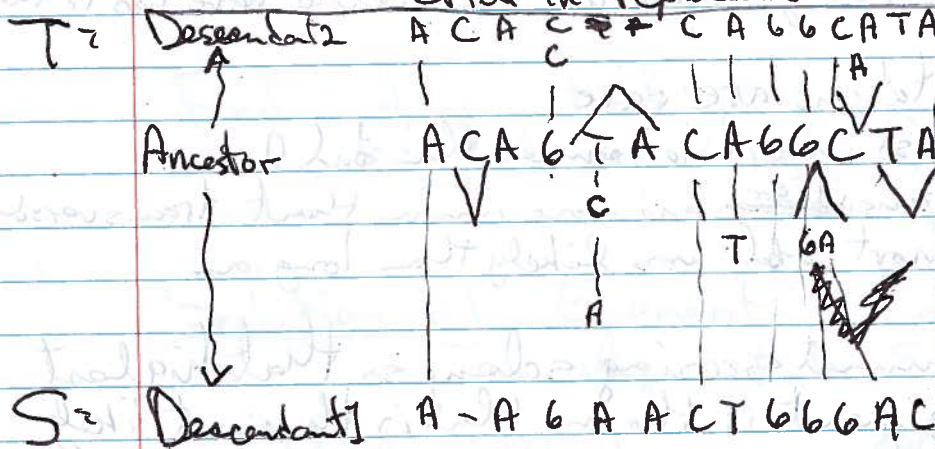


Sept 12 2016

# Sequence alignment

Mutation: Change in sequence due to ~~error in replication~~ error in replication or unrepaired DNA damage



Substitution: Change one nucleotide for another

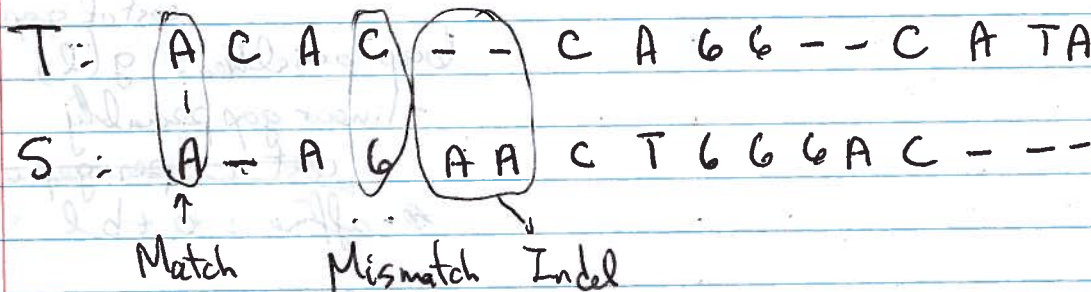
A ↔ G } transition  
C ↔ T } common

A ↔ C } transversion  
G ↔ T } (23 times less common)

Insertion: Insert 1 or more nucleotides

Deletion: Delete one or more nucleotides

Alignment of S, T is obtained by inserting gaps in S and T to obtain S' and T', so that in S' and T', nucleotides derived from same ancestor occur at same position. Gaps reflect nucleotides missing due to del. or insertion



Problem: Ancestor and mutation scenarios are generally not known.

# Idealized Pairwise Global Alignment Problem

Given:  $S, T$

Find: Alignment of  $S$  and  $T$  that is most likely to reflect the evolutionary scenario that would have led to them

- Key ideas:
- (1) Mutations are rare
  - (2) Subst. are more common than indel
  - (3) Transitions are more common than transversions
  - (4) Short indel more likely than long one

⇒ Define alignment scoring scheme so that highest scoring alignment is the one that is the most likely to be correct.

## Scoring a given alignment

$S$ : A C A G T C - - T A

$T$ : A T A - T A A A T A

1 -1 1 -2 1 -1 -2 -2 1 1 = -3

Score: Subst. score + indel score

Subst. cost matrix

A C G T

$$M = \begin{pmatrix} & A & C & G & T \\ A & +1 & -2 & -1 & -2 \\ C & -2 & +1 & -2 & -1 \\ G & & & & \\ T & & & & \end{pmatrix}$$

cost of gap of length  $l$   
Gap penalties:  $g(l)$

- linear gap penalty:  $b \cdot l = -5$

cost  $c$  per gap character

- affine:  $a + b \cdot l = -5 - 12$



# Pairwise seq. Aln. Problem (linear gap penalty)

Given:  $S_1, \dots, S_m$   
 $T_1, \dots, T_n$  sequences  
 $M$  subst. cost matrix  
 $c$

Find: Alignment of  $S, T$   
s.t. that  $\text{score}(\text{Aln}(S, T))$  is maximized

Solution #1: Enumerate and evaluate all possible alignments for  $S$  on  $T$ , report best

Problem: Way too slow  $O(2^{m+n})$

---

Solution #2: Needleman - Wunsch algo. (1970)  
 $S = AGCT$   
 $T = ACA$