP-score ( $O(n^4)$ ) via DP-algorithm  
if all n sequences have length m

## - Star Alignment Heuristics

① Pick a center sequence  $S_c$

↳ compute all  $O(n^2)$  pairwise alignments and select the sequence that has the largest similarity to all other sequences

② Align all remaining sequences to  $S_c$  using a pairwise sequence alignment algorithm

↳ "once a gap, always a gap" strategy

↳ gaps inserted in  $S_c$  cannot be removed again

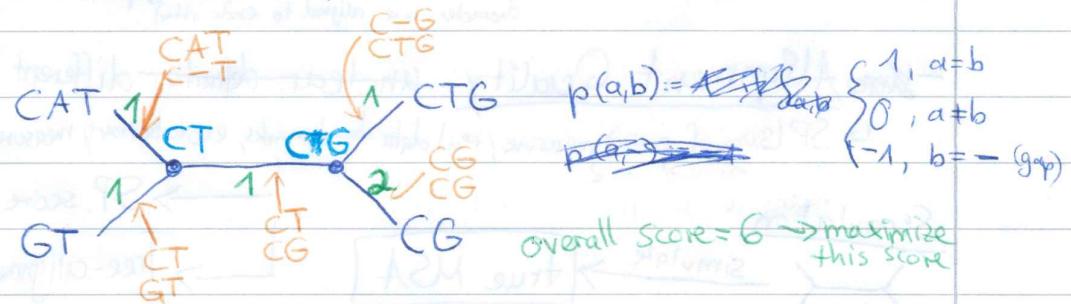
↳ in a merging step, shift right all previously done alignments if a new gap is introduced into  $S_c$  (insert gaps there also)

→ produces an MSA (multiple sequence alignment) whose SP score is < 2-optimum

- Tree Alignment: if an evolutionary tree for the sequences is available

↳ Find an assignment of sequences to the inner nodes such that the sum over the similarity scores on all branches is maximized

Example:

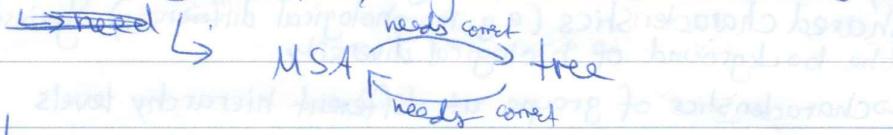


- NP-hard because we don't have the ancestral states

# Practical approaches for MSA

## ~~(Tree-Based Alignment)~~

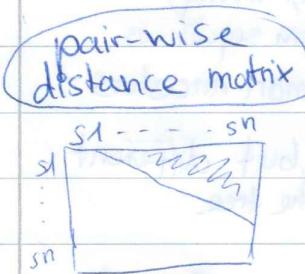
- Tree-Based Alignment is a hen and egg problem



↳ Simultaneous inference of tree and alignment

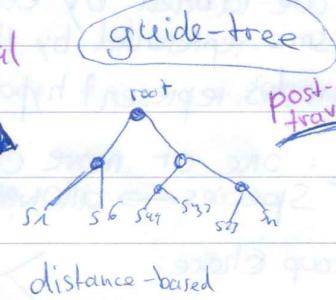
↳ NP-hard

### guide-tree based approach

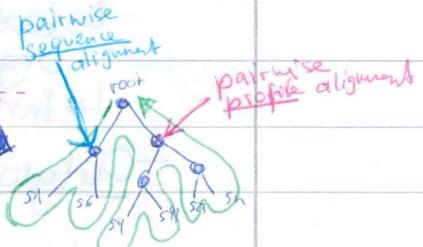


exactly  $O(n^2)$   
or approximate methods based on k-mers

hierarchical clustering  
(e.g. neighbor joining)



distance-based



bottom-up alignments:

- Sequence-sequence
- Sequence-profile
- profile-profile

## ↳ Progressive MSA!

Profile alignment: Merge alignments of the two descendant nodes

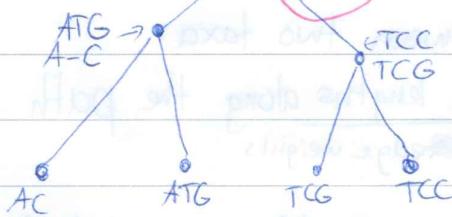
→ profile

profile alignment

Scoring  
weighted  
average over

all possibilities

leaves  
(leaves from left  
subtree  $\times$  leaves from  
right subtree)



- problem with progressive MSA: initial pair-wise

alignments are "frozen"  
↳ every error is propagated through to the final result

can't be corrected when new evidence emerges

## ↳ Iterative Progressive MSA:

↳ execute progressive MSA several times to refine the alignment,

e.g. MUSCLE Re-Finement: can return to previous alignments and change them

- Motif-based approaches: Find small motif (substring) common to all sequences

↳ shift sequences such that the motifs are "in alignment"

↳ align regions around these motifs using for instance progressive alignment

16/36

# Phylogenetics

Living

Taxonomy: classification of organisms by looking at shared characteristics (e.g. morphological differences) against the background of biological diversity

↳ characteristics of groups at different hierarchy levels

Phylogeny: refers to evolutionary history

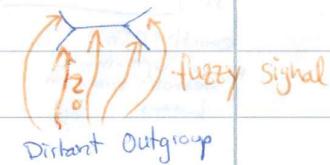
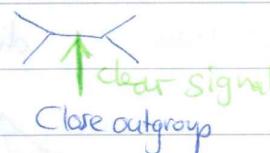
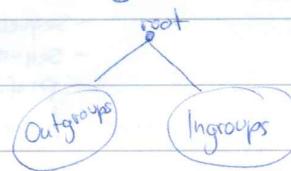
↳ unrooted strictly binary tree

↳ leafs are labeled by extant (currently living) organisms represented by their DNA/protein sequences

↳ inner nodes represent hypothetical common ancestors

- Outgroup: one or more closely related but different species → allows to root the tree

↳ Outgroup choice:



Taxon: denotes a subtree, one or more species that form a biological unit → unclear

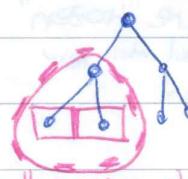
- in phylogenetics: leaf of a tree, a tree with n leaves has n taxa

Patristic distance between two taxa:

↳ sum of over branch lengths along the path in the tree  
↳ edge weights

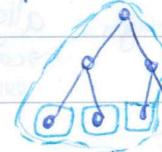
taking leafs:

monophyletic



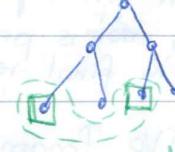
all species share a direct common ancestor, all species derived from that common ancestor are included

paraphyletic



all species share a common ancestor, not all species derived from that common ancestor are included

polyphyletic



species do not share an immediate common ancestor, excluding other members that would link them

# Applications of Phylogenetics

## - identifying unknown species

- ↳ phylogenetic placement for identifying anonymous sequences,
- ↳ phylogenetic species concept:  
↳ irreducible cluster, distinct from other such clusters, ~~with~~ with parental pattern of ancestry and descent
- Example: bird strikes
  - plane turbine
  - dead bird → which bird died?

## - viral outbreaks : trace back outbreak order

## - Snakebites

- ↳ antivenom same antivenom for closely related species
- interpret results from model organisms for other organisms

↳ species that is extensively studied / sequenced to understand particular biological phenomena, with the expectation that discoveries made for this one will provide insight to workings of other organisms

↳ selection criteria: easy to manipulate / grow / extract DNA, economical importance, less ethical problems, ...

↳ examples: Escherichia coli: gut bacterium, grows fast, cheap to culture

• Drosophila Melanogaster: fruit fly → breeds quickly

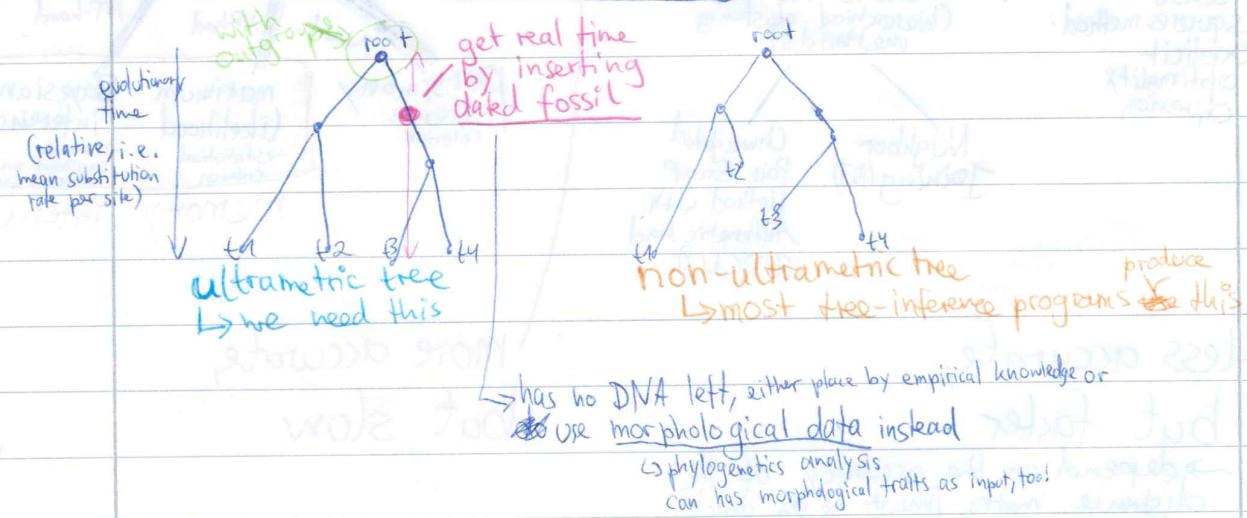
• Arabidopsis Thaliana: flowering plant → small genome

• mouse, rat, chimpanzee, ...

## - diversification rates: investigate extinction and speciation

## - divergence time estimates → for producing ultrametric trees

### ↳ needed in Tree Dating (dated trees)



18/36

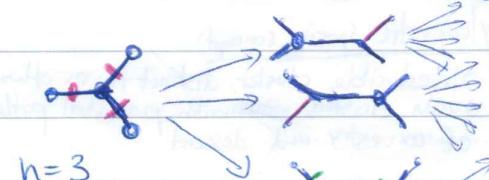
# Tree Inference

Tree search problem

Number of unrooted ~~binary~~ trees for  $n$  taxa: (unordered)

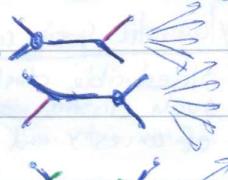
→ with node degrees ~~3~~ (unordered) → each vertex has either 1 or 3 neighbors

↳ trees that would be binary if we chose a root

 $n=3$ 

→ one tree

3 possible positions to insert new taxon

 $n=4$ 

→ 3 trees

5 possible positions to insert new taxon

$$\prod_{i=1}^{n-2} (2i-1)$$

→ die ersten  $n-2$  ungeraden Zahlen

↳ this number explodes



- number of rooted binary trees with  $n$  taxa

= number of unrooted ~~binary~~ trees with  $n+1$  taxa

- needs MSA, since alignment-free tree inference  
is typically less accurate → we have not established homology via a MSA

## Building Trees (Phylogeny reconstruction methods)

### distance-based methods

- use MSA to compute matrix of pair-wise distances
- build a tree using these distances

(least squares method:  
explicit optimality criterion)

NP-hard

heuristics  
(hierarchical clustering methods)

Neighbor Joining (NJ)

Unweighted Pair Group Method with Arithmetic Mean (UPGMA)

### character-based methods

- optimality criteria operate directly on MSA & tree
- take current tree topology & MSA to calculate a score tells us how well MSA data fits the tree

Parsony  
- discrete criterion

NP-hard

maximum likelihood  
→ statistical criterion

NP-hard

Bayesian inference  
→ integrate likelihood over entire tree space

Memory intensive

Less accurate,  
but faster

- depend on the accuracy of the distance matrix (must be an exact reflection of the tree)

More accurate,  
but slow

# Distance-based methods

19/36

## - Neighbor Joining : $O(n^3)$

Principle: Given distance matrix of size  $n \times n$ .

① Find Minimum

② Merge Taxa

③ Compute new distance matrix of size  $(n-1) \times (n-1)$

④ Goto ① if still more than 1 taxa there.

$O(n^2)$  space,  $O(n^3)$  time

③a: Compute average distances  $u_i = \sum_j D_{ij} \cdot \frac{1}{n-2}$

③b: Q-Matrix:  $D_{ij} - u_i - u_j \leftarrow$  find minimum  $i, j$ ,  
create ancestor & tree

$$\frac{1}{2}D_{ij} + \frac{1}{2}(u_i - u_j) \quad i \quad j \quad \frac{1}{2}D_{ij} + \frac{1}{2}(u_j - u_i)$$

③c: Update distance matrix:

$$D_{x,k} = (D_{ik} + D_{jk} - D_{ij}) / 2$$

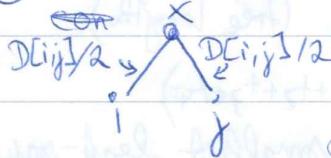
→ does not produce ultrametric trees.

## - UPGMA: Unweighted Pair Group Method with Arithmetic Mean: $O(n^2)$

→ can be used if we know that we have an ultrametric tree, usually not the case (for instance higher evolutionary pressure → difficult life conditions)

→ produces rooted, ultrametric trees → misleading results if trees not ultrametric  
↳ lineages with shorter lengths are grouped

① Find Minimum, merge nodes  $i, j$ ,



$k \in \{i, j\} \cup$

number of elements in node  $i$

$$D[(i,j),k] = \frac{n_i}{n_i + n_j} \cdot D[i,k] + \frac{n_j}{n_i + n_j} \cdot D[j,k]$$

② Update distances:  $D[(i,j),k] = \frac{n_i}{n_i + n_j} \cdot D[i,k] + \frac{n_j}{n_i + n_j} \cdot D[j,k]$

$O(n^2)$  via list of per-column (or per-row) minima

update list:  $O(n)$ , look for minimum:  $O(n)$

↳ don't need to update the entire matrix each time

20136

## Least Squares:

Given: fixed, fully binary, unrooted tree  $T$  with  $n$  taxa  
 • pair-wise distance matrix  $D$

Wanted: assign branch-lengths to the tree such that sum of differences between pair-w

$$Q = \left[ \sum_{i < j} (D_{ij} - d_{ij})^2 \right]^{\frac{1}{2}} \min \quad \square$$

squared

plain pair-wise distances      patristic (tree based) distances

We can also weight this if we want, e.g. with  $w_{ij} := \frac{1}{D_{ij}^2}$  or  $\frac{1}{d_{ij}^2}$

-  $Q$  can be minimized by taking derivative and solving system of linear equations in  $O(n^3)$   
 ↳ with regard to tree-like structure of  $Q$ :  $O(n^2)$  or even  $O(n)$

↳ NP-hard to find the minimum least-squares tree because there are too many possible tree topologies!

Minimum Evolution Method: similar to least squares

→ obtain branch lengths via least-squares method

→ minimize total branch length (tree length)

of the reconstructed tree ( $t_1 + t_2 + t_3 + t_4$ )

↳ do not return the tree with the ~~smallest~~ least-squares score, but the one with smallest total length out of these.

↳ NP-hard, too.

# Character-based methods I

Parsimony: explain the data with the least amount of mutations

Given:

- Multiple sequence alignment
- tree topology

Wanted: assignment to inner nodes such that the number of mutations on the tree is minimized

→ similar to tree alignment problem, but here we already have to multiple sequence alignment

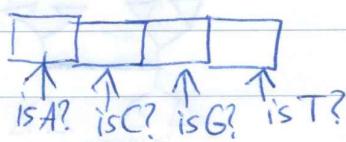
→ gaps (also called indels) are treated as undetermined characters (can be A,C,G or T; like "don't-care's" in TI lecture)

① Insert Virtual root somewhere in the tree → makes computation easier, score stays same  
 → we assume that evolution is time-reversible (occurred in the same way forward or backward in time)

② Post-order traversal to compute inner states

③ Compute score on a site-per-site basis  
 → assume that sites evolve independently

- Trick: A=1000, C=0100, G=0010, T=0001



bit-wise AND  
 ↳ if  $\neq 0000 \rightarrow$  intersection is non-empty  
 ↳ no mutation  
 ↳ else: mutation, count +1

ancestral states are bit-wise OR

Since there are  $n^2$  inner nodes when in tree  
 ↳ number of sites

↳ union  
 ↳ defer decision which letter to take to a later point

$O(n \cdot m)$  space & time-complexity,

but with small constant  
 dynamic programming algorithm

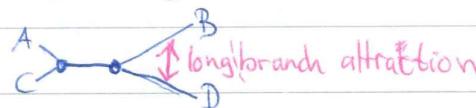
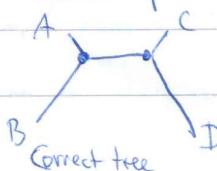
→ NP-hard if tree not given

→ returns discrete scores

↳ implementation trick:  
 AVX instructions & number of bits set to 0/1 as built-in operation

## Long Branch Attraction

↳ because parsimony tries to minimize the number of mutations it faces  
 some problems on trees with long branches



Wrong tree inferred by parsimony

22/36

# I. Phylogenetic tree search algorithms

## Tree-Search Algorithms

- Obtaining initial starting tree with  $n$  taxa

↳ NJ or UPGMA tree → only 1 tree

↳ random tree

↳ start with 3 randomly chosen sequences as leaves

↳ for each other sequence (randomly chosen) chose randomly a branch of the current tree to add it there

↳ stepwise addition algorithm (randomized)

↳ for each addition, chose branch that leads to best parsimony insertion score

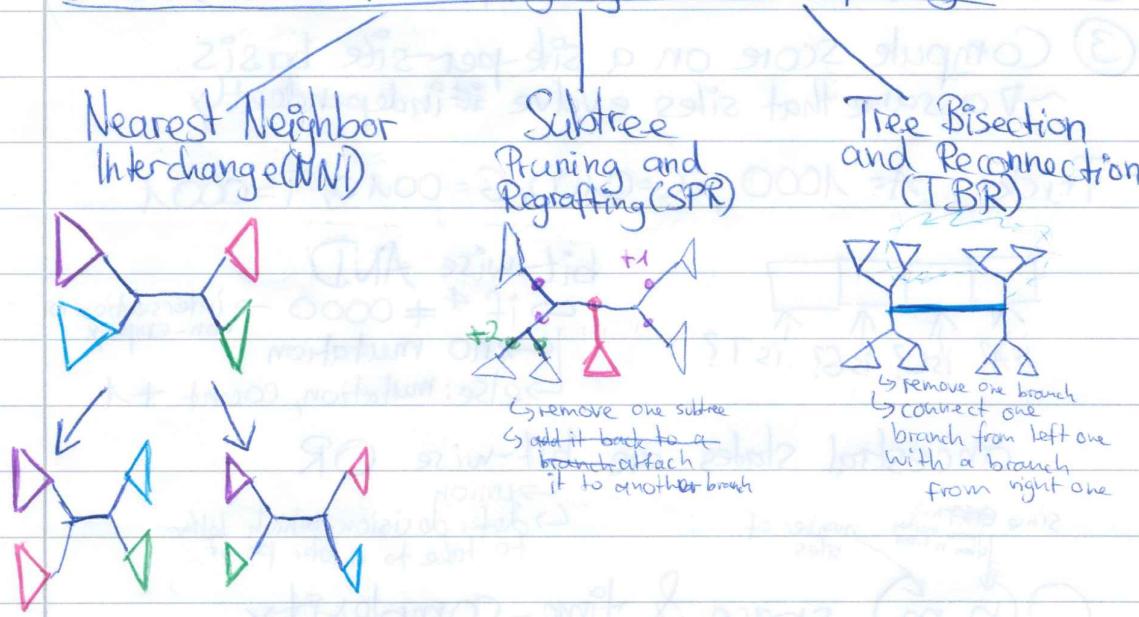
↳ addition order is randomized

↳ distinct starting trees are useful because of local maxima vs. global maxima

↳ reasonable starting trees: lead faster to a good optimality score (comprehensive) than random starting trees

## Tree Rearrangement

### (Basic Moves for changing a tree topology)



## Parsimonator Algorithm

↳ build randomized stepwise addition order parsimony tree

↳ apply SPR-moves to all current subtrees within rearrangement radius of 20

↳ greedy: if an SPR-move leads to a better score, keep it immediately

# Maximum Likelihood

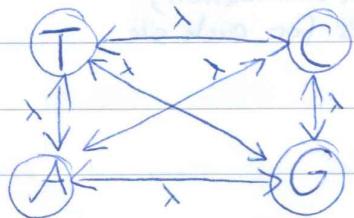
## Distances as stochastic processes: (Models of Sequence Evolution)

↳ We can have different ways to come from  $x$  to  $y$ :

$$\begin{array}{l} A \rightarrow A \\ A \rightarrow C \rightarrow G \rightarrow A \rightarrow C \rightarrow A \\ A \rightarrow T \rightarrow G \\ \vdots \end{array}$$

- Assumptions:
- Markov process: Probability of  $i \rightarrow j$  does not depend on what was before  $i$
  - Time-continuous: Substitutions can happen at any time
  - homogeneous: substitution probabilities do not change in different parts of the tree
  - stationary: Overall character frequencies  $\pi_i$  of the nucleotides are at equilibrium and remain constant

### Jukes-Cantor (JC)

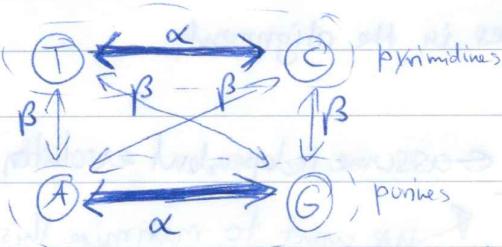


$$\lambda_{ij} = \lambda_{ji}$$

$$Q = \begin{pmatrix} -3\lambda & \lambda & \lambda & \lambda \\ \lambda & -3\lambda & \lambda & \lambda \\ \lambda & \lambda & -3\lambda & \lambda \\ \lambda & \lambda & \lambda & -3\lambda \end{pmatrix}$$

$$\pi = \left( \frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4} \right)$$

### Kimura (K)

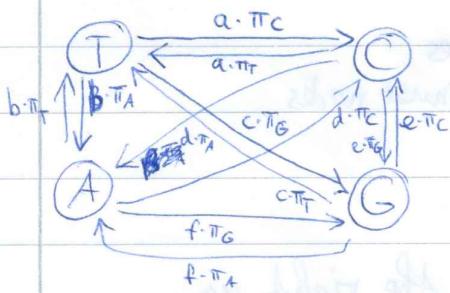


$$Q = \begin{pmatrix} -(\alpha+2\beta) & \alpha & \beta & \beta \\ \alpha & -(\alpha+2\beta) & \beta & \beta \\ \beta & \beta & -(\alpha+2\beta) & \alpha \\ \beta & \beta & \alpha & -(\alpha+2\beta) \end{pmatrix}$$

$$\alpha: \text{transition rate}$$

$$\beta: \text{transversion rate}$$

### General Time Reversible (GTR)



$$\pi = (\pi_A, \pi_C, \pi_G, \pi_T)$$

Time reversible:  $\pi_i \cdot \lambda_{ij} = \lambda_{ji} \cdot \pi_j$

Mutation in each direction is equally likely:  $p(i \rightarrow j) = p(j \rightarrow i)$

Likelihood computation does not depend on where the tree is rooted

$Q$ : Substitution rate matrix (Transition Rate Matrix)  
↳ Each row sums to zero → provides substitution rate between pairs of nucleotides for an instant time  $t$

transition matrix

$\pi$ : stationary distribution (equilibrium frequencies)  
↳ sum up to 1

$t$  is a measure of branch length, not time!

$$P(t) = e^{Q \cdot t}$$

$P_{i,j}(t) \triangleq$  probability to end in state  $j$  after time  $t$ , starting from state  $i$ ,  $\triangleq$  probability of

after time  $t$ , starting from state  $i$ ,  $\triangleq$  probability of



Trick:  $P(t) = e^{Q \cdot t} = e^{U \cdot \Lambda \cdot U^{-1} \cdot t} = U \cdot e^{\Lambda \cdot t} \cdot U^{-1} = U \cdot \Lambda^k \cdot U^{-1} \rightarrow$  easy to compute

$U$  is diagonalizable  
↳ eigenvalues of  $Q$  are eigenvalues of  $U$

↳ spectral decomposition

eigenvectors of  $Q$  as columns

14/36

# Topic 10: Phylogenetic Models

## Rate Heterogeneity

↳ different mutation rates at different sites

$$P(t) = e^{Q \cdot u_i \cdot t}$$

↳ draw these rates from a gamma distribution

↳ introduces 1 additional parameter (shape  $\alpha$ )

↳ draw 4 rates  $f$

↳ Gamma model: draw 4 rates from gamma distribution, for each site average over these 4 rates for each site

## Protein Substitution Models

↳ too many free parameters

↳ use precomputed empirical models

↳ count substitutions ~~in alignments~~

in alignments of closely related sequences (no  $A \rightarrow C$  happens)

## Obtain base frequencies:

↳ count the occurrences of states in the alignment

## Likelihood Function ← assume independent evolution

$$LH(T | D) = P(D | T)$$

$$\text{Tree data} = \prod_{S_i \text{ sites}} P(S_i; T)$$

$$\log LH(T | D) = \sum_{S_i \text{ sites}} \log(P(S_i; T))$$

Given: tree topology, branch lengths, leaves

↳ sum over all possibilities for inner nodes

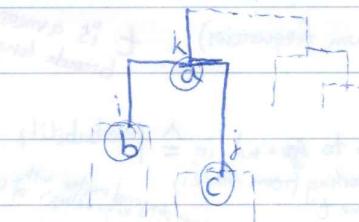
to get the ~~per~~ likelihood

↳ we want to maximize this  
 (find ~~per~~ best parameters  
 → Substitution rates,  
 equilibrium distribution,  
~~gamma~~ shape  
 branch lengths)

## Felsenstein-pruning-algorithm:

↳ move the summation signs to the right as far as possible (place them directly before the related term)

↳ dynamic programming algorithm, computes the partial likelihoods bottom-up by post-order-traversal of the tree



$$P(T | \text{node at } k \text{ is } a) = \sum_b P(T_i | \text{node } i \text{ is } b) \cdot P(k | i)$$

$$\cdot \sum_c P(T_j | \text{node } j \text{ is } c) \cdot P(l | j)$$

$$\text{Likelihood of sequence } A-U-A-U-A-U-A = \dots = S = [T]S = [T]S - [T]T$$

## Optimizing Parameters

- Continuous optimization
- (1) Optimize Parameters of the substitution model
  - (2) Optimize branch lengths of the tree
    - ↳ fix (1), do Newton-Raphson method
      - ↳ typically with iterative method, each branch is optimized separately since  $P(t) = U_1 A_1 e^{-t} U_1^{-1}$  usw.
  - (3) Search over all possible tree topologies
    - ↳ difficult (root & discrete optimization)
    - ↳ do tree rearrangement (NNI, SPR, TBR) → Iterated Local Search

## Discrete Operations on Trees

- tree sets

↳ trees with same parsimony score

↳ different trees for different genes / from different databases

↳ ...

## Phylogenetic bootstrap

↳ we want support values for each branch

↳ supported by  $\times\%$  of the data

- evolution: stochastic process that yields a multinomial

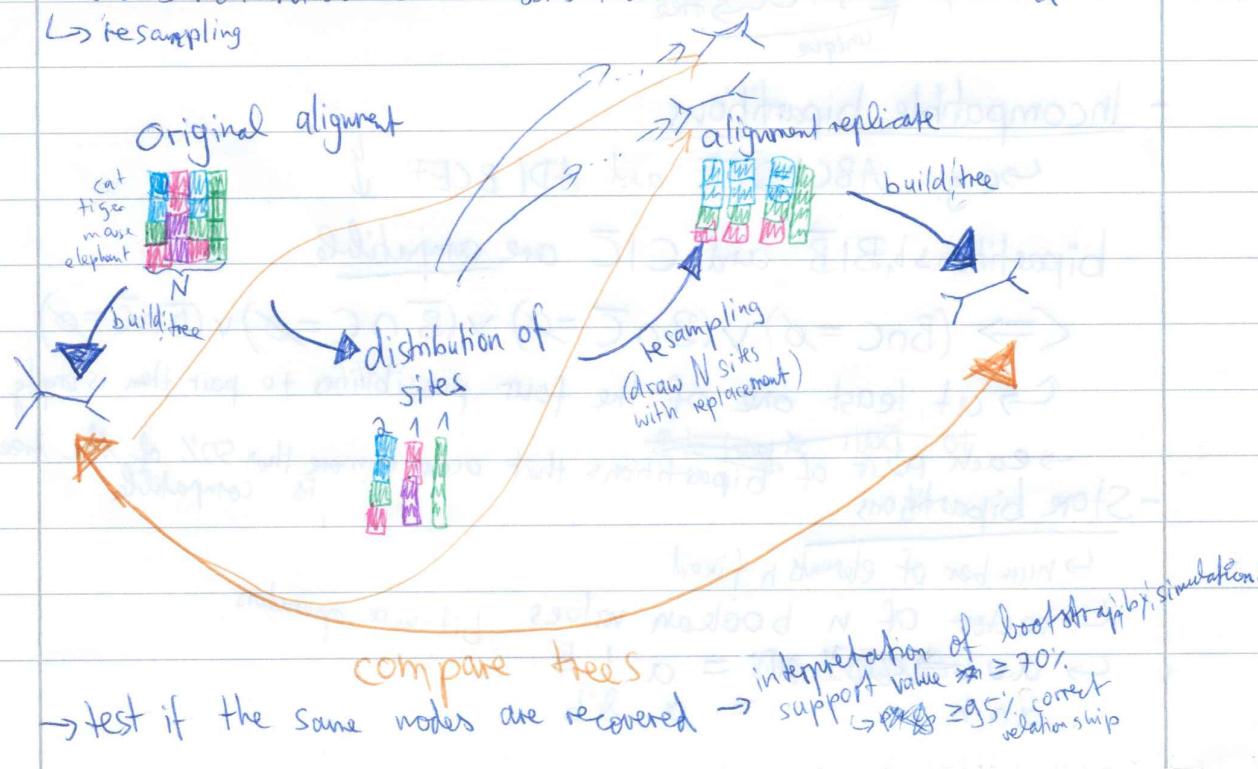
- distribution of evolutionary events (→ alignment sites)

- bootstrap new replicate alignments from original alignment

↳ alternative version of the original alignment

- does not introduce new sites that have not been observed

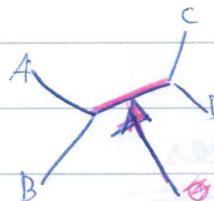
↳ resampling



16/36

## ML tree with bootstrap support values

- ① Infer ML tree on full alignment
- ② Create bootstrap tree set
- ③ extract bipartitions from ML tree and bootstrap tree set
- ④ annotate ML tree with relative frequency of its bipartitions in the bootstrap tree set



How many percent of the bootstrap trees do have this bipartition?

## Bipartitions

→ inner branches are bipartitions

AB|C|DEF → taxa A, B more closely related to each other than taxa CD, EF

→ outer branches are trivial bipartitions

A|B|C|D|E → contained in every tree with taxa ABCDE

→ not informative  
- set of bipartitions ↗ sometimes no tree  
    ↳ tree unique

### - Incompatible bipartitions

↳ e.g. ABC|DEF and AD|BCEF

bipartitions  $\overline{B}, \overline{B}|\overline{B}$  and  $\overline{C}, \overline{C}|\overline{C}$  are compatible

$$\Leftrightarrow (\overline{B} \cap C = \emptyset) \vee (\overline{B} \cap \overline{C} = \emptyset) \vee (\overline{B} \cap \overline{C} = \emptyset) \vee (\overline{B} \cap \overline{C} = \emptyset)$$

↳ at least one of the four possibilities to pair them is empty.

→ each pair of bipartitions that occur in more than 50% of the trees is compatible

### - Store bipartitions

↳ number of elements n fixed

↳ vector of n boolean values bit-wise operators

$$a \vee b = \text{bitwise OR} = a | b$$

$$a \wedge b = \text{bitwise AND} = a \& b$$

## Extract bipartitions from tree

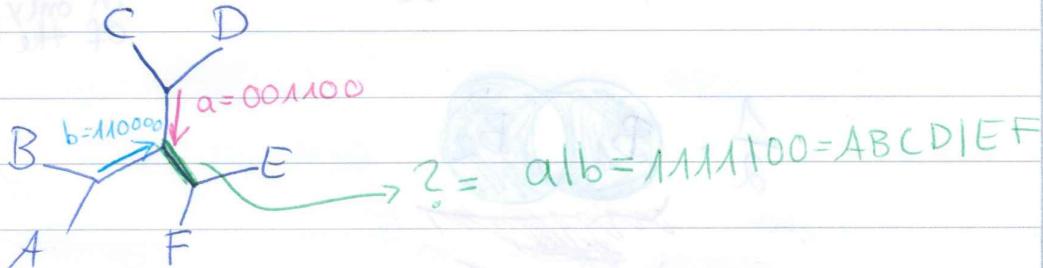
↳ via hashing → store bipartitions in hashtable

↳ post-order traversal

↳ compute bipartition from adjacent bipartitions

↳  $r[]$ : one random number per taxon

↳ hash: xor-value of all corresponding random numbers in array



$$\text{hash}(a) = r[C] \text{ xor } r[D]$$

$$\text{hash}(b) = r[A] \text{ xor } r[B]$$

$$\text{hash}(?) = \text{hash}(a) \text{ xor } \text{hash}(b)$$

Consensus Tree → extract bipartitions from tree set, determine, count percentage of occurrences

majority-rule consensus tree  
only bipartitions with freq > 50%

strict majority-rule consensus tree  
only bipartitions with freq = 100%

extended majority-rule consensus tree

- ① compute majority-rule consensus tree
- ② greedily refine tree with non-consensus bipartitions, starting with most frequent one:  
- check if can be added, if so, do so, else proceed to next one  
Optimal solution is NP-hard

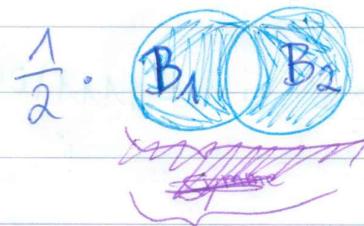
## Robinson-Foulds Distance (RF distance) between trees

Given:

- Unrooted tree  $T_1$  with bipartitions  $B_1$
- Unrooted tree  $T_2$  with bipartitions  $B_2$

$$RF(T_1, T_2) = \frac{|B_1 \cup B_2| - |B_1 \cap B_2|}{2}$$

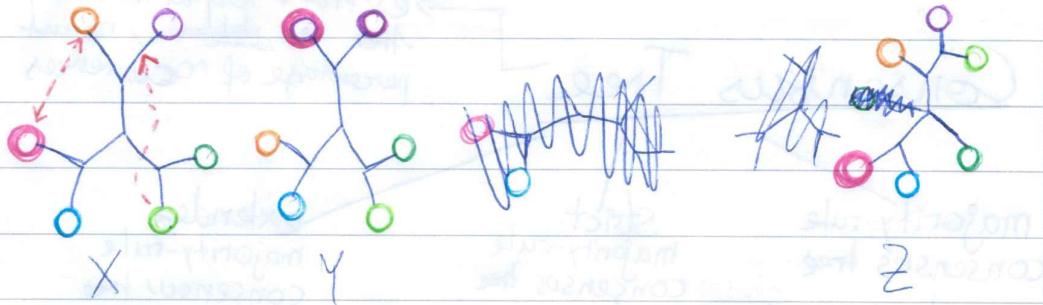
number  
of bipartitions  
that appear  
in only one  
of the trees



- RF-distance is a metric

- The RF-distance is not ultrametric:

↪ it does not hold that  $d(x, z) \leq \max(d(x, y), d(y, z))$



$$d(X, Y) = 2$$

$$d(X, Z) = 2$$

$$d(Y, Z) = 3$$

absolute RF-distance:  $\frac{B_1 \Delta B_2}{2}$

relative RF-distance:  $\frac{B_1 \Delta B_2}{2 \cdot (n-3)}$

# MC MC Methods

## - Difference between ML and MCMC

In common:- both re-evaluate the phylogenetic likelihood over and over again for different tree topologies, branch lengths and model parameters  
 ↳ Spend most of their total running time on likelihood calculations on trees

### Difference:



## Bayes Theorem

$$P(\text{Tree, Params} \mid \text{Alignment}) =$$

$$P(\text{Alignment} \mid \text{Tree, Params}) \cdot P(\text{Tree, Params})$$

$$\frac{P(\text{Alignment})}{\text{prior probability}}$$

Posterior Probability

marginal probability  
 → ouch!

Monte Carlo: pick random points, average over them, to app

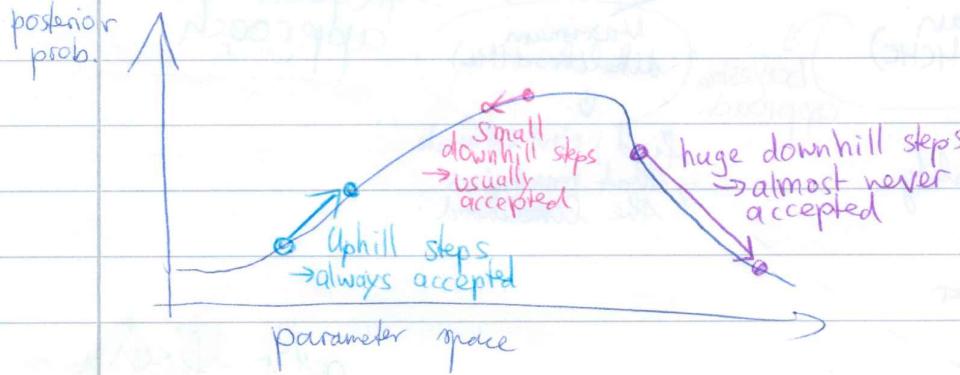
- Monte Carlo Methods: randomly sample data-points in the huge parameter space to approximate the interval

## -Markov Chain Monte Carlo

↳ we are interested in areas of high posterior probability → get there faster

↳ biased random walks → spend more time in interesting regions

## Robot Metaphor



accept/reject proposals: from point 1 to point 2

acceptance ratio depends on

$$\frac{P(\text{Point2}|\text{data})}{P(\text{Point1}|\text{data})}$$

ratio of posterior densities  
marginal probability cancels out

$$\frac{P(\text{point 2})}{P(\text{point 1})}$$

prior ratio

$$\frac{P(\text{data}|\text{point 2})}{P(\text{data}|\text{point 1})}$$

likelihood ratio

$$\frac{Q(\text{point 1}|\text{point 2})}{Q(\text{point 2}|\text{point 1})}$$

correction factor

hastings ratio

needed if proposal distribution is asymmetric

$Q \triangleq$  proposal probability

ensures that the probability going from point1 to point2 is the same as the probability going from point2 to point1, regardless of the proposal distribution

## Target distribution

posterior distribution we are trying to sample (integrate over)

## proposal distribution

↳ decides which point (how far/close) in the landscape to randomly go to/take next

good mixing: find a proposal distribution that requires few steps until MCMC chain converges

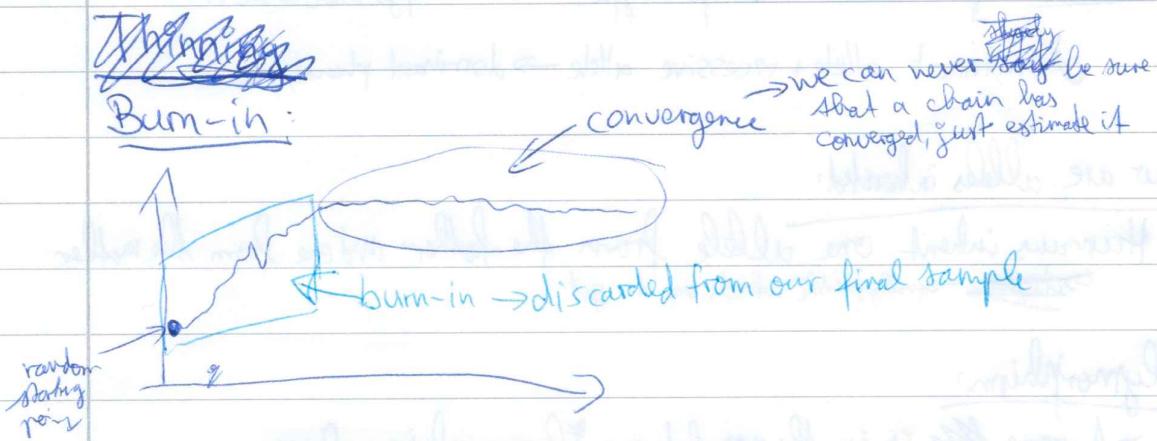
↳ balance pros and cons

→ very small steps; less rejections

→ larger steps; more rejections

# Metropolis-Coupled Markov Chain Monte Carlo

- ↳ multiple Markov Chains run simultaneously
- ↳ Cold Chain: Sees landscape as it is → ①
- ↳ Hot Chain: Sees a melted/flatter version of the landscape → several
  - moves more easily between peaks → larger steps
- only the cold chain samples
- cold robot can swap positions with the hot ones if the hot robots found some better places
- easier to escape local maxima



## MCMC in practice

- ① Start with random tree with random branch lengths and random substitution parameters
- ② Iterate using the following steps:
  - proposed change to the tree
  - Propose either
    - a new tree topology
    - a new branch length
    - ~~new~~ a new substitution parameter and re-calculate the likelihood
  - calculate acceptance ratio of the proposal, accept change or reject it
  - print the current tree with current branch lengths to file **only** every k (e.g. 1000) iterations
    - ↳ generates a sample from the chain, avoids writing TB's of files, makes samples be

Summarize samples: using means, histograms, credible intervals, Consensus Trees, ...

# Population Genetics

Biological

Evolution: Change of the frequency of occurrence of features of individuals of the population should be inherited

## 4 Evolutionary forces

↳ Genetic Drift: Chance (other than a random mutation)

↳ Migration

↳ Mutation

↳ Natural Selection: Response to some pressure  
(e.g. antibiotics, climate change)

Genotype: Full hereditary information

Phenotype: ~~actual~~ observed properties → expression

Dominant: produces a dominant phenotype → one copy is enough

Recessive: produces a recessive phenotype → both copies have to be it

Dominant allele + recessive allele → dominant phenotype

How are alleles inherited:

Humans inherit one allele from the father and one from the mother  
~~decided~~ - equally probable which one to get

Polymorphism:

- A gene ~~is~~ is in the population polymorphic when there exist multiple alleles

SNP: Single Nucleotide Polymorphisms

↳ polymorphic sites in the DNA of individuals  
(have more than one nucleotide state)

# Hardy-Weinberg Model

## Assumptions

- infinite population size
- random mating, no natural selection, no migration
- diploid population
- gene A with 2 alleles: A and a
- Current frequencies (at generation 0) of allele pairs defining the genotype

~~Aa~~

$$f(AA) = x$$

$$f(Aa) = 2y$$

$$f(aa) = z$$

- Current frequencies of gametes (single alleles) at generation 0:

$$f_0(A) = x + y$$

$$f_0(a) = z + y$$

- Sexes have same distribution of the 3 genotypes AA, Aa, aa

		Females	
		A(p)	a(q)
Punnett Square	A	AA	Aa
	a	Aa	aa

offspring genotype possibilities

	A(p)	AA ( $p^2$ )	Aa ( $2pq$ )
	a(q)	Aa ( $q \cdot p$ )	aa ( $q^2$ )

$$p^2 + 2pq + q^2 = 1$$

↳ if we filter out evolutionary forces ~~like~~ and assume infinite population size, the frequencies of the alleles stay the same.

↳ Hardy-Weinberg Equilibrium

↳ If ~~we~~ we see change in allele frequencies between parents and offspring → evolution happened

## Wright-Fisher model

- ↳ populations are of finite size  $2N$
- ↳ non-overlapping generations (e.g. annual plants)
- ↳ constant population per generation
- ↳ each individual from offspring generation picks a parent at random ~~from~~ from previous generation, uniformly distributed, each parent can be picked more than once
- ↳ each offspring inherits information from the parent

### → binomial random sampling

↳ probability to pick allele  $A$  as parent is  $\frac{\#A}{2N}$ .

↳ binomial distribution of alleles

$$\binom{n}{k} \cdot p^k \cdot (1-p)^{n-k}$$

Absorbing state:

↳ probability to escape this state is zero

→ if allele frequency is 0 or 1.

↳ fixation: probability to end up in  $f(A)=1$

$$\hookrightarrow = \frac{\#A}{2N}$$

Random genetic drift:

Change in allele frequencies over generations in finite populations due to re-sampling

Heterozygosity:

Probability that two alleles are different

$$Het_t = Het_{(t-1)} \cdot \left(1 - \frac{1}{2N}\right) \quad \hookrightarrow Het_t = Het_0 \cdot \left(1 - \frac{1}{2N}\right)^t$$

Probability  
that two randomly  
chosen alleles are  
different

↳ expected probability that an individual will be heterozygous → proportion of individuals heterozygous at a locus  
↳ measure for gene diversity from original proportion

$$\hookrightarrow -\frac{1}{2N} \cdot Het = \text{loss of heterozygosity per generation due to genetic drift}$$

## Positive selection:

if natural selection has favored allele A  
 ↳ reproduces more often

## Mutation-drift balance

↳ genetic drift removes polymorphisms

↳ mutation introduces polymorphisms

rate of evolution: the probability of a new mutation to arise in the population and to eventually become fixed

p: probability of mutation per generation and individual

2N individuals → 2Np mutations per generation

— probability that a mutation will be fixed is  $\frac{1}{2N}$

↳ rate at which a mutation will arise and fix is

$$\frac{1}{2N} \cdot 2Np = p$$

→ we want  $\Delta Het = 0$

gain in heterozygosity due to mutation:

$$\hookrightarrow 2p \cdot (1 - Het)$$

$\underbrace{2p}_{\text{Mutation rate}}$   $\underbrace{(1 - Het)}_{\text{Homozygosity}}$

at least one of the two mutations

$$Het_{t+1} = \underbrace{1 - \frac{1}{2N}}_{\text{initial}} \cdot Het_t + \underbrace{2p \cdot (1 - Het_t)}_{\text{due to mutation}}$$

## Phylogenetics vs

## Population genetics

Study of Polymorphisms  
in a population → why? how? how can we get away?

- ↳ given data (MSA of individuals of the same species)
- ↳ simulate evolution under different scenarios
- ↳ compare if one of the scenarios fits the Summary statistics (e.g. # SNPs) of our empirical dataset
  - help us to assess if positive selection occurred, help, to allow parameter estimation of population
- ↳ Evolution at different scales
  - ↳ less mutations → correction of MSA more important

↓ inheritance of a trait

How the frequency of alleles which control the trait change over time

