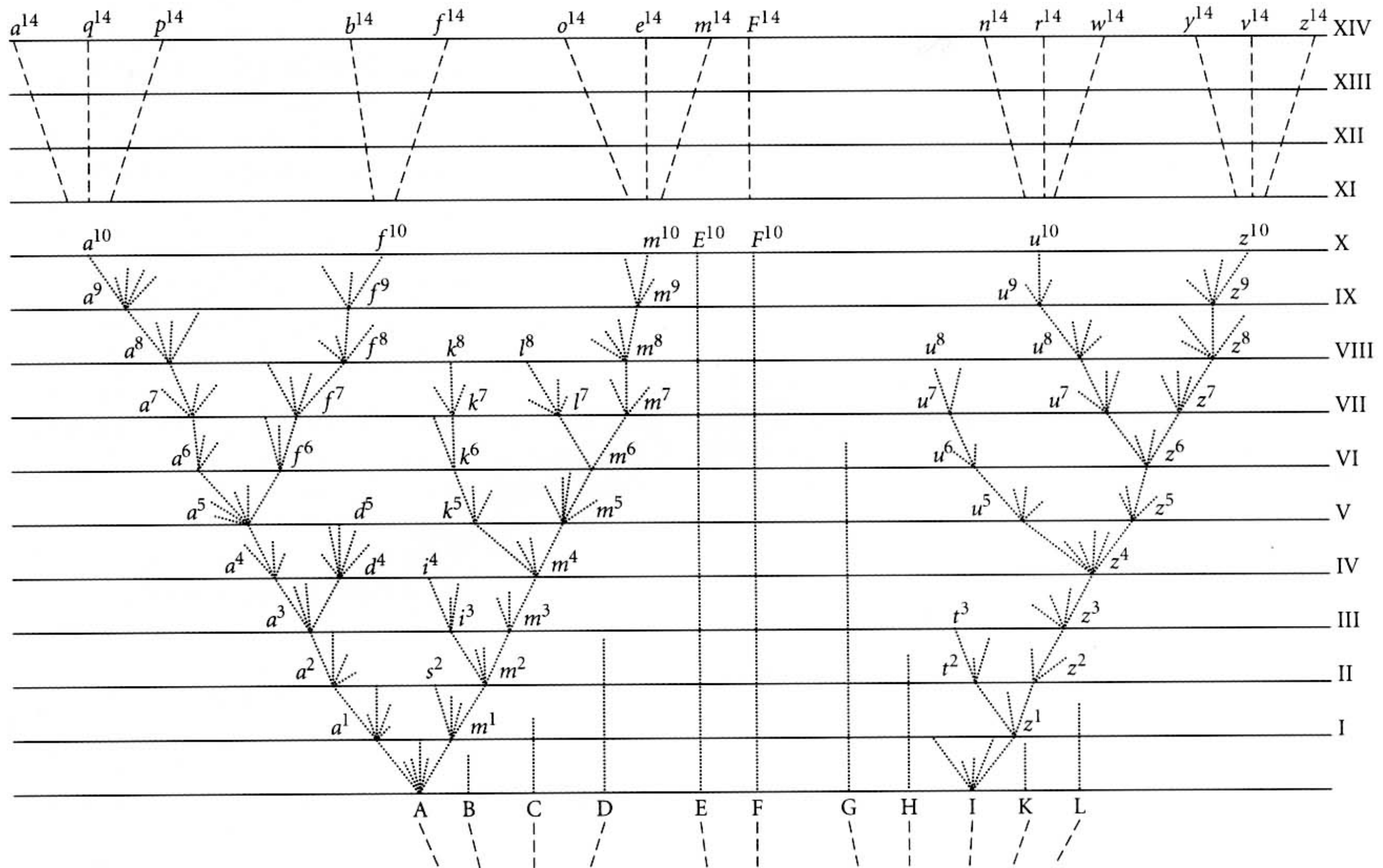# Week 3

Alignment, homology, evolution

# Charles Darwin (1809-1882)

- Avid Naturalist from youth
- Bad Student
- Quit Med School
- Tried Clergy
- Connections got him the Naturalist job on the Beagle

# The only figure in the Origin of Species

# Chapter 1: Variation under Domestication
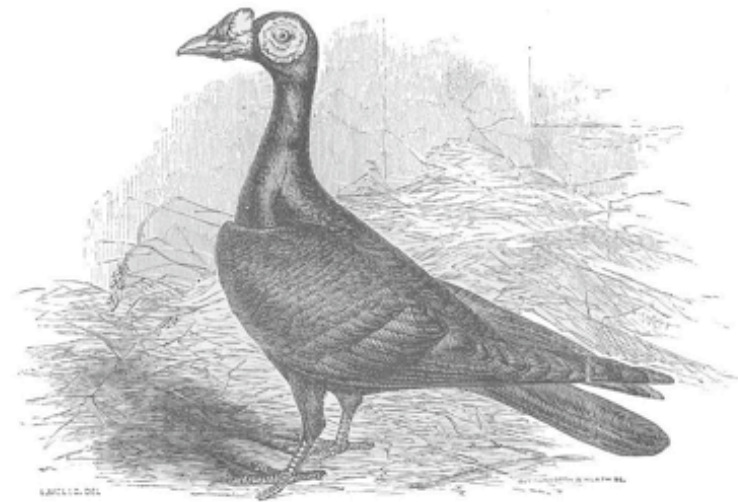


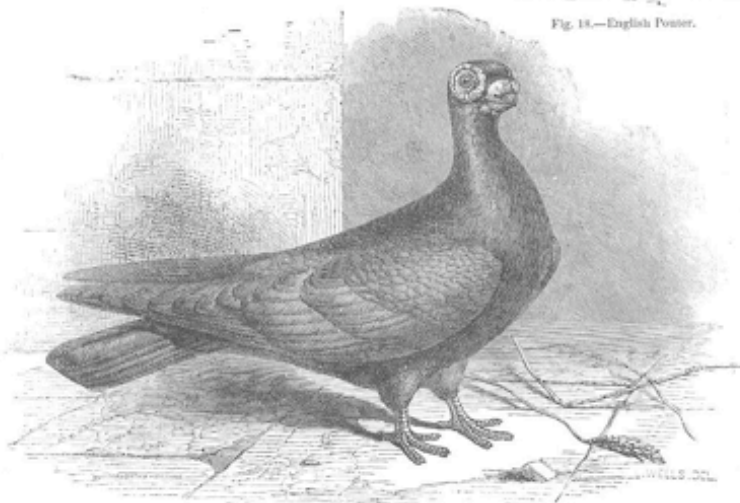Fig. 18.—English Pouter.

Fig. 19.—English Carrier.

Fig. 20.—English Barb.

Fig. 21.—The Rock Pigeon, or Columba livia.
The parent-form of all domesticated Pigeons.

Fig. 21.—English Fantail.

Fig. 22.—African Owl.

Fig. 23.—Short-faced English Tumbler.

# Domesticated *Brassica oleracea*
# (Cabbages)

Cauliflower

Kohlrabi

Broccoli

Turnips

Kale

Brussel Sprouts

# The only figure in the Origin of Species

# From Speciation to Phylogenies

Many rounds of speciation define a unique phylogenetic tree of relationships

Sequential speciation creates a hierarchy

Speciation is lineage splitting

# Phylogenies



a.k.a. evolutionary trees

# Phylogenies

**Genome Sequence**



gtcgagttat cccaacctgg

agtatactga tagcttccac

caaccaggtg ctctactagg

gagtcctggg cacagcccta

atatggtgag ctcatacgat

# Homology:
## similarity of structure due to descent from a common ancestor



Humerus

Radius

Ulna

Carpals

Metacarpals

Phalanges

**Human**

**Mole**

**Horse**

**Dolphin**

**Bat**

# Molecular Homology of Amino Acid positions in Proteins

| | |
|---|---|
| Bacteria | P L F D F A Y Q G F A R G – L E E D A E G L R A F A A M H K E L I V A S S Y S K N F G L Y N E R V G |
| Yeast | A L F D T A Y Q G F A T G D L D K D A Y A V R X X L S T V S P V F V C Q S F A K N A G M Y G E R V G |
| Alfalfa | P F F D S A Y Q G F A S G S L D A D A Q P V R L F V A D G G E L L V A Q S Y A K N M G L Y G E R V G |
| Chicken | P F F D S A Y Q G F A S G S L D K D A W A V R Y F V S E G F E L F C A Q S F S K N F G L Y N E R V G |
| Rat | P F F D S A Y Q G F A S G D L E K D A W A I R Y F V S E G F E L F C P Q S F S K N F G L Y N E R V G |
| Horse | P F F D S A Y Q G F A S G N L D R D A W A V R Y F V S E G F E L F C A Q S F S K N F G L Y N E R V G |
| Pig | P F F D S A Y Q G F A S G N L E K D A W A I R Y F V S E G F E L F C A Q S F S K N F G L Y N E R V G |
| Human | P F F D S A Y Q G F A S G N L E R D A W A I R Y F V S E G F E F F C A Q S F S K N F G L Y N E R V G |

**Figure 12–1**

*A comparison of eight organisms for a 50-amino-acid-long sequence of the enzyme aspartate transaminase. For the amino acid abbreviations, see Fig. 7-2 or Table 12-2. (Adapted from Benner et al.)*

# Deep homology of genetic code



Table 1.1 The standard genetic code

# Phylogeny of Primates based on DNA-DNA Hybridization



Sibley and Ahlquist 1984. Mol. Evol. 20:2-15

# So want to focus on homologous traits / characters

| | |
|---|---|
| *Sequence1* | -TCAGGA-TGAAC----- |
| *Sequence2* | ATCACGA-TGAACC--- |
| *Sequence3* | ATCAGGAATGAATCC-- |
| *Sequence4* | -TCACGATTGAATCGC- |
| *Sequence5* | -TCAGGAATGAATCGCM |

In genomes this means ALIGNMENT will be critical- which bases are the same bases?

# Phylogenies

**Genome Sequence**



gtcgagttat cccaacctgg

agtatactga tagcttccac

caaccaggtg ctctactagg

gagtcctggg cacagcccta

atatggtgag ctcatacgat

# Sequence Alignment

Key algorithms introduced starting in 1970

**A General Method Applicable to the Search for Similarities in the Amino Acid Sequence of Two Proteins**

SAUL B. NEEDLEMAN AND CHRISTIAN D. WUNSCH

*Department of Biochemistry, Northwestern University, and Nuclear Medicine Service, V. A. Research Hospital Chicago, Ill. 60611, U.S.A.*

*(Received 21 July 1969)*

**Identification of Common Molecular Subsequences**

The identification of maximally homologous subsequences among sets of long sequences is an important problem in molecular sequence analysis. The problem is straightforward only if one restricts consideration to contiguous subsequences (segments) containing no internal deletions or insertions. The more general problem has its solution in an extension of sequence metrics (Sellers 1974; Waterman *et al.*, 1976) developed to measure the minimum number of "events" required to convert one sequence into another.

**Global**: Needleman-Wunsch algorithm

**Local**: Smith-Waterman algorithm

# Sequence Alignment

key idea— edit distance

how many changes to make two words / sets the same?

PEAR
→ edit distance: 1,
swap 'P' for 'B'
BEAR

SHOT
→ edit distance: 1,
remove 'O'
SH-T

# Sequence Alignment

key idea— edit distance

how many changes to make two words / sets the same?

PEAR

BEAR

→ edit distance: 1,
swap 'P' for 'B'

SHOT

SH-T

→ edit distance: 1,
remove 'O'

**mismatch**

PEAR

BEAR

# Sequence Alignment

key idea— edit distance

how many changes to make two words / sets the same?

PEAR

⟶ edit distance: 1,
swap 'P' for 'B'

BEAR

SHOT

⟶ edit distance: 1,
remove 'O'

SH-T

**match**

PEAR

BEAR

# Sequence Alignment

key idea— edit distance

how many changes to make two words / sets the same?

PEAR
→ edit distance: 1,
swap 'P' for 'B'
BEAR

SHOT
→ edit distance: 1,
remove 'O'
SH-T

**gap**

PE-R

BEAR

# Sequence Alignment

key idea— edit distance

how many changes to make two words / sets the same?

PEAR  →  edit distance: 1,
BEAR        swap 'P' for 'B'

SHOT  →  edit distance: 1,
SH-T        remove 'O'

What about DNA?

TGTTACGG
GGTTGACTA

TG-TT-ACGG
-GGTTGACTA

edit distance: 5

TGTT-ACGG
GGTTGACTA

edit distance: 4

# Sequence Alignment

Let's align the these two sequences

# Dynamic Programming

## Tower of Hanoi Game



Goal: Get all disks in size order to last peg
Rules: One disk moves per turn. Smaller has to sit on larger

# Dynamic Programming

Tower of Hanoi Game



Can break up **big** problem into simple **smaller** problem

# Dynamic Programming

## Fibonacci sequence



Originally meant to model the growth in rabbit population sizes!

$$F_0 = 1; F_1 = 1$$

$$F_i = F_{i-1} + F_{i-2}$$

| F(0) | F(1) | F(2) | F(3) | F(4) | F(5) | F(6) |
|------|------|------|------|------|------|------|
| 1    | 1    | 2    | 3    | 5    | 8    | 13   |

# Sequence Alignment

So want to find minimum edit distance
how do we search for it?

seq_1 = "TACGGACGG"
seq_2 = "TAGACTA"

Next key idea

treat sequences as matrix

|   |   | T | A | C | G | G | A | C | G | G |
|---|---|---|---|---|---|---|---|---|---|---|
| T |   |   |   |   |   |   |   |   |   |   |
| A |   |   |   |   |   |   |   |   |   |   |
| G |   |   |   |   |   |   |   |   |   |   |
| A |   |   |   |   |   |   |   |   |   |   |
| C |   |   |   |   |   |   |   |   |   |   |
| T |   |   |   |   |   |   |   |   |   |   |
| A |   |   |   |   |   |   |   |   |   |   |

# Sequence Alignment

Global Alignment
Needleman-Wunsch

Assume scores as:

match = 4
mismatch = -3
gap = -2

1. initialize

Seq 1

|   |   | T | A | C | G | G | A | C | G | G |
|---|---|---|---|---|---|---|---|---|---|---|
|   | 0 | -2 | -4 | -6 | -8 | -10 | -12 | -14 | -16 | -18 |
| T | -2 | | | | | | | | | |
| A | -4 | | | | | | | | | |
| G | -6 | | | | | | | | | |
| A | -8 | | | | | | | | | |
| C | -10 | | | | | | | | | |
| T | -12 | | | | | | | | | |
| A | -14 | | | | | | | | | |

Seq 2

# Sequence Alignment

Global Alignment
Needleman-Wunsch

Assume scores as:

match = 4
mismatch = -3
gap = -2

1. initialize

Seq 1

|   |   | T | A | C | G | G | A | C | G | G |
|---|---|---|---|---|---|---|---|---|---|---|
|   | 0 | -2 | -4 | -6 | -8 | -10 | -12 | -14 | -16 | -18 |
| T | -2 | | | | | | | | | |
| A | -4 | | | | | | | | | |
| G | -6 | | | | | | | | | |
| A | -8 | | | | | | | | | |
| C | -10 | | | | | | | | | |
| T | -12 | | | | | | | | | |
| A | -14 | | | | | | | | | |

**gap**

Seq 2

# Sequence Alignment

## Global Alignment
## Needleman-Wunsch

Assume scores as:

match = 4
mismatch = -3
gap = -2

1. initialize
2. fill in table

consider value of x
3 ways to get there:
1) from previous (mis)match
2) gap in seq1
3) gap in seq2

Seq 1

| | | T | A | C | G | G | A | C | G | G |
|---|---|---|---|---|---|---|---|---|---|---|
| | 0 | -2 | -4 | -6 | -8 | -10 | -12 | -14 | -16 | -18 |
| T | -2 | x | | | | | | | | |
| A | -4 | | | | | | | | | |
| G | -6 | | | | | | | | | |
| A | -8 | | | | | | | | | |
| C | -10 | | | | | | | | | |
| T | -12 | | | | | | | | | |
| A | -14 | | | | | | | | | |

Seq 2

# Sequence Alignment

## Global Alignment
## Needleman-Wunsch

Assume scores as:

match = 4
mismatch = -3
gap = -2

$x_{i,j} = \max(\ x_{i-1,j} + gap,$
$\qquad\qquad x_{i,j-1} + gap,$
$\qquad\qquad x_{i-1,j-1} + (mis)match)$

Seq 1

| | | T | A | C | G | G | A | C | G | G |
|---|---|---|---|---|---|---|---|---|---|---|
| | 0 | -2 | -4 | -6 | -8 | -10 | -12 | -14 | -16 | -18 |
| T | -2 | 4 | | | | | | | | |
| A | -4 | | | | | | | | | |
| G | -6 | | | | | | | | | |
| A | -8 | | | | | | | | | |
| C | -10 | | | | | | | | | |
| T | -12 | | | | | | | | | |
| A | -14 | | | | | | | | | |

Seq 2

# Sequence Alignment

## Global Alignment
## Needleman-Wunsch

Assume scores as:

match = 4
mismatch = -3
gap = -2

x_i,j = max( x_i-1,j + gap,
         x_i,j-1 + gap,
         x_i-1,j-1 + (mis)match)

Seq 1

|   |   | T | A | C | G | G | A | C | G | G |
|---|---|---|---|---|---|---|---|---|---|---|
|   | 0 | -2 | -4 | -6 | -8 | -10 | -12 | -14 | -16 | -18 |
| T | -2 | 4 | 2 |   |   |   |   |   |   |   |
| A | -4 |   |   |   |   |   |   |   |   |   |
| G | -6 |   |   |   |   |   |   |   |   |   |
| A | -8 |   |   |   |   |   |   |   |   |   |
| C | -10 |   |   |   |   |   |   |   |   |   |
| T | -12 |   |   |   |   |   |   |   |   |   |
| A | -14 |   |   |   |   |   |   |   |   |   |

Seq 2

# Sequence Alignment

Global Alignment
Needleman-Wunsch

Assume scores as:

match = 4
mismatch = -3
gap = -2

x_i,j = max( x_i-1,j + gap,
        x_i,j-1 + gap,
        x_i-1,j-1 + (mis)match)

Seq 1

| | | T | A | C | G | G | A | C | G | G |
|---|---|---|---|---|---|---|---|---|---|---|
| | 0 | -2 | -4 | -6 | -8 | -10 | -12 | -14 | -16 | -18 |
| T | -2 | 4 | 2 | 0 | -2 | -4 | | | | |
| A | -4 | 2 | 8 | 6 | 4 | 2 | | | | |
| G | -6 | 0 | 6 | 5 | 10 | 8 | | | | |
| A | -8 | | | | | | | | | |
| C | -10 | | | | | | | | | |
| T | -12 | | | | | | | | | |
| A | -14 | | | | | | | | | |

Seq 2

# Sequence Alignment

Global Alignment
Needleman-Wunsch

Assume scores as:

match = 4
mismatch = -3
gap = -2

$x_{i,j} = \max(\ x_{i-1,j} + \text{gap},$
$\qquad\quad x_{i,j-1} + \text{gap},$
$\qquad\quad x_{i-1,j-1} + (\text{mis})\text{match})$

1. initialize
2. fill in table
3. traceback

Seq 1

|  |  | T | A | C | G | G | A | C | G | G |
|---|---|---|---|---|---|---|---|---|---|---|
|  | 0 | -2 | -4 | -6 | -8 | -10 | -12 | -14 | -16 | -18 |
| T | -2 | 4 | 2 | 0 | -2 | -4 | -6 | -8 | -10 | -12 |
| A | -4 | 2 | 8 | 6 | 4 | 2 | 0 | -2 | -4 | -6 |
| G | -6 | 0 | 6 | 5 | 10 | 8 | 6 | 4 | 2 | 0 |
| A | -8 | -2 | 4 | 3 | 8 | 7 | 12 | 10 | 8 | 6 |
| C | -10 | -4 | 2 | 8 | 6 | 5 | 10 | 16 | 14 | 12 |
| T | -12 | -6 | 0 | 6 | 5 | 3 | 8 | 14 | 13 | 11 |
| A | -14 | -8 | -2 | 4 | 3 | 2 | 7 | 12 | 11 | 10 |

Seq 2

# Sequence Alignment

### Global Alignment
### Needleman-Wunsch

traceback step

start at bottom right
work way back up
follow max score path

can read off alignment
backwards

Seq 1     **G**
Seq 2     **A**

Seq 1

|   |   | T | A | C | G | G | A | C | G | G |
|---|---|---|---|---|---|---|---|---|---|---|
|   | 0 | -2 | -4 | -6 | -8 | -10 | -12 | -14 | -16 | -18 |
| T | -2 | 4 | 2 | 0 | -2 | -4 | -6 | -8 | -10 | -12 |
| A | -4 | 2 | 8 | 6 | 4 | 2 | 0 | -2 | -4 | -6 |
| G | -6 | 0 | 6 | 5 | 10 | 8 | 6 | 4 | 2 | 0 |
| A | -8 | -2 | 4 | 3 | 8 | 7 | 12 | 10 | 8 | 6 |
| C | -10 | -4 | 2 | 8 | 6 | 5 | 10 | 16 | 14 | 12 |
| T | -12 | -6 | 0 | 6 | 5 | 3 | 8 | 14 | 13 | 11 |
| A | -14 | -8 | -2 | 4 | 3 | 2 | 7 | 12 | 11 | 10 |

Seq 2

# Sequence Alignment

## Global Alignment
## Needleman-Wunsch

traceback step

start at bottom right
work way back up
follow max score path

can read off alignment
backwards

Seq 1    **GG**

Seq 2    **TA**

Seq 1

Seq 2

|   |   | T | A | C | G | G | A | C | G | G |
|---|---|---|---|---|---|---|---|---|---|---|
|   | 0 | -2 | -4 | -6 | -8 | -10 | -12 | -14 | -16 | -18 |
| T | -2 | 4 | 2 | 0 | -2 | -4 | -6 | -8 | -10 | -12 |
| A | -4 | 2 | 8 | 6 | 4 | 2 | 0 | -2 | -4 | -6 |
| G | -6 | 0 | 6 | 5 | 10 | 8 | 6 | 4 | 2 | 0 |
| A | -8 | -2 | 4 | 3 | 8 | 7 | 12 | 10 | 8 | 6 |
| C | -10 | -4 | 2 | 8 | 6 | 5 | 10 | 16 | 14 | 12 |
| T | -12 | -6 | 0 | 6 | 5 | 3 | 8 | 14 | 13 | 11 |
| A | -14 | -8 | -2 | 4 | 3 | 2 | 7 | 12 | 11 | 10 |

# Sequence Alignment

### Global Alignment
### Needleman-Wunsch

Seq 1

traceback step

start at bottom right
work way back up
follow max score path

can read off alignment
backwards

Seq 1    `TACGGACGG`
Seq 2    `TA--GACTA`

Seq 2

|   |   | T | A | C | G | G | A | C | G | G |
|---|---|---|---|---|---|---|---|---|---|---|
|   | 0 | -2 | -4 | -6 | -8 | -10 | -12 | -14 | -16 | -18 |
| T | -2 | 4 | 2 | 0 | -2 | -4 | -6 | -8 | -10 | -12 |
| A | -4 | 2 | 8 | 6 | 4 | 2 | 0 | -2 | -4 | -6 |
| G | -6 | 0 | 6 | 5 | 10 | 8 | 6 | 4 | 2 | 0 |
| A | -8 | -2 | 4 | 3 | 8 | 7 | 12 | 10 | 8 | 6 |
| C | -10 | -4 | 2 | 8 | 6 | 5 | 10 | 16 | 14 | 12 |
| T | -12 | -6 | 0 | 6 | 5 | 3 | 8 | 14 | 13 | 11 |
| A | -14 | -8 | -2 | 4 | 3 | 2 | 7 | 12 | 11 | 10 |

# Sequence Alignment

Local Alignment
Smith-Waterman

No negative scores

traceback starts at highest
score

|   |   | T | A | C | G | G | A | C | G | G |
|---|---|---|---|---|---|---|---|---|---|---|
|   | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| T | 0 | 4 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| A | 0 | 2 | 8 | 6 | 4 | 2 | 0 | 0 | 0 | 0 |
| G | 0 | 0 | 6 | 5 | 10 | 8 | 6 | 4 | 2 | 0 |
| A | 0 | 0 | 4 | 3 | 8 | 7 | 12 | 10 | 8 | 6 |
| C | 0 | 0 | 2 | 8 | 6 | 5 | 10 | 16 | 14 | 12 |
| T | 0 | 0 | 0 | 6 | 5 | 3 | 8 | 14 | 13 | 11 |
| A | 0 | 0 | 0 | 4 | 3 | 2 | 7 | 12 | 11 | 10 |

# Sequence Alignment

## Local Alignment
## Smith-Waterman

No negative scores

traceback starts at highest score

terminates at first zero

```
TACGGAC
TA__GAC
```

|   |   | T | A | C | G | G | A | C | G | G |
|---|---|---|---|---|---|---|---|---|---|---|
|   |   | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| T | 0 | 4 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| A | 0 | 2 | 8 | 6 | 4 | 2 | 0 | 0 | 0 | 0 |
| G | 0 | 0 | 6 | 5 | 10 | 8 | 6 | 4 | 2 | 0 |
| A | 0 | 0 | 4 | 3 | 8 | 7 | 12 | 10 | 8 | 6 |
| C | 0 | 0 | 2 | 8 | 6 | 5 | 10 | 16 | 14 | 12 |
| T | 0 | 0 | 0 | 6 | 5 | 3 | 8 | 14 | 13 | 11 |
| A | 0 | 0 | 0 | 4 | 3 | 2 | 7 | 12 | 11 | 10 |

# Sequence Alignment

```
Seq1    TACGGACGG
Seq2    TAGACTA
```

Local Alignment
Smith-Waterman

```
TACGGAC
TA__GAC
```

Global Alignment
Needleman-Wunsch

```
TACGGACGG
TA__GACTA
```

in this case very similar, but not generally true
use local alignment when you want short, subset matches
global when whole sequence alignment wanted