

Transferts linguistiques des langues régionales vers le français, transferts technologiques du français vers les langues régionales

Ph. Boula de Mareüil (LIMSI, U. Paris Sud)
& D. Bernhard (LiLPa, U. Strasbourg)

Variation linguistique et Crowdsourcing :
Étudier la variation au 21ème siècle

23 octobre 2015



Le projet RESTAURE



- Projet financé par l'ANR
- Début : janvier 2015
- Durée : 42 mois
- Objectifs :
 - Acquisition et normalisation de ressources (corpus et lexiques)
 - Développement d'outils pour l'acquisition et l'analyse de corpus
 - Diffusion auprès du grand public
- Site web : <http://restaure.unistra.fr/>

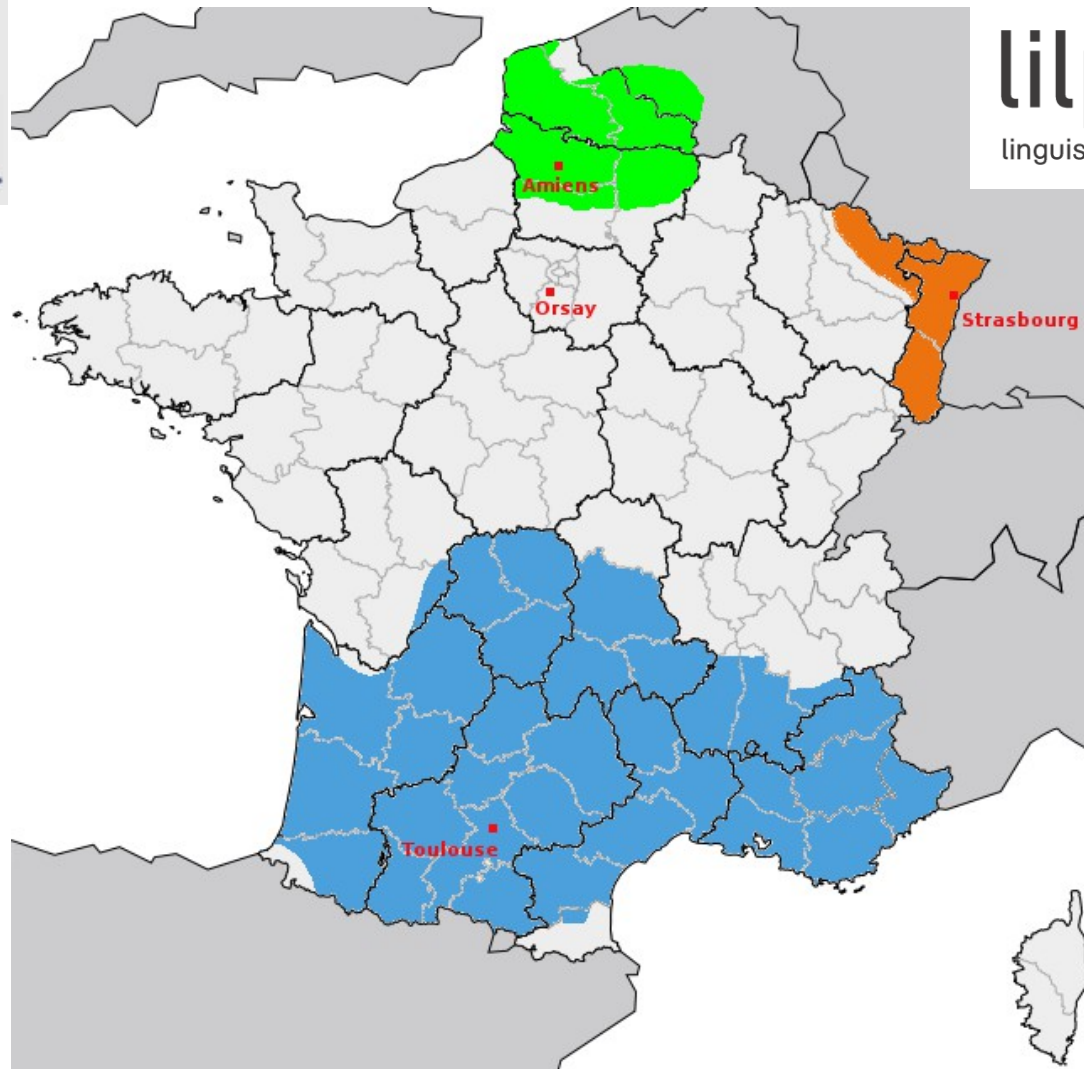
Partenaires et langues régionales



Picard
Resp : C. Rey



Occitan
Resp : M. Bras



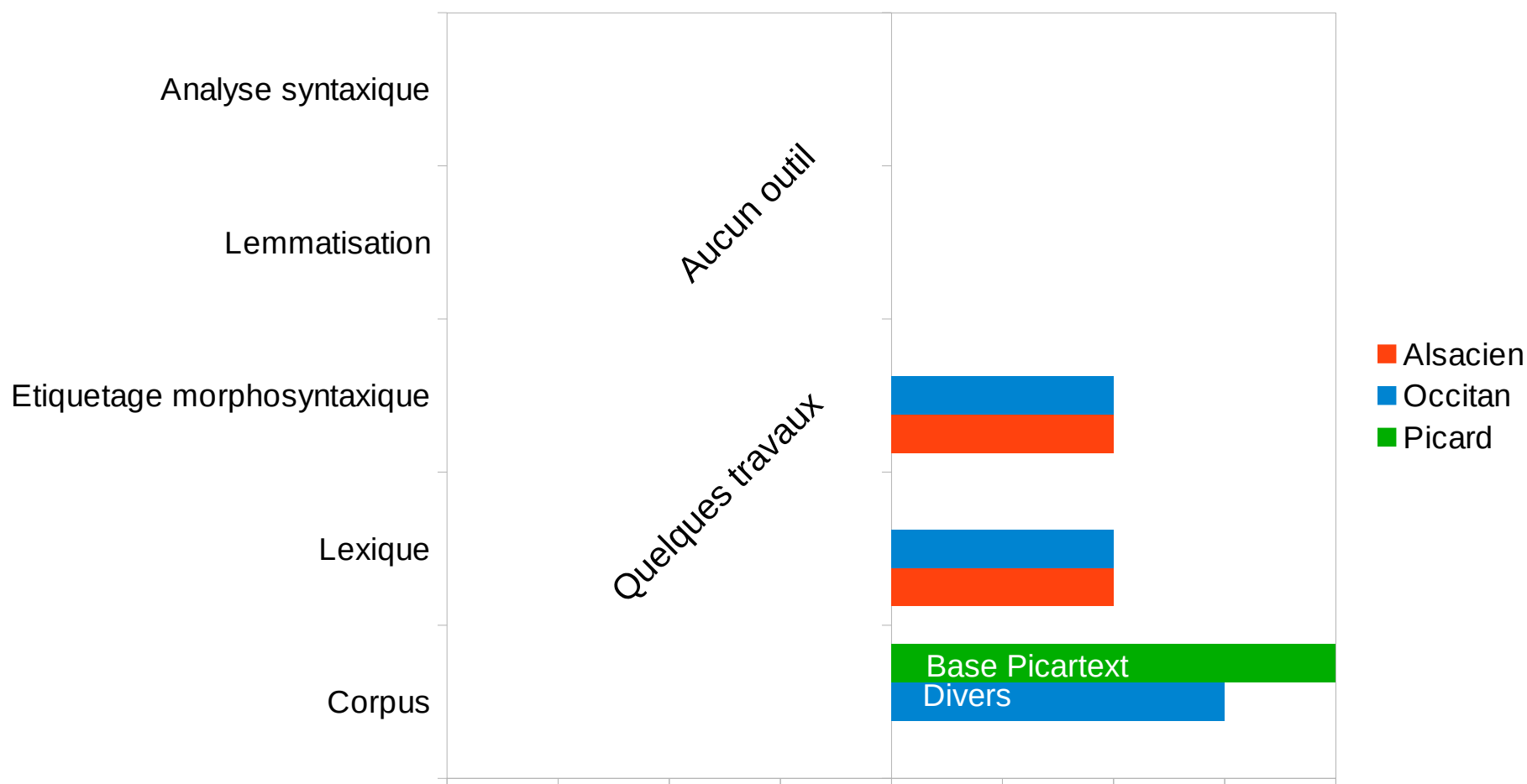
Alsacien
Resp : D. Bernhard



TAL
Resp. : A-L Ligozat

Motivations

Les langues régionales de France sont des langues peu dotées



Transferts technologiques depuis des langues "mieux dotées"

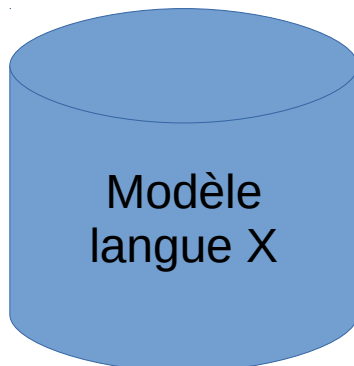
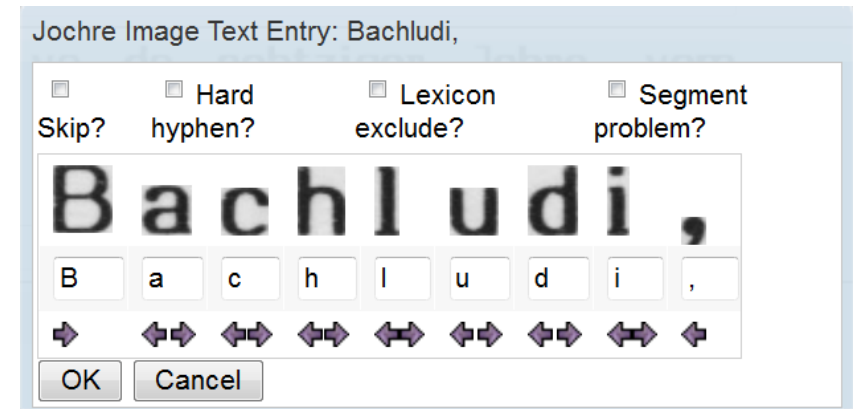
- Méthode : exploiter des outils et ressources disponibles pour des langues proches
 - Français pour le picard ?
 - Allemand pour l'alsacien ?
 - Catalan pour l'occitan ?
- Influence de l'application sur le choix de la langue "mieux dotée"
 - Exemple de la reconnaissance optique de caractères (OCR)

Principe de fonctionnement de l'OCR : exemple de Jochre (Assaf Urieli, CLLE-ERSS Toulouse et Joliciel informatique)

Segmentation des images en paragraphes, lignes, mots et « formes »



Mise en correspondance forme - caractère



Apprentissage d'un modèle spécifique à la langue :

- Extraction de "traits"
- Utilisation d'un algorithme d'apprentissage

Interface de correction en ligne de Jochre

Katz - Mi Sundgäu. Alemannische Gedichte in Sundgauer Mundart, Page 10

'S ROSELE

I

« Brruuungng ! » het's dunneret, un d'erschti
Tropfe sin käit, gross wie Tübeeier un chlitzerig
wie Silber, un wu si uf d'Acker un uf d'Wäg
üfprosslet hai, so het's ganz chleini Stäubwilchle
gäh.

As isch o alles brottrocke gsi.

Scho vier Wuche n e Hitz, ass mr hätt mächte
z'Wasser wärde, un e kä Tropfe Räge, un d'Lit
hai efange n all Obe d' Bärömäter bschäut, eb's
Quacksilber hoch säig, un sin vor de Wätterhisle
gschtange, un hai glüegt eb's Männle edder 's
Wible dusse n isch, un mr het efange nit meh
ghert gha ass : « Chinge, isch das aber e Hitz !
As geht eim jo no alles hi ! E güed Rägele jetze
wär gar nit z'zahle ! »

« Brruuungng ! » Un jetz isch's losgange ! Was
vom Himmel het obe n abe chänne !

Un do het Eini schnäll ihr Haiele uf dr Buckel
gnu, un dert het Ein schnäll d'Ross üsgschpanne,
— un drüfsgässe ! Un Galopp, was es gange n
isch, im Därfle züe !

Mr het nit gseh, ass Lit wu springe : d' Wiber
ihri roti edder bläui Räck iber e Chopf, un
d'Männer ihri Chäpf in dr Chittelchrage n ine-
zoge. — Un ab, was es gange n isch ; alles im
Därfle züe, edder enaime n unger e Bäum, edder
e n angere Ungerschlupf . . .

10

« Brruuungng ! » het's dunneret, un d'erschti
|«| B r r u u u n g n g ! » het's dunneret, un d'erschti

[[«|«| Brruuungng ! » het's dunneret, un d'erschti

Tropfe sin käit, gross wie Tübeeier un chlitzerig
T r o p f e s i n k ä i t , g r o s s w i e T ü b e e i e r u n c h l i t z e r i g

Tropfe sin käit, gross wie Tübeeier un chlitzerig

wie Silber, un wu si uf d'Acker un uf d'Wäg
w i e S i l b e r , u n w u s i u f d ' A c k e r u n u f d ' W ä g

wie Silber, un wu si uf d'Acker un uf d'Wäg

üfprosslet hai, so het's ganz chleini Stäubwilchle
gäh.
ü f p r o s s l e t h a i , s o h e t ' s g a n z c h l e i n i S t ä u b w i l c h l e

üfprosslet hai, so het's ganz chleini Stäubwilchle

As isch o alles brottrocke gsi.
Ä s i s c h o a l l e s b r o t t r o c k e g s i .

As Isch o alles brottrocke gsi.

Scho vier Wuche n e Hitz, ass mr hätt mächte
S c h o v i e r W u c h e n e H i t z , a s s m r h ä t t m ä c h t e

Scho vier Wuche n e Hitz, ass mr hätt mächte

OCR pour l'alsacien

Outils (open source) :

- **Tesseract** : open source depuis 2005 et développé par Google

- **Jochre** : développé par Assaf Urieli

Modèles existants

- français
 - allemand
 - ... (plus de 100 langues / typographies)
-
- yiddish et occitan

Dépasser les a priori...

- *"Il n'existe en effet qu'une seule définition scientifiquement correcte de la langue régionale en Alsace, ce sont les dialectes alsaciens dont **l'expression écrite est l'allemand**".*

Recteur Pierre Deyon, 1985

Dépasser les a priori...

- *"Il n'existe en effet qu'une seule définition scientifiquement correcte de la langue régionale en Alsace, ce sont les dialectes alsaciens dont l'**expression écrite est l'allemand**".*

Recteur Pierre Deyon, 1985

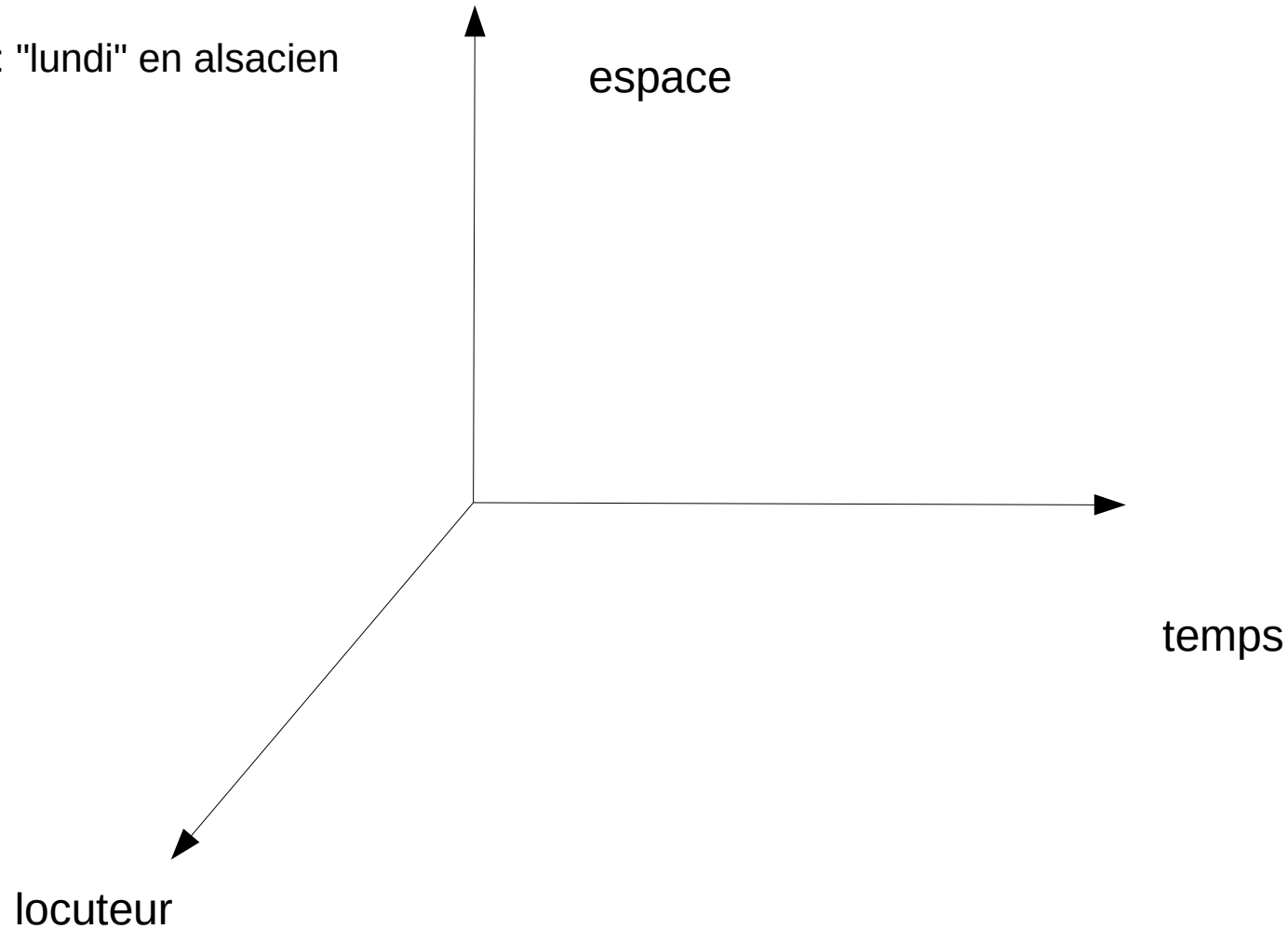
- Résultats pour Tesseract :

Modèle	Pourcentage d'erreurs (caractères)	Exactitude des mots (après normalisation)
allemand	2,33%	91,25%
français	2,42%	91,03%
français + allemand	1,68%	94,27%

- Pourquoi le modèle français fonctionne-t-il aussi bien ?

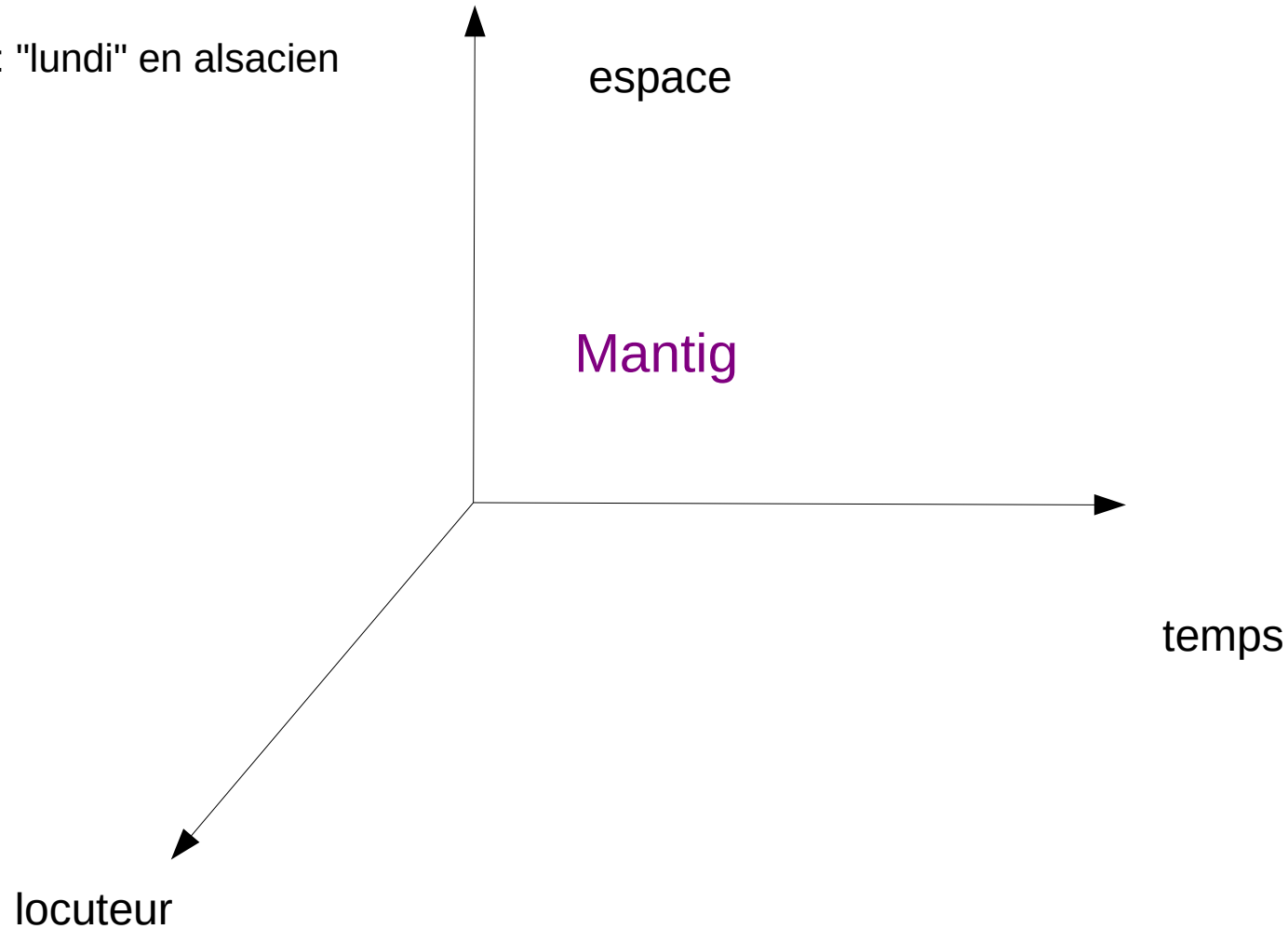
Un défi majeur : la variation graphique

Exemple : "lundi" en alsacien



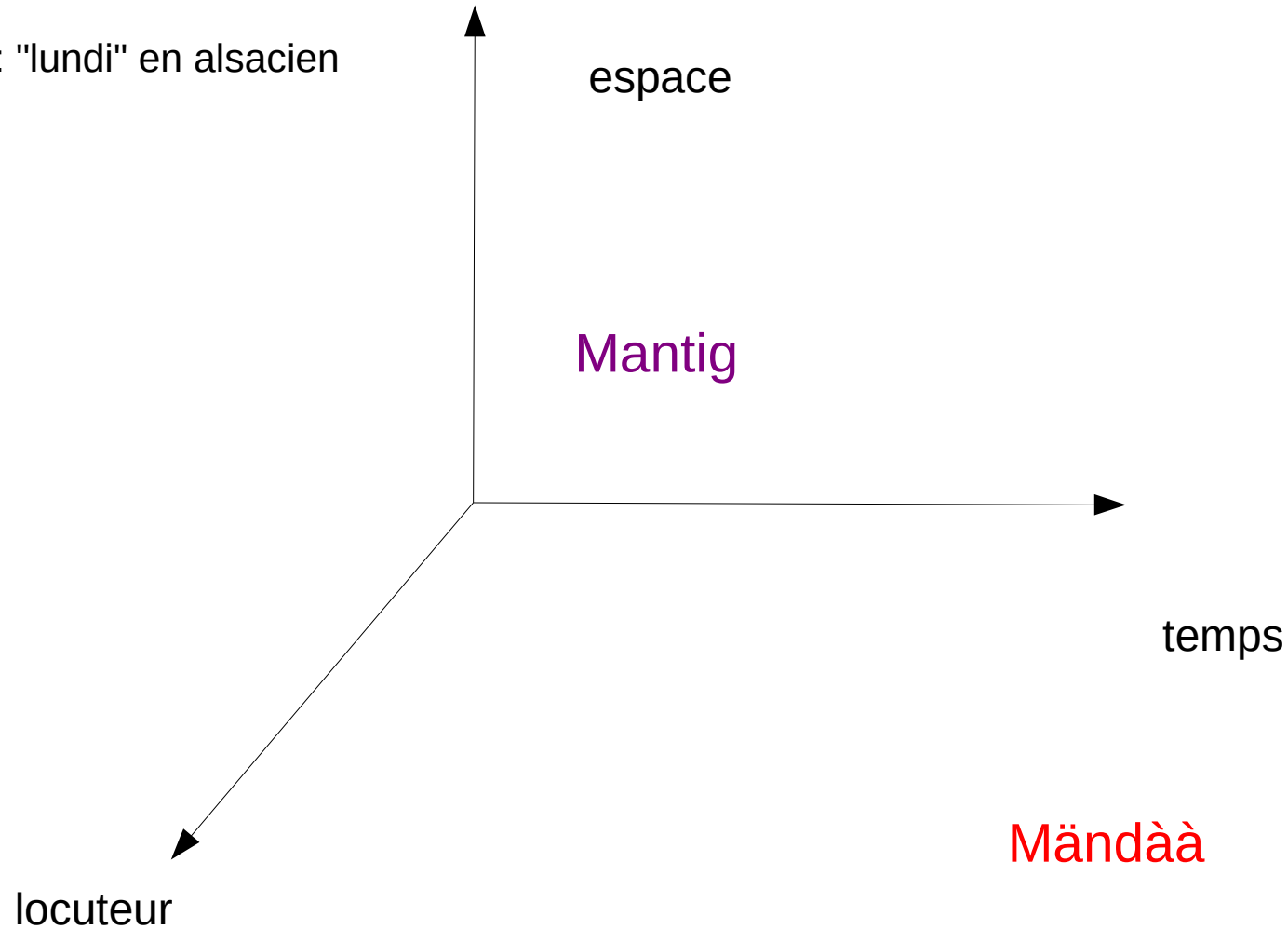
Un défi majeur : la variation graphique

Exemple : "lundi" en alsacien



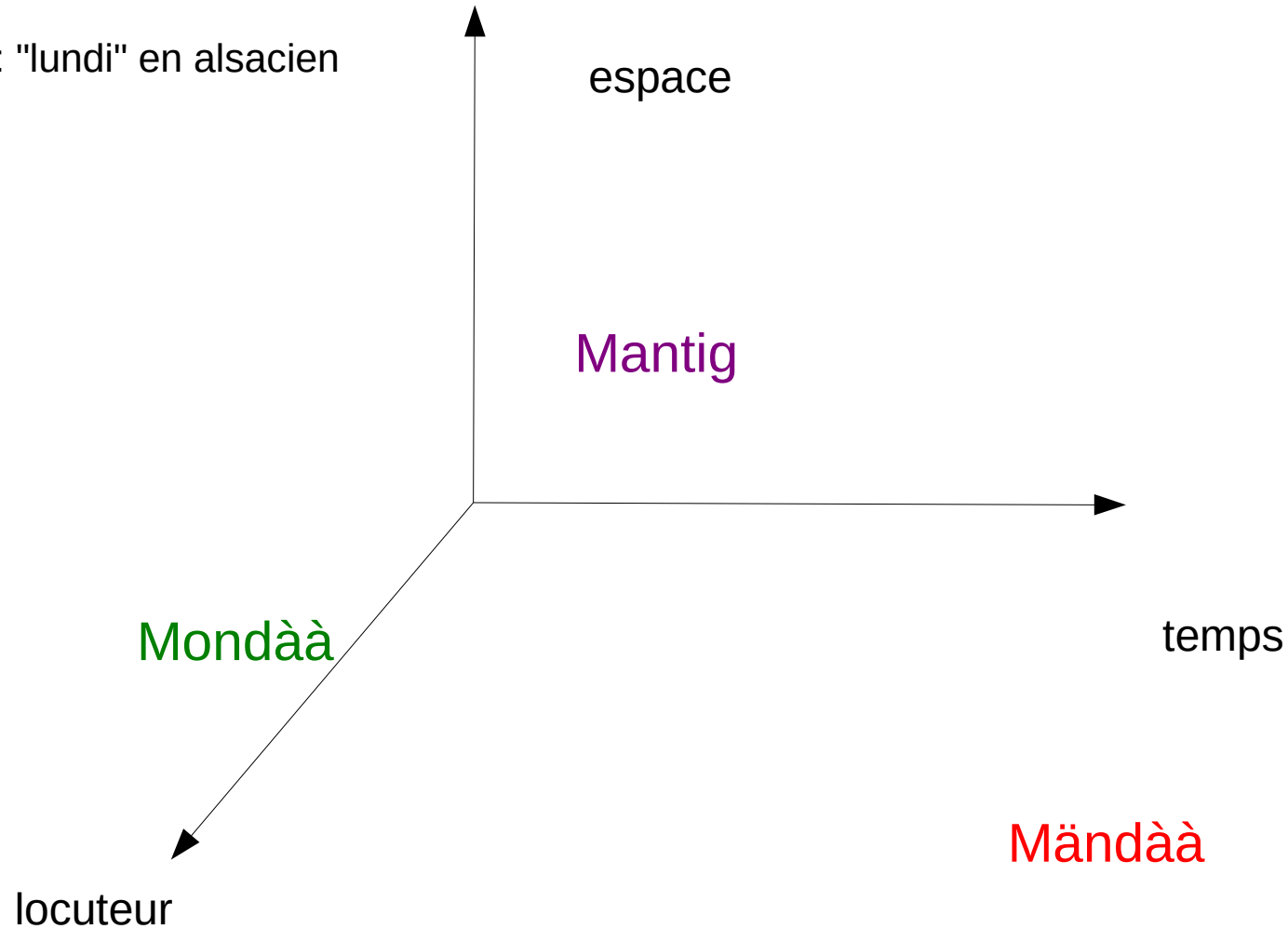
Un défi majeur : la variation graphique

Exemple : "lundi" en alsacien



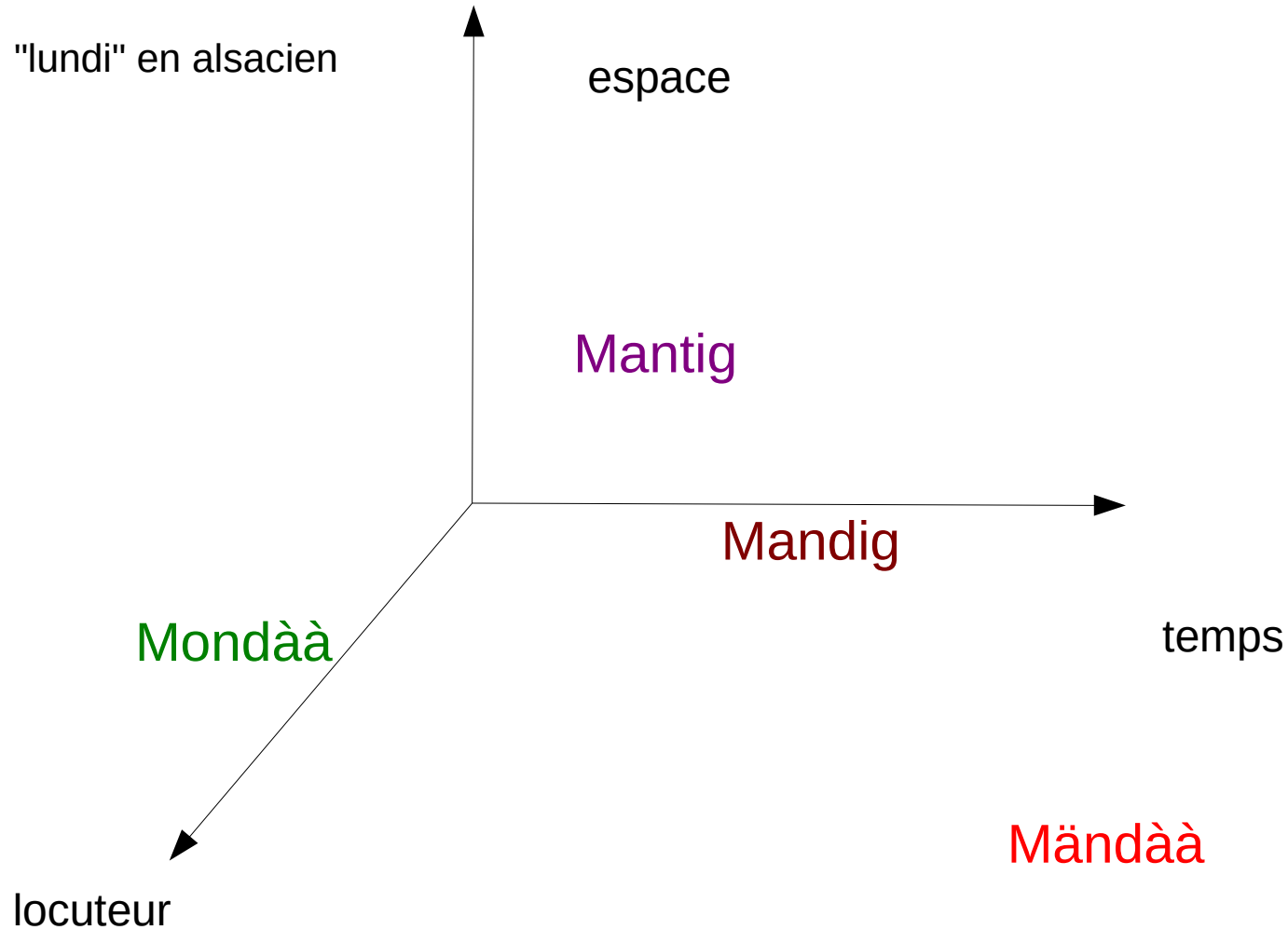
Un défi majeur : la variation graphique

Exemple : "lundi" en alsacien



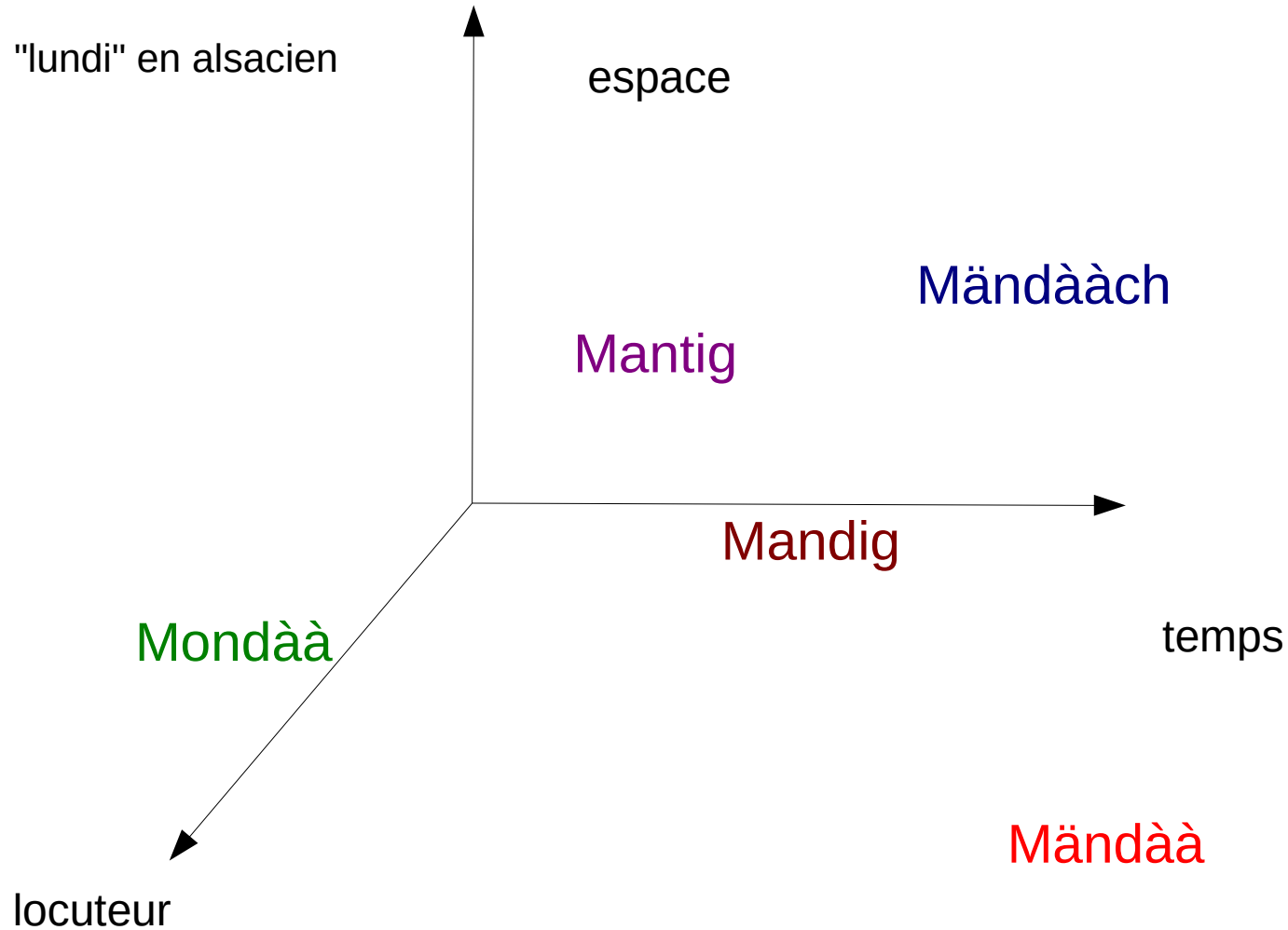
Un défi majeur : la variation graphique

Exemple : "lundi" en alsacien



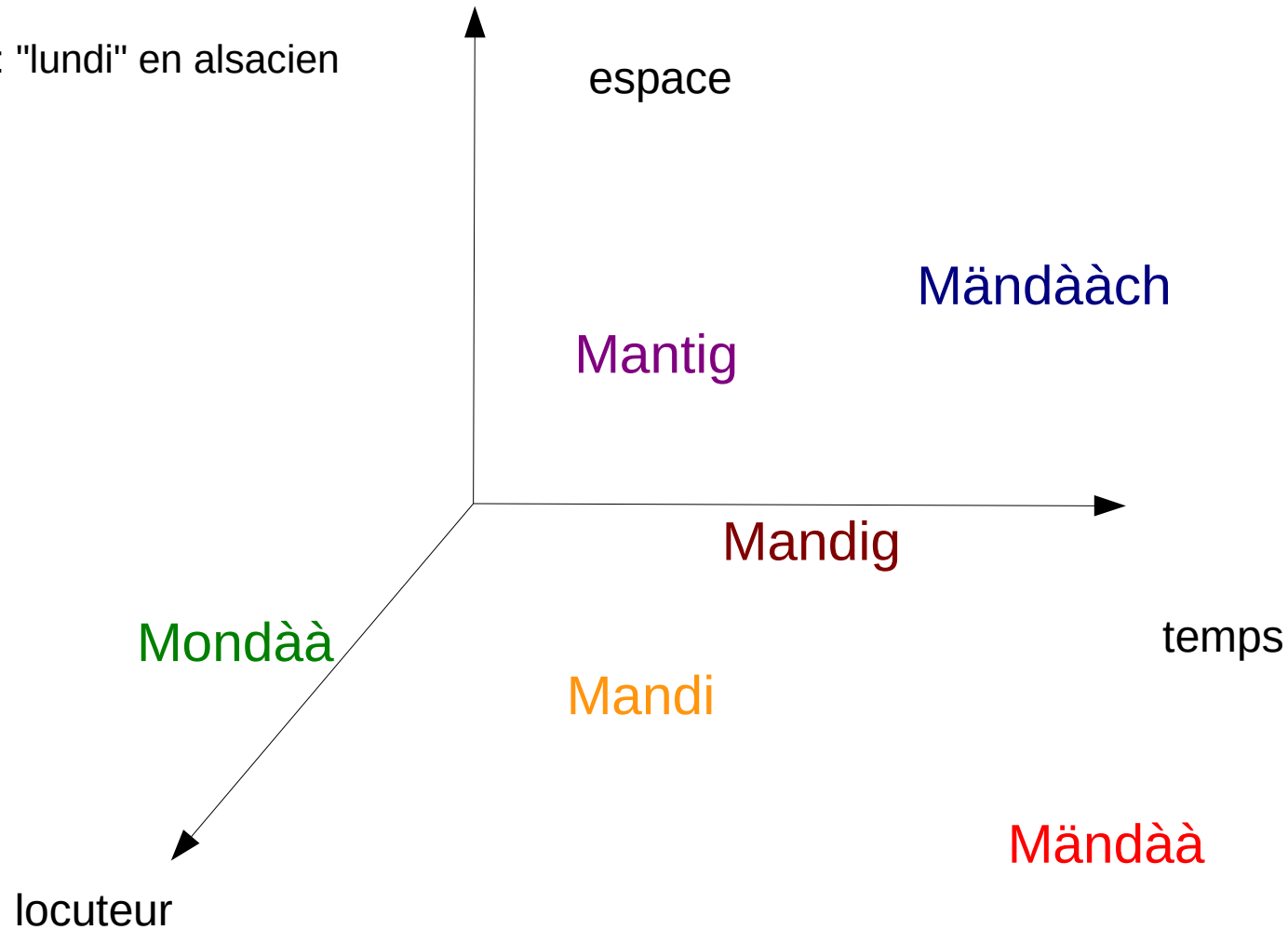
Un défi majeur : la variation graphique

Exemple : "lundi" en alsacien



Un défi majeur : la variation graphique

Exemple : "lundi" en alsacien



L'alsacien a du caractère

Caractères spécifiquement allemands	Total : 563
ä	419
Ä	5
ö	83
ß	56

Les caractères allemands couvrent **98,3 %** du corpus...

... les caractères français : **99 %**

Caractères spécifiquement français	Total : 960
à	651
À	22
â	75
ç	4
Ç	1
é	154
É	1
è	21
ê	1
î	23
ï	4
ô	2
ù	1

Exemple

Mit dàm G'spràch sin die zwei uf Altkilch ku, si hân nit g'wisst wie. Dr Lüwi isch glich in's "Reglis" ufe g'fahre, un het si Füehrwark dert îg'stellt. Vater un Grossvater han scho allewîl dert îkehrt, un o àr sàlber, wenn 'r îne g'fahre isch. D'Wirtschaft isch güet, un dr alte Stallknächt vertraüt g'si ; un mer het'm riehig derfe Ross un G'schirr avertraüe.

'S Bachludis
Mariannle, Charles
Zumstein, 1986

Modèle Tesseract

allemand

Mit d^äm G'spr^äch sin die zwei uf Altkilch ku, si h^än nit g'wisst wie. Dr Lüwi isch glich in's "Reglis" ufe g'fahre, un het si Füehrwark dert îg'stellt. Vater un Grossvater han scho allew^{'i}l dert ^{'i}kehrt, un o ^är s^älber, wenn 'r ⁱne g'fahre isch. D'Wirtschaft isch güet, un dr alte Stallkn^ächt vertraüt g'si ; un mer het'm riehig derfe Ross un G'schirr avertraüe.

français

Mit d^àm G'spr^àch sin die zwei uf Altkilch ku, si h^àn nit g'wisst wie. Dr Lüwi isch glich in's "Reglis" ufe g'fahre, un het si Füehrwark dert îg'stellt. Vater un Grossvater han scho allew^îl dert îkehrt, un o ^àr s^àlber, wenn 'r ^îne g'fahre isch. D'Wirtschaft isch güet, un dr alte Stallkn^àcht vertraüt g'si ; un mer het'm riehig derfe Ross un G'schirr avertraüe.

Perspectives

- Poursuite du travail sur l'OCR pour améliorer les résultats
 - Entraînement de modèles pour l'alsacien
 - Combinaison avec les modèles existants
- Prochaine application visée : annotation morpho-syntaxique
 - Définition de jeux d'étiquettes
 - Annotation de corpus pour l'évaluation
 - Constitution de lexiques morphosyntaxiques
 - Transfert technologique d'outils existants

Merci vielmols !