

# Mesurer la distance morphologique

Nabil HATHOUT, Basilio CALDERONE, Fabio MONTERMINI  
CLLE-ERSS, CNRS & UTM

Séminaire de l'UE TAL

12 décembre 2011

- Qu'est ce qu'une mesure de la distance morphologique ?
- À quoi ça peut servir ?
- Comment ça peut se calculer ?
- Comment comparer / évaluer ces mesures ?

# Distance morphologique

- Une mesure de la distance morphologique est une fonction qui associe à chaque couple de mots du lexique une valeur numérique.
- Deux mots sont en relation morphologique s'ils partagent en même temps une partie de leur forme et de leur sens.

- Il ne s'agit pas de distance au sens mathématique du terme (symétrie, séparation, inégalité triangulaire).
- Si  $w_1$  et  $w_2$  sont en relation morphologique,  $D(w_1, w_2) > 0$
- Si  $w_1$  et  $w_2$  n'ont aucune relation morphologique,  $D(w_1, w_2) = 0$
- Si  $w_1$  est construit à partir de  $w_2$  et  $w_2$  est construit à partir de  $w_3$ , alors  $D(w_1, w_2) < D(w_1, w_3)$

$$D(\text{incorrection\_N}, \text{incorrect\_A}) < D(\text{incorrection\_N}, \text{correct\_A})$$

- 1 Aspects linguistiques
- 2 Création d'un voisinage d'évaluation
- 3 Une morphologie sans sémantique
- 4 Distance d'édition de Levenshtein
- 5 Phonotactic Activation System (PHACTS)
- 6 Proxinettes
- 7 Entropie maximale
- 8 Comparaison des mesures

- Le lexique est une entité structurée.
- Chaque élément lexical est inséré dans un environnement (un voisinage).
- Par exemple, il fait partie de familles et de séries morphologiques.
- Il est possible de quantifier la ressemblance des membres de ces ensembles.

- En théorie, il serait possible de mesurer la relation de n'importe quelle unité du lexique avec n'importe quelle autre unité.
- En pratique, seules certaines relations sont pertinentes / intéressantes.
  - ▶ *plage / sable*
  - ▶ *admirable / sable*
  - ▶ *admirable / éligible*
  - ▶ *admirable / admiration*
  - ▶ *éligible / élection*

**Distance formelle (phonologique)** : en fonction des segments (phonèmes) partagés et de leur distribution.

**Distance sémantique** : en fonction des traits sémantiques partagés, des propriétés distributionnelles.

**Distance morphologique** : en fonction de la densité de l'environnement lexical, du nombre de voisins (des membres de la famille ou des sous-séries) partagés.

- Dans plusieurs domaines (phonologie, morphologie, syntaxe) on se sert de distinctions qui font référence à la distance entre deux mots :
  - ▶ **allomorphie** vs. **supplétion**
  - ▶ **transparent** vs. **non transparent**
- En général, on fait comme s'il s'agissait de distinctions qualitatives et qu'il était possible d'étendre à tout le lexique ce qu'on observe localement pour des couples de mots.



- Donner une métrique quantitative à des différences supposées qualitatives.
- On peut continuer de se servir des mêmes catégories mais en se basant sur des mesures objectives.

Quelques applications :

- Linguistique théorique (phonologie, morphologie, syntaxe, etc.) :
  - ▶ allomorphie, transparence, cas de déviation forme / sens (*anticancéreux*)
- Psycholinguistique :
  - ▶ transparence, taille / densité de la famille lexicale
- TAL :
  - ▶ les voisinages permettent de tenir compte des variations morphologiques dans les systèmes de RI ;
  - ▶ dégrossissage des relations morphologiques dans un lexique ;
  - ▶ réduire l'espace de recherche lorsqu'on cherche à connecter les mots qui sont morphologiquement apparentés.

# Vue d'ensemble

- 1 Aspects linguistiques
- 2 Création d'un voisinage d'évaluation**
- 3 Une morphologie sans sémantique
- 4 Distance d'édition de Levenshtein
- 5 Phonotactic Activation System (PHACTS)
- 6 Proxinette
- 7 Entropie maximale
- 8 Comparaison des mesures

- Comparer différentes mesures qui permettent d'**estimer la parenté morphologique** qui existe entre les mots du dictionnaire.
- La comparaison est réalisée relativement à un lexique de référence : **la partie anglaise de la base CELEX.**
- CELEX est la seule base de données morphologiques dérivationnelles largement diffusée et utilisée (anglais, allemand, néerlandais).

- CELEX contient des informations
  - ▶ phonologiques et graphémiques ;
  - ▶ **morphologiques** (flexionnelles et **dérivationnelles**) ;
  - ▶ catégorielles ;
  - ▶ fréquences par lemme et par forme.

## Quelques informations fournies par CELEX

acoustically/B/6/C/SA/((acoustic)[A],(ally)[B|A.])[B]/**derivation**

applicable/A/81/C/SA/((apply)[V],(able)[A|V.])[A]/**derivation**

applique/N/1/Z/S/((applique)[V])[N]/**conversion**

applejack/N/1/C/SS/((apple)[N],(jack)[N])[N]/**composition**

aerodrome/N/10/C/AA/((aero)[N|.x],(drome)[N|x.])[N]/**neoclassique**

- Comment calculer des distances entre les mots à partir de leurs **décompositions morphématisques** ?

# Graphe CELEX

- Les décompositions morphématiques sont analysées en utilisant Lex et Yacc.
- Construction d'un graphe CELEX
  - ▶ Les sommets sont les décompositions des lexèmes et des éléments de composition.
  - ▶ Les arcs relient les décompositions des lexèmes complexes aux décompositions de leurs constituants. Elles sont symétriques.

## Exemples d'arcs dans le graphe CELEX

acoustically \_B ↔ acoustic \_A

apologize \_V ↔ apology \_N

applicable \_A ↔ apply \_V

applique \_N ↔ applique \_V

applejack \_N ↔ apple \_N ; applejack \_N ↔ jack \_N

aerodrome \_N ↔ aero \_N.x ; aerodrome \_N ↔ drome \_N.x

# Familles dérivationnelles dans CELEX

- Le graphe est parcouru pour constituer des familles morphologiques.
- ① On calcule la fermeture symétrique et transitive des relations de dérivations et de conversion.
- ② Pas de fermeture transitive (relations transversales) pour les composés.
  - ▶ Un mot composé est dans les familles dérivationnelles de tous ses éléments de composition.
  - ▶ Un composant appartient à la famille dérivationnelle de chaque composé dont il est l'un des éléments.

*aerodrome* et *airplane* n'appartiennent pas à la même famille

**aeroplane\_N** ↔ **aerodrome\_N**, **airplane\_N**

**airplane\_N** ↛ **aerodrome\_N**

**aerodrome\_N** ↔ **aero\_N.x**; **aerodrome\_N** ↔ **drome\_N.x**

**aeroplane\_N** ↔ **aero\_N.x**; **aeroplane\_N** ↔ **plane\_N**

**airplane\_N** ↔ **air\_N**; **airplane\_N** ↔ **plane\_N**

# Familles dérivationnelles dans CELEX

## Distance simple entre les membres d'une famille

Si  $w_1$  et  $w_2$  appartiennent à la même famille, on peut définir  $D_f(w_1, w_2)$  comme le nombre d'arcs du plus court chemin entre  $w_1$  et  $w_2$ .

## Distance des éléments de la famille de *determinant* $_N$

determinable\_A=2, determinant\_A=2, determinant\_N=0, determinate\_A=2,  
determination\_N=2, determinative\_A=3, determinative\_N=4, determine\_V=1,  
determiner\_N=2, determinism\_N=2, deterministic\_A=3, indeterminable\_A=3,  
indeterminably\_B=4, indeterminacy\_N=4, indeterminate\_A=3, predetermination\_N=3,  
predetermination\_N=3, predetermine\_V=2, predeterminer\_N=3

## Distance des éléments de la famille de *fruit* $_N$

breadfruit\_N=1, fruit\_N=0, fruit\_V=1, fruitcake\_N=1, fruiterer\_N=1, fruitful\_A=1,  
fruitfully\_B=2, fruitfulness\_N=2, fruition\_N=1, fruitless\_A=1, fruitlessly\_B=2,  
fruitlessness\_N=2, fruity\_A=1, grapefruit\_N=1, unfruitful\_A=2



# Séries dérivationnelles dans CELEX

- Un mot appartient généralement à **plusieurs sous-séries**.

*indeterminably\_B* appartient aux sous-séries des adverbes en :

*in-V-ably* : incurably\_B, indisputably\_B, inexplanably\_B

*in-A-ly* : inhumanely\_B, incautiously\_B, incompetently\_B

*A-ly* : disapprovingly\_B, sincerely\_B, categorically\_B

- Les sous-séries d'une entrée peuvent être calculées à partir de sa décomposition en réalisant des **abstractions** sur ses constituants.

**Signatures des sous-séries de *indeterminably\_B***

décomposition = (((in)[A|.A],((determine)[V],(able)[A|V.])[A])[A],(ly)[B|A.])[B]

-

- (((in)[A|.A],(**@**,(able)[A|V.])[A])[A],(ly)[B|A.])[B]
- (((in)[A|.A],**@**)[A],(ly)[B|A.])[B]
- (**@**,(ly)[B|A.])[B]

## Distance simple entre un mot et les éléments de ses sous-séries

- 1 Ordonner les sous-séries de  $w$  par ordre décroissant du nombre de symboles de leurs signatures.
- 2 Si  $\sigma$  est une sous-série de  $w$ ,  $D_s(w, \sigma)$  est le rang de  $\sigma$  parmi les sous-série de  $w$ .
- 3 Si  $w'$  appartient aux sous-séries  $\sigma_1, \dots, \sigma_n$  de  $w$ ,

$$D_s(w, w') = \min_{i=1 \dots n} D_s(w, \sigma_i)$$

Comment combiner les mesures de distance morphologiques dans les familles et les sous-séries ?

- Proposer une mesure unique pour les membres des familles et des sous-séries morphologiques.
- L'idée est d'estimer la distance morphologique entre deux mots en fonction du nombre d'analogies dans lesquels ils apparaissent :

## Distance analogique

- 1 Rechercher toutes les analogies morphologiques qui s'établissent entre les mots du lexique.
- 2 Définir une distance  $D_r$  telle que

$$D_r(w_1, w_2) = \frac{|\{(x, y) \in L \times L / w_1 : w_2 = x : y\}|}{N(w_1)}$$

- ▶ où  $L$  est l'ensemble des mots du lexique
- ▶  $a : b = c : d$  est vrai s'il existe une analogie formelle entre les décompositions morphématiques de  $a$ ,  $b$ ,  $c$  et  $d$
- ▶  $N(w)$  est le nombre total d'analogies dans lesquelles apparaît  $w$ .

# Analogies morphologiques structurelles

*academic : academically = psychedelic : psychedelically*

|   |                              |                   |
|---|------------------------------|-------------------|
| € | ((academy)[N],(ic)[A N.])[A] | €                 |
| ( | ((academy)[N],(ic)[A N.])[A] | ,(ally)[B A.])[B] |

|   |                  |                   |
|---|------------------|-------------------|
| € | (psychedelic)[A] | €                 |
| ( | (psychedelic)[A] | ,(ally)[B A.])[B] |

- Calculer toutes les analogies qui s'établissent entre toutes les décompositions de CELEX est difficile à réaliser car de trop grande complexité computationnelle.

# Analogies morphologiques structurelles

- On peut collecter un grand nombre de ces analogies en associant à chaque couple une signature qui décrit les opérations d'édition qui permettent de passer d'une décomposition à l'autre.

*academically* : *academic*

D/['(')/[] ; D/['(ally)', '[B|A.]', ')', '[B]']/[]

- On utilise la `diff1ib` (l'algorithme du `diff` de Unix) en faisant varier le découpage des décompositions.
- Nous avons calculé 31 millions d'analogies pour les 41 120 entrées de CELEX.
- Les analogies dont l'une des relations est une conversion sont exclues.

## 20 voisins les plus proches de *avoidable*\_A

avoid\_V :0.445826 unavoidable\_A :0.422735 unavoidably\_B :0.113677

avoidance\_N :0.017762 favourable\_A :0.008881 comfortable\_A :0.008881

utterable\_A :0.007105 speakable\_A :0.005329 serviceable\_A :0.005329 readable\_A :0.005329

questionable\_A :0.005329 mistakable\_A :0.005329 inhabitable\_A :0.005329

impeachable\_A :0.005329 governable\_A :0.005329 forgettable\_A :0.005329

fathomable\_A :0.005329 endurable\_A :0.005329 desirable\_A :0.005329 deniable\_A :0.005329

believable\_A :0.005329 bearable\_A :0.005329 accountable\_A :0.005329

- L'évaluation des mesures est faite en termes de **Précision**, **Rappel** et **Fscore** pour les tailles de voisinages de 5, 10, 20, 50, 100, 200, 500 et 1000 voisins. Le Fscore est la moyenne harmonique de la Précision et du Rappel.

# Vue d'ensemble

- 1 Aspects linguistiques
- 2 Création d'un voisinage d'évaluation
- 3 Une morphologie sans sémantique**
- 4 Distance d'édition de Levenshtein
- 5 Phonotactic Activation System (PHACTS)
- 6 Proxinettes
- 7 Entropie maximale
- 8 Comparaison des mesures

## Autonomie de la morphologie

Dans les présentes simulations, l'effort de quantifier une mesure de 'morphologicité' est basé **exclusivement sur les aspects de 'surface' des mots**, sans l'adjonction d'informations sur la composante sémantico-contextuelle des mots (*blind-to-semantics framework*)

## Que reste sans la sémantique ?

- Analogies et relations purement formelles entre les mots
- Effets de fréquence (type et token)
- Aspects phonotactiques de la langue
- Structures paradigmatiques



## Motivations

- **Evidences comportementales et psycholinguistiques :**
  - ▶ Existence de effets morphologiques en 'word recognition' qui ne sont pas influencés par la transparence sémantique. (early visual recognition : Taft, 1994 ; Taft, 1979 ; Caramazza et al., 1988 ; Schreuder and Baayen, 1995 ; Meunier and Segui, 2002 ; Rastle and Davis, 2004)
  - ▶ Acquisition de la morphologie et maîtrise de la compétence morphologique (Pinker et Ullman, 2002 ; Pirrelli et al. 2011)
- **Evidences linguistiques-descriptives :**
  - ▶ Structure et organization du lexique : analyses de type 'stem space' (Bonami et Boyé, 2003 ; Boyé et Montermini, 2011)
  - ▶ Structures paradigmatiques émergents (Giraud et al. in press ; Albright, 2009 ; Calderone et Celata, 2011)

## Attention

En affirmant l'autonomie de la morphologie pour certains structures/niveaux linguistiques, nous ne sommes pas en train de caractériser la morphologie en termes de modularité

# Vue d'ensemble

- 1 Aspects linguistiques
- 2 Création d'un voisinage d'évaluation
- 3 Une morphologie sans sémantique
- 4 Distance d'édition de Levenshtein**
- 5 Phonotactic Activation System (PHACTS)
- 6 Proxinette
- 7 Entropie maximale
- 8 Comparaison des mesures

## Notion basique

La **distance de Levenshtein** mesure la similarité entre deux chaînes (séquences)  $S$  et  $T$  de caractères. Elle est égale au nombre minimal de caractères qu'il faut *supprimer*, *insérer* ou *remplacer* pour passer d'une chaîne à l'autre.

- Mesure de surface  $\rightarrow$  *matching* de caractères (des techniques d'alignement)
- Plusieurs versions : normalisation par nb. de caractères, différents coûts de passage, etc.

# Distance de Levenshtein

## Matrice de Levenshtein

- Si  $|S| = m$  et  $|T| = n$  il suffit construire une matrice de taille  $m + 1 * n + 1$
- $D(0,0) = 0$ ,  $D(i,0) = i$  où  $i = 1 \dots m$ ,  $D(0,j) = j$  où  $j = 1 \dots n$
- Remplir la matrice avec une itération des coûts :
  - ▶ Si  $S(i) == T(j)$  coût = 0, si no coût = 1
  - ▶  $D(i,j) = \min(D(i-1,j-1)+\text{coût}, D(i,j-1)+1, D(i-1,j)+1)$

|   |    | m  | e  | i | l | e | n | s | t | e | i  | n  |
|---|----|----|----|---|---|---|---|---|---|---|----|----|
| l | 0  | 1  | 2  | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
| e | 1  | 1  | 2  | 3 | 3 | 4 | 5 | 6 | 7 | 8 | 9  | 10 |
| v | 2  | 2  | 1  | 2 | 3 | 3 | 4 | 5 | 6 | 7 | 8  | 9  |
| e | 3  | 3  | 2  | 2 | 3 | 4 | 4 | 5 | 6 | 7 | 8  | 9  |
| n | 4  | 4  | 3  | 3 | 3 | 3 | 4 | 5 | 6 | 6 | 7  | 8  |
| s | 5  | 5  | 4  | 4 | 4 | 4 | 3 | 4 | 5 | 6 | 7  | 7  |
| h | 6  | 6  | 5  | 5 | 5 | 5 | 4 | 3 | 4 | 5 | 6  | 7  |
| t | 7  | 7  | 6  | 6 | 6 | 6 | 5 | 4 | 4 | 5 | 6  | 7  |
| e | 8  | 8  | 7  | 7 | 7 | 7 | 6 | 5 | 4 | 5 | 6  | 7  |
| i | 9  | 9  | 8  | 8 | 8 | 7 | 7 | 6 | 5 | 4 | 5  | 6  |
| n | 10 | 10 | 9  | 8 | 9 | 8 | 8 | 7 | 6 | 5 | 4  | 5  |
|   | 11 | 11 | 10 | 9 | 9 | 9 | 8 | 8 | 7 | 6 | 5  | 4  |

# Distance de Levenshtein

## Voisins pour *direction*\_N

*directions*\_N=0.947368421053, *directional*\_A=0.9, *diction*\_N=0.875,  
*misdirection*\_N=0.857142857143, *dissection*\_N=0.842105263158,  
*discretion*\_N=0.842105263158, *reaction*\_N=0.823529411765, *erection*\_N=0.823529411765,  
*director*\_N=0.823529411765, *imprecation*\_N=0.8, *distraction*\_N=0.8, ***direct***\_V=0.8,  
***direct***\_B=0.8, ***direct***\_A=0.8, *diffraction*\_N=0.8, *dereliction*\_N=0.8, *deprecation*\_N=0.8,  
*rejection*\_N=0.777777777778, *reduction*\_N=0.777777777778, *redaction*\_N=0.777777777778,  
*reception*\_N=0.777777777778, *injection*\_N=0.777777777778, *infection*\_N=0.777777777778,  
***directory***\_N=0.777777777778, ***directive***\_N=0.777777777778, *directive*\_A=0.777777777778,  
*digestion*\_N=0.777777777778, *dictation*\_N=0.777777777778, *detection*\_N=0.777777777778,  
*depiction*\_N=0.777777777778, *dejection*\_N=0.777777777778, *defection*\_N=0.777777777778,  
*deduction*\_N=0.777777777778, *deception*\_N=0.777777777778, *bisection*\_N=0.777777777778,  
*addiction*\_N=0.777777777778, *intersection*\_N=0.761904761905,  
*interjection*\_N=0.761904761905, *insurrection*\_N=0.761904761905,  
*indiscretion*\_N=0.761904761905, *imperfection*\_N=0.761904761905,  
*disinfection*\_N=0.761904761905

## Réflexion

- Il s'agit d'une distance très bonne pour trouver les mots qui partagent un maximum de séquences de lettres communes
  - Mais elle ne prend pas en compte la position du caractère modifié dans la chaîne : le fait qu'une lettre soit ajoutée/supprimée/modifiée au début, au milieu ou à la fin du mot n'a pas d'importance.
- 
- Pour la présente simulation :
    - ▶ Calcul de la distance de Levenshtein pour les 30.000 mots de CELEX
    - ▶ Repérage de  $n$ -voisins pour chaque mot de CELEX

# Vue d'ensemble

- 1 Aspects linguistiques
- 2 Création d'un voisinage d'évaluation
- 3 Une morphologie sans sémantique
- 4 Distance d'édition de Levenshtein
- 5 Phonotactic Activation System (PHACTS)**
- 6 Proxinettes
- 7 Entropie maximale
- 8 Comparaison des mesures



## Phonotaxe

- C'est quoi la phonotaxe d'une langue ?
- À quoi ça peut servir ?
- Comment peut-on calculer une matrice phonotactique ?
- Comment peut-on 'construire' une représentation de mot basée sur la phonotaxe ?

- C'est quoi la phonotaxe d'une langue ?
  - ▶ Combinaisons de phonèmes (ou caractères) attestés et récurrents dans la langue

## En français

Le phonème /o/ peut être transcrit par plusieurs séquences de graphèmes : < o >, ... < au >, ... < eau >, ... < ot > etc

Souvent le graphème utilisé dépend de plusieurs variables :

- Position de /o/ dans le mot
- Son contexte consonantique

Par exemple, en position centrale, /o/ est plus souvent transcrit comme < o > que < au > entre < b > et < r >, mais est bien souvent transcrit comme < au > (en comparaison à < o >) entre < p > et < v >. En position final /o/ est fréquemment transcrit < eau > après < r > ou < t >, mais il n'est jamais transcrit comme < eau > après < f >.

# Construire mots à partir de la phonotaxe d'une langue

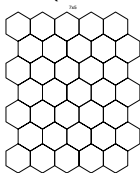
- Comment peut-on exploiter la connaissance phonotactique pour créer une représentation (psycholinguistiquement légitime) du mot (conçu comme une séquence temporelle d'unités phonologiques) ?
- Comment peut-on intégrer dans cette représentation les effets distributionnels de fréquence des séquences fortement attestées ( en termes de *type* et *token*) ?

**Inspiration** : cartes auto-organisatrices (Kohonen 2001), mémoires associatives (Hopfield 1982)

- Systèmes probabilistes d'organisation des données d'input
- Structure physique (en théorie) et computation en parallèle

# À la base de PHACTS : cartes adaptatives

- D'un point de vue architectural, les cartes adaptatives sont constituées d'une grille de noeuds ou neurones (le plus souvent uni- ou bidimensionnelle)
- Dans chaque noeud de la grille se trouve un neurone-vecteur (composants **random** au début, avant apprentissage)
- Chaque neurone est (en théorie) lié à un vecteur référent, responsable d'une position dans l'espace des données d'input
- L'idée de base est d'associer chaque donnée d'input à un (ou plusieurs) neurone sur la carte (après apprentissage)

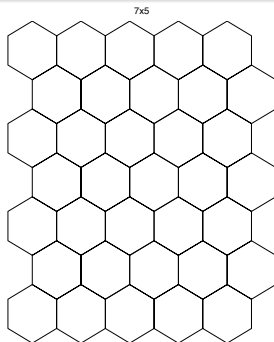


**NB** : Les données structurellement similaires doivent être représentées par des neurones adjacents sur la carte  $\Rightarrow$  *Clustering*

# À la base de PHACTS : cartes adaptatives

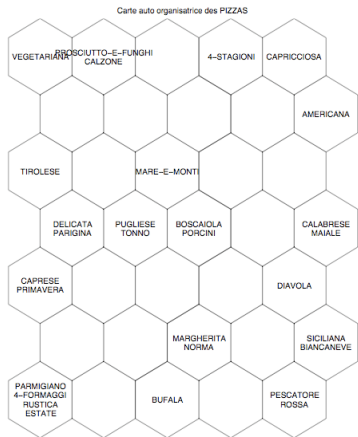
## Questions

- Mais quel type des données ?
- Mais quel type de clusters ?
- Mais quel type d'apprentissage ?
- Où est la phonotaxe ?





# Cartes adaptatives : un exemple italien



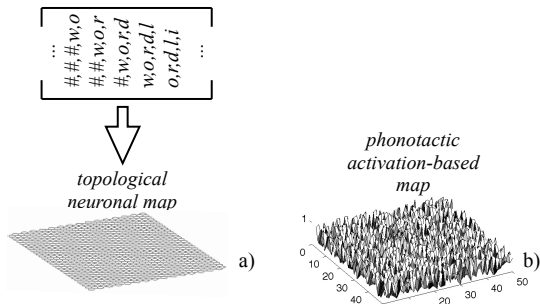
- Organisation topologique des pizzas
- En bas : pizzas *fromageuses*  
Au milieu : pizzas *charcuteries*  
En haut (gauche) : pizzas végétariennes
- La carte adaptative se déploie de façon à représenter un ensemble de données, et chaque neurone se spécialise pour représenter un groupe particulier des données selon les points communs qui les rassemblent
- Techniquement, la carte réalise une quantification vectorielle de l'espace des données



**Idée** : Obtenir une organisation topologique de la phonotaxe de la langue

- Pas ingrédients, mais *n-grams* transcrits de manière binaire (orthogonale) : chaque caractère est indépendant des autres
- Données : 30.000 mots de CELEX
- Dimension de la carte : 100x110
- Métrique d'apprentissage (vecteur-input vs vecteur-référent) : inner product

# PHACTS : l'architecture



- Importance de l'activation de chaque neurone : l'activation est liée à la fréquence de token du *n-gram* considéré
- Importance de la position de chaque neurone : la position est liée au type de *n-gram* considéré
- Importance de la position du *n-gram* dans les mots : les *n-grams* 'right-side' (fin de mot) sont les plus actifs

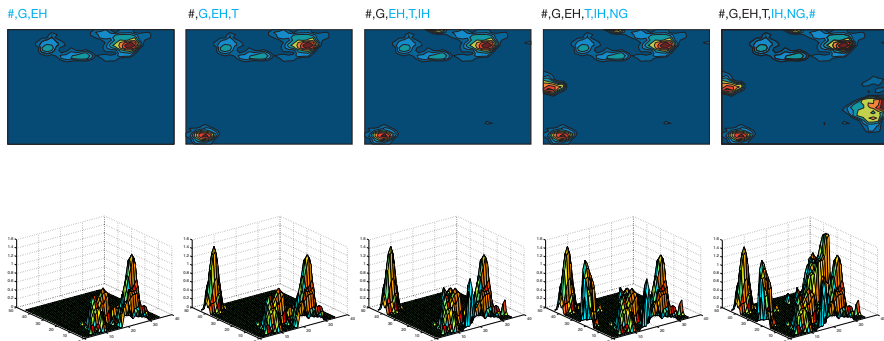
**Idée** : Reconstruire le mot à partir de cette carte phonotactique

- Reconstruire une dimension de mot en échantillonnant tous ses *n-grams* au moyen de la carte phonotactique
- La sommation des états d'activation déclenchés par des *n-grams* définit une représentation vectorielle où les mots sont recodés sur la base de la connaissance phonotactique préalablement acquise

$$F_{\text{phacts}}(\text{mot}) = \sum \Phi(\text{ngrams})$$

$\Phi$  = activations de la carte phonotactique

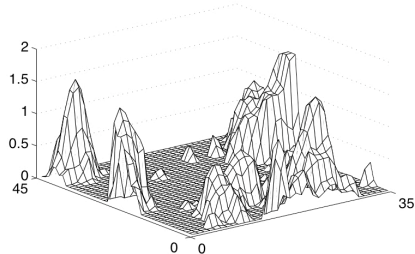
# PHACTS : l'architecture



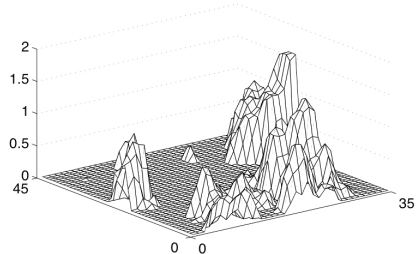
- Dans cette représentation il y a : effets de token (*n-grams* plus fréquents), effets de type (*n-grams* représentatifs de familles/séries de mots), importance de la position dans le mot (#\_ position initiale et \_# position finale)

# PHACTS : l'architecture

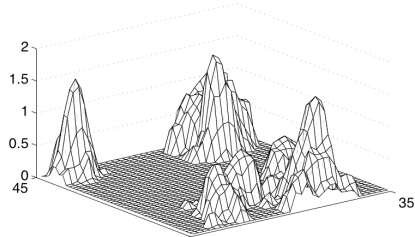
Accumulation output of the SOM for the word <getting>  
#G,EH,T,I,NG,# (ARPABET encoding)



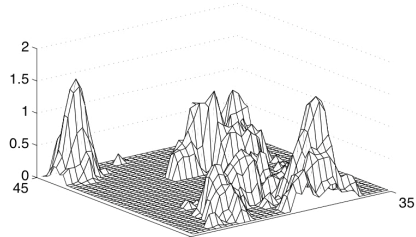
Accumulation output of the SOM for the word <being>  
#B,IY,IH,NG,# (ARPABET encoding)



Accumulation output of the SOM for the word <get>  
#G,EH,T,# (ARPABET encoding)



Accumulation output of the SOM for the word <gets>  
#G,EH,T,Z,# (ARPABET encoding)



# À PHACTS : distances morphologiques

## Voisins pour *direction* ( $1 - \cos(\text{mot}_1, \text{mot}_2)$ )

**directions** :0.0521546840 digestion :0.0583374380 vivisection :0.0664256310  
disaffection :0.0739004050 dissection :0.0805263970 discretion :0.0946790670  
**directional** :0.0977329500 rejection :0.1120468100 objection :0.1139618900  
defection :0.1144368500 diffraction :0.1167522700 erection :0.1199846300  
affection :0.1200441600 reflection :0.1204489500 deflection :0.1210436000  
election :0.1217490200 disruption :0.1222895200 recollection :0.1274858100  
correction :0.1275360100 reception :0.1319253200 misconception :0.1354825300  
indigestion :0.1355658600 deception :0.1356182300 differentiation :0.1376578000  
dissolution :0.1386147200 electrocution :0.1390508700 subjection :0.1391268500  
digression :0.1391603400 detection :0.1419080700 question :0.1428684900  
misrepresentation :0.1429052900 creation :0.1457034100 exception :0.1464719800  
disposition :0.1466483400 infection :0.1492515300 injection :0.1494845300  
resurrection :0.1498249400 inflection :0.1503713100 prediction :0.1511868000  
diminution :0.1543900400 accretion :0.1545539700 distraction :0.1565128300  
precaution :0.1567157400

# À PHACTS : distances morphologiques

- Pour la présente simulation :
  - ▶ Création d'une carte phonotactique à partir des *n-grams* de 30.000 mots de CELEX (pas d'information sémantique, seulement information de surface). Pas de catégories grammaticales (comme Proxinette par exemple)
  - ▶ Reconstruction des 30.000 mots comme sommation des états d'activation des *n-grams* dans la carte
  - ▶ Calcul d'une matrice de distance CELEX  $\times$  CELEX (30.000  $\times$  30.000)  $\rightarrow$  métrique *cosinus*

|  |       |   |   |   |     |
|--|-------|---|---|---|-----|
|  | CELEX |   |   |   |     |
|  | 0     |   |   |   |     |
|  |       | 0 |   |   |     |
|  |       |   | 0 |   |     |
|  |       |   |   | 0 |     |
|  |       |   |   |   | ... |

CELEX

# Vue d'ensemble

- 1 Aspects linguistiques
- 2 Création d'un voisinage d'évaluation
- 3 Une morphologie sans sémantique
- 4 Distance d'édition de Levenshtein
- 5 Phonotactic Activation System (PHACTS)
- 6 Proxnette**
- 7 Entropie maximale
- 8 Comparaison des mesures



- a. Deux formes qui partagent **à la fois** des propriétés sémantiques **et** formelles sont morphologiquement proches.
- b. La proximité morphologique de deux mots est d'autant plus grande que **le nombre** des propriétés sémantiques et formelles qu'ils partagent est grand.
- c. La proximité morphologique de deux mots est d'autant plus grande que les propriétés qu'ils partagent leur sont **spécifiques**, c'est-à-dire que peu d'autres mots en sont munis.

# Proximité formelle

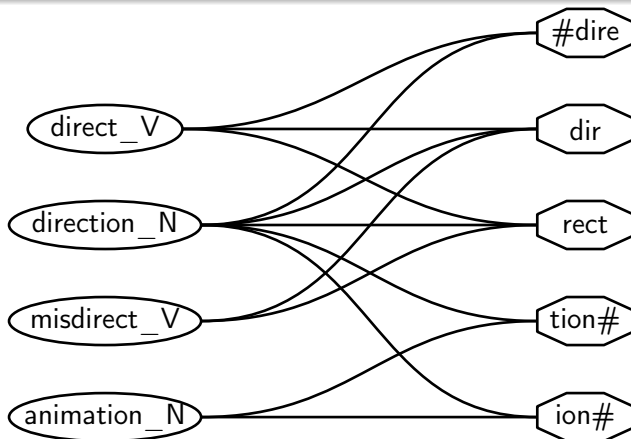
- Proximité est une mesure de **proximité formelle** : elle ne prend en compte que les propriétés formelles des entrées.
- Catégorie et sens ne sont pas utilisés.
- Les propriétés sémantiques et formelles qui servent à déterminer la proximité doivent être **maximalement redondantes** pour capter le plus grand nombre possible de similarités qui existent entre les mots.
- Les propriétés formelles d'un mot sont l'ensemble des séquences sonores qui apparaissent dans sa forme.

## Traits formels de *direction* \_ N

```
#direction#;  
#direction; direction#;  
#directio; direction; irection#;  
#directi; directio; irection; rection#;  
#direct; directi; irectio; rection; ection#;  
#direc; direct; irecti; rectio; ection; ction#;  
#dire; direc; irect; recti; ectio; ction; tion#;  
#dir; dire; irec; rect; ecti; ctio; tion; ion#;  
#di; dir; ire; rec; ect; cti; tio; ion; on#;
```

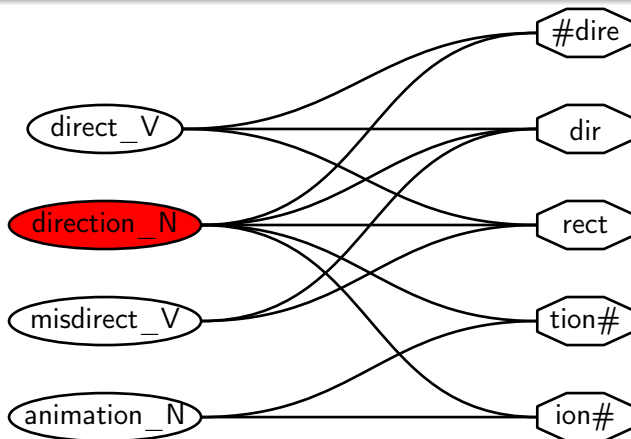
# Proxnette

On construit un graphe où chaque mot est relié à l'ensemble de ses traits phonémiques.

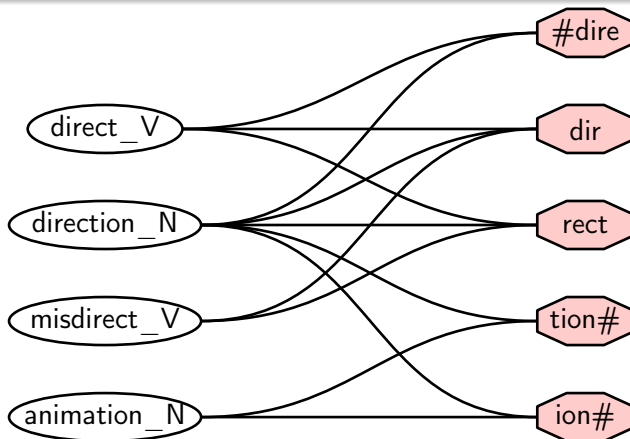


# Proxnette

Une activation est générée au niveau du sommet qui représente *direction\_N*.

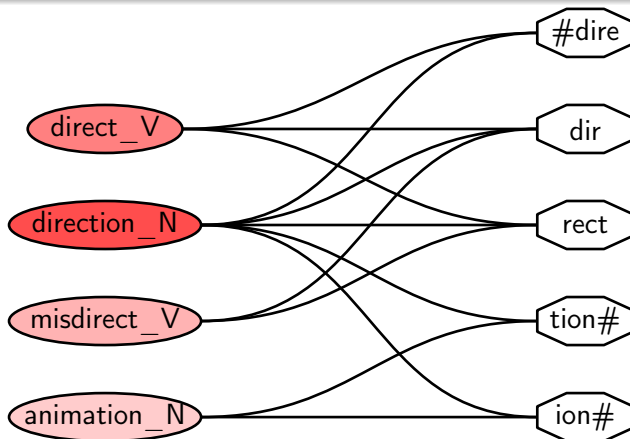


L'activation est propagée uniformément vers l'ensemble des traits de *direction\_N*.



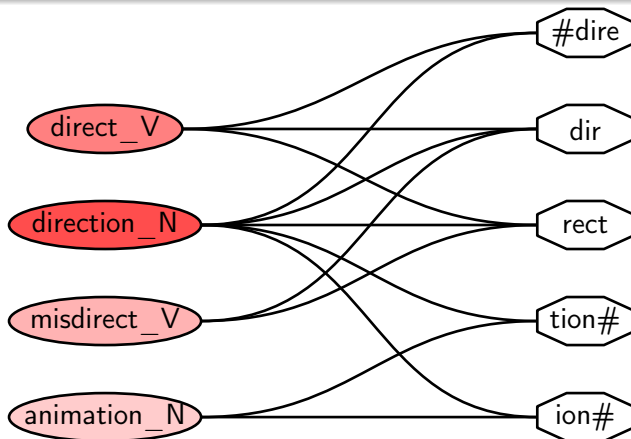
# Proxinette

Les activations sont propagées uniformément vers l'ensemble des mots qui partagent des traits phonémiques avec *direction\_N*.



# Proxinette

Le **niveau d'activation** sur chaque sommet est une estimation du **degré de parenté** des mots correspondants avec *direction\_N*.



## Voisins de *direction*\_N

**direction\_N** :1.0000000000 **directions\_N** :0.7359091020 **directional\_A** :0.7359091020  
**misdirection\_N** :0.7177257909 **directive\_N** :0.3035433276 **directive\_A** :0.3035433276  
**directory\_N** :0.1551731447 **directorship\_N** :0.1551731447 **directorate\_N** :0.1551731447  
**director\_N** :0.1551731447 **directness\_N** :0.1551731447 **directly\_B** :0.1551731447  
**direct\_V** :0.1551731447 **direct\_B** :0.1551731447 **direct\_A** :0.1551731447  
**resurrection\_N** :0.1424231938 **insurrection\_N** :0.1424231938 **erection\_N** :0.1424231938  
**correction\_N** :0.1424231938 **insurrectionist\_N** :0.0885551504  
**insurrectionism\_N** :0.0885551504 **redirect\_V** :0.0831217902 **misdirect\_V** :0.0831217902  
**indirectness\_N** :0.0831217902 **indirectly\_B** :0.0831217902 **indirect\_A** :0.0831217902  
**direfully\_B** :0.0527713005 **direful\_A** :0.0527713005 **dire\_A** :0.0527713005  
**rectification\_N** :0.0252537278 **recollection\_N** :0.0251417002 **dissection\_N** :0.0244677880  
**disinfection\_N** :0.0244677880 **disconnection\_N** :0.0244677880 **disaffection\_N** :0.0244677880  
**vivisection\_N** :0.0241078174 **subjection\_N** :0.0241078174 **selection\_N** :0.0241078174



# Vue d'ensemble

- 1 Aspects linguistiques
- 2 Création d'un voisinage d'évaluation
- 3 Une morphologie sans sémantique
- 4 Distance d'édition de Levenshtein
- 5 Phonotactic Activation System (PHACTS)
- 6 Proxinettes
- 7 Entropie maximale**
- 8 Comparaison des mesures

- Méthode proposée par De Pauw & Wagacha (2007)
- Mêmes traits formels que dans Proxinet :  $n$ -grammes de caractères où  $3 \leq n \leq L(w)$  où  $L(w)$  est la longueur du mot  $w$
- On ajoute en plus **un trait catégoriel** (= la catégorie)

## Classifieur csvLearner (Urieli)

Utilisation « créative » du classifieur :

- 1 On crée un modèle dans lequel il y a une catégorie différente pour chaque mot (autant de catégories qu'il y a de mots)
- 2 On calcule pour chaque mot la probabilité qu'il puisse appartenir à chacune des catégories du modèle.

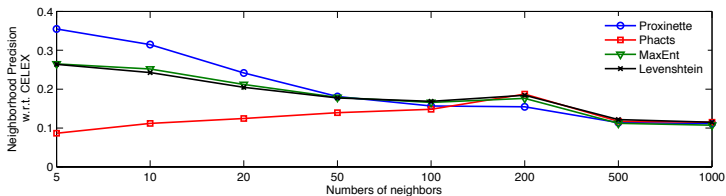
- On n'a utilisé que les mots dont la fréquence dans le corpus COBUILD est supérieure à 10.
- Le classifieur n'utilise que les traits qui caractérisent au moins 5 classes.

## Voisins de *direction*\_N

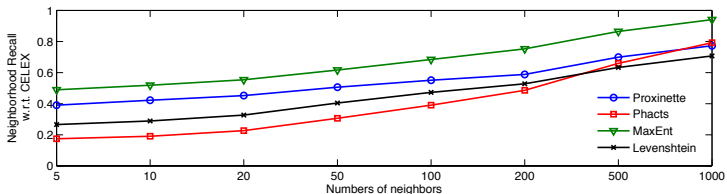
**direction**\_N :0.042159 **directions**\_N :0.015270 **directional**\_A :0.009194  
**directive**\_N :0.003306 **erection**\_N :0.003285 **resurrection**\_N :0.002327  
**correction**\_N :0.001961 **directness**\_N :0.001953 **direct**\_V :0.001840 **direct**\_A :0.001839  
**directly**\_B :0.001785 **director**\_N :0.001595 **directory**\_N :0.001474  
**directorship**\_N :0.001119 **director**ate\_N :0.000970 **dissection**\_N :0.000881  
**vivisection**\_N :0.000738 **inspection**\_N :0.000720 **rejection**\_N :0.000716 **detection**\_N :0.000703  
**defection**\_N :0.000701 **injection**\_N :0.000697 **connection**\_N :0.000683  
**interjection**\_N :0.000673 **disaffection**\_N :0.000672 **election**\_N :0.000669  
**intersection**\_N :0.000668 **infection**\_N :0.000655 **inflection**\_N :0.000644 **selection**\_N :0.000617  
**deflection**\_N :0.000615 **recollection**\_N :0.000586 **misdirect**\_V :0.000580  
**redirect**\_V :0.000574 **objection**\_N :0.000573 **reflection**\_N :0.000571 **affection**\_N :0.000570  
**indirect**\_A :0.000563 **projection**\_N :0.000545 **indirectly**\_B :0.000536  
**subjection**\_N :0.000529 **protection**\_N :0.000529 **perfection**\_N :0.000522  
**imperfection**\_N :0.000502 **collection**\_N :0.000489 **introspection**\_N :0.000443  
**sectional**\_A :0.000280 **diffraction**\_N :0.000267 **confectionery**\_N :0.000257

# Vue d'ensemble

- 1 Aspects linguistiques
- 2 Création d'un voisinage d'évaluation
- 3 Une morphologie sans sémantique
- 4 Distance d'édition de Levenshtein
- 5 Phonotactic Activation System (PHACTS)
- 6 Proxinette
- 7 Entropie maximale
- 8 Comparaison des mesures**

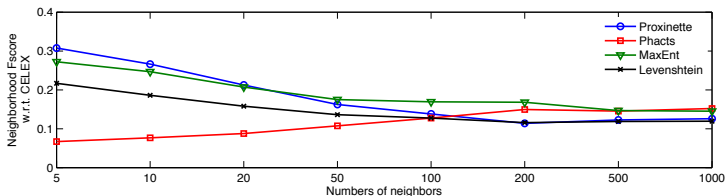


- Proxinet a le meilleur score jusqu'à 50 voisins, après il y a une convergence presque uniforme pour les 4 systèmes (0.13)
- Tendance un peu anormale pour Phacts : la seule ligne qui monte avant la convergence
  - ▶ Phacts est basé sur les *n*-grams de la langue et il développe une sensibilité statistique pour les *n*-grams très fréquents en fin de mot, comme *\_ion#*, *\_able#*, *\_al#*, *\_lly#*, (etc.) qui définissent les séries dérivationnelles de chaque mot (et il y a plus de familles dérivationnelles)



- MaxEnt a le meilleur rappel (0.96)

- ▶ Ce score est très probablement lié au fait que MaxEnt est capable de capter les correspondances régulières entre catégorie et marque caractéristique des séries
- ▶ De ce fait il est plus précis dans la sélection des éléments qui appartiennent à la série. C'est le seul qui est capable de privilégier les mots de la même catégorie



- Phacts est le meilleur au dessus de 500 et Proxinette le meilleur en dessous de 20.
- MaxEnt s'en sort très bien et fait mieux de Proxinette au dessus de 20 parce qu'il utilise les catégories ce qui lui permet d'être plus précis dans l'identification de la série
  - ▶ Bonne performance de Phacts : il n'utilise pas des traits de catégories (comme MaxEnt)

- Directions futures

- ▶ Analyses plus détaillée du comportement de chaque système par rapport les analogies de CELEX : voir et isoler les approximations pour les séries et les familles
- ▶ Pour Phacts : donner une transcription plus phonologique et pas seulement graphémique, avec un codage (binaire) des lieux/modes d'articulation
- ▶ Pour MaxEnt : essayer une apprentissage sans les marques des catégories grammaticales, pour voir si cet aspect est directement lié à sa performance
- ▶ Jouer avec des locuteurs : l'axe Y devrait avoir des valeurs de jugement fournies par des locuteurs sur des tâches de similarité morphologique/wordlikeness



# Distribution des valeurs : Phact vs. Proxnette

