**trivial**

**Hamming Distance:** $d_H$
Number of mismatches.

**Median String:**
$APA..A \to TT - 7$
find best hamming distance

$M = |P|$   Brute force: $O(m \cdot n)$
$n = |t|$   $KMP = O(m+n)$

**KMP Match:**
```
F = failureFunction(P)
i = 0
j = 0
while i < n
  if T[i] = P[j]
    if j = m-1
      return i-j
    else
      i = i+1
      j = j+1
  else
    if j > 0
      j = F[j-1]
    else
      i = i+1
      j = 0
return -1
```

**KMP failure:**
```
F[0] = 0
i = 1
j = 0
m = len(p)
while i < m
  if P[i] = P[j]
    F[i] = j+1
    i = i+1
    j = j+1
  else if j > 0
    j = F[j-1]
  else
    F[i] = 0
    i = i+1
```

**Suffix Trees:**
Collapsed keyword trees:

Let $s = abab$, suffix $(x,s)$
tree of $s$ is a
compressed tree of
all suffixes of $s = abab$
which are
$(abab$, $bab$, ...$)$ add $ to find

**UPGMA:**
$x_i \to C_i$ assignment (one for each x)

Find $C_i$ and $C_j$ such that $d_{ij}$ is min.

Let $C_k \in C_i \cup C_j$
Add vertex connecting
$C_i, C_j$ heighted $\frac{d_{ij}}{2}$
delete $C_i, C_j$
repeat until 1
cluster remains.

---

**Boyer Moore:**
**Bad Char:** Last occurrence of each char.

**Good Suffix 1:** if $t \in P$ matched with text, find last occurrence at $T$ left side.

**Good Suffix 2:** find $t'$ left side.

**Fingerprint:** $\theta$
$f(P) \Rightarrow O(m)$ $h = f \mod q$
**Rabin Karp:** $O(n)$
Hash pattern. Hash text's first m chars and compare. $1-m, 2-m \ldots$ if same, compare strings.

Select $q > m$, $q \to$ prime

**Finite Automata:** linear good $O(n)$
States: $q_1, 2, \ldots, m$ preprocess $\to O(m|\Sigma|)$
Memory $\to$ bad

**Bitap Algorithm:** Shift-And $O(mn)$
$M(i,j) = 1$ iff $P[1\ldots i] = T[j-i+1\ldots j]$
$m$ rows $\}$ M   Bitshift$(a) \to$ shift $a$ to right by adding 1 to first
$n$ cols $\}$   bit
array: 1 0 0 1 0 0 0 0 0 0 0
$M(j) = Bitshift(j-1) \wedge U(T(j))$
Define U for each char in alphabet.
if last row has 1: found match

**Pattern Matching:** $O(nm)$
**Suffix Trees:** $O(n)$   Let $N$ label n's
Multiple Pattern Match $\to$ brute force $O(km \cdot n)$
keyword tree $\to O(N)$ + Naive thread $\to O(N+nm)$
$\downarrow$   or
Construct using   +Aho-Corasick $\to O(N+m)$
patterns
and thread over text

**Weiner Algorithm** can do it in linear time. put these into tree.

**Aho-Corasick:** Search in keyword trees!
**Failure links:**
$L(v) \to$ the word that is the concatenation of the chars until node $v$ from 0.
$lp(v) \to$ longest suffix of $L(v)$ which is a prefix of some word starting from failure link from "0" node
$\to O(n)$ failure links

**Star Alignment:**
Sim $\to$ score of pairs
Compute sim for every pair $(i,j)$
$star\_score(i) = \sum_{j \in P} sim(i,j)$
$max(star score) \to$ center pattern
align according to it.

$d_{ic} + d_{jc} = D_{ij}$
$d_{ic} + d_{kc} = D_{ik}$
$+ d_{jc} + d_{kk} = D_{jk}$
Using Math:
$d_{ic} = \frac{D_{ij} + D_{jk} + D_{ik}}{2}$
can be calculated.

n-ease tree
has 2n-3 edges.
if n > 3 this may not be solvable!
not sure!

---

**Burrows-Wheeler Transformation**
b2|t2
Generate all rotations of word. Sort them. $\$ < $ all of alphabet
Output last column!
Last-to-first char map
$n^{th}$ "x" on left is also $n^{th}$ "x" char on right!

$occ(j, 'c') \Rightarrow$ number of occurrences of char 'c' until jth pos (inclusive) (table), $cnt('c')$ table contains occurrence number of each char, $rank('c')$ first occurrence of each char in first word

**Edit Distance:** $d_E$ [non-trivial]
Minimum number of elementary operations to get $S_2$ from $S_1$.

**Multiple Alignment:**
Frequencies for each column
$P_A, P_G, P_C, P_T, P_{(gap)}$ char col
$P_x = \frac{\# of A's}{\# of rows} \to P(j,i)$

$S(x,i) = \sum_j (\delta(x,y) \times p(y,i))$

char col on sequence

**Profile**

|   | 0 | -1 | -2 | -3 | -4 | -5 | -6 | -7 | -8 |
|---|---|----|----|----|----|----|----|----|----|
| C | -1 | 0.8 | -0.2 | -1.2 | -2.2 | -2 | -3 | -4 | -5 |
| A | -2 | -0.2 | 1.8 | 1.8 | 0.8 | -0.2 | -1.2 | -1 | -2 |
| G | -3 |  | -1.2 | 1.8 | 4.8 |  |  |  | -0.2 |
| G | -4 |  | -1.2 | 2.8 |  | 4.4 | 3.4 |  | 0.4 |
| T | -5 |  | -3.2 | -0.4 |  |  |  | 3.4 | 3.4 |
| A | -6 |  | -1.2 | -1.2 |  |  |  | 7.4 | 6.4 |
| C | -7 |  | -1.2 | -2.2 |  |  |  |  | 3.6 |

**Scoring Multiple Alignments:**
Multiple LCS: Match number only
```
A A A   match!
A P A
A A T
A T C
```

**Entropy:**
$-\sum P_x \log P_x$
$X = A, T, G, C$
if $P_x = 0$, then 0

worst entropy: A, T, G, C $-\frac{1}{4}\log(\frac{1}{4}) \cdot 4 = +2$
0.25 for all!

Sum all entropies to find alignment score (of cols)

**Evolutionary Trees:**
Leaves: existing species
Internal vertices: ancestors
Roots: oldest evolutionary ancestor

unrooted   rooted

Edges may have weight $\to$ # of mutations or time estimate of evolution process.
$d_{ij}(T) \to$ tree distance between $i$ and $j$
from 1 to 2: $d_{12}(T) = x + y + z$
$D_{ij} \to n \times n_q$ distance matrix edit

---

**Use gap-open opening gaps!** $O(nm)$ to fill

**Global Alignment (Pairwise)**
$$S_{i,j} = \max \begin{cases} S_{i-1,j-1} + match\_score & if\ v_i = w_i \\ S_{i-1,j-1} + mismatch\_score & if\ v_i \neq w_i \\ S_{i-1,j} + gap\_score \\ S_{i,j-1} + gap\_score \end{cases}$$

$$Sequence\ Identity = \frac{num\_matches}{len(alignment)} \times 100$$

**Local Alignment:**
Don't decrease scores under '0'
$$S_{i,j} = \max \begin{cases} S_{i-1,j-1} + \delta(v_i, w_i) \\ S_{i-1,j} + \delta(v_i, -) \\ S_{i,j-1} + \delta(-, w_i) \\ 0 \end{cases}$$

$$V(i,j) = \max \begin{cases} V(i-1,j-1) + S(S_x(i), j) \\ V(i-1,j) + \delta(S_i(i), '-') \\ V(i,j-1) + \delta('-', j) \end{cases}$$

$S_i \to$ sequence! not score!

| A | 1 | | 1 | | 1 | | 0.4 |
| C | 0.6 | 1 | 0.2 | 1 | | | |
| G | | | | | 1 | | 0.6 |
| T | 0.2 | | | | | | |
| - | 0.2 | | 0.8 | | | | |

**ClustalW**
similarity = frequency of exact matches

**Guide Tree ex.:**
$v_1 \ v_2 \ v_3 \ v_4$
| $v_1$ | - | | |
| $v_2$ | .11 | - | |
| $v_3$ | .87 | .25 | - |
| $v_4$ | .59 | .93 | .62 | - |
then

$v_1 \ v_3 \ v_4 \ v_2$
align 1-3, then 4, then with gaps and m's no lolos

**Sum of Pairs Score**
Calculate as pairwise score! with gaps and mismatches
for all pairs:
if 4 alignment then 6 probabilities

**Degenerate Triples: (DT)**
$D_{ij} + D_{jk} = D_{ik}$
j can be removed.
if there is no DT, create one by shortening edges.
hanging (connected at leaf)

# Change Problem:

$M = c_1 i_1 + c_2 i_2 + .$

where $i_1 + i_2 ...$ minimized.

$c_1, c_2 ...$ denominators.

Algorithm:

$$minCoins(M) = min \begin{cases} minCoins(M-c_i) + 1 \\ minCoins(m-c_n) + 1 \end{cases}$$

## Additive_Phylogeny(D)

if D is a 2×2 matrix:

  T = tree of a single edge of len $D_{1,2}$

  return T

if D is non-degenerate:

  $\delta$ = trimming parameter of matrix D

  for all $1 \leq i \neq j \leq n$

    $D_{ij} \leftarrow D_{ij} - 2\delta$

else

  $\delta = 0$

Find a triple i,j,k such that $D_{ij} + D_{jk} = D_{ik}$

x = $D_{ij}$

Remove jth row and jth column

T = Additive_Phylogeny(D)

Add a new vertex v to T at distance x from i to k

Add j back to T by creating edge (v,j) of len(?)

for every leaf l in T:

  if distance from l to v in the tree $\neq D_{l,i}$

    output "matrix is not additive"

    return

Extend all hanging edges by $\delta$

return T

---

## Four Point Condition:

i,j,k,l → four leaves of tree

Two of $D_{x,y}$'s are equal and third sum is smaller, four point cond. satisfied.

Theorem: n×n matrix is additive iff four point cond. satisfied for all quartets!

## Parsimony Score: non-weighted

ACCC
/ | \
ACCA  ACCG  ACCC
                / \
              ACCA  ATCC
ATCG  ATCC        2 / \
sum 6           ATCG  ACCG

sum5 → more parsimonius.

## Weighted parsimony: Sankoff Algorithm.

Begin from leaf.

leaf have score according to the character it have.

After: $s_k(parent) = min\{s_i(left\ child) + \delta_{i,k}\} + min\{s_j(right\ child) + \delta_{j,t}\}$

---

# Clustering:

**UPGMA**

## K-Means Clustering:

### Lloyd Algorithm:

Arbitrarily assign K cluster centers.

while cluster centers keep changing:

  Assign each data point to cluster center (nearest)

  Find center of gravity of each cluster and assign them as new cluster centers

{Another

## Greedy Algorithm:

Select arbitrary partition P into k clusters

while forever

  best change ⇐ 0

  for every C (cluster)

    for i not in C (element)

      if moving i into C reduces score:

        if cost(P) - cost $P_{i=c}$ > best change

          best change ⇐ cost P - cost $P_{i=c}$

          $i^* = i$

          $c^* = C$

  if best change > 0

    move $i^*$ to $c^*$

  else

    return P.

## Graph Clustering:

Clique Graph: A graph whose every vertex is connected to every other vertex.

## Fitch Algorithm (Small parsimony)

choose arbitrarily

if A:

if it has parent's label choose it else choose arbitrarily.

{A,C} {T,A}
{A,C}  {T,A}
A  C   T  A

---

# Types of Rearrangements:

Reversal
1 2 3 4 5 → 1 2 -5 -4 -3
6             6

Trans
1 2 3   →  1 2 6 4
4 5 6 location  5 3

Fusion
1 2 3 4  →  1 2 3 4 5
5 6         6
Fission

$\pi$: 1 2 3 4 5 6 7 >
$\rho(3,5)$ ↓

1 2 5 4 3 6 7 8
$\rho(5,6)$ ↓

1 2 5 4 6 3 7 8

if $\pi = $ 1 2 3 6 5 4
  already arranged.
prefix($\pi$) = 3
  increase it every step.
flip 6 5 (two bytes?)

GREEDYSORTING SGθ
P = 0
while S ≠ 0

## Distance Graph:

Threshold $\theta$.

if distance < $\theta$ correct vertices

## Break Point Reversal Search ($\pi$)

while $b(\pi) >> $

choose $\rho$ min($b(\pi \cdot \rho)$)

if $\pi_{i+1} = \pi_i \pm 1$

$\pi_i$ and $\pi_{i+1}$ are adjacent.

there is a break point between any adjacent that are not consecutive.

---

For A:

left
old & new

right
old & new

|   | A | T | G | C |
|---|---|---|---|---|
| A | 0 | 3 | 4 | 9 | sum
| T | 3 | 0 | 2 | 4 |
| G | 4 | 2 | 0 | 4 |
| C | 9 | 4 | 4 | 0 |

A T G C

A
/ \
(A)  (C)

0 ∞ ∞ ∞     ∞ ∞ 0 ∞
A T G C     A T G C

# BLAT:

Index database and find query but BLAST index query and search database

Blast: O(nm)

Dict: All words of len(w)

no gap, score > threshold returns.

Gapped blast.

## Pattern Hunter: Spaced seed:

if that well  but in pattern hunter
1111111...   1101001......
              18 elements
              there are 11, 1's

Group of Models

Extend each hit.