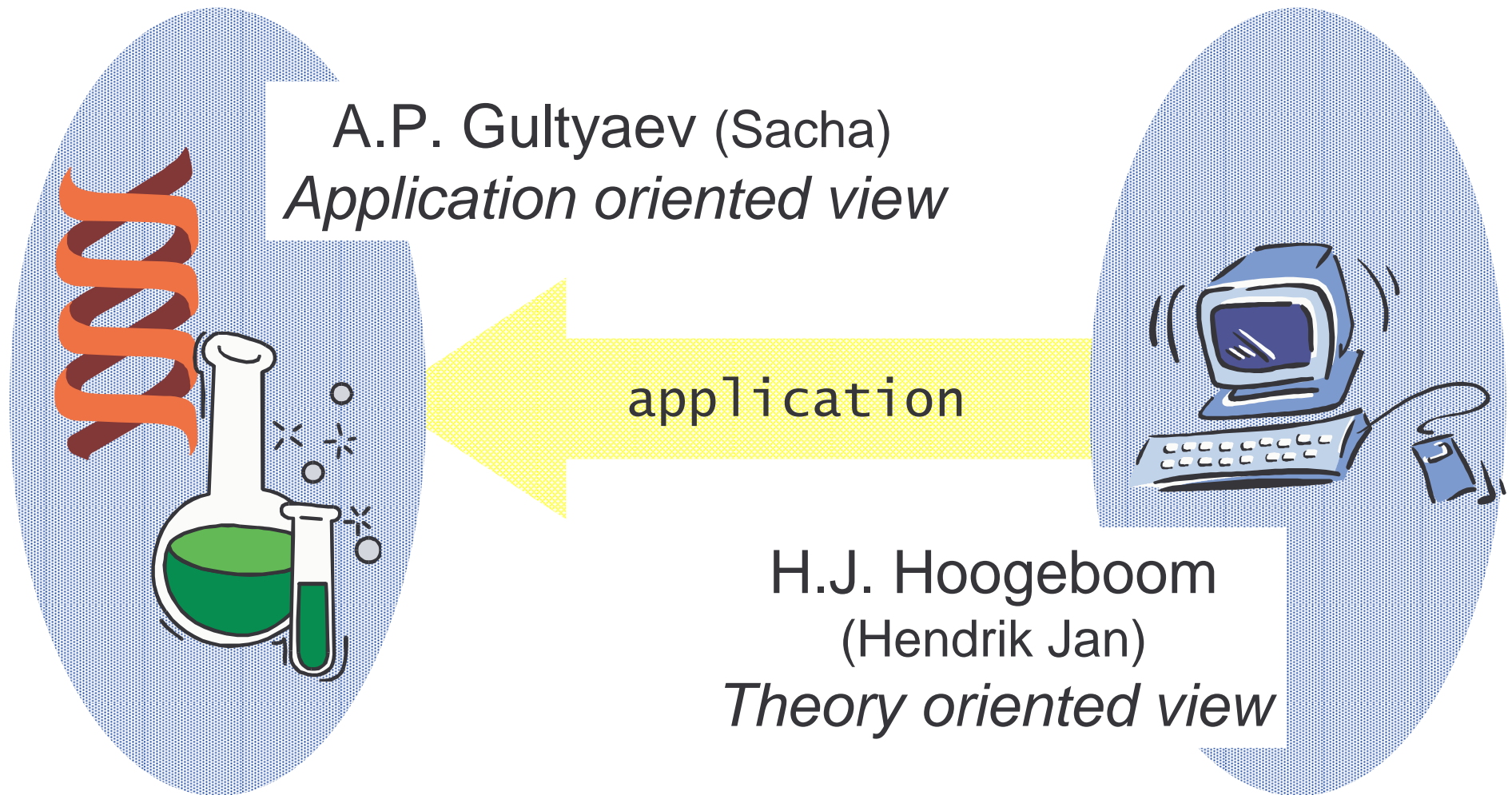


# Computational Biology



[www.liacs.nl/home/hoogeboo/mcb/](http://www.liacs.nl/home/hoogeboo/mcb/)

problem  $\Rightarrow$  model (eg. graph)

- known algorithms
- characterization

unprecise data

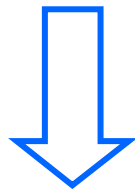
complexity

$\Rightarrow$  heuristics

what *is* the right answer?

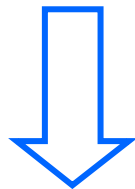
# central dogma

DNA



transcription  
& splicing

RNA



translation

protein  
*eiwit*

‘gene expression’

# two alphabets

DNA

bases

4 symbols

a c t g

RNA

a c u g

proteins

amino acids

20 symbols

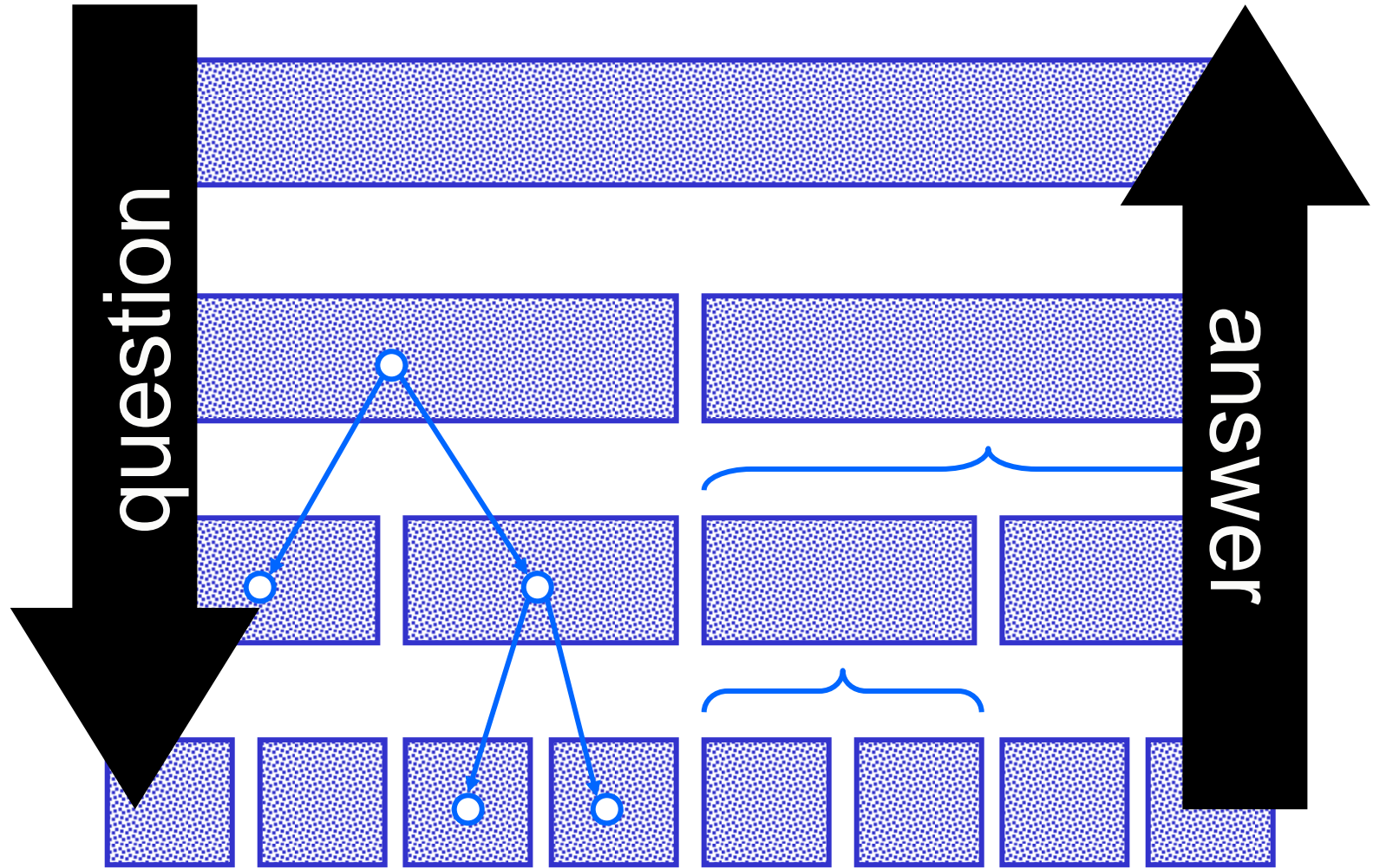
A R D N C

E Q G H I

L K M F P

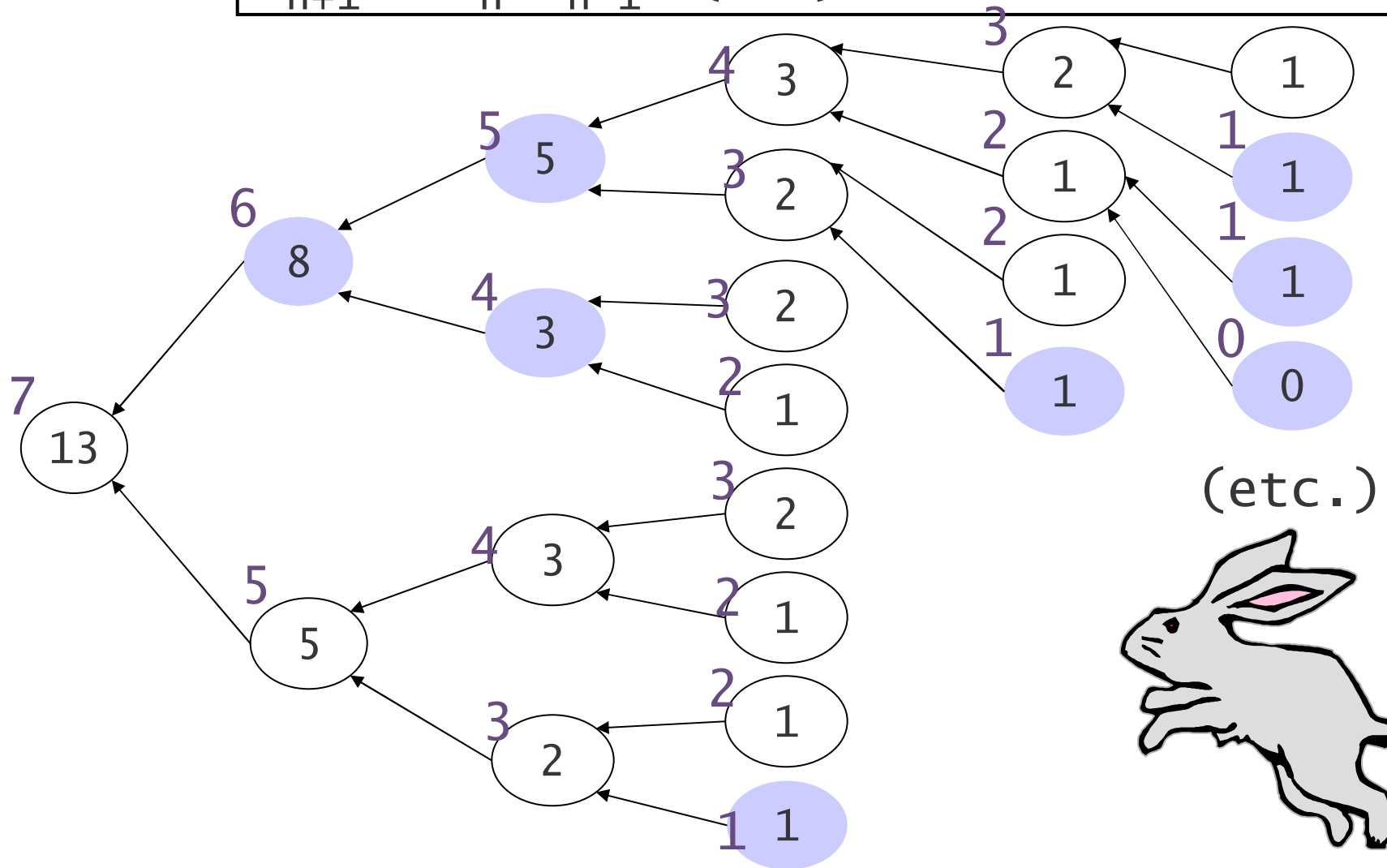
S T W Y V

recursion



# recursion: bad example

$$a_0 = 0 \quad a_1 = 1$$
$$a_{n+1} = a_n + a_{n-1} \quad (n \geq 1)$$



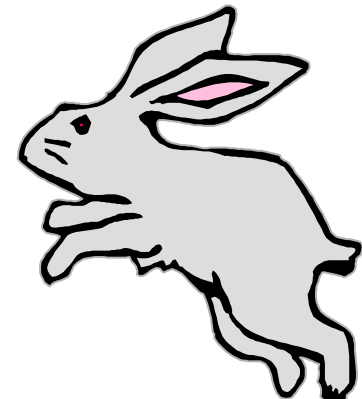
# recursion: dynamic programming

$$a_0 = 0 \quad a_1 = 1$$
$$a_{n+1} = a_n + a_{n-1} \quad (n \geq 1)$$

0, 1, 1, 2, 3, 5, 8, 13, 21, ...

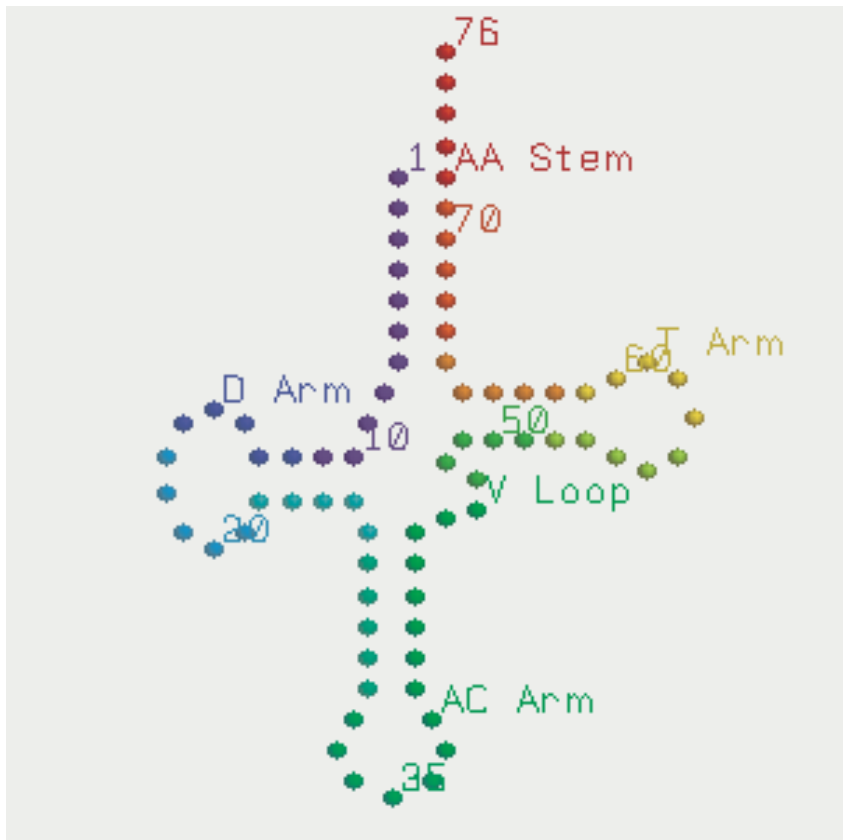
## dynamic programming

- *memoization*  
store results in table
- ( work bottom-up )

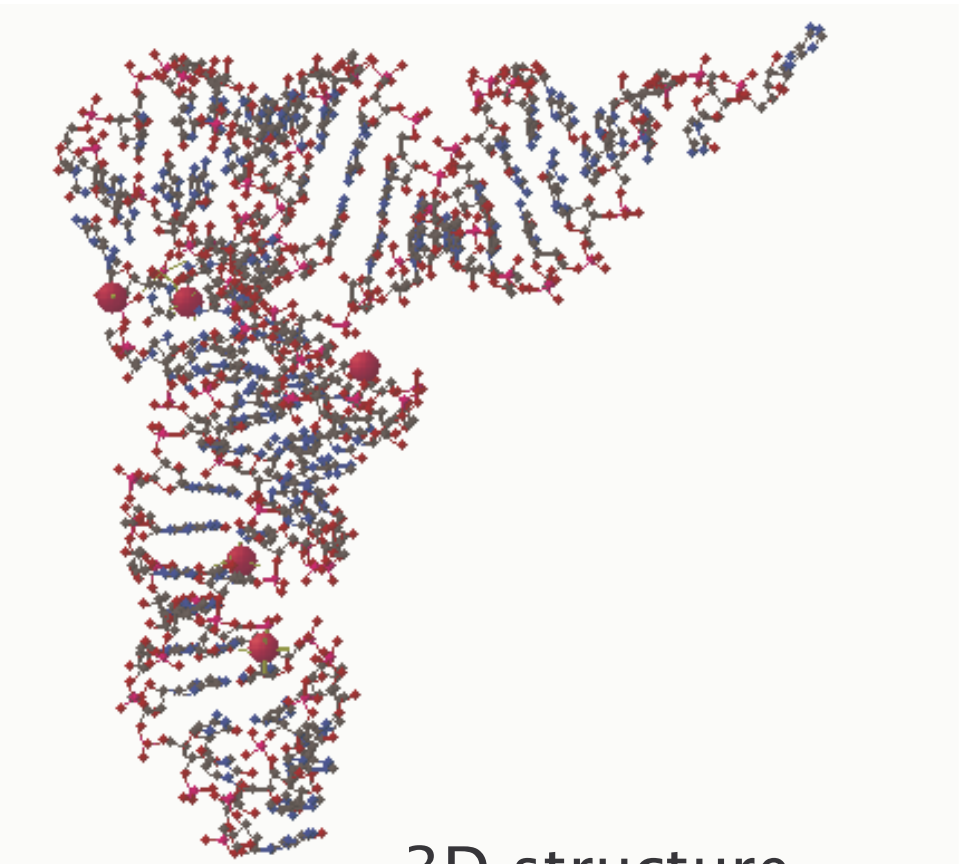


# structure of tRNA

sequence  $\Rightarrow$  structure  $\Rightarrow$  function



2D structure



3D structure



# sequence alignment

TCAGACGATTG  
TCGGAGCTG

TCAG-ACG-ATTG  
TC-GGA-GC-T-G

insertion  
deletion  
substitution

TCAGACGATTG  
TCGGAGC--TG

match  
mismatch  
gap

TCAG-ACGATTG  
TC-GGA-GCT-G

# sequence alignment

TCAGACGATTG  
TCGGAGCTG

TCAG-ACG-ATTG  
TC-GGA-GC-T-G

$$14-6=8$$

TCAGACGATTG  
TCGGAGC--TG

$$12-3-2=7$$

match +2  
mismatch -1  
gap -1

TCAG-ACGATTG  
TC-GGA-GCT-G

$$14-1-4=9$$

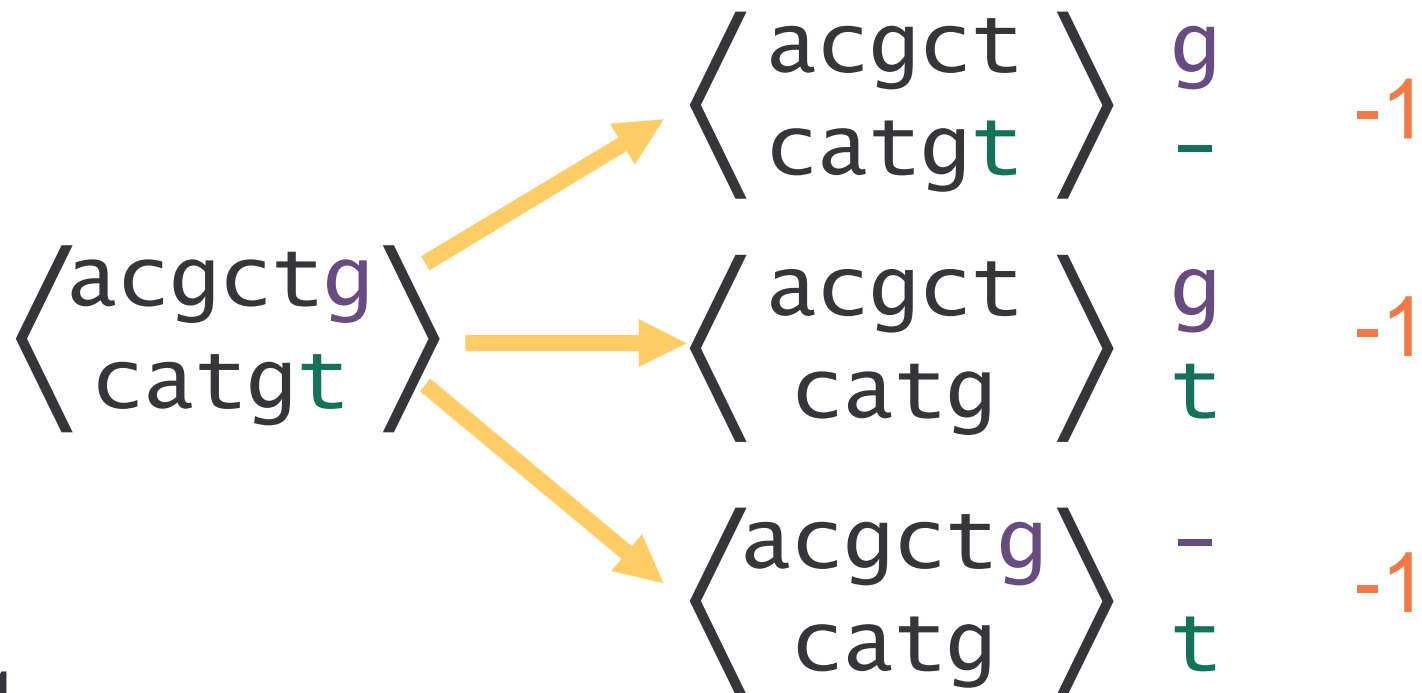


basic algorithm

## alignment

- recursive principle
- dynamic programming

# alignment: recursion



$$\sigma(-,x) = -1$$

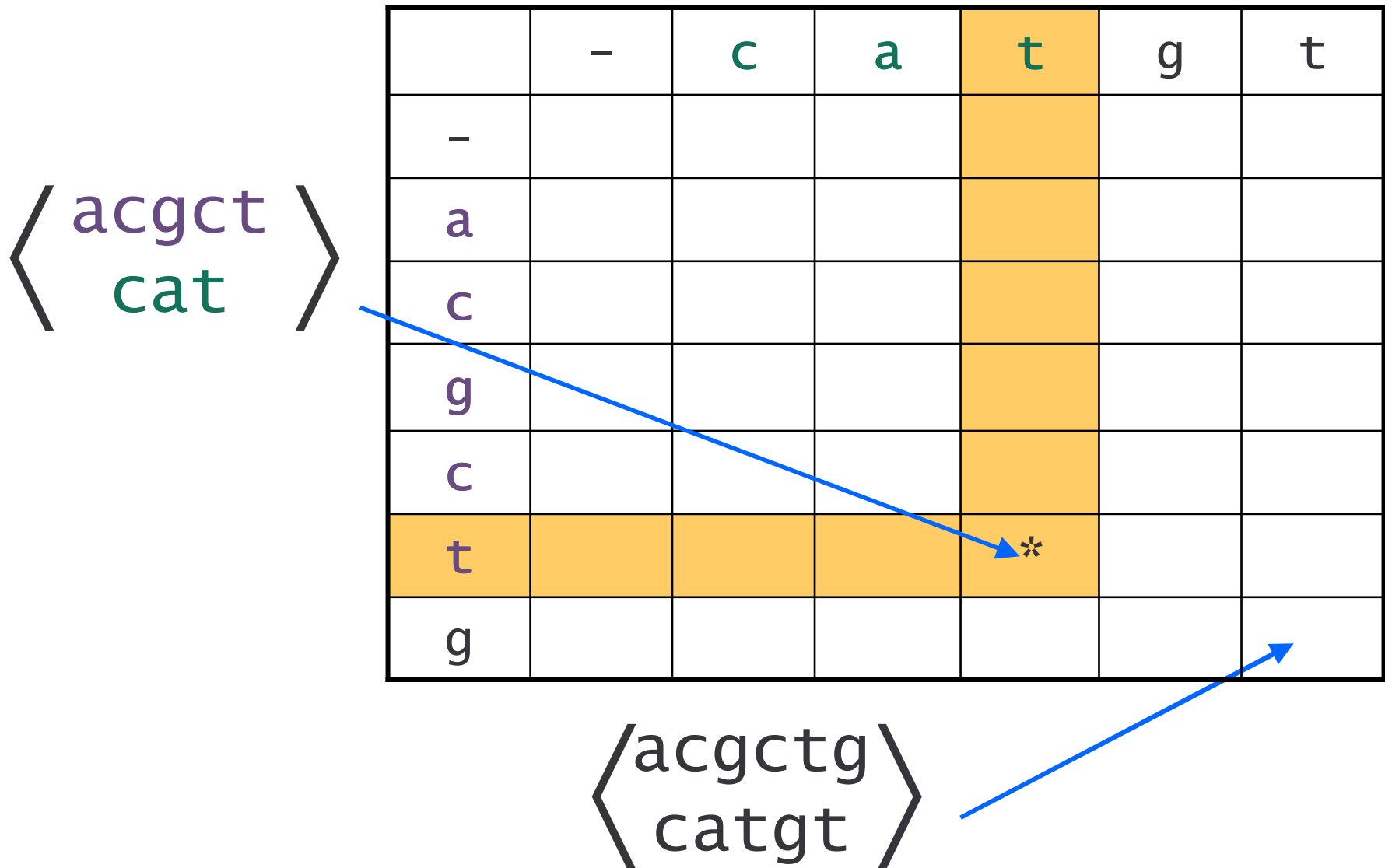
$$\sigma(x,-) = -1$$

$$\sigma(x,y) = -1$$

$$\sigma(x,x) = 2$$

'reward and punishment'

# dynamic programming



# dynamic programming

⟨ acgct  
cat ⟩

	-	c	a	t	g	t
-	0	-1	-2	-3	-4	-5
a	-1	-1	1	0	-1	-2
c	-2	1	0	0	-1	-2
g	-3	0	0	-1	2	1
c	-4	-1	-1	-1	1	1
t	-5	-2	-2	1	0	3
g	-6	-3	-3	0	3	2

⟨ acgctg  
catgt ⟩

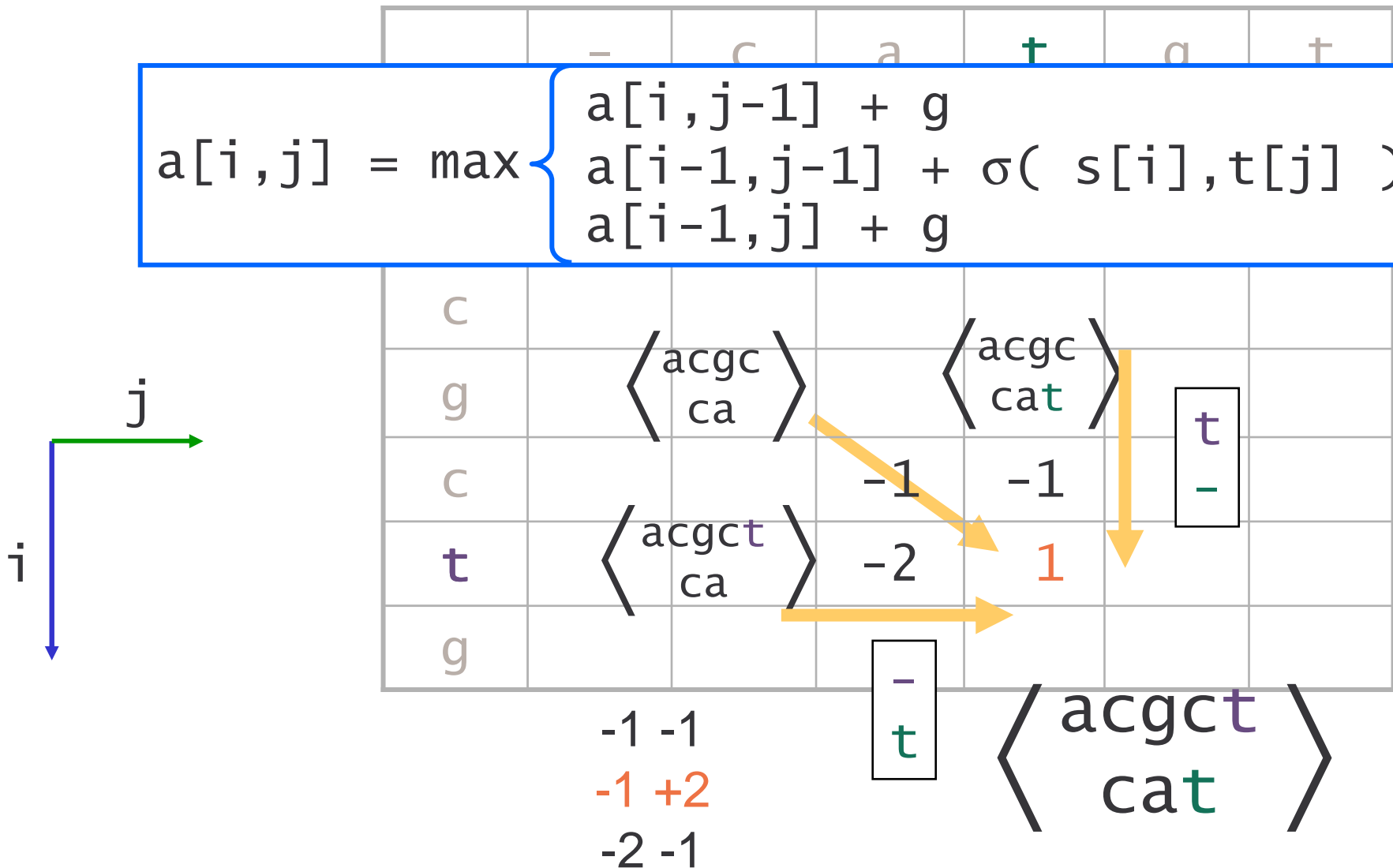
# initialization

< ---  
cat >

	-	c	a	t	g	t
-	0	-1	-2	-3	-4	-5
a	-1					
c	-2					
g	-3					
c	-4					
t	-5					
g	-6					

# dynamic programming

$$a[i, j] = \max \begin{cases} a[i, j-1] + g \\ a[i-1, j-1] + \sigma( s[i], t[j] ) \\ a[i-1, j] + g \end{cases}$$





# alignment

	-	c	a	t	g	t
-	0	-1	-2	-3	-4	-5
a	-1	-1	1	0	-1	-2
c	-2	1	0	0	-1	-2
g	-3	0	0	-1	2	1
c	-4	-1	-1	-1	1	1
t	-5	-2	-2	1	0	3
g	-6	-3	-3	0	3	2

⟨acgctg⟩  
⟨catgt⟩

# traceback

-acgctg

catg-t-

acgctg-

-ca-tgt

acgctg-

-c-atgt

⟨acgctg⟩  
⟨catgt⟩

	-	c	a	t	g	t
-	0	-1	-2	-3	-4	-5
a	-1	-1	1	0	-1	-2
c	-2	1	0	0	-1	-2
g	-3	0	0	-1	2	1
c	-4	-1	-1	-1	1	1
t	-5	-2	-2	1	0	3
g	-6	-3	-3	0	3	2

g  
-

-  
t

# Eddy S.R. (2004a)

```
/* Recursion: the heart of the DP algorithm.*/
/* Initialization. */

s[0][0] = 0;
for (i = 1; i <= M; i++)  S[i][0] = i * INDEL;
for (j = 1; j <= N; j++)  S[0][j] = j * INDEL;

/* Dynamic programming, global alignment (Needleman/wunsch) recursion. */

for (i = 1; i <= M; i++)
  for (j = 1; j <= N; j++)
  {
    /* case #1: i,j are aligned */
    if (x[i] == y[j]) S[i][j] = S[i-1][j-1] + MATCH;
    else               S[i][j] = S[i-1][j-1] + MISMATCH;

    sc = S[i-1][j] + INDEL;          /* case #2: i aligned to - */
    if (sc > S[i][j]) S[i][j] = sc;

    sc = S[i][j-1] + INDEL;         /* case #3: j aligned to - */
    if (sc > S[i][j]) S[i][j] = sc;
  }

/* The result (optimal alignment score) is now in S[M][N]. */
```

# scoring and parameter choice

A vs. A  
A vs. C

match M +2 +2  
mismatch m -1 -2  
gap g -1 -1

A vs. A-  
C vs. -C

AT vs. -AT  
TA vs. TA-

	A	C	G	T
A	91	-114	-31	-123
C		100	-125	-31
G			100	-114
T				91

gap 400 + 30k

ask your favourite molecular biologist !

# PAM250 Matrix

details → notes APG

<b>C</b>	12																			
<b>S</b>	0	2																		
<b>T</b>	-2	1	3																	
<b>P</b>	-3	1	0	6																
<b>A</b>	-2	1	1	1	2															
<b>G</b>	-3	1	0	-1	1	5														
<b>N</b>	-4	1	0	-1	0	0	2													
<b>D</b>	-5	0	0	-1	0	1	2	4												
<b>E</b>	-5	0	0	-1	0	0	1	3	4											
<b>Q</b>	-5	-1	-1	0	0	-1	1	2	2	4										
<b>H</b>	-3	-1	-1	0	-1	-2	2	1	1	3	6									
<b>R</b>	-4	0	-1	0	-2	-3	0	-1	-1	1	2	6								
<b>K</b>	-5	0	0	-1	-1	-2	1	0	0	1	0	3	5							
<b>M</b>	-5	-2	-1	-2	-1	-3	-2	-3	-2	-1	-2	0	0	6						
<b>I</b>	-2	-1	0	-2	-1	-3	-2	-2	-2	-2	-2	-2	-2	2	5					
<b>L</b>	-6	-3	-2	-3	-2	-4	-3	-4	-3	-2	-2	-3	-3	4	2	6				
<b>V</b>	-2	-1	0	-1	0	-1	-2	-2	-2	-2	-2	-2	-2	2	4	2	4			
<b>F</b>	-4	-3	-3	-5	-4	-5	-4	-6	-5	-5	-2	-4	-5	0	1	2	-1	9		
<b>Y</b>	0	-3	-3	-5	-3	-5	-2	-4	-4	-4	0	-4	-4	-2	-1	-1	-2	7	10	
<b>W</b>	-8	-2	-5	-6	-6	-7	-4	-7	-7	-5	-3	2	-3	-4	-5	-2	-6	0	0	17
	<b>C</b>	<b>S</b>	<b>T</b>	<b>P</b>	<b>A</b>	<b>G</b>	<b>N</b>	<b>D</b>	<b>E</b>	<b>Q</b>	<b>H</b>	<b>R</b>	<b>K</b>	<b>M</b>	<b>I</b>	<b>L</b>	<b>V</b>	<b>F</b>	<b>Y</b>	<b>W</b>

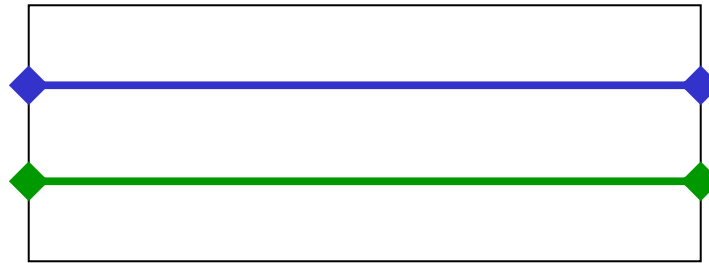
- biochemical properties
- mutation prob (evolution)

# alignment

CAGCACTTGGATTCTCGG  
CAGC-----G-T-----GG

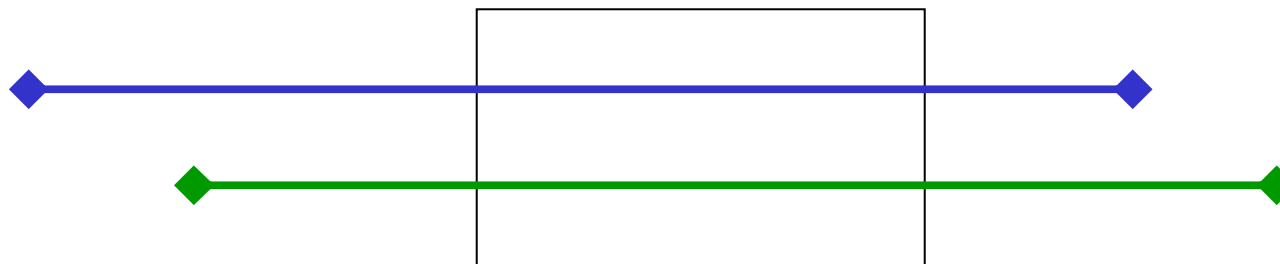
CAGCA-CTTGGATTCTCGG  
---CAGCGTGG-----

# variants of alignment



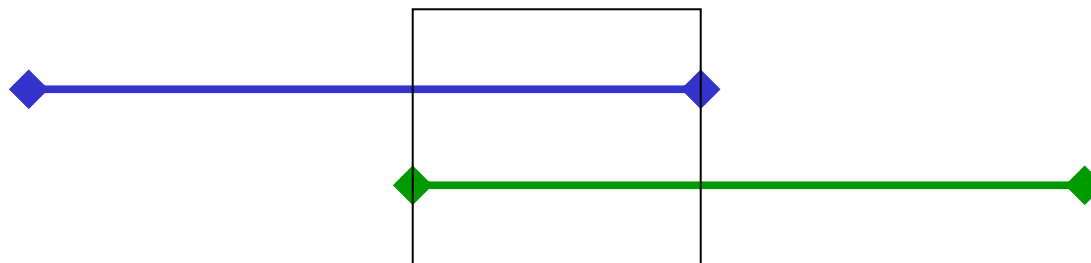
> global

*Needleman & Wunsch*

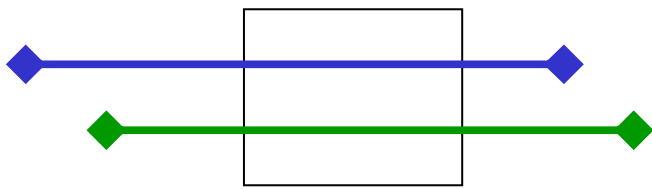


local

*Smith & Waterman*

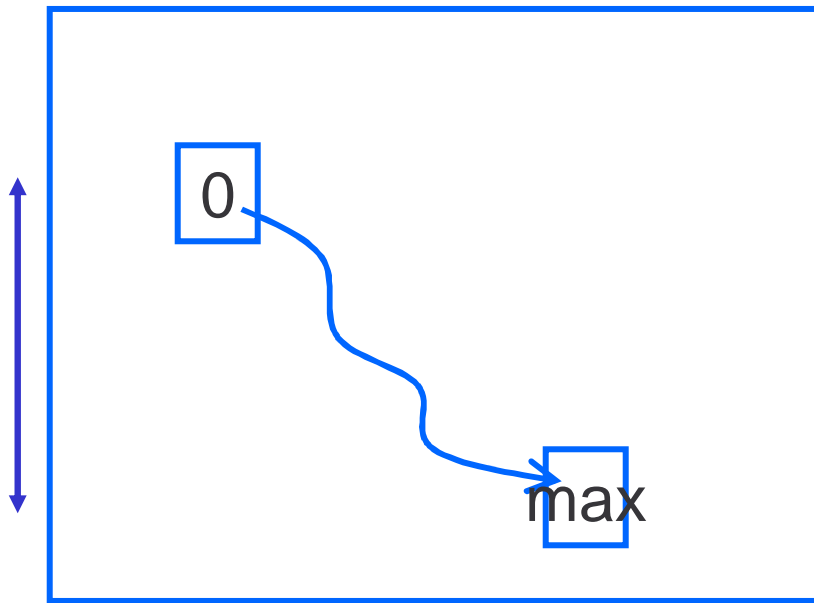


semi-global  
end



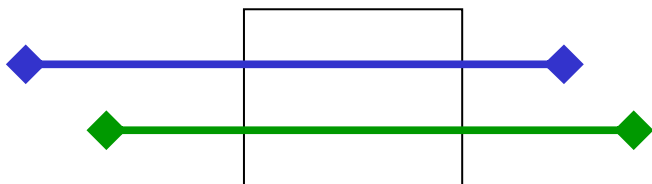
# local alignment

$$a[i, j] = \max \begin{cases} a[i, j-1] + g \\ a[i, j] + \sigma( s[i], t[j] ) \\ a[i-1, j] + g \\ 0 \end{cases}$$

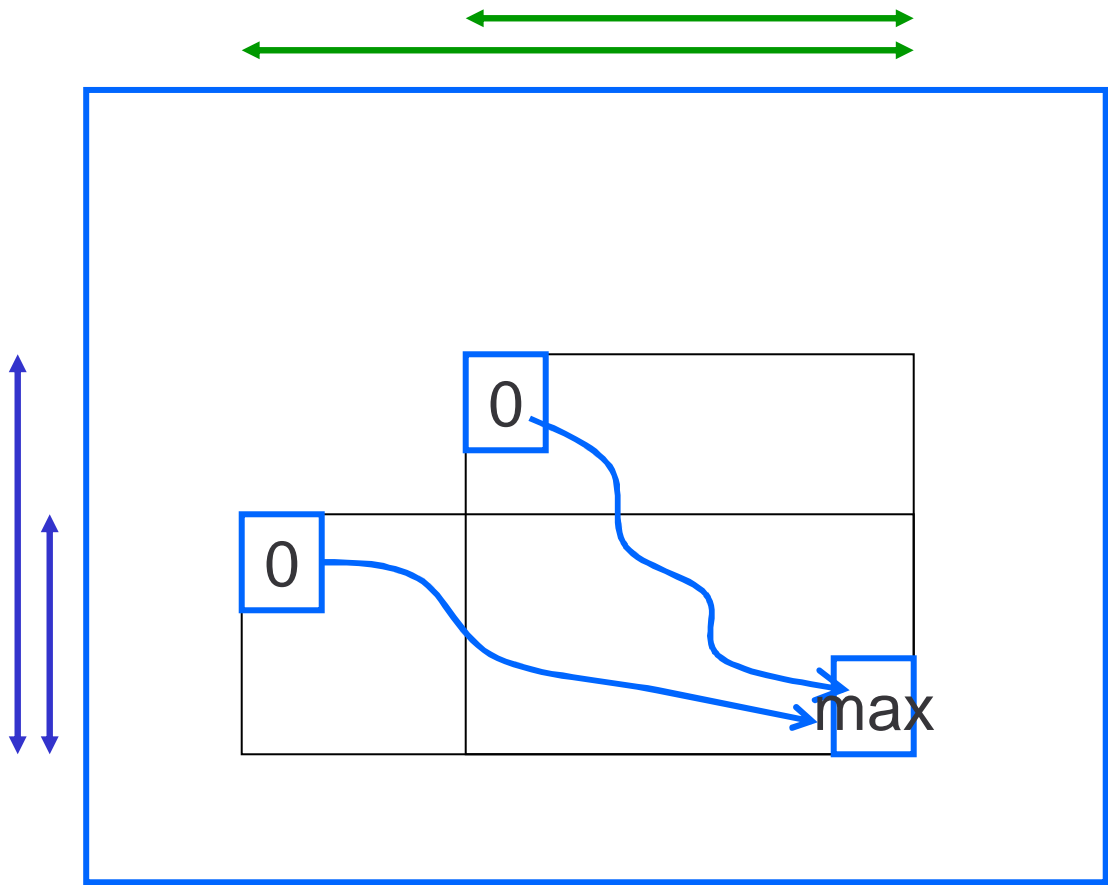


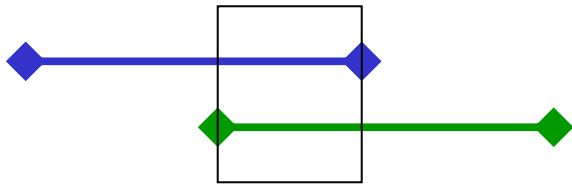
solution (*traceback*)  
from max to 0



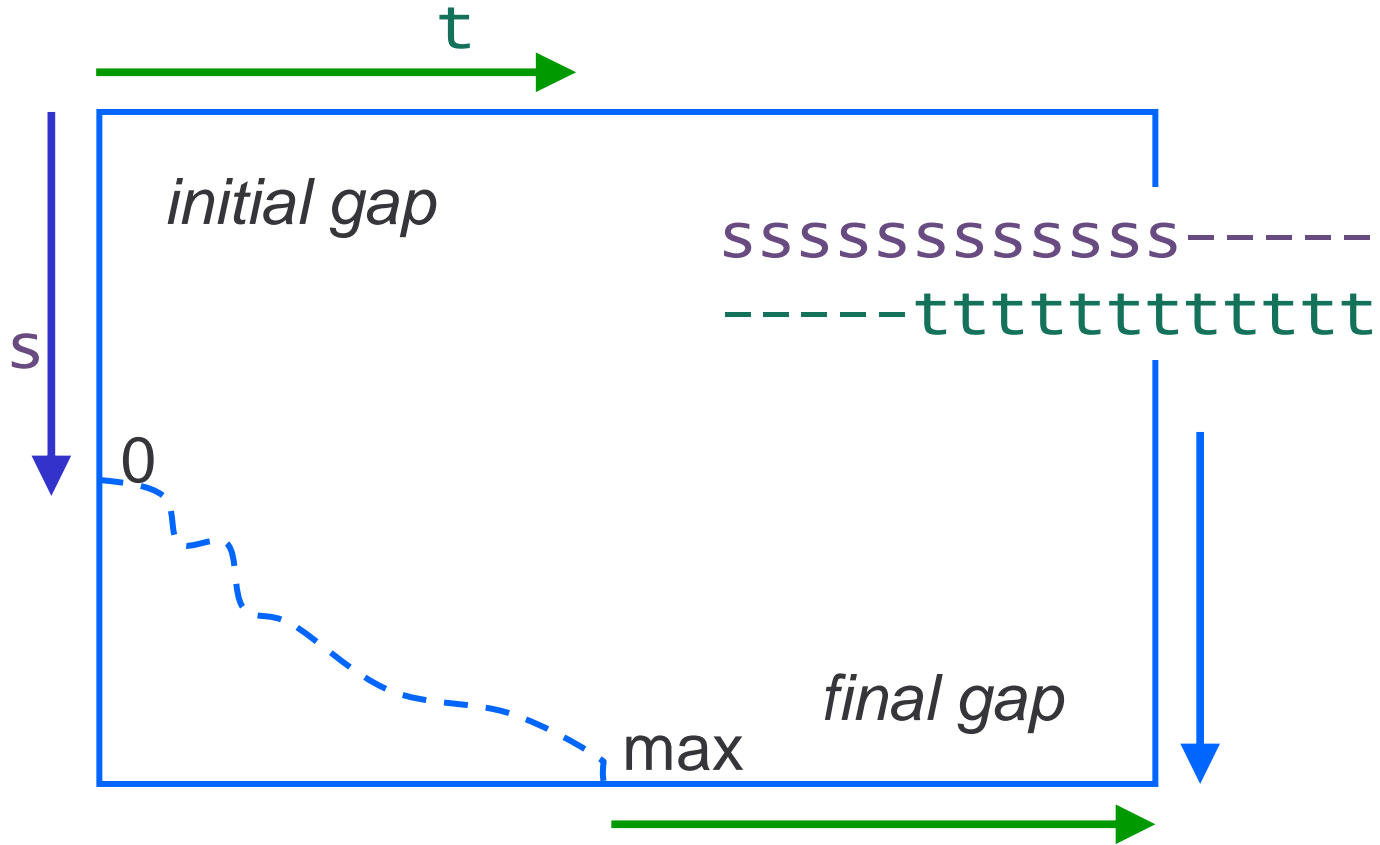


# local alignment





end alignment

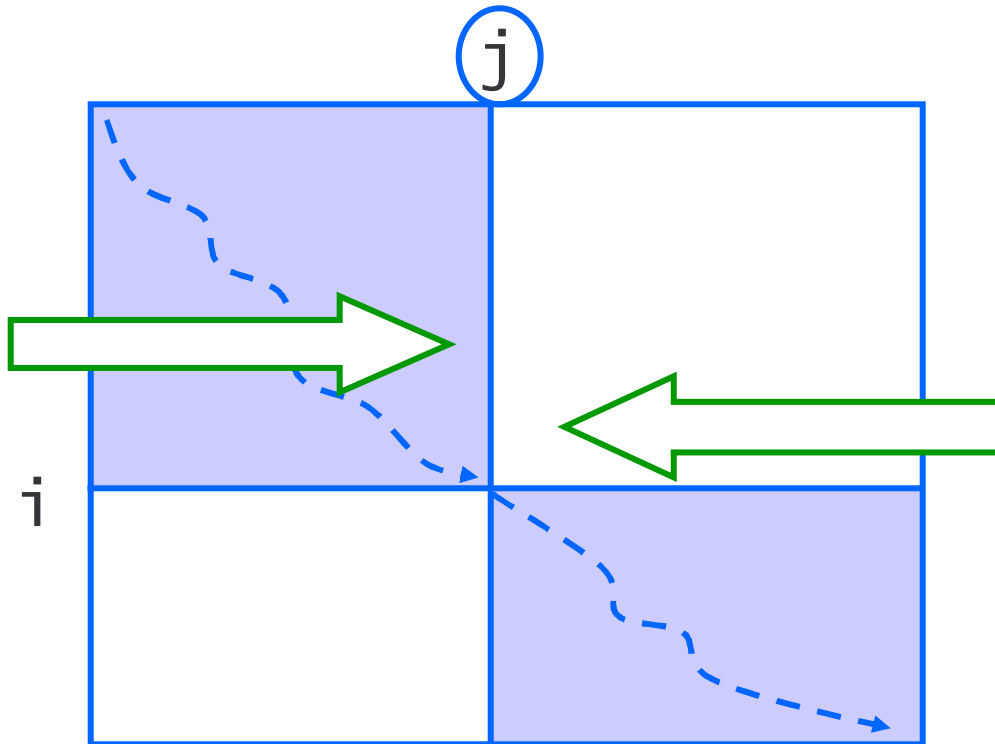


initialization with zero's  
max on border(s)

# saving space

$O(mn)$  both space & time

Hirschberg: linear space



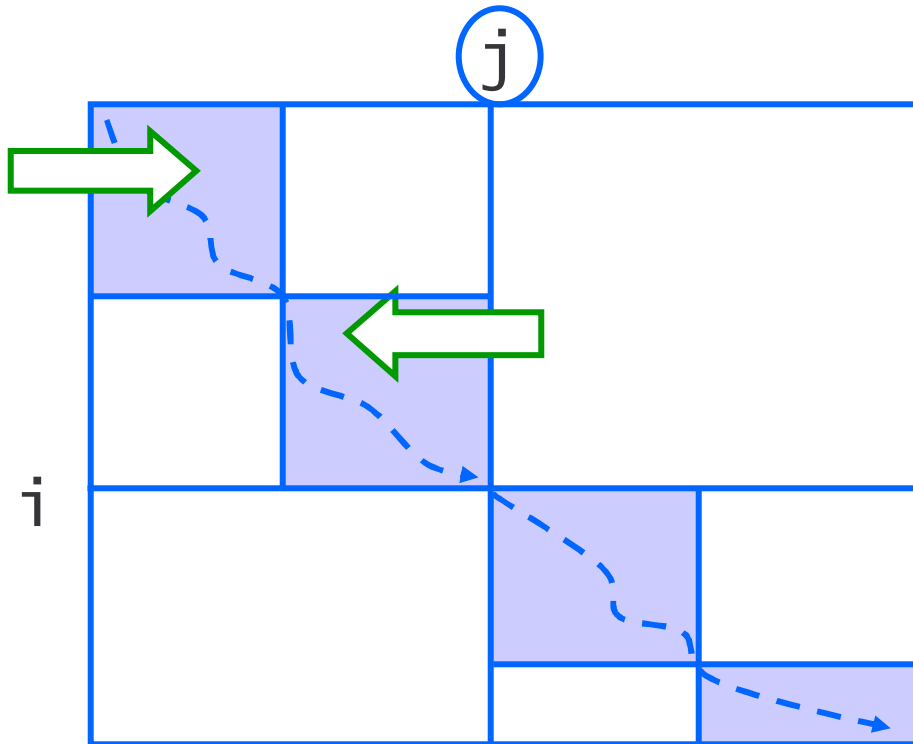
$$T(m,n) \leq 2mn$$

$$\begin{aligned} T(m,n) &\leq mn/2 + mn/2 \\ &\quad + T(i,n/2) + T(n-i,n/2) \\ &\leq mn + in + (m-i)n = 2mn \end{aligned}$$

# saving space

$O(mn)$  both space & time

Hirschberg: linear space



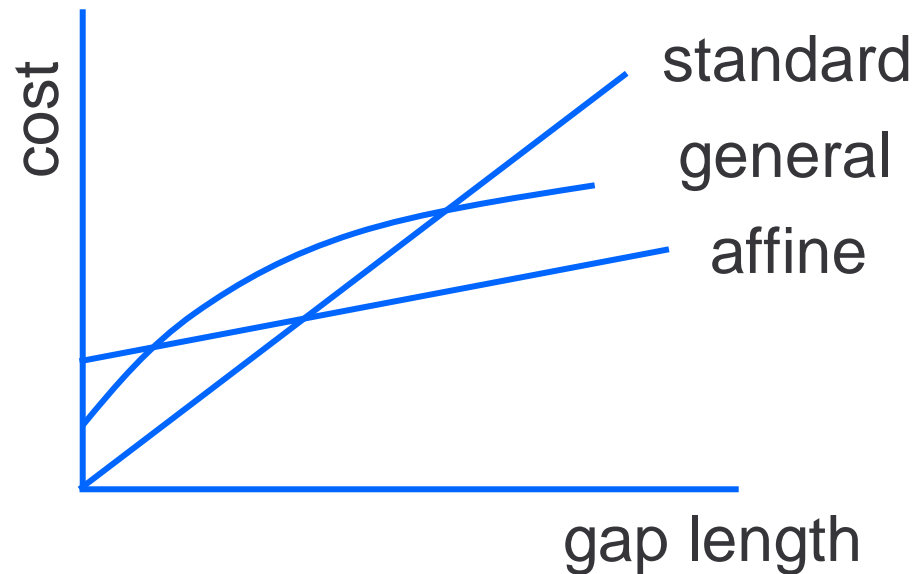
$$1 + 1/2 + 1/4 + \dots = 2$$

# general gap penalties

$g(k) = kg$  gap penalty

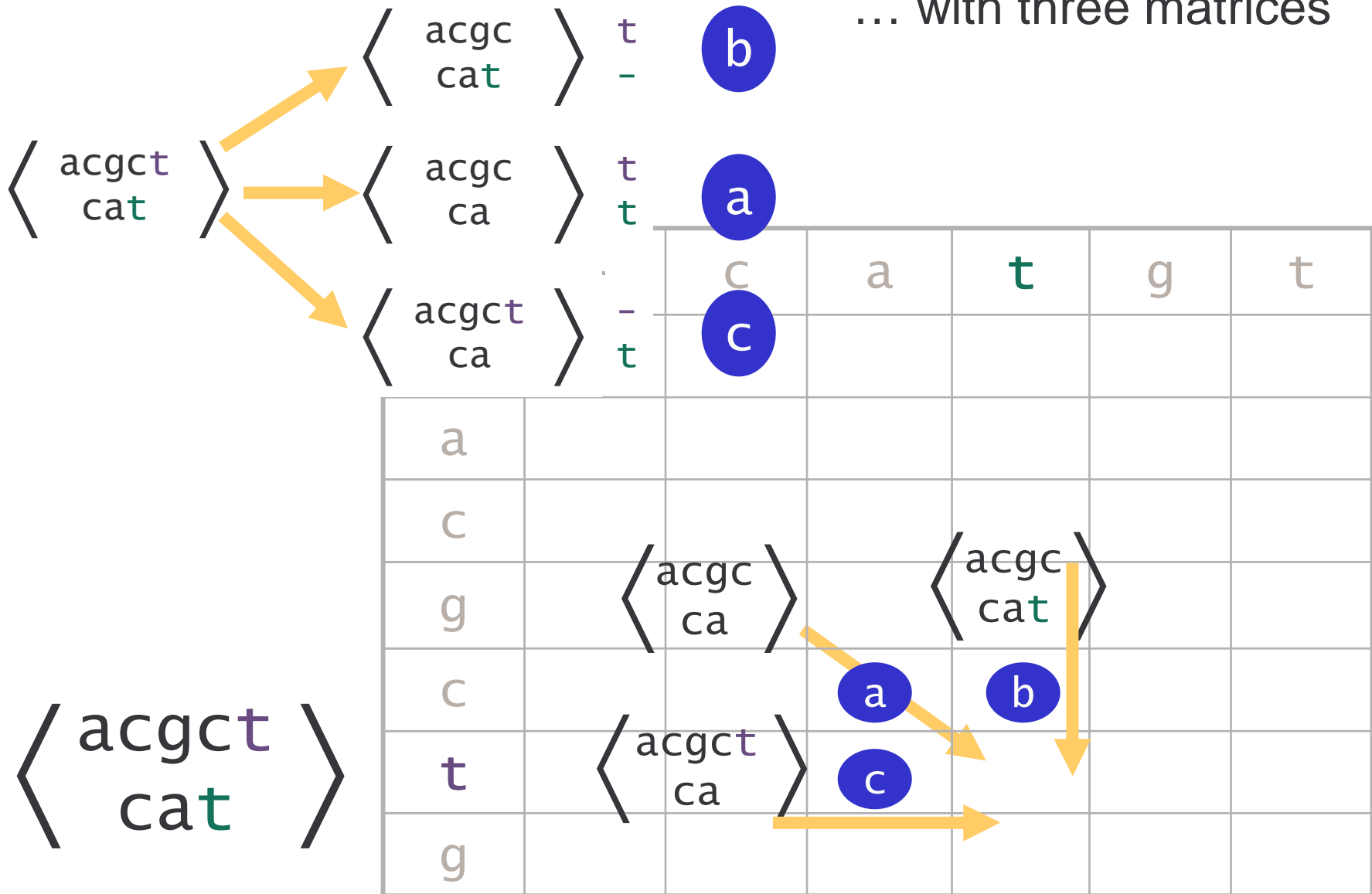
$g(k) = g + ke$  open gap, extend gap

$g(k)$  arbitrary



# dynamic programming

... with three matrices

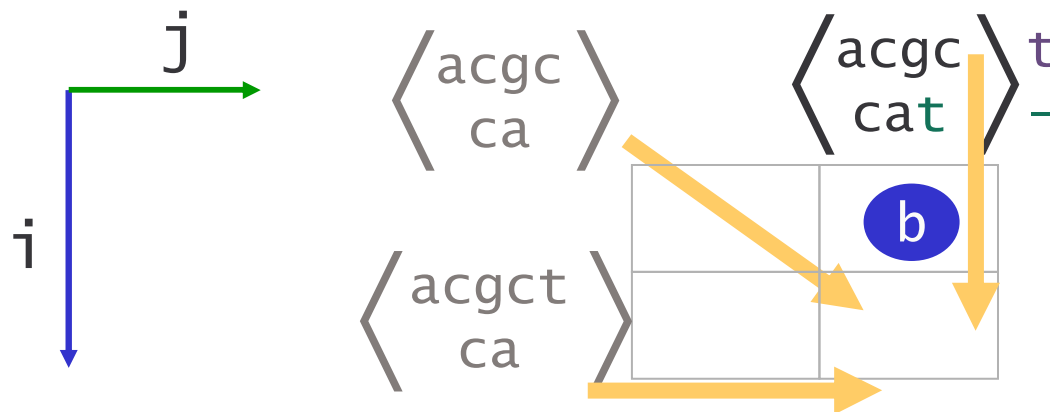


# affine gap penalties

dynamic programming  
with three matrices

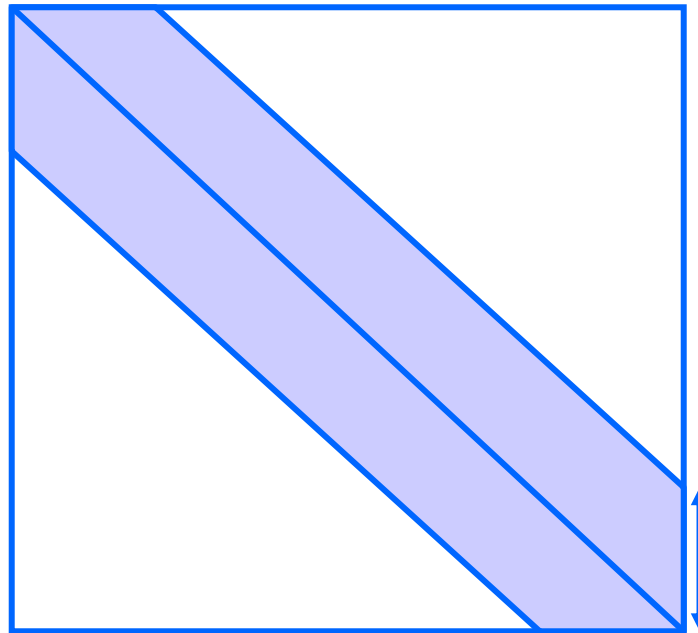
$g(k) = g + ke$       open gap, e extend gap

$$b[i, j] = \max \begin{cases} a[i-1, j] - g - e \\ b[i-1, j] - e \\ c[i-1, j] - g - e \end{cases} \quad \leftarrow \text{gap was already open}$$



$O(mn)$  vs.  
 $O(mn^2 + m^2n)$   
general gap penalty

# comparing similar sequences



$O(mn)$  *quadratic*  
both space & time  
reduce time

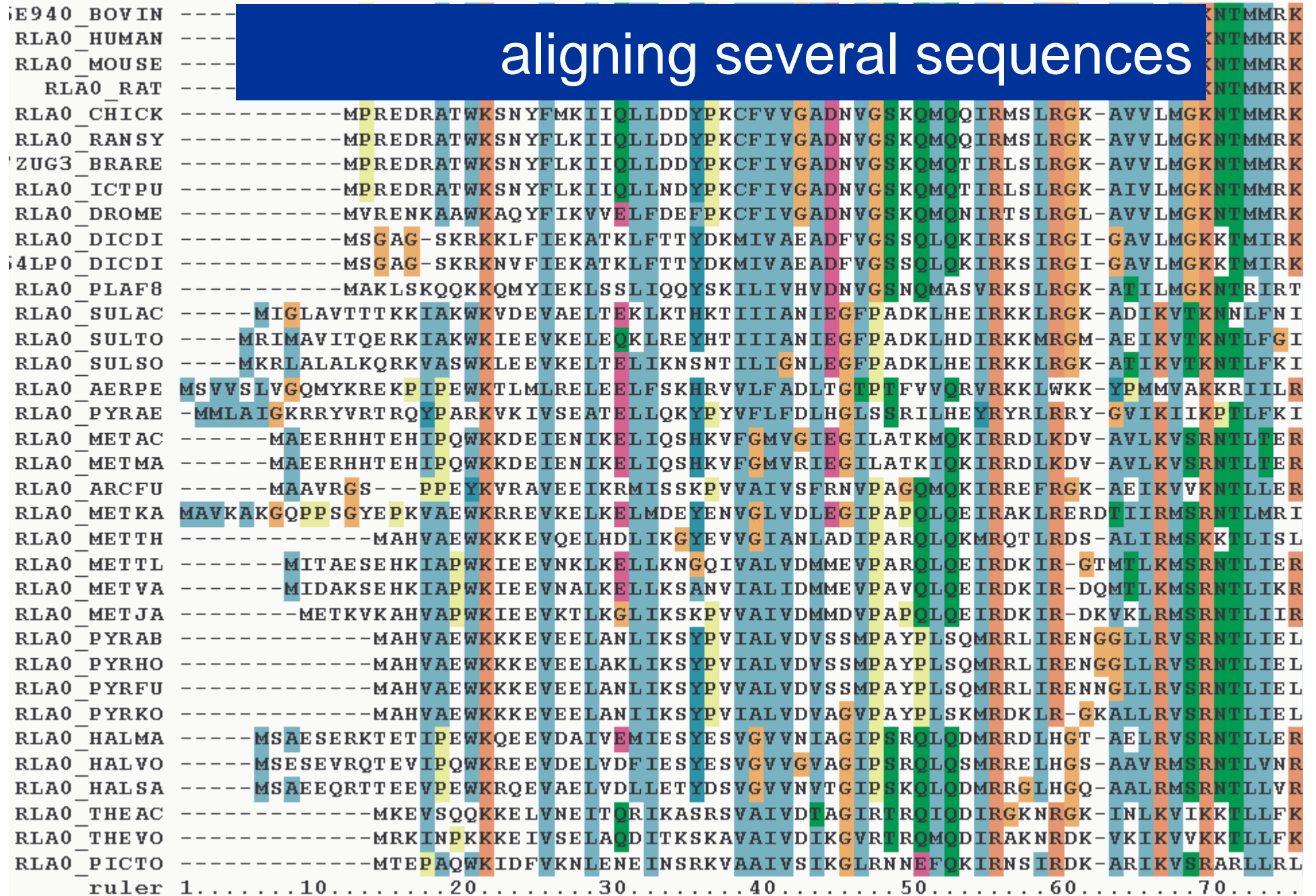
$k$  depends on  
'number of dashes'

banded alignment  
*linear time*

- heuristic
- set  $k$  based on upper bound for alignment
- exact  $\rightarrow$  while improvement do  $k = 2k$



# aligning several sequences



acidic ribosomal protein P0 homolog (L10E) encoded by the Rplp0 gene

# comparing multiple sequences

details → notes APG

how to score?

## SP sum-of-pairs

- pairs of symbols
- pairs of strings

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

## dynamic programming

possible ... but too much space & time

heuristic reduce search space (banded alignment)

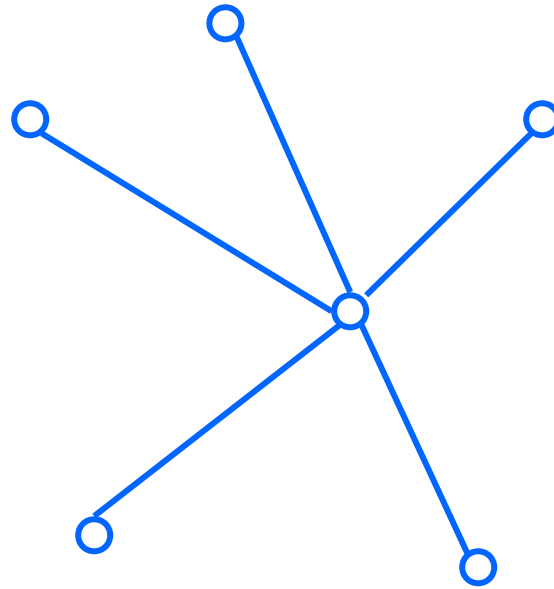
problem is NP-complete (# strings = dimension!)

## heuristic

star alignment / progressive alignment

hidden Markov model

# star alignment



string in center (which one?)  
compute pair-wise alignments  
extend pairs into multiple alignment  
“once a gap, always a gap”

$O(k^2m^2)$        $k$  sequences,  $m$  length

theory: within factor 2 of optimal alignment

# star alignment

1 ATTGCCATT  
 2 ATGGCCATT  
 3 ATCCAATTTT  
 4 ATCTTCTT  
 5 ACTGACC

	1	2	3	4	5	$\Sigma$
1		7	-2	0	-3	2
2			-2	0	-4	1
3				0	-7	-11
4					-3	-3
5						-17

1 ATTGCCATT  
 2 ATGGCCATT

1 ATTGCCATT--  
 3 ATC-CAATTTT

1 ATTGCCATT  
 4 ATCTTC-TT

1 ATTGCCATT  
 5 ACTGACC--

1 ATTGCCATT--

2 ATGGCCATT--

3 ATC-CAATTTT

1 ATTGCCATT--

2 ATGGCCATT--

3 ATC-CAATTTT

4 ATCTTC-TT--

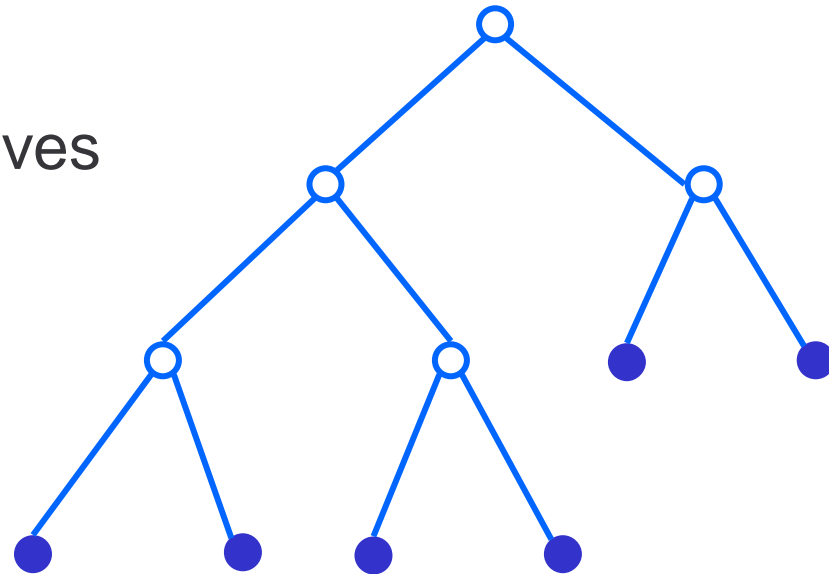
5 ACTGACC----

# progressive alignment

*guide tree*

bottom-up

sequences in leaves



- tree known ← evolutionary relation
- build tree ← clustering algo's ClustalW

inner nodes

combine alignments (profiles variant basic algo)



problem solved !?

## **much too slow**

- long strings
- huge databases

## **heuristics**

- **FASTA** along diagonal
- **BLAST** minimal close match

## **multiple alignment**

(several strings)

NP complete ... exponential

## statistics

- 23 pairs of chromosomes
- $3.1 \times 10^9$  nucleotide bases
- average 3000 bases / gene
  - dystrophin 2.4 million
- 30,000 - 40,000 genes
- protein variants 1 million (splicing)
- 99.9% exactly the same in humans

# Basic Local Alignment Search Tools

query

query word (W=3)

GSVEDTTGSQSLAALLNKCKT **PQG** QRLVNQWIKQPLMDKNRIEERLNL

neighbourhood  
words

<b>PQG</b>	18
PEG	15
PRG	14
PMG	13
PQA	12
PQN	12

threshold

SLAALLNKCKT **PQG** QRLVNQWIKQPLMDKNRIEERLNL  
 TLASVLDCTVT **PMG** SRMLKRWLHMPVRDTRVLLERQQT

subject

high-scoring segment pair (HSP)

parameters:  
 W word size  
 T threshold  
 S score  
 E expected