



HAL
open science

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 Šmilauer. 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 multiplicité 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 phonologique-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 n'est pas présente pour un locuteur natif au 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 le genre « en dur » dans son vocabulaire ou essayer de trouver et d'appliquer des règles basées sur la forme qui permettraient 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. À 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 vocale du lemme ou la terminaison consonnantique du radical). Ces marqueurs permettent avec un degré variable de certitude d'identifier le genre des substantifs : le nominatif singulier terminé par « -a » désigne sans ambiguïté un substantif neutre, le nominatif singulier terminé par « -s » 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 féminin, mais également féminin, voir par exemple « stůl » (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 « stůli » d'après le type de déclinaison consonnantique féminin correspondant, au lieu de « stůlu » 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 diagrammatisme 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 apprenants francophones et le rôle de l'ambiguïté des marqueurs formels des substantifs. Nous supposons que les erreurs de genre les plus fréquentes devaient être celles qui ont lieu 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 automatisant l'attribution du genre à un substantif par le 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.

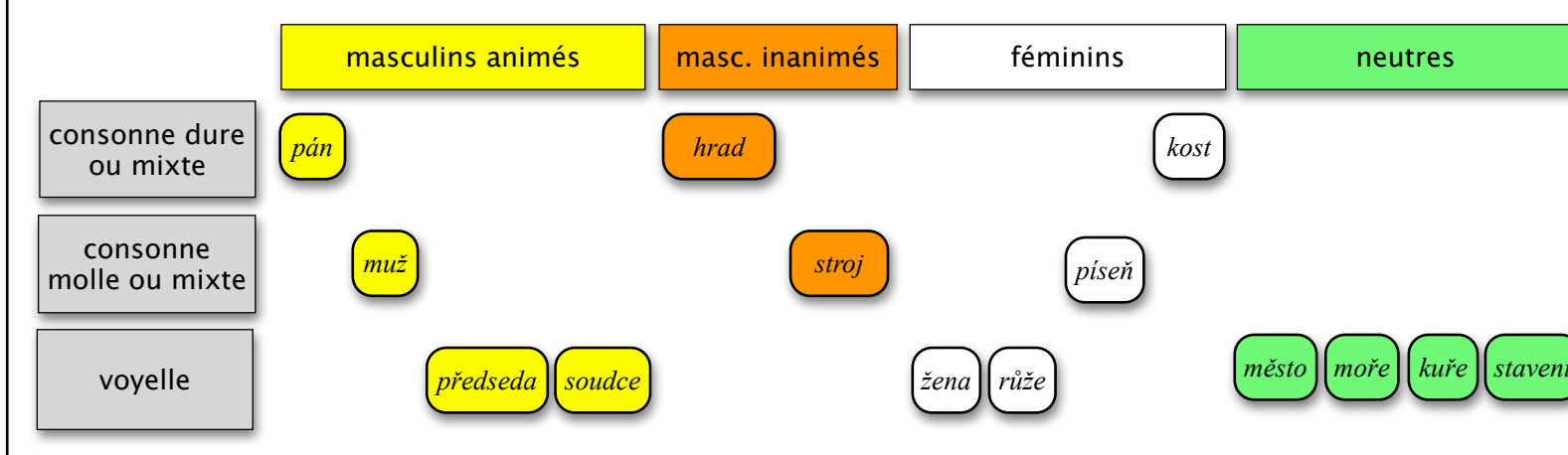
Déclinaison et genre grammatical en tchèque

- langue slave occidentale
- flexion nominale très riche : sept cas (nominatif, génitif, datif, accusatif, vocatif, locatif, instrumental), quatre genres (masculin animé et inanimé, féminin, neutre), deux nombres (singulier, pluriel) avec des résidus du duel
- nombreux paradigmes de déclinaison avec des sous-types et des exceptions.
- des variantes (doubles, triples) de registre et/ou de fonction.
- alternances vocales et consonnantiques du radical lors de la déclinaison.
- alternance vocale avec une opposition phonologique de longueur (quantité vocale).
- système consonnantique relativement complexe distinguant trois grandes classes de consonnes :
 - consonnes dures : h, ch, k, g, r, d, t, n
 - consonnes molles : b, f, l, m, p, s, v, z
 - consonnes mixtes : c, č, ř, š, ž, x, y
- grâce à la flexion, l'ordre des mots est géré par la structure informationnelle de la phrase.

Exemple d'un paradigme neutre	Exemple d'un paradigme féminin	Exemple d'un paradigme masculin
nominatif: kůř génitif: kůř datif: kůři accusatif: kůř locatif: kůři instrumental: kůřemi	nominatif: kůř génitif: kůř datif: kůři accusatif: kůř vocatif: kůři locatif: kůři instrumental: kůřemi	nominatif: kůř génitif: kůř datif: kůři accusatif: kůř vocatif: kůři locatif: kůři instrumental: kůřemi

Paradigmes de la déclinaison nominale tchèque et le genre grammatical

Les grammaires traditionnelles reconnaissent 14 types de déclinaison nominale. Le type paradigmatique d'un substantif (l'ensemble de désinences s'attachant à son radical) est donné par son genre grammatical et par la nature morpho-phonologique de la périphérie droite (terminaison) de son lemme. Chaque type est représenté par un mot modèle.



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 souvent la forme du lemme qui est déterminante pour l'identification du type paradigmatique (il existe notamment des différences entre les types féminins « -ě » et « -y » et il ne serait pas utile de les préciser en contexte ici).

Erreur de genre sur CETLEF.fr

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 flexive 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 erreur.

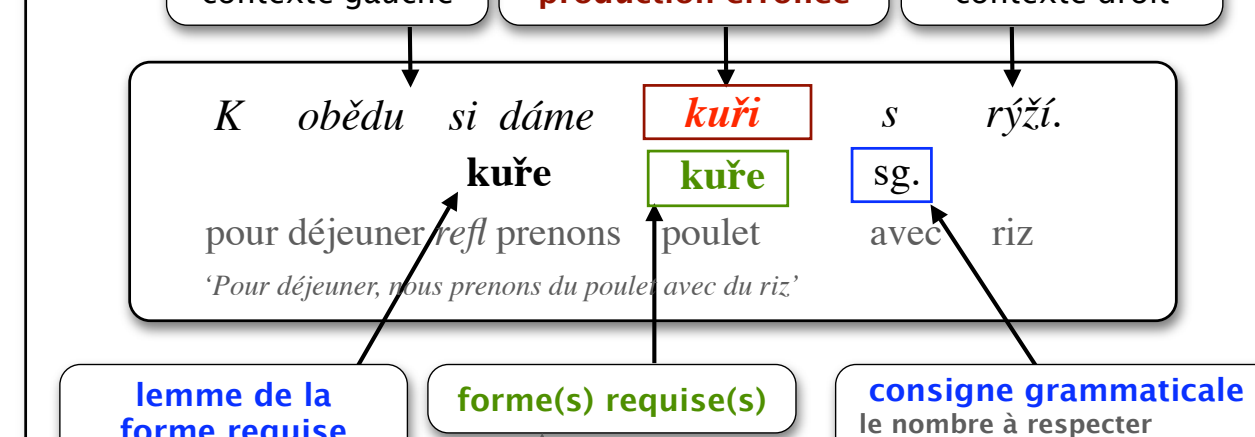
Une erreur de déclinaison est considérée non pas comme un phénomène aléatoire mais comme le résultat d'une activité succombant à des règles d'ordre linguistique et cognitif. Le diagnostic automatique des erreurs, implémenté sur cette plate-forme, est basé sur cette hypothèse. Il a nécessité l'élaboration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

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.

Tâche de déclinaison au sein d'un exercice à trous

Exemple : « K obědu si dáme (kůř, sg.) s rýží. »



Annotation morphologique de la forme requise

(1) lemme de la forme requise, (2) catégorie lexicale, (3) type morphologique (4) paradigme, (5) cas, (6) genre, (7) nombre, (8) alternance éventuelle

Exemple : kůř | kůř | subst | N | ž | acc | sg | n | sans

Remarque : L'annotation est assignée automatiquement à la forme requise à l'aide d'un générateur morphologique rudimentaire et vérifiée par l'auteur des exercices.

Une tâche de déclinaison présentée sur CETLEF.fr

Jana má hodného muže.

K obědu si dáme kůř s rýží.

Dobře znám své kůř, sg. poulet.

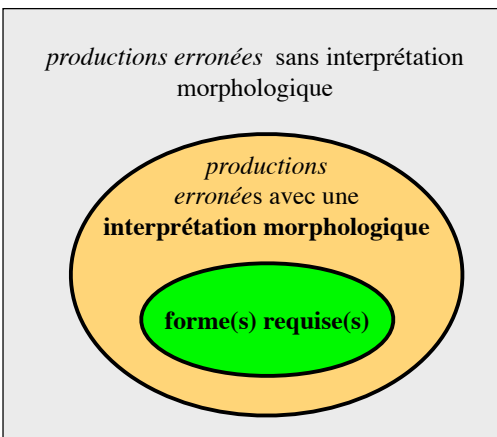
Večer Honza nepracuje, ale poslouchá.

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 requise par :

- emploi d'une désinence inappropriée appartenant (a) au paradigme de la forme requise, (b) appartenant à un paradigme différent
- simple concaténation du radical et d'une désinence sans la réalisation d'une alternance obligatoire
- erreurs d'ordre graphique : erreurs de diacritique, de casse

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



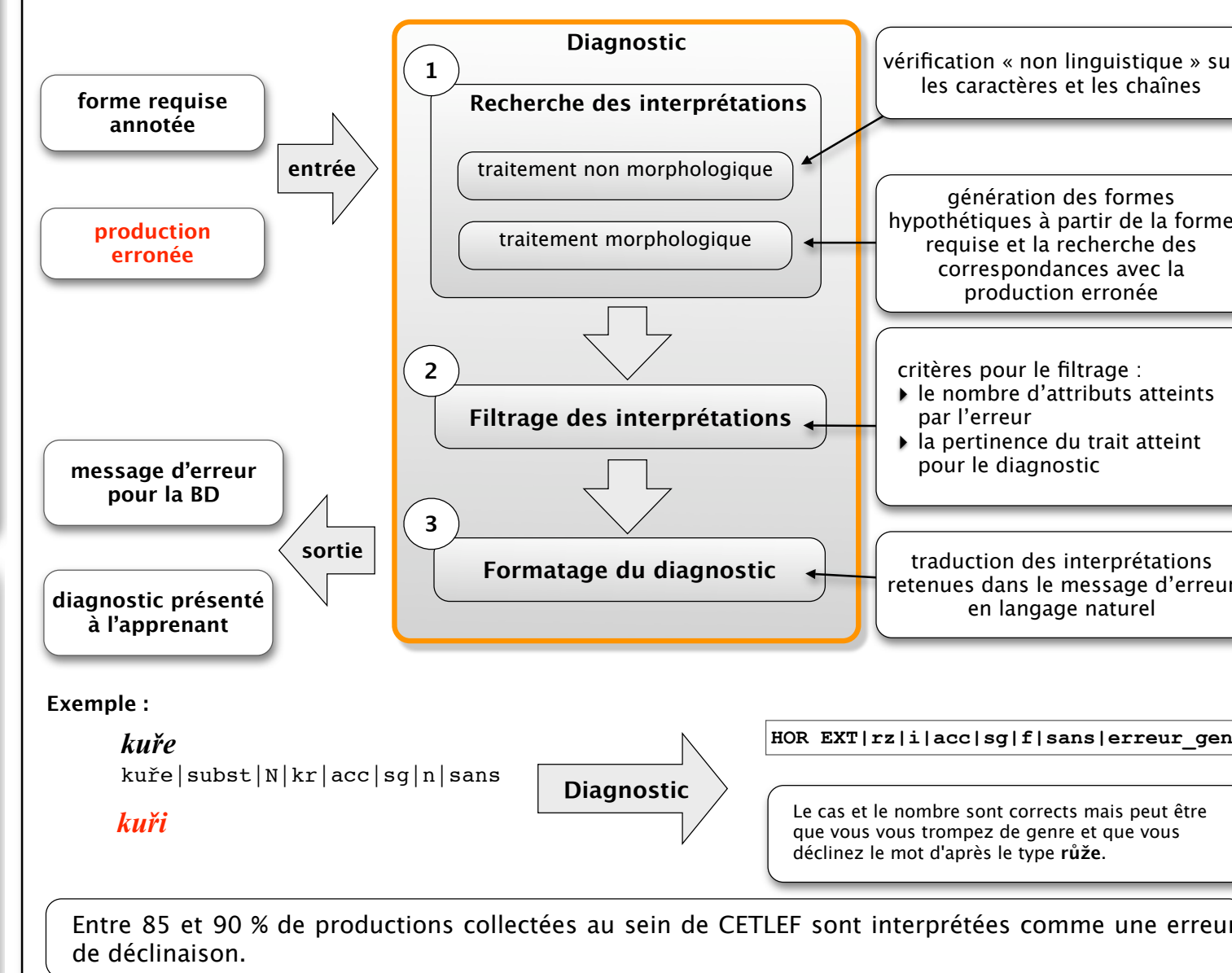
Typologie des erreurs d'après l'attribut atteint et l'erreur de genre

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.

attribut	type d'erreur	Exemple de l'interprétation d'une production erronée en tant qu'erreur de genre :	forme requise	production erronée
cas	erreur de cas	kůř	kůř	kůři
num	erreur de nombre	kůř	kůř	kůři
gen	erreur de genre	kůř	kůř	kůři
alt	erreur d'alternance	kůř	kůř	kůři
pdgm	erreur de type paradigmatique	kůř	kůř	kůři
pdgm	erreur de sous-type paradigmatique	kůř	kůř	kůři
tagMorph	erreur de type morphologique	kůř	kůř	kůři
dia	erreur de diacritique	kůř	kůř	kůři
casse	erreur de casse	kůř	kůř	kůři

Comme erreur de genre est donc considérée toute forme qui peut être interprétée comme ayant une désinence exprimant le cas et le nombre corrects, mais qui a ces fonctions au sein d'un paradigme de genre différent.

Diagnostic automatique des erreurs



Entre 85 et 90 % de productions collectées au sein de CETLEF sont interprétées comme une erreur de déclinaison.

Étude quantitative des marqueurs formels du genre

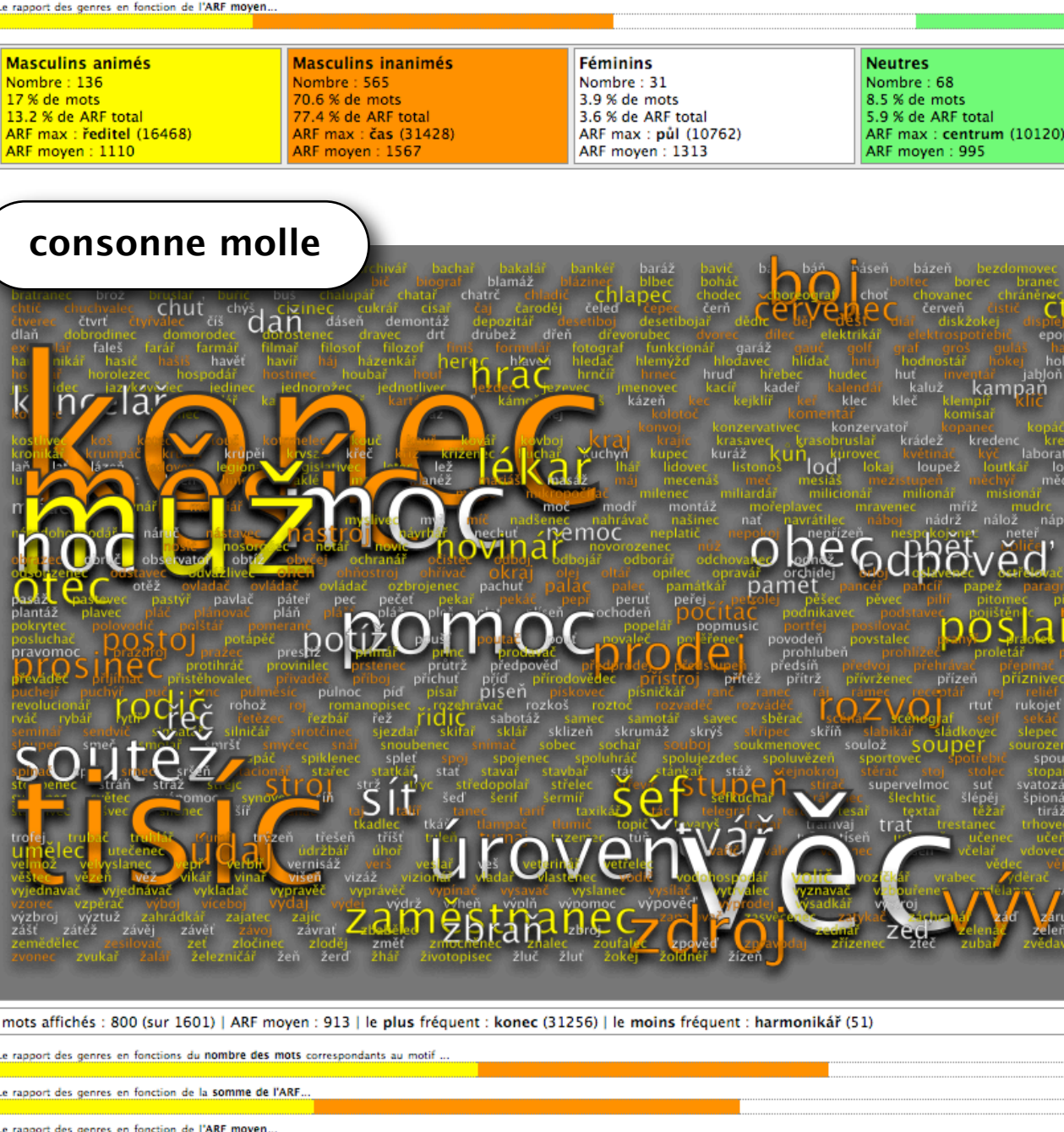
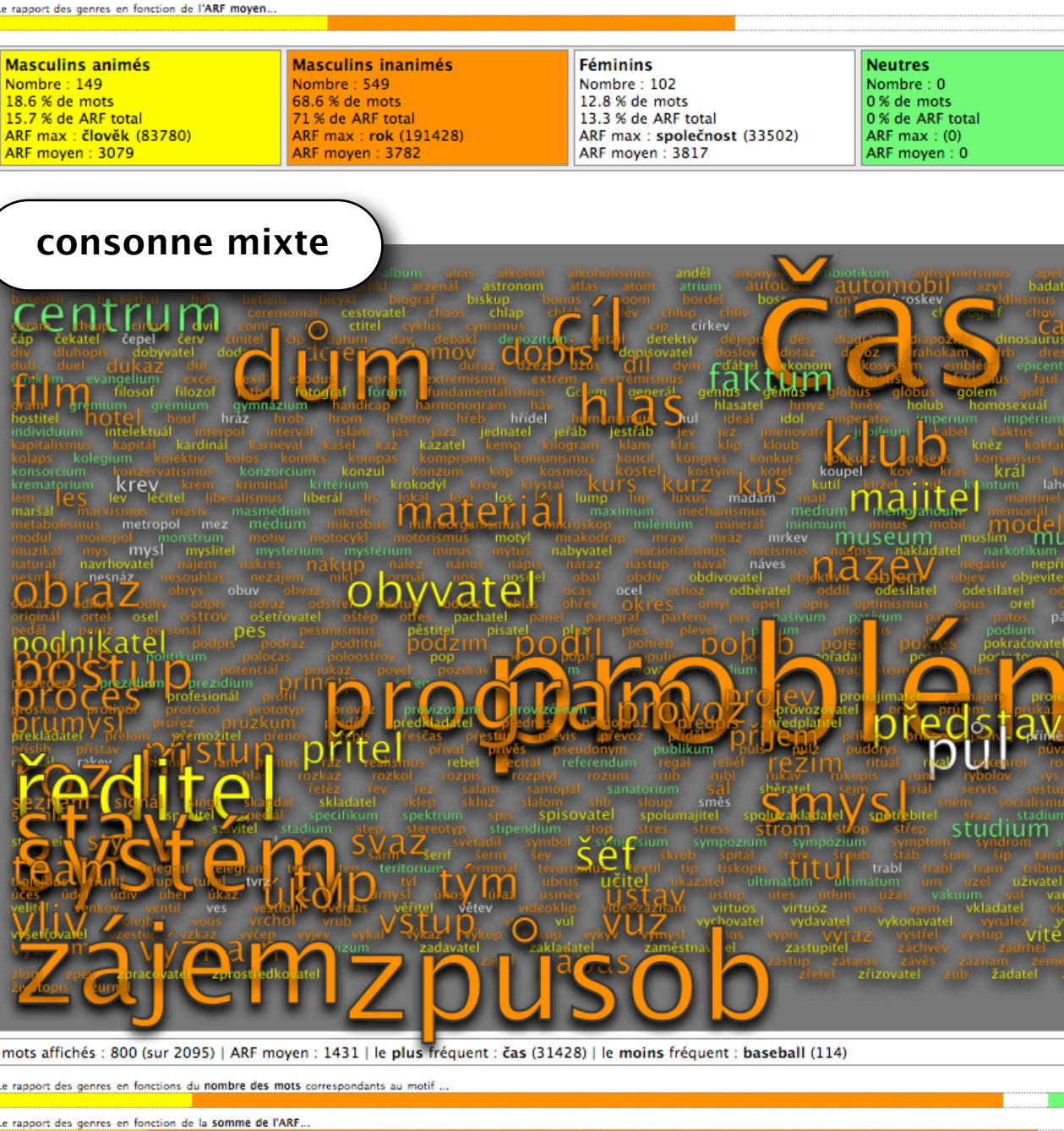
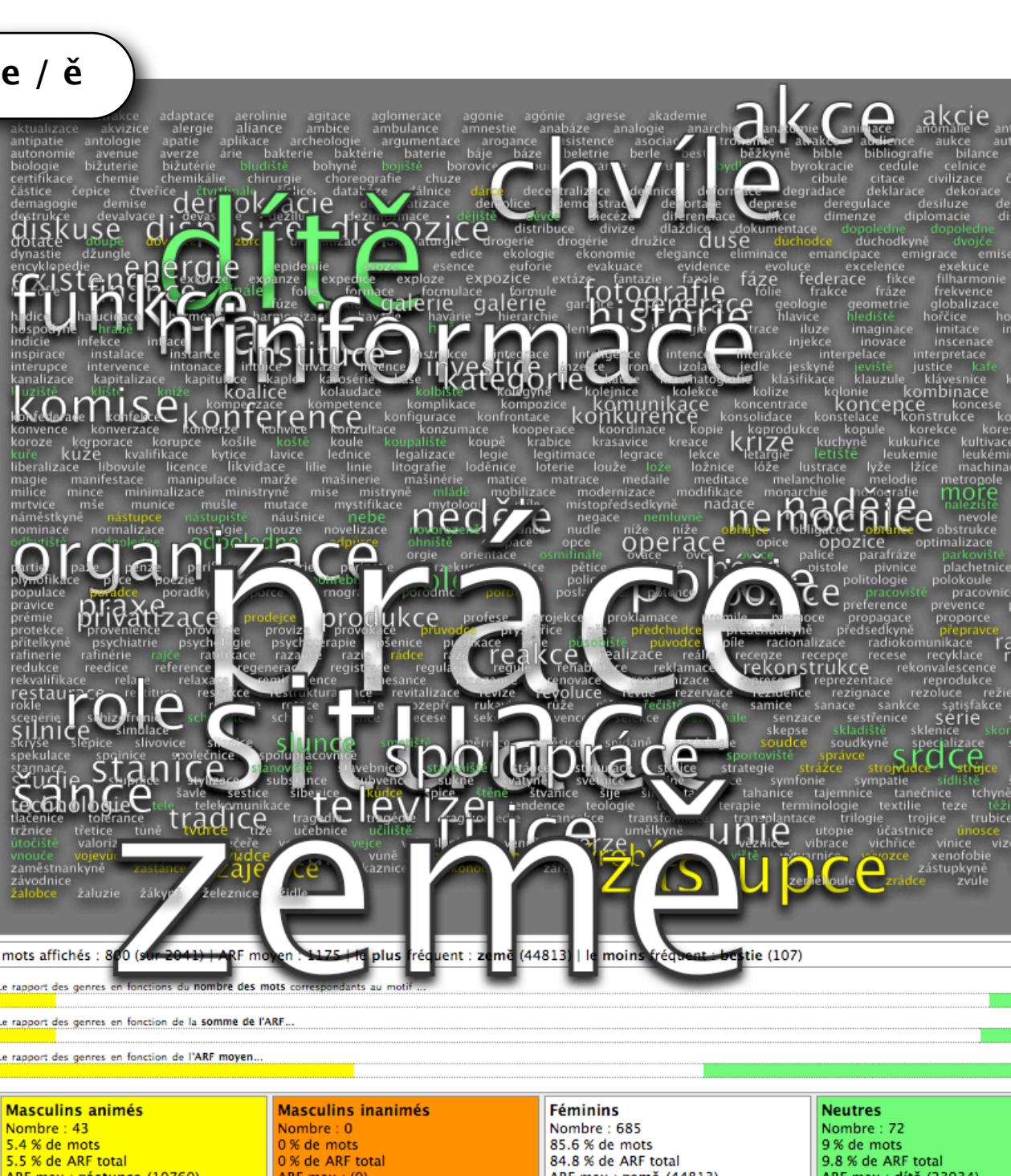
La définition d'un marqueur formel est en principe arbitraire – le marqueur ultime d'un lexème serait sa forme entière, ce qui revient à même que de construire un lexique.

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 vocaux (o, i, a, e/ě) et 3 classes de marqueurs consonnantiques (consonne dure, consonne mixte, consonne molle).

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 Tchéco (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 (masculin animé : jaunes ; masculin inanimé : oranges ; féminins : blanc ; neutres : verts) et une taille en fonction de sa fréquence moyenne réduite (ARF - Average Reduced Frequency, cf. Hlavčová & Savičková 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.

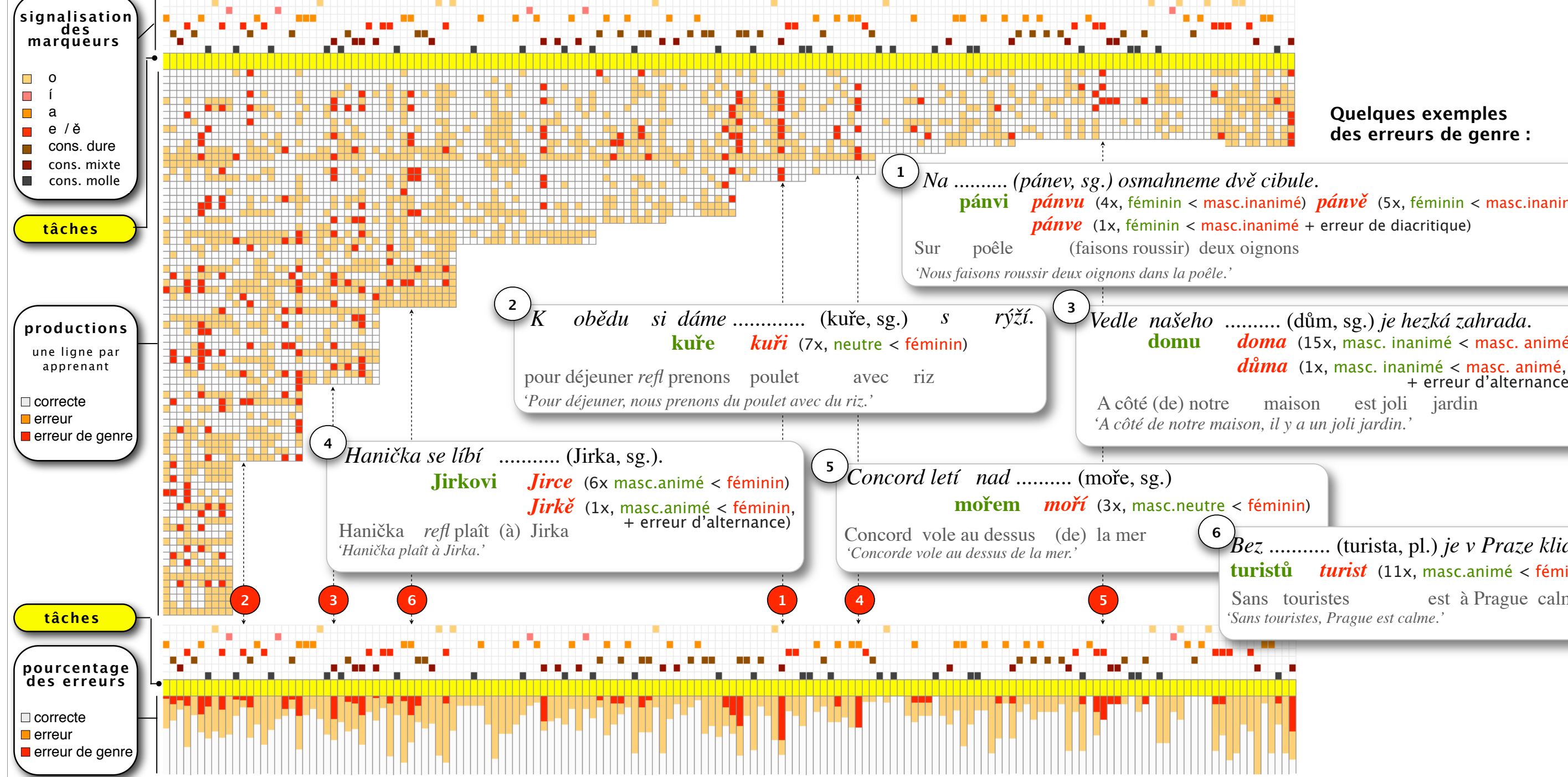
Il serait possible d'établir des marqueurs à l'aide des segments plus importants, notamment ceux qui peuvent être identifiés comme des suffixes de dérivation. Néanmoins, cet affinement, qui va au détriment de la maniableté des règles par l'apprenant, n'apporte pas forcément une baisse significative de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

À côté des critères formels, un moyen possible de lever l'ambiguïté d'un substantif est son analyse par rapport au genre sémantique (nature) 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.

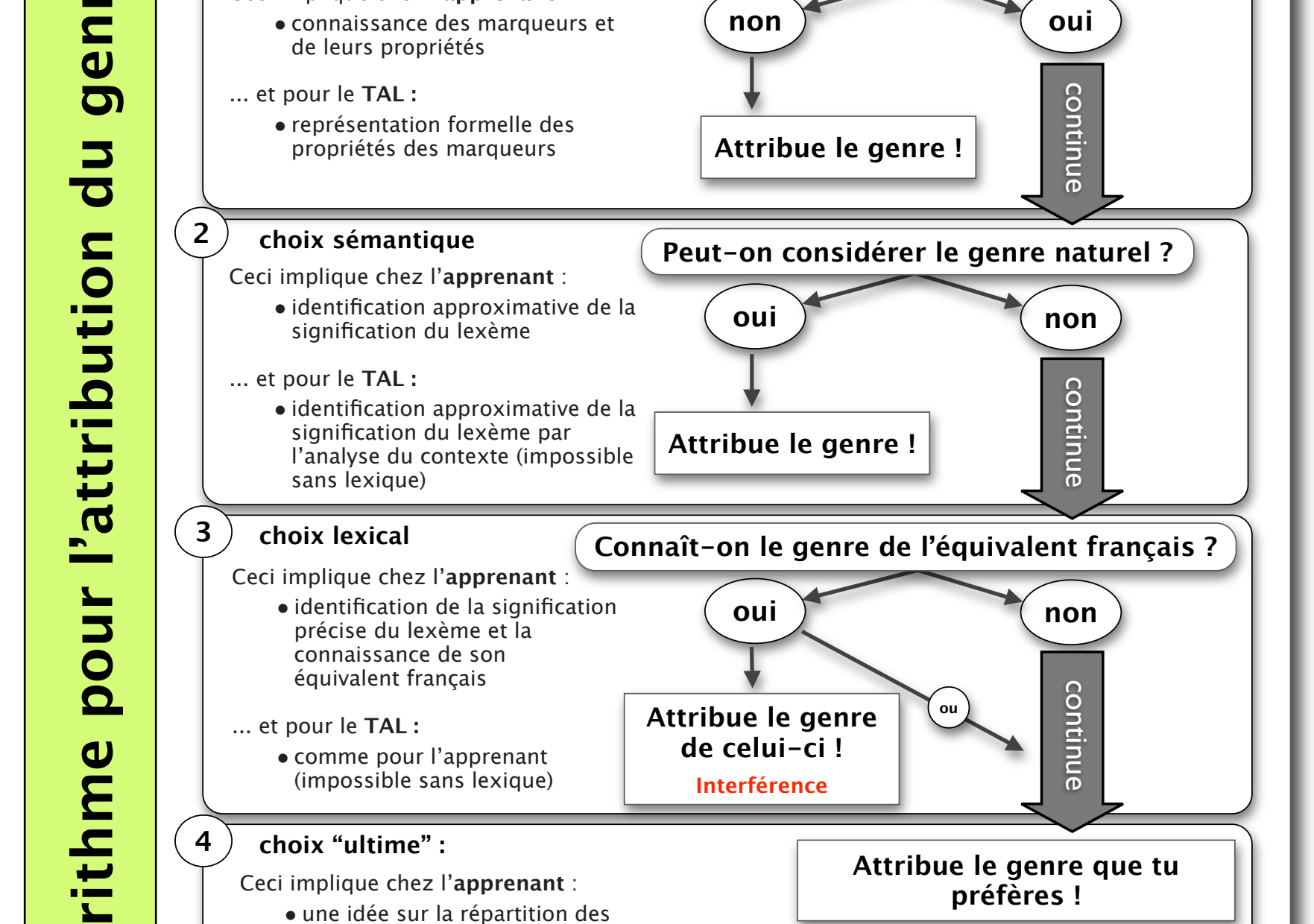


Erreurs de genre dans les productions recueillies

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. À l'aide de cette affiche, possible grâce à la base de données et au diagnostic des erreurs, il est aisé d'observer les différentes occurrences des erreurs de genre et leur relation avec la nature du marqueur formel du lemme de la production requise. Avec plus de données recueillies dans le futur, il sera possible de mieux dégager des tendances générales.



Algorithme pour l'attribution du genre



Bibliographie

- 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(4):453-467.
- Čermák, F. (1997). *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.
- Hlavčová, J. & Savičková, E. (2002). *Measures of word commonness*. *Journal of Quantitative Linguistics*, 9(3): 215-213.
- Klanten, R. éditeur (2009). *Data Flow. Design graphique et visualisation d'information*. Thames & Hudson, Paris.
- Sgall, P., Hajičová, E. & Panevová, J. (1986). *The Meaning of the Sentence in its Semantic and Pragmatic Aspects*. D. Reidel Publishing Company, Praha - Amsterdam.
- Šmilauer, I. (2008). *Acquisition du tchèque par les francophones : analyse automatique des erreurs de déclinaison*. *The Prague Bulletin of Mathematical Linguistics* 90:33-56.