# VERB ARGUMENT BROWSER

### ARGUMENT FRAMES IN THE HUNGARIAN NATIONAL CORPUS

Bálint Sass
sass.balint@nytud.mta.hu

Research Institute for Linguistics, Hungarian Academy of Sciences

Corpus resources for quantitative and psycholinguistic analysis
Eger, 2-3 June 2014

## PREVIEW

- Verb Argument Browser – a specific corpus query tool for investigating argument structure of verbs.
- original version:
    - for Hungarian
      based on the old version of the Hungarian National Corpus
    - not just arguments
      *all* NP and PP dependents of verbs
      subjects, objects, complements, adjuncts included
    - investigating
      verb subcategorization frames, institutionalized phrases,
      light verb constructions, idiomatic verbal expressions,
      figures of speech . . .
      common property: verb + some NP/PP dependents
    - *examples*
- language independence

1. SENTENCE MODEL

2. VERB PHRASE CONSTRUCTIONS AS COLLOCATIONS

3. USAGE & EXAMPLES

4. APPLICATIONS

5. LANGUAGE INDEPENDENCE

## SENTENCE MODEL

- *Basic unit:* simple sentence or clause.

> A      lány      váll-at              von.
> the    girl      shoulder-ACC    pull.
> 'The girl shrugs her shoulder.'

- clause = verb + set of NP/PP dependents → verb frame

> *verb*=von      *NOM*=lány      *ACC*=váll
> *verb*=shrug    *SUBJ*=girl     *OBJ*=shoulder

- *Dependent types*: defined . . .

    – syntactically:      word order      (in English)
    – morphologically:    case markers    (in Hungarian)

## SENTENCE MODEL

- *Basic unit:* simple sentence or clause.

| A | lány | váll-at | von. |
|---|------|---------|------|
| the | girl | shoulder-ACC | pull. |

'The girl shrugs her shoulder.'

- clause = verb + set of NP/PP dependents → verb frame

| *verb*=von | *NOM*=lány | *ACC*=váll |
|------------|------------|------------|
| *verb*=shrug | *SUBJ*=girl | *OBJ*=shoulder |

- *Dependent types*: defined ...

    – syntactically:     word order    (in English)
    – morphologically:  case markers  (in Hungarian)

## DEPENDENT TYPES

- in Hungarian: 20 different case markers
  in English: usually prepositions

| case marker | case | abbr. | English |
|---|---|---|---|
| -∅ | nominative | NOM | word order |
| -t | accusative | ACC | word order |
| -bAn | inessive | INE | *in*-phrase |
| -rÓl | delative | DEL | *from*-phrase[1] |
| -bÓl | elative | ELA | *from*-phrase[2] |
| . . . | | | |

## EXAMPLES

Az    emberek    az    időjárás-ról    beszélnek.
the    people    the    weather-DEL    talk.
'People talk about the weather.'

*verb*=beszél    *NOM*=ember    *DEL*=időjárás
*verb*=talk    *SUBJ*=people    *ABOUT*=weather

Péter    fél    az    ismeretlen-től.
Peter    fear    the    unknown-ABL.
'Peter fears of the unknown.'

*verb*=fél    *NOM*=Péter    *ABL*=ismeretlen
*verb*=fear    *SUBJ*=Peter    *OF*=unknown

## EXAMPLES

| Az | emberek | az | időjárás-ról | beszélnek. |
|----|---------|----|--------------|------------|
| the | people | the | weather-DEL | talk. |

'People talk about the weather.'

| *verb*=beszél | *NOM*=ember | *DEL*=időjárás |
|---------------|-------------|-----------------|
| *verb*=talk | *SUBJ*=people | *ABOUT*=weather |

| Péter | fél | az | ismeretlen-től. |
|-------|-----|----|-----------------|
| Peter | fear | the | unknown-ABL. |

'Peter fears of the unknown.'

| *verb*=fél | *NOM*=Péter | *ABL*=ismeretlen |
|------------|-------------|------------------|
| *verb*=fear | *SUBJ*=Peter | *OF*=unknown |

## EXAMPLES

| Az | emberek | az | időjárás-ról | beszélnek. |
|----|---------|-----|--------------|-----------|
| the | people | the | weather-DEL | talk. |

'People talk about the weather.'

| *verb*=beszél | *NOM*=ember | *DEL*=időjárás |
|---------------|-------------|----------------|
| *verb*=talk | *SUBJ*=people | *ABOUT*=weather |

| Péter | fél | az | ismeretlen-től. |
|-------|-----|-----|-----------------|
| Peter | fear | the | unknown-ABL. |

'Peter fears of the unknown.'

| *verb*=fél | *NOM*=Péter | *ABL*=ismeretlen |
|------------|-------------|------------------|
| *verb*=fear | *SUBJ*=Peter | *OF*=unknown |

## FIXED AND FREE DEPENDENTS

| Hogy | jöttek | lét-re | | az | első | csillagok? |
|------|--------|--------|---|----|------|-----------|
| how | came | existence-SUB | | the | first | stars? |

'How the first stars came into existence?'

| *verb*=jön | *SUB*=lét | *NOM*=csillag |
|-----------|-----------|---------------|
| *verb*=come | *INTO*=existence | *SUBJ*=star |

- *fixed dependent:*
  cannot change the content word
  without changing the meaning of the VPC

- *free dependent:*
  can change the content word
  without changing the meaning of the VPC

## FIXED AND FREE DEPENDENTS

Hogy   jöttek   lét-re          az     első   csillagok?
how     came     existence-SUB   the    first  stars?
'How the first stars came into existence?'

*verb*=jön        *SUB*=lét            *NOM*=csillag
*verb*=come       *INTO*=existence     *SUBJ*=star

- *fixed dependent:*
  cannot change the content word
  without changing the meaning of the VPC

- *free dependent:*
  can change the content word
  without changing the meaning of the VPC

## MULTI-WORD VERBS

- *multi-word verb:* verb stem + some fixed dependent(s)

> lét-re           jön
> existence-SUB   come
> 'come into existence'

Multi-word verbs have. . .

- separate meaning
- own argument structure

> rész-t      vesz   bAn
> part-ACC    take   INE
> 'take part in SOMETHING'

Typical units to be investigated using the VAB.

SENTENCE MODEL

> sentence = verb + set of dependents
> dependent = type + content word

i.e.

| *verb*=jön | *SUB*=lét | *NOM*=csillag |
| *verb*=come | *INTO*=existence | *SUBJ*=star |

In this way we can investigate VPCs independently from the particular word order in which they appear in the corpus.

## CORPUS PREPARATION

*Input:* Hungarian National Corpus
(POS-tagged and disambiguated)

- clause boundary detection
  – regexps based on conjunction and punctuation patterns
- verb normalization
  – e.g. separated verbal prefixes attached
- noun phrase chunking
  $\rightarrow$ case and lemma of the head of dependent phrases

$\rightarrow$ representation according to the model

## VERB PHRASE CONSTRUCTIONS AS COLLOCATIONS

A specific kind of VPCs:

'take part in SOMETHING'

- fixed dependent (object) + free dependent (*in*-phrase)
- multi-word verb with argument structure

These kind of expressions are

1. subcategorization frames
   *and*
2. collocations
   at the same time.

The idea behind the VAB is: **treat VPCs as collocations.**

# VERB PHRASE CONSTRUCTIONS AS COLLOCATIONS

We search for collocations in the space of these structures:

| | | |
|---|---|---|
| *verb*=jön | *SUB*=lét | *NOM*=csillag |
| *verb*=come | *INTO*=existence | *SUBJ*=star |

### IDEA

Apply an association measure (designed for bigrams) taking ...

- the content word of a particular dependent – as one unit,
- *all other* parts of the verb frame – as the other unit

of the collocation.

## VERB PHRASE CONSTRUCTIONS AS COLLOCATIONS

We search for collocations in the space of these structures:

| | | |
|---|---|---|
| *verb*=jön | *SUB*=lét | *NOM*=? |
| *verb*=come | *INTO*=existence | *SUBJ*=? |

### IDEA

Apply an association measure (designed for bigrams) taking ...

- the content word of a particular dependent – as one unit,
- *all other* parts of the verb frame – as the other unit

of the collocation.

If we choose the subject's content word as the first unit,

## VERB PHRASE CONSTRUCTIONS AS COLLOCATIONS

We search for collocations in the space of these structures:

| | | |
|---|---|---|
| *verb*=jön | *SUB*=lét | *NOM*=? |
| *verb*=come | *INTO*=existence | *SUBJ*=? |

### IDEA

Apply an association measure (designed for bigrams) taking . . .

- the content word of a particular dependent – as one unit,
- *all other* parts of the verb frame – as the other unit

of the collocation.

If we choose the subject's content word as the first unit,
we will query the most important or usual subjects of this
construction, namely *what* is used to come into existence.

## VERB PHRASE CONSTRUCTIONS AS COLLOCATIONS

The Verb Argument Browser can answer the following typical research question:

- What are the salient words which can appear as a particular dependent of a given verb frame?

- What are the most important collocates of a given verb (or verb frame) as a particular dependent?

Association measure: *salience* (adjusted mutual information)

$$S(x, y) = log_2 f(x) \cdot log_2 N \frac{f(x, y)}{f(x) \cdot f(y)}$$

## VERB PHRASE CONSTRUCTIONS AS COLLOCATIONS

Consequence:

The Verb Argument Browser can treat not just a single word but a whole verb frame (a verb together with some arguments) as one unit in collocation extraction.

It can collect . . .

- salient subjects of a verb,
- salient objects of a given verb–subject pair,
- salient locatives of a given verb–subject–object triplet . . .

# USAGE

- *corpus:* Hungarian National Corpus (187 million words)
- *response times:* a few seconds

# USAGE

- *corpus:* Hungarian National Corpus (187 million words)
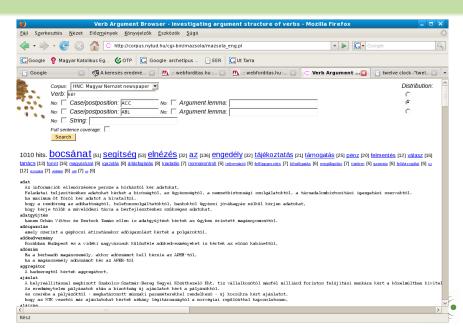- *response times:* a few seconds

# USAGE

- *corpus:* Hungarian National Corpus (187 million words)
- *response times:* a few seconds

# USAGE

- *corpus:* Hungarian National Corpus (187 million words)
- *response times:* a few seconds

# USAGE

- *corpus:* Hungarian National Corpus (187 million words)
- *response times:* a few seconds

Verb Argument Browser - investigating argument structure of verbs - Mozilla Firefox

Fájl   Szerkesztés   Nézet   Előzmények   Könyvjelzők   Eszközök   Súgó

http://corpus.nytud.hu/cgi-bin/mazsola/mazsola_eng.pl

Google   Magyar Katolikus Eg...   OTP   Google: archetípus ...   EER   Út Tárra

Google   A keresés eredmé...   .:: webforditas.hu :...   .:: webforditas.hu :...   Verb Argument ...   twelve clock -"twel...

Corpus: HNC: Magyar Nemzet newspaper

Verb: kér

No: ☐   Case/postposition: ACC          No: ☐   Argument lemma:
No: ☐   Case/postposition: ABL          No: ☐   Argument lemma:
No: ☐   String:
Full sentence coverage: ☐
Search

Distribution:

1010 hits. **bocsánat** [51] **segítség** [53] **elnézés** [32] **az** [136] **engedély** [32] **tájékoztatás** [21] **támogatás** [25] **pénz** [20] **felmentés** [12] **válasz** [16] **tanács** [13] forint [16] magyarázat [9] igazolás [8] állásfoglalás [8] kiadatás [7] normakontroll [6] információ [9] felfüggesztés [7] kihallgatás [6] megállapítás [7] türelem [6] garancia [6] felülvizsgálat [5] ez [12] vizsgálat [7] vélemén [6] ami [7] az [6]

adat
    Az információk ellenőrzésére persze a kórházától kér adatokat,
    Feladatai teljesítéséhez adatokat kérhet a bíróságtól, az ügyészségtől, a nemzetbiztonsági szolgálatoktól, a társadalombiztosítási igazgatási szervektől,
    ha maximum őt főről kér adatot a hivataltól.
    hogy a rendőrség az adóhatóságtól, telefonszolgáltatóktól, bankoktól ügyészi jóváhagyás nélkül kérjen adatokat,
    hogy kérje tőlük a művelődési tárca a bérfejlesztéshez szükséges adatokat.
adatgyűjtés
    hanem Orbán Viktor és Deutsch Tamás ellen is adatgyűjtést kértek az ügyben érintett magánnyomozótól.
adóigazolás
    amely szerint a gépkocsi átíratásakor adóigazolást kértek a polgároktól.
adókedvezmény
    Korábban Budapest és a vidéki nagyvárosok különféle adókedvezményeket is kértek az előző kabinettól.
adószám
    Ha a bérbeadó magánszemély, akkor adószámot kell kérnie az APEH-tól,
    ha a magánszemély adószámot kér az APEH-tól
aggregátor
    A hadseregtől kértek aggregátort,
ajánlat
    A helyreállítással megbízott Szabolcs-Szatmár-Bereg Megyei Közútkezelő Kht. tíz vállalkozótól másfél milliárd forintos felújítási munkára kért a közelmúltban kivitel
    Az eredménytelen pályázatok után a bizottság új ajánlatot kért a pályázóktól.
    én cserébe a pályázóktól - meghatározott műszaki paraméterekkel rendelkező - új locsírka kért ajánlatot.
    hogy az MTK vezetői már ajánlatokat kértek néhány légitársaságtól a norvégiai repülőúttal kapcsolatosan,
aláírás

Kész

## QUERY: 2 FREE DEPENDENTS

kér    -t       -tól
ask    ACC   ABL
'ask for SOMETHING from SOMEBODY'

*verb*=kér    *ABL*=?    *ACC*=?
*verb*=ask    *from*=?    *for*=?

## QUERY: 2 FREE DEPENDENTS

kér    -t      -tól
ask   ACC   ABL
'ask for SOMETHING from SOMEBODY'

*verb*=kér    *ABL*=?    *ACC*=?
*verb*=ask    *from*=?   *for*=?

*Result:* (Most salient objects:)

- bocsánat – 'forgiveness'
- segítség – 'help'
- elnézés – also 'forgiveness'
- engedély – 'permission'
- . . .

# QUERY: 1 FIXED + 1 FREE DEPENDENT

| | | |
|---|---|---|
| vesz | figyelem-bA | -t |
| take | consideration-ILL | ACC |

'take SOMETHING into consideration'

| | | |
|---|---|---|
| *verb*=vesz | *ILL*=figyelem | *ACC=?* |
| *verb*=take | *INTO*=consideration | *OBJ=?* |

## QUERY: 1 FIXED + 1 FREE DEPENDENT

vesz    figyelem-bA         -t
take    consideration-ILL    ACC
'take SOMETHING into consideration'

*verb*=vesz    *ILL*=figyelem          *ACC*=?
*verb*=take    *INTO*=consideration    *OBJ*=?

*Result:* (Most salient direct objects:)

- szempont – 'point of view'
- érdek – 'interest'
- vélemény – 'opinion'
- . . .

# A TRICK: QUERYING THE VERB

## A SIMPLE QUERY

ad      -t
give    ACC
'give SOMETHING'

*verb*=ad      *ACC*=?
*verb*=give    *OBJ*=?

## A SIMPLE QUERY

> ad      -t
> give   ACC
> 'give SOMETHING'

> *verb*=ad       *ACC*=?
> *verb*=give     *OBJ*=?

*Result:* (Most salient direct objects:)

- hang – 'voice' → 'to give voice to SOMETHING'
- hír – 'news' → to give news ∼ 'to report'
- igaz – 'true' → to give true ∼ 'to take sy's side'
- . . .

→ **multi-word verbs**

## ANOTHER SIMPLE QUERY

üt       ∅
beat   NOM
'SOMETHING beats'

*verb*=üt       *NOM*=?
*verb*=beat   *SUBJ*=?

## ANOTHER SIMPLE QUERY

üt      ∅
beat    NOM
'SOMETHING beats'

*verb*=üt      *NOM*=?
*verb*=beat    *SUBJ*=?

*Result:* (Some salient subjects:)

- óra – 'clock' → 'The clock strikes twelve.'
- forint → 10 Ft beat his palm. ∼ 'He receives 10 Ft.'
- kő – 'stone' → Üsse kő!
  – Let a stone beat it! ∼ 'It does not matter.'
- . . .

→ **multi-word verbs, figures of speech**

## COLLECTING MWVS

Important property of the Verb Argument Browser:

For any specific dependent, the tool provides constructions
where *this dependent is fixed*, if there is any such construction.
(e.g. light verb constructions, idiomatic verbal expressions, figures of speech)

+ 'take' + into → 'consideration', 'account' . . .
– 'eat' + OBJ → just some kinds of food
  in this case we obtain frequent words with literal meaning,
  often forming a semantically coherent class

## COLLECTING MWVS

Important property of the Verb Argument Browser:

> For any specific dependent, the tool provides constructions
> where *this dependent is fixed*, if there is any such construction.
> (e.g. light verb constructions, idiomatic verbal expressions, figures of speech)

- + 'take' + into → 'consideration', 'account' ...
- – 'eat' + OBJ → just some kinds of food
  in this case we obtain frequent words with literal meaning,
  often forming a semantically coherent class

> *Warning:* VPCs with fixed position(s) are frequent.

They are not to be ignored, it is necessary to deal with them.

E.g. when doing something with verbs, do not forget about multi-word verbs.

## CORPUS SIZE ↔ ANNOTATION RICHNESS

- compared to. . .
  – large raw or POS-tagged corpora ($\sim$ big data)
  – small syntactically annotated corpora ($\sim$ rich information)

- a VAB uses (and works well with) corpora which are
  big enough + have "some" syntactic information annotated

- using this approach
  corpus-driven information can be gathered
  about some *"higher level"* phenomena
  (i.e. the predicate-argument structure in our case)
  based on querying a quite *large* piece of text

$\rightarrow$ a big corpus (size $\gg$ treebanks)
  with *some shallow* syntactic annotation
  can be a valuable resource. :)

1 SENTENCE MODEL

2 VERB PHRASE CONSTRUCTIONS AS COLLOCATIONS

3 USAGE & EXAMPLES

4 APPLICATIONS

5 LANGUAGE INDEPENDENCE

## APPLICATIONS

Past...

- lexical database development of a Hungarian to English machine translation system – handling MWVs
- searching for MWVs to include them into the Hungarian WordNet
- lexicography

Future...

- language teaching
- determining the frequency of particular VPCs to be used in experiments in psycholinguistic research
- linguistic research
  - studying verb synonyms
  - classifying verbs based on argument structure similarities
  - studying selectional preferences of verbs

## LANGUAGE INDEPENDENCE

*Claim:* the approach is language independent.

- To extend the methodology to other languages all we need is a corpus represented according to our sentence model.

- Can we create such a representation?

- Essentially, the representation relies on the very fact that there are *some kind of predicate-argument structure* in the languages.

- All we should do is segment the text into sentences/clauses (containing a verb and its dependents) and specify the relationship between the verb and the particular dependents.

## A PILOT VAB FOR SERBIAN

Just to show that the approach works. :)

*Serbian:* dependents are defined by
case markers *or* case marker + preposition combinations.

- much smaller corpus (Intera)
- much simpler (pre)processing
    1. clause boundary detection = just split at punctuations
    2. verb identification = take the last verb + attach *se* if occurs
    3. noun phrase chunking = extract PPs according to this
       simple pattern: *a preposition + possible not-nouns + a noun*
    4. no case information:
       all NPs without preposition
       go to a big class (called ANYCASE)

## REPRESENTATION OF A SERBIAN CLAUSE

Example clause from the corpus:

| Svako | ima | pravo | na | rad. |
|---|---|---|---|---|
| Everyone | has | right | to/for | work. |

'Everyone has the right to work.'

Representation of the example clause:

| *verb*=imati | *ANYCASE*=pravo | *na*=rad |
|---|---|---|
| *verb*=have | *OBJ*=right | *to/for*=work |

# EXAMPLE: A MULTIWORD VERB

Query:
'imati pravo na' ('have right to/for')



184 hits. **sloboda** [27] **zaštita** [23] **naknada** [13] jezik [7] odsustvo [5] lek [4] poštovanje [4] podrška [4] život [4]

freedom, protection, compensation, language …

# EXAMPLE: DISCOVERING MULTIWORD VERBS

Query:
'ići u' ('go in')

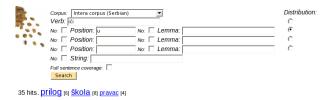

35 hits. prilog [6] škola [8] pravac [4]

benefit, school, direction

# EXAMPLE: DISCOVERING MULTIWORD VERBS

Query:
'ići u' ('go in')



35 hits. **prilog** [6] <u>škola</u> [8] <u>pravac</u> [4]

benefit, school, direction

> *prilog* does not fit into this little semantic class.
> → This phenomenon is a good indicator of being a MWV!
>
> 'ići u prilog' is a MWV.
> Meaning: ∼ *support* (?) – literally: 'go in benefit' (?)

## LANGUAGE INDEPENDENCE

The methodology can be extended to other languages,
and a fully functioning VAB can be created
if a shallow parsed, adequately processed corpus is available.

AVAILABILITY

Available for you:

- Hungarian version:

```
http://corpus.nytud.hu/vab
```

- Serbian version:

```
http://corpus.nytud.hu/vabs
```

username: eger; password: vab

AVAILABILITY

Available for you:

- Hungarian version:

  http://corpus.nytud.hu/vab

- Serbian version:

  http://corpus.nytud.hu/vabs

  username: eger; password: vab

  Thank you for your attention!