

Comparaison de séquences 2 à 2

Alain Denise

Bioinformatique

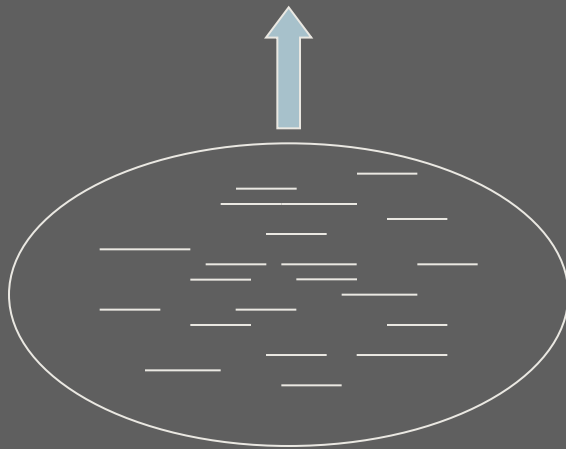
LRI Orsay

Université Paris-Sud 11

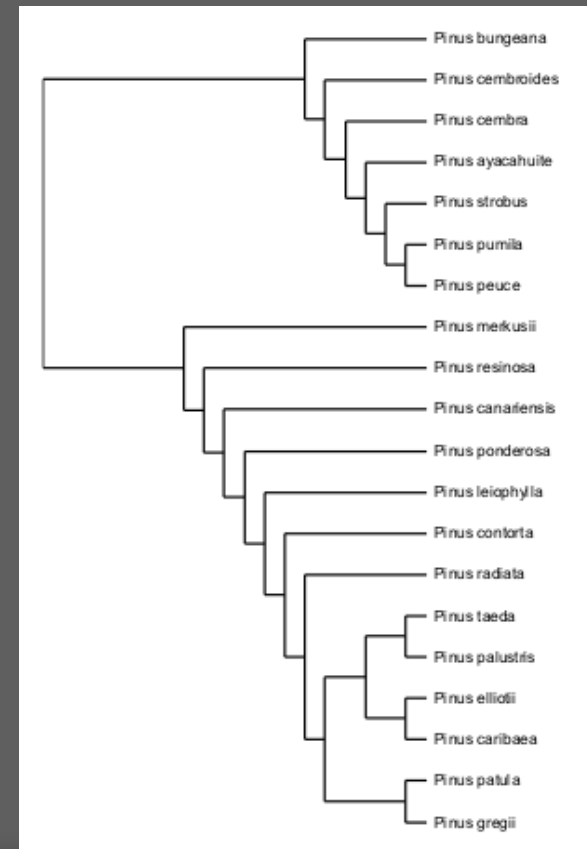
Motivations

⇒ Détection de gènes :

Génome nouveau



⇒ Phylogénie :



(Alignements multiples)

tfe1_mouse	~ERRMANNARERVRVDINEAFRELQRMQLHLKSDKAQTKLLILQAVQVILGLE
tfe2_human	~ERRVANNARERLVRVDINEAFKELQRMQLHLNSEKPKTKLLILHQAVSVILNLE
myod_human	ADRRKAATMRERRRLSKVNEAFETLKRCTSSNP..NQLRPKVEILRNAIKYIEGLQA
myod_mouse	ADRRKAATMRERRRLSKVNEAFETLKRCTSSNP..NQLRPKVEILRNAIKYIEGLQA
id1_human	~LYDMNGCYSRLKELVPTLPQNRKVSQVEILQHVIDYIRDQ
id1_mouse	~LYDMNGCYSRLKELVPTLPQNRKVSQVEILQHVIDYIRDQ
id1_rat	~LYDMNGCYSRLKELVPTLPQNRKVSQVEILQHVIDYIRDQ
ndf1_human	~RRMKANARERNRMHGLNAALDNLRKVVPCYSKTQKLSKIETLRLAKNYIWAQ
ndf1_mesau	~RRMKANARERNRMHGLNAALDNLRKVVPCYSKTQKLSKIETLRLAKNYIWAQ
ndf1_mouse	~RRMKANARERNRMHGLNAALDNLRKVVPCYSKTQKLSKIETLRLAKNYIWAQ
ndf2_human	~RRKANARERNRMHGLNAALDNLRKVVPCYSKTQKLSKIETLRLAKNYIWAQ
scl_human	~RRIFTNSRERWRQNVNGAFARLKLIPTHPPDKKLSKNEILRLAMKYINFL
scl_mouse	~RRIFTNSRERWRQNVNGAFARLKLIPTHPPDKKLSKNEILRLAMKYINFL
ly11_human	~RRVFTNSRERWRQNVNGAFARLKLIPTHPPDRKLSKNEVLRAMKYI
ly11_mouse	~RRVFTNSRERWRQNVNGAFARLKLIPTHPPDRKLSKNEVLRAMKYI
twst_human	~RVMANVRERRTSLNEAFAALRKIIPTLPSDKLSKIQTCLKLAARYIDFL
twst_mouse	~RVMANVRERRTSLNEAFAALRKIIPTLPSDKLSKIQTCLKLAARYIDFL
twst_xenla	~RVMANVRERRTSLNEAFSSLRKIIPTLPSDKLSKIQTCLKLASRYIDFL
max_chick	ADKRAHHNALERKRRDHIKDSFHSRLRDSVPSLGG.EKASRAQILDKATEYIYMR
max_human	ADKRAHHNALERKRRDHIKDSFHSRLRDSVPSLGG.EKASRAQILDKATEYIYMR
max_mouse	ADKRAHHNALERKRRDHIKDSFHSRLRDSVPSLGG.EKASRAQILDKATEYIYMR
max_rat	ADKRAHHNALERKRRDHIKDSFHSRLRDSVPSLGG.EKASRAQILDKATEYIYMR
max_xenla	ADKRAHHNALERKRRDHIKDSFHSRLRDSVPSLGG.EKASRAQILDKATEYIYMR
myc1_human	~KRKNHNFLERKRRNDLRSRFLALRDVPTLASCSPKPKVVILSKALEYLQAL
myc1_mouse	~KRKNHNFLERKRRNDLRSRFLALRDVPTLASCSPKPKVVILSKALEYLQAL
myc_human	~KRRTHNVLERRRMELKRSFFALRDQIPELENNEKAPKVVILKATAYILSVQA
myc_mouse	~DKRRTHNVLERRRMELKRSFFALRDQIPELENNEKAPKVVILKATAYILSVQA
tfe3_human	~KKDNHNLIERRRRNFINDRIKELQTLIPKSSDPEMRWNNKGTILKASVDYIRKL
tfe3_mouse	~KKDNHNLIERRRRNFINDRIKELQTLIPKSSDPEMRWNNKGTILKASVDYIRKL
sre1_human	~EKRTAHNAIEKRYRSSINDKIIELKDLVVGTT..EAKLNKSAVLRKAIDYIRFL
sre2_human	~ERRTTHNIIIEKRYRSSINDKIIELKDLVMGTT..DAKNHKSQVLRKAIDYIKYL

Distance d'édition

Deux séquences $v = v_1v_2\dots v_n$ et $w = w_1w_2\dots w_m$

Opérations d'édition :

- $\text{ins}(x,i)$
- $\text{suppr}(x,i)$
- $\text{subs}(x,y,i)$

CHAT - $\text{suppr}(C,1) \rightarrow$ **HAT** - $\text{subs}(H,R,1) \rightarrow$ **RAT**

Distance d'édition

- Chaque modification a un poids, dépendant de l'opération et des lettres en cause.
- **Distance d'édition** entre v et w : poids minimal d'une suite d'opérations permettant de transformer v en w .

CHAT - $\text{suppr}(C,1) \rightarrow$ HAT - $\text{subs}(H,R,1) \rightarrow$ RAT

Alignement

X : alphabet des séquences

$$X' = X \cup \{-\}$$

$$X'' = X' \setminus \{-,-\}$$

2 morphismes de monoïdes :

$$F_1: X''^* \rightarrow X^*$$

$$(x,y) \rightarrow x \text{ si } x \neq - \\ \varepsilon \text{ sinon}$$

$$F_2: X''^* \rightarrow X^*$$

$$(x,y) \rightarrow y \text{ si } x \neq - \\ \varepsilon \text{ sinon}$$

Un alignement de $v = v_1v_2\dots v_n$ et $w = w_1w_2\dots w_m$ est un mot A de X'' tel que $F_1(A) = v$ et $F_2(A) = w$.

$v = \text{aatca}$

$w = \text{agca}$

aa-tca

agc-a-

OK

aatc-a

-agc-a

~~OK~~

aatca

-agca

OK

Score d'alignement

- ⇒ Le score d'un alignement est la somme des scores lettre à lettre.
- ⇒ On cherche l'alignement de score maximal.

Scoring matrix based on large set of distantly related blocks: **Blosum62**

10	1	6	6	4	8	2	6	6	9	2	4	4	4	5	6	6	7	1	3	%
A	C	D	E	F	G	H	I	K	L	M	N	P	Q	R	S	T	V	W	Y	
8	0	-4	-2	-4	0	-4	-2	-2	-2	-2	-4	-2	-2	-2	2	0	0	-6	-4	A
	18	-6	-8	-4	-6	-6	-2	-6	-2	-2	-6	-6	-6	-6	-2	-2	-2	-4	-4	C
		12	4	-6	-2	-2	-6	-2	-8	-6	2	-2	0	-4	0	-2	-6	-8	-6	D
			10	-6	-4	0	-6	2	-6	-4	0	-2	4	0	0	-2	-4	-6	-4	E
				12	-6	-2	0	-6	0	-6	-8	-6	-6	-4	-4	-2	2	6	6	F
					12	-4	-8	-4	-8	-6	0	-4	-4	-4	0	-4	-6	-4	-6	G
						16	-6	-2	-6	-4	2	-4	0	0	-2	-4	-6	-4	4	H
							8	-6	4	2	-6	-6	-6	-6	-4	-2	6	-6	-2	I
								10	-4	-2	0	-2	2	4	0	-2	-4	-6	-4	K
									8	4	-6	-6	-4	-4	-4	-2	2	-4	-2	L
										10	-4	-4	0	-2	-2	2	-2	-2	-2	M
											12	-4	0	0	2	0	-6	-8	-4	N
												14	-2	-4	-2	-2	-4	-8	-6	P
													10	2	0	-2	-4	-4	-2	Q
														10	-2	-2	-6	-6	-4	R
															8	2	-4	-6	-4	S
																10	0	-4	-4	T
																	8	-6	-2	V
																		22	4	W
																			14	Y

Score d'alignement

A C G T

A 1 0 0 0

C 0 1 0 0

G 0 0 1 0

T 0 0 0 1

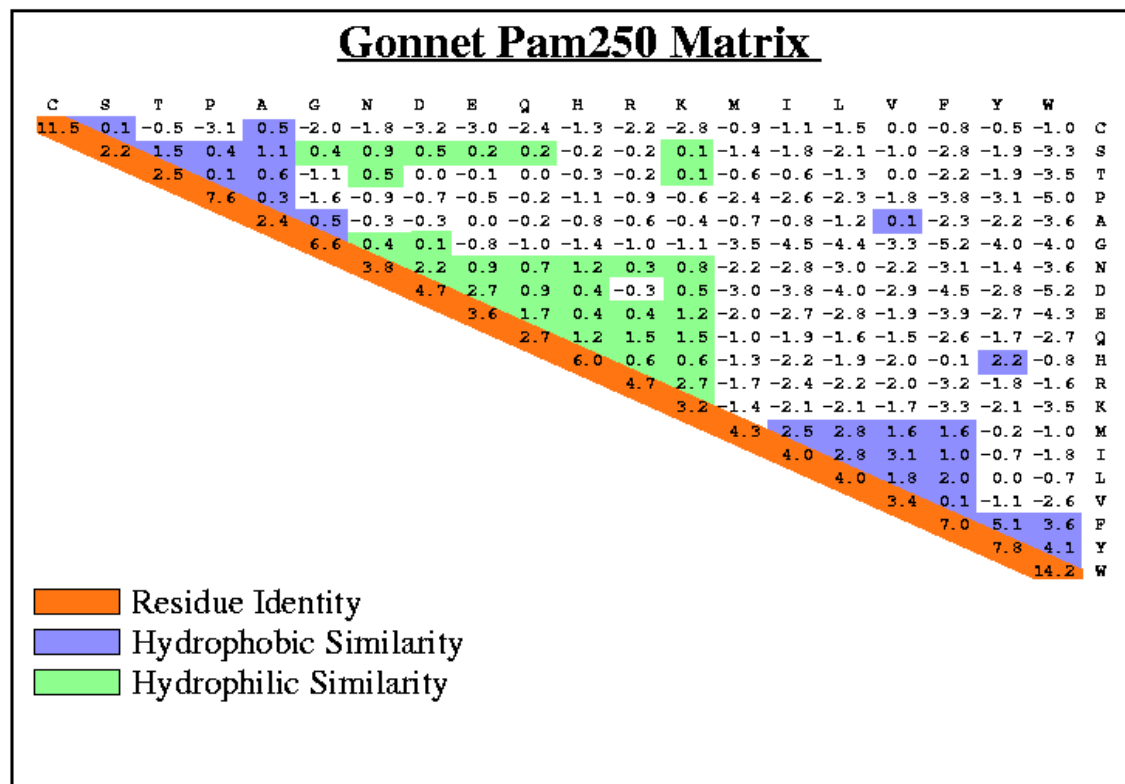
A C G T

A 3 0 1 0

C 0 3 0 1

G 1 0 3 0

T 0 0 1 3



Un alignement de deux protéines

seq 1: 3 EHYYSSEKPSVKSNNKQTSFRLRNKDFTFSTDSGVFSKKEVDFGSRLLLIDSFEEPEVEGGI 62
| | | | | . | . | | | . | | | | | | | . | | | | | . |

seq 2: 2 GHYYSREPNVPLKTKQIDVCIRGYCFKFKFITASGVFSFGKLDRGTELLIENM-ILKPDWKI 60

seq 1: 63 **LDVGCYGF**IGLSLASDFKDRTIHMIDVNERAVELSNENAEQNGITNVKIYQSDLFSNVD 122
| | | | | | | | | | . | | | | | | | . | | . | . . | . | .

seq 2: 61 **LDLGCYGV**IGI-VASRFVNYVV-MTDINKRAVQIARKNIKINGVKNAEVRLGNLYEPVE 118

seq 1: 123 SAQTFASILT**NPP**IRAGKKVVHAI FEKSAEHLKASGELWIVIQKKQGAPSAIEKLEELFD 182
. | | | | | | | | | . . | | | | | | | | | . . . |

seq 2: 119 -GEKFSIIT**NPP**VHAGKDILREIVINAPNYLHDGGMLQLVIKTKLGAKFIKDLMKDTFT 177

seq 1: 183 EVSVVQKKKGY 193
| | | | |

seq 2: 178 EVVELAKGSY 188

Equivalence édition - alignement

AACTA-CGAT

A-GTACCGTT

suppr(A,2)

subs(C,G,2)

ins(C,5)

subs(A,T,8)

Algorithme de Needleman et Wunsch (1970), Gotoh (1982)

$$V = V_1 V_2 \dots V_n$$

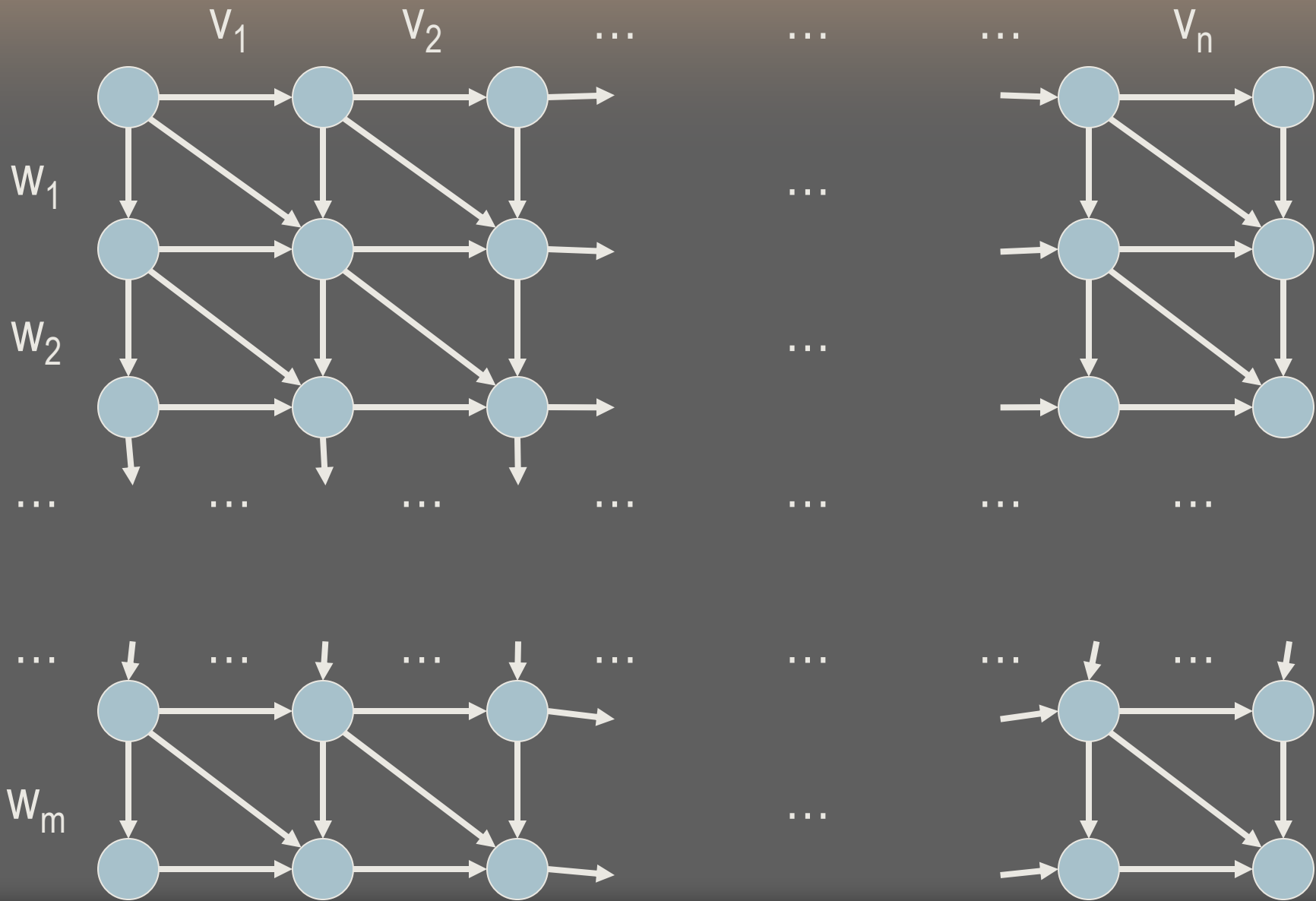
$$W = W_1 W_2 \dots W_m$$

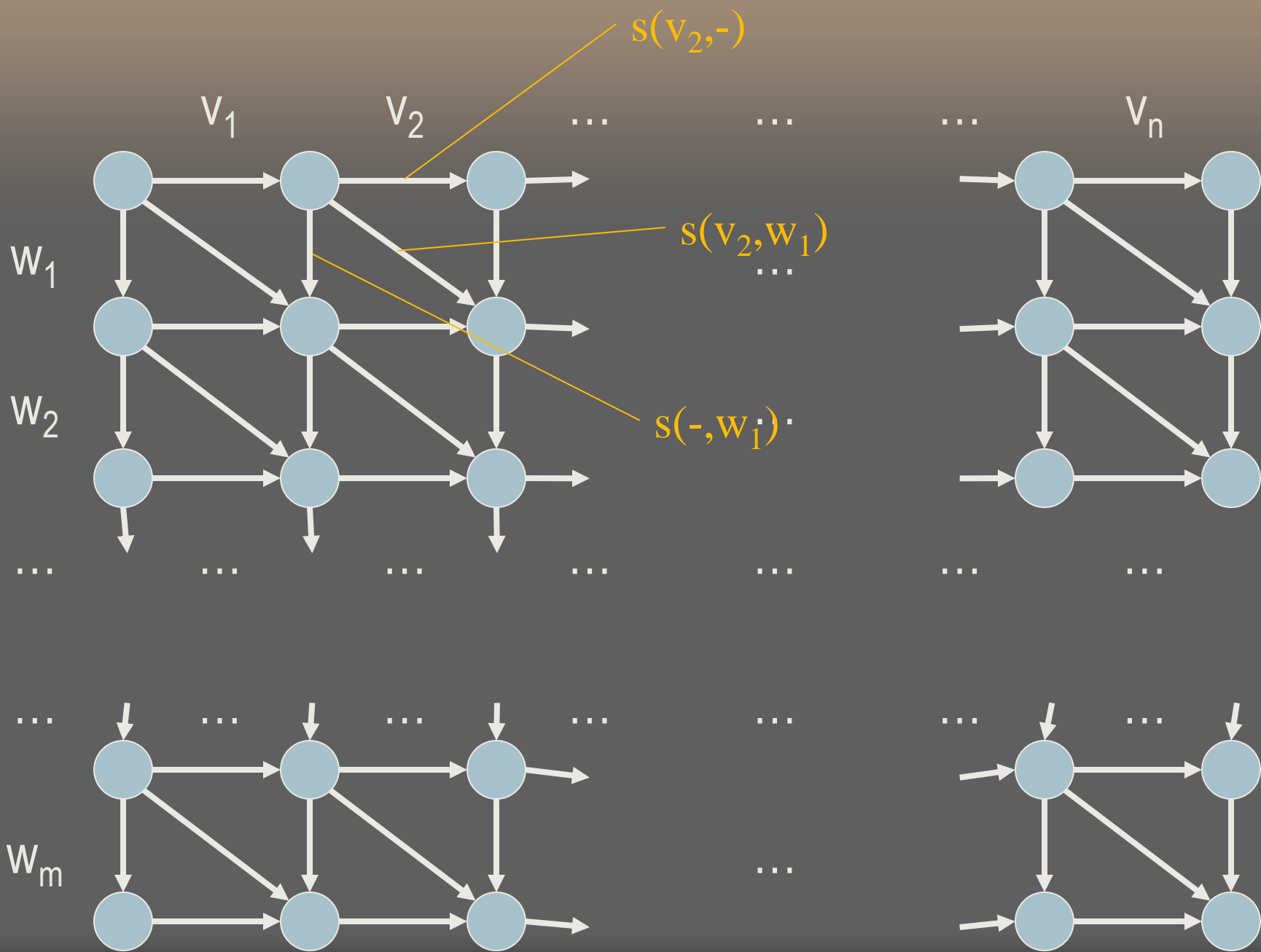
$s(x,y)$: score de substitution de x en y

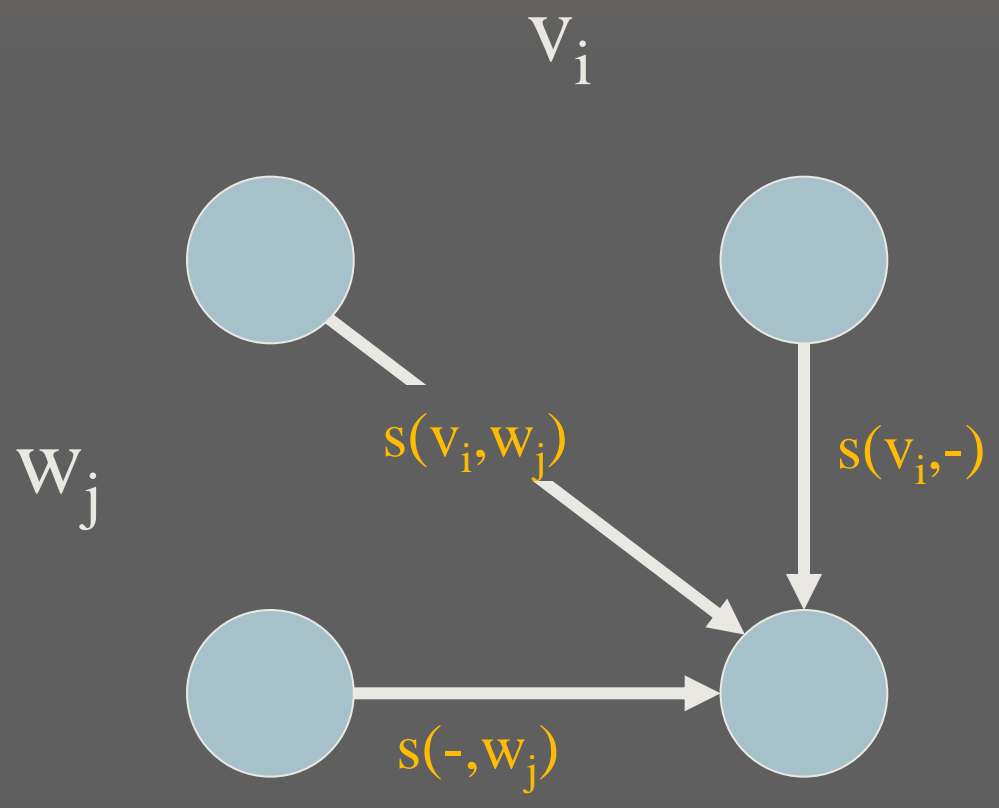
$s(x,-)$: score de suppression de x

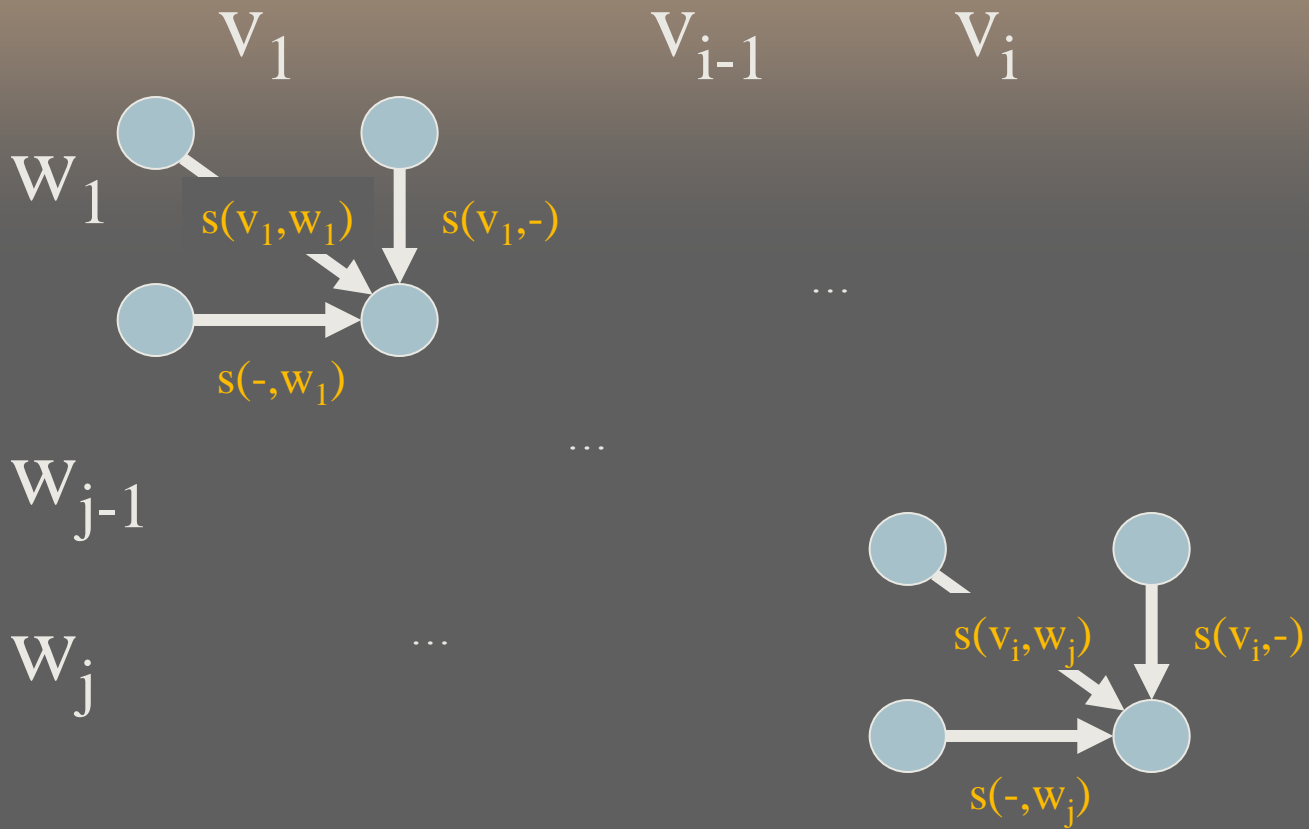
$s(-,y)$: score d'insertion de y

$S(v,w)$: score d'alignement de v et w









$$S(v_1 \dots v_i, w_1 \dots w_j) = \text{Max} \{$$

$$S(v_1 \dots v_{i-1}, w_1 \dots w_{j-1}) + s(v_i, w_j)$$

$$S(v_1 \dots v_{i-1}, w_1 \dots w_j) + s(v_i, -)$$

$$S(v_1 \dots v_i, w_1 \dots w_{j-1}) + s(-, w_j)$$

$$\}$$

W I N D O X S

L
I
N
E
S

0	-1	-2	-3	-4	-5	-6	-7
-1							
-2							
-3							
-4							
-5							

$$s(x,y) = 2 \text{ si } x=y \\ -1 \text{ sinon}$$

$$s(x,-) = s(-,x) = -1$$

W I N D O X S

L	0	-1	-2	-3	-4	-5	-6	-7
I	-1	-1						
N	-2							
U	-3							
X	-4							
	-5							

$$s(x,y) = 2 \text{ si } x=y \\ -1 \text{ sinon}$$

$$s(x,-) = s(-,x) = -1$$

W I N D O X S

L	0	-1	-2	-3	-4	-5	-6	-7
I	-1	-1	-2					
N	-2							
U	-3							
X	-4							
	-5							

$$s(x,y) = 2 \text{ si } x=y \\ -1 \text{ sinon}$$

$$s(x,-) = s(-,x) = -1$$

W I N D O X S

L
I
N
E
S

0	-1	-2	-3	-4	-5	-6	-7
-1	-1	-2	-3	-4	-5	-6	-7
-2	-2						
-3							
-4							
-5							

$$s(x,y) = 2 \text{ si } x=y$$

$$-1 \text{ sinon}$$

$$s(x,-) = s(-,x) = -1$$

W I N D O X S

L
I
N
E
S

0	-1	-2	-3	-4	-5	-6	-7
-1	-1	-2	-3	-4	-5	-6	-7
-2	-2	1					
-3							
-4							
-5							

$$s(x,y) = \begin{cases} 2 & \text{si } x=y \\ -1 & \text{sinon} \end{cases}$$

$$s(x,-) = s(-,x) = -1$$

W I N D O X S

L
I
N
E
S

0	-1	-2	-3	-4	-5	-6	-7
-1	-1	-2	-3	-4	-5	-6	-7
-2	-2	1 → 0					
-3							
-4							
-5							

$$s(x,y) = 2 \text{ si } x=y$$

$$-1 \text{ sinon}$$

$$s(x,-) = s(-,x) = -1$$

W I N D O X S

L
I
N
E
S

0	-1	-2	-3	-4	-5	-6	-7
-1	-1	-2	-3	-4	-5	-6	-7
-2	-2	1	0	-1	-2	-3	-4
-3	-3	0	3	2	1	0	-1
-4	-4	-1	2	2	1	0	-1
-5	-5	-2	1	1	1	3	2

$$s(x,y) = 2 \text{ si } x=y$$

$$-1 \text{ sinon}$$

$$s(x,-) = s(-,x) = -1$$

W I N D O X S

L
I
N
E
S

0	-1	-2	-3	-4	-5	-6	-7
-1	-1	-2	-3	-4	-5	-6	-7
-2	-2	1	0	-1	-2	-3	-4
-3	-3	0	3	2	1	0	-1
-4	-4	-1	2	2	1	0	-1
-5	-5	-2	1	1	1	3	2

$$s(x,y) = 2 \text{ si } x=y$$

$$-1 \text{ sinon}$$

$$s(x,-) = s(-,x) = -1$$

W I N D O X S

L
I
N
E
S

0	-1	-2	-3	-4	-5	-6	-7
-1	-1	-2	-3	-4	-5	-6	-7
-2	-2	1	0	-1	-2	-3	-4
-3	-3	0	3	2	1	0	-1
-4	-4	-1	2	2	1	0	-1
-5	-5	-2	1	1	1	3	2

$$s(x,y) = 2 \text{ si } x=y$$

$$-1 \text{ sinon}$$

$$s(x,-) = s(-,x) = -1$$

W I N D O X S

L
I
N
E
S

0	-1	-2	-3	-4	-5	-6	-7
-1	-1	-2	-3	-4	-5	-6	-7
-2	-2	1	0	-1	-2	-3	-4
-3	-3	0	3	2	1	0	-1
-4	-4	-1	2	2	1	0	-1
-5	-5	-2	1	1	1	3	2

$$s(x,y) = 2 \text{ si } x=y$$

$$-1 \text{ sinon}$$

$$s(x,-) = s(-,x) = -1$$

W I N D O X S

L
I
N
E
S

0	-1	-2	-3	-4	-5	-6	-7
-1	-1	-2	-3	-4	-5	-6	-7
-2	-2	1	0	-1	-2	-3	-4
-3	-3	0	3	2	1	0	-1
-4	-4	-1	2	2	1	0	-1
-5	-5	-2	1	1	1	3	2

$$s(x,y) = 2 \text{ si } x=y$$

$$-1 \text{ sinon}$$

$$s(x,-) = s(-,x) = -1$$

W I N D O X S

L
I
N
U
X

0	-1	-2	-3	-4	-5	-6	-7
-1	-1	-2	-3	-4	-5	-6	-7
-2	-2	1	0	-1	-2	-3	-4
-3	-3	0	3	2	1	0	-1
-4	-4	-1	2	2	1	0	-1
-5	-5	-2	1	1	1	3	2

$$s(x,y) = 2 \text{ si } x=y$$

$$-1 \text{ sinon}$$

$$s(x,-) = s(-,x) = -1$$

W I N D O X S

L
I
N
U
X

0	-1	-2	-3	-4	-5	-6	-7
-1	-1	-2	-3	-4	-5	-6	-7
-2	-2	1	0	-1	-2	-3	-4
-3	-3	0	3	2	1	0	-1
-4	-4	-1	2	2	1	0	-1
-5	-5	-2	1	1	1	3	2

$$s(x,y) = 2 \text{ si } x=y$$

$$-1 \text{ sinon}$$

$$s(x,-) = s(-,x) = -1$$

W I N D O X S

L
I
N
E
S

0	-1	-2	-3	-4	-5	-6	-7
-1	-1	-2	-3	-4	-5	-6	-7
-2	-2	1	0	-1	-2	-3	-4
-3	-3	0	3	2	1	0	-1
-4	-4	-1	2	2	1	0	-1
-5	-5	-2	1	1	1	3	2

$$s(x,y) = 2 \text{ si } x=y$$

$$-1 \text{ sinon}$$

$$s(x,-) = s(-,x) = -1$$

W I N D O X S

L
I
N
U
X

0	-1	-2	-3	-4	-5	-6	-7
-1	-1	-2	-3	-4	-5	-6	-7
-2	-2	1	0	-1	-2	-3	-4
-3	-3	0	3	2	1	0	-1
-4	-4	-1	2	2	1	0	-1
-5	-5	-2	1	1	1	3	2

W I N D O X S
L I N - U X -

W I N D O X S

L
I
N
U
X

0	-1	-2	-3	-4	-5	-6	-7
-1	-1	-2	-3	-4	-5	-6	-7
-2	-2	1	0	-1	-2	-3	-4
-3	-3	0	3	2	1	0	-1
-4	-4	-1	2	2	1	0	-1
-5	-5	-2	1	1	1	3	2

W I N D O X S

W I N D O X S

L I N - U X -

L I N U - X -

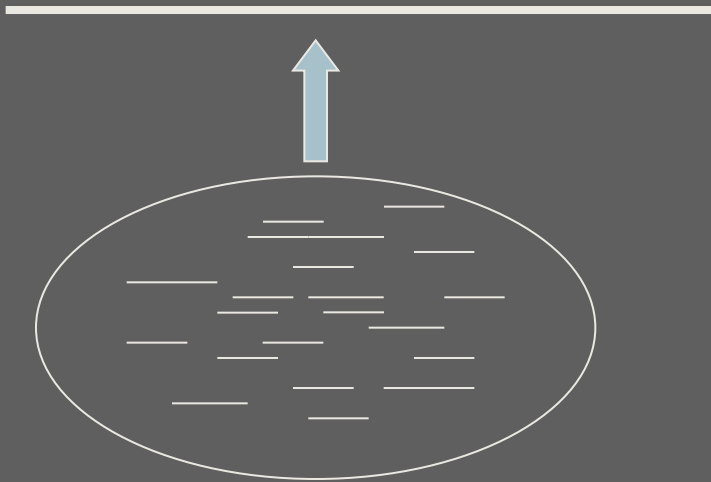
Complexité

En espace : $m \quad n$

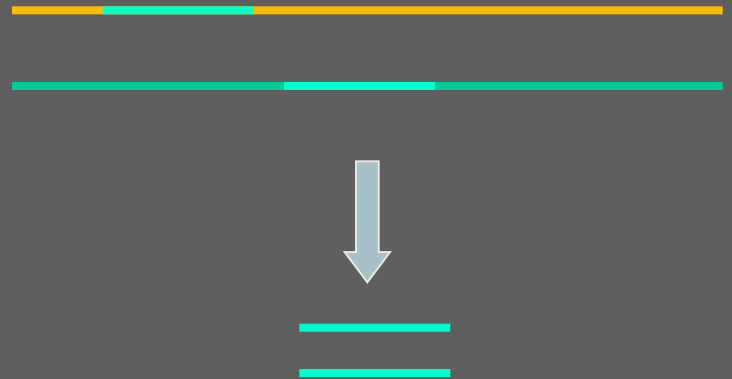
En temps : $m \quad n$

Alignement local (Smith et Waterman 1981)

Aligner des petites séquences
sur des grandes :



Trouver des similarités locales :



$s(v_2, -)$

V_1

V_2

...

...

...

V_n

W_1

$s(v_2, w_1)$

...

W_2

$s(-, w_1)$

...

...

...

...

...

...

...

...

...

...

...

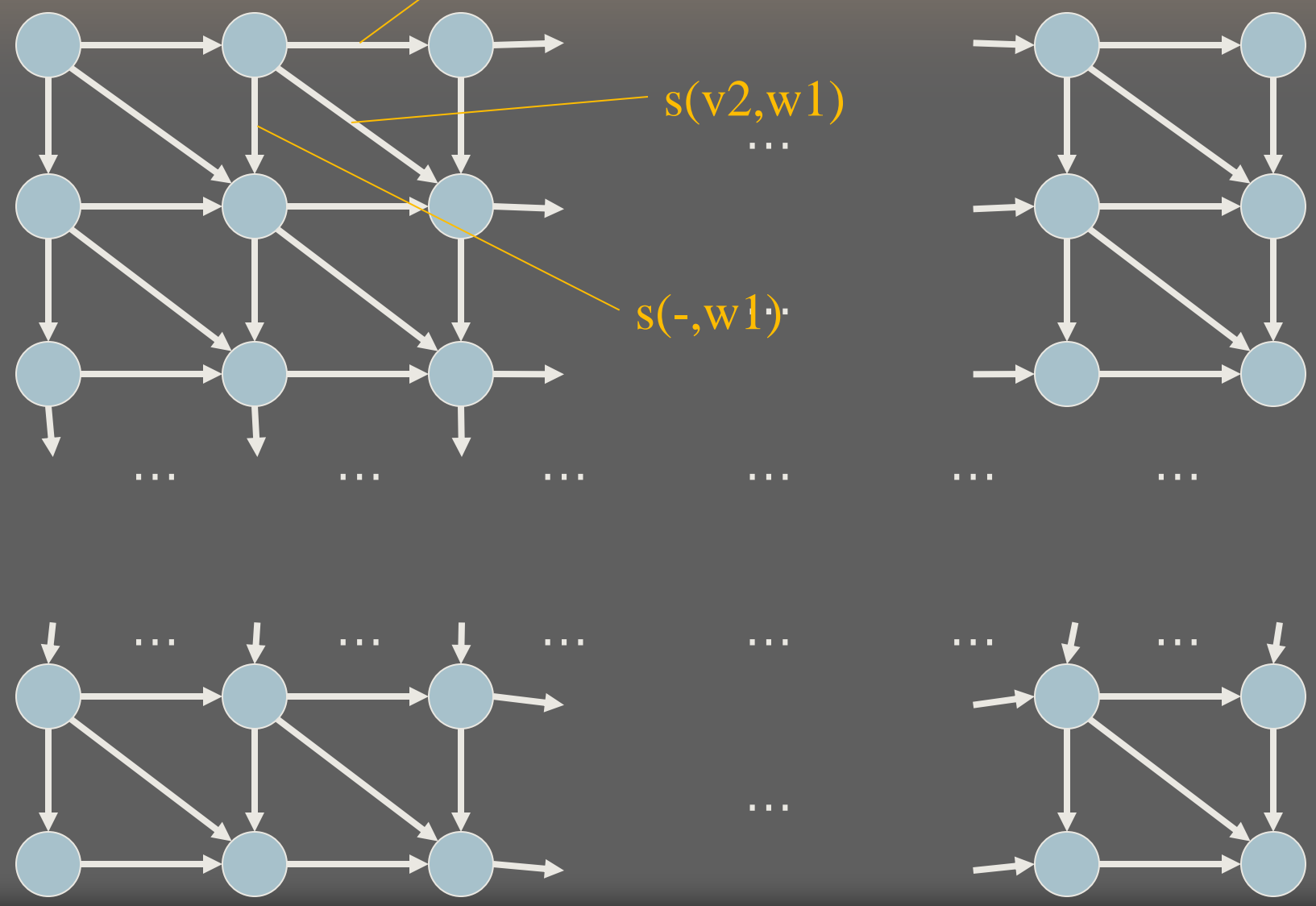
...

...

...

W_m

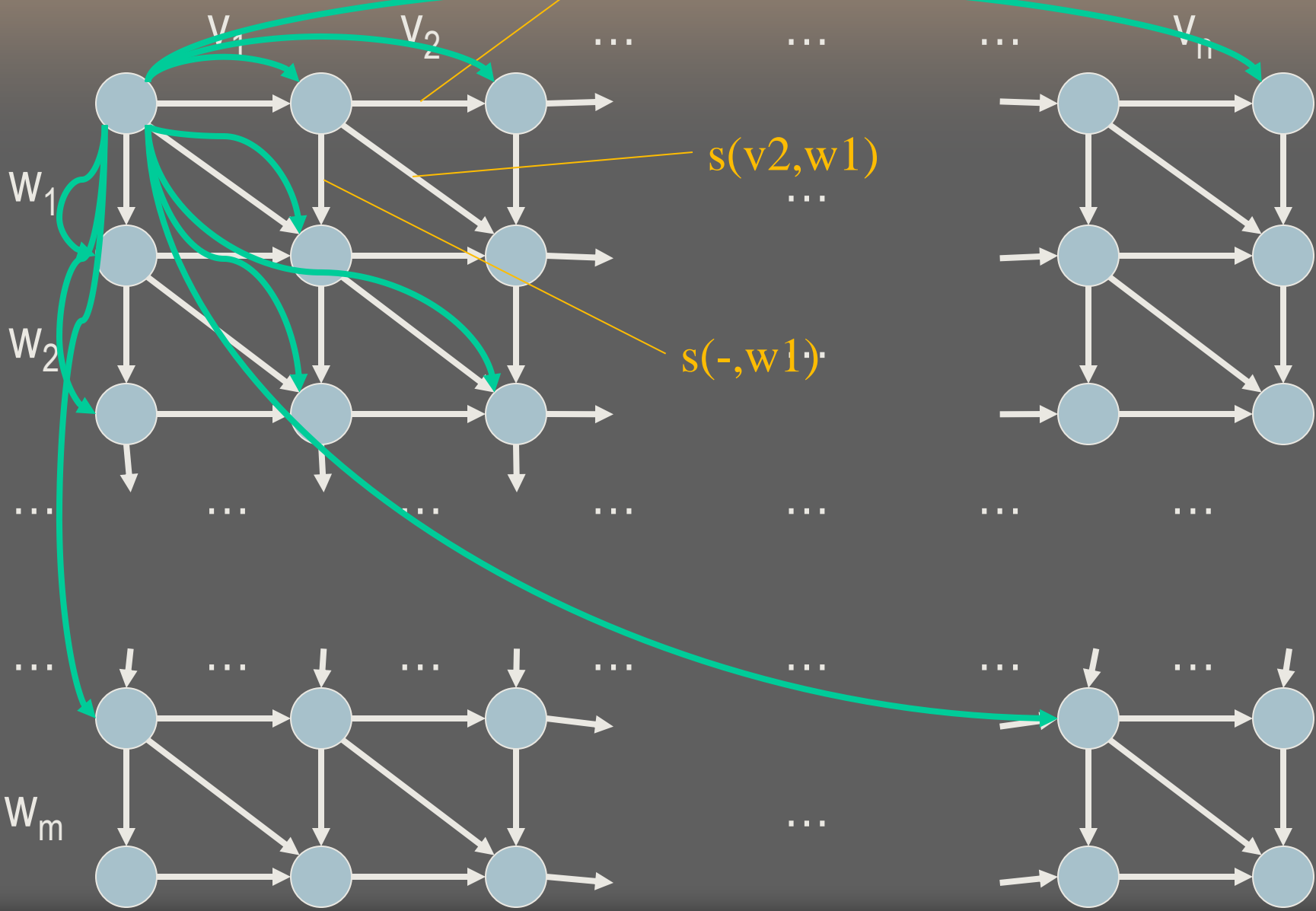
...



$s(v2,-)$

$s(v2,w1)$

$s(-,w1)$

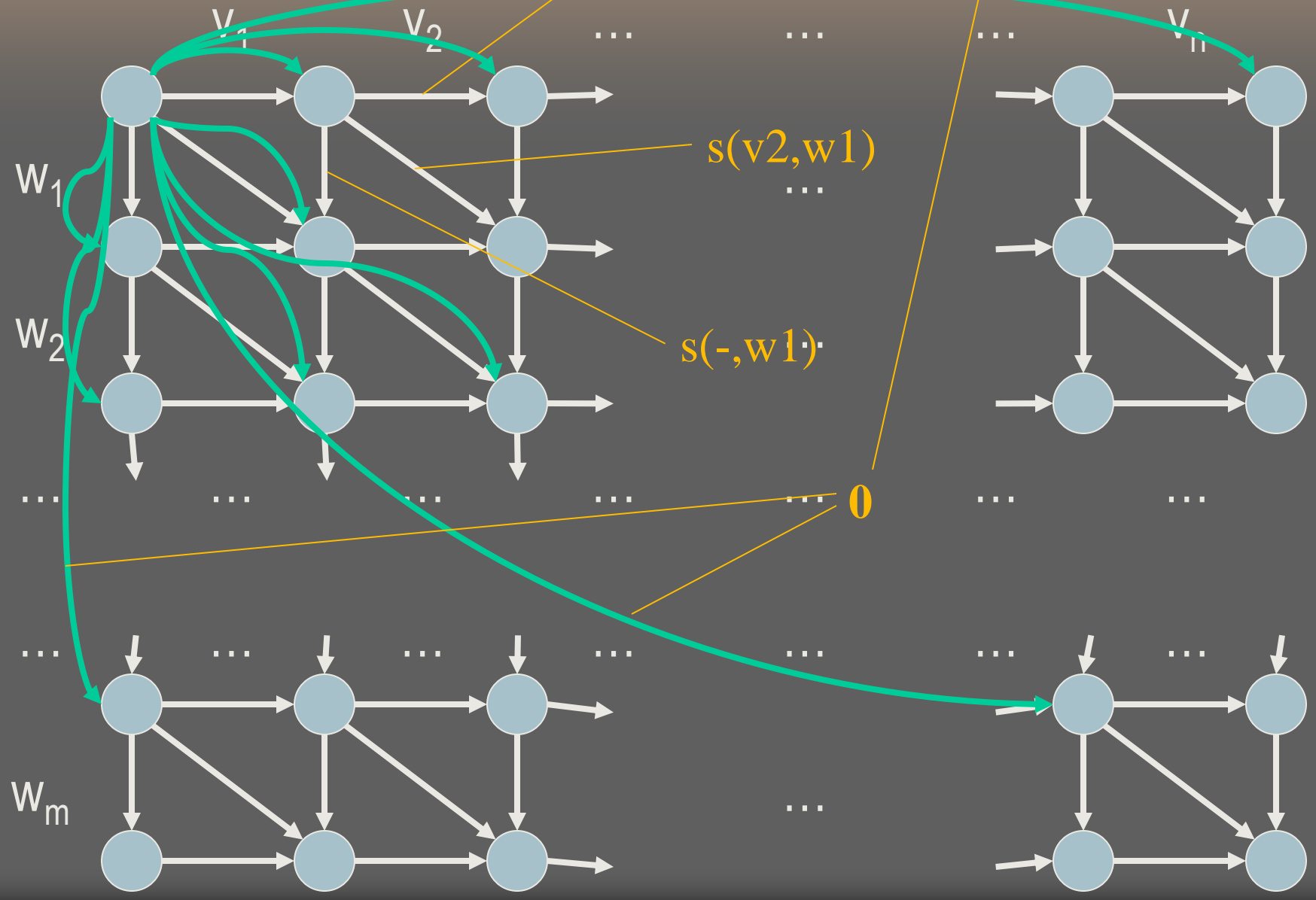


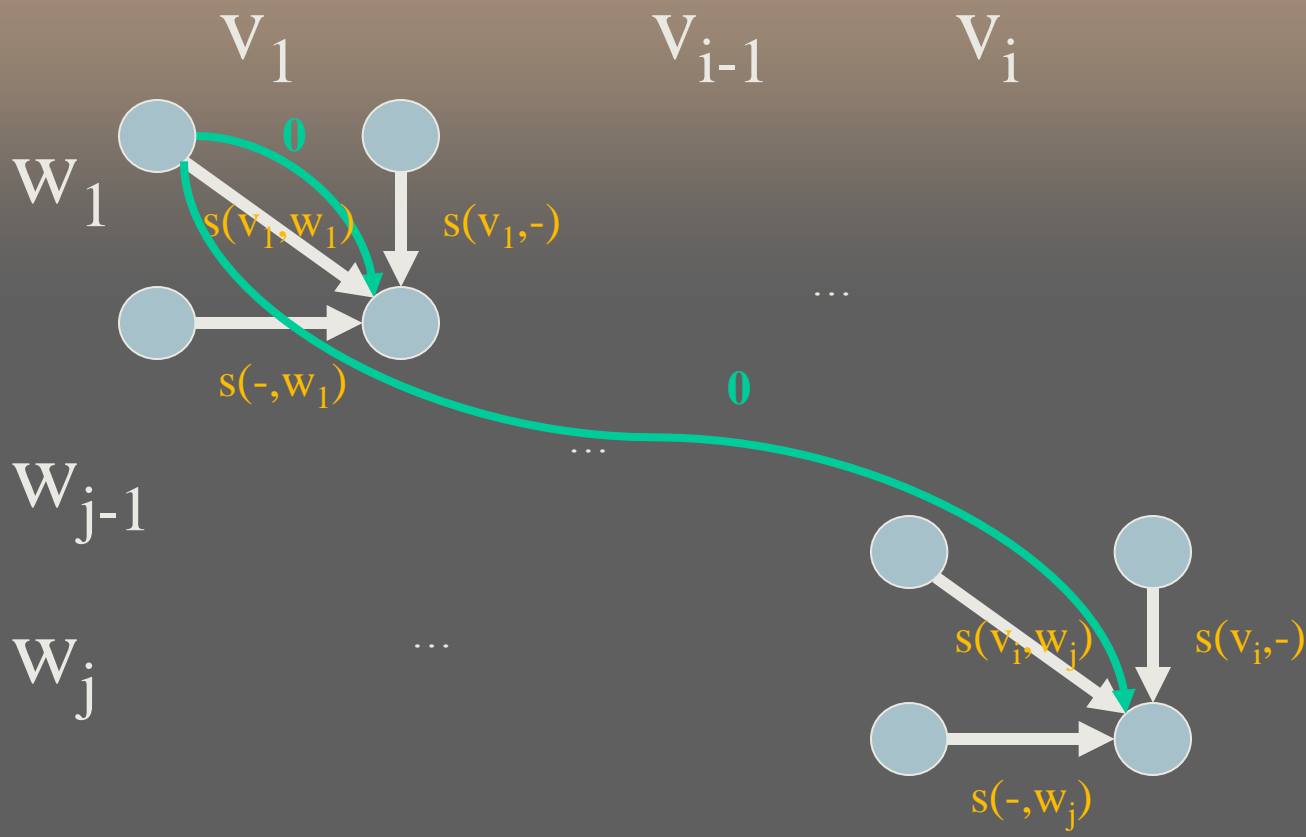
$s(v2,-)$

$s(v2,w1)$

$s(-,w1)$

θ





$$S(v_1 \dots v_i, w_1 \dots w_j) = \text{Max} \{$$

$$\begin{aligned}
 & 0 \\
 & S(v_1 \dots v_{i-1}, w_1 \dots w_{j-1}) + s(v_i, w_j) \\
 & S(v_1 \dots v_{i-1}, w_1 \dots w_j) + s(v_i, -) \\
 & S(v_1 \dots v_i, w_1 \dots w_{j-1}) + s(-, w_j)
 \end{aligned}$$

$$\}$$

W I N D O X S

L
I
N
E
S

0	0	0	0	0	0	0	0
0							
0							
0							
0							
0							

$$s(x,y) = 2 \text{ si } x=y \\ -1 \text{ sinon}$$

$$s(x,-) = s(-,x) = -1$$

W I N D O X S

L	0	0	0	0	0	0	0
I	0	0					
N	0						
U	0						
X	0						

$$s(x,y) = 2 \text{ si } x=y$$

$$-1 \text{ sinon}$$

$$s(x,-) = s(-,x) = -1$$

W I N D O X S

L

I

N

U

X

0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0
0	0	2					
0							
0							
0							

$$s(x,y) = 2 \text{ si } x=y \\ -1 \text{ sinon}$$

$$s(x,-) = s(-,x) = -1$$

W I N D O X S

L
I
N
E
U
X

0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0
0	0	2	1	0	0	0	0
0	0	1	4	3	2	1	0
0	0	0	3	3	2	1	0
0	0	0	2	2	2	4	3

$$s(x,y) = 2 \text{ si } x=y$$

$$-1 \text{ sinon}$$

$$s(x,-) = s(-,x) = -1$$

W I N D O X S

L
I
N
E
U
X

0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0
0	0	2	1	0	0	0	0
0	0	1	4	3	2	1	0
0	0	0	3	3	2	1	0
0	0	0	2	2	2	4	3

$$s(x,y) = 2 \text{ si } x=y$$

$$-1 \text{ sinon}$$

$$s(x,-) = s(-,x) = -1$$

W I N D O X S

L
I
N
E
S

0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0
0	0	2	1	0	0	0	0
0	0	1	4	3	2	1	0
0	0	0	3	3	2	1	0
0	0	0	2	2	2	4	3

$$s(x,y) = 2 \text{ si } x=y$$

$$-1 \text{ sinon}$$

$$s(x,-) = s(-,x) = -1$$

W I N D O X S

L
I
N
U
X

0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0
0	0	2	1	0	0	0	0
0	0	1	4	3	2	1	0
0	0	0	3	3	2	1	0
0	0	0	2	2	2	4	3

I N
I N

I N D O X
I N - U X

I N D O X
I N U - X