

Genre des substantifs en tchèque: l'ambiguïté de ses marqueurs formels et le diagnostic automatique des erreurs du point de vue de son acquisition par les apprenants francophones

Ivan Šmilauer

▶ To cite this version:

Ivan Śmilauer. Genre des substantifs en tchèque: l'ambiguïté de ses marqueurs formels et le diagnostic automatique des erreurs du point de vue de son acquisition par les apprenants francophones. Colloque des doctorants et des jeunes chercheurs (COLDOC09) "Ambiguïté dans les sciences du langage", Jun 2009, Nanterre, France. hal-01375644

HAL Id: hal-01375644 https://inalco.hal.science/hal-01375644

Submitted on 20 Jan 2017

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Genre des substantifs en tchèque

Ivan ŠMILAUER LaLIC-CERTAL, INALCO (Paris) smilauer@cetlef.fr

session poster COLDOC'09

l'ambiguïté de ses marqueurs formels du point de vue de son acquisition par les apprenants francophones et le diagnostic automatique des erreurs

Ambiguïté morphologique du point de vue d'un apprenant de LE

Les effets de l'ambiguïté, définie dans un modèle linguistique stratificationnel par la multitude de fonctions prises par un élément de niveau n au niveau n+1 (cf. Sgall et al. 1986) peuvent être également étudiés dans les productions langagières des apprenants d'une langue étrangère. Nous l'illustrons par l'étude des phénomènes d'ambiguïté entre le niveau phonologico-graphémique et le niveau morphologique du tchèque qui résultent en différentes erreurs de genre dans les productions des apprenants francophones.

Dans la plupart des cas, la relation entre un substantif et son genre grammatical ne présente pour un locuteur natif aucun caractère d'ambiguïté car cette information fait partie de son stock « naturel » de connaissances lexicales. La situation est différente pour un apprenant étranger confronté à un substantif au genre inconnu. Celui-ci se trouve devant deux solutions : rechercher et intégrer l'information sur genre « en dur » dans son vocabulaire ou essayer de trouver et d'appliquer des règles basées sur la forme qui permettrait d'attribuer le genre sans recours à un lexique. En se servant de cette seconde technique, l'apprenant se trouve nécessairement heurté aux limites des règles établies dont la surgénéralisation produit des erreurs.

C'est ainsi que nous considérons une erreur de genre comme le résultat de l'attribution d'une fonction inadéquate à une forme ambiguë par rapport aux règles dont se sert l'apprenant. De ce point de vue, nous estimons que les substantifs sont plus ou moins ambigus par rapport au genre.

Le tchèque possède quatre valeurs pour la catégorie grammaticale de genre : masculin animé et inanimé, féminin, neutre. A l'aide des critères formels et sans recours à une connaissance d'ordre lexicale, il est possible de définir des règles pour l'attribution du genre en fonction des marqueurs situés à la périphérie droite des lexèmes (la désinence casuelle vocalique du lemme ou la terminaison consonantique du radical). Ces marqueurs permettent avec un degré variable de certitude d'identifier le genre des substantifs : le nominatif singulier terminé par "-o" désigne sans ambiguïté un substantif neutre, le nominatif singulier terminé par "-a" désigne le plus souvent un féminin, mais aussi, marginalement, un masculin etc.

Dans une production d'apprenant, une erreur de genre peut se manifester (a) par une **erreur d'accord** des éléments dépendants de ce substantif, (b) par l'attribution d'un type paradigmatique erroné qui se manifeste lors de la déclinaison du substantif en question. Par exemple, le substantif masculin "stůl" (table) peut être interprété par l'apprenant en tant que féminin car un substantif terminé par la consonne "-l" peut être masculin, mais également féminin, voir par exemple "sůl" (sel). Ainsi, cette analyse erronée peut mener à des productions du type (a), par exemple "velká stůl", qui pourrait être traduit littéralement comme « grand table » ; ou à des erreurs de type (b), par exemple le génitif singulier "stoli" d'après le type de déclinaison consonantique féminin correspondant, au lieu de "stolu" d'après le type correct masculin. C'est sur ce deuxième type d'erreur que nous nous concentrons ici.

L'étude de l'ambiguïté des marqueurs formels du genre est intéressante également du point de vue pédagogique car cette propriété est en relation directe avec le degré de diagrammaticité du substantif (caractère indexical de sa construction formelle par rapport à sa signification) qui peut faciliter le traitement cognitif du lexème par l'apprenant.

Nous apportons ici des éléments pour étudier, dans un cadre expérimental spécifique (exercices de déclinaison sur une plate-forme ELAO), les erreurs de genre produites par des apprenant francophones et le rôle de l'ambiguïté des marqueurs formels des substantifs. Nous supposons que les erreurs de genre les plus fréquentes devraient être celles qui occurrent dans les substantifs formellement ambigus par rapport à la catégorie du genre.

En nous basant sur une étude quantitative, nous présentons d'abord l'ambiguïté des marqueurs de genre en tchèque. Ensuite, nous introduisons l'application CETLEF.fr (un outil d'enseignement du tchèque assisté par ordinateur avec un diagnostic automatique des erreurs) qui sert pour la collecte et l'analyse des productions langagières des apprenants.

Finalement, nous présentons les données recueillies à l'aide de CETLEF.fr et nous proposons un algorithme modélisant l'attribution du genre à un substantif par les choix d'ordres formels, sémantiques et lexicaux effectués par l'apprenant pendant l'analyse de la forme au genre inconnu. Cet algorithme pourrait servir à des fins pédagogiques mais également pour l'affinement du diagnostic automatique existant au sein de CETLEF.fr.



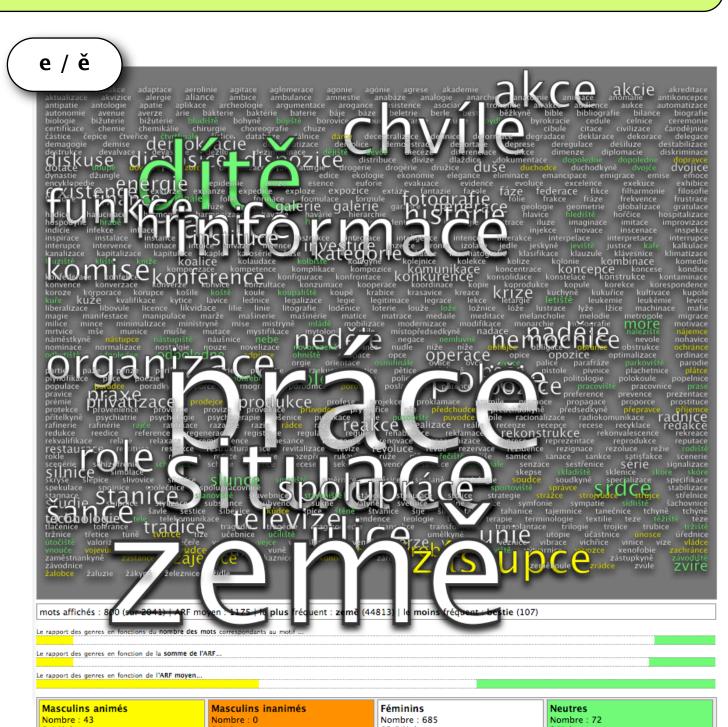
La définition d'un marqueur formel est en principe arbitraire – le marqueur ultime d'un lexème serait sa forme entière, ce qui revient à même que de construire un lexique.

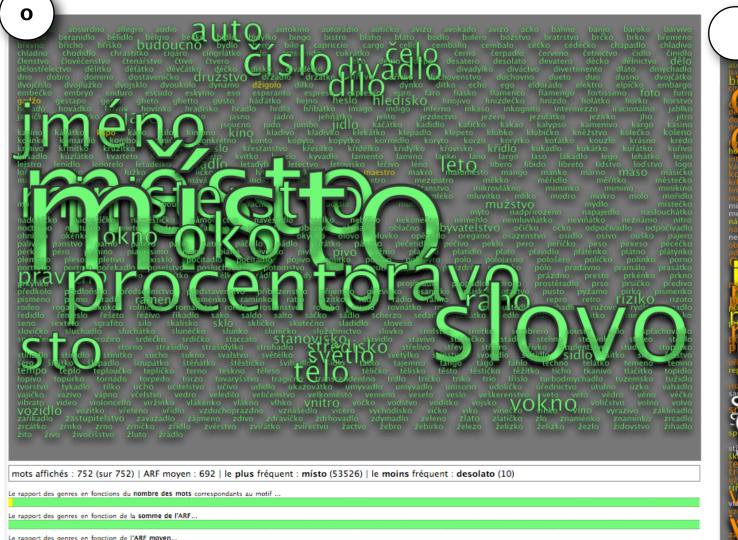
Etant donné le souci pour l'économie des moyens employés dans les règles, nous avons défini 7 marqueurs différents, couvrant uniquement le dernier phonème/graphème du mot : 4 marqueurs vocaliques (o, í, a, e/ě) et 3 classes de marqueurs consonantiques (consonne dure, consonne mixte, consonne

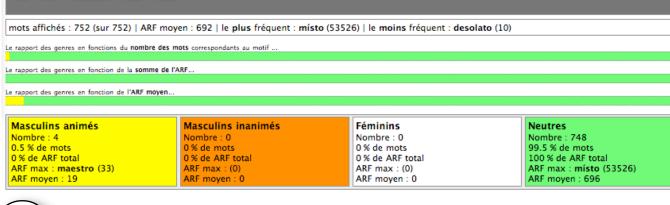
Pour chaque marqueur, nous avons recherché les substantifs correspondants dans un lexique annoté contenant les 21 986 substantifs tchèques les plus fréquents dans le Corpus National Tchèque (cf. www.korpus.cz, Čermák 1997, Čermák & Křen 2004, 2005). Pour une meilleure visualisation des données qui permet une appréciation générale de l'ambiguïté des marqueurs, nous avons attribué à chaque mot une couleur en fonction de son genre (masculins animés : jaunes ; masculins inanimés : oranges ; féminins : blanc ; neutres : verts) et une taille en fonction de sa fréquence moyenne réduite (ARF - Average Reduced Frequency, cf. Hlaváčová and Savický 2002). Seulement les premiers 800 substantifs les plus fréquents sont affichés. Le rapport entre les différents genres pour un marqueur donné offre une possibilité de mesurer le taux de son ambiguïté.

Il serait possible d'établir des marqueurs à l'aide des segments plus importants, notamment ceux qui peuvent être identifiés comme des suffixes de dérivation. Néanmoins, cet affinement, qui va au détriment de la maniabilité des règles par l'apprenant, n'apporte pas forcément une baisse significative de leur ambiguïté.

A côté des critères formels, un moyen possible de lever l'ambiguïté d'un substantif est son analyse par rapport au genre sémantique (naturel) de l'objet désigné. Dans le cas des personnes et des animaux mâles et femelles, le genre grammatical et sémantique sont identiques, ainsi que pour les enfants humains et animaux (neutre). Dans le cas des objets, cette relation est complètement arbitraire. Une fois identifiée, la différence entre les masculins animés et inanimés est gérée par des critères du même ordre.











mots affichés : 800 (sur 4809) | ARF moyen : 2612 | le plus fréquent : doba (61603) | le moins fréquent : konzerva (393) Le rapport des genres en fonctions du nombre des mots correspondants au motif Le rapport des genres en fonction de l'ARF moyen Nombre: 761 95.1 % de mots

96.6 % de ARF total

ARF max : doba (61603)

2.8 % de ARF total

ARF max : předseda (1421



mots affichés : 800 (sur 1601) | ARF moyen : 913 | le plus fréquent : konec (31256) | le moins fréquent : harmonikář (51

0.5 % de mots 7 % de ARF total

Nombre: 221

27.6 % de mots

35.5 % de ARF total ARF max : věc (26962)

ARF moyen: 1172

0 % de mots 0 % de ARF total

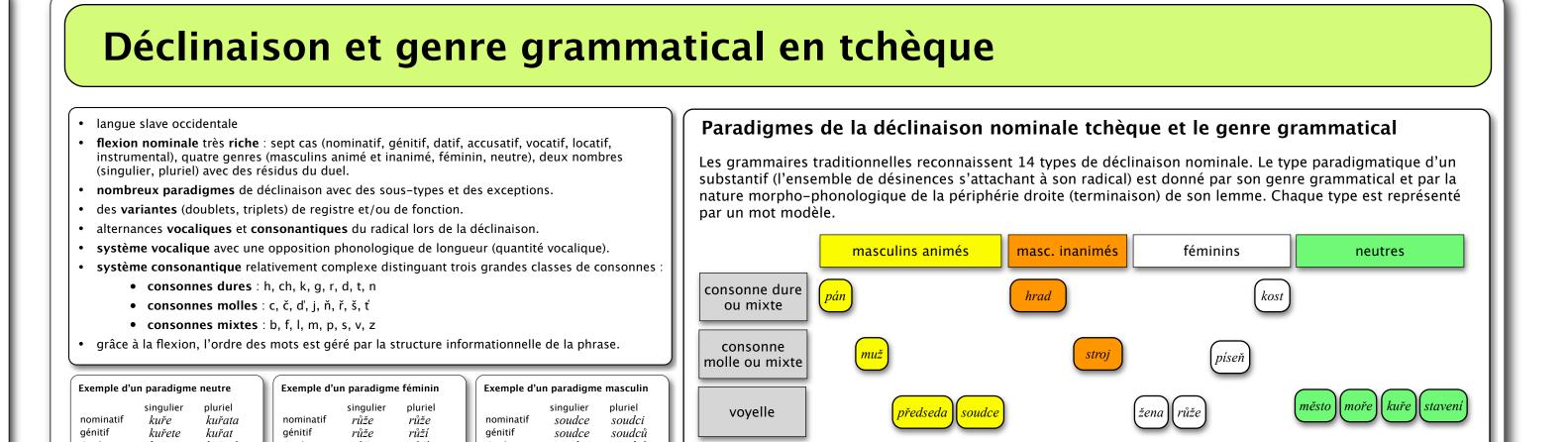
e rapport des genres en fonction de l'ARF moyen.

27.5 % de ARF tota

ARF max : **muž** (2213

5 % de ARF total

RF max : **téma** (7208 RF moyen : 2217



Erreur de genre sur CETLEF.fr

soudce

soudce

CETLEF.fr est un outil d'enseignement de langue assisté par ordinateur (ELAO) proposant des exercices de déclinaison tchèque : la tâche de l'apprenant est de créer la forme fléchie d'un lemme en fonction de son contexte syntaxique au sein d'une proposition donnée. Une erreur commise dans une telle tâche est appelée erreur de déclinaison. Erreur de genre est un des types d'une telle

accusatif

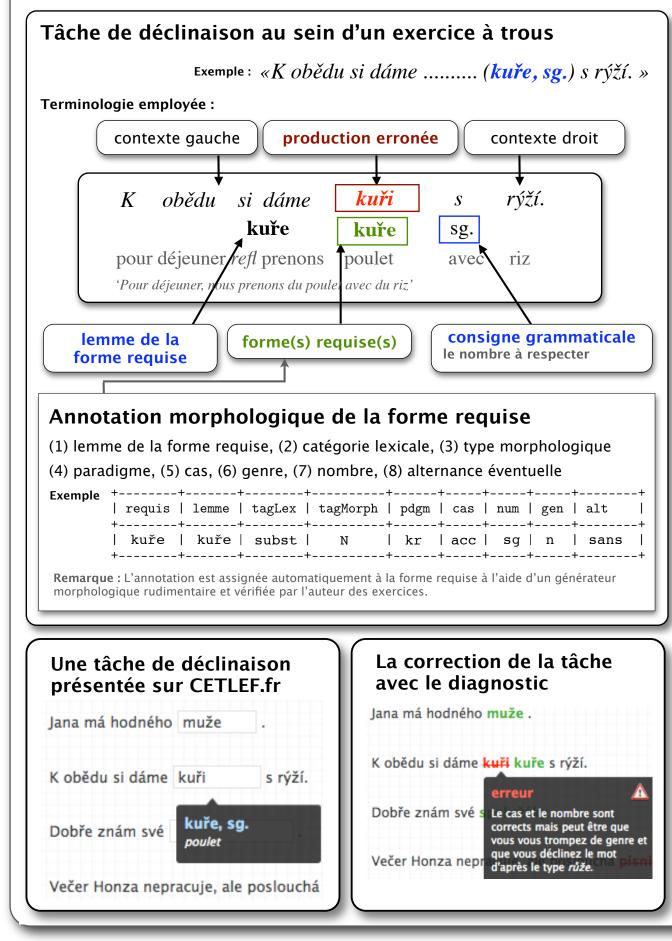
accusatif

alternances morphématiques.

Une erreur de déclinaison est considérée non pas comme un phénomène aléatoire mais comme le résultat d'une activité succombant à des règles d'ordre linguistique et cognitif. Le diagnostic automatique des erreurs, implémenté sur cette plate-forme, est basé sur cette hypothèse. Il a nécessité l'élaboration d'un modèle formel spécifique de la déclinaison contenant un classement des types paradigmatiques et des règles pour la réalisation des

CETLEF.fr illustre les possibilités d'un riche modèle morphologique et des techniques de TAL employées dans un outil d'enseignement de langue assisté par ordinateur, voir par exemple (Heift & Schulze 2007).

Du point de vue d'une recherche sur l'acquisition de langue étrangère, CETLEF.fr sert comme un outil pour la compilation d'un corpus d'erreurs. Par rapport aux productions libres, l'analyse des erreurs recueillies au sein des exercices grammaticaux permet un meilleur contrôle sur le volume de données pertinentes, car les productions doivent nécessairement contenir les phénomènes qui ont été établis comme l'objet de l'investigation.



Définition d'une erreur de déclinaison Comme erreur de déclinaison est considérée toute production différente de la (les) forme(s) requise(s) ayant une interprétation morphologique. Une interprétation morphologique peut être assignée à une production erronée, si elle correspond à une des formes hypothétiques générées à partir du radical de la forme emploi d'une **désinence inappropriée** appartenant (a) au paradigme de la forme requise, (b) appartenant à un paradigme

attribut | type d'erreur

simple concaténation du radical et d'une désinence sans la réalisation d'une alternance obligatoire • erreurs d'ordre **graphique** : erreurs de diacritique, de casse

productions erronées sans interprétation morphologique erronées avec une nterprétation morphologiq

production erronée

kuři

forme requise

kuře

HOR EXT|rz|i|acc|sg|f|sans|erreur gen

Le cas et le nombre sont corrects mais peut être que vous vous trompez de genre et que vous

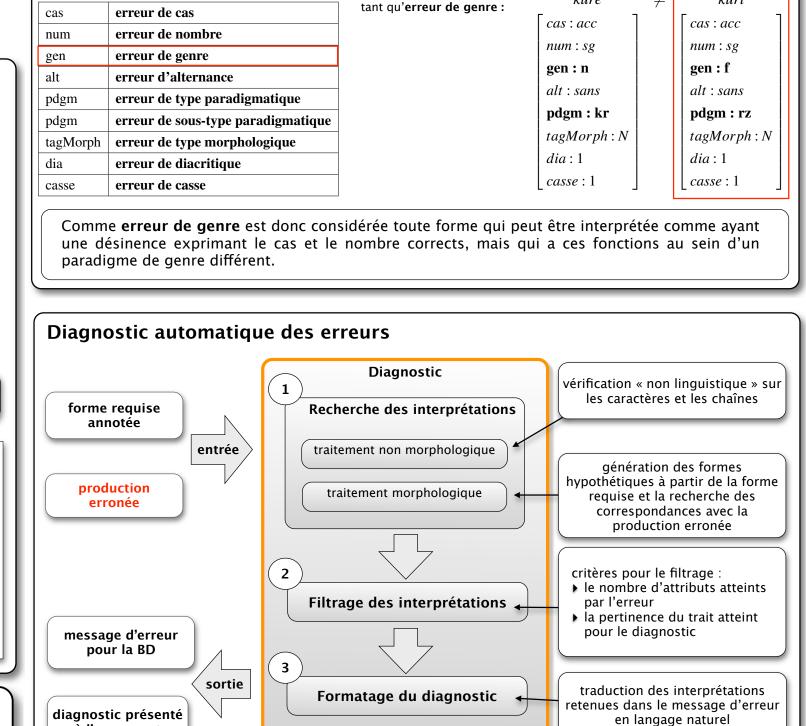
déclinez le mot d'après le type růže

Typologie des erreurs d'après l'attribut atteint et l'erreur de genre La typologie formelle des erreurs de déclinaison est basée sur les différences entre les valeurs des attributs morphologiques de la forme requise et celles de la forme hypothétique correspondant à la production erronée.

Exemple de l'interprétation

. plusieurs interprétations d'une forme erronée sont possibles, notamment à cause de l'homonymie des désinences.

Remarque: Au sein de chaque type, il existe de nombreux sous-types, définis par les différences par rapport au mot modèle. Cependant, ces différences touchent rarement la forme du lemme qui est déterminante pour l'identification du type paradigmatique (il existe notamment des



Diagnostic

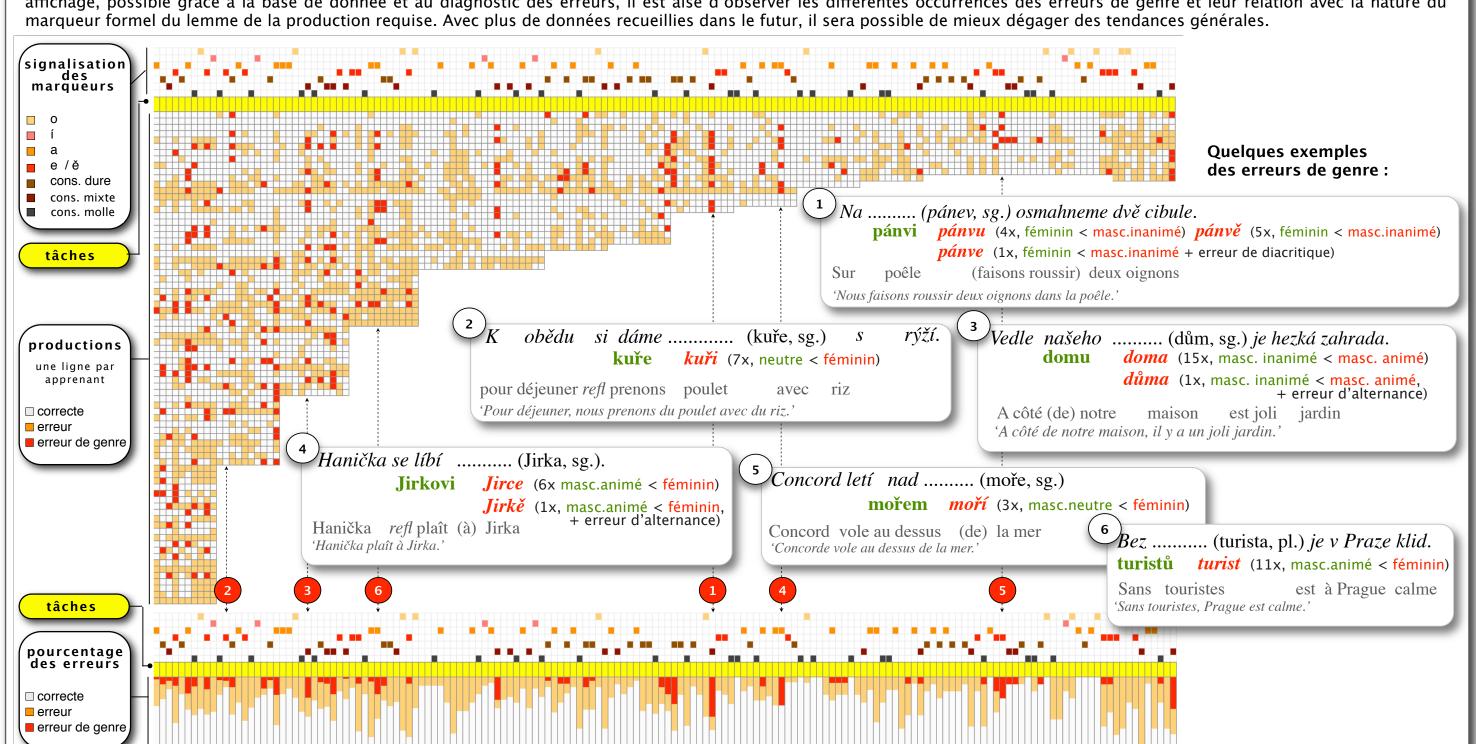
Entre 85 et 90 % de productions collectées au sein de CETLEF sont interprétées comme une erreui

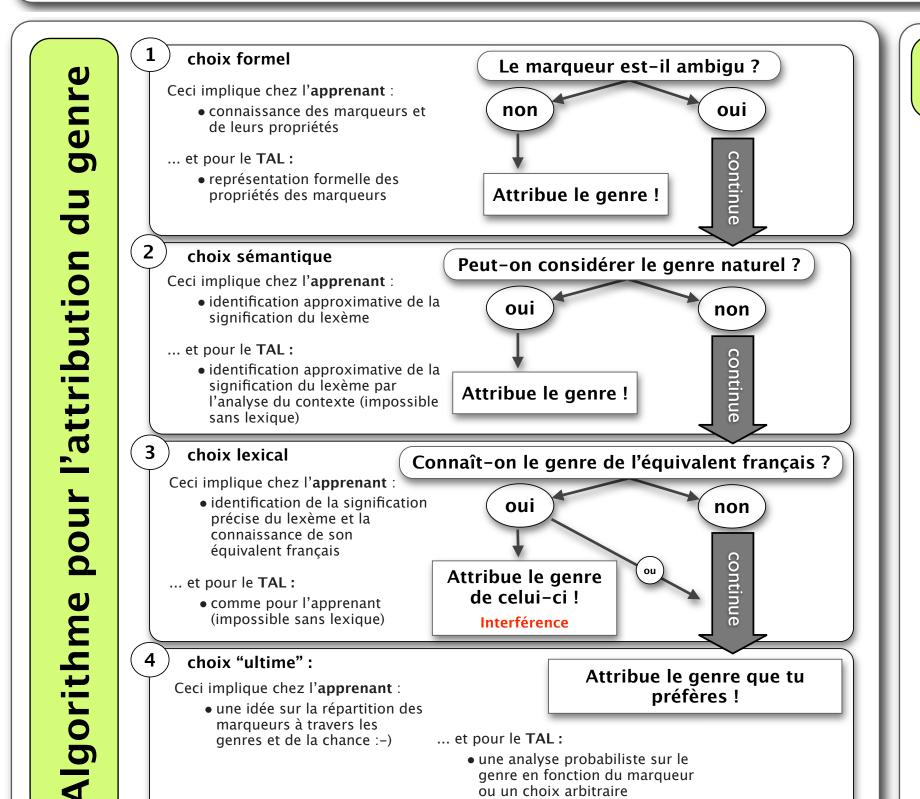


Le graphique ci-dessous présente les erreurs de genres commises par les apprenants dans 162 tâches de déclinaison au sein de 16 exercices sur CETLEF.fr. A l'aide de cette affichage, possible grâce à la base de donnée et au diagnostic des erreurs, il est aisé d'observer les différentes occurrences des erreurs de genre et leur relation avec la nature du

à l'apprenant

kuře|subst|N|kr|acc|sg|n|sans





Bibliographie

Čermák, F. (1997).

Hudson, Paris.

- Allerton, D. J., Tschichold, C. et Wieser, J., éditeurs (2005). Linguistics, Language Learning and Language Teaching. Schwabe, Basel.
- Čermák, F. & Křen, M. (2004). Frekvenční slovník češtiny. Nakladatelství Lidové Noviny, Praha. Čermák, F. & Křen, M. (2005). New generation corpus-based frequency dictionaries: The case of czech. International Journal of Corpus Linguistics, 10:453-467.
- Czech national corpus: A case in many contexts. International Journal of Corpus Linguistics, 2(2):181-197. Heift, T. & Schulze, M. (2007). Errors and Intelligence in Computer- Assisted Language Learning:
- Parsers and Pedagogues. Routledge, UK. Hlaváčová, J. & Savický, P. (2002). Measures of word commonness. Journal of Quantitative Linguistics, 9(3):
- Klanten, R. éditeur (2009)
- Sgall, P., Hajičová, E. & Panevová, J. (1986). The Meaning of the Sentence in its Semantic and Pragmatic Aspects. D. Reidel Publishing Company, Praha - Amsterdam.

Data Flow. Design graphique et visualisation d'information. Thames &

Acquisition du tchèque par les francophones : analyse automatique des erreurs de déclinaison. The Prague Bulletin of Mathematical Linguistics