

4.xi.

AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCGGTAGGAT
AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGATGCG

1. direct optimization

2. summary

Direct optimization (DO)

AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGAT
AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGATGCG

ACG

ACGTGC

CT

ACT

AACA

?

?

?

?

GA	R
GT	K
GC	S
GAT	D
GAC	V
GTC	B
ATC	H
AT	W
AC	M
CT	Y
ACGT	N
ACGT-	X

sequences UNEQUAL in length

Sequence alignment

AACGGTTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTT
AACGGTTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTT

ONLY LOCAL
comparison

BLAST, comparison of 2 sequences

MULTIPLE sequence alignment (MSA)

GLOBAL
comparison

mostly done PROGRESSIVELY

e.g. CLUSTAL

computationally
demanding

1) pairwise sequence comparison

2) pairwise distance matrix

3) **guide** tree formed based on this matrix

4) **pair-wise alignment** performed based on this tree
connecting seqs. of all the terminals included

in later developed programs (e.g. muscle, MAFFT)

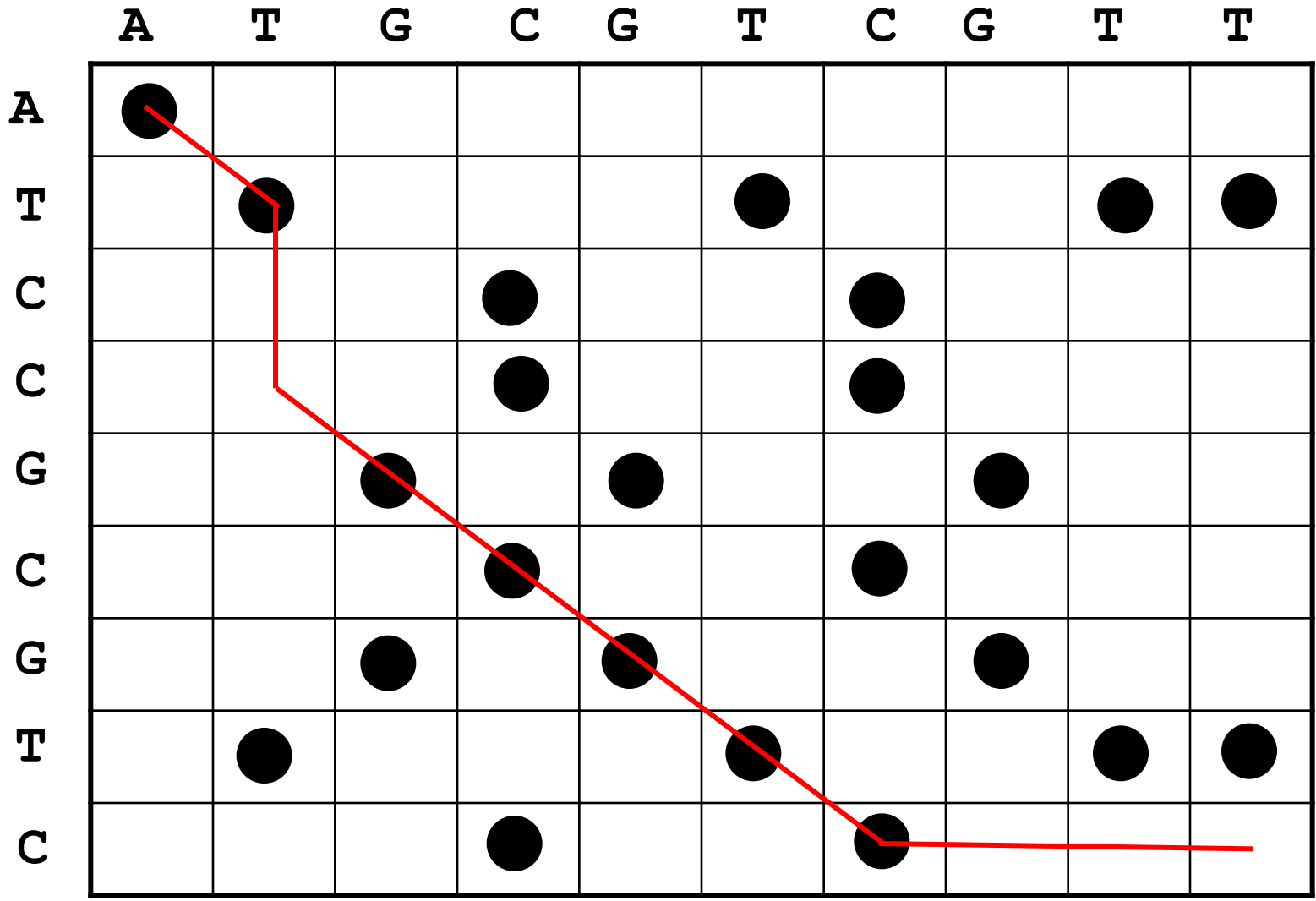
repeated cycles of these operations, with final round of
edge refinement using sum of pairs

Pair-wise alignment

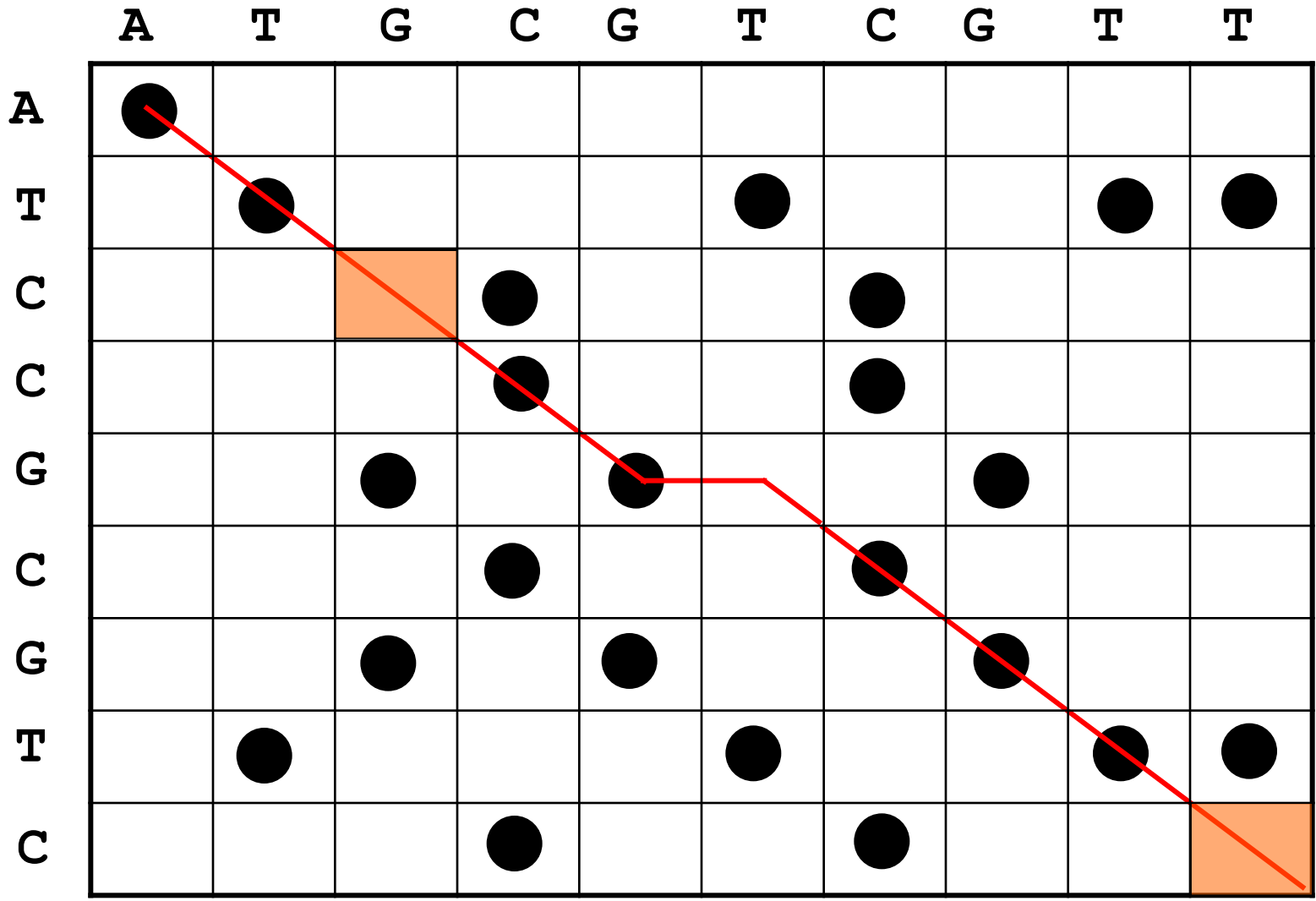
```
AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCGGTAGGAT  
AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGATGCG
```

dot-plot method

simple visualization of matches



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



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

Needleman-Wunsch algorithm

AACGGTTAAGGTACGGAGAATTAGGCAACCCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCGGTAGGAT

Needleman, S.B. & Wunsch, C.D. 1970. A general method applicable to the search of similarities in the amino acid sequence of two proteins. *J. Mol. Biol.* 48: 443-453.

4 basic steps:

1) laying out the alignment matrix

2 sequences define the axes

2) initiation of the matrix

3) wave-front update of the matrix elements

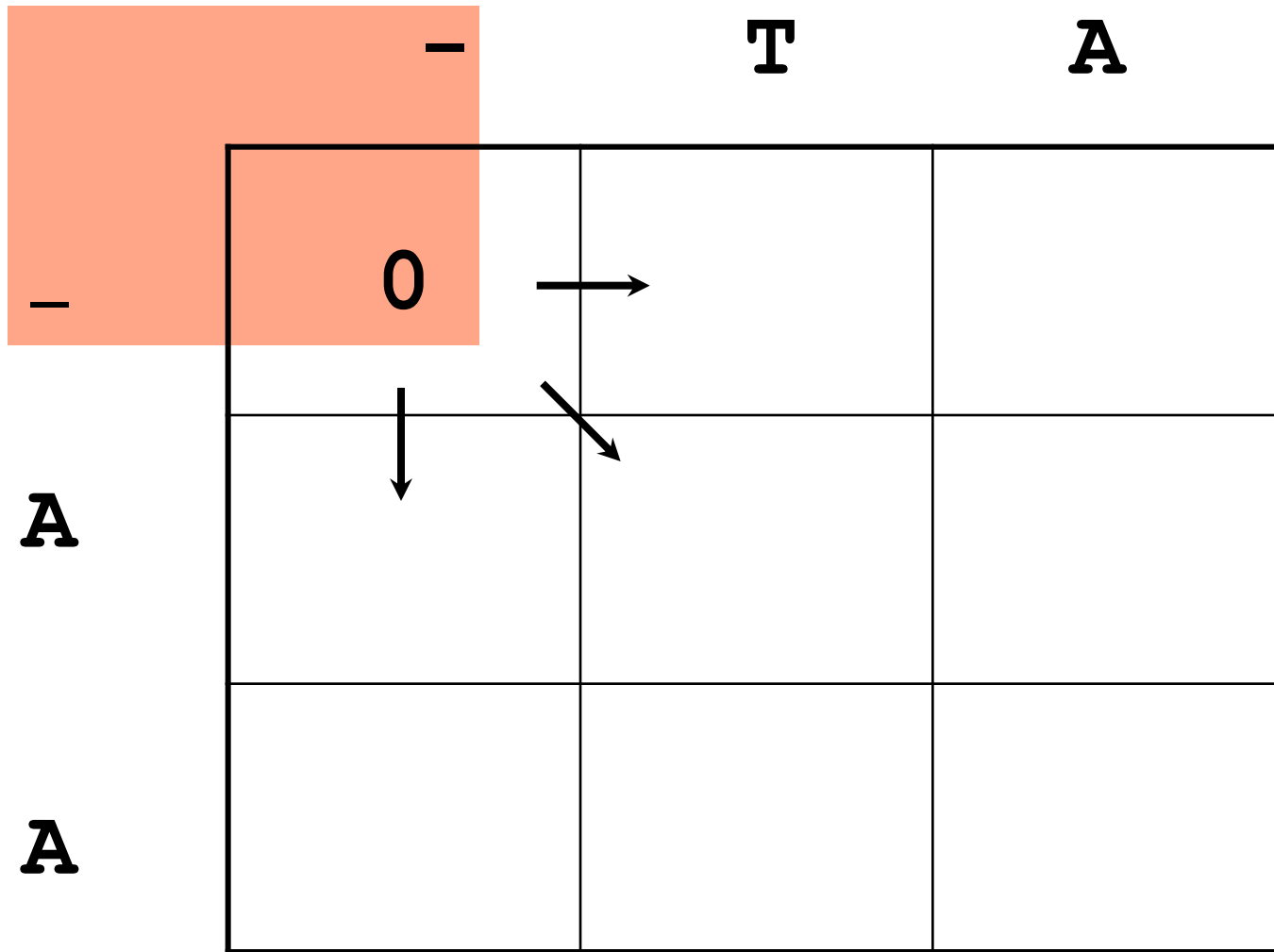
4) trace back

- T A A A T T G C A

-	0									
A		1	0	0	0	1	1	1	1	0
A		1	0	0	0	1	1	1	1	0
T		0	1	1	1	0	0	1	1	1
T		0	1	1	1	0	0	1	1	1
T		0	1	1	1	0	0	1	1	1
G		1	1	1	1	1	1	0	1	1
G		1	1	1	1	1	1	0	1	1
G		1	1	1	1	1	1	0	1	1
C		1	1	1	1	1	1	1	0	1
C		1	1	1	1	1	1	1	0	1
A		1	0	0	0	1	1	1	1	0

- T A A A T T G C A

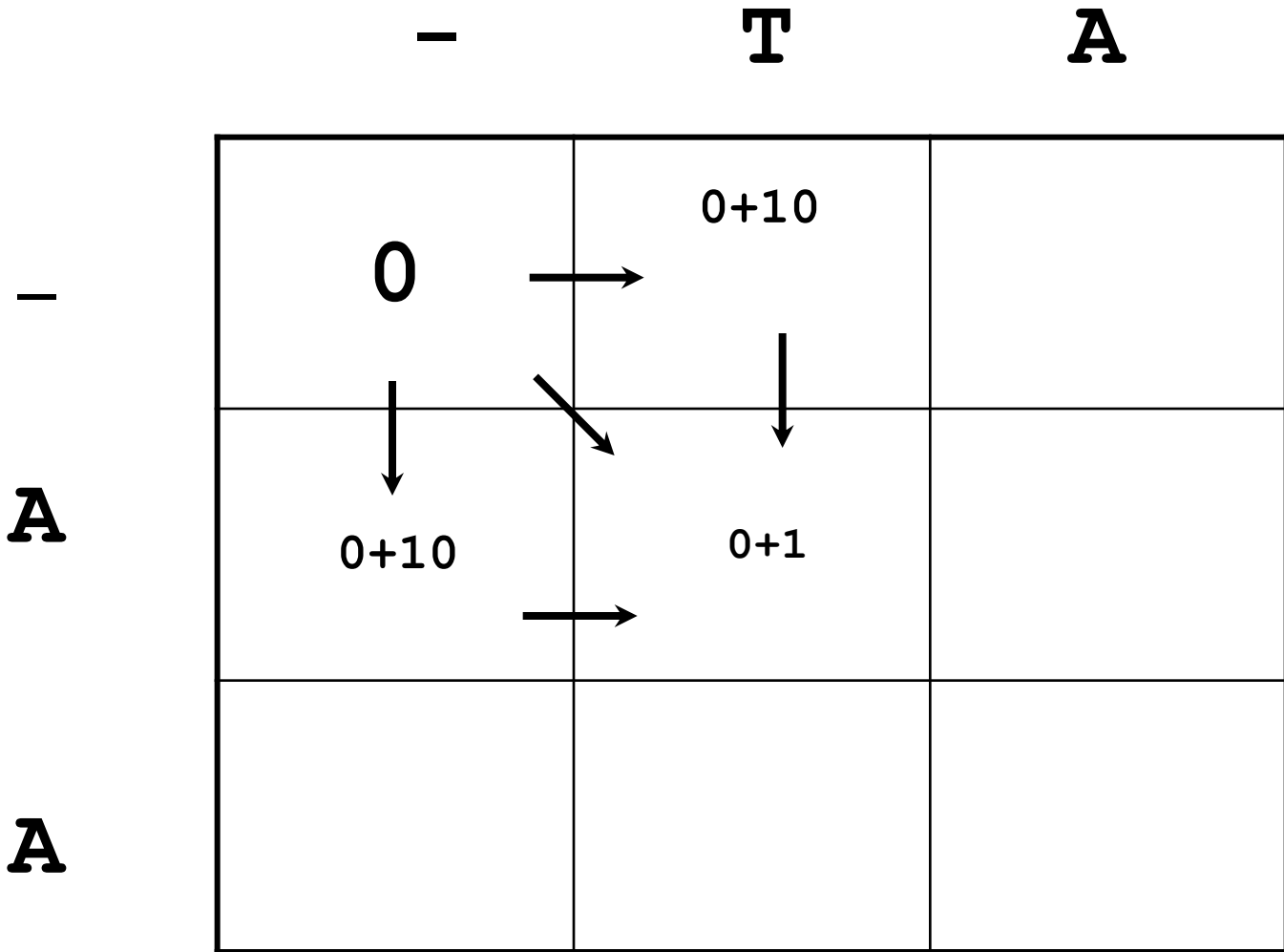
-	0								
A	1	0	0	0	1	1	1	1	0
A	1	0	0	0	1	1	1	1	0
T	0	1	1	1	0	0	1	1	1
T	0	1	1	1	0	0	1	1	1
T	0	1	1	1	0	0	1	1	1
G	1	1	1	1	1	1	0	1	1
G	1	1	1	1	1	1	0	1	1
G	1	1	1	1	1	1	0	1	1
C	1	1	1	1	1	1	1	0	1
C	1	1	1	1	1	1	1	0	1
A	1	0	0	0	1	1	1	1	0



match = 0

substitution = 1

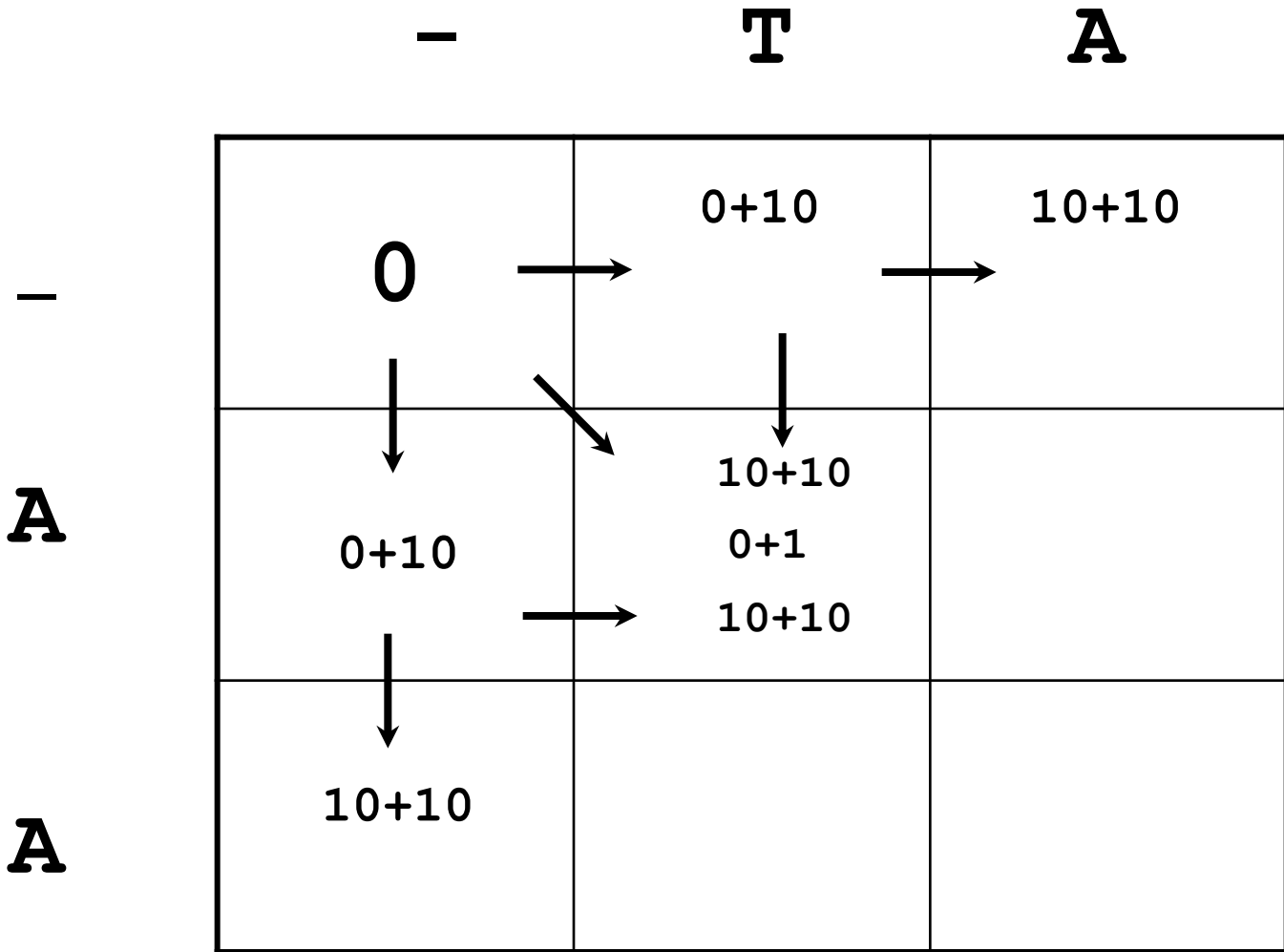
gap penalty = 10



match = 0

substitution = 1

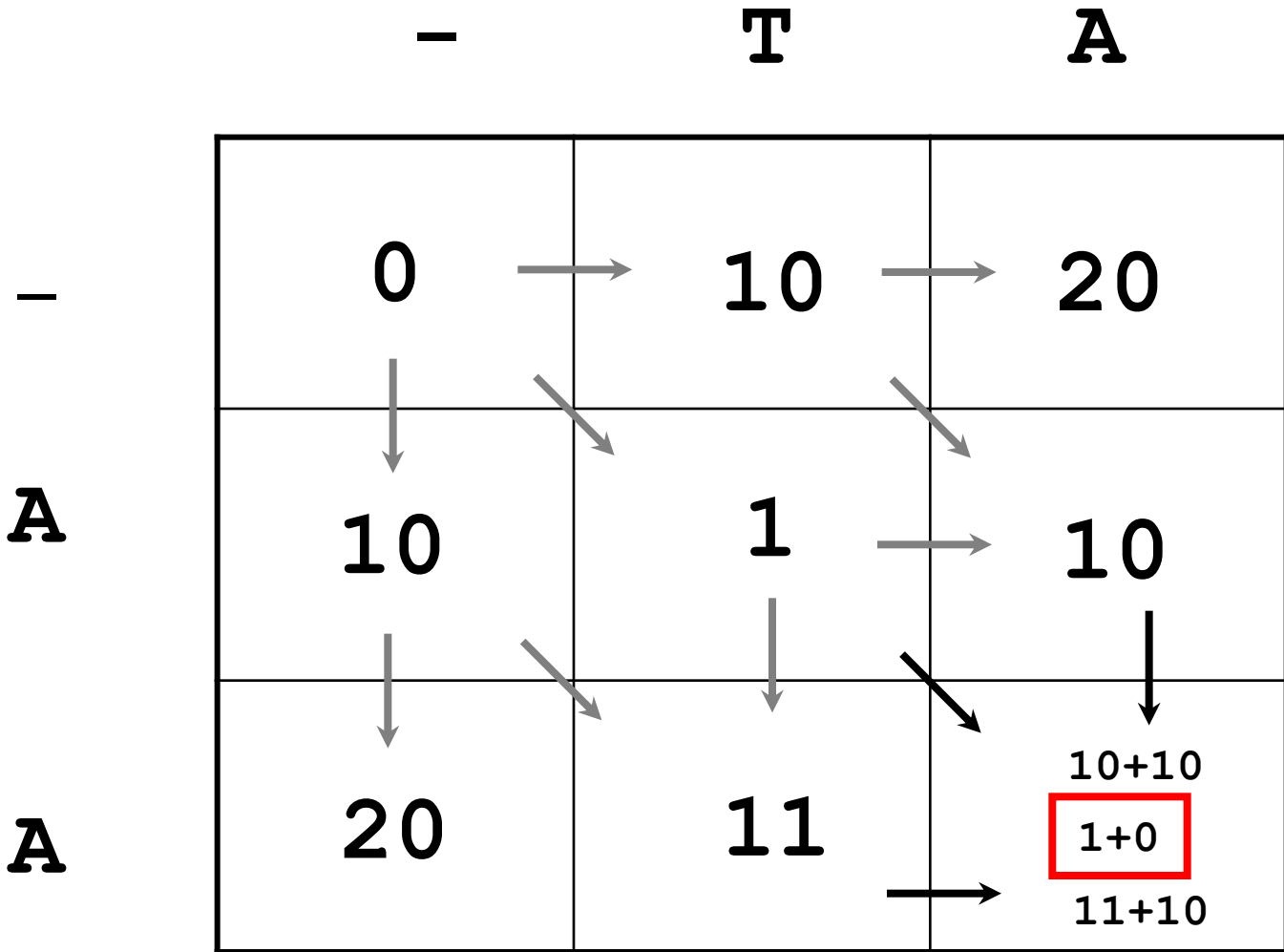
gap penalty = 10

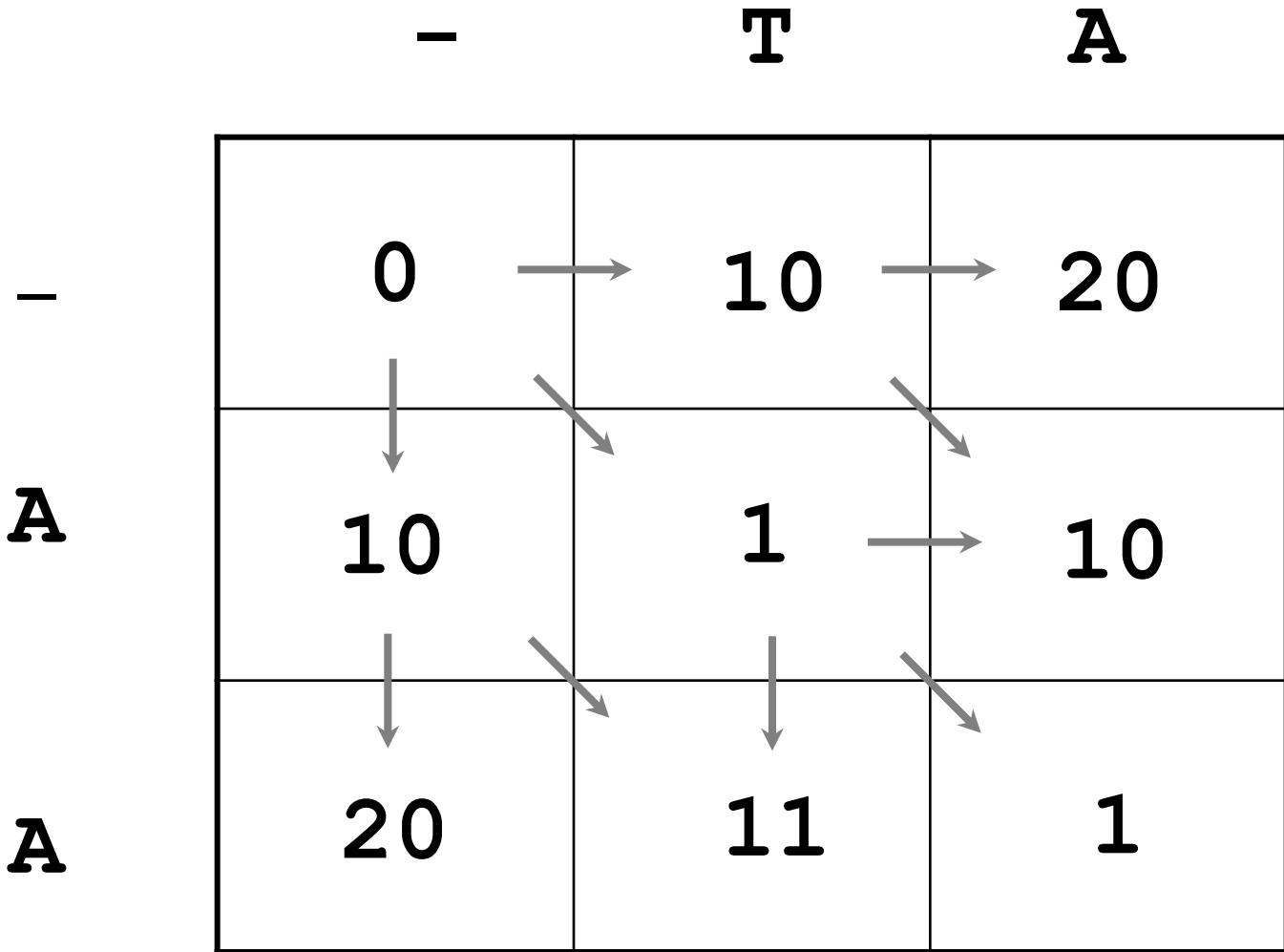


match = 0

substitution = 1

gap penalty = 10





- T A A A T T G C A

-	0	→10	→20	→30	→40	→50	→60	→70	→80	→90
A	↓10	1	10	→20	→30	→40	→50	→60	→70	→80
A	↓20	↓11	2	10	→20	→30	→40	→50	→60	→70
T	↓30	↓20	↓11	2	11	20	→30	→40	→50	→60
T	↓40	↓30	↓21	↓12	3	11	→20	→30	→40	→50
T	↓50	↓40	↓31	↓22	↓13	3	11	→21	→31	→41
G	↓60	↓50	↓41	↓32	↓23	↓13	4	11	→21	→31
G	↓70	↓60	↓51	↓42	↓33	↓23	↓14	4	12	→22
G	↓80	↓70	↓61	↓52	↓43	↓33	↓24	↓14	5	13
C	↓90	↓80	↓71	↓62	↓53	↓43	↓34	↓24	↓14	6
C	↓100	↓90	↓81	↓72	↓63	↓53	↓44	↓34	↓24	15
A	↓110	↓100	↓90	↓81	↓72	↓63	↓54	↓44	↓34	24

Needleman-Wunsch algorithm

```
AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGAT  
AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGATG
```

4 basic steps:

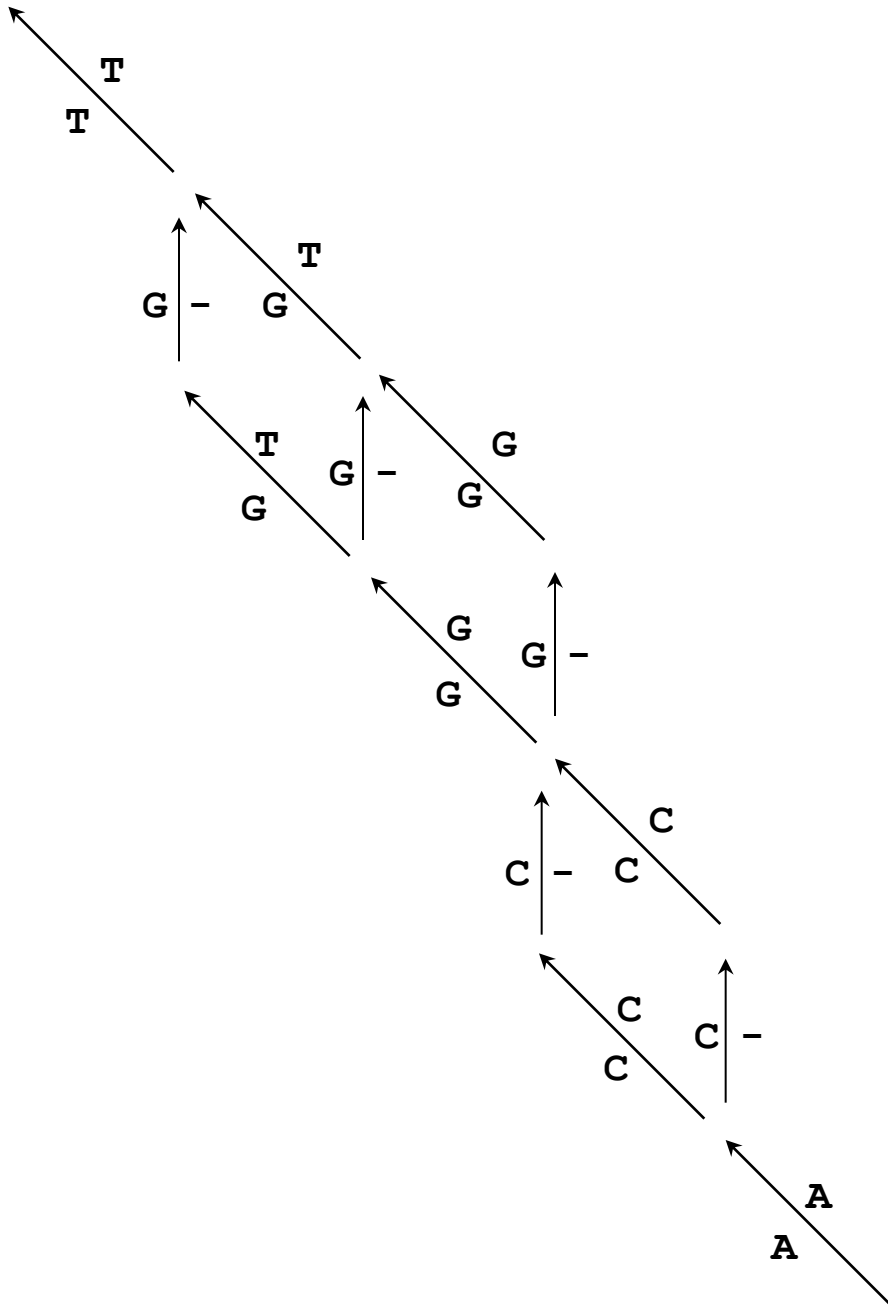
- 1) laying out the alignment matrix
- 2) initiation of the matrix
- 3) wave-front update of the matrix elements
- 4) trace back

- T A A A T T G C A

-	0	10	20	30	40	50	60	70	80	90
A	10	1	10	20	30	40	50	60	70	80
A	20	11	2	10	20	30	40	50	60	70
T	30	20	11	2	11	20	30	40	50	60
T	40	30	21	12	3	11	20	30	40	50
T	50	40	31	22	13	3	11	21	31	41
G	60	50	41	32	23	13	4	11	21	31
G	70	60	51	42	33	23	14	4	12	22
G	80	70	61	52	43	33	24	14	5	13
C	90	80	71	62	53	43	34	24	14	6
C	100	90	81	72	63	53	44	34	24	15
A	110	100	90	81	72	63	54	44	34	24

- T A A A T T G C A

-	0	10	20	30	40	50	60	70	80	90
A	10	1	10	20	30	40	50	60	70	80
A	20	11	2	10	20	30	40	50	60	70
T	30	20	11	2	11	20	30	40	50	60
T	40	30	21	12	3	11	20	30	40	50
T	50	40	31	22	13	3	11	21	31	41
G	60	50	41	32	23	13	4	11	21	31
G	70	60	51	42	33	23	14	4	12	22
G	80	70	61	52	43	33	24	14	5	13
C	90	80	71	62	53	43	34	24	14	6
C	100	90	81	72	63	53	44	34	24	15
A	110	100	90	81	72	63	54	44	34	24



TTG-C-A
TGGGCCA

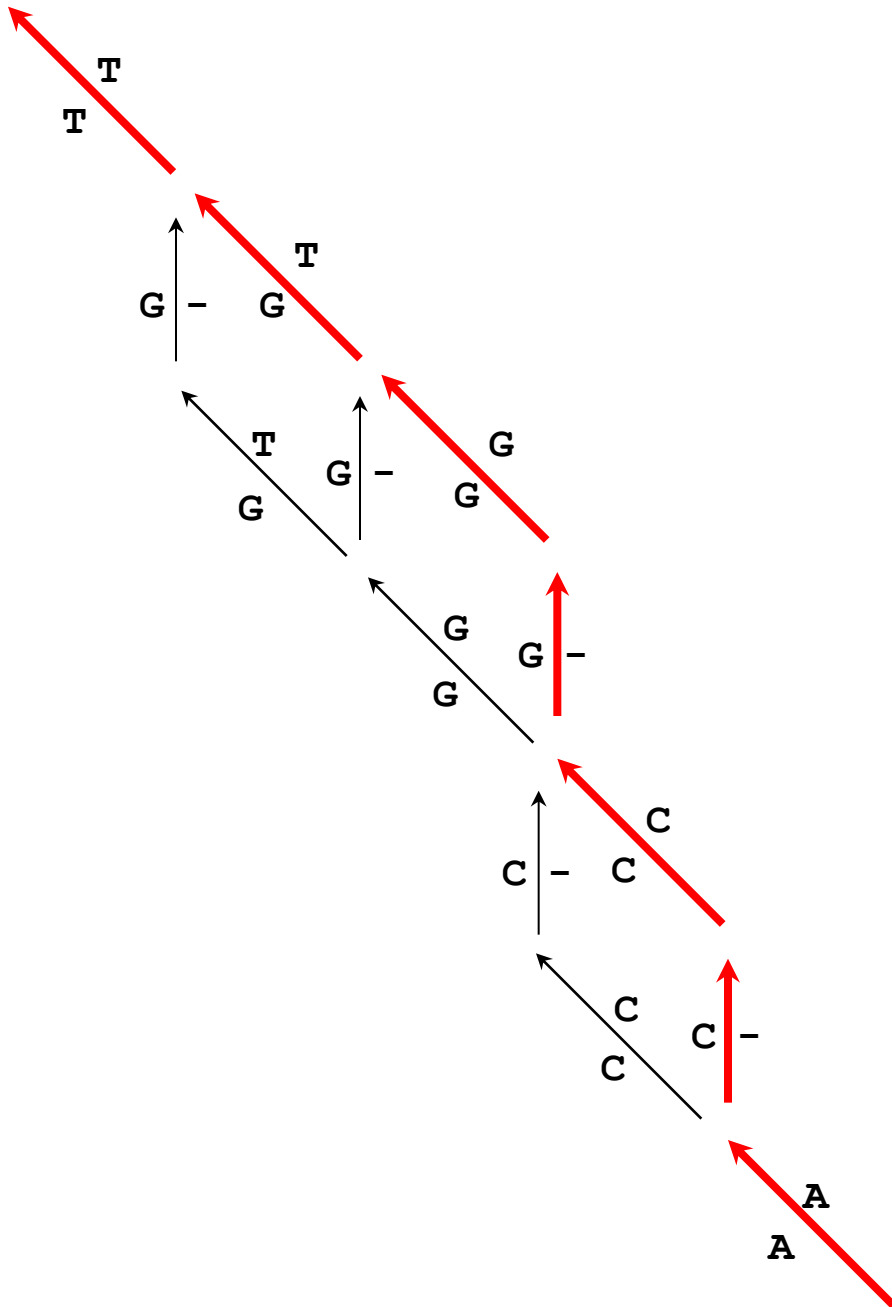
TTG--CA
TGGGCCA

TT-G-CA
TGGGCCA

TT-GC-A
TGGGCCA

T-TGC-A
TGGGCCA

T-TG-CA
TGGGCCA



TTG-C-A
TGGGCCA

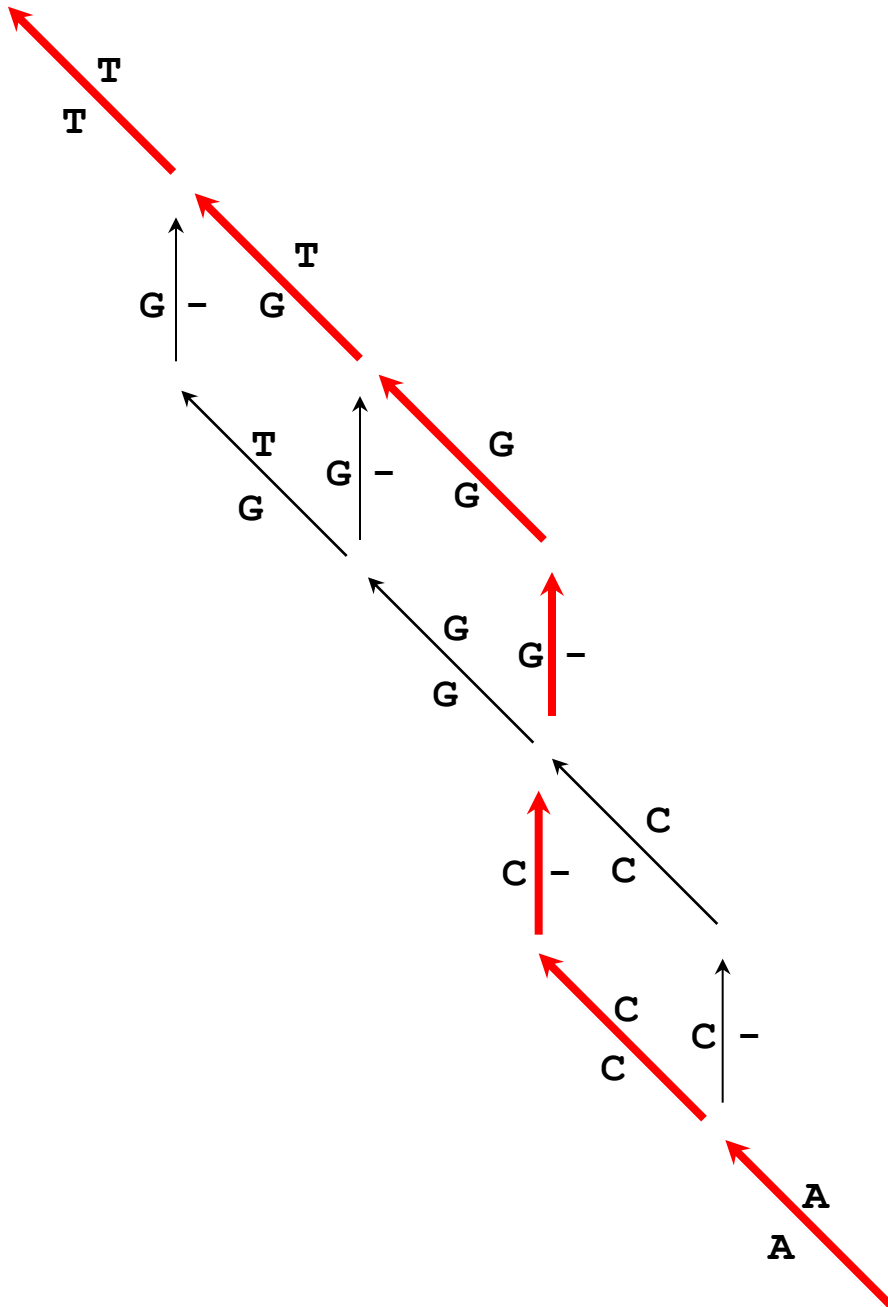
TTG--CA
TGGGCCA

TT-G-CA
TGGGCCA

TT-GC-A
TGGGCCA

T-TGC-A
TGGGCCA

T-TG-CA
TGGGCCA



TTG-C-A
TGGGCCA

TTG--CA
TGGGCCA

TT-G-CA
TGGGCCA

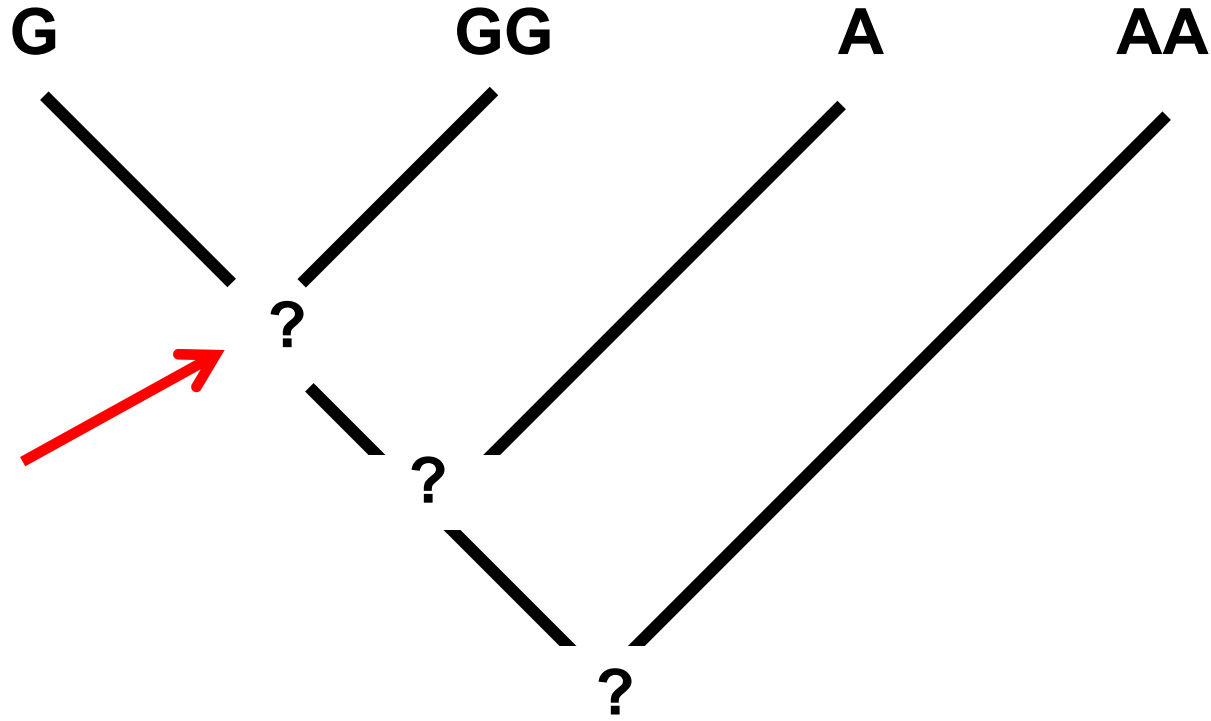
TT-GC-A
TGGGCCA

T-TGC-A
TGGGCCA

T-TG-CA
TGGGCCA

Direct optimization (DO)

AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGAT
AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGATGCG



	A	C	G	T	-	M	R	W	(A)	S	Y	(C)	K	(G)	(T)	V	H	(M)	D	(R)	(W)	B	(S)	(Y)	(K)	N	(V)	(H)	(B)	(D)	X	
A	0	2	1	2	3	0	0	0	0	1	2	2	1	1	2	0	0	0	0	0	0	1	1	2	1	0	0	0	1	0	0	
C	0	2	1	3	0	2	1	2	0	0	0	0	1	2	1	0	0	0	1	2	1	0	0	0	1	0	0	0	0	1	0	
G	0	2	3	1	0	1	1	1	0	2	2	2	0	0	2	0	1	1	0	0	1	0	0	2	0	0	0	1	0	0	0	
T	0	3	1	2	0	2	1	0	1	0	1	1	0	2	0	1	0	1	0	2	0	0	1	0	0	0	1	0	0	0	0	
-	0	3	3	3	0	3	3	0	3	3	0	3	3	0	3	3	0	3	0	0	3	0	0	0	3	0	0	0	0	0	3	
M	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	
R	0	0	0	0	0	0	2	2	0	0	2	2	0	0	2	0	0	0	0	0	0	0	0	2	0	0	0	0	0	0	0	
W	0	1	0	1	0	1	0	0	1	0	1	0	1	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	
(A)	GA					R			0	1	2	0	1	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0
S									0	0	0	0	0	0	1	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0
Y	GT					K					0	0	0	2	0	0	0	0	0	2	0	0	0	0	0	0	0	0	0	0	0	
(C)											0	0	1	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	
K	GC					S							0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	
(G)													0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
(T)	GAT					D									0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
V															0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
H															0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
(M)	GTC					B									0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
D															0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
(R)	ATC					H									0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
(W)															0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
B	AT					W															0	0	0	0	0	0	0	0	0	0	0	
(S)																					0	0	0	0	0	0	0	0	0	0	0	
(Y)	AC					M																		0	0	0	0	0	0	0	0	
(K)																								0	0	0	0	0	0	0	0	
N	CT					Y																		0	0	0	0	0	0	0	0	
(V)	ACGT					N																		0	0	0	0	0	0	0	0	
(H)																									0	0	0	0	0	0	0	
(B)	A-					(A)																								0	0	
(D)	AC-					(M)																								0	0	
X																															0	
																															
	N-					X																										

transitions 1
transversions 2
indels 3

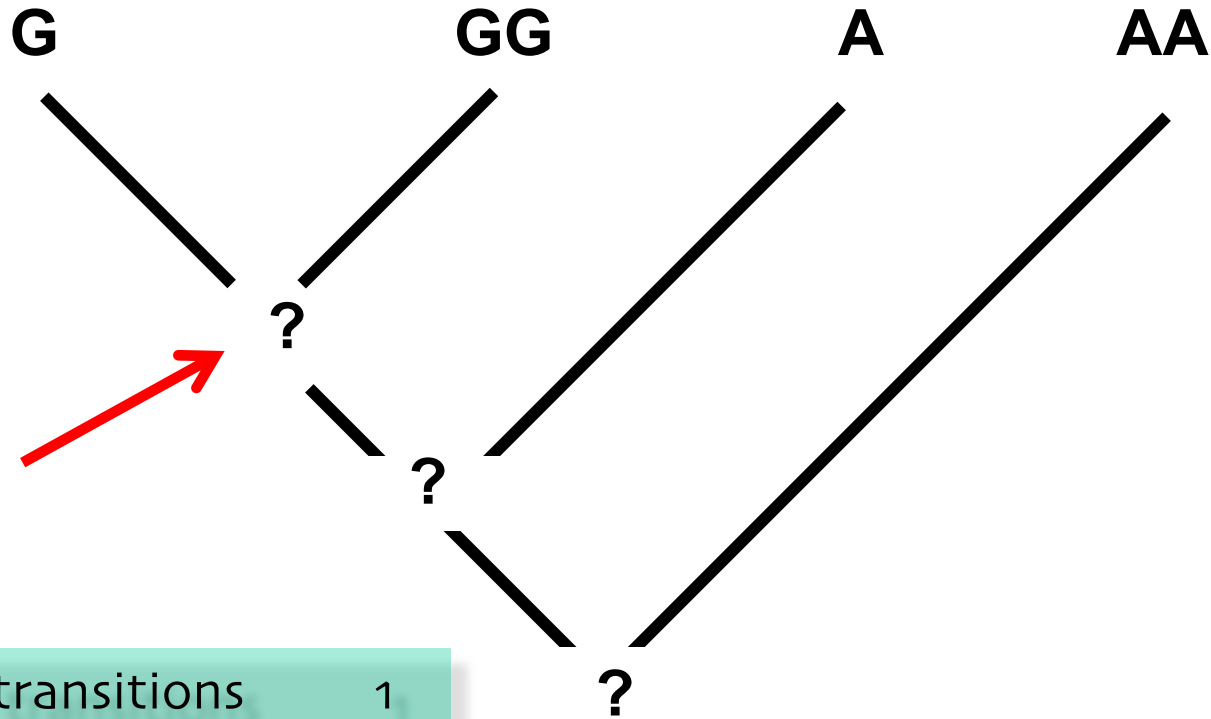
$$31 \times 31 = 961 - 31 = 930 / 2 = 465$$

	A	C	G	T	-	M	R	W	(A)	S	Y	(C)
A	0	2	1	2	3	0	0	0	0	1	2	2
C		0	2	1	3	0	2	1	2	0	0	0
G	GA	R	0	2	3	1	0	1	1	0	2	2
T	GT	K		0	3	1	2	0	2	1	0	1
-	GC	S			0	3	3	3	0	3	3	0
M	GAT	D										
R	GAC	V				0	0	0	0	0	0	0
W	GTC	B										
(A)	ATC	H					0	0	0	0	2	2
S	AT	W						0	0	0	1	0
Y	AC	M							0	1	2	0
(C)	CT	Y								0	0	0
	ACGT	N										
	A-	(A)										
	AC-	(M)									0	0
											
	N-	X										0

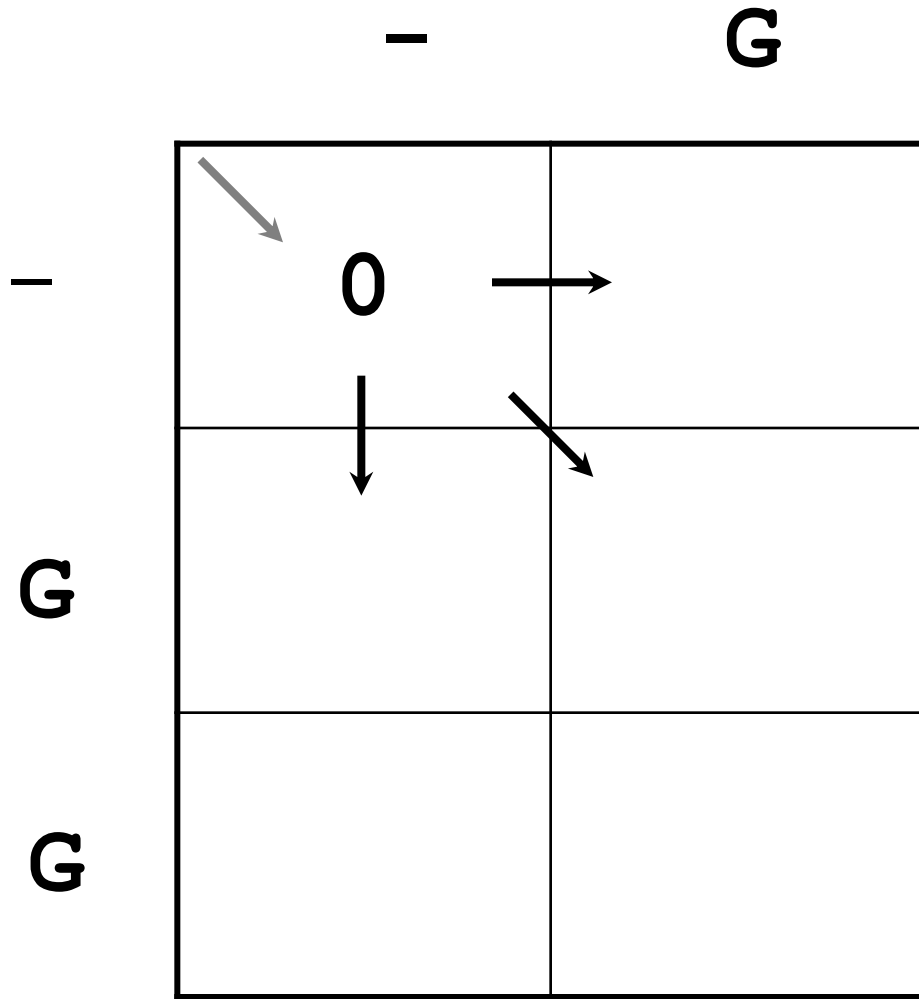
transitions 1
transversions 2
indels 3

Direct optimization (DO)

AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGAT
AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGATGCG



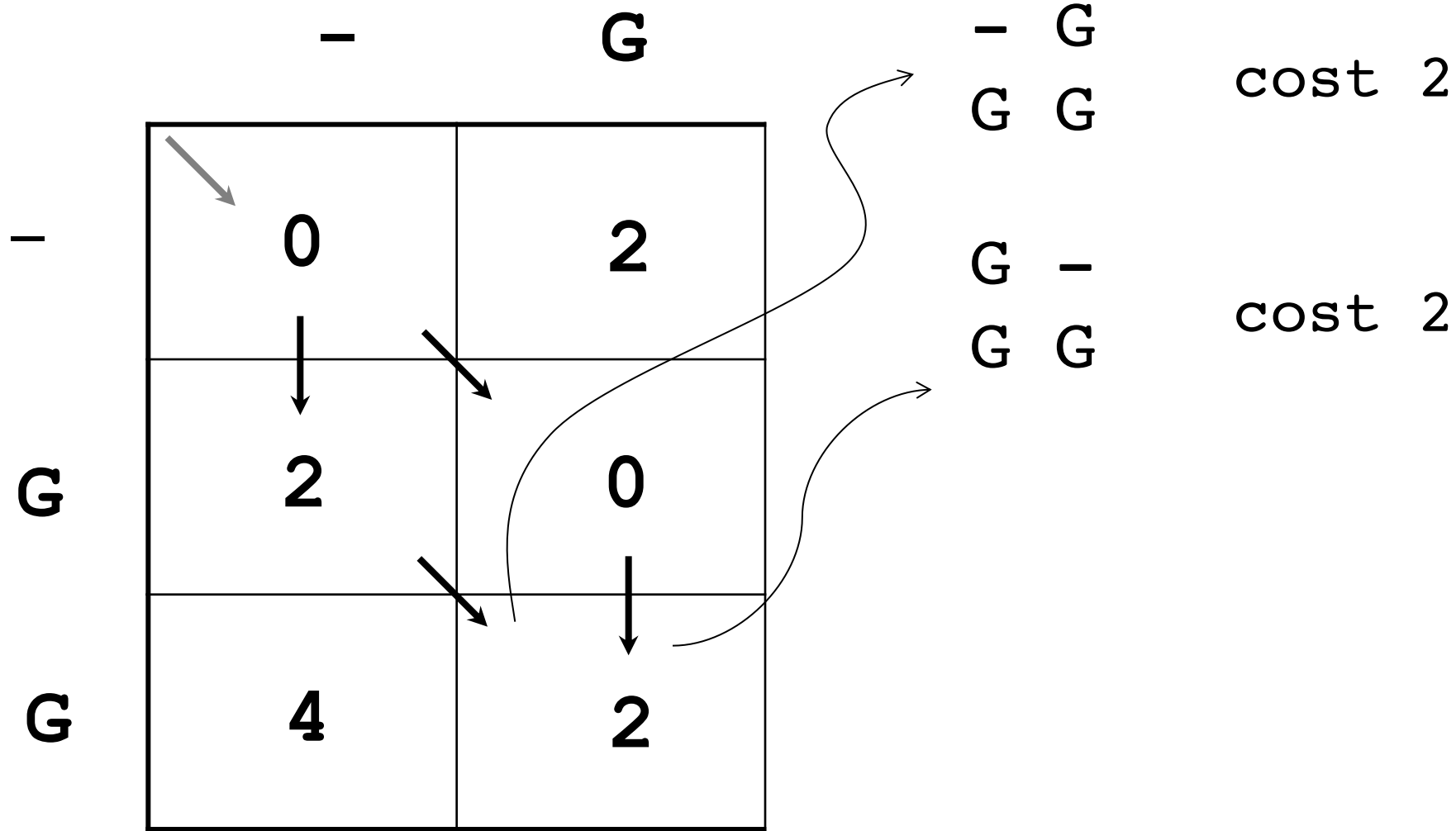
transitions	1
transversions	1
indels	2



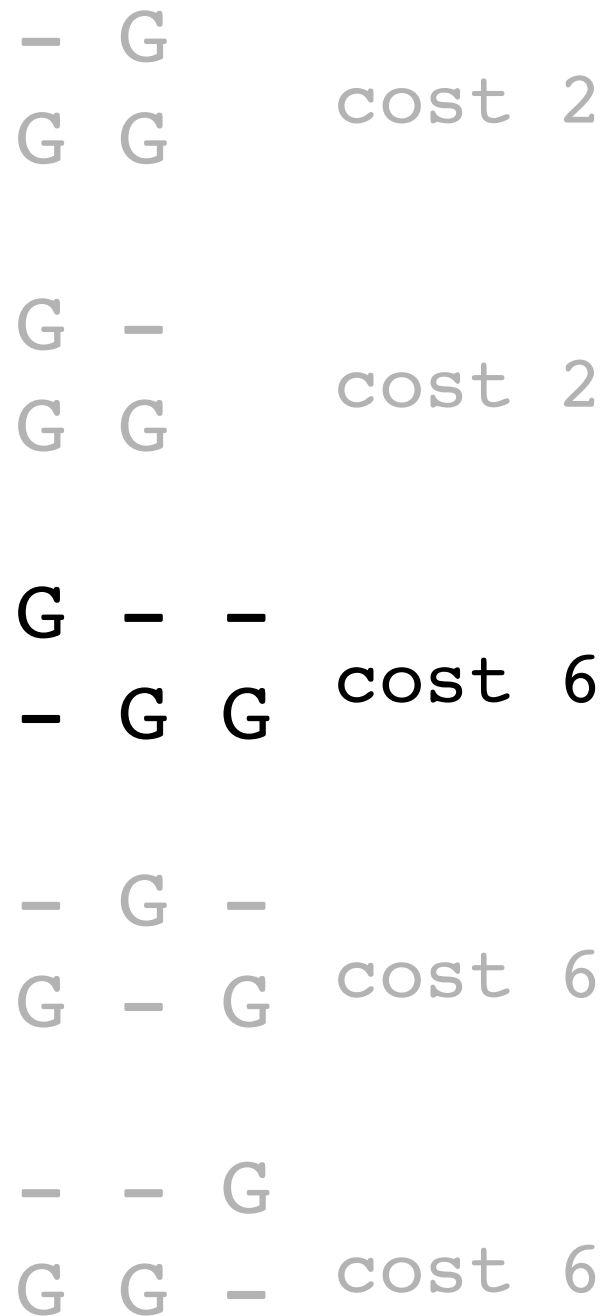
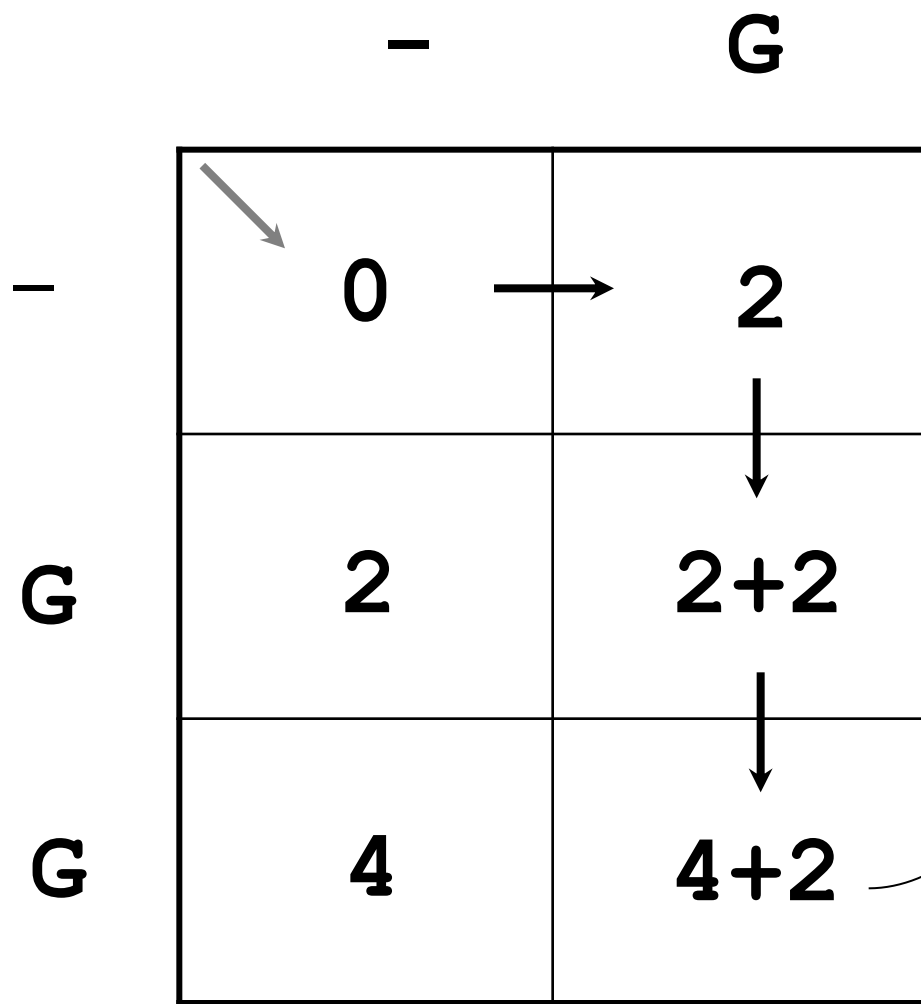
transitions 1

transversions 1

indels 2



transitions	1
transversions	1
indels	2



transitions	1
transversions	1
indels	2

Direct optimization (DO)

AACGGTTAAGGTACGGAGAATTAGGC
 AACGGTTAAGGTACGGAGAATTAGGC

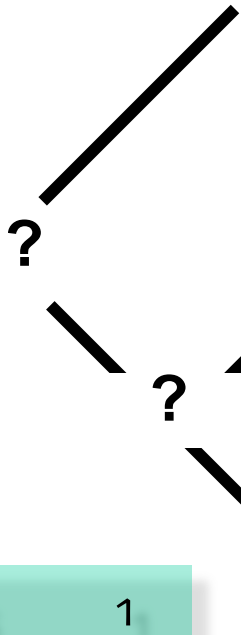
RULE 1: if both terminals share the character state this is also marked for their common ancestor (intersection, \cap)

G **GG**

A

AA

RULE 2: if terminals have different character states (intersection, $\cap = \emptyset$) mark their union (\cup) for their common ancestor



- G

cost 2

G G

G -

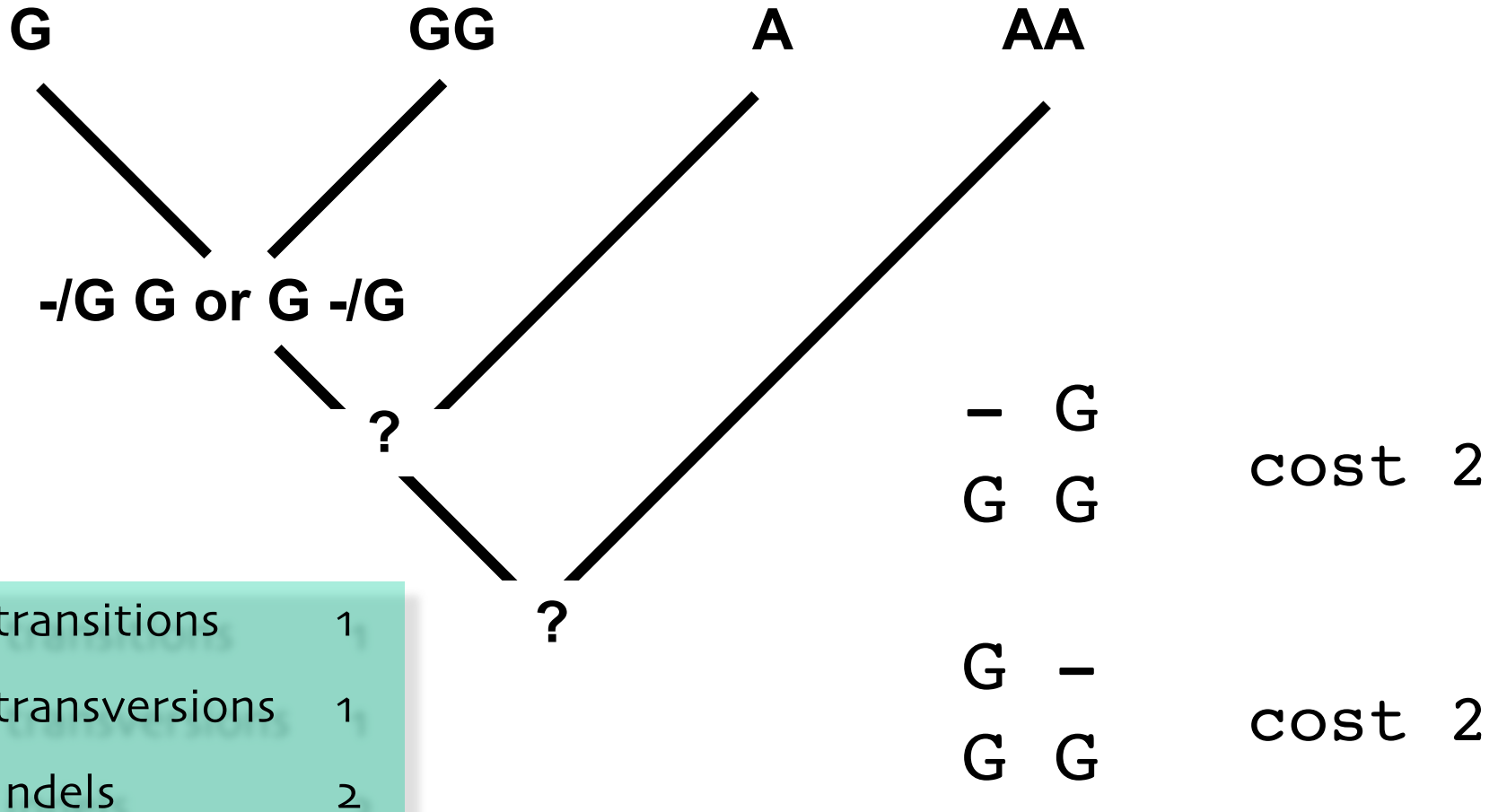
cost 2

G G

transitions	1
transversions	1
indels	2

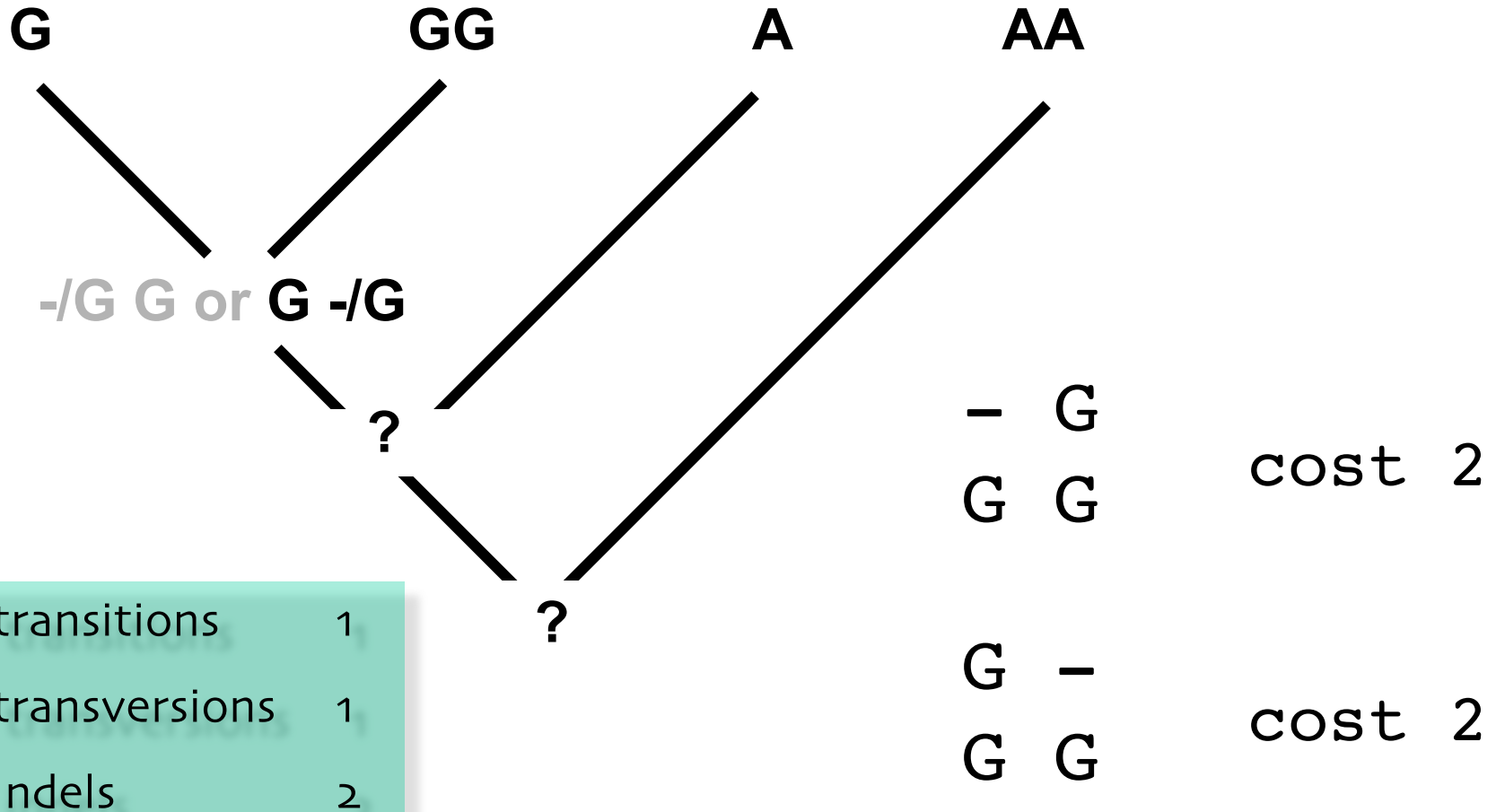
Direct optimization (DO)

AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGAT
 AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGATGCG



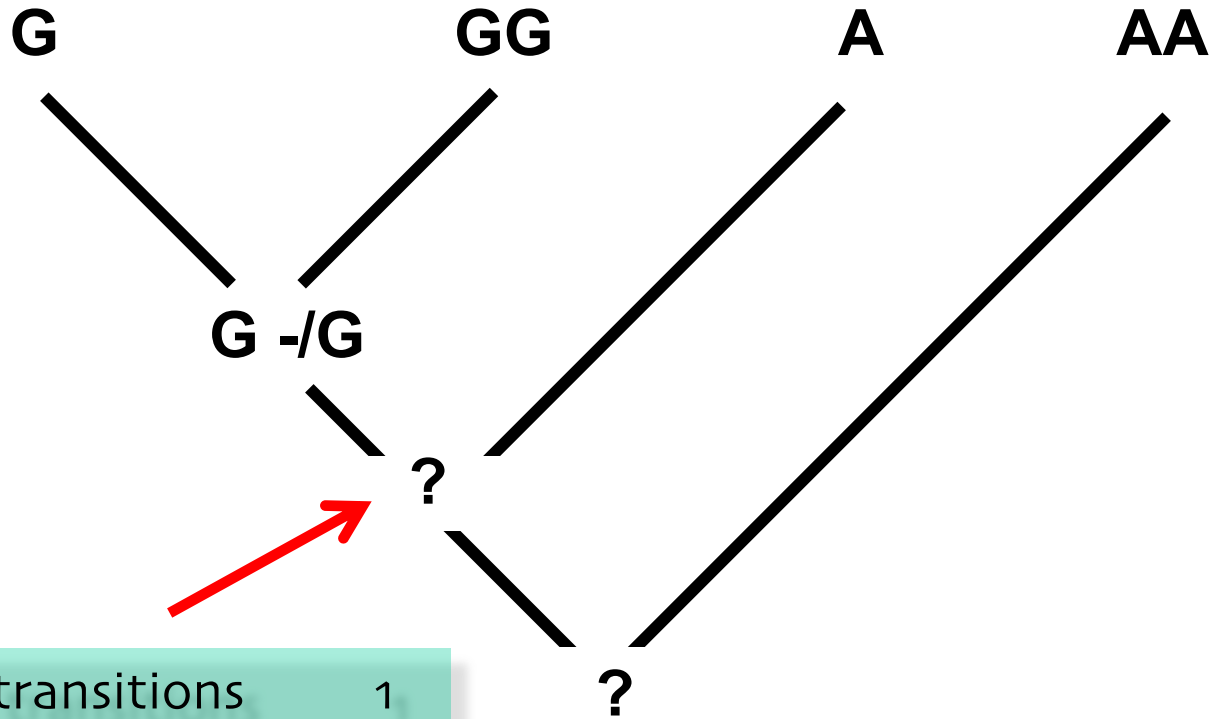
Direct optimization (DO)

AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGAT
 AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGATGCG



Direct optimization (DO)

AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGAT
AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGATGCG



transitions	1
transversions	1
indels	2

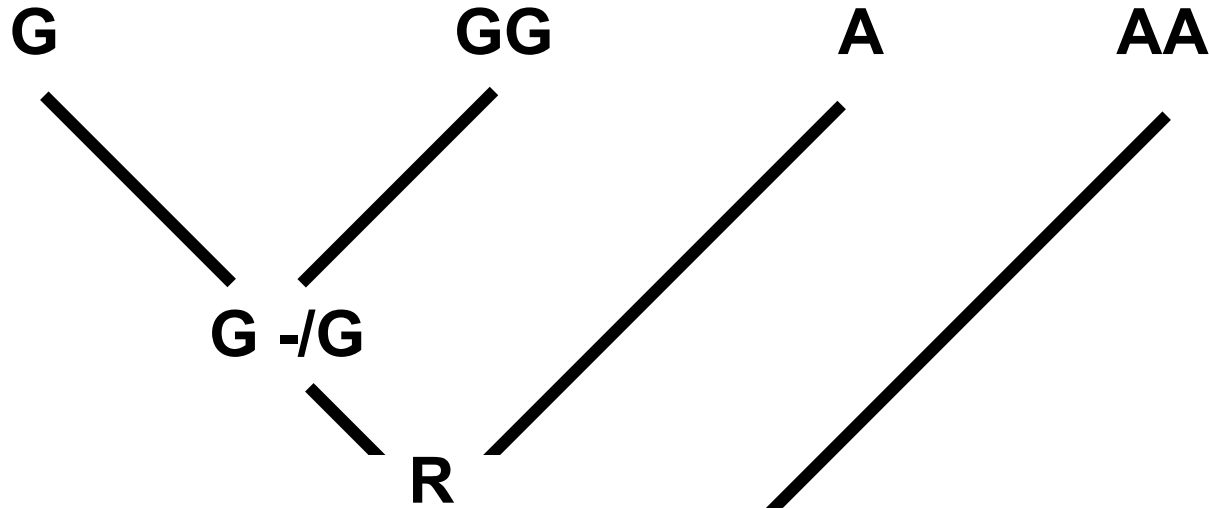
		-	A
-	0		2
G	2		1
-/G	2		1

- A - cost 1
- G - /G cost 1
- A cost 3
- G - /G cost 3
- A - - cost 4
- G - /G cost 4
- A - cost 4
- G - - /G cost 4
- - A cost 4
- G - /G - cost 4

Direct optimization (DO)

AACGGTTAAGGTACGGAGAATTAGGC
AACGGTTAAGGTACGGAGAATTAGGC

RULE 1: if both terminals share the character state this is also marked for their common ancestor (intersection, \cap)



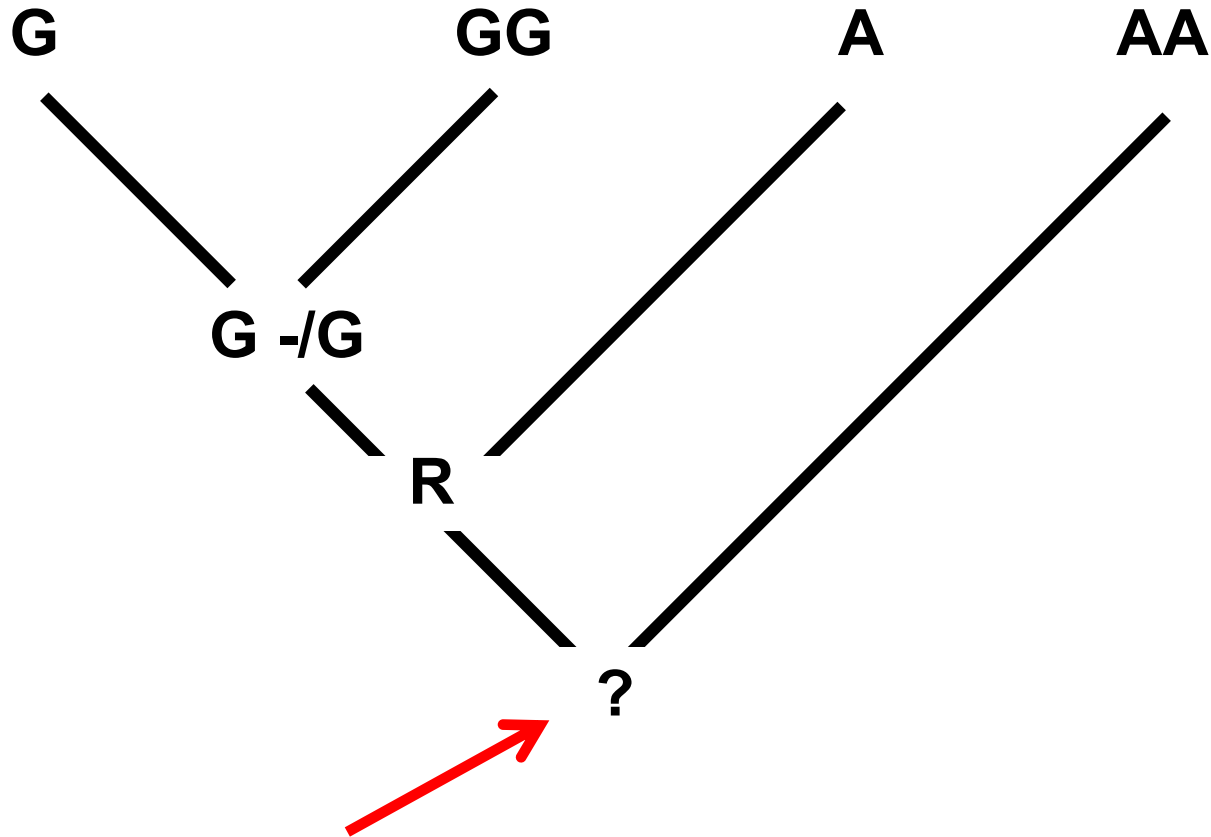
GA	R
GT	K
GC	S
GAT	D
GAC	V
GTC	B
ATC	H
AT	W
AC	M
CT	Y
ACGT	N
ACGT-	X

RULE 2: if terminals have different character states (intersection, $\cap = \emptyset$) mark their union (\cup) for their common ancestor

A -
G -/G cost 1

Direct optimization (DO)

AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGAT
AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGATGCG



	-	R
-	0	2
A	2	0
A	4	2

R - cost 2
A A

- R cost 2
A A

R - - cost 6
- A A

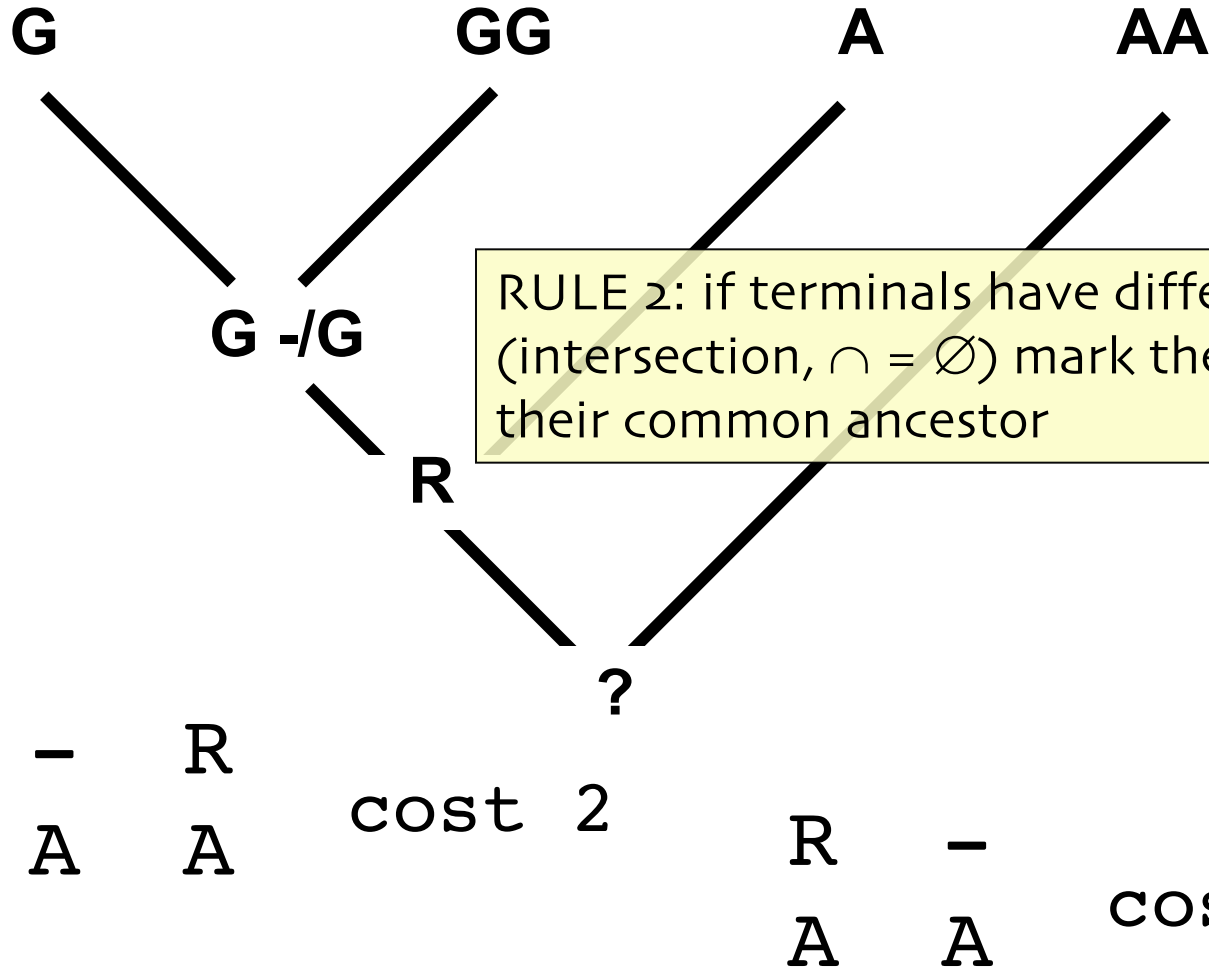
- R - cost 6
A - A

- - R cost 6
A A -

Direct optimization (DO)

AACGGTTAAGGTACGGAGAATTAGGC
AACGGTTAAGGTACGGAGAATTAGGC

RULE 1: if both terminals share the character state this is also marked for their common ancestor (intersection, \cap)

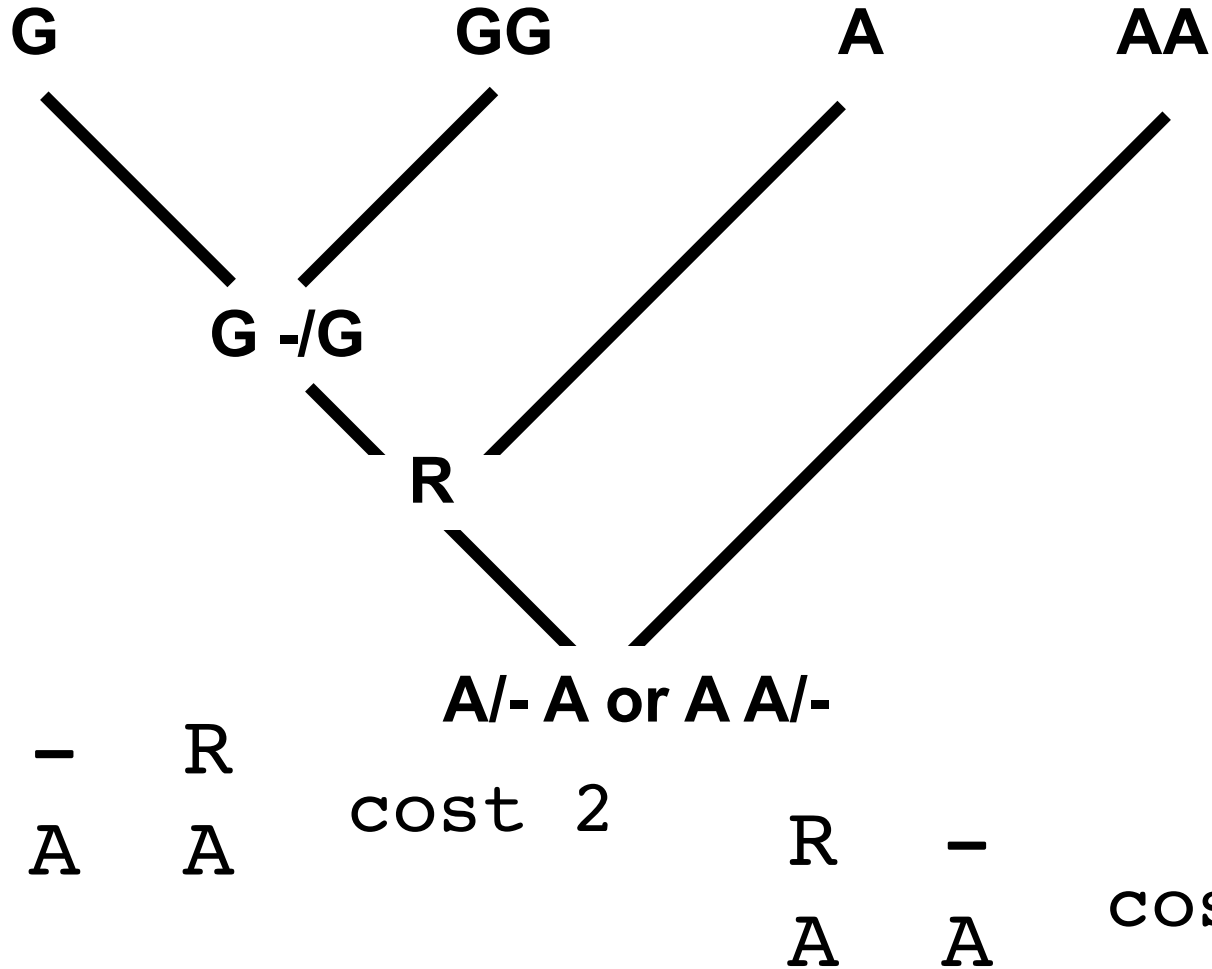


RULE 2: if terminals have different character states (intersection, $\cap = \emptyset$) mark their union (\cup) for their common ancestor

GA	R
GG	S
GAT	D
GAC	V
GTC	B
ATC	H
AT	W
AC	M
CT	Y
ACGT	N
ACGT-	X

Direct optimization (DO)

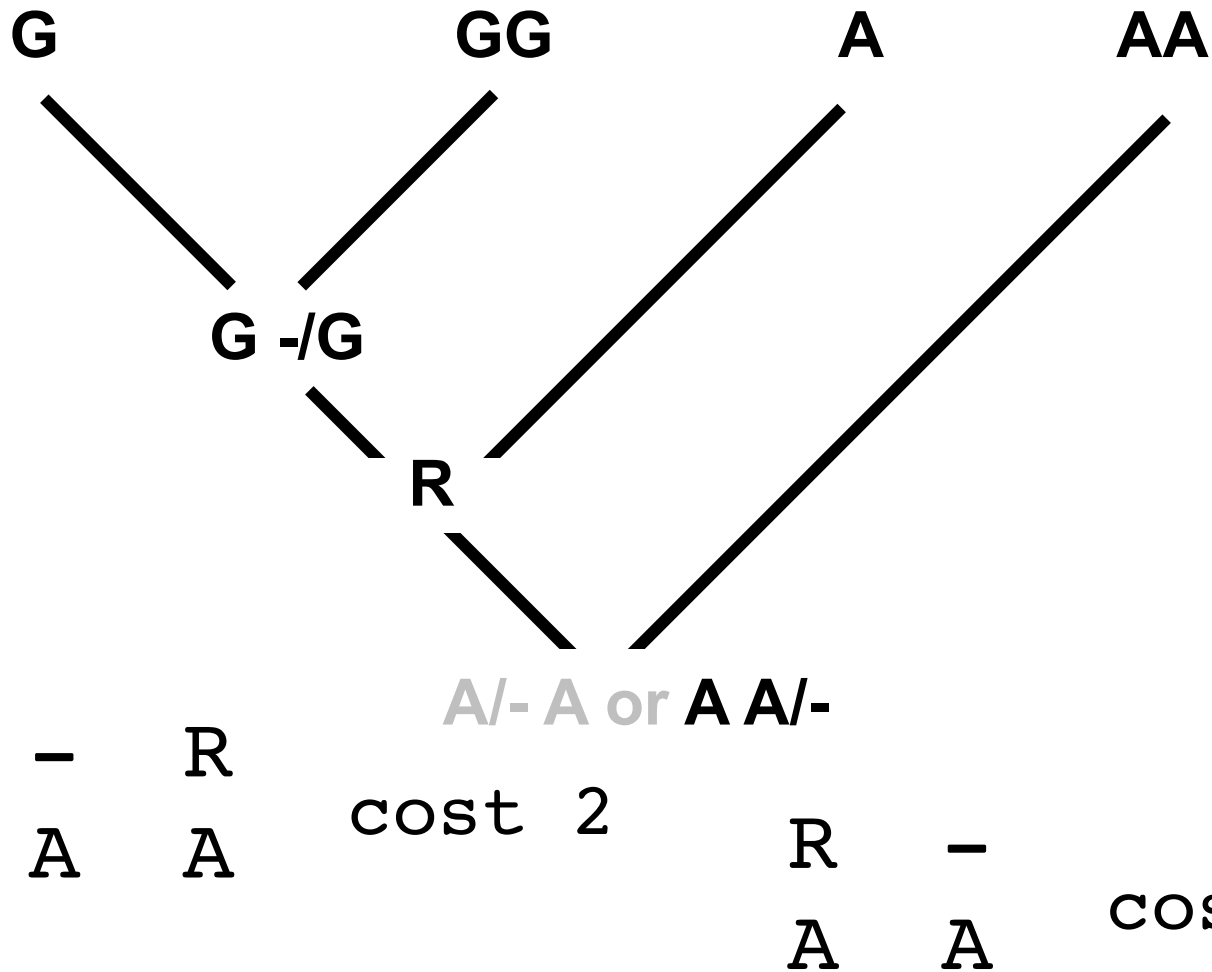
AACGGTTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGAT
 AACGGTTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGATGCG



GA	R
GT	K
GC	S
GAT	D
GAC	V
GTC	B
ATC	H
AT	W
AC	M
CT	Y
ACGT	N
ACGT-	X

Direct optimization (DO)

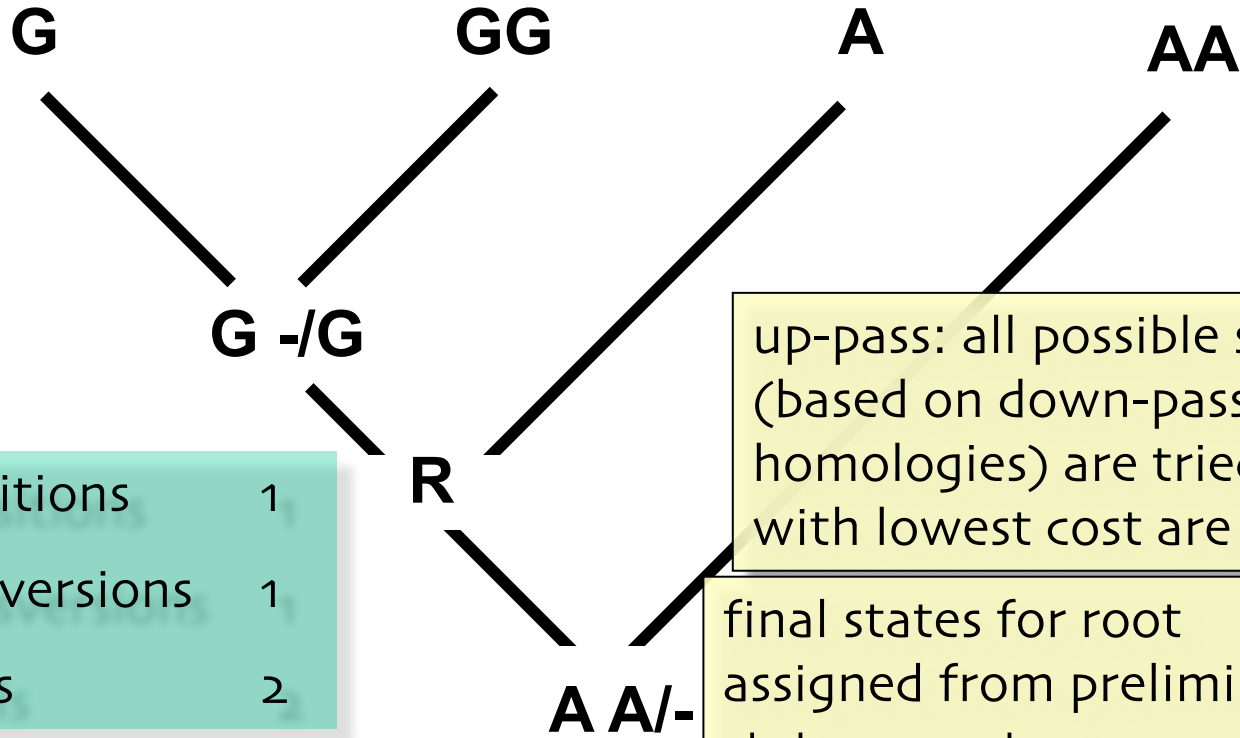
AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGATGCG
 AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGATGCG



GA	R
GT	K
GC	S
GAT	D
GAC	V
GTC	B
ATC	H
AT	W
AC	M
CT	Y
ACGT	N
ACGT-	X

Direct optimization (DO)

AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGAT
 AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGATGCG



up-pass: all possible states (based on down-pass homologies) are tried, those with lowest cost are kept

final states for root assigned from preliminary states or out-group

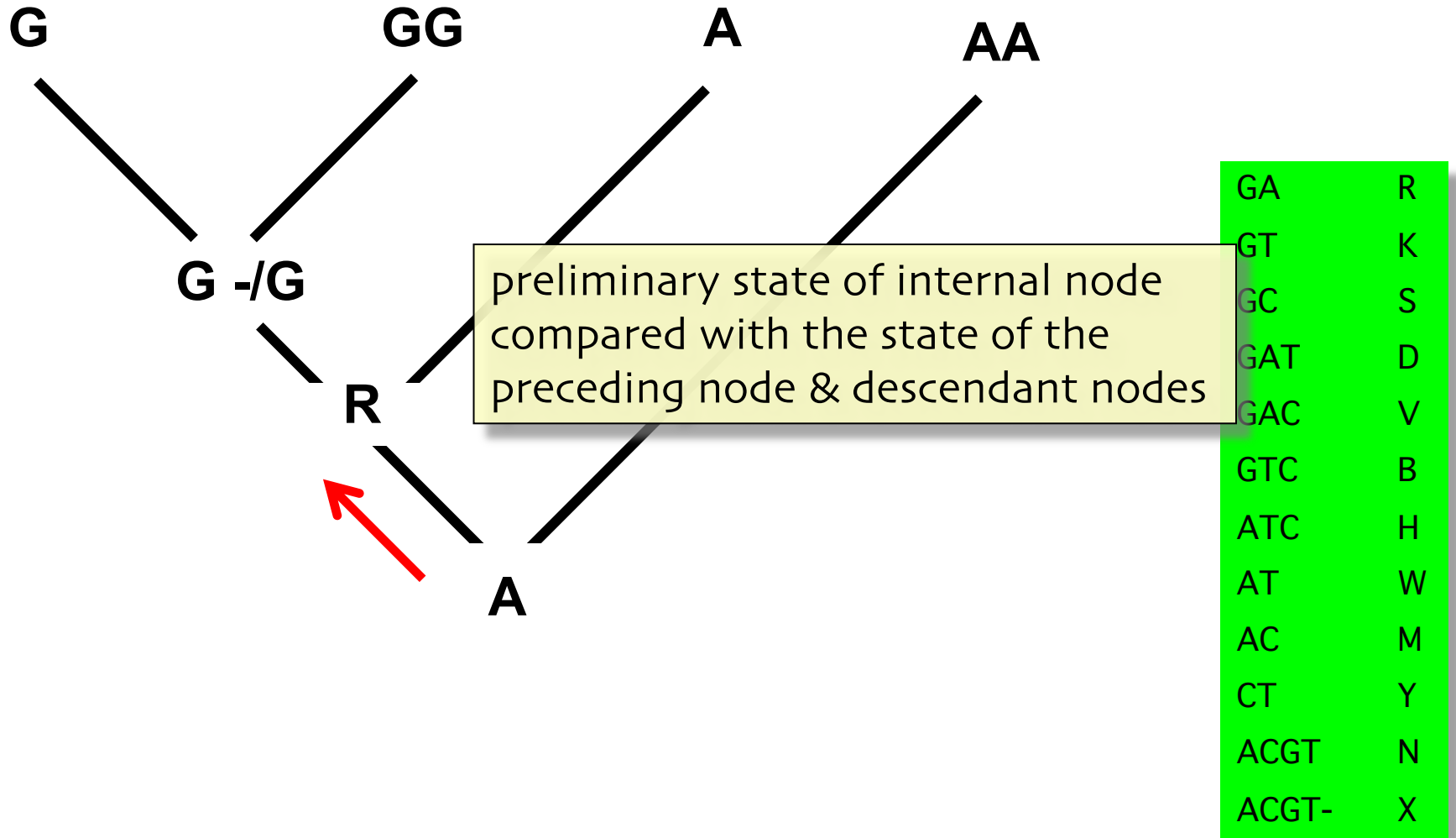
transitions	1
transversions	1
indels	2

GA	R
GT	K
GC	S
GAT	D
GAC	V
GTC	B
ATC	H
AT	W
AC	M
CT	Y
ACGT	N
ACGT-	X

total cost 2+1+2=5

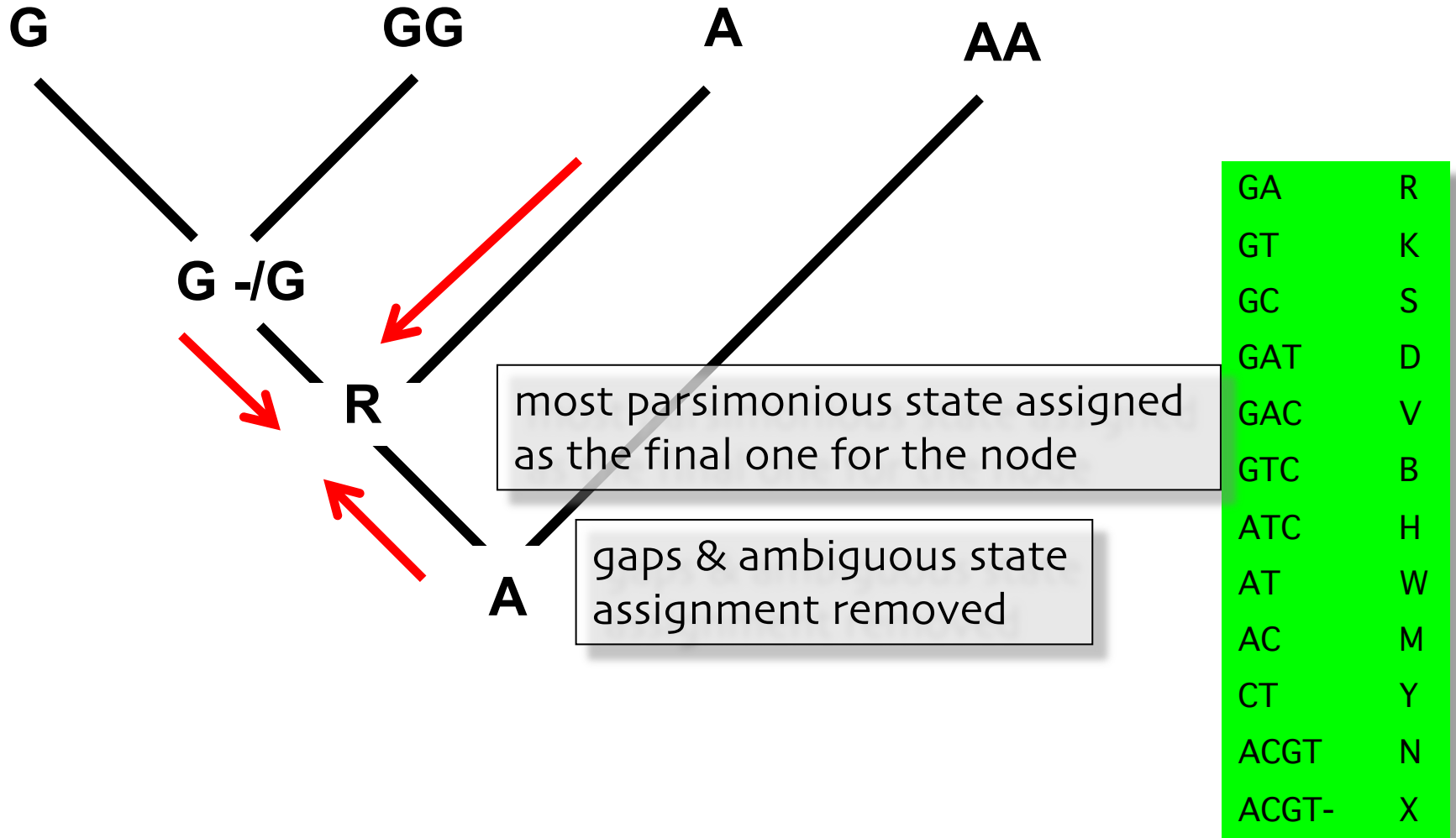
Direct optimization (DO)

AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGAT
AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGATGCG



Direct optimization (DO)

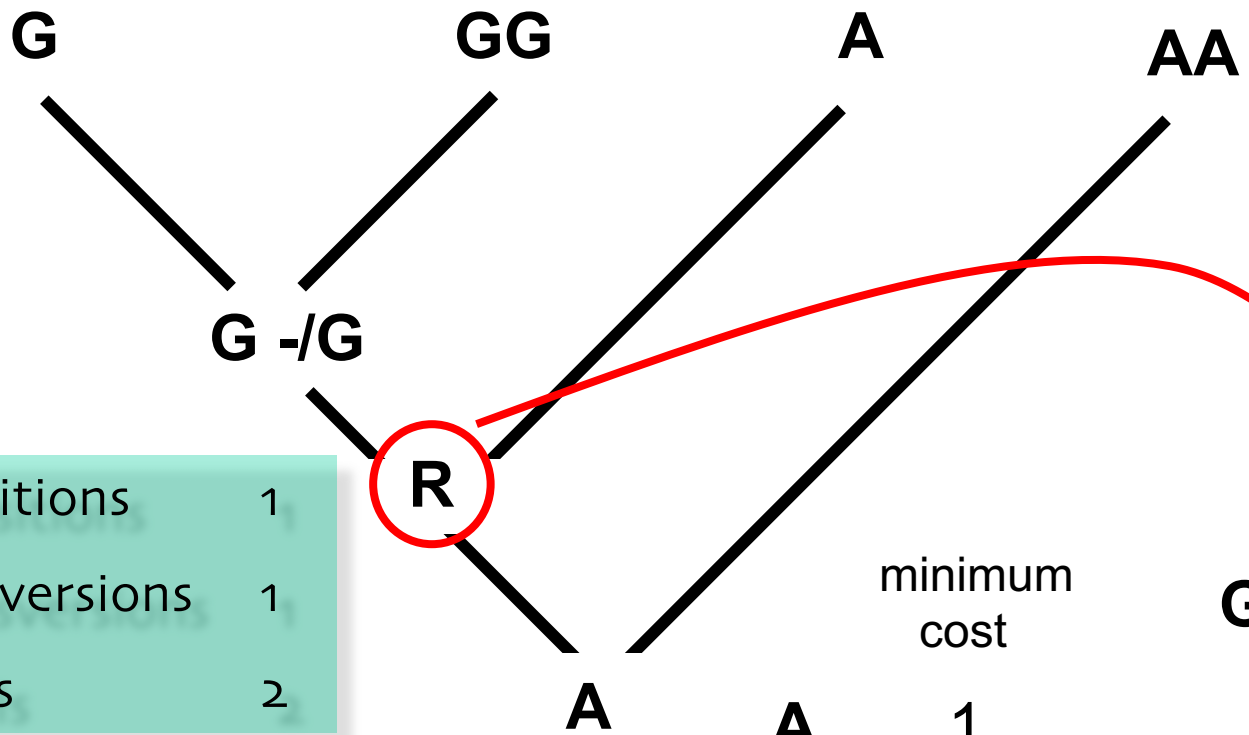
```
AACGGTTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGATAGGAT  
AACGGTTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGATGCG
```



Direct optimization (DO)

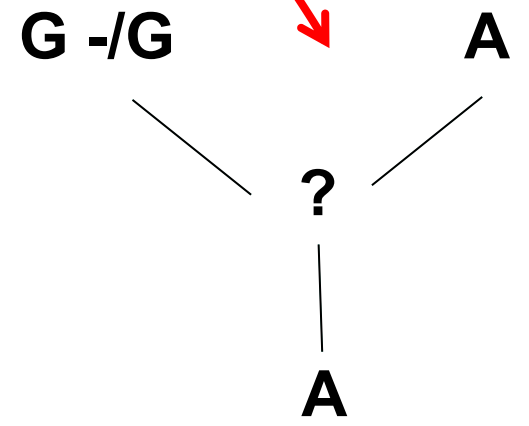
AACGGTTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGG
 AACGGTTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGG

GA	R
GT	K
GC	S
GAT	D
GAC	V
GTC	B
ATC	H
AT	W
AC	M
CT	Y
ACGT	N
ACGT-	X



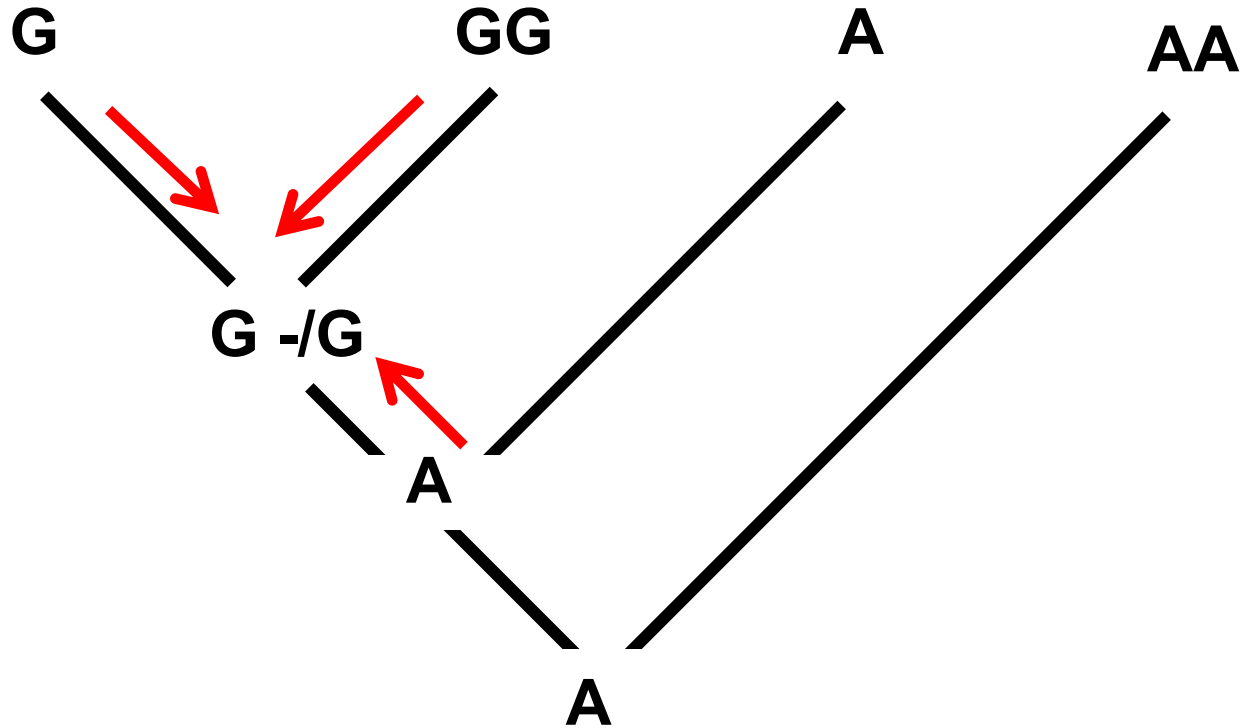
transitions	1
transversions	1
indels	2

	minimum cost
A	1
G	2
GG	6



Direct optimization (DO)

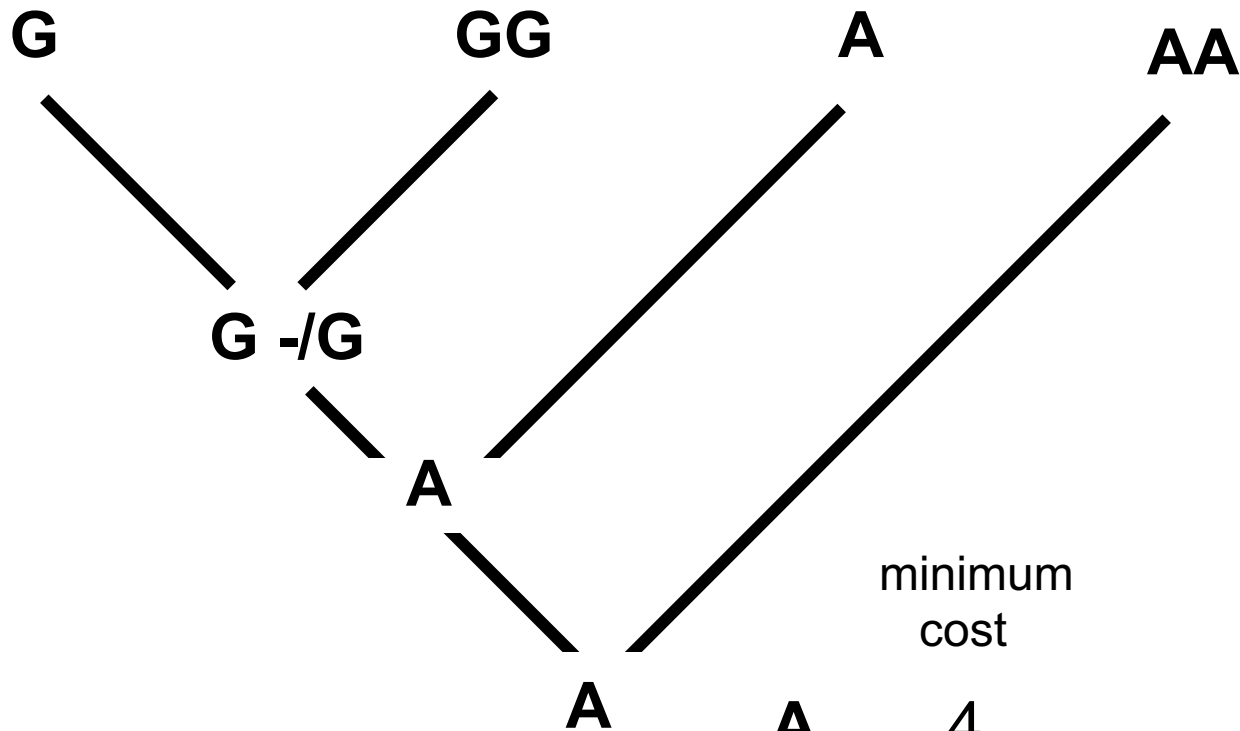
AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGAT
AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGATGCG



GA	R
GT	K
GC	S
GAT	D
GAC	V
GTC	B
ATC	H
AT	W
AC	M
CT	Y
ACGT	N
ACGT-	X

Direct optimization (DO)

AACGGTTTAAGGTACGGAGAATTAGGCCAACCCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCCAACCCGGTAGGAT
AACGGTTTAAGGTACGGAGAATTAGGCCAACCCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCCAACCCCTAGGATGC

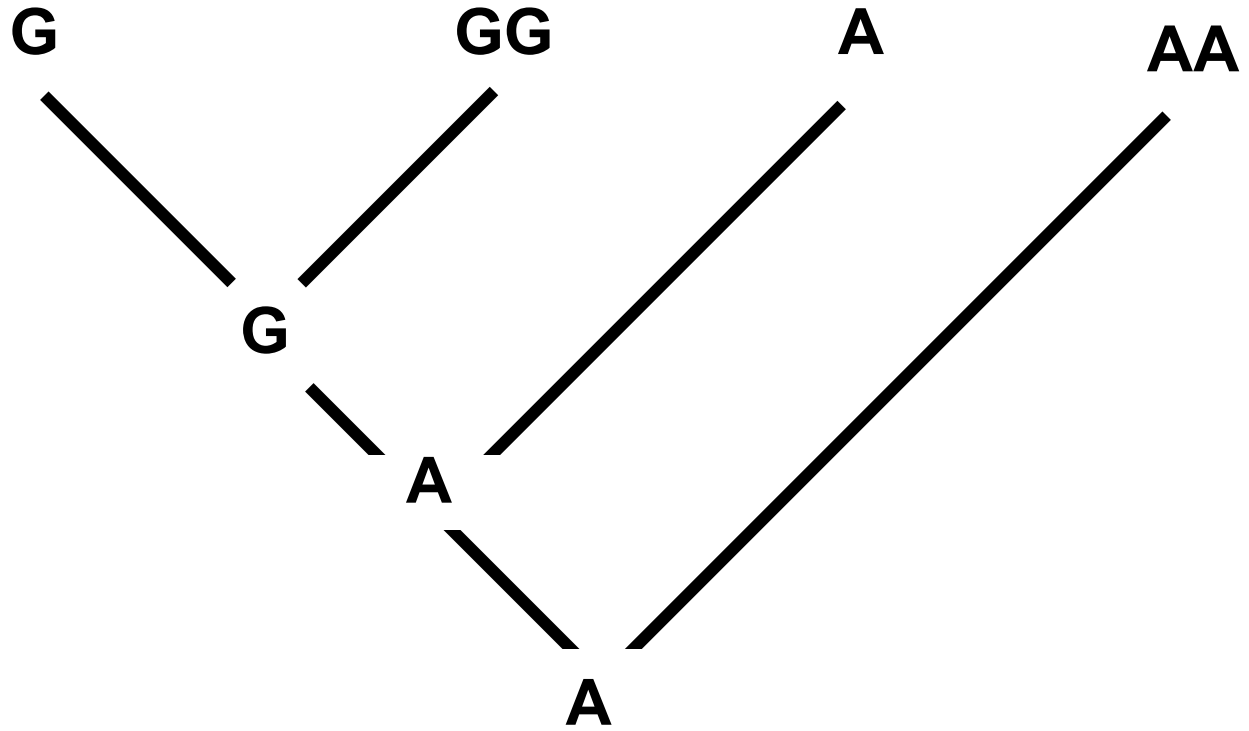


	minimum cost
A	4
G	3
GG	5

GA	R
GT	K
GC	S
GAT	D
GAC	V
GTC	B
ATC	H
AT	W
AC	M
CT	Y
ACGT	N
ACGT-	X

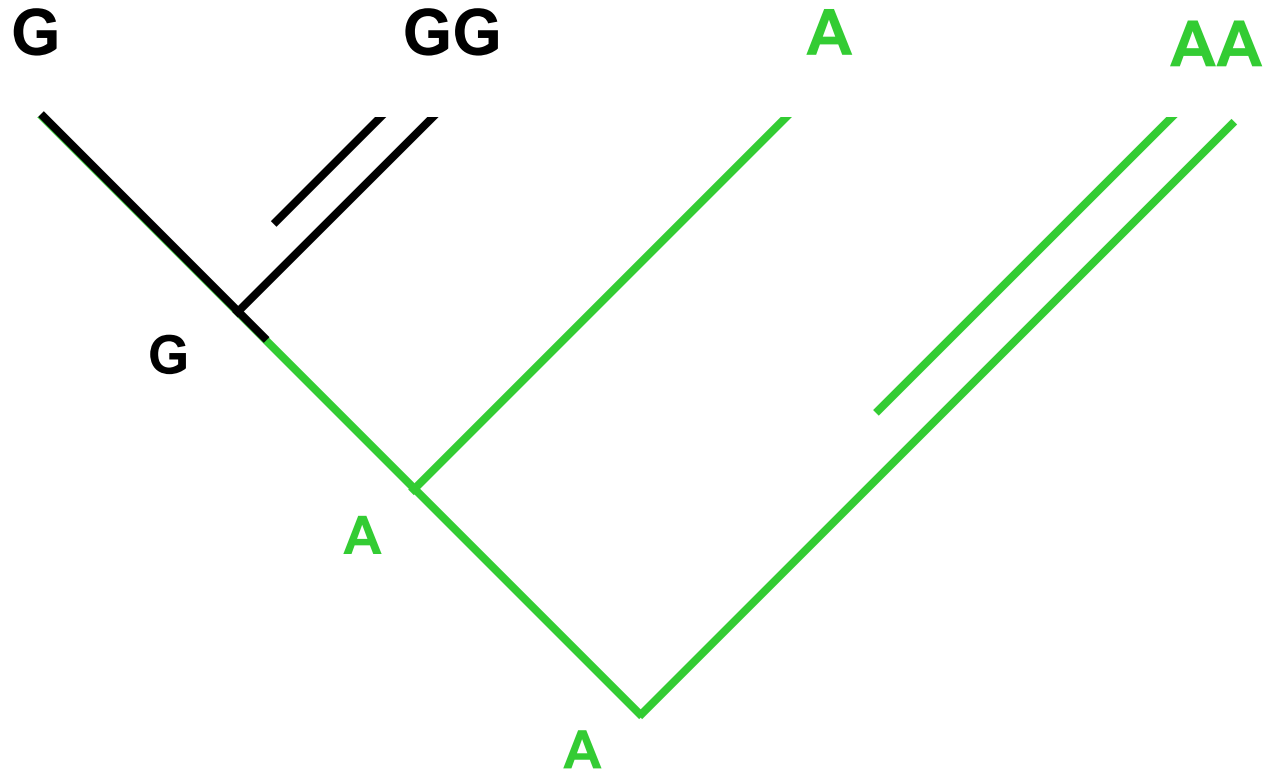
Direct optimization (DO)

AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGAT
AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGATGCG



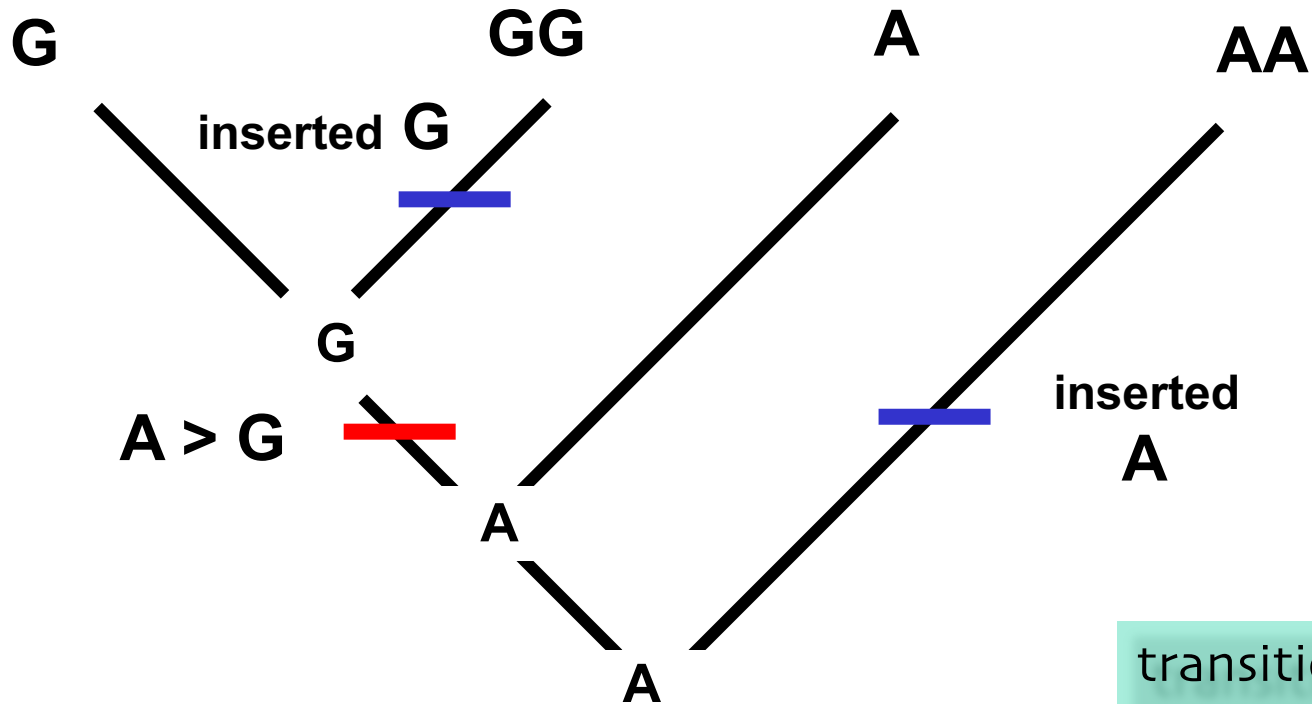
Direct optimization (DO)

AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGAT
AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGATGCG



Direct optimization (DO)

AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGAT
AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGATGCG

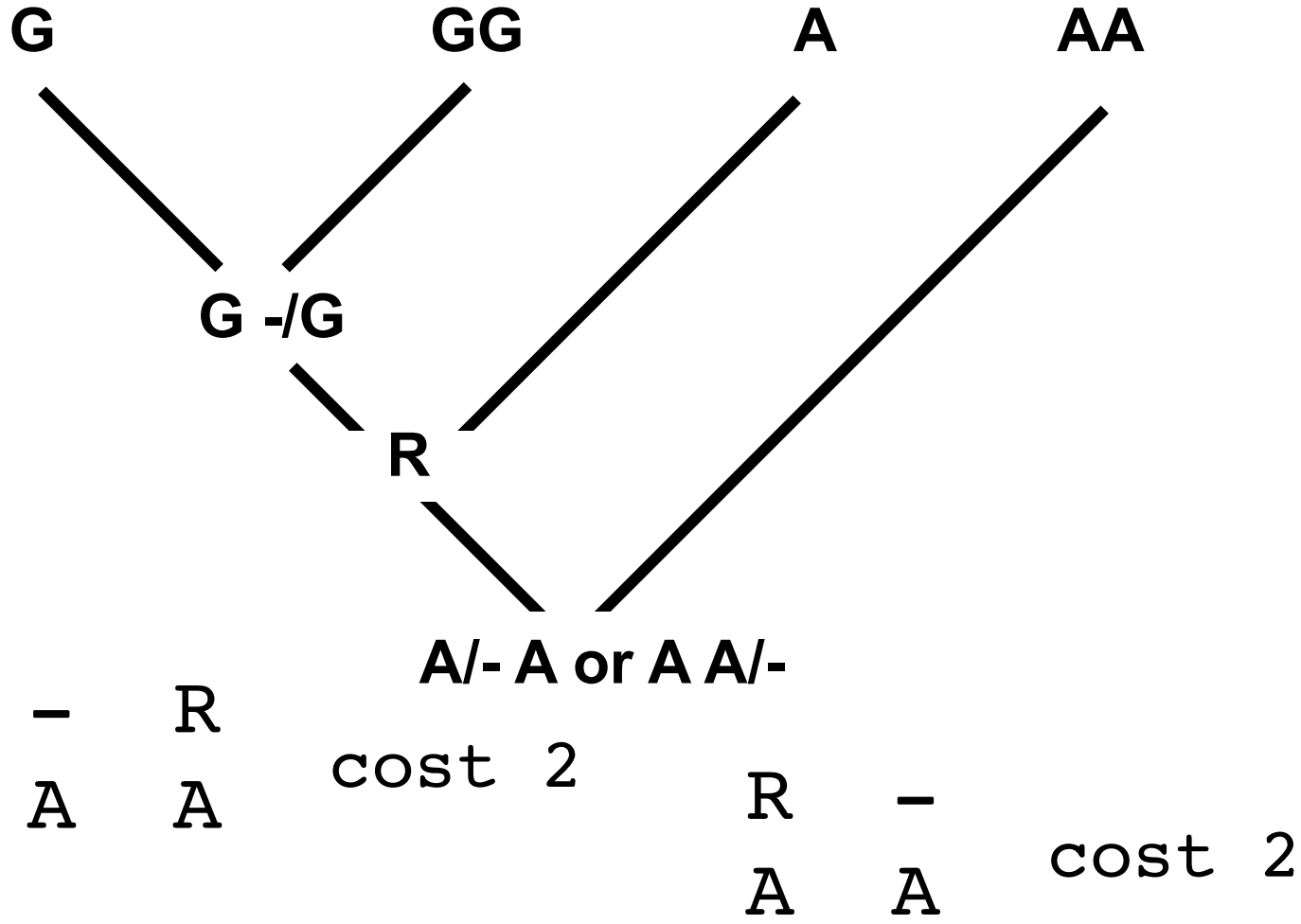


transitions	1
transversions	1
indels	2

total cost 2+1+2=5

Direct optimization (DO)

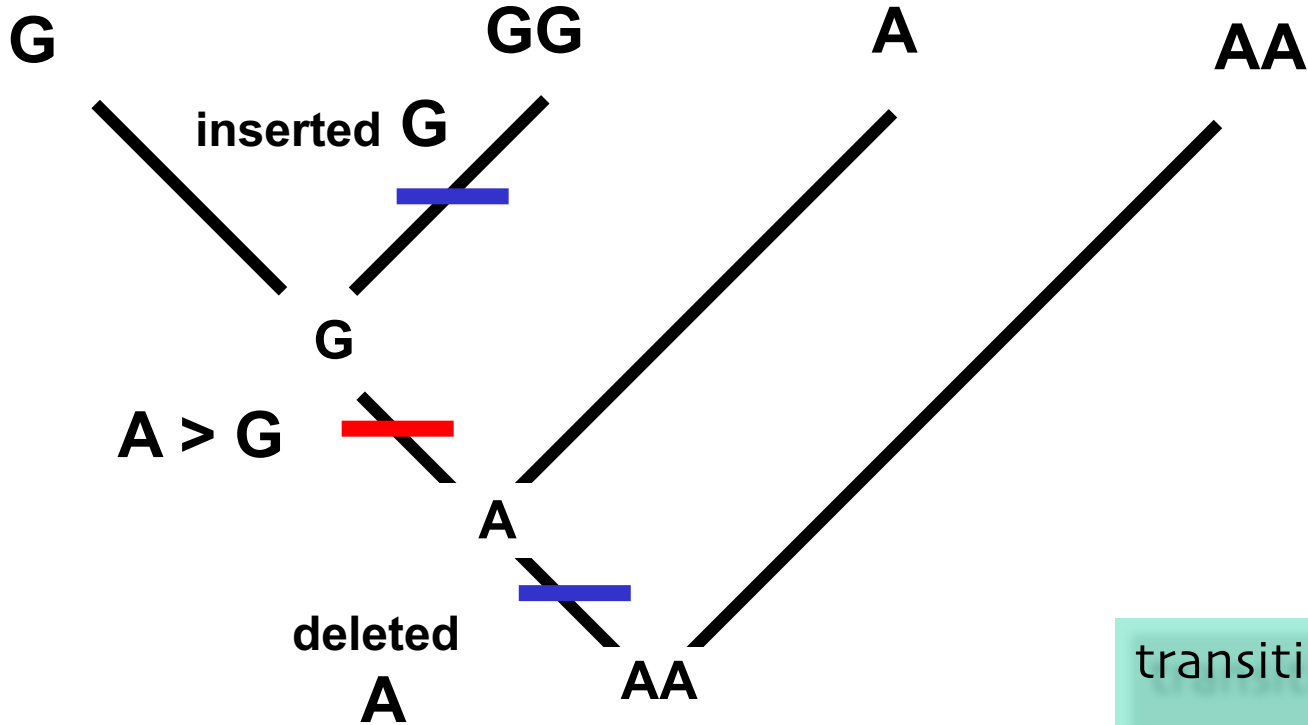
AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGAT
 AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGATGCG



GA	R
GT	K
GC	S
GAT	D
GAC	V
GTC	B
ATC	H
AT	W
AC	M
CT	Y
ACGT	N
ACGT-	X

Direct optimization (DO)

AACGGTTAAGGTACGGAGAATTAGGCAACCCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCGGTAGGAT
AACGGTTAAGGTACGGAGAATTAGGCAACCCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCCTAGGATGCG

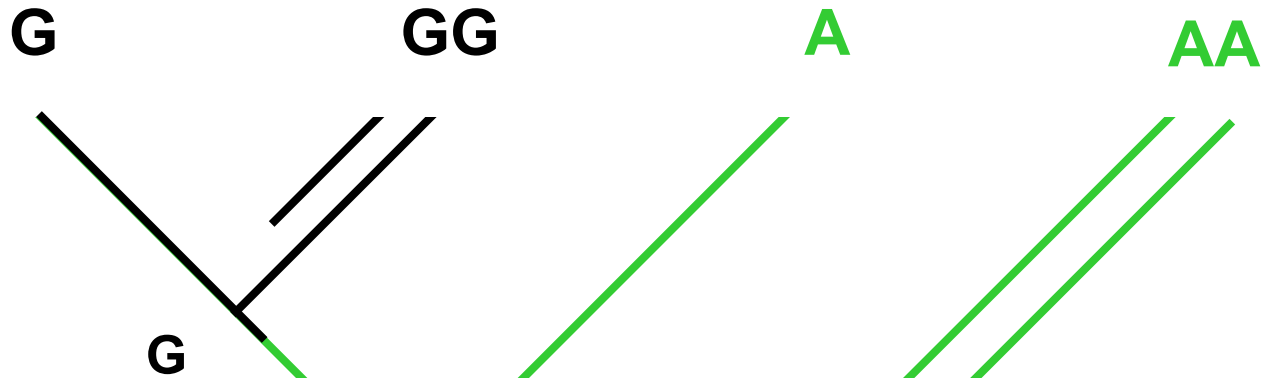


total cost 2+1+2=5

transitions	1
transversions	1
indels	2

Direct optimization (DO)

AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGAT
AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGATGCG



...sequence data do not present themselves in neat packages....

...homologies are **TOPOLOGY SPECIFIC**...

Wheeler, W. 2006. Dynamic homology and the likelihood criterion. *Cladistics* 22: 157-170.

AA

ADDITIONAL INFORMATION:

AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGAT
AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGATGCG

Phillips, A.J. 2006. Homology assessment and molecular sequence alignment. *J. Biomedical Informatics* 39: 18-33.

Wheeler, W.C. 2002. "Optimization alignment: down, up, error, and improvements." Pp. 55-69. in R. Desalle et al. eds. *Techniques in Molecular Systematics and Evolution*. Birkhäuser Verlag, Basel Switzerland.

Wheeler, W. 1996. Optimization alignment: the end of multiple sequence alignment in phylogenetics? *Cladistics* 12: 1-9.

Algorithm 10.6: DirectOptimizationFirstPass

Data: Input strings X and Y of lengths $|X|$ and $|Y|$

Data: Element cost matrix σ' of elements nearest to all pairs of elements in Σ and λ (indel) based on element cost matrix σ of Algorithm 8.1

Result: Median cost.

Initialize first row and column of matrices;

$direction[0][0] \leftarrow \swarrow;$

$cost[0][0] \leftarrow 0;$

$length[0][0] \leftarrow 0;$

for $i = 1$ **to** $|X|$ **do**

$cost[i][0] \leftarrow cost[i-1][0] + \sigma'_{X_i,\lambda};$

$direction[i][0] \leftarrow \rightarrow;$

$length[i][0] \leftarrow length[i-1][0] + 1;$

end

for $j = 1$ **to** $|Y|$ **do**

$cost[0][j] \leftarrow cost[0][j-1] + \sigma'_{Y_j,\lambda};$

$direction[0][j] \leftarrow \downarrow;$

$length[0][j] \leftarrow length[0][j-1] + 1;$

end

Update remainder of matrices $cost$, $direction$, and $length$;

for $i = 1$ **to** $|X|$ **do**

for $j = 1$ **to** $|Y|$ **do**

$ins \leftarrow cost[i-1][j] + \sigma'_{X_i,\lambda};$

$del \leftarrow cost[i][j-1] + \sigma'_{Y_j,\lambda};$

$sub \leftarrow cost[i-1][j-1] + \sigma'_{X_i,Y_j};$

$cost[i][j] \leftarrow \min(ins, del, sub);$

if $cost[i][j] = ins$ **then**

$direction[i][j] \leftarrow \rightarrow;$

$length[i][j] \leftarrow length[i-1][j] + 1;$

else if $cost[i][j] = del$ **then**

$direction[i][j] \leftarrow \downarrow;$

$length[i][j] \leftarrow length[i][j-1] + 1;$

else

$direction[i][j] \leftarrow \swarrow;$

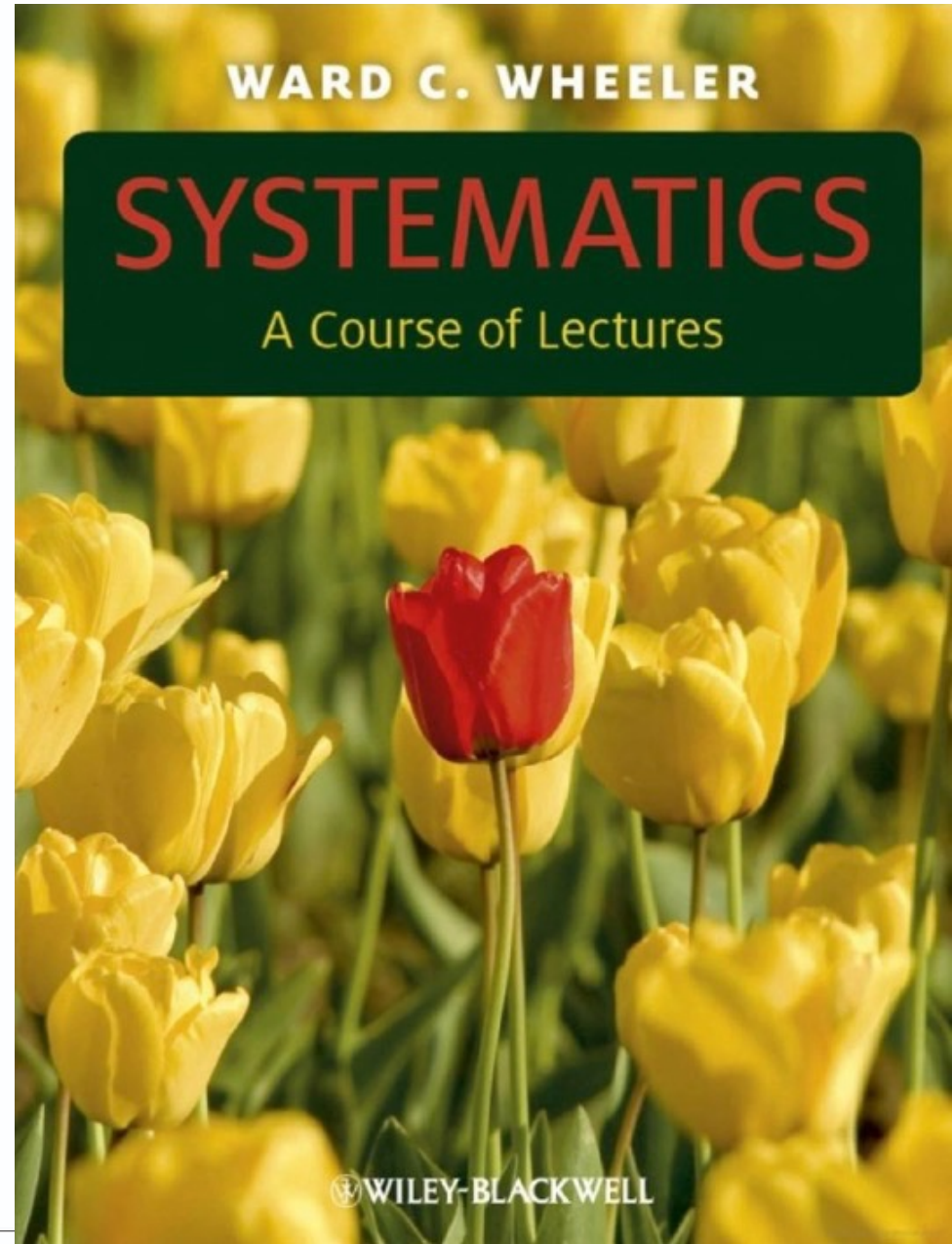
$length[i][j] \leftarrow length[i-1][j-1] + 1;$

end

end

end

return $cost[|X|][|Y|]$



Homology - static or dynamic?

AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGAT
AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGATGCG

conventional alignment

direct optimization

terminal	1	TATACTTT
	2	----CT-C
	3	TA--C----

POSITIONS of nucleotides
FIXED before phylogenetic
analysis

Homology - static or dynamic?

AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGAT
 AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGATGCG

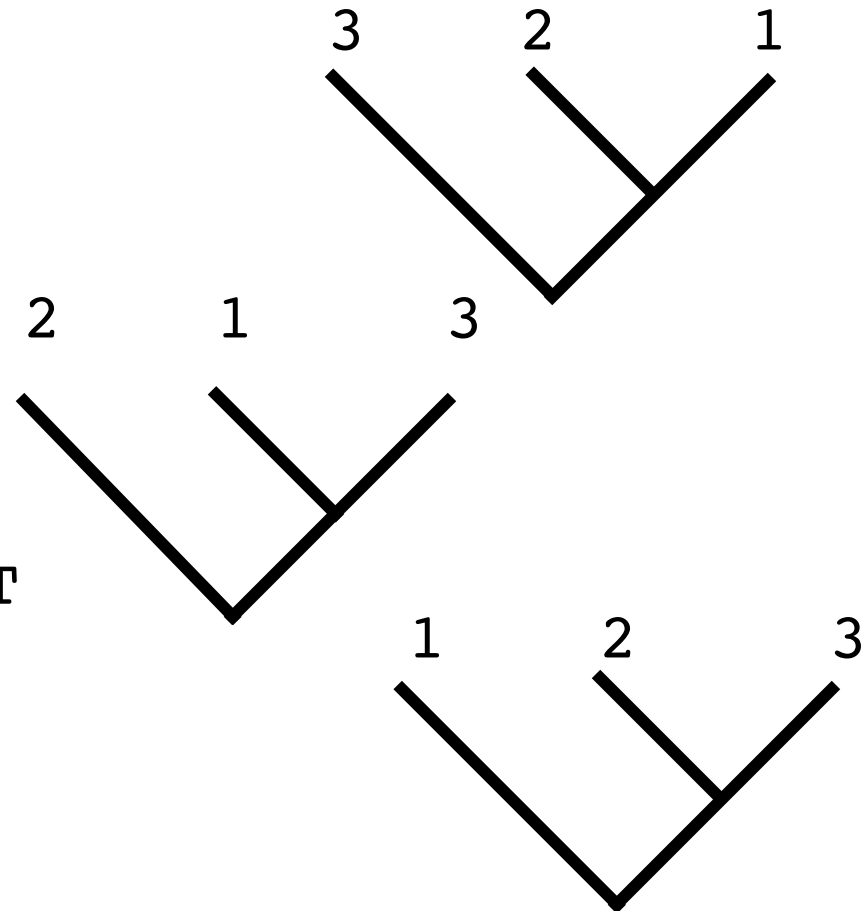
conventional alignment

terminal 1 TATACTTT
 2 ----CT-C
 3 TAC

positions of nucleotides ARE
*FREE to change during
 phylogenetic analysis*

terminal 1 TATACTTT
 2 CTC
 3 TAC

direct optimization



Number of possible alignments

AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGAT
 AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGATGCG

$f(n, m)$ for $1 \leq n \leq 10; 1 \leq m \leq 5$

$n \backslash m$	2	3	4	<u>5</u>
1	3	13	75	541
2	13	409	23917	2244361
3	63	16081	10681263	14638756721
4	321	699121	5552351121	117629959485121
5	1683	32193253	3147728203035	1.05×10^{18}
<u>10</u>	8097453	9850349744182729	3.32×10^{26}	<u>1.35×10^{38}</u>

Slowinski (1998)

Finding THE optimal alignment for even a small number of short sequences is **IMPOSSIBLE**

T B(T)

3	1
4	3
5	15
6	105
7	945
8	10 395
9	135 135
10	2×10^6
15	8×10^{12}
20	2×10^{20}
50	3×10^{74}

$$B(n) = (2n-5)!!$$

in DO these
2 NP complete
problems are coupled

computationally **VERY**
DEMANDING
only heuristic "solutions"

Direct optimization (DO)

AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGAT
AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGATG

do we search homology
between nucleotides

AAATCGCGGATT

base to base homology

AACTCCCGGAGT

AAATCGC-GAGT

or between STRETCHES of nucleotides

fragment homology

AAATCGCGGATT

AACTCCCGGAGT

AAATCGCGGAGT

Direct optimization

AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGAT
AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGATG

base to base homology

POSITION of single nucleotide = character

different nucleotides (A C G T) plus alignment gaps =
character states

AAATCGCGGATT

AACTCCGGAGT

AACTCGC-GAGT

12 characters

A C G T & - 5 character states

Fixed states optimization (FSO)

Wheeler, W.C. 1999. Cladistics 15: 379-385

fragment homology

stretch of DNA = character

homologous stretches in different terminals
= character states

AAATCGCGGATTT

AACTCGCGGAGTT

AACTCGCGGAGTT

1 character

number of character states $x \leq n$

n = number of terminals studied

Direct optimization (DO)

AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGAT
AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGATGCG

we try to figure out homologies at the level alphabets

P A N T S	P - A N T S
A N T S	- - A N T S
P L A N T	P L A N T S
R A T S	R - A - T S
B A T S	B - A - T S

Fixed states optimization (FSO)

AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGAT
AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGATGCG

we do NOT try to figure out homologies between alphabets

instead we are comparing whole WORDS

each stretch to be compared potentially different character state number of terminals to be compared = n

number of possible characters states = x

$$x \leq n$$

PANTS

DUNGAREES

TROUSERS

SLACKS

JEANS

Fixed states optimization (FSO)

```
AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGAT  
AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGATG
```

- 1) **pair-wise** comparison of sequences
- 2) cost of transformation determined between each pair
(= cost between different character STATES)
this is based on the nucleotide composition of each stretch
- 3) each stretch (= character state) located on internal nodes of trees
- 4) the tree with lowest overall cost is chosen

Fixed states optimization (FSO)

AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGAT
AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGATGCG

Advantage: large number of character states
origin of exactly same kind of sequence in 2 unrelated
terminals is VERY unlikely in DO only 5 states

A C G T -

Disadvantage: subjectivity in determination of
characters, how long, what kind of stretches?

in DO this is NOT a problem, nucleotides
CANNOT be reduced any further

METRICITY is required between character states

Fixed state optimization (FSO)

AACGGTTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGAT
AACGGTTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGATGCG

AATTT

ATTAA

TT

TAATT

terminals with numerous missing entries might have SAME EDIT cost with MANY other terminals

this might lead to violation of triangle inequality

$$d(x,y) \leq d(x,z) + d(z,y)$$

determination of adjacency of character states mathematically IMPOSSIBLE

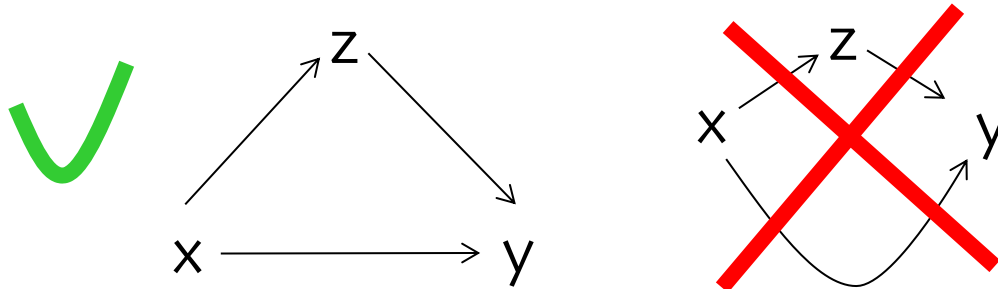
...optimization will lead to inconsistent results & thus characters states have to fulfill the following requirements:

$$\begin{aligned}\forall x & d(x,x) = 0 \\ \forall x, y; x \neq y & d(x, y) > 0 \\ \forall x, y & d(x,y) = d(y,x) \\ \forall x, y, z & d(x,y) \leq d(x, z) + d(z, y)\end{aligned}$$

x, y, z character states
d change from ch. state to another

i.e.

- 1) change from one ch. state to another > 0
- 2) change $x \rightarrow y = y \rightarrow x$ (changes symmetrical)
- 3) change $x \rightarrow y \leq$ change $x \rightarrow z +$ change $z \rightarrow y$
(triangle inequality)



SUMMARY

AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGAT
AACGGTTAAGGTACGGAGAATTAGGCAACCCTAGGATGCGATGCGCAGAGTTAGGTACGGAGAATTAGGCAACCCTAGGATGCG

number of possible alignments is **VAST**

problems of multiple alignment & tree construction
are **inextricably linked**

dynamic homology enables more parsimonious
interpretation of sequence level data

DO computationally demanding

new algorithms & techniques are continuously
developed

HOME "QUIZ"

Why, in this age of whole genome sequences, it is still relevant to study & use also morphological characters in phylogenetic analyses

Why does cutting of continuous sequence stretches within CONSERVED lead into huge speed-up in DO analyses?

come up with arguments for the 1st demo Mon 14.xi.