

# Multiple Sequence Alignment - main points

---

- The goal of MSA is to introduce gaps such that residues in the same column are homologous (all residues in the column descended from a residue in their common ancestor).
- The problem is recast as:
  - reward matches (+ scores)
  - penalize rare substitutions (- scores),
  - penalize gaps (- scores),
  - try to find an alignment that maximizes the total score
- pairwise alignment is tractable
- MSA is usually done progressively
- progressive alignment algorithms are heuristic, and do not optimize an evolutionary defensible criterion

# Multiple Sequence Alignment tools

---

- clustal variants are popular, but not very reliable.
- simultaneous inference of MSA and tree is the most defensible (but computationally demanding)
- Promising tools for MSA (roughly in order of computational tractability):
  1. Simultaneous MSA + Trees (Handel, BAliPhy, BEAST, AliFritz...)
  2. FSA (fast statistical alignment); Infernal (for rRNA); Prank
  3. MAFFT, Muscle, ProbCons
- Iterative “meta-solutions” (e.g. SATè ) allow MSA uncertainty to be incorporated in tree inference.
- GBlocks (and similar tools) cull ambiguously aligned regions.

human	KRSV
chimp	KRV
orang	KPRV

human

chimp

orangutan

KRSV

KRV

KPRV

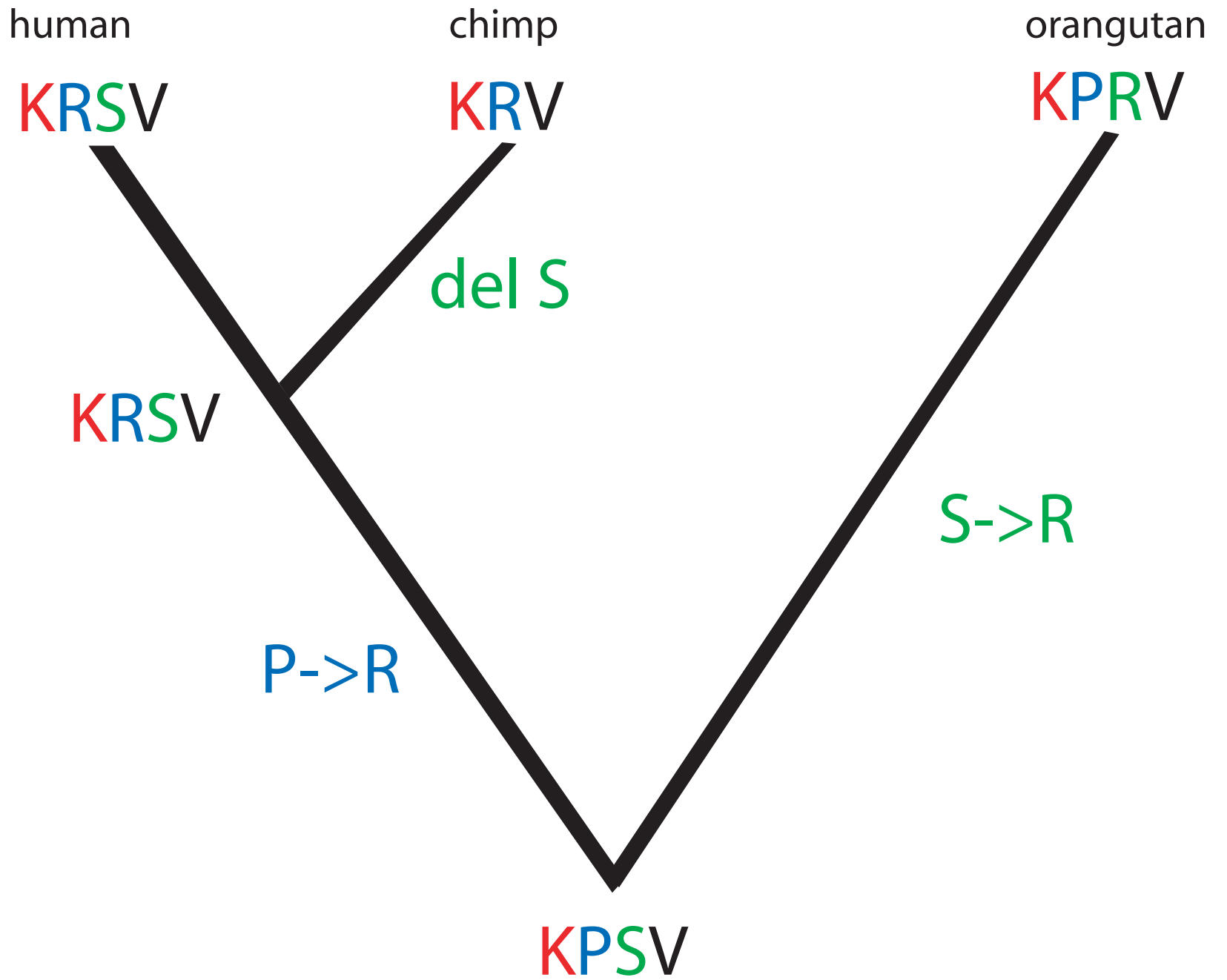
del S

KRSV

S->R

P->R

KPSV



human KRSV  
chimp KRV  
gorilla KSV  
orang KPRV

How should we align these sequences?

human	KRSV		human	KRSV
chimp	KR-V	OR	chimp	K-RV
gorilla	KS-V		gorilla	K-SV
orang	KPRV		orang	KPRV

## Pairwise alignment

---

Gap penalties and a substitution matrix imply a score for any alignment. Pairwise alignment involves finding the alignment that maximizes this score.

- substitution matrices assign positive values to matches or similar substitutions (for example Leucine→Isoleucine).
- unlikely substitutions receive negative scores
- gaps are rare and are heavily penalized (given large negative values).

## Scoring an alignment. Simplest case

---

Costs:

Match	1
Mismatch	0
Gap	-5

Alignment:

<i>Pongo</i>	V	D	E	V	G	G	E	L	G	R	L	F	V	V	P	T	Q
<i>Gorilla</i>	V	E	V	A	G	D	L	G	R	L	L	I	V	Y	P	S	R
<b>Score</b>	1	0	0	0	1	0	0	0	0	0	1	0	1	0	1	0	0

Total score = 5

## Scoring an different alignment. Simplest case

---

Match 1  
Mismatch 0  
Gap -5

<i>Pongo</i>	V	D	E	V	G	G	E	L	G	R	L	-	F	V	V	P	T	Q
<i>Gorilla</i>	V	-	E	V	A	G	D	L	G	R	L	L	I	V	Y	P	S	R
<b>Score</b>	1	-5	1	1	0	1	0	1	1	1	1	-5	0	1	0	1	0	0

Total score = 0



# BLOSUM 62 Substitution matrix

---

	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
A	4																			
R	-1	5																		
N	-2	0	6																	
D	-2	-2	1	6																
C	0	-3	-3	-3	9															
Q	-1	1	0	0	-3	5														
E	-1	0	0	2	-4	2	5													
G	0	-2	0	-1	-3	-2	-2	6												
H	-2	0	1	-1	-3	0	0	-2	8											
I	-1	-3	-3	-3	-1	-3	-3	-4	-3	4										
L	-1	-2	-3	-4	-1	-2	-3	-4	-3	2	4									
K	-1	2	0	-1	-3	1	1	-2	-1	-3	-2	5								
M	-1	-1	-2	-3	-1	0	-2	-3	-2	1	2	-1	5							
F	-2	-3	-3	-3	-2	-3	-3	-3	-1	0	0	-3	0	6						
P	-1	-2	-2	-1	-3	-1	-1	-2	-2	-3	-3	-1	-2	-4	7					
S	1	-1	1	0	-1	0	0	0	-1	-2	-2	0	-1	-2	-1	4				
T	0	-1	0	-1	-1	-1	-1	-2	-2	-1	-1	-1	-1	-2	-1	1	5			
W	-3	-3	-4	-4	-2	-2	-3	-2	-2	-3	-2	-3	-1	1	-4	-3	-2	11		
Y	-2	-2	-2	-3	-2	-1	-2	-3	2	-1	-1	-2	-1	3	-3	-2	-2	2	7	
V	0	-3	-3	-3	-1	-2	-2	-3	-3	3	1	-2	1	-1	-2	-2	0	-3	-1	4
	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V

## Scoring an alignment with the BLOSUM 62 matrix

---

<i>Pongo</i>	V	D	E	V	G	G	E	L	G	R	L	F	V	V	P	T	Q
<i>Gorilla</i>	V	E	V	A	G	D	L	G	R	L	L	I	V	Y	P	S	R
<b>Score</b>	4	2	-2	0	6	-6	-3	-4	-2	-2	4	0	4	-1	7	4	1

The score for the alignment is

$$D_{ij} = \sum_k d_{ij}^{(k)}$$

If  $i$  indicates *Pongo* and  $j$  indicates *Gorilla*

$$D_{ij} = 12$$

## Scoring an alignment with gaps

---

If the GP is -8:

<i>Pongo</i>	V	D	E	V	G	G	E	L	G	R	L	-	F	V	V	P	T	Q
<i>Gorilla</i>	V	-	E	V	A	G	D	L	G	R	L	L	I	V	Y	P	S	R
<b>Score</b>	4	-8	5	5	0	6	2	4	6	5	4	-8	0	4	-1	7	4	1

By introducing gaps we have improved the score:

$$D_{ij} = 40$$

## Gap Penalties

---

Gaps are penalized more heavily than substitutions to avoid alignments like this:

<i>Pongo</i>	VDEVGGE-LGRLFVVPTQ
<i>Gorilla</i>	VDEVGG-WLGRLFVVPTQ

## Gap Penalties

---

Because multiple residues are often inserted or deleted at the same time, *affine gap penalties* are often used:

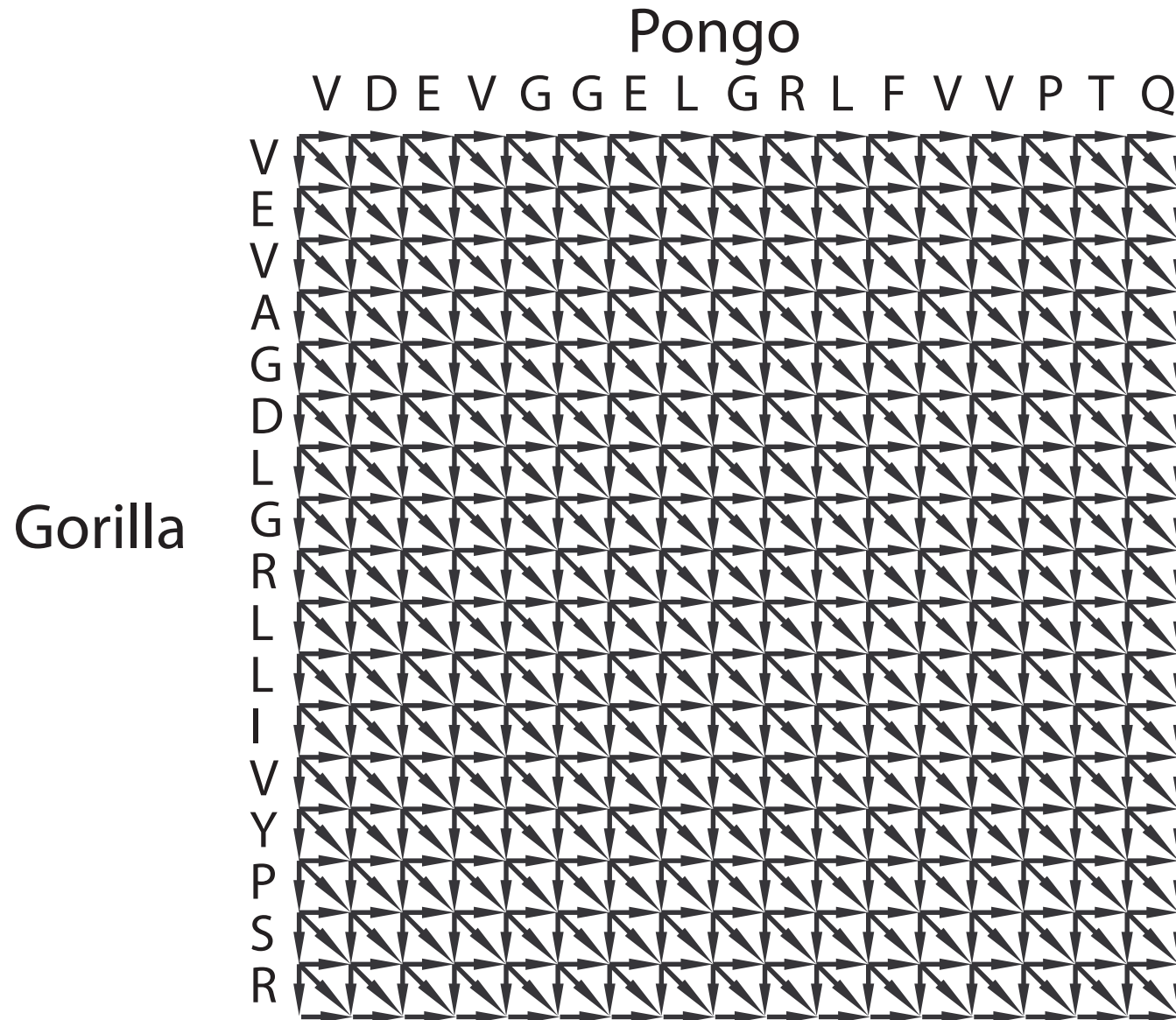
$$GP = GO + lGE$$

where:

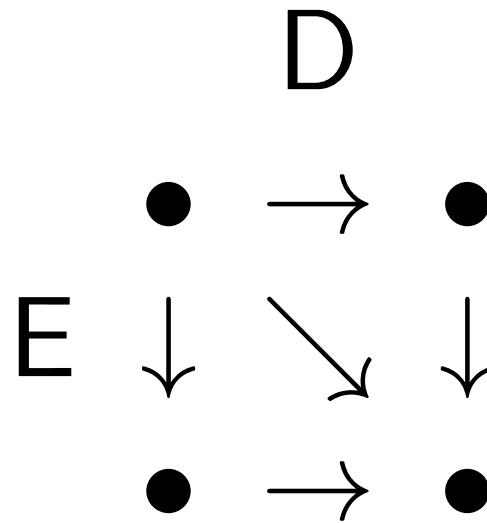
- GP is the gap penalty.
- GO is the “gap-opening penalty”
- GE is the “gap-extension penalty”
- $l$  is the length of the gap

# Finding an optimal alignment

---

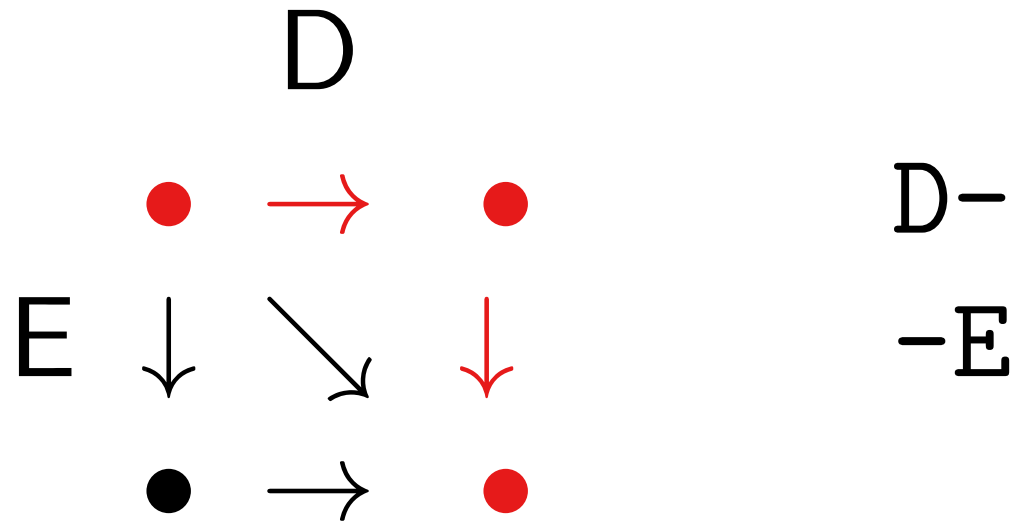


# Aligning two sequences, each with length = 1



# Alignment 1

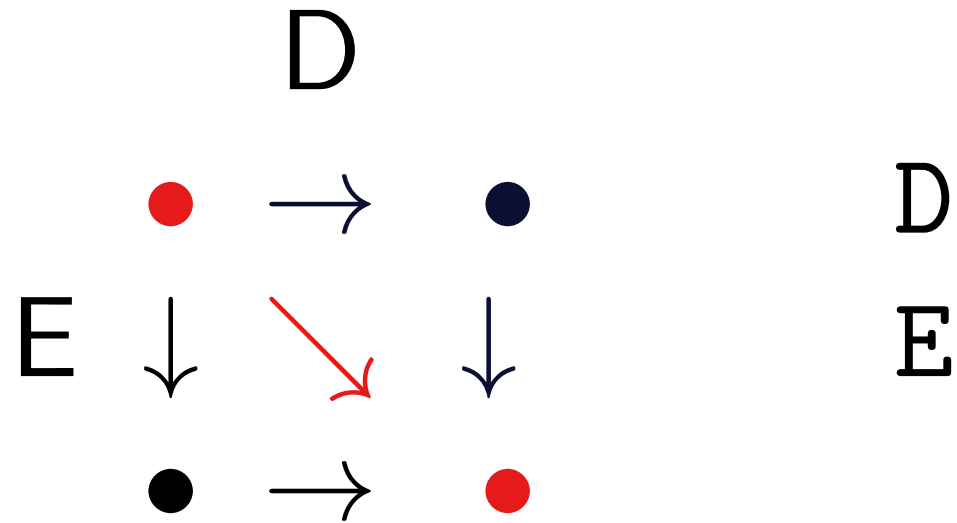
---





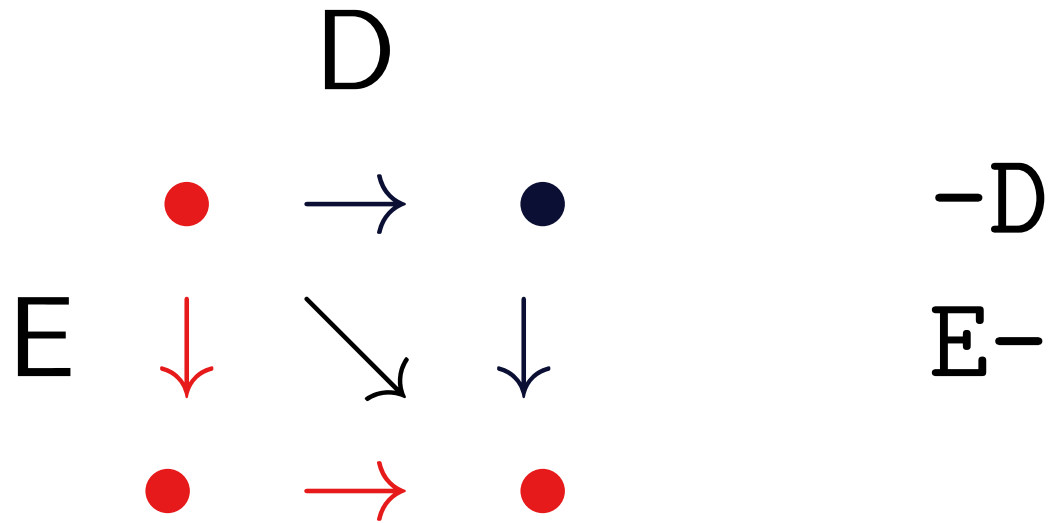
# Alignment 2

---

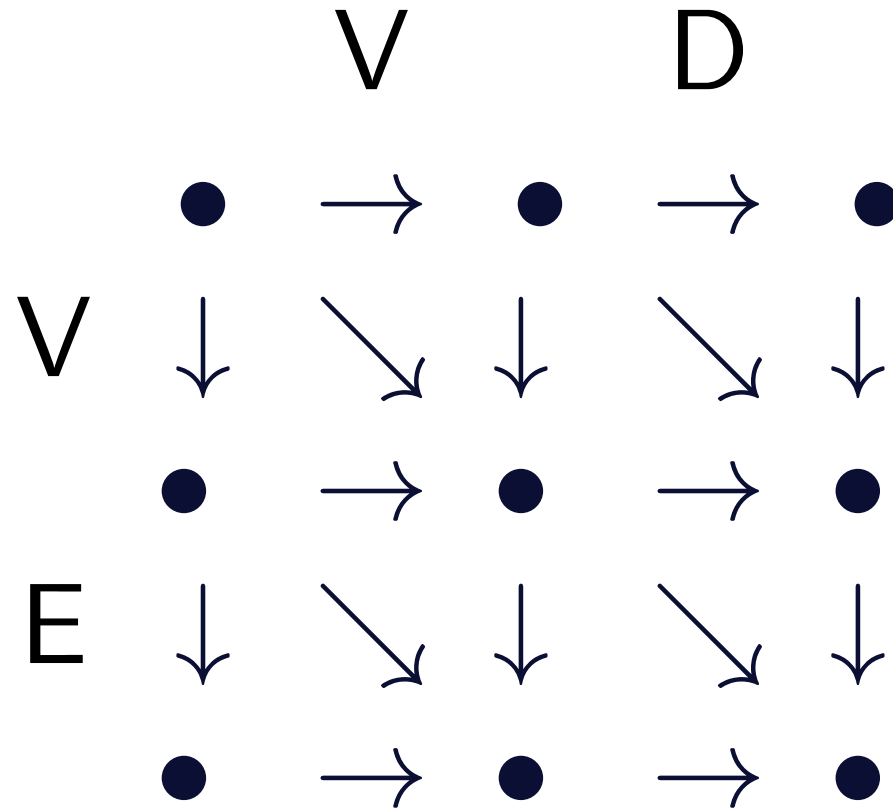


# Alignment 3

---

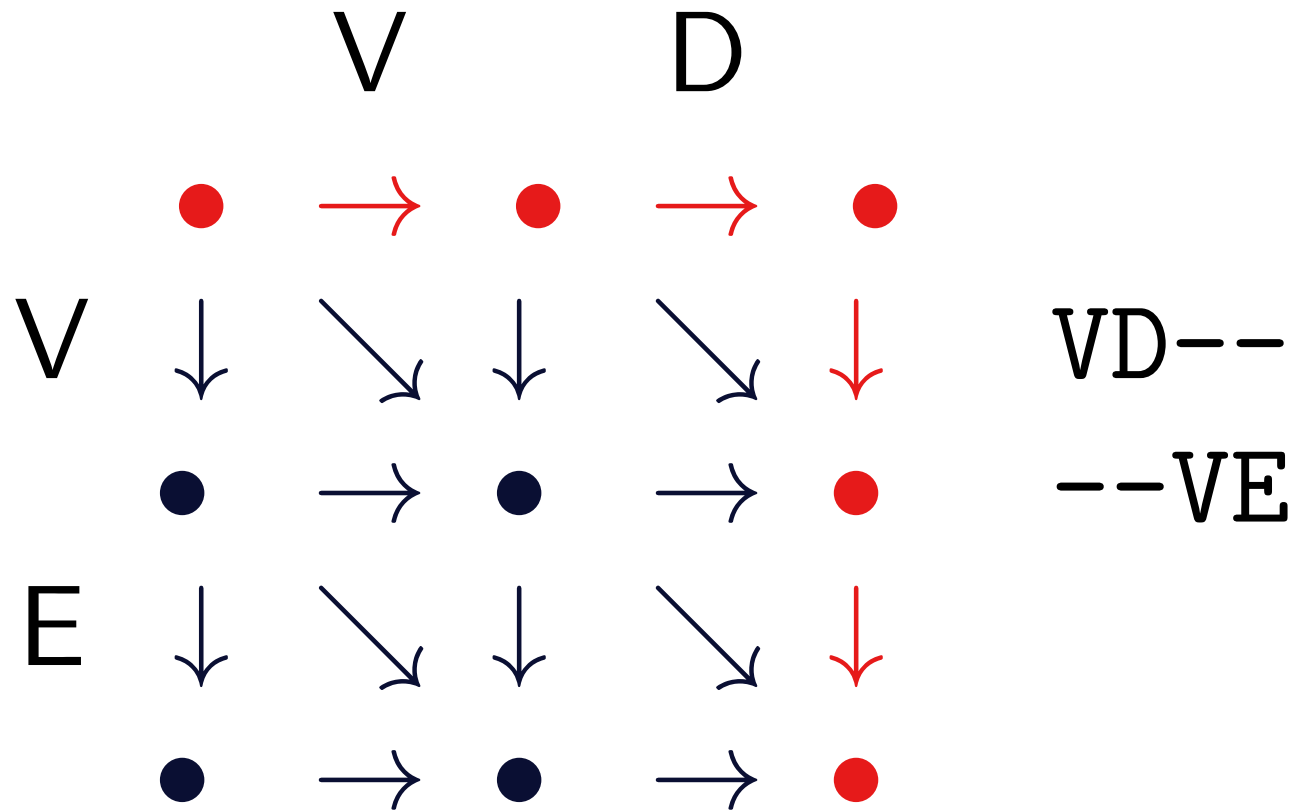


## Longer sequences – up to 2 amino acids!



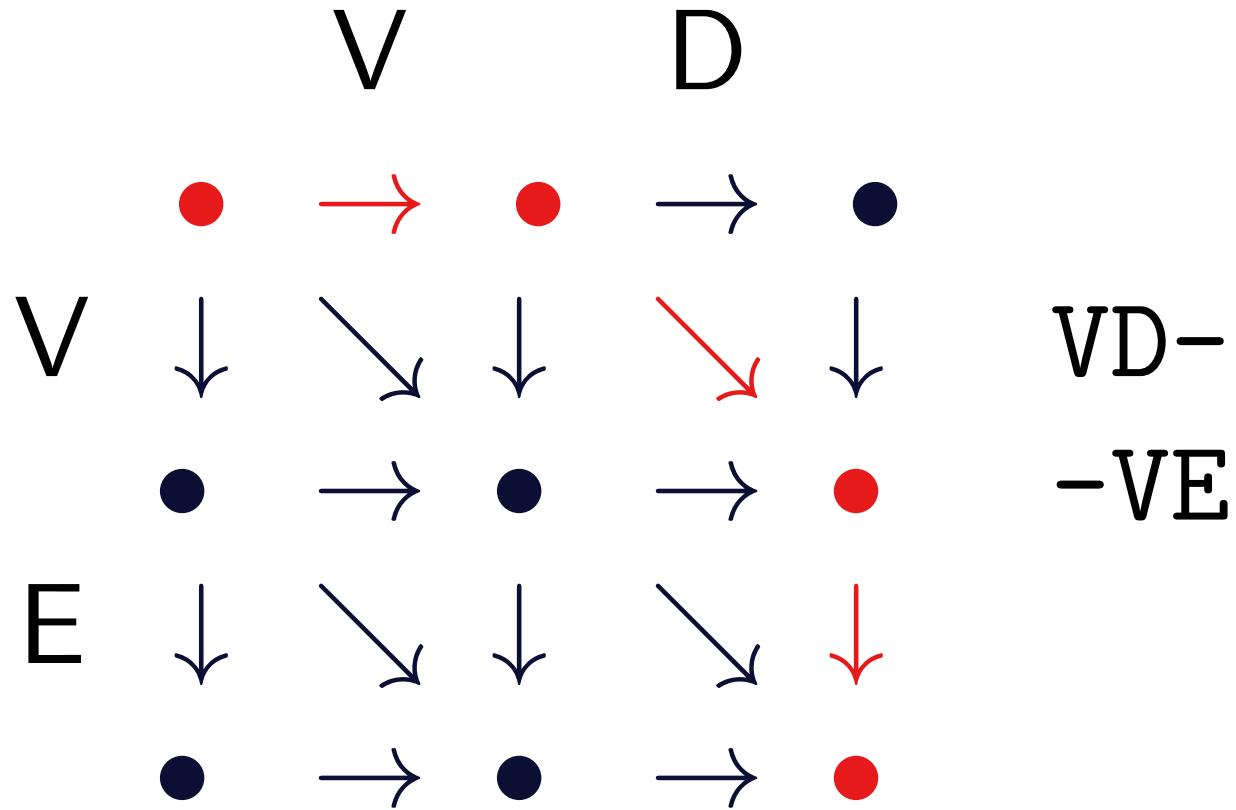
# Alignment 1

---



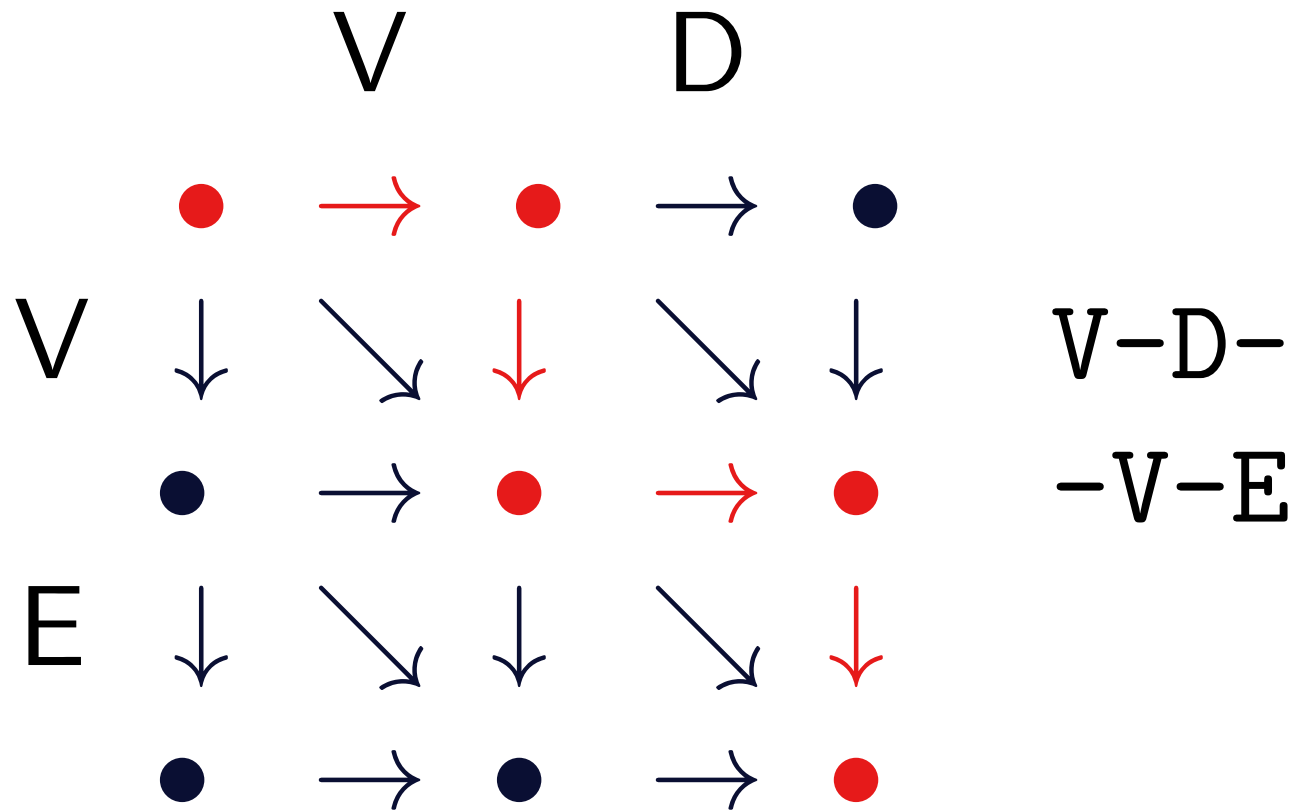
# Alignment 2

---



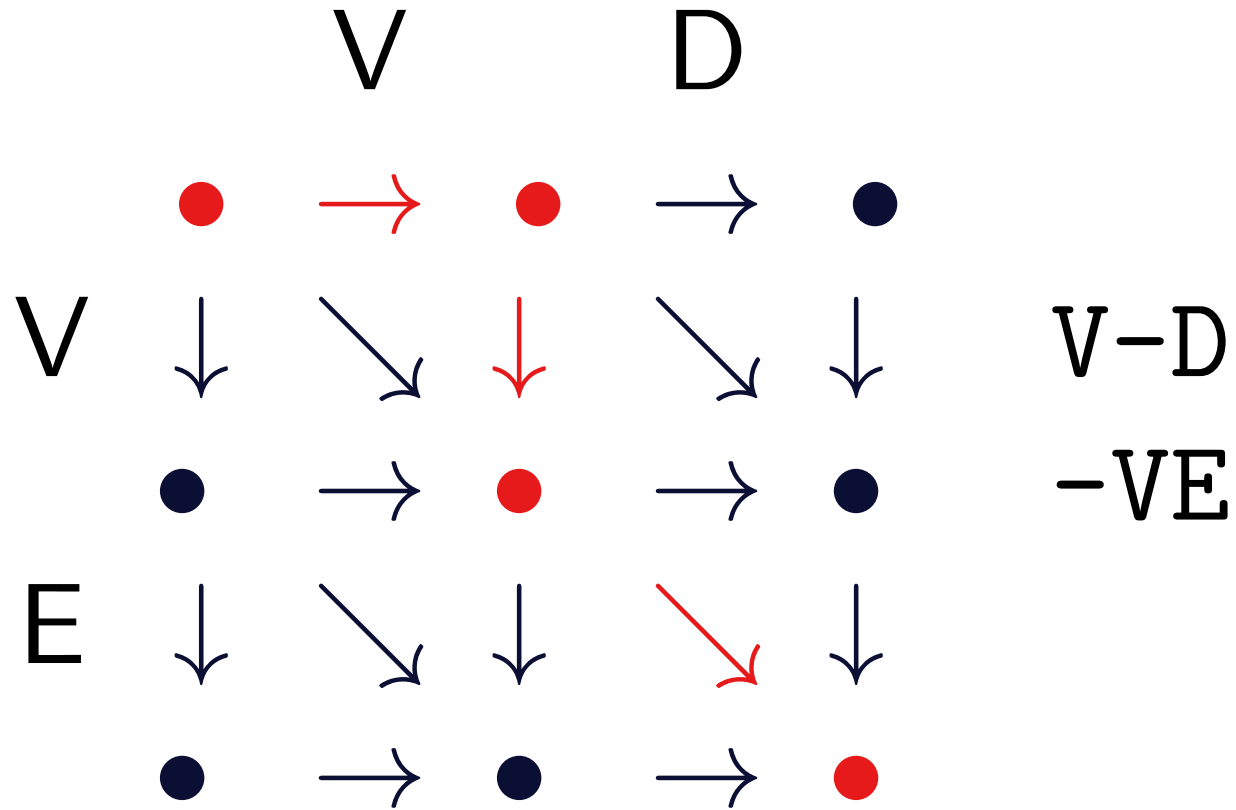
# Alignment 3

---



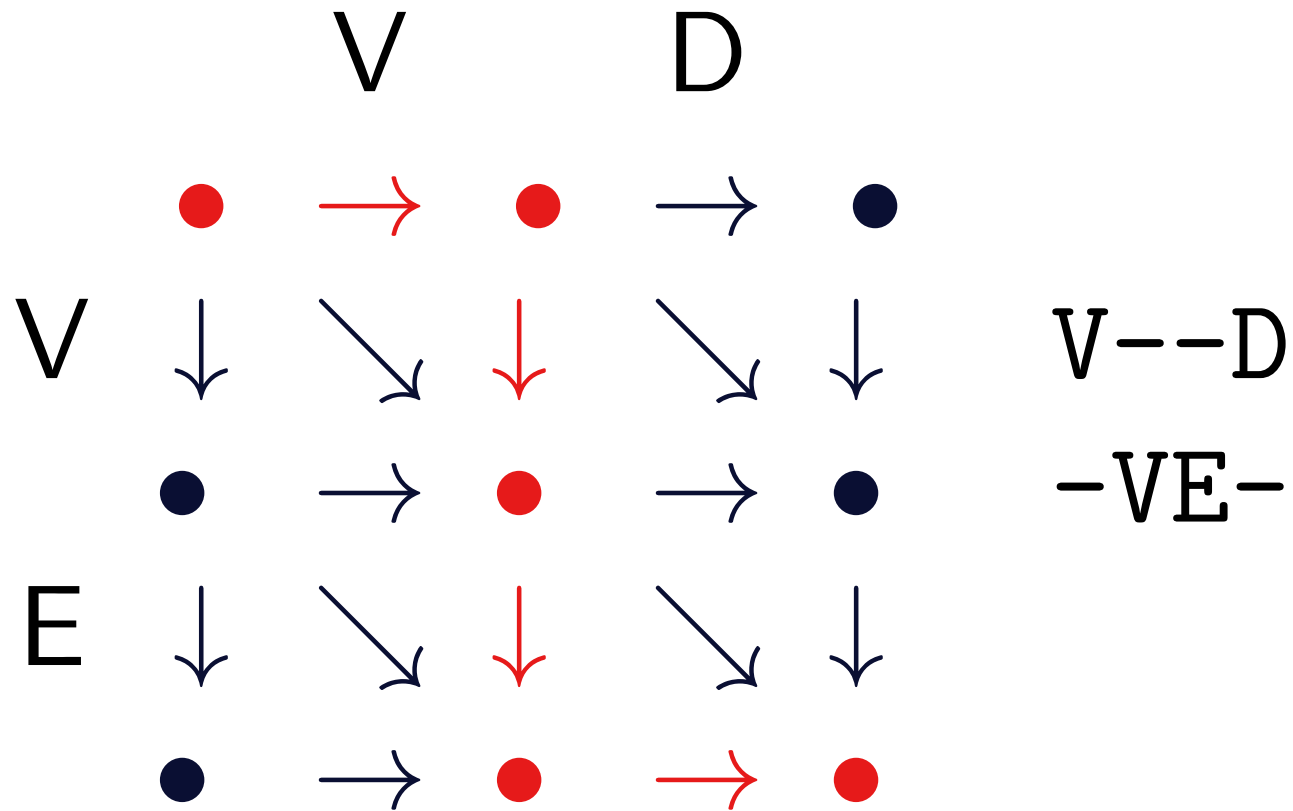
# Alignment 4

---



# Alignment 5

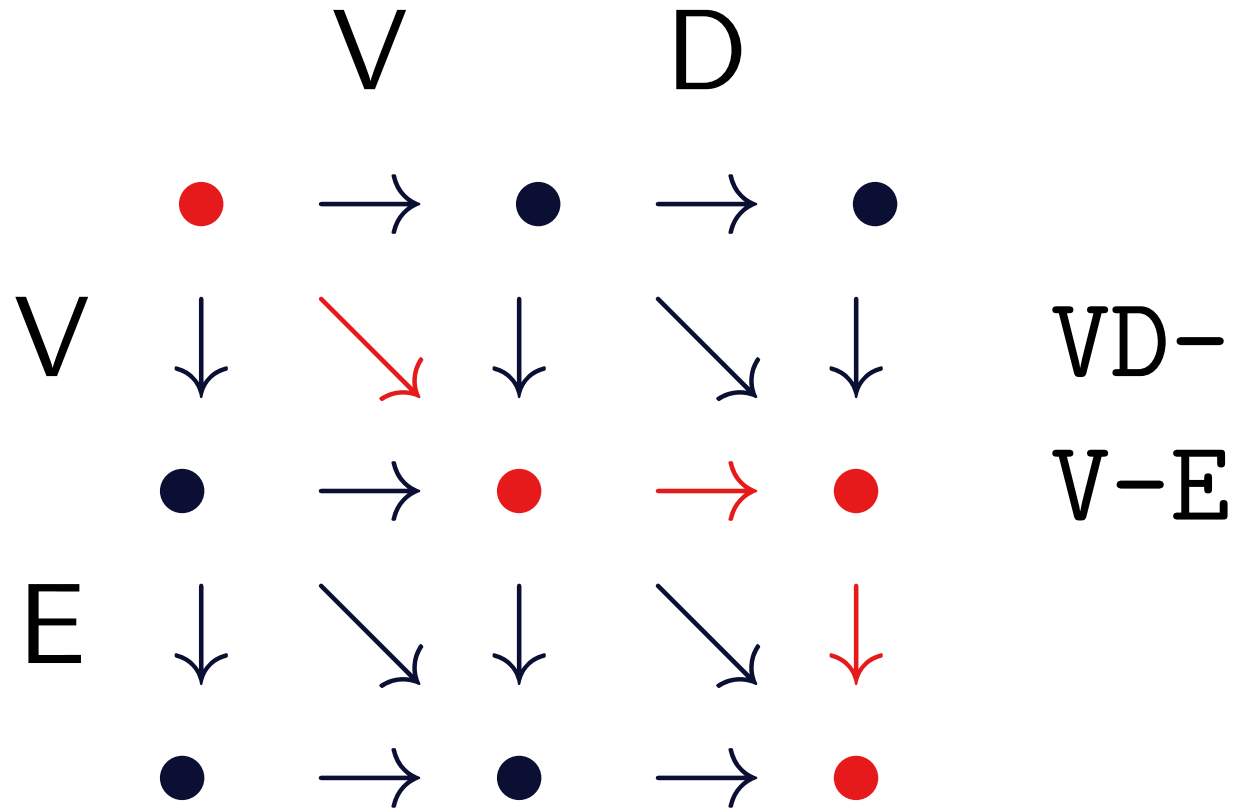
---





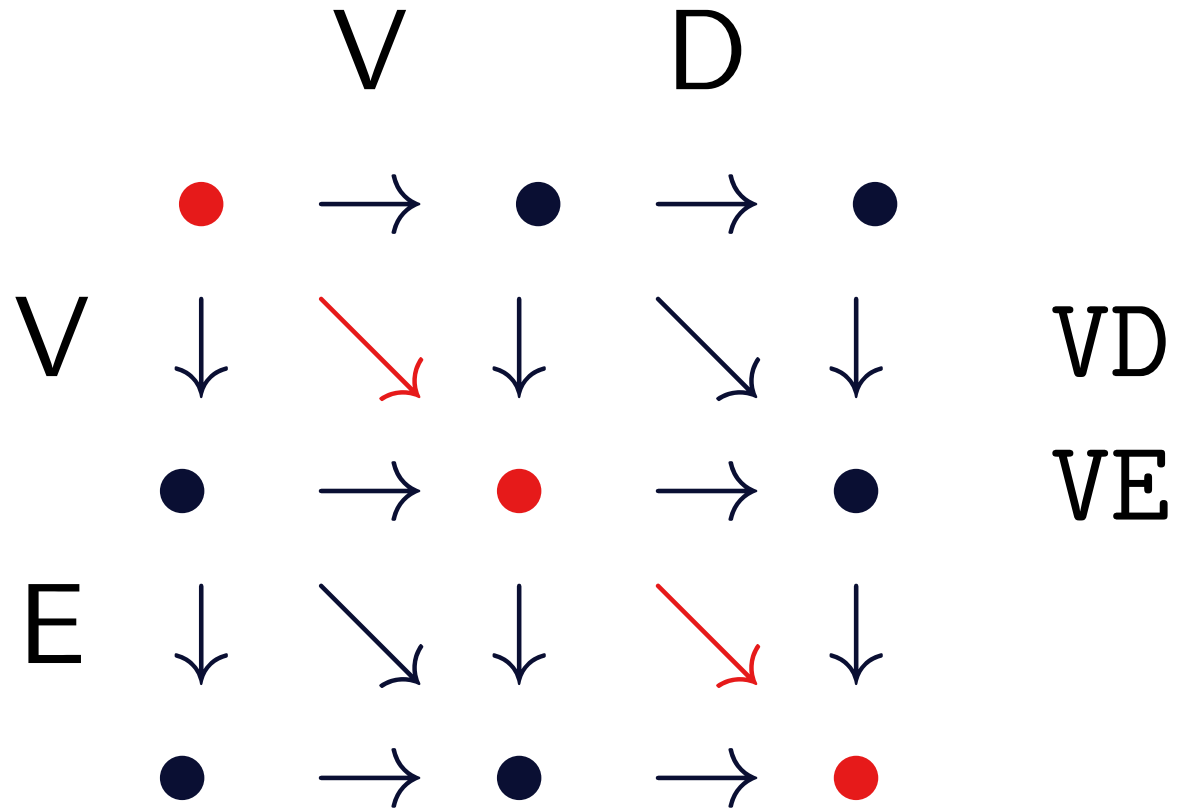
# Alignment 6

---



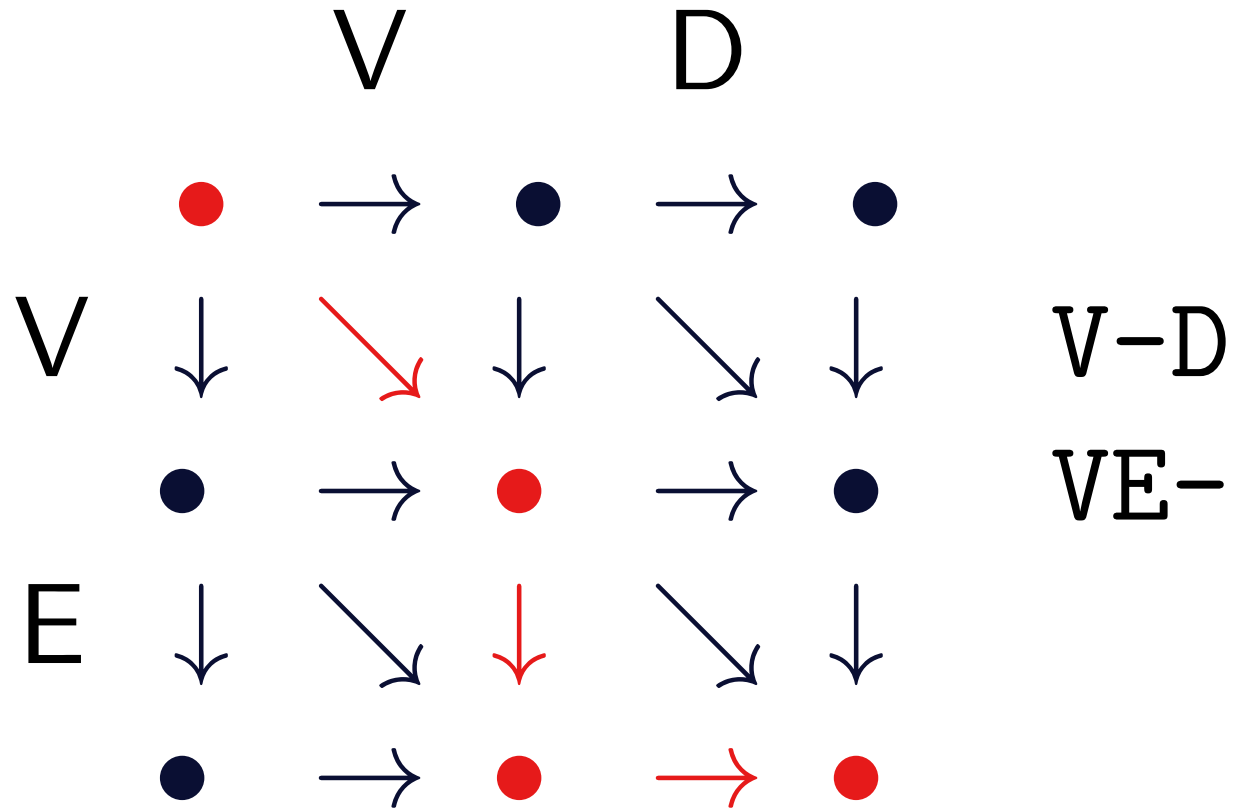
# Alignment 7

---



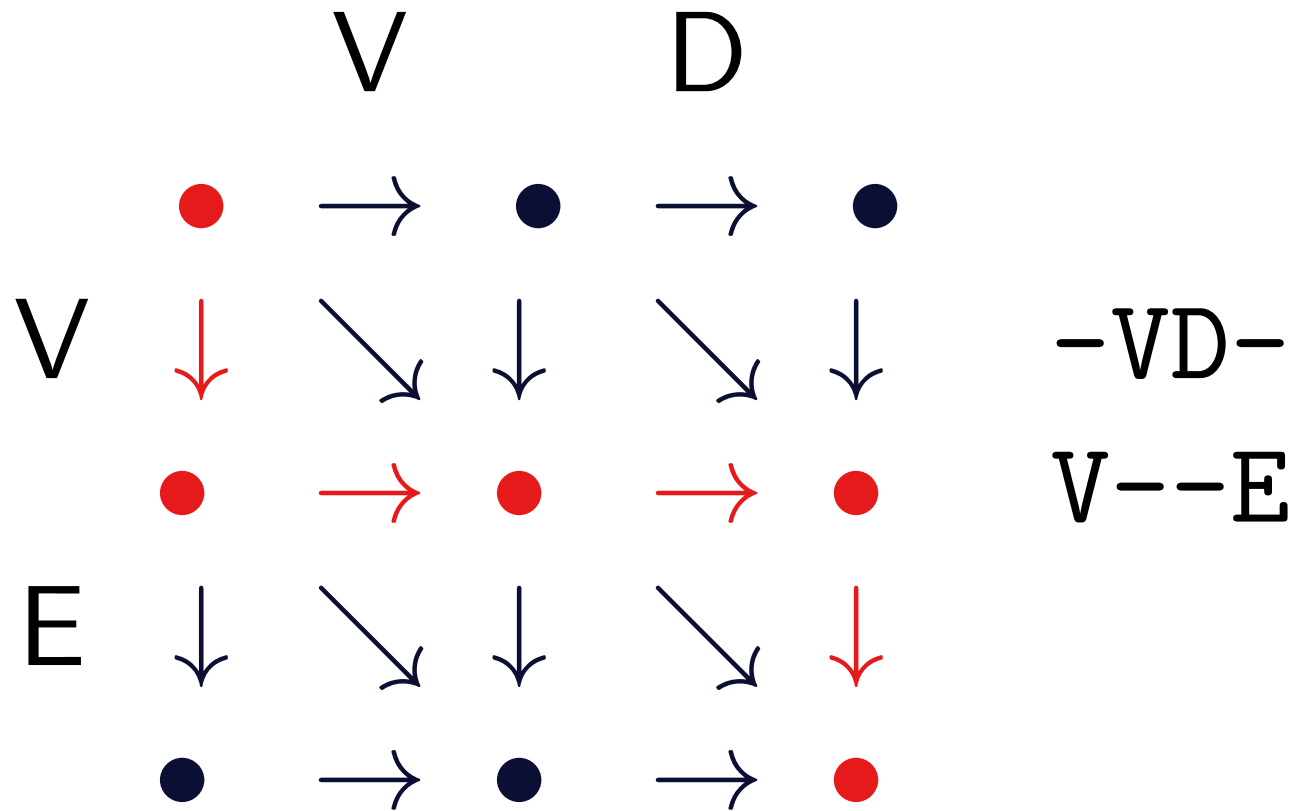
# Alignment 8

---



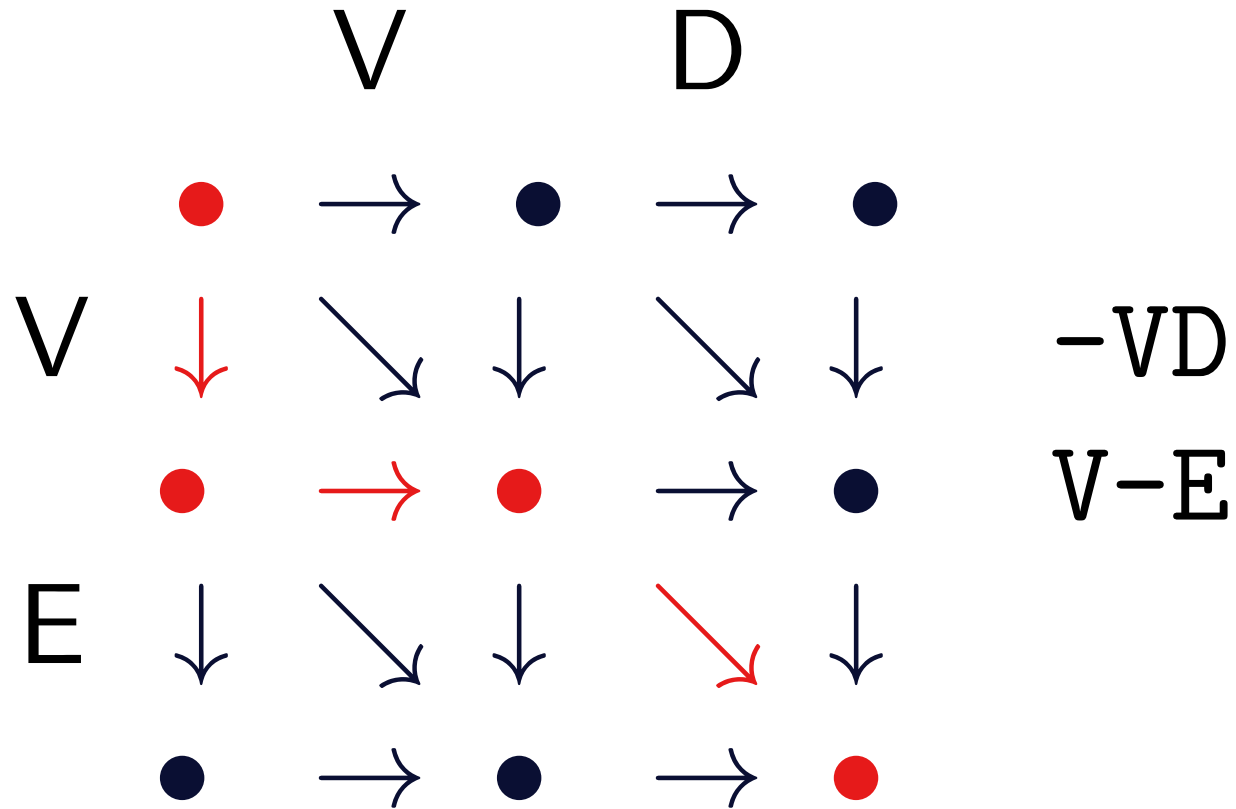
# Alignment 9

---



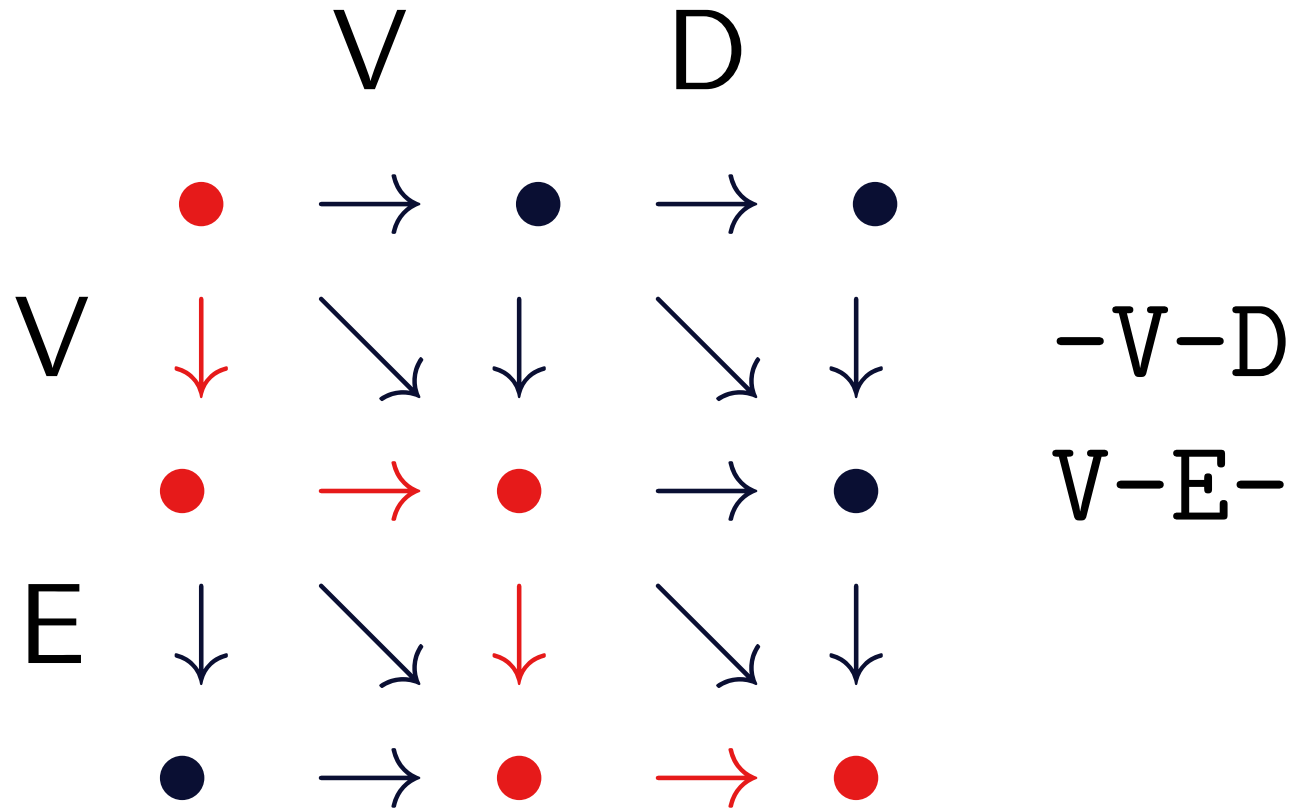
# Alignment 10

---



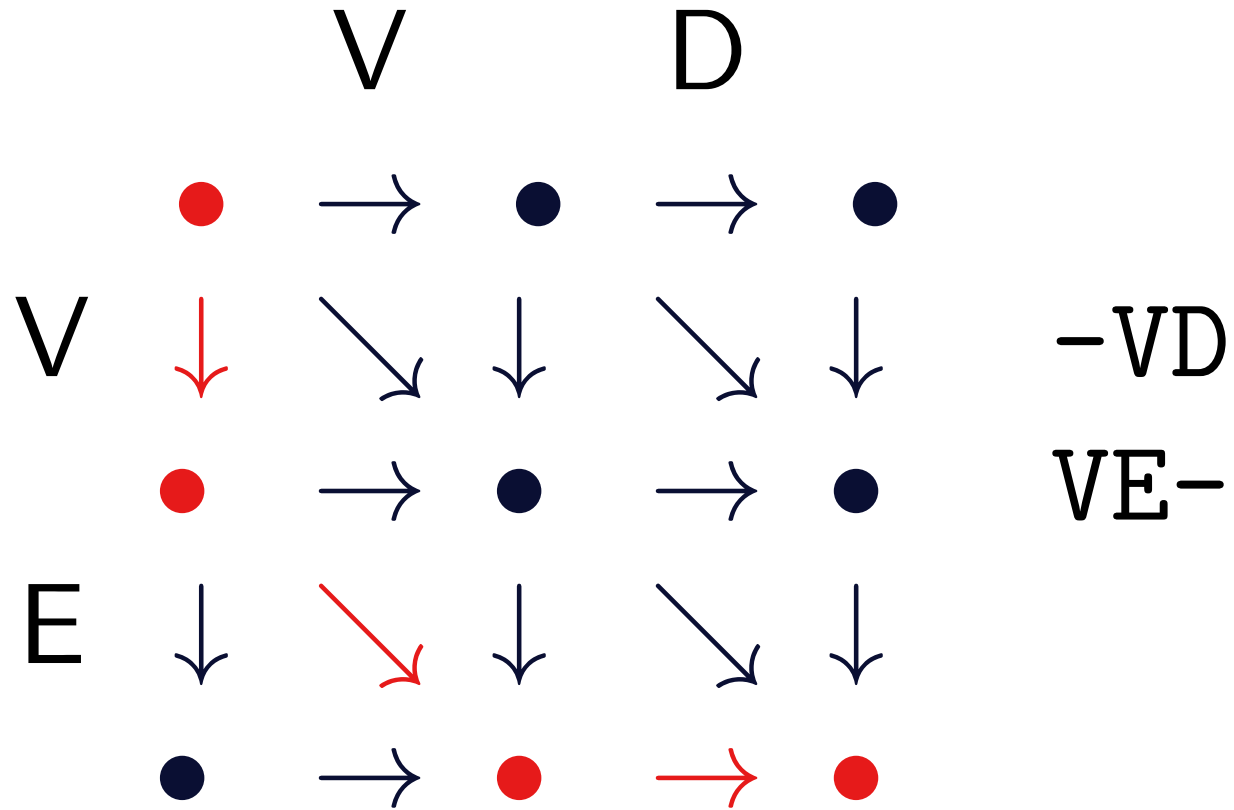
# Alignment 11

---



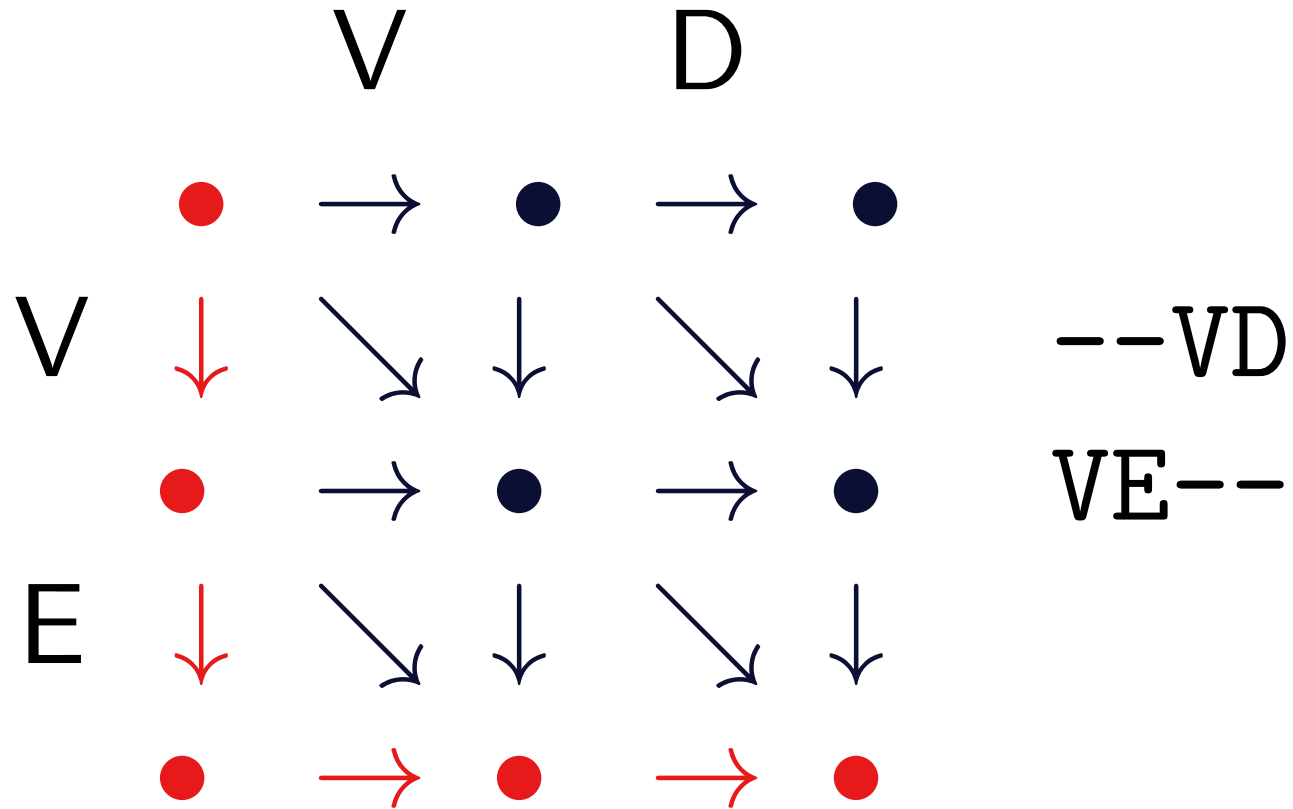
# Alignment 12

---



# Alignment 13

---



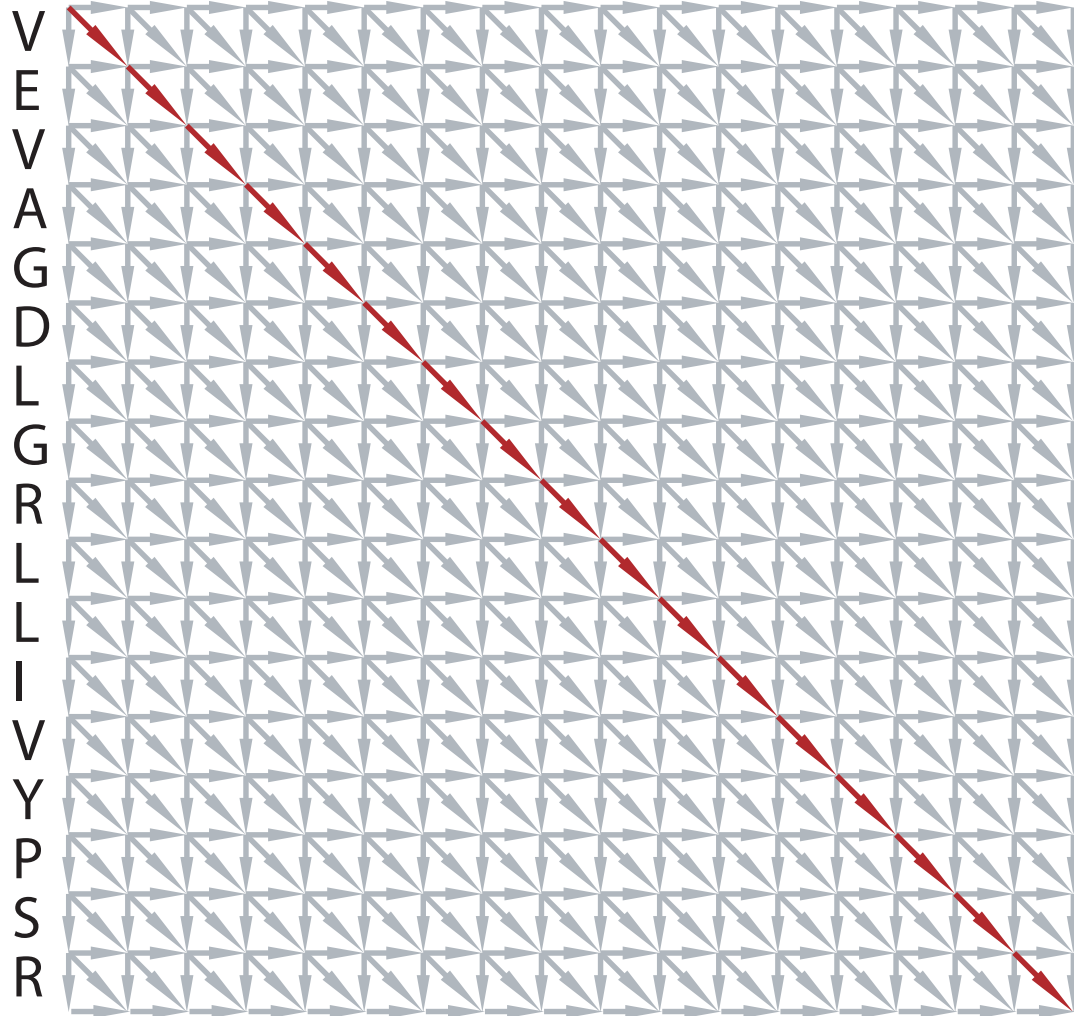


<i>Pongo</i>	V	D	E	V	G	G	E	L	G	R	L	F	V	V	P	T	Q
<i>Gorilla</i>	V	E	V	A	G	D	L	G	R	L	L	I	V	Y	P	S	R
<b>Score</b>	4	2	-2	0	6	-6	-3	-4	-2	-2	4	0	4	-1	7	4	1

## Pongo

V D E V G G E L G R L F V V P T Q

Gorilla





length	Seq # 1	length	Seq # 2	# alignments
1		1		3
2		2		13
3		3		63
4		4		321
5		5		1,683
6		6		8,989
7		7		48,639
8		8		265,729
9		9		1,462,563
⋮		⋮		⋮
17		17		1,425,834,724,419

## Needleman-Wunsch algorithm (paraphrased)

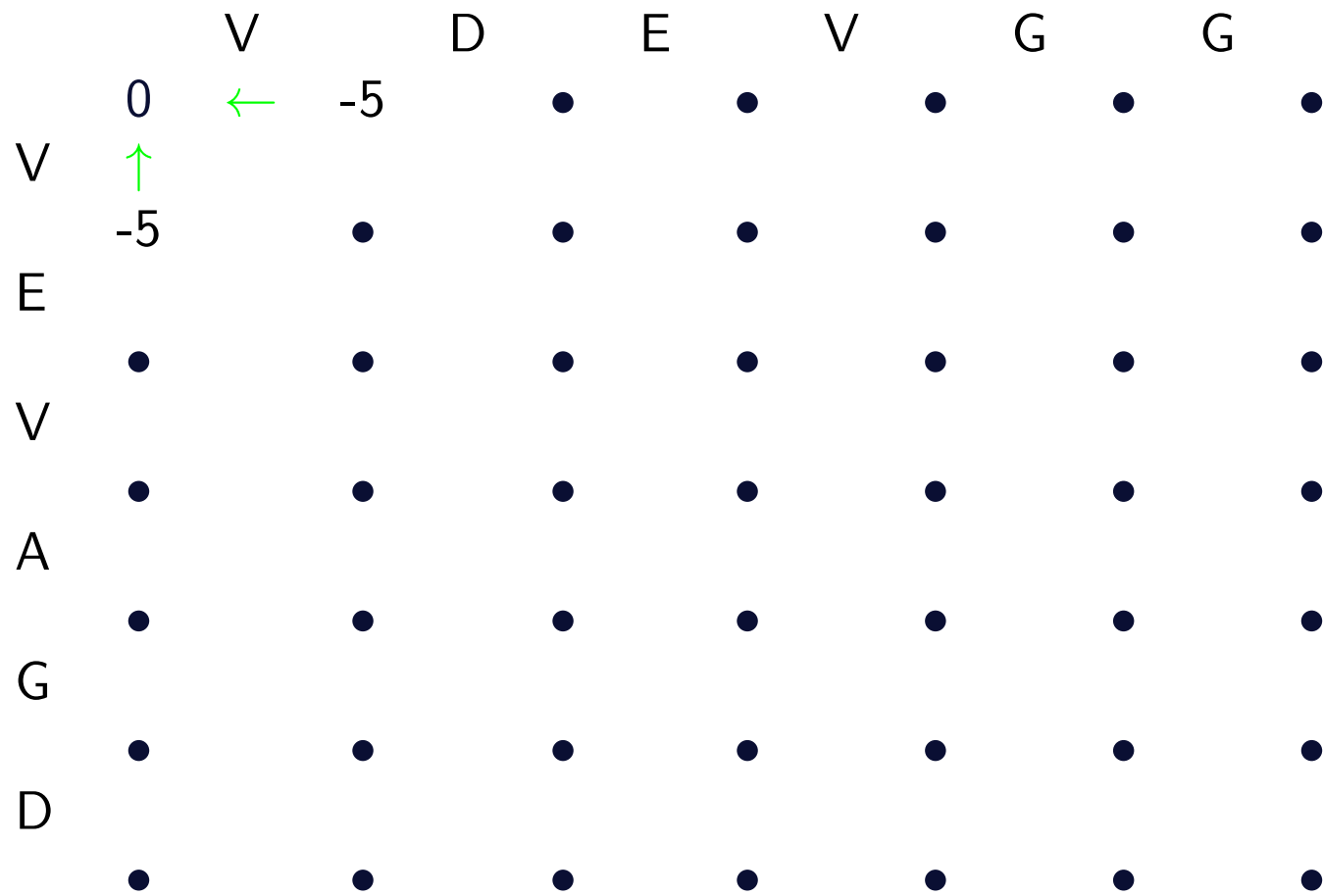
- Work from the top left (beginning of both sequences)
- For each cell store the highest score possible for that cell and a “back” pointer to tell point to the previous step in the best path
- When you reach the lower right corner, you know the optimal score and the back pointers tell you the alignment.

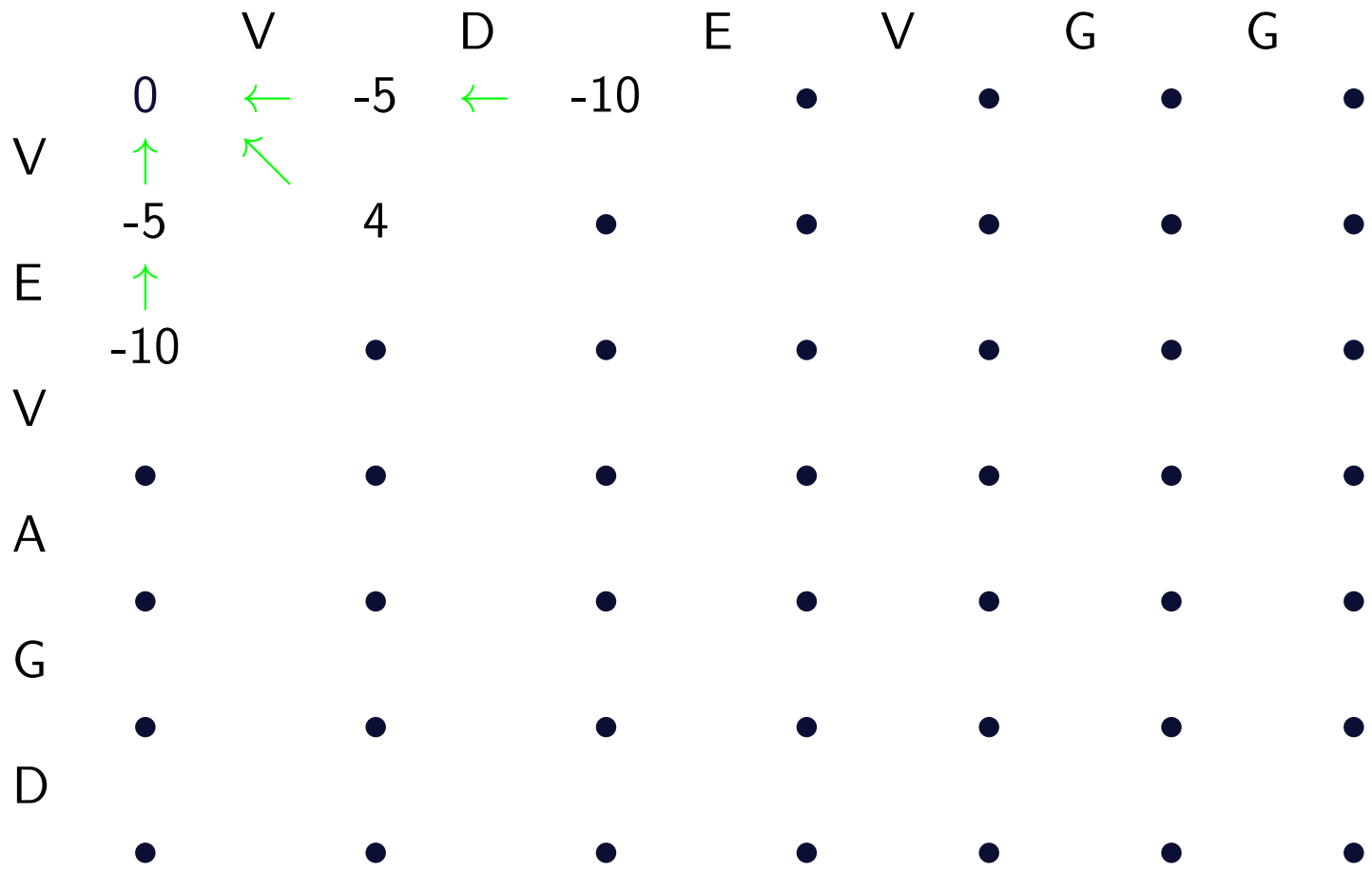
The highest-score calculation at each cell only depends on its the cell's three possible previous neighbors.

If one sequence is length  $N$ , and the other is length  $M$ , then Needleman-Wunsch only takes  $\approx 6NM$  calculations.

But there are a **much** larger number of possible alignments.







		V	D	E	V	G	G
V	0	←	←	←			
E	-5	↖	←				
V	-10						
A	-15						
G							
D							

V  
E  
V  
A  
G  
D

0  
-5  
-10  
-15

←  
↖

-5  
4  
-1  
•  
•  
•  
•

←

-10  
-1  
•  
•  
•  
•

←

-15  
•  
•  
•  
•  
•  
•

V

•  
•  
•  
•  
•  
•  
•

G

•  
•  
•  
•  
•  
•  
•

G

•  
•  
•  
•  
•  
•  
•





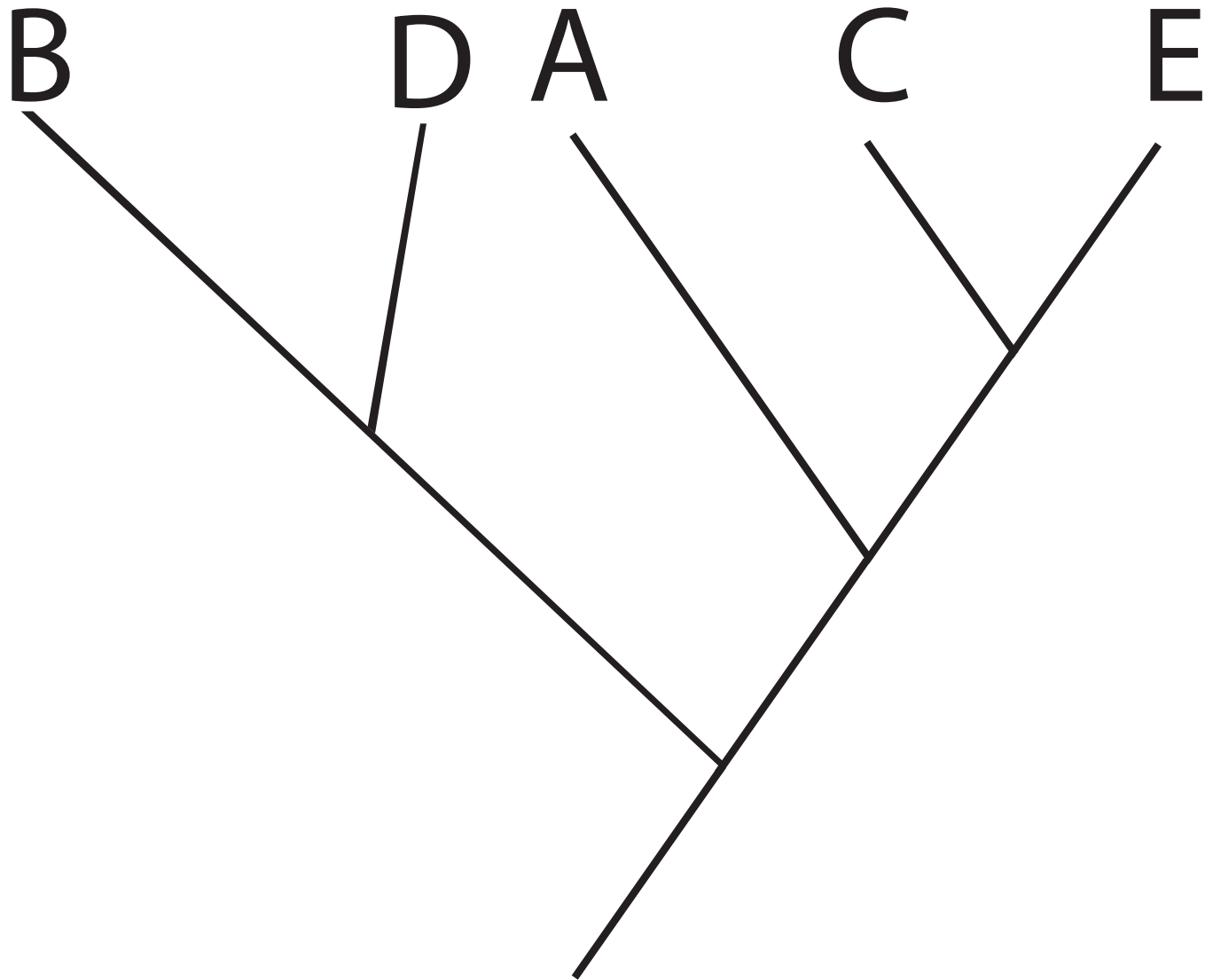


		V	D	E	V	G	G	E	L	G	R
V	0	← -5	← -10	← -15	← -20	← -25	← -30	← -35	← -40	← -45	← -50
	↑	↖			↖						
	-5	4	← -1	← -6	-11	← -16	← -21	← -26	← -31	← -36	← -41
E	↑	↑	↖					↖			
	-10	-1	6	↖	4	← -1	← -6	← -11	-16	← -21	← -26
V	↑	↖	↑	↖		↖					
	-15	-6	1	4	8	← 3	← -2	← -7	← -12	← -17	← -22
A	↑	↑	↑	↖		↖				↖	
	-20	-11	-4	0	4	8	3	← -2	← -7	-12	← -17
G	↑	↑	↑	↑	↑	↖				↖	
	-25	-16	-9	-5	-1	10	14	← 9	← 4	-1	← -6
D	↑	↑	↖	↖	↑	↑	↖				
	-30	-21	-10	-7	-6	5	9	16	← 11	← 6	← 1
L	↑	↑	↑	↑	↖	↑	↑	↑	↖		
	-35	-26	-15	-12	-6	0	4	11	20	← 15	← 10
G	↑	↑	↑	↖	↑	↖		↑	↑	↖	
	-40	-31	-20	-17	-11	0	6	6	15	26	← 21
R	↑	↑	↑	↖	↑	↑	↑	↖	↑	↑	↖
	-45	-36	-25	-20	-16	-5	1	6	10	21	31
L	↑	↑	↑	↑	↖	↑	↑	↑	↖	↑	↑
	-50	-41	-30	-25	-19	-10	-4	1	10	16	26
L	↑	↑	↑	↑	↖	↑	↑	↑	↖	↑	↑
	-55	-46	-35	-30	-24	-15	-9	-4	5	11	21
I	↑	↑	↑	↑	↖	↑	↑	↑	↑	↑	↑
	-60	-51	-40	-35	-27	-20	-14	-9	0	6	16
V	↑	↖	↑	↑	↖	↑	↑	↑	↑	↑	↑



# Aligning multiple sequences

---



# Progressive alignment

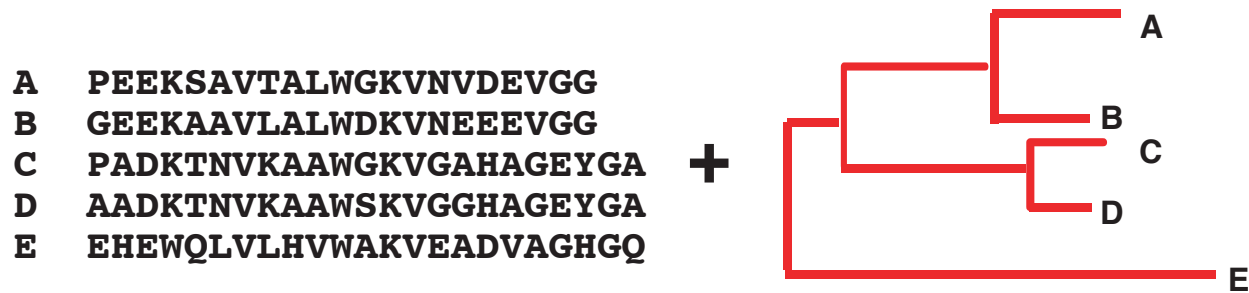
---

Devised by Feng and Doolittle 1987 and Higgins and Sharp, 1988. An approximate method for producing multiple sequence alignments using a guide tree.

- Perform pairwise alignments to produce a distance matrix
- Produce a guide tree from the distances
- Use the guide tree to specify the ordering used for aligning sequences, closest to furthest.

		<b>pairwise alignment</b>							
A	PEEKSAVTALWGKVVNDEVGG	→	A	-					
B	GEEKAAVLALWDKVNNEEVGG		B	.17	-				
C	PADKTNVKAAWGKVGGAHAGEYGA		C	.59	.60	-			
D	AADKTNVKAAWSKVGGHAGEYGA		D	.59	.59	.13	-		
E	EHEWQLVLHVWAKVEADVAGHGQ		E	.77	.77	.75	.75	-	

↓ **tree inference**



↓ **alignment stage**

A	PEEKSAVTALWGKVN--VDEVGG
B	GEEKAAVLALWDKVN--EEEVGG
C	PADKTNVKAAWGKVGGAHAGEYGA
D	AADKTNVKAAWSKVGGHAGEYGA
E	EHEWQLVLHVWAKVEADVAGHGQ

## **Alignment stage of progressive alignments**

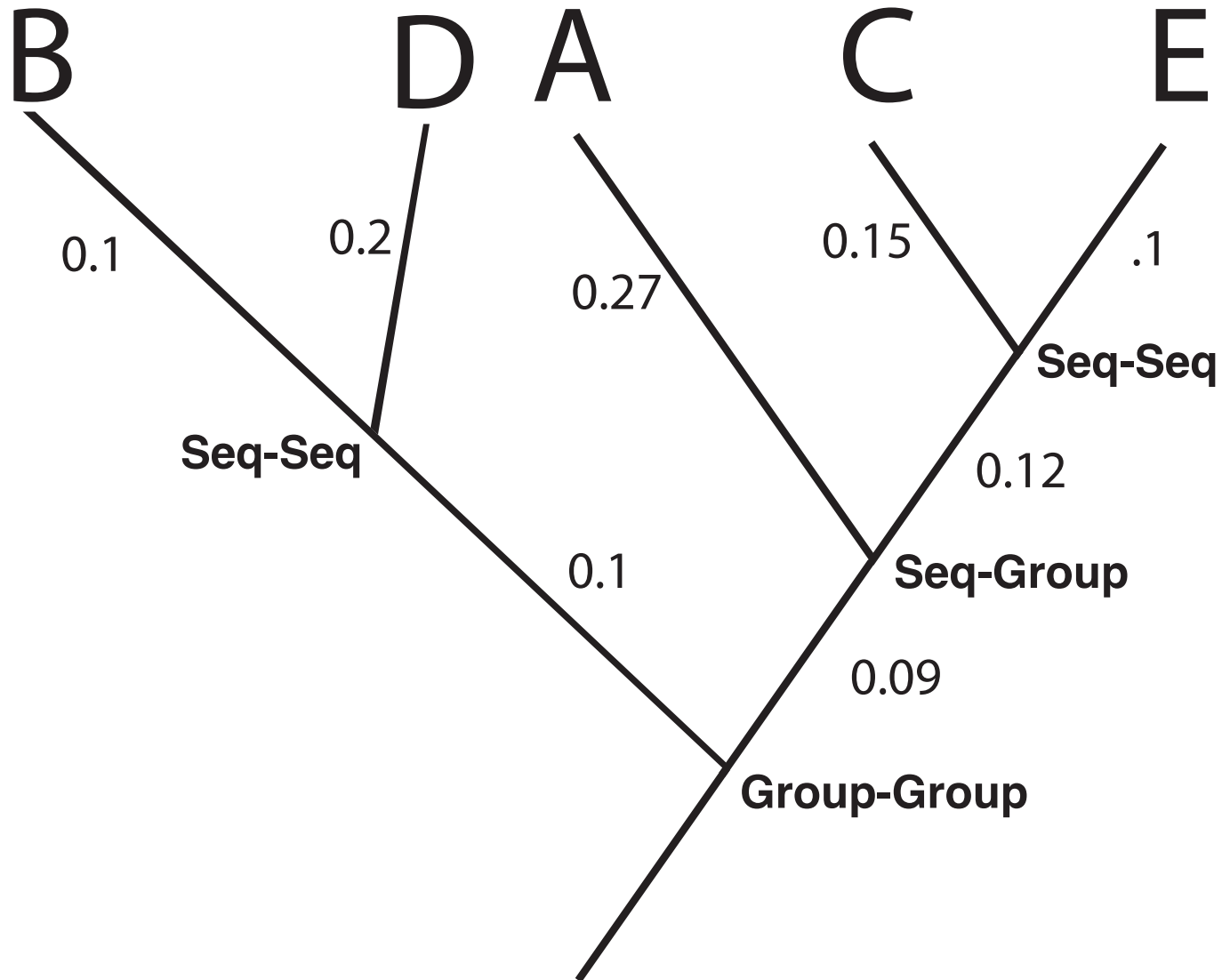
Sequences of clades become grouped into “profiles” as the algorithm descends the tree. The next youngest internal nodes is selected at each step to create a new profile. Alignment at each step involves

- Sequence-Sequence
- Sequence-Profile
- Profile-Profile



# Aligning multiple sequences

---



# Profile to Profile alignment

---



## Profile to profile alignments

---

Adding a gap to a profile means that every member of that group of sequences gets a gap at that position of the sequence.

Usually the scores for each edge in the Needleman-Wünsch graph are calculated using a “sum of pairs” scoring system.

`clustal W`<sup>1</sup> uses weights assigned to each sequence in a profile group to downweight closely related sequences so that they are not overrepresented.

---

<sup>1</sup>Thompson, Higgins, and Gibson. **Nuc. Acids. Res.** 1994

Profile 1			Profile 2		
Seq	weight	AA	Seq	weight	AA
taxon A	0.3	V	taxon B	0.15	V
taxon C	0.24	A	taxon D	0.25	M
taxon E	0.19	I			

$$\begin{aligned}
D_{P1,P2} &= \frac{\sum_i \sum_j w_i w_j d_{ij}}{n_i n_j} \\
&= \frac{1}{6} [d(V, V)w_A w_B + d(V, M)w_A w_D + d(A, V)w_C w_B \dots \\
&= \dots d(A, M)w_C w_D + d(I, V)w_E w_B + d(I, M)w_E w_D] \\
&= \frac{1}{6} (4 \times 0.3 \times 0.15 + 1 \times 0.3 \times 0.25 + 0 \times 0.24 \times 0.15 \dots \\
&= \dots -1 \times 0.24 \times 0.25 + 3 \times 0.19 \times 0.15 + 1 \times 0.19 \times 0.1 \\
&= 1.46225
\end{aligned}$$

# Dealing with alignment ambiguity<sup>2</sup>

(a)

	X			Y						Z		
	1	2	3	4	5	6	7	8	9	1	1	1
										0	1	2
Outgroup	T	A	G	A	G	C	A	C	T	C	A	G
Taxon A	T	A	G	A	G	C	A	C	T	C	A	G
Taxon B	T	A	G	T	G	A	A	G	C	C	A	G
Taxon C	T	A	G	T	G	A	A	G	C	C	A	G
Taxon D	T	A	G		A	G	C			C	A	G
Taxon E	T	A	G		A	G	C			C	A	G

(b)

	X			Y						Z		
	1	2	3	4	5	6	7	8	9	1	1	1
										0	1	2
Outgroup	T	A	G	A	G	C	A	C	T	C	A	G
Taxon A	T	A	G	A	G	C	A	C	T	C	A	G
Taxon B	T	A	G	T	G	A	A	G	C	C	A	G
Taxon C	T	A	G	T	G	A	A	G	C	C	A	G
Taxon D	T	A	G	A	G	C	-	-	-	C	A	G
Taxon E	T	A	G	A	G	C	-	-	-	C	A	G

(c)

	X			Y						Z		
	1	2	3	4	5	6	7	8	9	1	1	1
										0	1	2
Outgroup	T	A	G	A	G	C	A	C	T	C	A	G
Taxon A	T	A	G	A	G	C	A	C	T	C	A	G
Taxon B	T	A	G	T	G	A	A	G	C	C	A	G
Taxon C	T	A	G	T	G	A	A	G	C	C	A	G
Taxon D	T	A	G	-	-	-	A	G	C	C	A	G
Taxon E	T	A	G	-	-	-	A	G	C	C	A	G

<sup>2</sup>from M. S. Y. Lee, *TREE*, 2001

# Dealing with alignment ambiguity<sup>3</sup> - deletion

(a)

	X			Y						Z				
	1	2	3	4	5	6	7	8	9	1	1	1		
										0	1	2		
Outgroup	T	A	G	A	G	C	A	C	T	C	A	G		
Taxon A	T	A	G	A	G	C	A	C	T	C	A	G		
Taxon B	T	A	G	T	G	A	A	G	C	C	A	G		
Taxon C	T	A	G	T	G	A	A	G	C	C	A	G		
Taxon D	T	A	G			A	G	C				C	A	G
Taxon E	T	A	G			A	G	C				C	A	G

	X			Z		
	1	2	3	1	1	1
				0	1	2
Outgroup	T	A	G	C	A	G
Taxon A	T	A	G	C	A	G
Taxon B	T	A	G	C	A	G
Taxon C	T	A	G	C	A	G
Taxon D	T	A	G	C	A	G
Taxon E	T	A	G	C	A	G

	X			Y						Z		
	1	2	3	4	5	6	7	8	9	1	1	1
										0	1	2
Outgroup	T	A	G	A	G	C	A	C	T	C	A	G
Taxon A	T	A	G	A	G	C	A	C	T	C	A	G
Taxon B	T	A	G	T	G	A	A	G	C	C	A	G
Taxon C	T	A	G	T	G	A	A	G	C	C	A	G
Taxon D	T	A	G	?	?	?	-	-	-	C	A	G
Taxon E	T	A	G	?	?	?	-	-	-	C	A	G

<sup>3</sup>from M. S. Y. Lee, *TREE*, 2001

# Dealing with alignment ambiguity<sup>4</sup>

Elision method (Wheeler, 1995) involves simply concatenating matrices.

	X			Y						Z			X			Y						Z											
	1	2	3	4	5	6	7	8	9	1	1	1	1	2	3	4	5	6	7	8	9	1	1	1	0	1	2						
Outgroup	T	A	G	A	G	C	A	C	T	C	A	G	T	A	G	A	G	C	A	C	T	C	A	G	T	A	G	A	G	C	C	A	G
Taxon A	T	A	G	A	G	C	A	C	T	C	A	G	T	A	G	A	G	C	A	C	T	C	A	G	T	A	G	A	G	C	C	A	G
Taxon B	T	A	G	T	G	A	A	G	C	C	A	G	T	A	G	T	G	A	A	G	C	C	A	G	T	A	G	T	G	A	A	G	C
Taxon C	T	A	G	T	G	A	A	G	C	C	A	G	T	A	G	T	G	A	A	G	C	C	A	G	T	A	G	T	G	A	A	G	C
Taxon D	T	A	G	-	-	-	A	G	C	C	A	G	T	A	G	A	G	C	-	-	-	C	A	G	T	A	G	A	G	C	-	-	-
Taxon E	T	A	G	-	-	-	A	G	C	C	A	G	T	A	G	A	G	C	-	-	-	C	A	G	T	A	G	A	G	C	-	-	-

<sup>4</sup>from M. S. Y. Lee, *TREE*, 2001

## Simultaneous tree inference and alignment

- Ideally we would address uncertainty in both types of inference at the same time
- Allows for application of statistical models to improve inference and assessments of reliability
- Just now becoming feasible: POY (Wheeler, Gladstein, Laet, 2002), Hande1 (Holmes and Bruno, 2001), BA1iPhy (Redelings and Suchard, 2005), and BEAST(Lunter *et al.*, 2005, Drummond and Rambaut, 2003). SATe (Liu *et al* 2009; Yu and Holder software).



# References

---