

# POS Comes with Parsing: a Refined Word Categorisation Method

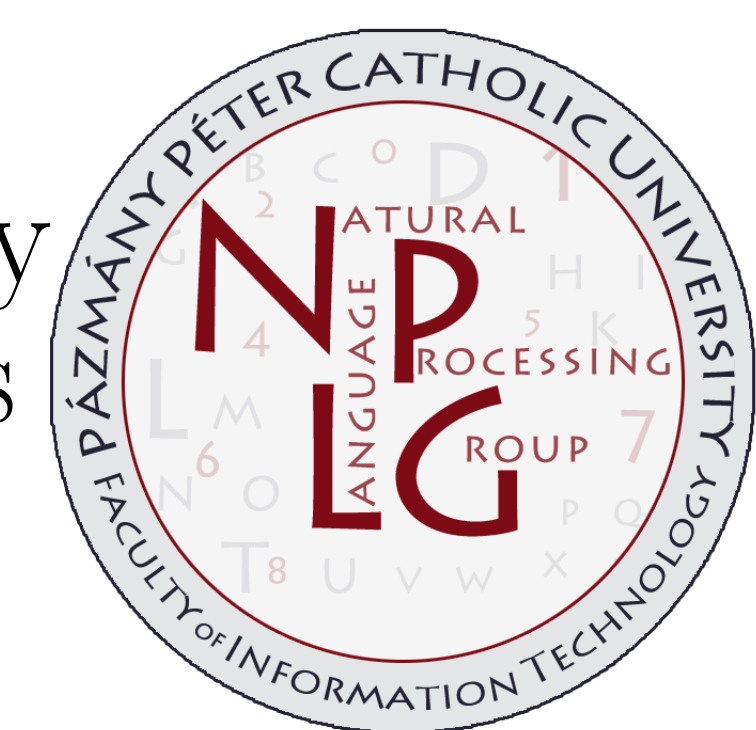
Balázs Indig<sup>1,2</sup> & Noémi Vadász<sup>1,3</sup>

<sup>1</sup>MTA–PPKE Hungarian Language Technology Research Group, Budapest, Hungary

<sup>2</sup>Pázmány Péter Catholic University, Faculty of Information Technology and Bionics

<sup>3</sup>Pázmány Péter Catholic University, Faculty of Humanities and Social Sciences

{indig.balazs,vadasz.noemi}@itk.ppke.hu

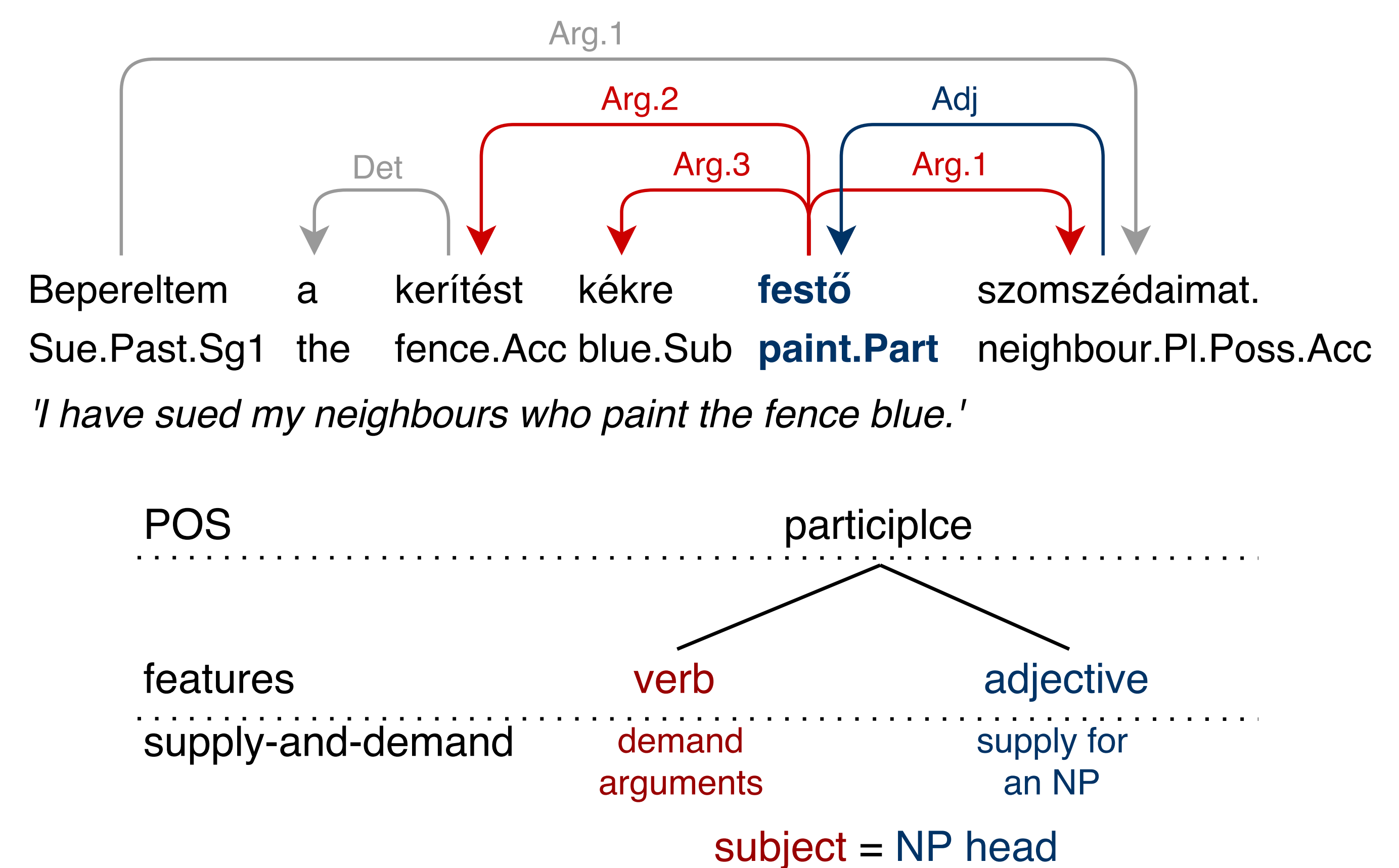


## Introduction

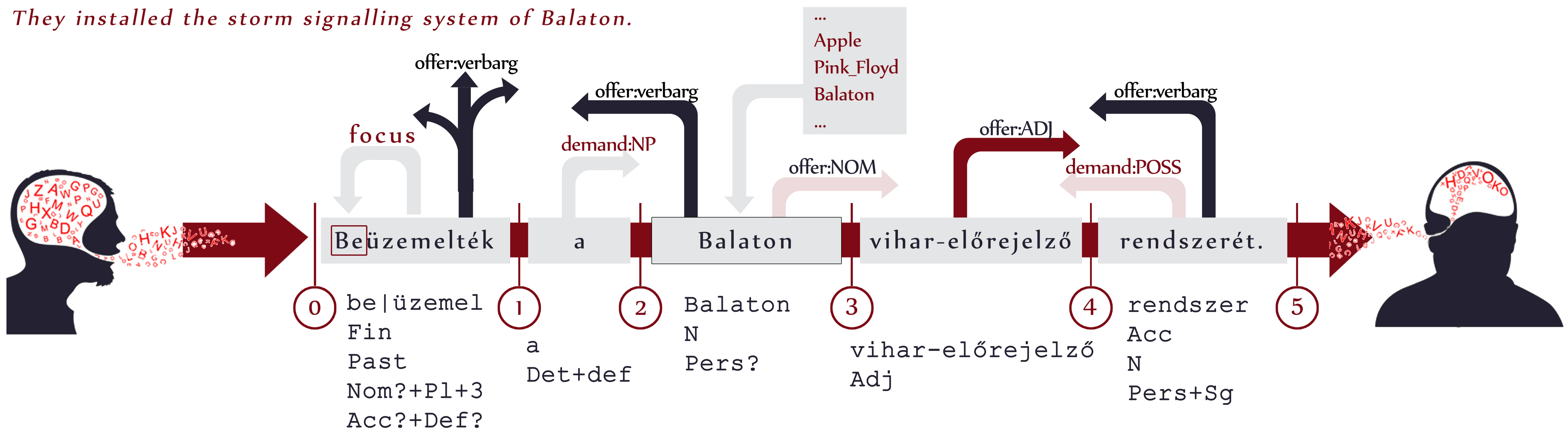
ANAGRAMMA [1, 3] is a **psycholinguistically motivated** parser:

- **left-to-right, word-by-word** approach: as humans do
- **performance-based**: only existing constructions are preferred
- **supply-and-demand framework** [2] based on features:
  - based on lexical representation and morpho-syntactic information
  - like Word-Expert-Parsing [4]: features instead of POS tags
  - *main and subcategories*
    - \* based on syntactic and semantic features (statistically supported)
    - \* categories can start different processes during the parsing
    - \* the main category inherits its processes to the subcategories
    - \* we can mix features of classical categories within one word

## Participle = Verb & Adjective



*They installed the storm signalling system of Balaton.*



## Detached Preverbs or Adverbs

We gathered statistics from two different corpora to model Hungarian.

- Hungarian Gigaword Corpus (HGC), 709 million tokens: both edited and unedited texts from different domains
- InfoRadio Corpus (IRC), 1.953 419 tokens: edited texts of uniform domains, short political news

### Positions that can be considered real preverbs:

Frequent preverbs	-2	0	+1	+2	+3
<i>meg, ki, be, le, fel, föl, el, át, rá</i>	0.5%	58.5%	40%	1%	0.01%

**98.5% of the preverbs appear on the verb (0) or immediately after (+1) it. The others (+4–) are not real preverbs.**

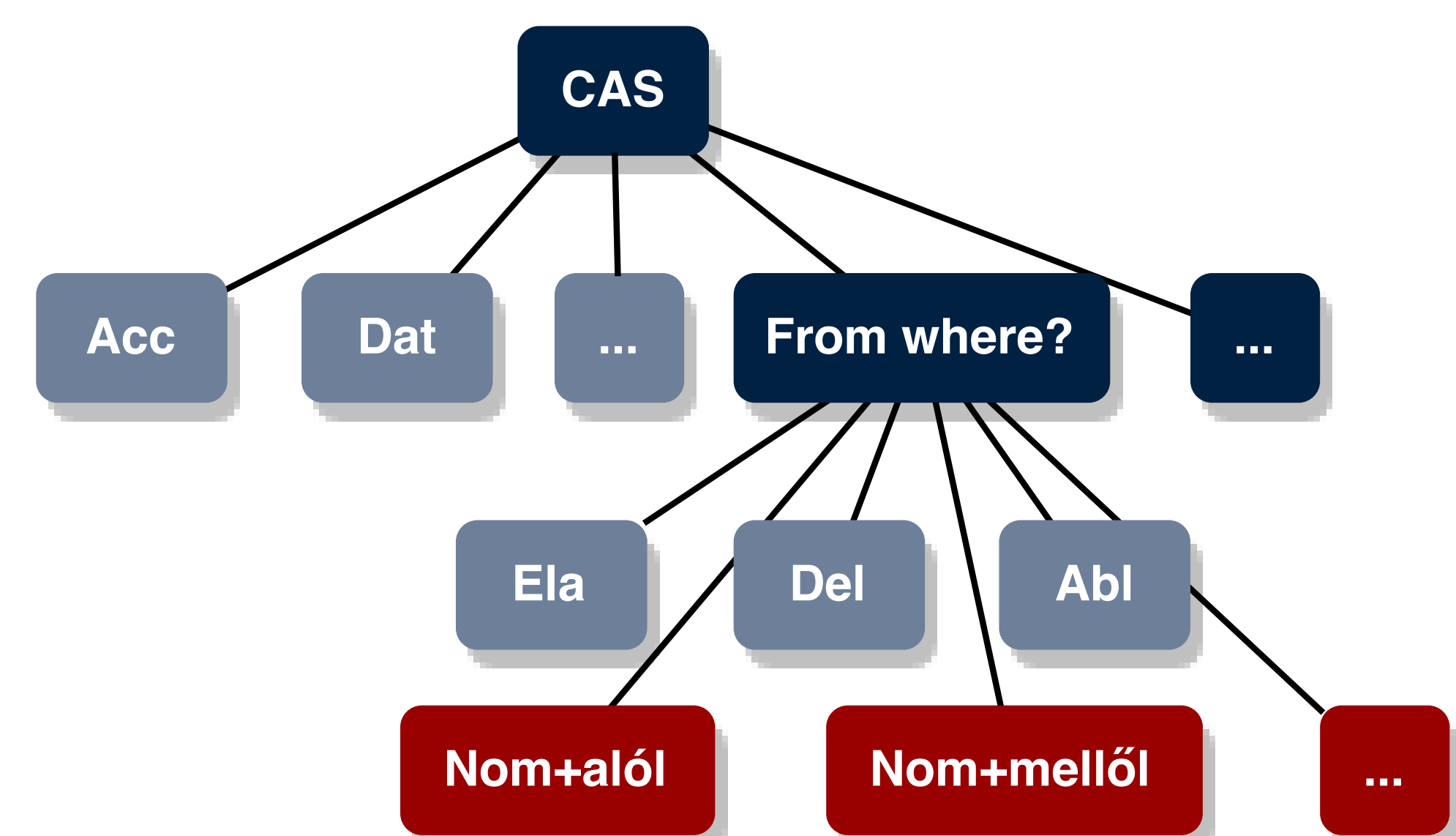
V.FIN	+1	+2	+3	+4	+5	+6	+7
HGC	7.527.308	163.993	5.126	1.193	267	101	27
IRC	23.552	220	-	-	-	-	-
HGC%	97.78%	2.13%	0.0666%	0.015%	0.003%	0.001%	3.5e-4%
IRC%	99.999%	0.001%	-	-	-	-	-

**Csábítson** téged a rettyezáti nagy barna medve **oda** ahova akarod.  
**Lure** you the from Rettyezát big brown bear **there** where you want to  
 'Let the big brown bear from Rettyezát **lure** you **wherever** you want to.'

## Postpositions as special case markers

*Honnan jöttél?* 'Where do you come from?'

case or PP	abbrev.	Hungarian example	translation
elativus	Ela	<i>A börtön<b>ből</b>.</i>	'From the prison.'
delativus	Del	<i>A tárgyalás<b>ról</b>.</i>	'From the trial.'
ablativus	Abl	<i>Az ügyvéd<b>től</b>.</i>	'From the lawyer.'
postposition	Nom + <b>alól</b>	<i>A hegy <b>alól</b>.</i>	'From <b>under</b> the hill.'
postposition	Nom + <b>mellől</b>	<i>A Tisza <b>mellől</b>.</i>	'From <b>near</b> the lake Tisza.'
postposition	Nom + <b>mögül</b>	<i>A rács <b>mögül</b>.</i>	'From <b>behind</b> the bars.'



## The hierarchical system of Hungarian cases and postpositions

## Conclusion

We illustrate how the dissected features that build the classical POS categories start supplies and demands, still our method conforms to the existing theories.

[1] Balázs Indig and Gábor Prószék. Magyar szövegek pszicholingvisztikai indíttatású elemzése számítógéppel. *Alkalmazott nyelvtudomány*, 15(1-2):29–44, 2015.

[2] Balázs Indig and Noémi Vadász. Windows in human parsing – how far can a preverb go? In Tadić Marko and Božo Bekavac, editors, *Proceedings of the Tenth International Conference on Natural Language Processing (HrTAL2016)*, Dubrovnik, Croatia, sept 2016. (Accepted, In press).

[3] Gábor Prószék, Balázs Indig, and Noémi Vadász. Performanciaalapú elemző magyar szövegek számítógépes megértéséhez. In Kas Bence, editor, *"Szavad ne feledd!": Tanulmányok Bánréti Zoltán tiszteletére*, pages 223–232. MTA Nyelvtudományi Intézet, Budapest, 2016.

[4] Steven L. Small. Word expert parsing. In *Proceedings of the 17th Annual Meeting on Association for Computational Linguistics*, ACL '79, pages 9–13, Stroudsburg, PA, USA, 1979. Association for Computational Linguistics.