

Expressions polylexicales dans la linguistique computationnelle: on n'est pas sorti de l'auberge

Agata Savary

Université Paris-Saclay, France

Université de Varsovie 18 mai 2022

(exposés similaires ont été donnés le 18 janvier 2019 et le 16 avril 2020)

Expressions polylexicales

Qu'y a-t-il 'quel est le problème' de spécial avec les expressions mises en exergue 'mises en évidence'?

*Si vous **avez** tant **besoin** de **couper l'herbe sous le pied** de quelqu'un, je vous proposerais de **vous en prendre** au **rédacteur-en-chef**, Monsieur **Jean-Marc Petit**.*

Expressions polylexicales

Définition

Combinaisons de **plusieurs mots** qui possèdent des **propriétés irrégulières** au niveau du lexique, de la grammaire, de la sémantique, etc.

Exemples de propriétés irrégulières

- **Sémantique non-compositionnelles**: le sens global n'est pas déductible de manière régulière à partir des sens des composants, et des liens syntaxiques qui les relient.
 - *couper l'herbe sous le pied de quelqu'un* 'empêcher quelqu'un de réussir'
 - *s'en prendre à quelqu'un* 'prendre quelqu'un pour cible, lui attribuer un faute'
- Déterminant zéro
 - *avoir besoin*, **avoir un besoin*, *avoir un besoin important*
- Figement lexical:
 - *rédacteur-en-chef*, **journaliste-en-chef*, **rédacteur-en-leader*
- Figement syntaxique:
 - *Jean-Marc Petit*, **Marc-Jean Petit*

Figement

Une EP est moins flexible (variable) qu'une construction régulière de la même structure syntaxique.

Construction régulière	Expression polylexicale	Propriété
<i>livre bleu</i> ≈ ¹ <i>livre cyan</i> ≈ <i>bouquin bleu</i>	<i>cordon bleu</i> vs. <i>#cordon cyan</i> vs. <i>#corde bleue</i>	Figement lexical
<i>manger des salades</i> ≈ <i>manger une salade</i>	<i>raconter des salades</i> vs. <i>#raconter une salade</i>	Figement morphologique
<i>il a fait la soupe</i> ≈ <i>la soupe a été faite par lui</i>	<i>il a fait la tête</i> vs. <i>#la tête a été faite par lui</i>	Figement syntaxique
<i>les patates sont cuites</i> ≈ <i>nous avons cuit les patates</i>	<i>les carottes sont cuites</i> vs. <i>#on a cuit les carottes</i>	

¹'≈': glissement de sens prévisible à partir du remplacement lexical

Expressions polylexicales – un vrai *bric-à-brac* ‘amas d'objets hétéroclites’

- mots composés
 - *cordon bleu* ‘excellent cuisinier’, *pomme de terre*
- termes complexes
 - *mémoire vive* ‘RAM de l'ordinateur’
- entités nommées polylexicales
 - *Mer Morte*
- constructions à verbe support
 - *avoir besoin*, *rendre visite*
- idiomes
 - *ne pas être sorti de l'auberge* ‘ne pas avoir terminé avec des problèmes’
- proverbes
 - *qui sème le vent récolte la tempête* ‘celui qui incite à la violence doit s'attendre à de fâcheuses conséquences’

Elles se rencontrent *tous les quatre matins* 'très souvent'

Fréquence des EP dans les textes

Jusqu'à 20% des mots d'un texte appartiennent à des expressions polylexicales.

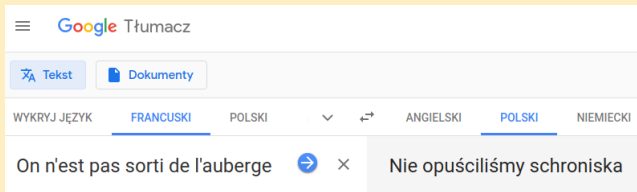
Si vous *avez* tant *besoin* de *couper l'herbe sous le pied* de quelqu'un, je vous proposerais de *vous en prendre* au *rédacteur-en-chef*, Monsieur *Jean-Marc Petit*.

- Ici: 17 composants d'EP sur 30 mot du texte → 57%

Les *casse-pieds* 'qui ennuient/dérangent' de la linguistique computationnelle

Du fait de la non-compositionnalité, les EP constituent un défi pour des tâches du traitement sémantique.

Traduction automatique



Traductions mot-à-mot ne captent pas le sens idiomatique.

Recherche d'information (RI)

- La tâche: pour une requête donnée (un ou plusieurs mots), trouver automatiquement tous les documents pertinents
- Moteurs de recherche: Google Search, Bing, Baidu, Yahoo!, ...
- Technique:
 - Tous les textes sont lemmatisés et **indexés**, i.e. représentés comme listes des mots qu'ils contiennent
 - Exemple: *les cornes et les pieds du taureau* → {corne - 1, du - 1, et - 1, le - 2, pied - 1, taureau - 1}
 - Chaque mot d'une requête d'un utilisateur est recherchée dans l'indexe. Les textes contenant les bons mots sont pondérés et retournés en réponse.
- Problème avec les EP:
 - Un texte contient *prendre le taureau par les cornes* 'faire face aux difficultés plutôt que les fuir'
 - La requête d'un utilisateur contient *cornes d'un taureau*
 - Ce texte n'est pas pertinent, mais sera probablement proposé



Fouille d'opinion

- La tâche: prédire automatiquement la valence (positive, neutre ou négative) de l'opinion exprimée par un texte
- Exemples:
 - *Je soutiens les gilets jaunes. Je suis très respectueux de leur cause. Leur mouvement est compréhensible.*
 - *Il ne faut pas accepter ce geste. Rien ne le justifie ! Cela s'appelle une agression.*
 - *Tout cela va mal finir. La violence est partout dans les paroles et dans les actes.*
- Technique:
 - Les mots simples sont annotés avec une valence élémentaire: *respectueux* → 1, *agression* → -2, *justifier* → 1
 - Des règles locales modifient la valence élémentaire:
 - *très, grand* doublent la valence; *très respectueux* → $1*2 = 2$; *grande agression* → $-2*2 = -4$
 - La négation inverse la valence: *rien ne le justifie* → $-1*1 = -1$

Fouille d'opinion – problèmes avec les EP

Texte

Valence
calculéeValence
réelle

bras₀ d'honneur₁ 'geste impoli qui marque la reprobation'

aimer₂ à la folie₋₂ 'aimer énormément'

boire₀ comme un trou₀ 'être alcoolique'

avoir un coup₀ de foudre₋₁ 'tomber amoureux d'un seul coup'

avoir la pêche₀ 'être en forme'

ne pas être dans son assiette₀ 'ne pas être en forme'

la moutarde₀ lui monte₀ au nez₀ 'il se met en colère'

avoir un poil₀ dans la main₀ 'être paresseux'

tourner₀ au vinaigre₀ 's'orienter vers la dispute'

raconter des salades₀ 'dire n'importe quoi'

en faire tout un fromage₀ 'exagérer'

faire la tête₀ 'bouder'

chercher la petite bête₀ 'regarder trop aux détails'

ne pas être sorti₀ de l'auberge₀ 'continuer à avoir des difficultés'

les carottes₀ sont cuites₀ 'la situation est compromise'



Fouille d'opinion – problèmes avec les EP

Texte

	Valence calculée	Valence réelle
<i>bras₀ d'honneur₁</i> 'geste impoli qui marque la reprobation'	1	-2
<i>aimer₂ à la folie₋₂</i> 'aimer énormément'	0	4
<i>boire₀ comme un trou₀</i> 'être alcoolique'	0	-2
<i>avoir un coup₀ de foudre₋₁</i> 'tomber amoureux d'un seul coup'	0	2
<i>avoir la pêche₀</i> 'être en forme'	0	1
<i>ne pas être dans son assiette₀</i> 'ne pas être en forme'	0	-1
<i>la moutarde₀ lui monte₀ au nez₀</i> 'il se met en colère'	0	-1
<i>avoir un poil₀ dans la main₀</i> 'être paresseux'	0	-1
<i>tourner₀ au vinaigre₀</i> 's'orienter vers la dispute'	0	-1
<i>raconter des salades₀</i> 'dire n'importe quoi'	0	-1
<i>en faire tout un fromage₀</i> 'exagérer'	0	-1
<i>faire la tête₀</i> 'bouder'	0	-1
<i>chercher la petite bête₀</i> 'regarder trop aux détails'	0	-1
<i>ne pas être sorti₀ de l'auberge₀</i> 'continuer à avoir des difficultés'	0	-1
<i>les carottes₀ sont cuites₀</i> 'la situation est compromise'	0	-1



Solutions

- Identifier automatiquement les EP dans le texte.
- Leur appliquer des traitements spéciaux:
 - traduction automatique
 - reformuler une EP avant la traduction
 - *il boit comme un trou* → *il est alcoolique* → *jest pijakiem*
 - recherche d'information
 - ne pas rajouter les composants d'une EP dans l'indexe du texte
 - rajouter l'expression en entier
 - *Le débat au parlement a tourné au vinaigre* → {au - 1, avoir - 1, débat - 1, le - 1, parlement - 2, tourner au vinaigre - 1}
 - fouille d'opinion
 - attribuer une valence à une EP en entier
 - *avoir un [coup de foudre]₂*

Rôle des linguistes

Créer des ressources linguistiques:

- lexiques
- corpus annotés

Annotation d'EP en corpus

The screenshot displays the FoLIA Linguistic Annotation Tool interface in a Chromium browser. The main window shows a corpus of French text with various linguistic annotations. On the left, there is a sidebar with a 'Perspective' dropdown set to 'Sentence', a 'page' indicator set to '2', and a 'Selector' dropdown set to 'Automatic (deepest)'. Below this is a 'Legend • Entity' section with a list of entities and their corresponding colors: EP-4.1-LEX (green), EN-2-ORG.Final (purple), EP-4.3-DET (red), EP-3-IRREG (yellow), EN-1-PERS.Final (cyan), EP-6.2-CL (pink), and EP-1-CRAN (orange). The main text area shows several paragraphs of text with annotations above and below the words. Annotations include entity codes in parentheses with arrows, such as (EP-4.1-LEX) above 'à', (EN-2-ORG.Final) above 'Union', and (EP-3-IRREG) above 'peut-être'. The text is numbered on the left margin from 176 to 186. The text includes phrases like 'Je voudrais rappeler à cet égard, qu'il y a quelques semaines, 80 000 jeunes des de les pays de l'Union européenne ont participé à un concours pour la recherche d'une devise pour l'Europe et que la devise qui a été finalement retenue par un grand jury a été "L'unité dans la diversité."', 'Je dois avouer que cela n'est peut-être pas génial, mais c'est plus intéressant qu'il n'y paraît parce que cela me semble répondre au à le sentiment très profond de beaucoup de citoyens de nos pays.', 'Enfin, vous avez rappelé, Monsieur le Président, les valeurs auxquelles à lesquelles vous teniez, et qui sont à la base de l'intégration européenne.', 'Vous avez aussi évoqué le souhait de ne pas perdre de vue la solidarité sociale, dans le contexte de la globalisation.', 'Là encore, il me semble que vous rejoignez parfaitement les objectifs de notre Parlement européen.', 'Je vous souhaite bonne chance ainsi qu'à toutes les autorités slovènes qui participent aux à les négociations.', 'Nous espérons vivement que ces négociations aboutiront dans les délais prévus.', 'Bonne chance, Monsieur le Président, et nous vous remercions encore de votre présence et de votre intervention.', '(La séance solennelle est close à 12h30)', 'Monsieur le Président, il devait y avoir un débat sur la violence dans le football.', and 'Les événements de la nuit dernière à Copenhague soulignent à quel point il est important que le Parlement européen tienne ce débat, comme il a décidé de le faire plus tôt dans la semaine.'

Corpus PARSEME d'expressions polylexicales verbales (EPV)

Coopération internationale

- réseau scientifique européen PARSEME: 26 équipes nationales
- terminologie et méthodologie unifiées
- corpus de 26 langues, 8 millions de words, 113 mille EP annotées

Familles de langues

- **Balto-slaves:** bulgare (BG), croate (HR), lituanien (LT), polonais (PL), serbe (RS), slovène (SL), tchèque (CS)
- **Germaniques:** allemand (DE), anglais (EN), suédois (SV)
- **Romanes:** espagnol (ES), français (FR), italien (IT), portugais brésilien (PT), roumain (RO)
- **Autres:** arabe (AR), basque (EU), chinois (ZH), farsi (FA), grec (EL), hébreu (HE), hindi (HI), hongrois (HU), irlandais (GA), maltais (MT), turc (TR)

Typologie

Catégories universelles (valables dans toutes les langues)

- idiomes verbaux (verbal idioms : **VID**)

ne pas être sorti de l'auberge

- constructions à verbe support (light-verb constructions : **LVCs**) - le nom exprime l'action/état, le verbe apporte seulement la morphologie (temps, mode, personne, etc.)

rendre visite (LVC.full)

entraîner la mort (LVC.cause)

Catégories quasi-universelles (valables dans plusieurs langues)

- verbes intrinsèquement pronominaux (inherently reflexive verbs : **IRVs**)

s'apercevoir

- verbes à particule (verb-particle constructions : **VPCs**)

(EN) *to do in* 'tuer' (VPC.full)

(EN) *to eat up* 'manger tout' (VPC.semi)

- constructions multi-verbales (multi-verb constructions : **MVCs**) – en langues asiatiques

(HI) *kar le-na* 'do take.INF' ⇒ 'to do something (for one's own benefit)'

Guide d'annotation – *metez la main à la pâte* 'travaillez vous-mêmes'

Tu feras d'une pierre deux coups, si tu te comportes bien.

Tu feras d'une pierre deux coups

- Suivre l' ▸ arbre de décision
- S.1 [1HEAD] (YES): *fera* est la seule tête verbale
- S.2 [1DEP] (NO): *d'une pierre* et *deux coups* sont les deux dépendants lexicalisés
- VID.1 [CRAN] (NO): tous les composants existent comme mots indépendants
- VID.2 [LEX] (YES): *#tu feras de deux pierres quatre coup*
- Résultat: **VID**

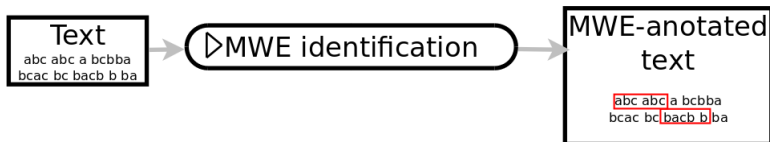
Guide d'annotation

Tu feras d'une pierre deux coups, si tu te comportes bien.

si tu te comportes bien

- S.1 [1HEAD] (YES): *comportes* est la seule tête verbale
- S.2 [1DEP] (YES): *te* est le seul dépendant lexicalisé de *comportes*
- S.3 [LEX-SUBJ] (NO): le sujet *tu* n'est pas lexicalisé
- S.4 [CATEG] (reflexive clitic): *te* est un pronom réflexif
- IRV.1 [INHERENT] (NON): *comporte* peut apparaître tout seul
- IRV.2 [DIFF-SENSE] (YES): *se comporter* a un sens très différent que *comporter*
- Résultat: **IRV**

Tâche d'identification automatique d'EPV



C'est *tip-top* 'excellent' ou ça *ne casse pas trois pattes à un canard* 'c'est médiocre' ?

Evaluation d'une annotation automatique

- **Vrais (V)** – entités existantes dans le texte (annotées par un humain)
- **Positifs (P)** – entités identifiées par le système
- **Vrais positifs (VP)** – entités correctement identifiées par le système

Mesures d'évaluation

- **Précision:** $P = \frac{|VP|}{|P|}$
- **Rappel:** $R = \frac{|VP|}{|V|}$
- **F-mesure:** $F = \frac{2 * P * R}{|P| + |R|}$

C'est tip-top ou ça ne casse pas trois pattes à un canard ?

Vrais (annotés par un humain)

Si vous **avez** tant **besoin** de **couper l'herbe sous le pied** de quelqu'un, je vous proposerais de **vous en prendre** au **rédacteur-en-chef**, Monsieur **Jean-Marc Petit**.

Positifs (identifiés par le système)

Si vous avez tant besoin de **couper l'herbe sous le pied** de quelqu'un, je **vous proposerais** de **vous en prendre** au rédacteur-en-chef, Monsieur **Jean-Marc Petit**.

- $|V| = ?$
- $|P| = ?$
- $|VP| = ?$
- $P = \frac{|VP|}{|P|} = ?$
- $R = \frac{|VP|}{|V|} = ?$
- $F = \frac{2 * P * R}{|P| + |R|} = ?$

C'est tip-top ou ça ne casse pas trois pattes à un canard ?

Vrais

Si vous **avez** tant **besoin** de **couper l'herbe sous le pied** de quelqu'un, je vous proposerais de **vous en prendre** au **rédacteur-en-chef**, Monsieur **Jean-Marc Petit**.

Positifs

Si vous avez tant besoin de **couper l'herbe sous le pied** de quelqu'un, je **vous proposerais** de **vous en prendre** au **rédacteur-en-chef**, Monsieur **Jean-Marc Petit**.

- $|V| = 5$

- $|P| = 4$

- $|VP| = 2$

- $P = \frac{2}{4} = 0,5$

- $R = \frac{2}{5} = 0,4$

- $F = \frac{2*P*R}{|P|+|R|} = \frac{2*0,5*0,4}{0,5+0,4} = \frac{0,4}{0,9} = 0,44$

Campagne d'évaluation PARSEME (édition 1.2 en 2020)

- 7 équipes, 9 systèmes d'identification d'EPV, 14 langues
- Le corpus PARSEME pour chaque langue est divisé en 2 parties
 - (grand) corpus d'entraînement (TRAIN)
 - (petit) corpus d'évaluation en deux versions: annotée (TEST) et non-annotée (TEST.blind)
- Les systèmes automatiques s'entraînent sur le TRAIN
- Ils sont appliqués au TEST.blind
- Les organisateurs comparent leurs résultats au TEST et calculent P, R et F

On a fait du chemin 'on a progressé', mais
on n'est pas sorti de l'auberge 'on n'en a pas fini avec
des difficultés'

Meilleurs résultats pour toutes les langues

Système	P	R	F
MTLB-STRUCT (édition 1.2 2020)	0,71	0,69	0,70
Seen2Seen (édition 1.2 2020)	0,76	0,58	0,66
Meilleur système dans l'édition 1.1 (2018)	0,66	0,51	0,58

Meilleurs résultats pour le français

Système	P	R	F1
TRAVIS-mono (édition 1.2 2020)	0,82	0,83	0,82
MTLB-STRUCT (édition 1.2 2020)	0,80	0,78	0,79
Meilleur système dans l'édition 1.1 (2018)	0,72	0,53	0,61

Point barre! 'tout a été dit'

Conclusions

- Les EP sont fascinantes!
 - Elles font passer le message de manière succincte et efficace
 - Elles cachent des traces de l'histoire, des stéréotypes, des connotations surprenantes
 - Elles peuvent être très drôles
- Les EP sont difficiles
 - Elles sont difficiles à apprendre pour les locuteurs non natifs
 - Elles témoignent du degré de maîtrise d'une langue
 - Elles se comportent différemment que des combinaisons régulières de mots
 - Elles sont difficiles à identifier, analyser et traduire automatiquement
- Les EP très fréquentes
 - Elles couvrent jusqu'à 20% de mots d'un texte en langage naturel