



CENTRE D'ÉTUDES FRANCO-RUSSE

ЦЕНТР ФРАНКО-РОССИЙСКИХ  
ИССЛЕДОВАНИЙ



# New Russian resources for the linguistic software Nooj

Improving linguistics resources  
and adding semantic tags  
for the Russian language  
*for Max Silberztein's Nooj software*

**Zagreb**  
**Nooj Conference**  
**June 2020**

Vincent BÉNET  
INALCO  
**CEFR-Moscow**

# Brief summary about Russian resources for Nooj

Zaliznyak dictionary

inverted dictionary

most useful dictionary for  
Russian grammar

Zagreb Conference Nooj  
June 2020

чудотворец	мо	5*a	
крюкотворец	мо	5*a	
языкотворец	мо	5*a	
песнотворец	мо	5*a	
миротворец	мо	5*a	
стихотворец	мо	5*a	
горец	мо	5*a	
черногорец	мо	5*a	
эквадорец	мо	5*a	
корец	м	5*b	
черноморец	мо	5*a	
поморец	мо	5*a	
североморец	мо	5*a	
шпорец	м	5*b	[// шпорца]
торец	м	5*b	
консерваторец	мо	5*a	
алюторец	мо	5*a	
шорец	мо	5*a	
хитрец	мо	5b	
острец	м	5b	
кострец	м	5b	
нутрец	м	5b	(болезнь лошадей);
	мо	5b	(лошадь, страдающая этой болезнью)
огурец	м	5*b	
сырец	м	5*b	
эльзасец	мо	5*a	
парнавец	мо	5*a	
запавец	м	5*a	
овсец	м	5b	
песец	мо	5*b	(в т. ч. о мехе)
тунивец	мо	5*a	
писец	мо	5*b	
живописец	мо	5*a	
борзописец	мо	5*a	
самописец	м	5*a	
стенписец	мо	5*a	
иконописец	мо	5*a	
баснописец	мо	5*a	
скорописец	мо	5*a	
летописец	мо	5*a	(составитель летописи);
	м	5*a	(летопись)
высотописец	м	5*a	



## Brief summary about Russian resources for Nooj

- Linguistic resources for Russian language  
**Zaliznyak's dictionary**

**-Very complete grammar information with some unformal unclassified added annotations**

**-All forms are included in the dictionary,  
BUT NO information about derivation (prefixes or suffixes)**

# Brief summary about Russian resources for Nooj

• Structure of Russian language :

**PREF (s) + ROOT + SUFF (s) + DES + SUF2**

	pi		t'	
na	pi	t	ok	
na	pi		t'	sja
na	pi	va	t'	sja

# How can we improve Nooj Russian resources?

PREFIX + ROOT + SUFFIX + DESINENCE

Verbal prefixation and derivation

Выпить *pf* – выпивать *ipf*  
vy-pi-t' → vy-pi-va-t' *to finish drinking*

ПИТЬ

pi-t'

*To drink*

напиться *pf* – напиваться *ipf*  
na-pi-t'-sja → na-pi-va-t'-sja *to get /be drunk*

пропить *pf* – пропивать *ipf*  
pro-pi-t' → pro-pi-va-t' *to spend sthg in drinking*

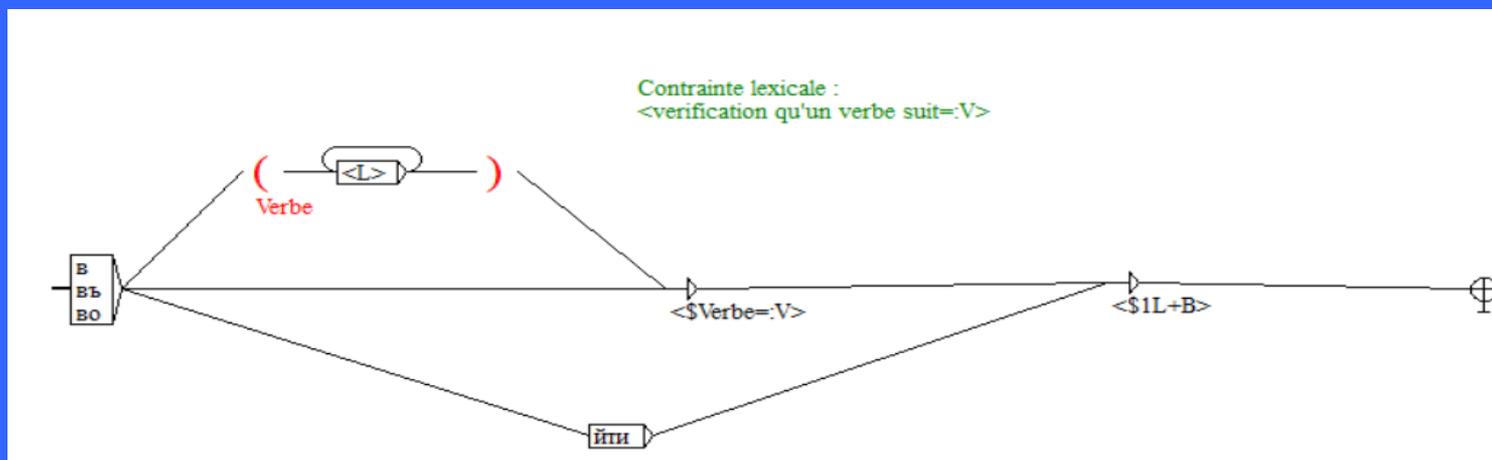
# How can we improve Nooj Russian resources?

## Implementing verbal prefixes and suffixes

More information in the dictionary?      More derivational grammars .nom

ПИТЬ      ПИТЬ... + DRV= «ва» + DRV=«ся» + DRV=  
pi-t'      «вы» + DRV = «на» + DRV=«за» + DRV=«про»

Grammar to recognize the verbal prefix В (meaning = into)



# Improving lexical resources

## Adding semantic resources to the main dictionary

Simple classification « A word and an idea »  
House, sports, body, ...

Semantic tags / annotation  
of the Russian national Corpus

Semantic dictionary of Tuzov  
(160 000 words)

# Improving lexical resources

## Adding semantic resources to the main dictionary

Nooj properties.def

A\_Sem = Animal | Color | ( App

N\_Sem = Hum | Prof | Parents | Body

Conc | Abstr | Org | Text |

Animal | Food | Health | Arts | Lit | Music | Sports

Topo | Country | River | City | Mount| Lake |

Posit | Time | Color ;

ADV\_Sem = Time |Topo | Modal;

V\_Sem = Color | Topo | Posit |Modal;

# Improving lexical resources

## Semantic resources added to the main dictionary

**Prof = 900**

**Parent/ Kinship = 160 items**

**Forenames = 2280**

**Animal = 370**

**Food = 280 (Liquid = 25 )**

**Body = 285**

**Health = 175**

**Arts = 65**

**Lit = 40**

**Music = 155**

**Sport = 65**

**Topo = 40**

**Country = 180**

**River = 15**

**City = 175**

**Mount = 5**

**Lake = 5**

**Posit = 25**

**Time = 135**

**Modal = 15**

**Color = 275**

# Improving lexical resources

Added semantic resources to the main dictionary

## Morphology or semantics ?

ходить, V+Mvt+Indet+ipf+intrLX=ходить

идти, V+Mvt+Det+ipf+intrLX=идти

входить, V+Mvt+Pvb+ipf+intrLX=ходить

войти, V+Mvt+Pvb+pf+intrLX=идти

выходить, V+Mvt+Pvb+ipf+intrLX=ходить

приезжать, V+Mvt+Pvb+ipf+intrLX=акать

# Working with Nooj Searching for « verbs of motion »

The screenshot shows the Nooj software interface. A dialog box titled "Locate a pattern in Tchexov\_Dama" is open, showing the search pattern "<V+Mvt+Det>". The search results are displayed in a table with columns "Before", "Seq.", and "After".

Before	Seq.	After
блондинка, в берете; за нею	бежал	белый шпиц. И потом он
неделю. Помолчали немного. - Время	идет	быстро, а между тем здесь
все равно, куда бы ни	идти	, о чем ни говорить. Они
и по ней от луны	шла	золотая полоса. Говорили о том
лошадях, и он провожал ее.	Ехали	целый день. Когда она садилась
огонь. Уже начались морозы. Когда	идет	первый снег, в первый день
Сергеевна не снилась ему, а	шла	за ним всюду, как тень
говорила: - Тебе, Дмитрий, совсем не	идет	роль фата. Однажды ночью, выходя
-то чепуха, и уйти и	бежать	нельзя, точно сидишь в сумасшедшем
банк надоел, не хотелось никуда	идти	, ни о чем говорить. В
-то старушка, а за нею	бежал	знакомый белый шпиц. Гуров хотел
афиша с очень крупными буквами:	шла	в первый раз "Гейша". Он
он - за ней, и оба	шли	бестолково, по коридорам, по лестница
вас всем святым, умоляю... Сюда	идут	! По лестнице снизу вверх кто
лестнице снизу вверх кто-то	шел	. - Вы должны уехать... - продолжала
С. и говорила мужу, что	едет	посоветоваться с профессором насчет
знал об этом. Однажды он	шел	к ней таким образом в
и не застал). С ним	шла	его дочь, которую хотелось ему
градуса тепла, а между тем	идет	снег, - говорил Гуров дочери. - Но



# The Tuzov dictionary

## Improving Nooj ressources

The semantic dictionary by V.Tuzov (Saint Petersburg university) is a most complete electronic dictionary of 190 000 entries, with full morphological and semantical annotation.

- 1. Introduction. The task of Syntactic-Semantic Analysis of Russian Texts.
- 2. Syntactic and semantic analyzers.
- 3. Main principles of V.A Tuzov's theory.
- 4. Sentence analysis.
- 5. The detection of links between words.
- 6. Examples.
- 7. Conclusions.

# The Tuzov dictionary

## Improving Nooj ressources

1<sup>st</sup> job (almost done) : make it complete in Excel format and ready for Nooj annotation

2<sup>nd</sup> make it compatible with the Zaliznyak's dictionary (morphological tagging)

# The Tuzov dictionary

## Improving Nooj ressources

There are 1526 different morphological types  
( that are the same as Zaliznyak's categories !)

There are 7452 semantical different types that can be  
reduced to 1638, using the category table.  
*(the types have been translated to English)*



# The Tuzov dictionary

## Improving Nooj ressources

\$100/	Life+/ /	47
\$100/0	Life+Existence-Inexistence	240
\$100/1	Life+Perception	56
\$100/11	Life+Perception+Viewing	299
\$100/111	Life+Perception+Viewing+Sighted_Blind	38
\$100/112	Life+Perception+Viewing+Far-viewing_Myope	6
\$100/113	Life+Perception+Viewing+Perceptible_Unperceptible	68
\$100/114	Life+Perception+Viewing+Visible_Uninvisible	15
\$100/12	Life+Perception+Hearing	100
\$100/13	Life+Perception+Smelling	44
\$100/14	Life+Perception+Touching	8
\$100/2	Life+Nutrition	599
\$100/3	Life+Respiration	42
\$100/31	Life+Respiration+Breathe	13
\$100/4	Life+Development	84
\$100/5	Life+Activities	18
\$101/1	Life+Food	239

How many tags?

Semantic tag      Signification

Number of words

# The Tuzov dictionary

## Improving Nooj ressources

кормильщик	кищълмрок	n%~пища	\$101/1	{м3о 96}
кормитель	ьлетимрок	n%~пища	\$101/1	{м2о 316}
кормить	ьтимрок	n%~кормление	\$100/2	{г4н 6391}
кормиться	ясьтимрок	n%~кормление	\$100/2	{г4н 10252}
кормление	еинелмрок		\$100/2	{с7 610}
кормленщик	кищнелмрок		\$12413402	{м3о 96}
кормленный	йынелмрок	n%~кормление	\$100/2	{п1 1387}
кормный	йынмрок	n%~кормление	\$100/2	{п1 36}
кормовой	йовомрок	n%~корма	\$121324/2	{п1 553}
кормовой	йовомрок	n%~кормление	\$100/2	{п1 553}
кормодобывание	еинавыбодомрок	n%~добыча	\$1527	{с7 610}
кормодобывающий	йищюавыбодомрок	n%~добыча	\$1527	{п4 7665}
кормозаготовительный	йыньлетивотогазомрок	n%~заготовка	\$15231	{п1 428}
кормозаготовка	аквотогазомрок	n%~заготовка	\$15231	{ж3 168}

Entry

Entry2

Commentary

Semantic tag

Grammatical tag

# The Nooj semantic dictionary

## Improving Nooj resources

быть	п%~будущее	\$160/010/055	{г16сн 7788}
быть	п%~время	\$165	{г16сн 7788}
быть	п%~действие	\$155	{г16сн 7788}
быть	п%~зависимость	\$1241/4123/135	{г16сн 7788}
быть	п%~звание	\$1241/125	{г16сн 7788}
быть	п%~количество	\$12/135	{г16сн 7788}
быть	п%~местопребывание	\$12/35	{г16сн 7788}
быть	п%~нахождение	\$12/35	{г16сн 7788}
быть	п%~носка	\$1525325	{г16сн 7788}
быть	п%~отношения	\$1/42325	{г16сн 7788}
быть	п%~приобретение	\$15310/0/055	{г16сн 7788}
быть	п%~происхождение	\$1241/175	{г16сн 7788}
быть	п%~психика	\$1241/45	{г16сн 7788}
быть	п%~свойство	\$1/10/005	{г16сн 7788}
быть	п%~состояние	\$11135	{г16сн 7788}
быть	п%~существование	\$111015	{г16сн 7788}
быть	п%~свойство	\$1/10/00	{г16сн 7788}

# The Nooj semantic dictionary

## Improving Nooj resources

говор	\$14402	Knowledge+Communication+Information+Speech
говорение	\$14402	Knowledge+Communication+Information+Speech
говорильный	\$14402	Knowledge+Communication+Information+Speech
говорильня	\$14402	Knowledge+Communication+Information+Speech
говорить	\$14402	Knowledge+Communication+Information+Speech
говорить	\$1526031	Action+Work+Intellectual+Commandment+Order
говориться	\$14402105	Knowledge+Communication+Information+Speech+Изложение+Высказывание
сказать	\$14402105	Knowledge+Communication+Information+Speech+Изложение+Высказывание
сказать	\$1526031	Action+Work+Intellectual+Commandment+Order
сказаться	\$1/37	Reason

# Improving semantic resources for Nooj

All the semantics tag are translated into English

## First step:

- one morphological dictionary (already existing)
- one semantical dictionary (soon available), both at Nooj format (soon available)

## Second step:

more derivational and syntactic grammars to provide a full linguistic analysis





# Improving Russian resources for Nooj

Thank you for your attention

[vincent.benet@cnrs.fr](mailto:vincent.benet@cnrs.fr)

[vincent.benet@inalco.fr](mailto:vincent.benet@inalco.fr)

