



**felipe@rali** Dernière modification January 25, 2015

# Introduction à l'apprentissage analogique

**Philippe Langlais**

`felipe@iro.umontreal.ca`

# Plan

## Apprentissage Analogique

Analogie

Principe

## Quelques réalisations

### Sous le capot

Définitions de l'analogie (formelle)

Solveurs d'équations

Recherche des analogies

Filtrer le bruit

### Applications

Traduction de mots inconnus

Traduction de termes

Translittération

## Apprentissage Analogique

Analogie

Principe

### Quelques réalisations

#### Sous le capot

Définitions de l'analogie (formelle)

Solveurs d'équations

Recherche des analogies

Filtrer le bruit

#### Applications

Traduction de mots inconnus

Traduction de termes

Translittération

# Analogie

$[x : y :: z : t]$

*“x est à y ce que z est à t”*

# Analogies sémantiques

[Christian : church :: Muslim : ?]

# Analogies sémantiques

[Christian : church :: Muslim : mosque]

# Analogies sémantiques

[Christian : church :: Muslim : mosque]

[3 : 4 :: 6 : ?]

# Analogies sémantiques

[Christian : church :: Muslim : mosque]

[3 : 4 :: 6 : 8]



# Analogies sémantiques

[Christian : church :: Muslim : mosque]

[3 : 4 :: 6 : 8]

[Paris : Vélib :: Montréal : ?]

# Analogies sémantiques

[Christian : church :: Muslim : mosque]

[3 : 4 :: 6 : 8]

[Paris : Vélib :: Montréal : Bixi]

# Analogies sémantiques

[Christian : church :: Muslim : mosque]

[3 : 4 :: 6 : 8]

[Paris : Vélib :: Montréal : Bixi]



# Analogies sémantiques

[Christian : church :: Muslim : mosque]

[3 : 4 :: 6 : 8]

[Paris : Vélib :: Montréal : Bixi]



# Analogies formelles

Exemples pris de (Lepage, 1998)

honor	:	hon	ōrem	::	ōrātor	:	ōrātōrem
reader	:	un	readable	::	doer	:	undoable
lang	:	läng	ste	::	scharf	:	schärfste
répression	:	répression	naire	::	réaction	:	réactionnaire
croys	:	cré	ons	::	mont royal	:	montréal
keras	:	menger	askan	::	kena	:	mengenakan

obj(big)&            :    obj(small)⊂obj(big) ::    obj(big)&            :    ?  
obj=circle            & obj=circle            obj=square

? = obj(small)⊂obj(big) & obj=square

# Apprentissage analogique

Principe d'après (Pirrelli & Yvon, 1999)

- ▶ un ensemble d'observations  $\mathcal{L} = \{\langle s, t \rangle\}$
- ▶ une observation incomplète,  $u = \langle s, ? \rangle$

**Problem:** Prédire les traits manquants de  $u$

- 1. Construire  $\mathcal{E}_{\mathcal{I}}(u)$ :**  
 $\{(x, y, z) \in \mathcal{L}^3 \mid [I(x) : I(y) :: I(z) : I(u)]\}$
- 2. Construire  $\mathcal{E}_{\mathcal{O}}(u)$ :**  
 $\{t \mid [O(x) : O(y) :: O(z) : t], \forall (x, y, z) \in \mathcal{E}_{\mathcal{I}}(u)\}$
- 3. Sélectionner les candidats de  $\mathcal{E}_{\mathcal{O}}(u)$ .**

where  $I(\langle s, t \rangle) \equiv s$  and  $O(\langle s, t \rangle) \equiv t$ .

## Illustration: parsing

$$\mathcal{L} \equiv \left\{ \begin{array}{l} (\text{S}(\text{Pr}(\text{elle}).\text{Vb}(\text{parle})), \text{ she talks}), \\ (\text{S}(\text{Pr}(\text{elle}).\text{Vb}(\text{mange})), \text{ she eats}), \\ (\text{S}(\text{Pr}(\text{elles}).\text{Vb}(\text{mangent})), \text{ they eat}) \end{array} \right\}$$

## Illustration: parsing

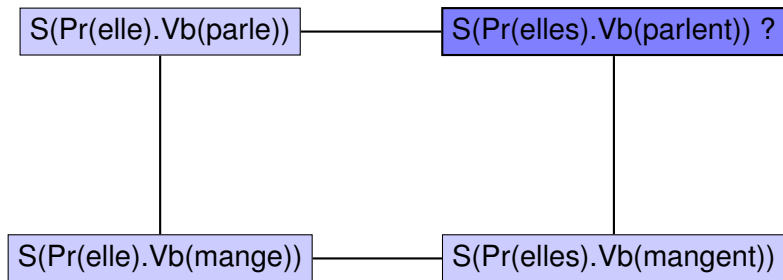
$$\mathcal{L} \equiv \left\{ \begin{array}{l} (\text{S}(\text{Pr}(\text{elle}).\text{Vb}(\text{parle})), \text{ she talks}), \\ (\text{S}(\text{Pr}(\text{elle}).\text{Vb}(\text{mange})), \text{ she eats}), \\ (\text{S}(\text{Pr}(\text{elles}).\text{Vb}(\text{mangent})), \text{ they eat}) \end{array} \right\}$$

S(Pr(elles).Vb(parlent)) ?



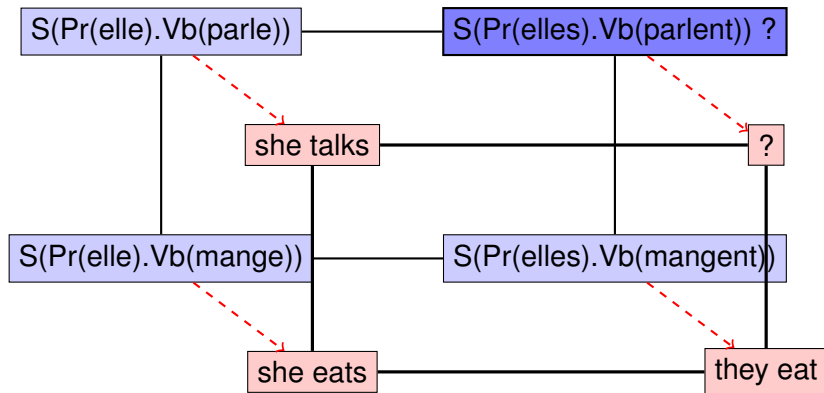
# Illustration: parsing

$$\mathcal{L} \equiv \left\{ \begin{array}{l} (\text{S}(\text{Pr}(\text{elle}).\text{Vb}(\text{parle})), \text{ she talks}), \\ (\text{S}(\text{Pr}(\text{elle}).\text{Vb}(\text{mange})), \text{ she eats}), \\ (\text{S}(\text{Pr}(\text{elles}).\text{Vb}(\text{mangent})), \text{ they eat}) \end{array} \right\}$$



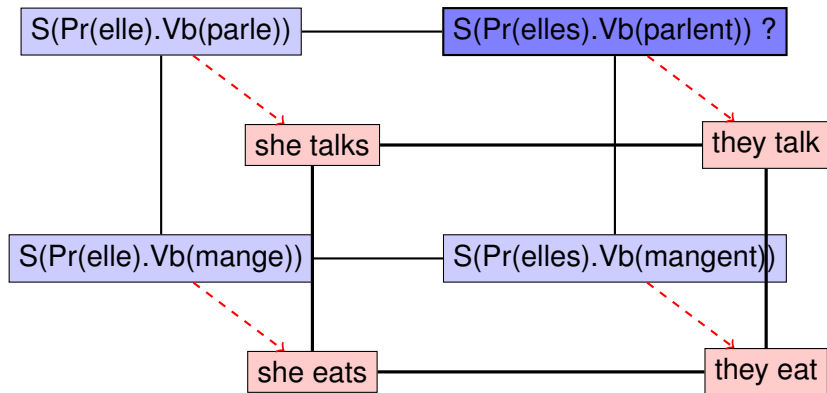
# Illustration: parsing

$$\mathcal{L} \equiv \left\{ \begin{array}{l} (\text{S}(\text{Pr}(\text{elle}).\text{Vb}(\text{parle})), \text{she talks}), \\ (\text{S}(\text{Pr}(\text{elle}).\text{Vb}(\text{mange})), \text{she eats}), \\ (\text{S}(\text{Pr}(\text{elles}).\text{Vb}(\text{mangent})), \text{they eat}) \end{array} \right\}$$



# Illustration: parsing

$$\mathcal{L} \equiv \left\{ \begin{array}{l} (S(\text{Pr}(\text{elle}).\text{Vb}(\text{parle})), \text{she talks}), \\ (S(\text{Pr}(\text{elle}).\text{Vb}(\text{mange})), \text{she eats}), \\ (S(\text{Pr}(\text{elles}).\text{Vb}(\text{mangent})), \text{they eat}) \end{array} \right\}$$



# Illustration: traduction de mots inconnus

Step-1: calcul de  $\mathcal{E}_J$  (futilité)

futilité?

# Illustration: traduction de mots inconnus

Step-1: calcul de  $\mathcal{E}_{\mathcal{J}}$  (futilité)

faillites

faillite

**futilité?**

futilités

[faillites : faillite :: futilités : **futilité**]

# Illustration: traduction de mots inconnus

Step-1: calcul de  $\mathcal{E}_f$  (futilité)

faillites

futile

faillite

**futilité?**

futilités

lucide

lucidité

[lucide : lucidité :: futile : **futilité**]

# Illustration: traduction de mots inconnus

Step-1: calcul de  $\mathcal{E}_f$  (futilité)

faillites

futile

faillite

hostilités

**futilité?**

hostiles

futilités

lucide

lucidité

[hostiles : hostilités :: futile : **futilité**]

# Illustration: traduction de mots inconnus

Step-1: calcul de  $\mathcal{E}_J$  (futilité)

faillites

réalités

faillite

futilités

futilités

lucide

réalité

**futilité?**

lucidité

futile

hostilités

hostiles

[réalités : réalité :: futilités : **futilité**]



# Illustration: traduction de mots inconnus

Step-1: calcul de  $\mathcal{E}_J$  (futilité)

faillites	brutale		futile
réalités	brutalité	réalité	mutilé
faillite		bestialité	hostilités
futilités	<b>futilité?</b>		natalité
facile		facilités	hostiles
futilités	bestiale	timide	natale
lucide	faillite		maille
	timidité	lucidité	

2944 équations considérées, 84 valides

# Illustration: traduction de mots inconnus

## Step-1: calcul de $\mathcal{E}_{\mathcal{J}}$ (futilité)

⟨activités, activité, futilités⟩ ⟨fatale, futile, fatalité⟩ ⟨timide, timidité, futile⟩  
⟨faciles, facilités, futile⟩ ⟨cupide, cupidité, futile⟩ ⟨utilisés, futilités, utilisé⟩  
⟨humide, humidité, futile⟩ ⟨brutalités, brutalité, futilités⟩ ⟨multipliés, multiplié, futilités⟩  
⟨qualités, qualité, futilités⟩ ⟨totale, totalité, futile⟩ ⟨utilisés, utilisé, futilités⟩  
⟨mutilés, mutilé, futilités⟩ ⟨félicités, futilités, félicité⟩ ⟨active, activité, futile⟩  
⟨mature, maturité, futile⟩ ⟨unités, futilités, unité⟩ ⟨habilités, habilité, futilités⟩  
⟨fragile, fragilité, futile⟩ ⟨subtile, subtilité, futile⟩ ⟨mutilés, futilités, mutilé⟩  
⟨autorités, autorité, futilités⟩ ⟨vitale, vitalité, futile⟩ ⟨autorisés, autorisé, futilités⟩  
⟨maille, faillite, mutilé⟩ ⟨mute, mutilé, futile⟩ ⟨facultés, faculté, futilités⟩  
⟨félicités, félicité, futilités⟩ ⟨utile, futile, utilité⟩ ⟨facilités, facilité, futilités⟩  
⟨rurale, ruralité, futile⟩ ⟨brutalités, futilités, brutalité⟩ ⟨finales, finalités, futile⟩  
⟨subtiles, subtilités, futile⟩ ⟨spatiale, spatialité, futile⟩ ⟨visités, visité, futilités⟩  
⟨réalités, réalité, futilités⟩ ⟨pénale, pénalité, futile⟩ ⟨brutales, brutalités, futile⟩  
⟨habilités, futilités, habilité⟩ ⟨hostilités, futilités, hostilité⟩

(+ 42 autres ...)

# Illustration: traduction de mots inconnus

## Step-2: calcul de $\mathcal{E}_O$ (futilité)

[faillites : faillite :: futilités : **futilité**]

faillites	↔	<b>bankruptcies</b> , bankruptcy
faillite	↔	collapse, bust, insolvency, ruin, complaining, bankrupt, business, <b>bankruptcy</b> , wall
futilités	↔	<b>trivialities</b>

[bankruptcies : bankruptcy :: trivialities : ?] { **triviality** }

# Illustration: traduction de mots inconnus

## Step-2: calcul de $E_O(\text{futilité})$

[lucide : lucidité :: futile : **futilité**]

lucide	↔	clear-thinking, blind, refocussed, lucid, far-sighted
lucidité	↔	well-informed, meticulously, moderation, penetrating, lucidity, clear-headed, clear-sighted
futilité	↔	meaningless, futile

[lucid : lucidity :: meaningless : ?]                    { meaninglessness }  
[lucid : lucidity :: futile : ?]                    { futiityle, futileity, futilyte }  
[far-sighted : clear-sighted :: futile : ?]                    { fucleile,  
cleutile, clefuile, cluteile, culetile, cfuleile, ... }

# Illustration: traduction de mots inconnus

## Step-3: sélection depuis $\mathcal{E}_O$ (futilité)

1124  $\neq$  formes générées, 120 au moins 2 fois

(trivialities,27)	(meaninglitiesss,3)	(meaniitnglyss,2)
(triviality,14)	(meaningiltyss,2)	(superfluous,2)
(futile,9)	(futiilty,2)	(meaningitilesss,2)
(futilitye,9)	(applicatitivialitis,2)	(meaningitlyss,2)
(meaningless,9)	(futiitiles,2)	(meaninglesity,2)
(meaninglityess,8)	(futiitly,2)	(meaninglessness,2)
(trivialitie,6)	(futiityl,2)	(meaninglityes,2)
(futility,4)	(futileity,2)	(meaniniglytyss,2)
(meaninglityss,4)	(futileessne,2)	(meaninitglyss,2)
(futilities,3)	(high-triviality,2)	(meaninitiglesss,2)
(meaninglesit,3)	(ltbbyirivialitig,2)	(mfaninglecilits,2)

...

# Illustration: traduction de mots inconnus

Step-3: sélection depuis  $\mathcal{E}_O(\text{futilité})$

1124  $\neq$  formes générées, 120 au moins 2 fois

(trivialities,27)	(meaninglitiesss,3)	(meaniitnglyss,2)
(triviality,14)	(meaningiltyss,2)	(superfluous,2)
(futile,9)	(futiilty,2)	(meaningitilesss,2)
(futilitye,9)	(applicatitrivialitis,2)	(meaningitlyss,2)
(meaningless,9)	(futiitiles,2)	(meaninglesity,2)
(meaninglityess,8)	(futiitly,2)	(meaninglessness,2)
(trivialitie,6)	(futiityl,2)	(meaninglityes,2)
(futility,4)	(futileity,2)	(meaniniglytyss,2)
(meaninglityss,4)	(futileessne,2)	(meaninitglyss,2)
(futilities,3)	(high-triviality,2)	(meaninitiglesss,2)
(meaninglesit,3)	(ltbbyirivialitig,2)	(mfaninglecilits,2)

...

# Illustration: traduction de mots inconnus

Step-3: sélection depuis  $\mathcal{E}_O$ (futilité)

## 13 solutions conservées

(trivialities,27) (triviality,14) (futile,9)  
(meaningless,9) (futility,4) (meaninglessness,3)  
(superfluous,2) (unwieldy,2) (unnecessary,2)  
(uselessness,2) (trivially,1) (tie,1) (trivial,1)

# Quelques idées fausses

- ▶ analogie  $\neq$  similarité ( $[x : y :: x : y]$ )
- ▶ analogie identifiées dans chaque espace (source et cible)  
**indépendamment**

biais inductif: homomorphisme

- ▶ ressemblances avec les k plus proches voisins, mais:
  - ▶ aucune distance
  - ▶ approche réversible
  - ▶ espace de sortie peut être structuré
- ▶ une analogie **fortuite** n'invalide pas l'apprentissage !  
[suddenly, she appeared : she appeared all of a sudden :: immediately, she appeared : she appeared all of a immediate]



# Pourquoi est-ce séduisant en TALN?

- ▶ Un corpus est plus qu'une séquence de **mots**
- ▶ Vue paradigmatisée des données:
  - ▶ [funny : funniest :: lucky : luckiest]
  - ▶ [suddenly, she appeared : she appeared suddenly :: today, she appeared : she appeared today]
  - ▶ [This guy drinks too much : This boat sank :: These guys drink too much : These boats sank]
  - ▶ [Elle est émue : Il est ému :: Elle est touchée : Il est touché]
- ▶ Aucune connaissance à utiliser pour définir la correspondance entre les deux espaces

Lire [Pirrelli & Yvon, 1999]

## Apprentissage Analogique

Analogie

Principe

## Quelques réalisations

### Sous le capot

Définitions de l'analogie (formelle)

Solveurs d'équations

Recherche des analogies

Filtrer le bruit

### Applications

Traduction de mots inconnus

Traduction de termes

Translittération

# Analogies sémantiques et tests SAT

(Turney & Littman, 2005; Turney, 2006)

stem	mason : stone
a)	teacher : chalk
b)	carpenter : wood
c)	soldier : gun
d)	photograph : camera
e)	book : word

<http://www.sadlier-oxford.com/phonics/analogies/analogiesx.htm>

# Analogies sémantiques et tests SAT

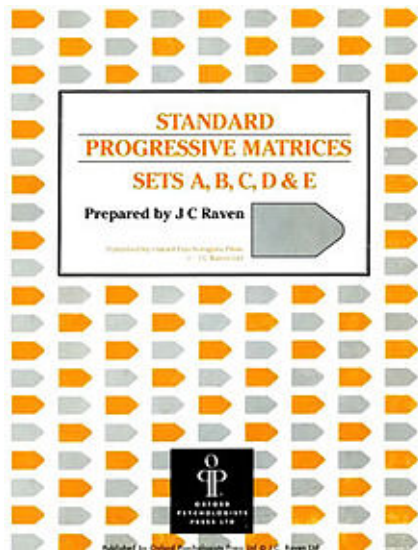
(Turney & Littman, 2005; Turney, 2006)

stem	mason : stone
a)	teacher : chalk
b)	carpenter : wood
c)	soldier : gun
d)	photograph : camera
e)	book : word

<http://www.sadlier-oxford.com/phonics/analogies/analogiesx.htm>

# Tests de Raven

Voir les slides de Correa, Prade & Richard



# ALEPH (Lepage & Denoual, 2005)

**Src:** *It floated across the river*

:

::

: It floated  
across the  
river

# ALEPH (Lepage & Denoual, 2005)

**Src:** *It floated across the river*

They swam in the sea : They swam across the river :: It floated in the sea : It floated across the river

# ALEPH (Lepage & Denoual, 2005)

**Src:** *It floated across the river*

They swam in the sea	:	They swam across the river	::	It floated in the sea	:	It floated across the river
↕		↕		↕		↕
Nadaron en el mar	:	Atraversaron el rio nadando	::	Flotó en el mar	:	x



# ALEPH (Lepage & Denoual, 2005)

Src: *It floated across the river*

They swam in the sea	:	They swam across the river	::	It floated in the sea	:	It floated across the river
↕		↕		↕		↕
Nadaron en el mar	:	Atraversaron el rio nadando	::	Flotó en el mar	:	x

Tgt: *Atraversó el rio flotando*

# Prolog est de retour !

## % bitext

translation( $s_1, \hat{s}_1$ ).

translation( $s_2, \hat{s}_2$ ).

...

translation( $s_n, \hat{s}_n$ ).

## % ALEPH

translation( $D, \hat{D}$ ) :-

translation( $A, \hat{A}$ )

translation( $B, \hat{B}$ )

translation( $C, \hat{C}$ )

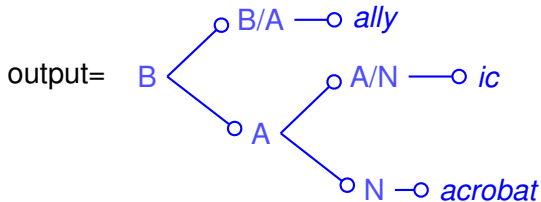
analogie( $A, B, C, D$ )

analogie( $\hat{A}, \hat{B}, \hat{C}, \hat{D}$ )

assert(translation( $D, \hat{D}$ )).

## (Stroppa & Yvon, 2005)

- ▶ **Task 1:** mot inconnu → traits morpho-syntaxiques  
input = *replying*  
output = {*reply*, V-PP}
- ▶ **Task 2:** mot inconnu → analyse morphologique  
input= *acrobatically*



## (Moreau et al, 2007)

▶  $m_1 = \text{désinstaller}$

$m_2 = m_1 - \text{dés}_p + \text{ré}_p - \text{er}_s + \text{ation}_s$

▶  $m_1 = \text{désydrater}$

$m_2 = \text{réinstallation}$

## (Moreau et al, 2007)

▶  $m_1 = \text{désinstaller}$

$m_2 = m_1 - \text{dés}_p + \text{ré}_p - \text{er}_s + \text{ation}_s$

▶  $m_1 = \text{désydrater}$

$m_2 = \text{réinstallation}$

$m_2 = \text{réhydratation}$

# (Moreau et al, 2007)

▶  $m_1 =$  désinstaller

$m_2 =$  réinstallation

$m_2 = m_1 - \text{dés}_p + \text{ré}_p - \text{er}_s + \text{ation}_s$

▶  $m_1 =$  déshydrater

## ENTRAÎNEMENT

1. soit  $d$  un document de la collection  $C$
2. soit les mots de  $d$  qui partagent une grande sous-chaîne ( $\geq 7$  car.)
3. mémoriser les correspondances
4. aller en 1

## EXPANSION DE REQUÊTE

1. Soit  $m$  un mot de la requête  $r$
2. Appliquer les règles à  $m \Rightarrow S$
3. Ajouter à  $r$  les mots de  $S \in \text{voc}(C)$

▶ **Ex:** *pollution des eaux souterraines*

↪ *pollution des eaux souterraines polluants dépollution  
anti-pollution pollutions polluées polluant eau souterraine  
souterrains souterrain*

Gains consistants pour des collections de petite taille dans 6 langues (30 requêtes)

## Apprentissage Analogique

Analogie

Principe

## Quelques réalisations

### Sous le capot

Définitions de l'analogie (formelle)

Solveurs d'équations

Recherche des analogies

Filtrer le bruit

### Applications

Traduction de mots inconnus

Traduction de termes

Translittération

# Définitions de l'analogie formelle

(Pirrelli & Yvon, 1999)

- ▶ Symboles:

$$[x : y :: z : t] \iff \text{or } \begin{cases} x = y \text{ and } z = t \\ x = z \text{ and } y = t \end{cases}$$

[Noun : Noun :: Verb : Verb] , [Noun : Verb :: Noun : Verb]

- ▶ Chaînes:

$$[x : y :: z : t] \iff \text{or } \begin{cases} x = bc, y = bd, z = ac, t = ad \\ x = bc, y = ac, z = bd, t = ad \end{cases}$$

[dream : dreamer :: eat : eater] , [steal : ceal :: stage : cage]

- ▶ S'étend aux ensembles de valeurs



# Définition de l'analogie formelle

(Lepage, 1998)

- ▶ chaînes:

$$[x : y :: z : t] \Rightarrow \begin{cases} \sigma(y, t) & = -|x| + |y| + \sigma(x, z) \\ \sigma(z, t) & = -|x| + |z| + \sigma(x, y) \\ \sigma(x, y, z, t) & = -|x| + \sigma(x, y) + \sigma(x, z) \\ |t|_a & = -|x|_a + |y|_a + |z|_a \forall a \end{cases}$$

[believer : unbelievable :: dreamer : undreamable]

- ▶ Formalisé par un algorithme

# Définition de l'analogie formelle

(Stroppa & Yvon, 2005)

Pour tout  $(x, y, z, t) \in \Sigma^{*4}$ ,  $[x : y :: z : t]$  **ssi** on peut trouver une **factorisation**  $(f_x, f_y, f_z, f_t) \in (\Sigma^{*d})^4$  telle que  $\forall i \in [1, d]$ :

$$(f_y^{(i)}, f_z^{(i)}) \in \left\{ (f_x^{(i)}, f_t^{(i)}), (f_t^{(i)}, f_x^{(i)}) \right\}$$

- ▶ le plus petit  $d$  pour lequel cela est vrai est appelé le **degré**
- ▶  $f_x^{(i)}, f_y^{(i)}, f_z^{(i)}$  and  $f_t^{(i)}$  sont appelés les **facteurs**

x ≡	th	is	guy	ε	dr	inks	too much
y ≡	th	is	boat	ε	s	inks	ε
z ≡	th	ese	guy	s	dr	unk	too much
t ≡	th	ese	boat	s	s	unk	ε

# Deux solveurs génériques

[Lepage, 1998]

Introduction Apprentissage Analogique Expériences Discussion

[even : usual = unevenly : ?]

4	4	4	4	4	N<	4	4	3	3	2	1	0	0	0
3	3	3	3	3	E	3	3	3	2	1	0	0	0	0
2	2	2	2	2	V	2	2	2	1	0	0	0	0	0
1	1	1	1	1	E	1	1	1	0	0	0	0	0	0
0	0	0	0	0	A	0	0	0	0	0	0	0	0	0
L	A	U	S	U	B	C	U	N	E	V	E	N	L	Y

△

{N,L,Y, [SUB,DEL], MOVE C + COPY Y }

	E	V	E	N
U	S	U	A	L

△

E	V	E	N				
U	N	E	V	E	N	L	Y

△

[Stroppa & Yvon, 2005]

Introduction Apprentissage Analogique Expériences Discussion

édition • dictateur \ éditeur

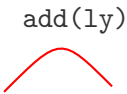
△

# Intuition

even :    evenly =    x :

# Intuition

add(1y)  
even : evenly = x :



# Intuition

add(ly)  
even : evenly =

add(ly)  
x : xly

# Intuition

he drinks : they drank = *he sinks*

# Intuition

?

he drinks : they drank = *he sinks*



# Intuition

?

he drinks : they drank =

?

*he sinks* they sank

# [even : usual = unevenly : ?]

4	4	4	4	4	4	N	4	4	3	3	2	1	0	0	0
3	3	3	3	3	3	E	3	3	3	2	1	0	0	0	0
2	2	2	2	2	2	V	2	2	2	1	0	0	0	0	0
1	1	1	1	1	1	E	1	1	1	0	0	0	0	0	0
0	0	0	0	0	0	A	0	0	0	0	0	0	0	0	0
L	A	U	S	U	B		C	U	N	E	V	E	N	L	Y

⟨N,L,Y, [SUB,DEL], MOVE C + COPY Y⟩

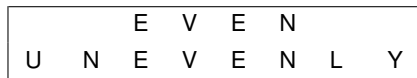
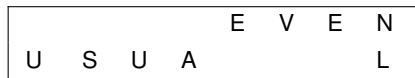
				E	V	E	N
U	S	U	A				L

		E	V	E	N		
U	N	E	V	E	N	L	Y

# [even : usual = unevenly : ?]

4	4	4	4	4	4	N	4	4	3	3	2	1	0	0	0
3	3	3	3	3	3	E	3	3	3	2	1	0	0	0	0
2	2	2	2	2	2	V	2	2	2	1	0	0	0	0	0
1	1	1	1	1	1	E	1	1	1	0	0	0	0	0	0
0	0	0	0	0	0	A	0	0	0	0	0	0	0	0	0
L	A	U	S	U	B		C	U	N	E	V	E	N	L	Y
△						△									

⟨N,L,L, [SUB,DEL], MOVE C + COPY L⟩



Y

# [even : usual = unevenly : ?]

<b>4</b>	4	4	4	4	4	N	<	4	4	3	3	2	1	<b>0</b>	<b>0</b>	<b>0</b>
3	<b>3</b>	3	3	3	3	E		3	3	3	2	1	<b>0</b>	0	0	0
2	<b>2</b>	2	2	2	2	V		2	2	2	1	<b>0</b>	0	0	0	0
1	<b>1</b>	1	1	1	1	E		1	1	1	<b>0</b>	0	0	0	0	0
0	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>A</b>		<b>0</b>	<b>0</b>	<b>0</b>	0	0	0	0	0	0
L	A	U	S	U	<b>B</b>			<b>C</b>	U	N	E	V	E	N	L	Y

<N,L,N, [SUB,<=>], MOVE A,B,C + COPY L >

										E	V	E	N
U	S	U	A										L

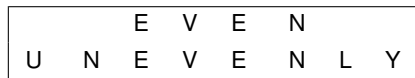
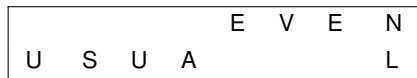
										E	V	E	N
U	N	E	V	E	N	L	Y						

LY

# [even : usual = unevenly : ?]

4	4	4	4	4	4	N	4	4	3	3	2	1	0	0	0
3	3	3	3	3	3	E<	3	3	3	2	1	0	0	0	0
2	2	2	2	2	2	V	2	2	2	1	0	0	0	0	0
1	1	1	1	1	1	E	1	1	1	0	0	0	0	0	0
0	0	0	0	0	0	A	0	0	0	0	0	0	0	0	0
L	A	U	S	U	B		C	U	N	E	V	E	N	L	Y

⟨E,A,E, [INS,<=>], MOVE A,C ⟩

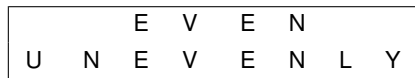
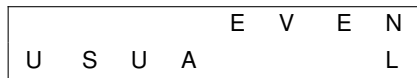


LLY

# [even : usual = unevenly : ?]

4	4	4	4	4	4	N	4	4	3	3	2	1	0	0	0
3	3	3	3	3	3	E	3	3	3	2	1	0	0	0	0
2	2	2	2	2	2	V<	2	2	2	1	0	0	0	0	0
1	1	1	1	1	1	E	1	1	1	0	0	0	0	0	0
0	0	0	0	0	0	A	0	0	0	0	0	0	0	0	0
L	A	U	S	U	B		C	U	N	E	V	E	N	L	Y

$\langle V, A, V, [INS, \langle = \rangle], MOVE A, C \rangle$



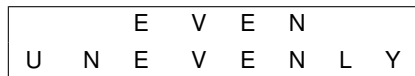
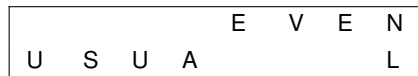
LLY

# [even : usual = unevenly : ?]

4	4	4	4	4	4	N	4	4	3	3	2	1	0	0	0
3	3	3	3	3	3	E	3	3	3	2	1	0	0	0	0
2	2	2	2	2	2	V	2	2	2	1	0	0	0	0	0
1	1	1	1	1	1	E<	1	1	1	0	0	0	0	0	0
0	0	0	0	0	0	A	0	0	0	0	0	0	0	0	0
L	A	U	S	U	B		C	U	N	E	V	E	N	L	Y



⟨E,A,E, [INS,<=>], MOVE A,C ⟩

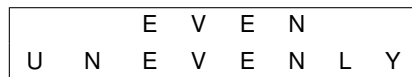
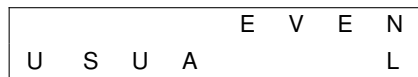


LLY

# [even : usual = unevenly : ?]

4	4	4	4	4	4	N	4	4	3	3	2	1	0	0	0
3	3	3	3	3	3	E	3	3	3	2	1	0	0	0	0
2	2	2	2	2	2	V	2	2	2	1	0	0	0	0	0
1	1	1	1	1	1	E	1	1	1	0	0	0	0	0	0
0	0	0	0	0	0	A $\triangleleft$	0	0	0	0	0	0	0	0	0
L	A	U	S	U	B		C	U	N	E	V	E	N	L	Y

$\langle \epsilon, A, N, [\text{DEL}, \text{DEL}], \text{MOVE } C + \text{COPY } N \rangle$



LLY



# [even : usual = unevenly : ?]

4	4	4	4	4	4	N	4	4	3	3	2	1	0	0	0
3	3	3	3	3	3	E	3	3	3	2	1	0	0	0	0
2	2	2	2	2	2	V	2	2	2	1	0	0	0	0	0
1	1	1	1	1	1	E	1	1	1	0	0	0	0	0	0
0	0	0	0	0	0	A $\triangleleft$	0	0	0	0	0	0	0	0	0
L	A	U	S	U	B		C	U	N	E	V	E	N	L	Y

$\langle \epsilon, A, U, [\text{DEL}, \text{DEL}], \text{MOVE } C + \text{COPY } U \rangle$

				E	V	E	N
U	S	U	A				L

		E	V	E	N		
U	N	E	V	E	N	L	Y

NLLY

# [even : usual = unevenly : ?]

4	4	4	4	4	4	N	4	4	3	3	2	1	0	0	0
3	3	3	3	3	3	E	3	3	3	2	1	0	0	0	0
2	2	2	2	2	2	V	2	2	2	1	0	0	0	0	0
1	1	1	1	1	1	E	1	1	1	0	0	0	0	0	0
0	0	0	0	0	0	A $\triangleleft$	0	0	0	0	0	0	0	0	0
L	A	U	S	U	B		C	U	N	E	V	E	N	L	Y

$\langle \epsilon, A, \epsilon, [\text{DEL}, \bullet], \text{MOVE B} + \text{COPY A} \rangle$

				E	V	E	N
U	S	U	A				L

		E	V	E	N		
U	N	E	V	E	N	L	Y

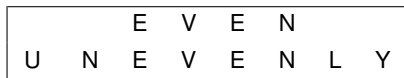
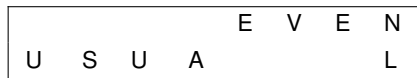
UNLLY

# [even : usual = unevenly : ?]

4	4	4	4	4	4	N	4	4	3	3	2	1	0	0	0
3	3	3	3	3	3	E	3	3	3	2	1	0	0	0	0
2	2	2	2	2	2	V	2	2	2	1	0	0	0	0	0
1	1	1	1	1	1	E	1	1	1	0	0	0	0	0	0
0	0	0	0	0	0	A $\triangleleft$	0	0	0	0	0	0	0	0	0
L	A	U	S	U	B		C	U	N	E	V	E	N	L	Y



$\langle \varepsilon, U, \varepsilon, [\text{DEL}, \bullet], \text{MOVE B} + \text{COPY U} \rangle$



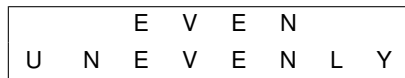
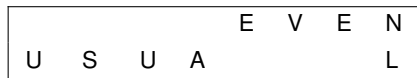
AUNLLY

# [even : usual = unevenly : ?]

4	4	4	4	4	4	N	4	4	3	3	2	1	0	0	0
3	3	3	3	3	3	E	3	3	3	2	1	0	0	0	0
2	2	2	2	2	2	V	2	2	2	1	0	0	0	0	0
1	1	1	1	1	1	E	1	1	1	0	0	0	0	0	0
0	0	0	0	0	0	A $\triangleleft$	0	0	0	0	0	0	0	0	0
L	A	U	S	U	B		C	U	N	E	V	E	N	L	Y



$\langle \epsilon, S, \epsilon, [\text{DEL}, \bullet], \text{MOVE B} + \text{COPY S} \rangle$



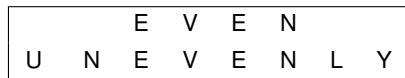
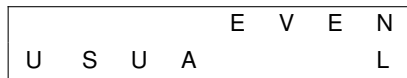
UAUNLLY

# [even : usual = unevenly : ?]

4	4	4	4	4	4	N	4	4	3	3	2	1	0	0	0
3	3	3	3	3	3	E	3	3	3	2	1	0	0	0	0
2	2	2	2	2	2	V	2	2	2	1	0	0	0	0	0
1	1	1	1	1	1	E	1	1	1	0	0	0	0	0	0
0	0	0	0	0	0	A $\triangleleft$	0	0	0	0	0	0	0	0	0
L	A	U	S	U	B		C	U	N	E	V	E	N	L	Y



$\langle \epsilon, U, \epsilon, [\text{DEL}, \bullet], \text{MOVE B} + \text{COPY U} \rangle$



SUAUNLLY

# [even : usual = unevenly : ?]

4	4	4	4	4	4	N	4	4	3	3	2	1	0	0	0
3	3	3	3	3	3	E	3	3	3	2	1	0	0	0	0
2	2	2	2	2	2	V	2	2	2	1	0	0	0	0	0
1	1	1	1	1	1	E	1	1	1	0	0	0	0	0	0
0	0	0	0	0	0	A◀	0	0	0	0	0	0	0	0	0
L	A	U	S	U	B		C	U	N	E	V	E	N	L	Y

△                      △  
 ⟨ $\epsilon, \epsilon, \epsilon, [\bullet, \bullet], \text{STOP}$ ⟩

						E	V	E	N
U	S	U	A						L



						E	V	E	N
U	N	E	V	E	N	L	Y		



USUAUNLLY

[even : usual = unevenly : ?]

4	4	4	4	4	4	N	4	4	3	3	2	1	0	0	0
3	3	3	3	3	3	E	3	3	3	2	1	0	0	0	0
2	2	2	2	2	2	V	2	2	2	1	0	0	0	0	0
1	1	1	1	1	1	E	1	1	1	0	0	0	0	0	0
0	0	0	0	0	0	A◀	0	0	0	0	0	0	0	0	0
L	A	U	S	U	B		C	U	N	E	V	E	N	L	Y



E				V	E	N
	U	S	U	A		L



		E	V	E	N		
U	N	E	V	E	N	L	Y



UNUSUALLY

# [even : usual = unevenly : ? ]

15 solutions (681 synchronisations):

<i>uunsually</i>	<i>usuaunlly</i>	<i>usuunalyl</i>
<b>unusually</b>	<i>usunually</i>	<i>uunslyual</i>
<i>unulyusual</i>	<i>uunsualyl</i>	<i>unuslyual</i>
<i>unusualyl</i>	<i>usunualyl</i>	<i>uunsulyal</i>
<i>unusulyal</i>	<i>usunulyal</i>	<i>usuunally</i>

(72 solutions selon la définition de [Stroppa & Yvon, 2005])



$$[x : y :: z : t] \iff t \in (y \bullet z) \setminus x$$

(Yvon & al., 2004)

- ▶ **shuffle**  $a \bullet b$  lire séquences dans  $a$  et  $b$  de gauche à droite, en autorisant de changer de chaîne

*spondyondontilalgias* et

*ondspndonlaltitigia*  $\in$  *spondylalgia*  $\circ$  *ondontitis*

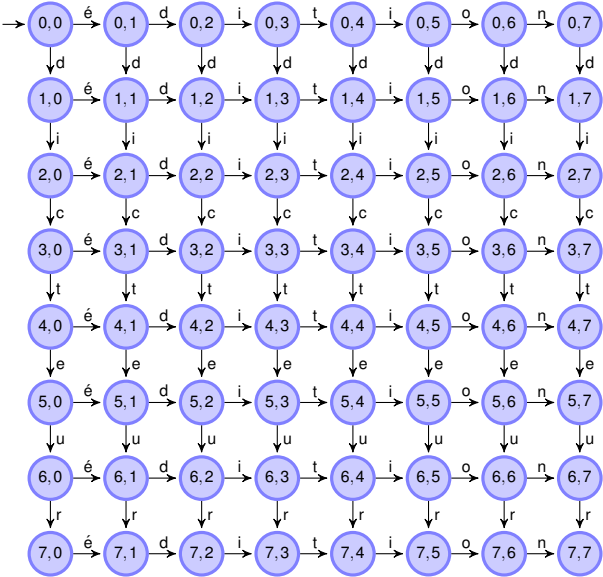
- ▶ **complément**  $a \setminus b$  chaînes obtenues en retirant la sous-chaîne  $b$  dans  $a$

*spondylitis*  $\in$  *spondyondontilalgias*  $\setminus$  *ondontalgia*

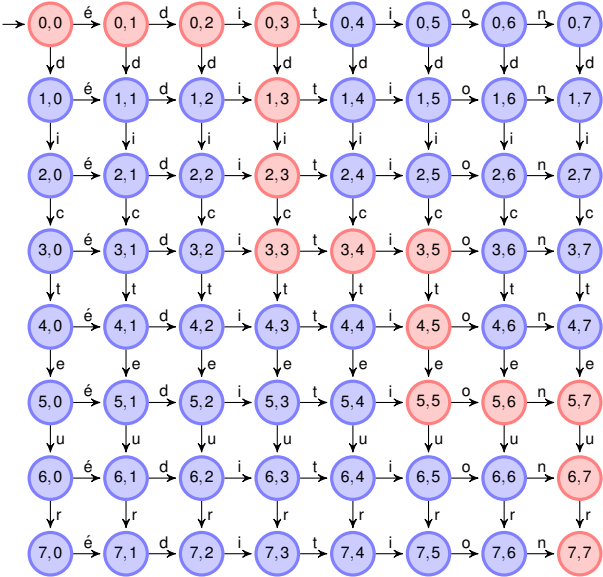
*spydoniltis*  $\in$  *spondyondontilalgias*  $\setminus$  *ondontalgia*

$\{(y \bullet z) \setminus x\}$  est un langage **régulier**

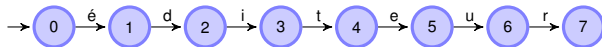
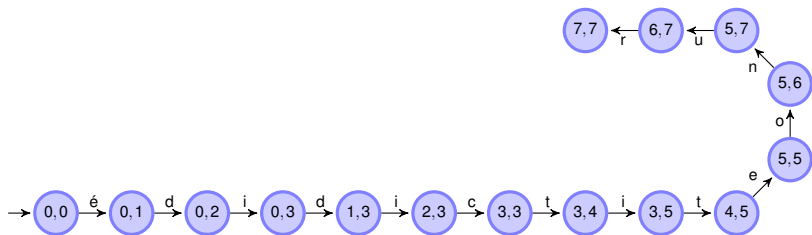
# [éditeur : édition :: dictateur : ?]



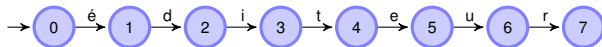
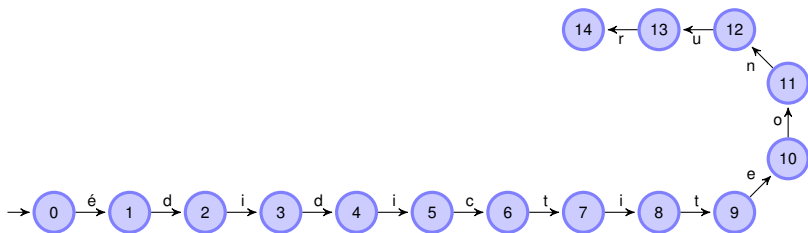
# [éditeur : édition :: dicteur : ?]



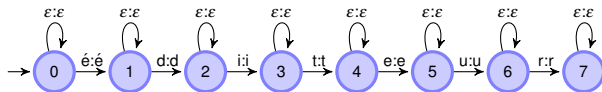
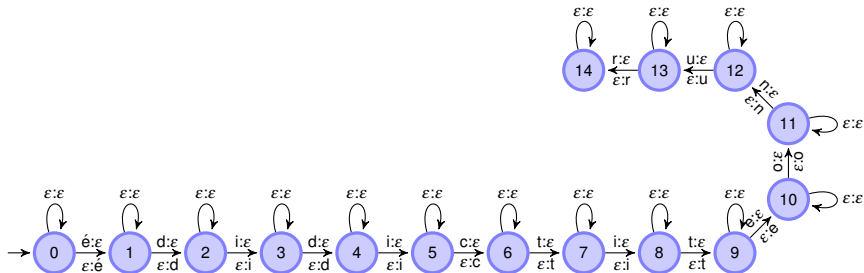
# édition • dicteur \ éditeur



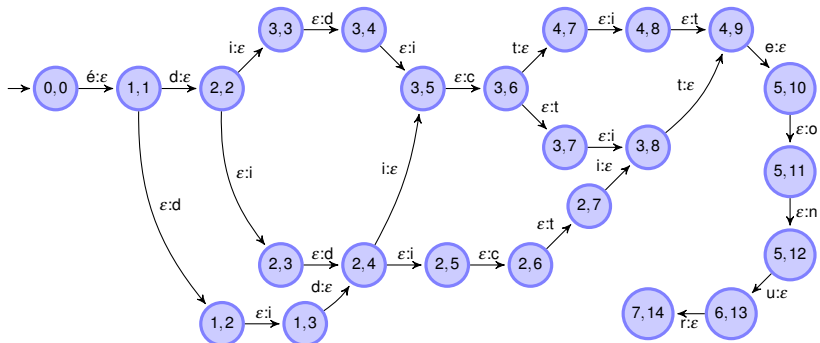
# édition • dicteur \ éditeur



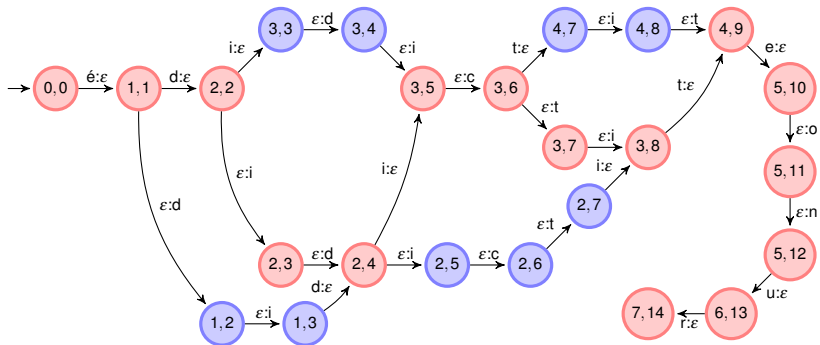
# édition • dicteur \ éditeur



# édition • dicteur \ éditeur



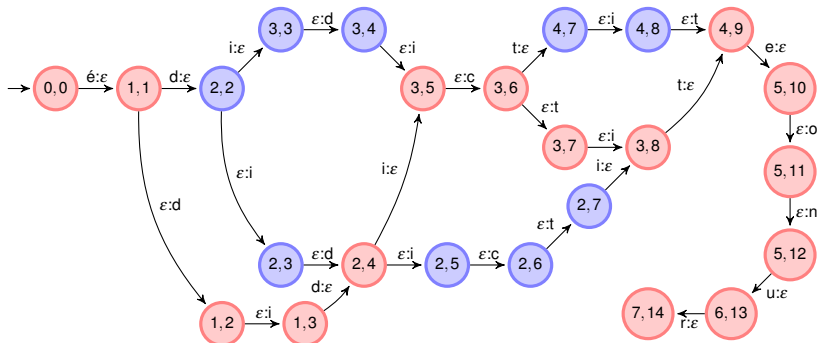
# édition • dicteur \ éditeur



idction

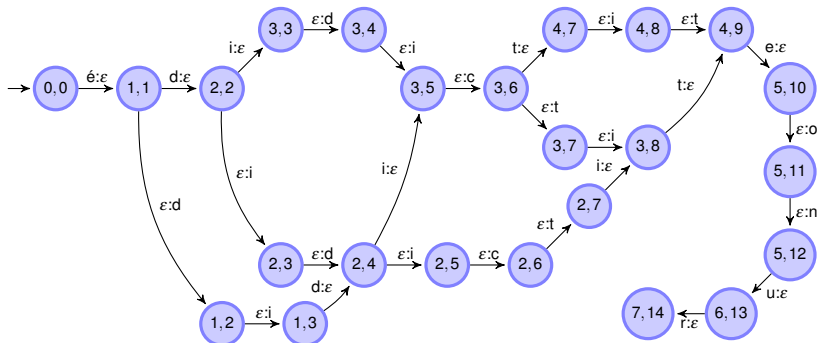


# édition • dicteur \ éditeur



idction , diction

# édition • dicteur \ éditeur



idiction , diction , idciton, diicton, diciton

# [éditeur : édition :: dicteur : ?]

**110 sol.**  $f(\text{diction}) = 5180, f(\text{ditcion}) = 1155, \dots, f(\text{tiondic}) = 1$

diction ditcion dcition diticon diciton diicton diitcon idction  
ditiocn diciotn cdition dticion dtiicon diicotn diioctn icdtion  
dctiion idtcion idticon itdicon ditionc dicient diitocn idcition  
dciiton dtiioen idicton iditcon diicont diioctn diionct ditoicn  
idtioen itdcion tdicion dtioicn itdioen tdiicon dioicn icdicion  
dciioen cdtiion cdtiion icdtion idciotn itidcon tidicon diitocn  
dtiionc idicotn idioctn cdiiton iditocn itdoicn itidocn tdiioen  
idtione dioient dioinct dionict ditoinc ditonic itdione dciiont  
dtioinc dtionic idciont idtoicn itioden itodien tdioicn tidioen  
tidoicn tiodien cdiioen icdiotn idoiotn idoiotn idoiotn idoiotn idoiotn  
idionct itidocn iditocn itdoicn itdonic itiodnc tdiionc dciiont  
icdiont idoiotn idoiotn idoiotn idoiotn idoiotn idoiotn idoiotn idoiotn  
iodnict iondict itionc itodinc itodnic itondic tdiioinc tdionic  
tidione tidoicn tidonic tiodinc tiodnic tiondic

# À propos de ces solveurs

- ▶ 0, 1 ou plusieurs solutions (souvent inutiles) à une équation analogique

# À propos de ces solveurs

- ▶ 0, 1 ou plusieurs solutions (souvent inutiles) à une équation analogique
- ▶ solutions(lepage)  $\subseteq$  solutions(Stroppa&Yvon)

# À propos de ces solveurs

- ▶ 0, 1 ou plusieurs solutions (souvent inutiles) à une équation analogique
- ▶ solutions(lepage)  $\subseteq$  solutions(Stroppa&Yvon)
- ▶ les deux solveurs ont une complexité exponentielle
  - ▶ implémentation par échantillonnage
  - ▶ aucune garantie d'obtenir une solution lorsqu'il en existe une

# Exemple

[even : usual :: unevenly : ?]

$\rho$	$nb$	solutions		
20	12	usuaunlly (3)	unusually (2)	usunually (2)
100	34	unusually (6)	usuaunlly (6)	uunusually (4)
1000	67	unusually (57)	uunusually (23)	usuunally (19)
2000	72	unusually (130)	uunusually (77)	usunually (43)

- ▶  $\rho$ : fréquence d'échantillonnage
- ▶  $nb$ : nombre de solutions produites

# Exemple

[this guy drinks too much : this boat sinks :: those guys drink too much : ?]

$\rho = 20$        $nb = 8$   
 $t = 0.0003$     $rang = \phi$

thos\_boate\_sinks (2)

tho\_boatse\_sinks (2)

thoatse\_\_sinks (2)

$\rho = 100$        $nb = 28$   
 $t = 0.001$       $rang = 13$

thoatse\_\_sinks (2)

tho\_boatse\_sinks (2)

those\_sboat\_sink (2)

$\rho = 1000$        $nb = 28$   
 $t = 0.009$       $rang = 2$

those\_boat\_ssink (5)

those\_boats\_sink (5)

thoes\_tboa\_sinks (5)

$\rho = 10^6$        $nb = 19796$   
 $t = 3.82$        $rang = 10$

thoes\_boat\_sinks (2550)

thoses\_boat\_sink (1037)

those\_boat\_ssink (999)

- ▶  $\rho$ : fréquence d'échantillonnage
- ▶  $nb$ : nombre de solutions produites
- ▶  $rang$ : rang de la solution attendue ( $\phi$  veut dire "pas trouvé")
- ▶  $t$ : temps mis par le solver (sec.)



# Exemple

[this guy drinks too much : this boat sinks :: those guys drink too much : ?]

---

$\rho = 20$	$nb = 8$
$t = 0.0003$	$rang = \phi$

---

thos\_boate\_sinks (2)  
tho\_boatse\_sinks (2)  
thoboatse\_\_sinks (2)

---

---

$\rho = 100$	$nb = 28$
$t = 0.001$	$rang = 13$

---

thoboatse\_\_sinks (2)  
tho\_boatse\_sinks (2)  
those\_sboat\_sink (2)

---

---

$\rho = 1000$	$nb = 28$
$t = 0.009$	$rang = 2$

---

those\_boat\_ssink (5)  
those\_boats\_sink (5)  
thoes\_tboa\_sinks (5)

---

---

$\rho = 10^6$	$nb = 19796$
$t = 3.82$	$rang = 10$

---

thoes\_boat\_sinks (2550)  
thoses\_boat\_sink (1037)  
those\_boat\_ssink (999)

---

- ▶  $\rho$ : fréquence d'échantillonnage
- ▶  $nb$ : nombre de solutions produites
- ▶  $rang$ : rang de la solution attendue ( $\phi$  veut dire "pas trouvé")
- ▶  $t$ : temps mis par le solver (sec.)

# Exemple

[this guy drinks too much : this boat sinks :: those guys drink too much : ?]

---

$\rho = 20$	$nb = 8$
$t = 0.0003$	$rang = \phi$

---

thos\_boate\_sinks (2)  
tho\_boatse\_sinks (2)  
thoboatse\_\_sinks (2)

---

---

$\rho = 100$	$nb = 28$
$t = 0.001$	$rang = 13$

---

thoboatse\_\_sinks (2)  
tho\_boatse\_sinks (2)  
those\_sboat\_sink (2)

---

---

$\rho = 1000$	$nb = 28$
$t = 0.009$	$rang = 2$

---

those\_boat\_ssink (5)  
those\_boats\_sink (5)  
thoes\_tboa\_sinks (5)

---

---

$\rho = 10^6$	$nb = 19796$
$t = 3.82$	$rang = 10$

---

thoes\_boat\_sinks (2550)  
thoses\_boat\_sink (1037)  
those\_boat\_ssink (999)

---

- ▶  $\rho$ : fréquence d'échantillonnage
- ▶  $nb$ : nombre de solutions produites
- ▶  $rang$ : rang de la solution attendue ( $\phi$  veut dire "pas trouvé")
- ▶  $t$ : temps mis par le solver (sec.)

# Exemple

[this guy drinks too much : this boat sinks :: those guys drink too much : ?]

---

$\rho = 20$	$nb = 8$
$t = 0.0003$	$rang = \phi$
<hr/>	
thos_boate_sinks (2)	
tho_boatse_sinks (2)	
thoboatse__sinks (2)	

---

---

$\rho = 1000$	$nb = 28$
$t = 0.009$	$rang = 2$
<hr/>	
those_boat_ssink (5)	
those_boats_sink (5)	
thoes_tboa_sinks (5)	

---

---

$\rho = 100$	$nb = 28$
$t = 0.001$	$rang = 13$
<hr/>	
thoboatse__sinks (2)	
tho_boatse_sinks (2)	
those_sboat_sink (2)	

---

---

$\rho = 10^6$	$nb = 19796$
$t = 3.82$	$rang = 10$
<hr/>	
thoes_boat_sinks (2550)	
thoses_boat_sink (1037)	
those_boat_ssink (999)	

---

- ▶  $\rho$ : fréquence d'échantillonnage
- ▶  $nb$ : nombre de solutions produites
- ▶  $rang$ : rang de la solution attendue ( $\phi$  veut dire "pas trouvé")
- ▶  $t$ : temps mis par le solver (sec.)

# Recherche des analogies

1.  $\mathcal{E}_{\mathcal{I}}(u) \equiv \{(s, v, w) \in \mathcal{L}^3 \mid [I(s) : I(v) :: I(w) : I(u)]\}$

- ▶ Cubique en  $\mathcal{I}$  ...
- ▶ Nombre quadratique de résolutions d'équations [Lepage & Denoual, 2005]:
  1. soit  $(x, y) \in \mathcal{I}^2$
  2. résoudre  $[y : x :: t : ?]$
  3. garder les solutions  $z$  qui appartiennent à  $\mathcal{I}$   
 $\Rightarrow$  triplets  $(x, y, z)$

Car:

$$[x : y :: z : t] \Leftrightarrow [y : x :: t : z]$$

- ▶ Trop couteux cependant pour des espaces sources modérés  
 $\Rightarrow$  échantillonnage de  $(x, y)$

# Recherche des analogies

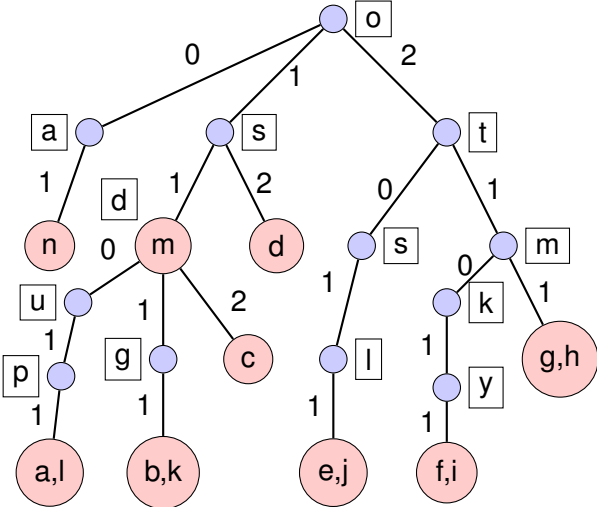
Peut découvrir des **analogies potentielles** en un temps linéaire, grâce à:

$$[x : y :: z : t] \Rightarrow |x|_c + |t|_c = |y|_c + |z|_c \quad \forall c \in \mathcal{A}$$

1. soit  $x \in \mathcal{I}$
2. chercher les paires  $(y, z)$  vérifiant la propriété sur les comptes
3. vérifier les **véritables analogies**  
(algorithme en  $O(|x| \times |y| \times |z| \times |t|)$ )

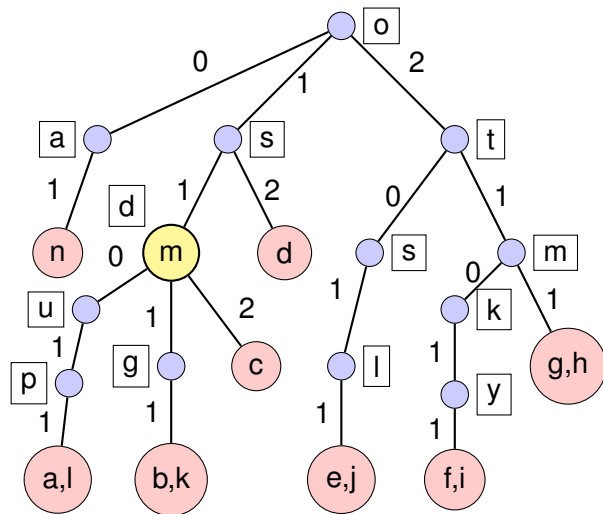
L'espace d'entrée doit être organisé (**tree-count**).

# Un tree-count



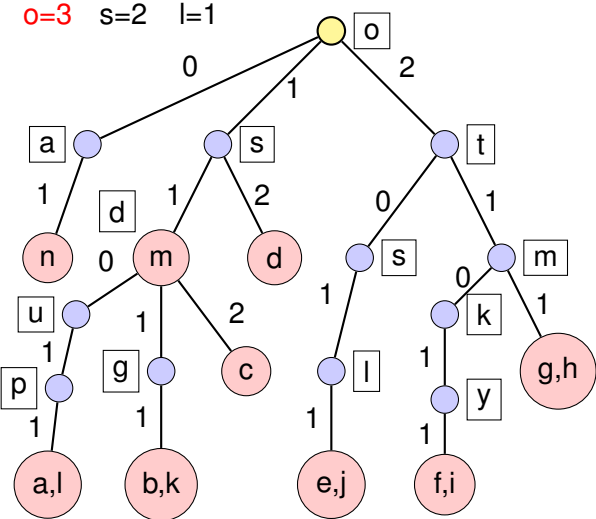
a	soup
b	gods
c	odds
d	sos
e	solo
f	tokyo
g	moot
h	moto
i	kyoto
j	oslo
k	dogs
l	opus
m	os
n	a

# Un tree-count



a	soup
b	gods
c	odds
d	sos
e	solo
f	tokyo
g	moot
h	moto
i	kyoto
j	oslo
k	dogs
l	opus
m	os
n	a

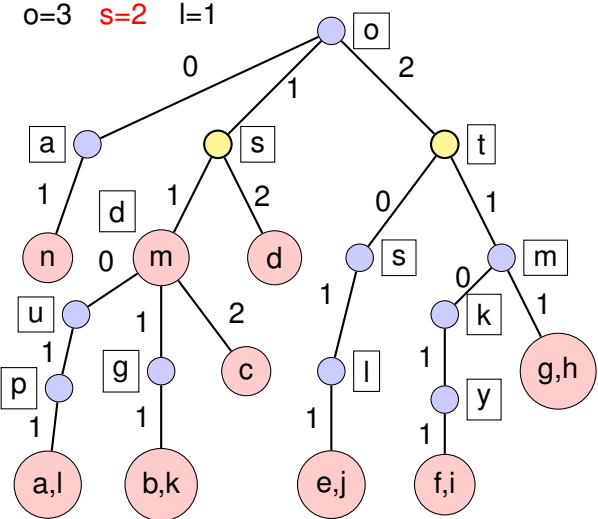
# Un tree-count



a	soup
b	gods
c	odds
d	sos
e	solo
f	tokyo
g	moot
h	moto
i	kyoto
j	oslo
k	dogs
l	opus
m	os
n	a



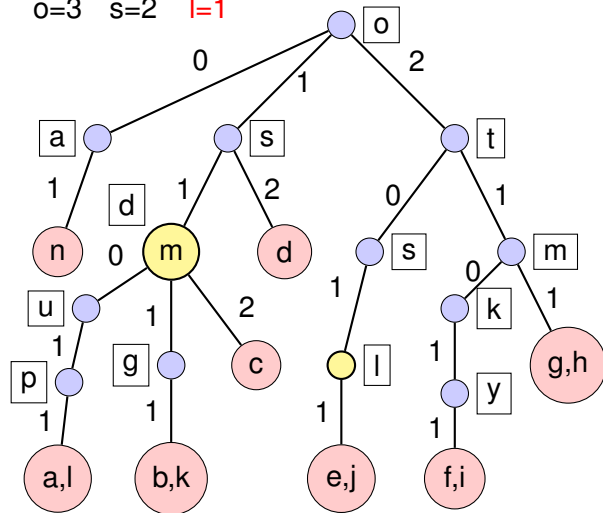
# Un tree-count



a	soup
b	gods
c	odds
d	sos
e	solo
f	tokyo
g	moot
h	moto
i	kyoto
j	oslo
k	dogs
l	opus
m	os
n	a

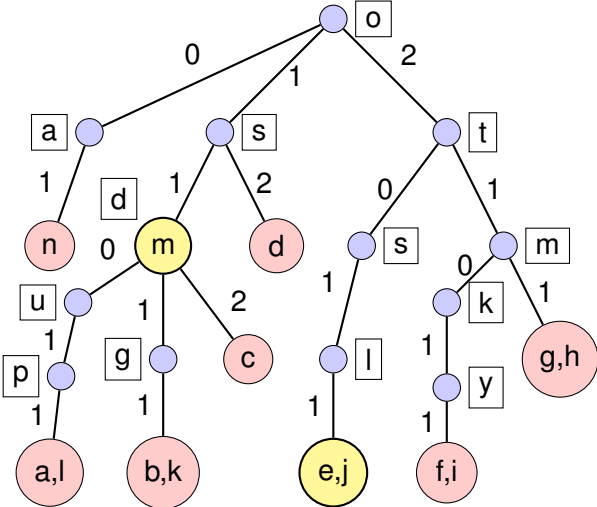
# Un tree-count

$o=3$   $s=2$   $l=1$



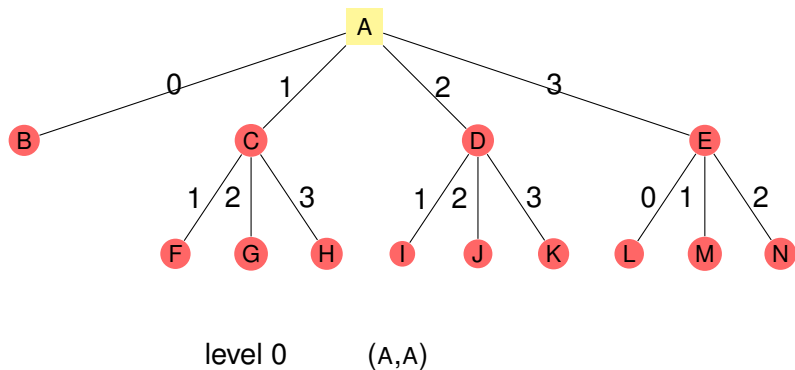
a	soup
b	gods
c	odds
d	sos
e	solo
f	tokyo
g	moot
h	moto
i	kyoto
j	oslo
k	dogs
l	opus
m	os
n	a

# Un tree-count

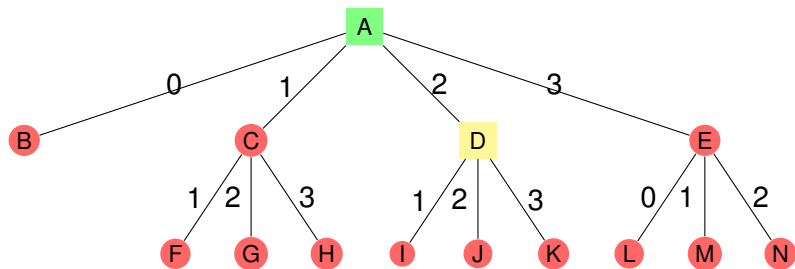


a	soup
b	gods
c	odds
d	sos
e	solo
f	tokyo
g	moot
h	moto
i	kyoto
j	oslo
k	dogs
l	opus
m	os
n	a

# Recherche des analogies potentielles

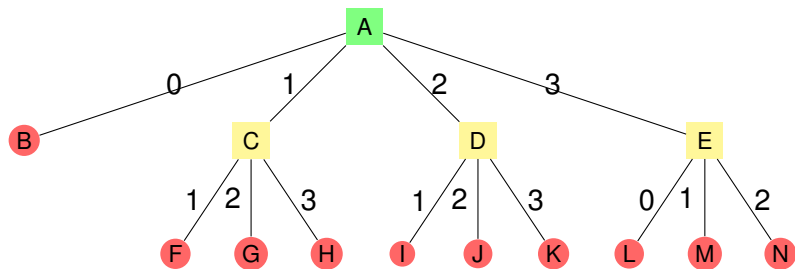


# Recherche des analogies potentielles



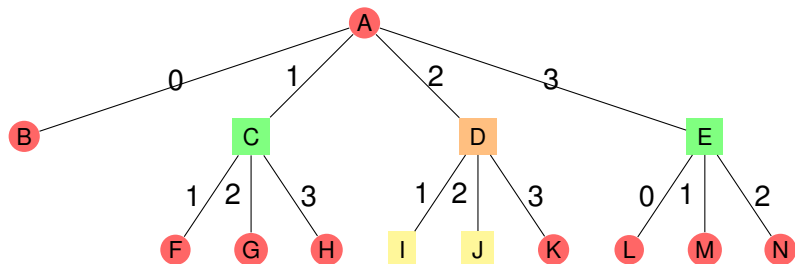
level 0		(A,A)
level 1	4	(D,D)
level 2	3	

# Recherche des analogies potentielles



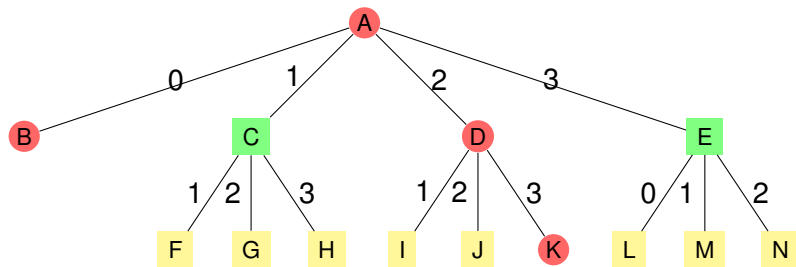
level 0		(A,A)
level 1	4	(D,D) (C,E)
level 2	3	

# Recherche des analogies potentielles



level 0		(A,A)
level 1	4	( <b>D,D</b> ) (C,E)
level 2	3	(I,J)

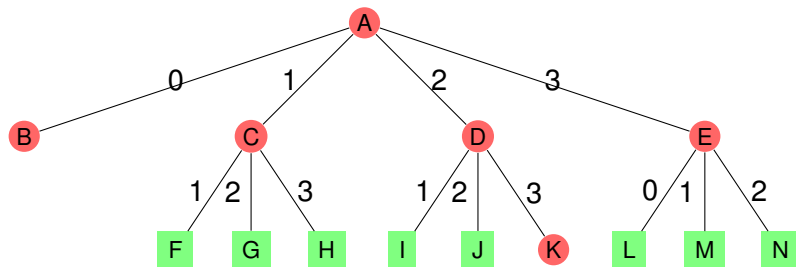
# Recherche des analogies potentielles



level 0		(A,A)
level 1	4	(D,D) <b>(C,E)</b>
level 2	3	(I,J) (G,M) (H,L) (F,N)



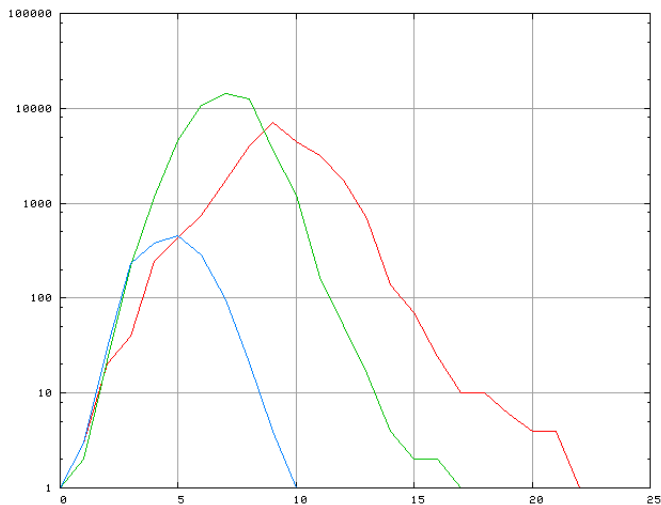
# Recherche des analogies potentielles



level 0      (A,A)  
level 1    4    (D,D) (C,E)  
level 2    3    (I,J) (G,M) (H,L) (F,N)  
...

# Recherche dans un tree-count

une figure sans légende ...



# A tree-count

- ▶ espace d'entrée: 11 317 717 formes
- ▶ en moyenne sur 1 000 recherches:

ratio	time (ms)	<i>frontier</i>	nodes
1/1000	5.5e-05	38	6.8
1/100	0.0003	150	6.3
1/10	0.003	1082	6.6
1/5	0.0055	1655	6.5
1/1	0.02	3921	5.8

- ▶ mémoire et temps  $\sim$  linéaire avec la taille de l'espace d'entrée

# Vérification d'une analogie

(Stroppa, 2005)

```
a(i,j,k,l) ← false , if i, j, k or l < 0
for i ← 0 to |x| do
  for j ← 0 to |y| do
    for k ← 0 to |z| do
      for l ← 0 to |t| do
        if i = j = k = l then
          a(i,j,k,l) ← true
        else
          a(i,j,k,l) ←
            or {
              a(i-1,j-1,k,l) ∧ x[i] = y[j]
              a(i-1,j,k-1,l) ∧ x[i] = z[k]
              a(i,j-1,k,l-1) ∧ t[l] = y[j]
              a(i,j,k-1,l-1) ∧ t[l] = z[k]
            }
```

**return** a(|x|, |y|, |z|, |t|)

# Vérification d'une analogie

- ▶ trop long ...
- ▶ vérifier d'abord ( $\sim 35\%$  d'économie):

$$\begin{aligned} [x : y :: z : t] &\Rightarrow \\ &(x[1] \in \{y[1], z[1]\}) \vee (t[1] \in \{y[1], z[1]\}) \\ &(x[\$] \in \{y[\$], z[\$]\}) \vee (t[\$] \in \{y[\$], z[\$]\}) \end{aligned}$$

# Espaces d'entrée très grands

- ▶  $0.02 \text{ ms} \times 10\text{M} = 200 \text{ seconds}$   
(sans même compter la vérification d'analogie)  
⇒ échantillonnage des  $x$ -forms ...
  
- ▶ Step 1:
  1. **sample**  $x$ -forms  $\in \mathcal{I}$
  2. search for all the pairs  $(y, z)$  vérifiant la propriété des comptes
  3. vérifier les **véritables analogies**  
(algorithme en  $o(|x| \times |y| \times |z| \times |t|)$ )

# Espaces d'entrée très grands

- ▶ Échantillonnage des  $x$ -forms

- ▶ [une pomme verte : des pommes vertes :: une voiture rouge : des voitures rouges]

- ▶ [dream : dreamer :: dreams : dreamers]

- ▶ [This guy drinks too much : This boat sinks :: These guys drank too much : These boats sank]

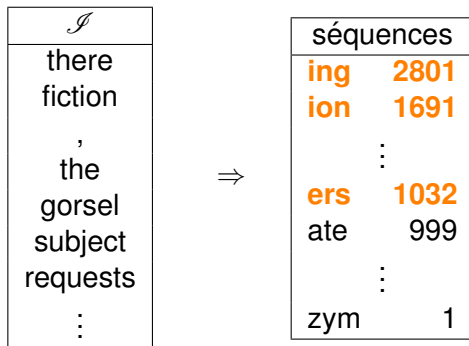
- ▶ l'approche vectorielle semble adaptée

# Vector space sampling

$\mathcal{I}$
there
fiction
,
the
gorsel
subject
requests
⋮



# Vector space sampling



# Vector space sampling

$\mathcal{I}$
there
fiction
,
the
gorsel
subject
requests
⋮

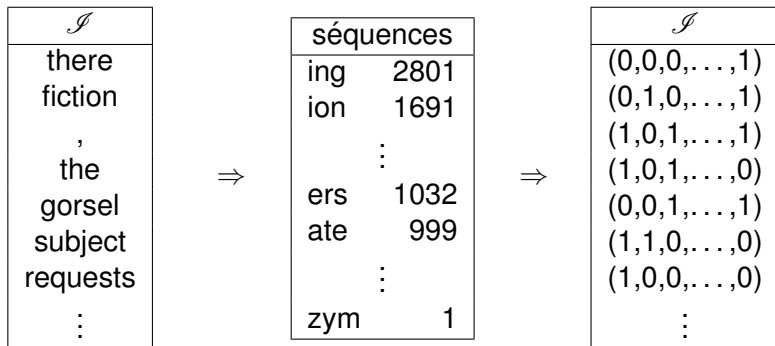
⇒

séquences	
ing	2801
ion	1691
	⋮
ers	1032
ate	999
	⋮
zym	1

⇒

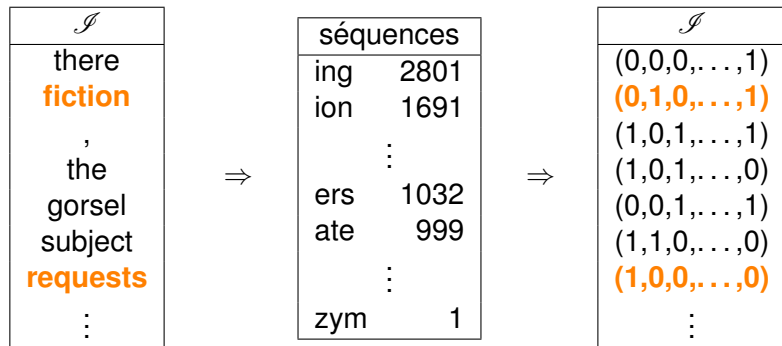
$\mathcal{I}$
(0,0,0,...,1)
(0,1,0,...,1)
(1,0,1,...,1)
(1,0,1,...,0)
(0,0,1,...,1)
(1,1,0,...,0)
(1,0,0,...,0)
⋮

# Vector space sampling



- ▶  $t \equiv$  question

# Vector space sampling

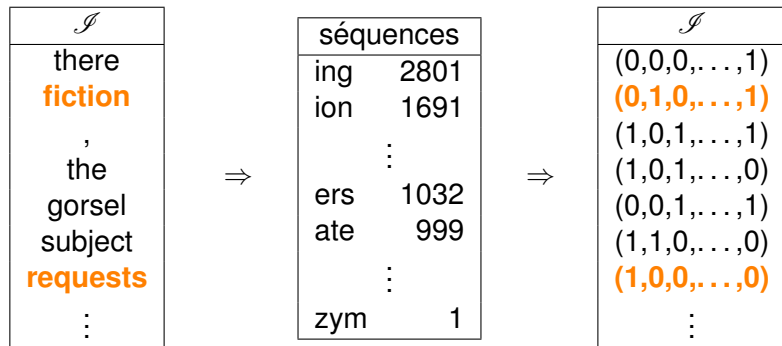


▶  $t \equiv (0,1,0,\dots,0)$

▶ **similarité cosinus** souvent utilisée :  $\text{similarité}(A,B) = \cos(\theta) =$

$$\frac{A \cdot B}{\|A\| \cdot \|B\|} = \frac{\sum_i A_i \times B_i}{\sqrt{\sum_i A_i^2} \times \sqrt{\sum_i B_i^2}}$$

# Vector space sampling



▶  $t \equiv (0,1,0,\dots,0) \Rightarrow$  request,requests, query, quest, fiction, ...

▶ **similarité cosinus** souvent utilisée :  $\text{similarité}(A,B) = \cos(\theta) =$

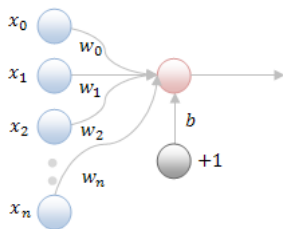
$$\frac{A \cdot B}{\|A\| \cdot \|B\|} = \frac{\sum_i A_i \times B_i}{\sqrt{\sum_i A_i^2} \times \sqrt{\sum_i B_i^2}}$$

## Gérer le bruit (step 3)

- ▶ Step 2 peut facilement générer des **millions** de formes candidates . . .
- ▶ Plusieurs stratégies proposées:
  - ▶ filtrer par la fréquence [Lepage & Denoual, 2005]
  - ▶ éliminer les formes non vues dans un grand corpus de formes cibles [Langlais & Patry, 2007]
  - ▶ éliminer les formes contenant un n-gramme de caractère non vu à l'entraînement [Lepage, 2007]
  - ▶ apprendre à reconnaître les analogies "utiles" [Langlais et al., 2008]

# Classification des solutions (bonne/mauvaise)

Perceptron [Rosenblatt, 1958]



- ▶  $x_j$  input (valeur réelle)
- ▶  $w_j$  poids synaptique (valeur réelle)
- ▶  $b$  biais (peut être représenté par un poids supplémentaire et une entrée fixe)
- ▶  $h = \sum_j x_j \times w_j + b$  est le signal arrivant au neurone
- ▶ l'activation du neurone est déterminée par une fonction d'activation  $t \equiv f(h)$  (ex:  $f = \text{sign}$ )

- ▶ soit  $(x_i, y_i)_{i \in [1, L]}$  un corpus d'entraînement où  $y_i$  (ici à valeur dans  $\{-1, +1\}$ ) est la réponse (supervision) et  $x_i \in \mathbb{R}^{n+1}$
- ▶ apprendre les  $w_j$  peut se faire **en ligne** en ajustant itérativement les poids une fois chaque observation  $x_i$  rencontrée:  
$$\forall j \in [0, n], w_j \leftarrow w_j + \eta (y_i - t_i) \cdot x_j$$
où  $\eta$  est le **learning rate**

résout des problèmes linéairement séparables

# Classification des solutions (bonne/mauvaise)

## Voted-perceptron [Freund:99]

- ▶ Entraînement en ligne:

**Require:**  $\{(\mathbf{x}_i, y_i)\}_{i \in [1, L]}$  où  $y_i \in \{+1, -1\}$  et  $\mathbf{x}_i \in \mathbb{R}^n$

**Ensure:** a pool of  $K$  perceptrons  $\{(\mathbf{v}_k, c_k)\}_{k \in [1, K]}$ ,  $\mathbf{v}_k \in \mathbb{R}^n$ ,  $c_k$   
entier

$k, c_1 \leftarrow 0$

$\mathbf{v}_1 \leftarrow [0 \dots 0]^T$

**for all epoch do**

**for all**  $i \in [1, L]$  **do**

$\hat{y} \leftarrow \text{sign}(\mathbf{v}_k \cdot \mathbf{x}_i)$

**if**  $\hat{y} \neq y_i$  **then**

$\mathbf{v}_{k+1} \leftarrow \mathbf{v}_k + y_i \mathbf{x}_i$

$c_{k+1} \leftarrow 1$

$k \leftarrow k + 1$

**else**

$c_k \leftarrow c_k + 1$

Garanties de convergence (même dans les cas non linéairement séparables)



# Classification des solutions (bonne/mauvaise)

Voted-perceptron [Freund:99]

- ▶ Test:

$$\hat{y} = \text{sign} \left( \sum_k c_k \cdot \text{sign}(\mathbf{v}_k \cdot \mathbf{x}) \right)$$

- ▶ requiert la sauvegarde de  $K$  perceptrons et de leur coefficient
- ▶ calcul en  $O(K \times n)$ . On peut également ne retenir qu'un sous-ensemble des perceptrons (à l'extrême un seul, par exemple celui de plus fort  $c_k$ ).

# Traits

- ▶ soit une paire d'analogies (source,cible) associant une forme  $t$  à sa solution  $s$ :

$$(a \equiv [I(x) : I(y) :: I(z) : t], \hat{a} \equiv [O(x) : O(y) :: O(z) : s])$$

- ▶ exemples de traits:
  - ▶ degré des analogies source ( $a$ ) et cible ( $\hat{a}$ )

[miracle : miraculeux :: fable : fabuleux] car:

x	mi	rac	$\epsilon$	le	$\epsilon$
y	mi	rac	u	le	ux
z	$\epsilon$	fab	$\epsilon$	le	$\epsilon$
t	$\epsilon$	fab	u	le	ux

# Traits

- ▶ soit une paire d'analogies (source,cible) associant une forme  $t$  à sa solution  $s$ :

$$(a \equiv [I(x) : I(y) :: I(z) : t], \hat{a} \equiv [O(x) : O(y) :: O(z) : s])$$

- ▶ exemples de traits:
  - ▶ **degré** des analogies source ( $a$ ) et cible ( $\hat{a}$ )

[miracle : miraculeux :: fable : fabuleux] mais aussi:

x	mirac	le
y	mirac	uleux
z	fab	le
t	fab	uleux

---

degré 2 (nombre minimum de facteurs)

# Traits

- ▶ soit une paire d'analogies (source,cible) associant une forme  $t$  à sa solution  $s$ :

$$(a \equiv [I(x) : I(y) :: I(z) : t], \hat{a} \equiv [O(x) : O(y) :: O(z) : s])$$

- ▶ exemples de traits:
  - ▶ degré des analogies source ( $a$ ) et cible ( $\hat{a}$ )

[miracle : miraculeux :: fable : fabuleux] mais aussi:

x	mirac	le
y	mirac	uleux
z	fab	le
t	fab	uleux
<hr/>		
cofacteurs	(mirac,fab)	(le,uleux)

# Traits

- ▶ soit une paire d'analogies (source,cible) associant une forme  $t$  à sa solution  $s$ :

$$(a \equiv [I(x) : I(y) :: I(z) : t], \hat{a} \equiv [O(x) : O(y) :: O(z) : s])$$

- ▶ exemples de traits:
  - ▶ degré des analogies source ( $a$ ) et cible ( $\hat{a}$ )

[miracle : miraculeux :: fable : ueufaxbl] car

x	mirac	$\epsilon$	l	$\epsilon$	$\epsilon$	$\epsilon$	$\epsilon$	e
y	mirac	u	l	eu	$\epsilon$	x	$\epsilon$	$\epsilon$
z	$\epsilon$	$\epsilon$	$\epsilon$	$\epsilon$	fa	$\epsilon$	bl	e
t	$\epsilon$	u	$\epsilon$	eu	fa	x	bl	$\epsilon$

# Traits

- ▶ soit une paire d'analogies (source,cible) associant une forme  $t$  à sa solution  $s$ :

$$(a \equiv [I(x) : I(y) :: I(z) : t], \hat{a} \equiv [O(x) : O(y) :: O(z) : s])$$

- ▶ exemples de traits:
  - ▶ degré des analogies source ( $a$ ) et cible ( $\hat{a}$ )

	[miracle : miraculeux :: fable : ueufaxbl] car							
x	mirac	ε	l	ε	ε	ε	ε	e
y	mirac	u	l	eu	ε	x	ε	ε
z	ε	ε	ε	ε	fa	ε	bl	e
t	ε	u	ε	eu	fa	x	bl	ε

- ▶ fréquence de  $s$ , son rang
- ▶ probabilité minimale, maximale, moyenne (etc.) de  $s$  selon un modèle n-gramme au niveau des caractères
- ▶ traits sur la présence de cofacteurs,
- ▶ etc.

## Apprentissage Analogique

Analogie

Principe

## Quelques réalisations

### Sous le capot

Définitions de l'analogie (formelle)

Solveurs d'équations

Recherche des analogies

Filtrer le bruit

## Applications

Traduction de mots inconnus

Traduction de termes

Translittération

# Traduction de mots inconnus

Ressource: corpus WMT'06 [Koehn et Monz, 2006]

UNKS dans le corpus de test (French → English):

- ▶ 20% de noms propres
- ▶ 12% d'expressions numériques (années, pages, etc.)
- ▶ 8% de mots composés
- ▶ 4% de mots d'emprunt (latin, grec)

~ **50% de mots *normaux***



# Protocole

▶  $\mathcal{L}_T$ : **lexique souche**

*Modèles IBM 2 ( $s2t \cap t2s$ ) entraînés sur les  $T$  premières paires de phrases de TRAIN.*

↪  $T$ : 5 000, 10 000, 100 000, 200 000, and 500 000

▶  $\mathcal{L}_{ref}$ : **lexique de référence**

*Modèles IBM 2 ( $s2t \cap t2s$ ) entraînés sur TRAIN.*

▶ Entrées dans un modèle IBM (anglais/roumain,  $p > 0.01$ ):

previous	[1162]
	<i>al</i> (0.155) <i>anterioare</i> (0.06113) <i>precedente</i> (0.03609) <i>se</i> (0.02009) <i>domnule</i> (0.01758) <i>cei</i> (0.01053) <i>anterior</i> (0.01005)
previously	[763]
	<i>persoane</i> (0.0216) <i>anterior</i> (0.01674) <i>puternic</i> (0.0158) <i>form</i> (0.01496) <i>deputai</i> (0.01341) <i>adecvat</i> (0.01315) <i>drepturilor</i> (0.01305) <i>sau</i> (0.01186) <i>la</i> (0.01115) <i>lor</i> (0.01087) <i>zona</i> (0.01042) <i>de-a</i> (0.01037) <i>suplimentare</i> (0.01003)

# Protocole

- ▶ Traduire les mots de TEST ne contenant pas de chiffre et présents dans  $L_{ref}$ , **mais** inconnus de  $L_T$ .
  
- ▶ Deux mesures:
  - rappel** % de mots UNKS avec une traduction correcte selon  $L_{ref}$
  - précision** % nombre de mots avec une traduction correcte selon  $L_{ref}$

# Approches de base (baselines)

$X$   $\equiv$  mot inconnu

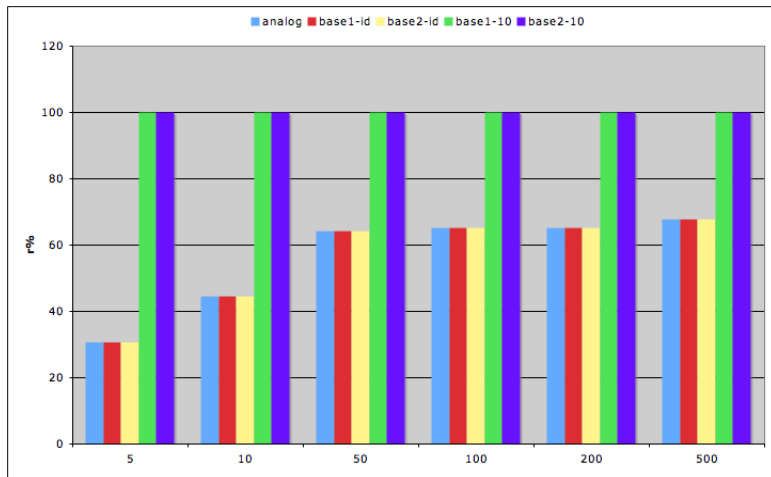
$$\mathbf{B1} \quad \hat{T} = \operatorname{argmin}_{T \in \mathcal{O}} \operatorname{edit-dist}(T, X)$$

*signalaient*  $\rightarrow$  *signalling, signalled, salient, ...*

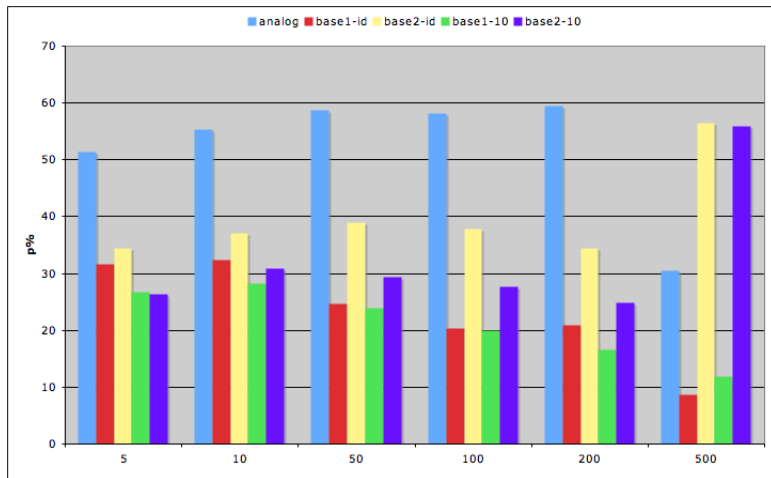
$$\mathbf{B2} \quad \hat{T} = \operatorname{argmax}_{S \in V(X)} \operatorname{proj}_{\mathcal{L}}(S) \text{ où} \\ V(X) = \operatorname{argmin}_{S \in \mathcal{S}} \operatorname{edit-dist}(X, S)$$

*signalaient*  $\rightarrow$  *signalement, signeraient, signalent, ...*  
 $\rightarrow$  *indicates, sign, signalling, ...*

# Recall (French $\rightarrow$ English)



# Précision (French → English)



# Cas de $T = 500\,000$ (French $\rightarrow$ English)

- ▶ Only a few words seen in  $\mathcal{L}_{ref}$ , but not in  $\mathcal{L}_{500\,000}$ 
  - ▶ 34 words (in-domain), 1.2 translations sanctioned by  $\mathcal{L}_{ref}$
  - ▶  $\sim 1/3$  of the reference translations are wrong

- ▶ exemples:

- ▶ **perquisitions:**

A	<i>seizures</i> (76) <i>seizure</i> (17) <i>raid</i> (1)
$\mathcal{L}_{ref}$	house-to-house

- ▶ **non-discriminatoire:**

A	<i>non-discriminatory</i> (72) <i>non-discrimination</i> (47) <i>nondiscrimination</i> (24) <i>nondiscriminatory</i> (23) <i>outlawing</i> (20) <i>race</i> (20) <i>prosperity</i> (20) <i>antidiscrimination</i> (20) <i>anti-discrimination</i> (18) <i>discrimination</i> (17)
$\mathcal{L}_{ref}$	affirmative

# Évaluation manuelle

contournant	(49 candidats)
A	◇ <b>(circumventing,55)</b> (undermining,20) (evading,19) (circumvented,17) (overturning,16) (circumvent,15) (circumvention,15) (bypass,13) (evade,13) (skirt,12) . . .
$\mathcal{L}_{ref}$	◇ <b>skirting, bypassing</b> , by-pass, overcoming

- ▶ 80% des **mots ordinaires** reçoivent une traduction valide en première position par analogie
- ▶ 35% dans le cas de B2

# Impact en traduction (vers l'anglais)

- ▶ Ajout dans le modèle de traduction de la première traduction produite par analogie pour chaque mot inconnu (lexique souche =  $L_{ref}$ )
- ▶ Évalué sur les phrases contenant **au moins** un mot inconnu

	French		Spanish		German	
	WER	BLEU	WER	BLEU	WER	BLEU
base	61.8	22.74	54.0	27.00	69.9	18.15
+B2	61.8	22.72	54.2	26.89	70.3	18.05
+A	<b>61.6</b>	<b>22.90</b>	<b>53.7</b>	<b>27.27</b>	<b>69.7</b>	<b>18.30</b>
sentences	387		452		814	

↪ améliorations stables mais non-significatives



# Traduction de segments inconnus

French → English

- ▶  $\mathcal{L}_{ref}, \mathcal{L}_T$ : phrase table

<b>expulsent</b> ◇ expelling (36) expel (31) are expelling (23) are expel (10)
--

<b>focaliserai</b> ◇ focus (10) focus solely (9) concentrate all (9) (will focus (9) will placing (9)
---

<b>dépasseront</b> ◇ will exceed (4) (exceed (3) will be exceed (3) we go beyond (2) will be exceeding (2)
--

<b>non-réussite de</b> ◇ lack of success for (4) lack of success of (4) lack of success (4)
---

que vous <b>subissez</b> ◇ you are experiencing (2)
---

- ▶ précision: 55%
- ▶ rappel: 10% (trop de filtrage pendant step-1)

# Enrichissement d'une *phrase table*

- ▶ TRAIN: table de segments entraînée sur WMT'06 (fr2en)

```
comprehensive safety ||| sécurité globale ||| 0.111111 0.000569718
comprehensive security ||| sécurité globale ||| 0.222222 0.000789918
general safety ||| sécurité globale ||| 0.111111 0.000161561
global safety ||| sécurité globale ||| 0.111111 0.00126761
global security ||| sécurité globale ||| 0.444444 0.00175755
advocate a modern comprehensive security ||| favorables à un concept de sécurité globale moderne |
```

- ▶ Sélecteur (step-3): aucun

# Exemples

- ▶ **âgées à leur sort. [1079]**  
*old to die . (57) old on their own . (56) old in the lurch . (53) old to their fate . (41) very old to die . (35) ...*
- ▶ **' acquis soient transposées [3610]**  
*acquired are transposed (47) acquired be transposed with (38) acquired will be transposed (37) ...*
- ▶ **a caractérisé la réunification allemande [3655]**  
*has characterised of german reunification, (24) has characterised german reunification (20) ...*
- ▶ **acceptables , sans mettre en [9985]**  
*acceptable without calling into (23) acceptable, without calling into (21) ...*
- ▶ **a été discutée et [406223]**  
*were debated and (151) was discussed this and (133) was discussed thi and (123) has been discussed and has (119)...*

# Traduction de termes médicaux

- ▶ but: comparer l'apprentissage analogique à la traduction statistique par segments (Koehn et al., 2003) où l'unité de base est le caractère (sinon c'est trop mauvais)
- ▶ paires de langues variées: RU-EN, FI-EN, SW-EN, ES-EN, FR-EN
- ▶ un besoin réel

# Corpora

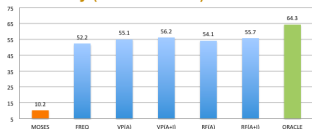
	TRAIN			TEST	DEV
	NB	AVG. LENGTH	NB	OOV	OOV
MESH (Medical Subject Headings thesaurus)					
FI	19787	19.3	1000	65.0	63.8
FR	17230	21.5	1000	35.8	36.8
RU	21407	38.5	1000	42.3	45.1
ES	19201	21.5	1000	37.4	34.9
SW	17090	17.3	1000	69.3	70.0
MEDDRA (Medical Drug Regulatory Activities thesaurus)					
ES	65276	34.6	1000	7.1	7.1

# Résultats

## Accuracy (MESH EN-FI)



## Accuracy (MEDDRA ES-EN)



- ATRANS  $\geq$  PB-SMT
  - ✓ Similar tendencies for other directions
  - ✓ Cascading ATRANS and PB-SMT improves drastically over each system
- Rescoring (RF) at a small advantage
- ATRANS  $\gg$  Moses
  - ➔ Similar observation for EN-ES
  - ➔ ATRANS (much) less silent
- Rescoring (RF) > classification (VP)

# Exemples

SW aikakauslehdet aiheena

REF periodicals as topic

ANA periodicals as topic

SMT timenancylages, topic

SW alfasalpaajat

REF adrenergic alpha-antagonists

ANA adrenergic alpha-antagonists

SMT alphablockers

SP instituciones de atención ambulatoria

REF ambulatory care facilities

ANA ambulatory care facilities

SMT institutions, attention ambulatory

FI märkivä kilpirauhastulehdus

REF thyroiditis, suppurative

ANA thyroiditis suppurativa

SMT rativa thyroid glandorum

FR malformations de la machoïre

REF jaw abnormalities

ANA jaw congenital abnormalities

SMT malformations jawory

FI rasva-alkoholit

REF fatty alcohols

ANA lipid alcohols

SMT fatty-alcohols

# Transliteration

NEWS 2009 EN-CH [NEWS09]

	TRAIN	DEV	TEST
examples	31 961	2 896	2 896
symbols	26	26	26
symbols	370	275	283

examples	
Emission	埃米申
Blagrove	布格夫
Aposhian	阿波希安

NEWS 2009 evaluation script:

- ▶ **ACC: accuracy of the first solution**
- ▶  $F_1$ : partial credit proportional to the longest subsequence between the 1st candidate and the reference
- ▶ MRR: Mean Reciprocal Rank



# System Tested

- ▶ Pure analogical devices:
  - ▶ ana-freq  $a_1$   
solutions ranked by frequency
  - ▶ ana-svm<sub>A</sub>  $a_2$   
solutions ranked by an SVM trained on analogical-based features
- ▶ Moses  $m$   
default setting (word  $\equiv$  character) ?
- ▶ Hybrids:
  - ▶ ana-svm<sub>A+M</sub>  $am_1$   
solutions ranked by an SVM trained on analogical **and** moses-based features
  - ▶  $am_2$  : cascading  $am_1$  and  $m$   
trust  $AM_1$  whenever a solution is produced,  $m$  otherwise

# Results

Configuration		ACC	$F_1$	MRR	<i>rank</i>
$a_1$	ana-freq	<b>56.6</b>	79.1	63.0	16
$a_2$	ana-svm <sub>A</sub>	<b>58.0</b>	80.0	58.8	15
m	moses	<b>66.6</b>	85.9	66.6	6
$am_1$	ana-svm <sub>A+M</sub>	<b>63.4</b>	82.0	64.1	10
$am_2$	casc( $am_1, m$ )	<b>68.5</b>	86.9	69.0	4
	last NEWS 2009	<b>19.9</b>	60.6	22.9	23
	first NEWS 2009	<b>73.1</b>	89.5	81.2	1

1- None of the systems tested outperform the best system at NEWS 2009

# Results

Configuration		ACC	F <sub>1</sub>	MRR	rank
a <sub>1</sub>	ana-freq	<b>56.6</b>	79.1	63.0	16
a <sub>2</sub>	ana-svm <sub>A</sub>	<b>58.0</b>	80.0	58.8	15
m	moses	<b>66.6</b>	85.9	66.6	6
am <sub>1</sub>	ana-svm <sub>A+M</sub>	<b>63.4</b>	82.0	64.1	10
am <sub>2</sub>	casc(am <sub>1</sub> ,m)	<b>68.5</b>	86.9	69.0	4
	last NEWS 2009	<b>19.9</b>	60.6	22.9	23
	first NEWS 2009	<b>73.1</b>	89.5	81.2	1

2- The best system we tested would have ranked 4th

# Results

Configuration		ACC	F <sub>1</sub>	MRR	rank
a <sub>1</sub>	ana-freq	<b>56.6</b>	79.1	63.0	16
a <sub>2</sub>	ana-svm <sub>A</sub>	<b>58.0</b>	80.0	58.8	15
m	moses	<b>66.6</b>	85.9	66.6	6
am <sub>1</sub>	ana-svm <sub>A+M</sub>	<b>63.4</b>	82.0	64.1	10
am <sub>2</sub>	casc(am <sub>1</sub> ,m)	<b>68.5</b>	86.9	69.0	4
	last NEWS 2009	<b>19.9</b>	60.6	22.9	23
	first NEWS 2009	<b>73.1</b>	89.5	81.2	1

3- Learning to recognize good solutions is preferable to ranking by frequency

# Results

Configuration		ACC	F <sub>1</sub>	MRR	rank
a <sub>1</sub>	ana-freq	<b>56.6</b>	79.1	63.0	16
a <sub>2</sub>	ana-svm <sub>A</sub>	<b>58.0</b>	80.0	58.8	15
m	moses	<b>66.6</b>	85.9	66.6	6
am <sub>1</sub>	ana-svm <sub>A+M</sub>	<b>63.4</b>	82.0	64.1	10
am <sub>2</sub>	casc(am <sub>1</sub> ,m)	<b>68.5</b>	86.9	69.0	4
	last NEWS 2009	<b>19.9</b>	60.6	22.9	23
	first NEWS 2009	<b>73.1</b>	89.5	81.2	1

4- Moses outperforms pure analogical devices  
(silence rate of analogical devices: 3.7%)

# Results

Configuration		ACC	F <sub>1</sub>	MRR	rank
a <sub>1</sub>	ana-freq	<b>56.6</b>	79.1	63.0	16
a <sub>2</sub>	ana-svm <sub>A</sub>	<b>58.0</b>	80.0	58.8	15
m	moses	<b>66.6</b>	85.9	66.6	6
am <sub>1</sub>	ana-svm <sub>A+M</sub>	<b>63.4</b>	82.0	64.1	10
am <sub>2</sub>	casc(am <sub>1</sub> ,m)	<b>68.5</b>	86.9	69.0	4
	last NEWS 2009	<b>19.9</b>	60.6	22.9	23
	first NEWS 2009	<b>73.1</b>	89.5	81.2	1

5- Adding a few Moses-based features helps classification a lot

# Results

Configuration		ACC	F <sub>1</sub>	MRR	rank
a <sub>1</sub>	ana-freq	<b>56.6</b>	79.1	63.0	16
a <sub>2</sub>	ana-svm <sub>A</sub>	<b>58.0</b>	80.0	58.8	15
m	moses	<b>66.6</b>	85.9	66.6	6
am <sub>1</sub>	ana-svm <sub>A+M</sub>	<b>63.4</b>	82.0	64.1	10
am <sub>2</sub>	casc(am <sub>1</sub> ,m)	<b>68.5</b>	86.9	69.0	4
	last NEWS 2009	<b>19.9</b>	60.6	22.9	23
	first NEWS 2009	<b>73.1</b>	89.5	81.2	1

6- Higher precision of the analogical device over Moses