United Institute of Informatics Problems of the National Academy of Sciences of Belarus

Processing of Quantitative Expressions with Units of Measurement in Scientific Texts as Applied to Belarusian and Russian Text-to-Speech Synthesis

Alena Skopinava, Yuras Hetsevich, Boris Lobanov

Scopes of Application

- Corpora and Database Management Systems, Libraries, Information Retrieval Systems:
- to formulate extended search queries,
- locate specific expressions on the Internet,
- support automatic text annotation and summarization;

Publishing Institutions:

- to automatically locate specified lists of expressions with MUs,
- classify resulting expressions as SI units, their derivatives or units out of the SI,
- check quickly if the extended names of units are used correctly;

Text-to-Speech Synthesis Systems:

 to generate orthographically correct texts and their tonal and prosodic peculiarities;

The main goal is

to develop algorithms and linguistic resources in order to identify, classify and generate measurement units (MU) and quantitative expressions (QE) with them on the material of hand-crafted text corpora for Belarusian and Russian; + to prove its importance for correct intonational marking.

SI derived unit of speed with a multiple prefix

Цягнік рухаўся з хуткасцю 200 км/г у Бекасава. Поезд двигался со скоростью 200 км/ч в Бекасово. 'The train was running at the speed of 200 km/h to Bekasovo.'

дзвесце кіламетраў у гадзіну двести километров в час 'two hundred kilometres per hour'

Three sets of algorithms for BE and RU:

- Identification & Classification of QE with MU according to the SI
- Identification & Classification of QE with MU according to word formation
- Processing of QE with MU into orthographical words

Difficulties

❖ The language-dependent origin;

```
Bel. " на" = Rus. "час" = Eng. "hour" = Ger. "Stunde"...
```

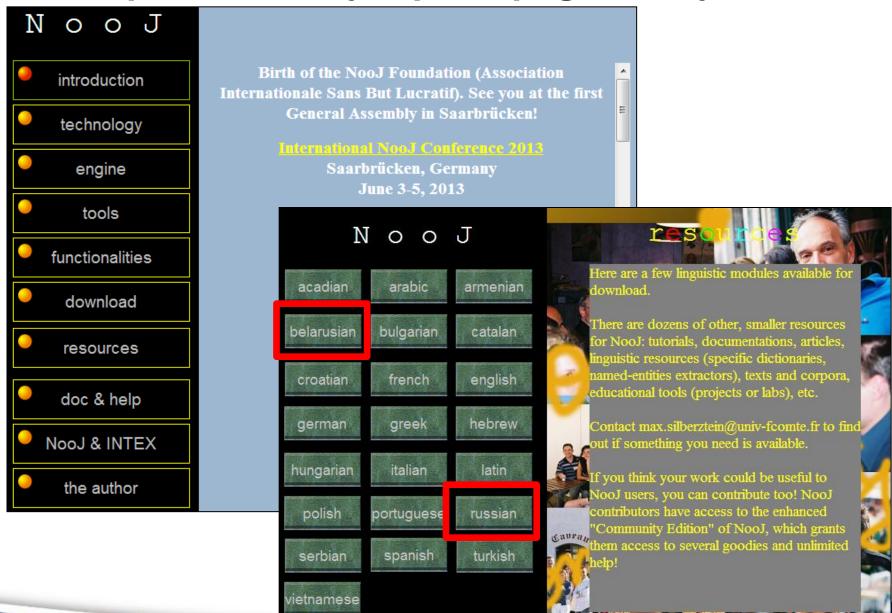
- A great variety of numeral quantifiers and names of units, both in writing and formation;
- intricate agreement:

25 метраў '25 meters', 21 метр_ '21 meter', 23 метры '23 meters'

synonymy of written forms:

```
2000 метраў '2000 meters' = 2000 м '2000 m' = 2 10<sup>3</sup> м '2 10<sup>3</sup> m' = 2 кіламетры '2 kilometers' = 2 км '2 km' = ...
```

http://www.nooj4nlp.net/pages/nooj.html



Scientific-technical text corpora

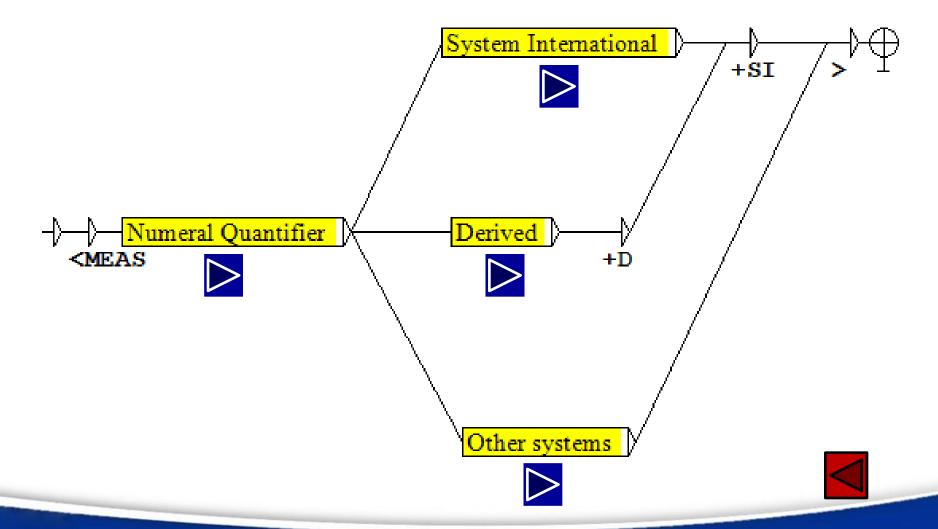
BE

	Size		Кампаніяй DigitalGlobe эксплуатуецца КА высокага разрашэння
Kosmas_1_bel	53887	11.07.2012	QuickBird-2, які быў выведзены на арбіту вышынёй 450 км у 2001 г.
Kosmas_2_bel	72228	11.07.2012	Забяспечвае атрыманне панхраматычных малюнкаў з разрашэннем
			0,64 м і мультыспектральных з разрашэннем 2,44 м у паласе захопу
			16,6 км. Час актыўнага функцыянавання – 7 гадоў.
Kosmas_5_bel	51073		Францыя валодае двума КА SPOT (SPOT-4 і SPOT-5). Спадарожнік
			SPOT-4 функцыянуе з 1998 г. і забяспечвае атрыманне здымкаў з

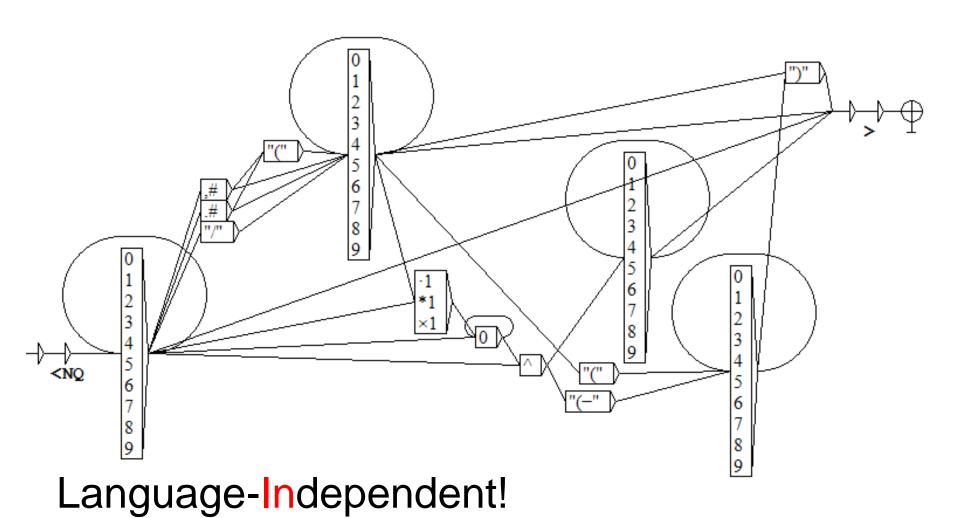
RU

File Name БИОЛОГИЯ	эксперимента единственный выживший на глубине 2 см
БИОТЕХНОЛОГИЯ	экземпляр имел 6 мм в длину и 3 - в ширину, тогда как самые
БИОФИЗИКА БИОХИМИЯ	крупные из 50 особей, выживших на глубине 20 см, достигали
БОТАНИКА	лишь 1.3 мм в длину и 0.77 мм в ширину. Следует отметить,
ВУЛКАНОЛОГИЯ ГЕНЕТИКА	что в каждой рамке находилось более тысячи семян.
ГЕОГРАФИЯ_пер	Первоначально проросли 13 % семян на глубине 1 см и 60% -

1. Identification & Classification of QE with MU according to the SI



Identification of Numeral Quantifiers



Results of NQ Identification

BE

Before	Seq.	After
электратэхнічнай камісіяй) ІЕС	60027	ужываецца пазначэнне Mbit
проста Mb). 1 мегабіт =	1000^2	біт = 10^6 біт = 1000000 біт
Mb). 1 мегабіт = 1000^2 біт =	10^6	біт = 1000000 біт. Дзесятков
Напрыклад: 1/6 = 0,166666 =	0,1(6)	; 1/7 = 0,1428571428 = 0,(14
0,1(6); 1/7 = 0,1428571428 =	0,(142857)	

RU

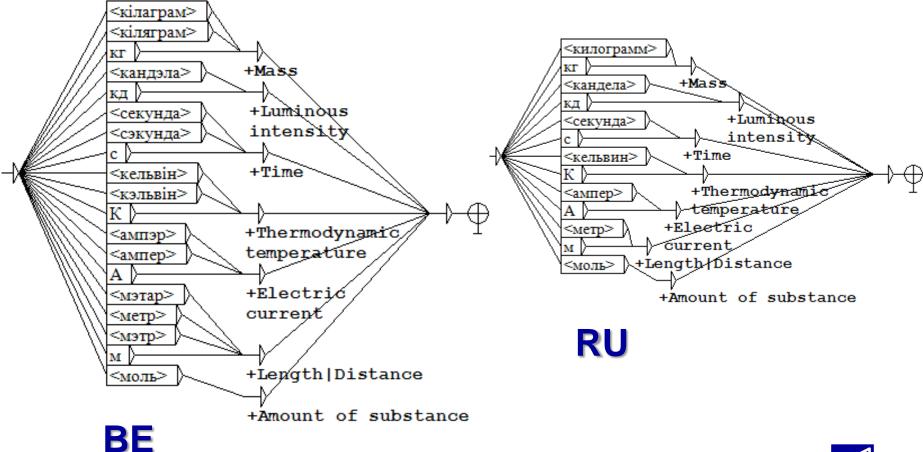
Before	Seq.	After
автомагистралях - не более	110	км/ч, на
двумя осями;	18,75	метра для сочлененного
в среднем составляет	5·10^(-5)	Тл, а на
на экваторе (широта 0°) —	3,1·10^(-5)	Тл. 5.Ом — единица
бомбардировке Хиросимы: около	6.10^13	Дж. Энергия фотона

EN

Before	Seq.	After
is equal to	6.24150974×10 ¹ 8	eV (electronvolts). 1 joule
is equal to	2.3901×10 ⁽⁻⁴⁾	kcal (thermochemical kilod
defined as exactly	0.0254	m, and the
defined as exactly	453.59237	g. Also a
are equivalent to	1/100	. An integer such

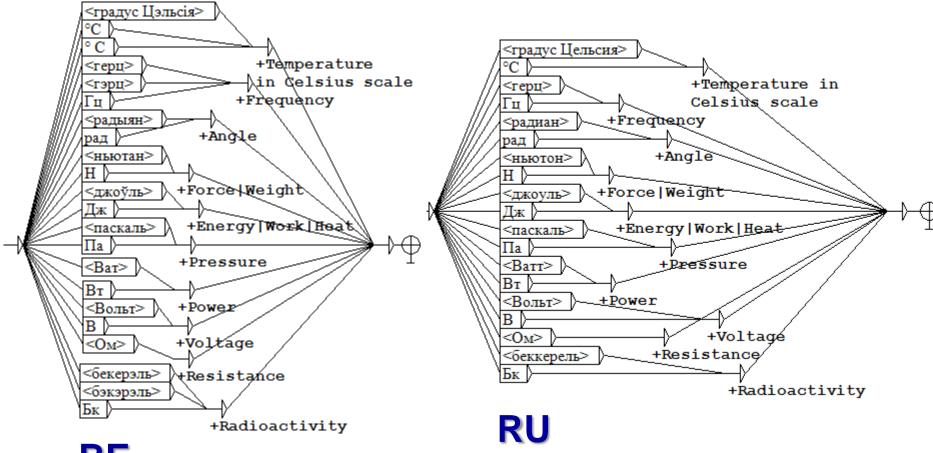


Identification of the SI basic MU





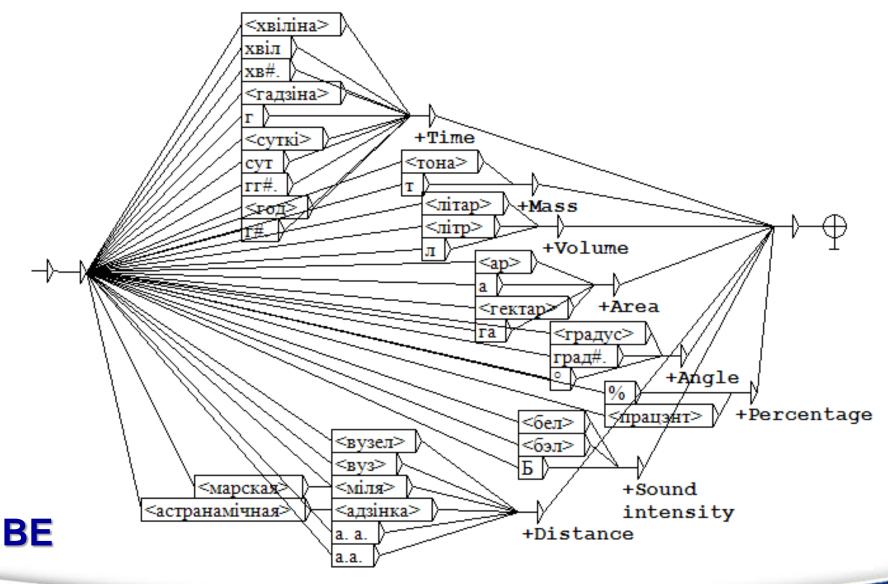
Identification of the SI derived MU



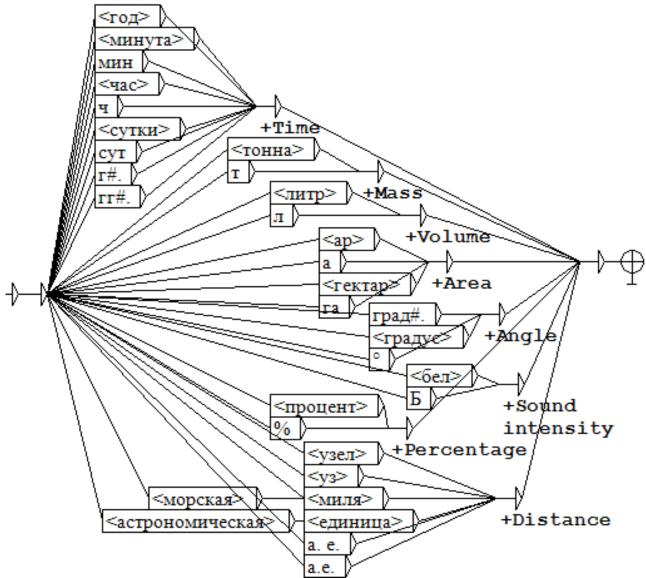
BE



Identification of MU out of the SI



Identification of MU out of the SI







Results of identification of QE with MU according to the SI for Bel & Rus

<MEAS>

BE

Before	Seq.	After
		(бач.), 5
IЭB, 2-30 кэв,	0,1 Гц/ <meas+frequency+d+si></meas+frequency+d+si>	-300 кгц,
	8 т/ <meas+mass></meas+mass>	, выведз
ання Зямлі. У	2005 г./ <meas+time></meas+time>	Іран зда
хвілін і ўхілам	74 градусы/ <meas+angle></meas+angle>	. Затым

1m 0,1 Hz 8 t year 2005 74 degrees

EN

RU

Before	Seq.	After
температуре	109 K/ <meas+thermodynamic td="" temperature+<=""><td></td></meas+thermodynamic>	
гимостью ок.		. Железныє
рии (а спустя	33 года/ <meas+time></meas+time>	– и его сын
э превышало	5°/ <meas+angle></meas+angle>), а потом на
ратуре выше	600° C/ <meas+temperature celsius="" in="" scal<="" td=""><td>, а халькоге</td></meas+temperature>	, а халькоге

109 K 200 000 I 33 years 5 600 C

Results of identification of QE with MU according to the SI for Bel & Rus

 $= <\!\!MEAS + SI + D\!\!>$

Before	Seq.	After
МЭВ, 2-200 МЭВ, 2-30 кэв,	0,1 Гц	-300 кгц, 0-50 кгц; перыядычнасць
ткавай атмасферы складаў	400 Па	, працягласць плазмавага імпульсу
гэтага тэрмомэтра пры	0 °C	, і літара, якая
у дыяпазоне ад -	50 °C	да +200 °C. Залежнасьць
ад −50 °C да +	200 °C	. Залежнасьць супору ад
у дыяпазоне ад -	260 °C	да +1100 °C. Залежнасьць

	•
0,1	Hz
400	Pa
0	C
50	C
200	C
260	C

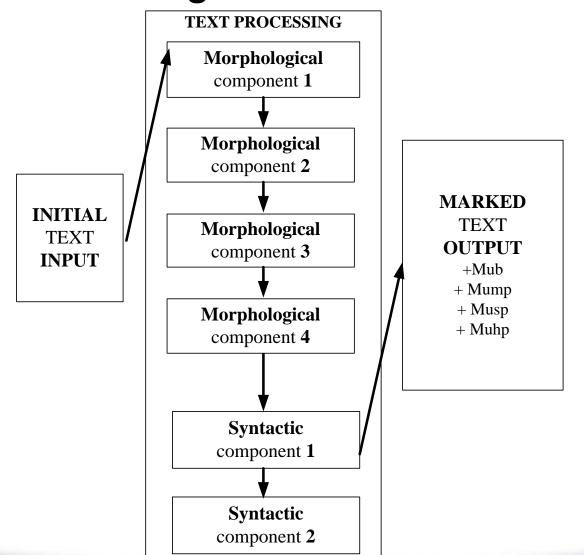
EN

RU

Before	Seq.	After
МэВ, 2–200 МэВ, 2–30 кэВ,	0,1 Гц	–300 кГц, 0–50 кГц
лазер накачки мощностью	25 Вт	возбуждает лазер
красителях выходной мощность		, который и дает
всех металлов теплоемкостью:	16,44 Дж	/(моль К) для
К) для ⊟-Ве,	30,0 Дж	/(моль К) для
С она составляет	209,3 Вт	/ (м К), что

0,1 Hz 25 W 4,25 W 16,44 J 30,0 J 209,3 W

2. Identification & Classification of QE with MU according to Word Formation





Excerpts from the NooJ dictionaries of basic stems

BE RU

Б,ABBREVIATION+Mbase
байт,NOUN+FLX=БАЙТ+s2+UNAMB+Base
бекерэль,NOUN+FLX=ABAЛЬ+s6+UNAMB+Base
біт,NOUN+FLX=БАЙТ+s2+UNAMB+Base
В,ABBREVIATION+Mbase
Вт,ABBREVIATION+Mbase
ват,NOUN+FLX=БАЙТ+s2+UNAMB+Base
вольт,NOUN+FLX=БАЙТ+s2+UNAMB+Base
г,ABBREVIATION+Mbase
га,ABBREVIATION+Mbase

гектар.NOUN+FLX=ГЕКТАР+s5+UNAMB+Base герц,NOUN+FLX=AMПЕР+s2+UNAMB+Base год,NOUN+FLX=ГОД+sN+UNAMB+Base град,ABBREVIATION+Mbase грам,NOUN+FLX=ГРАМ+s3+UNAMB+Base Гц,ABBREVIATION+Mbase гц,ABBREVIATION+Mbase

ампер, NOUN+FLX=AЛТЫH+s4+UNAMB+Base A.ABBREVIATION+Mbase байт.NOUN+FLX=АБАЖУР+s2+UNAMB+Base бит,NOUN+FLX=АБАЖУР+s2+UNAMB+Base **5.ABBREVIATION+Mbase** ватт, NOUN+FLX=AЛТЫH+s2+UNAMB+Base BT.ABBREVIATION+Mbase вольт.NOUN+FLX=AЛТЫH+s2+UNAMB+Base B,ABBREVIATION+Mbase гектар, NOUN+FLX=AБAЖУР+s5+UNAMB+Base га, ABBREVIATION+Mbase герц.NOUN+FLX=ГЕРЦ+s2+UNAMB+Base Гц.ABBREVIATION+Mbase год,NOUN+FLX=ГОД+sN+UNAMB+Base r,ABBREVIATION+Mbase rr.ABBREVIATION+Mbase град.ABBREVIATION+Mbase рамм,NOUN+FLX=АНГСТРЕМ+s3+UNAMB+Base

Describing Inflectional Classes...

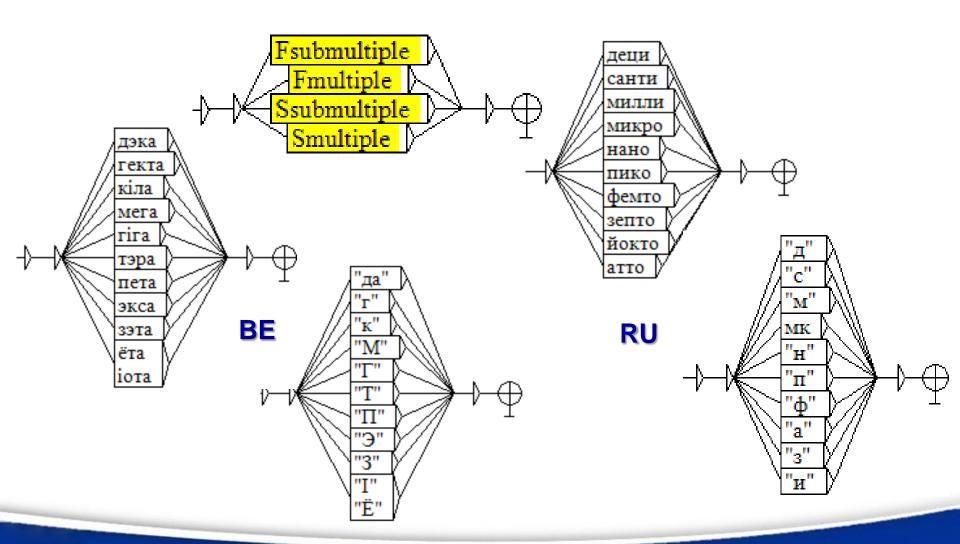
BE FEKTAP =

- <E>/Accusative+Common+Inanimate+Masculine
- + <E>/Common+Inanimate+Masculine+Nominative
- + <E>a/Common+Genitive+Inanimate+Masculine
- + <E>aÿ/Common+Genitive+Inanimate+Masculine+Plural
- + <E>aм/Common+Inanimate+Instrumental+Masculine
- + <E>aм/Common+Dative+Inanimate+Masculine+Plural
- + <E>aмi/Common+Inanimate+Instrumental+Masculine+Plural
- + <E>ax/Common+Inanimate+Masculine+Plural+Prepositional
- + <E>ы/Common+Inanimate+Masculine+Prepositional
- + <E>y/Common+Dative+Inanimate+Masculine
- + <E>ы/Accusative+Common+Inanimate+Masculine+Plural
- + <E>ы/Common+Inanimate+Masculine+Nominative+Plural;

RU AHFCTPEM =

- <E>/Common+Genetive+Inanimate+Masculine+Plural
- + <E>/Accusative+Common+Inanimate+Masculine+Singular
- + <E>/Common+Inanimate+Masculine+Nominative+Singular
- + <E>a/Common+Genetive+Inanimate+Masculine+Singular
- + <E>ам/Common+Dative+Inanimate+Masculine+Plural
- + <E>ами/Common+Inanimate+Instrumental+Masculine+Plural
- + <E>ax/Common+Inanimate+Masculine+Plural+Prepositional
- + <E>e/Common+Inanimate+Masculine+Prepositional+Singular
- + <E>oB/Common+Genetive+Inanimate+Masculine+Plural
- + <E>ом/Common+Inanimate+Instrumental+Masculine+Singular
- + <E>y/Common+Dative+Inanimate+Masculine+Singular
- + <E>ы/Accusative+Common+Inanimate+Masculine+Plural
- + <E>ы/Common+Inanimate+Masculine+Nominative+Plural;

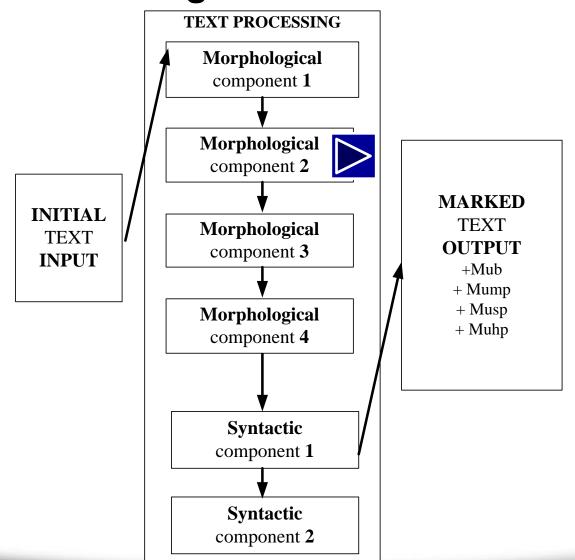
Databases of Prefixes according to the SI



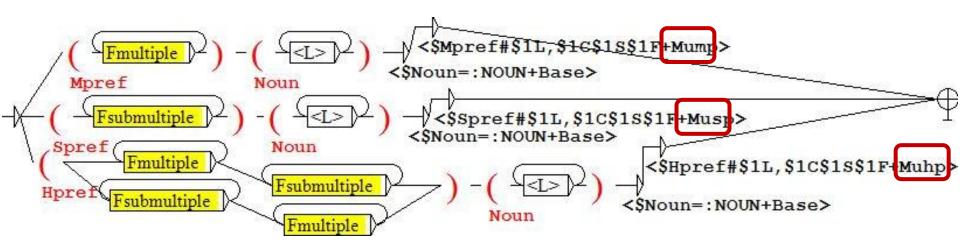
Word-formative classification of MU:

- 1. with full-form stems and without prefixes метр, Герц, Ом (eng. meter, Hertz, Ohm)
- 2. with full-form stems and full-form prefixes нанафарады, міліампер (eng. nanofarads, milliampere)
- 3. with full-form stems and shortened prefixes кБайт (eng. Kbyte)
- 4. with shortened stems and without prefixes Дж, га, Па (eng. J, ha, Pa)
- 5. with shortened stems and shortened prefixes км, дл, гПа (eng. km, dL, hPa)

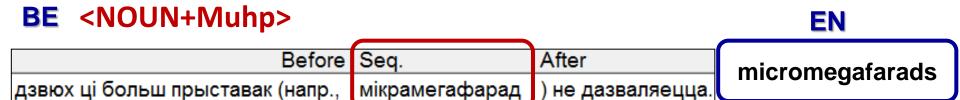
2. Identification & Classification of QE with MU according to Word Formation



M2 identifies <u>full-stem</u> MU with <u>full</u> multiple and/or submultiple prefixes



M2-operation results



RU <**NOUN+Mump>**

на расстоянии несколько сотен	километров	. Первый вариант ударного
кусков может достигать нескольких	килограммов	. Куски брони поражают
излучения мощностью в сотни	мегаватт	. Проблема в том
составляет уже десятки тысяч	мегагерц	, что соответствует волнам
своих жестких дисках тысячи	гигабайт	информации, третьи подкл
нергии лазерного излучения порядка	мегаджоуля	(106 Дж) и кпд

Before Sea.

kilometres kilograms megawatt megahertz gigabyte megajoule

Examples of annotated word forms

ВЕ
дэкалітрамі.

317

дэкалітр, NOUN+Meaning=Common
+Animation=Inanimate
+Case=Instrumental
+Gender=Masculine+Number=Plural
+s2+Meas=Base+Mump

```
наносекундами

□

наносекунда,NOUN+ProperCommon=Common

+Gender=Feminine+Animation=Inanimate

+Case=Instrumental+Number=Plural

+s4+Meas=Base+Musp

>
```

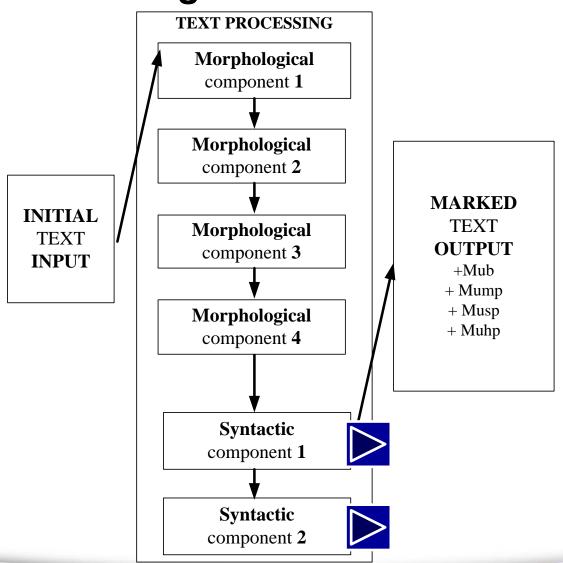
RU

There are <u>no</u> word forms "дэкалітрамі" (Eng. decaliters) "наносекундами" (Eng. nanoseconds) <u>in the resource dictionary</u>.

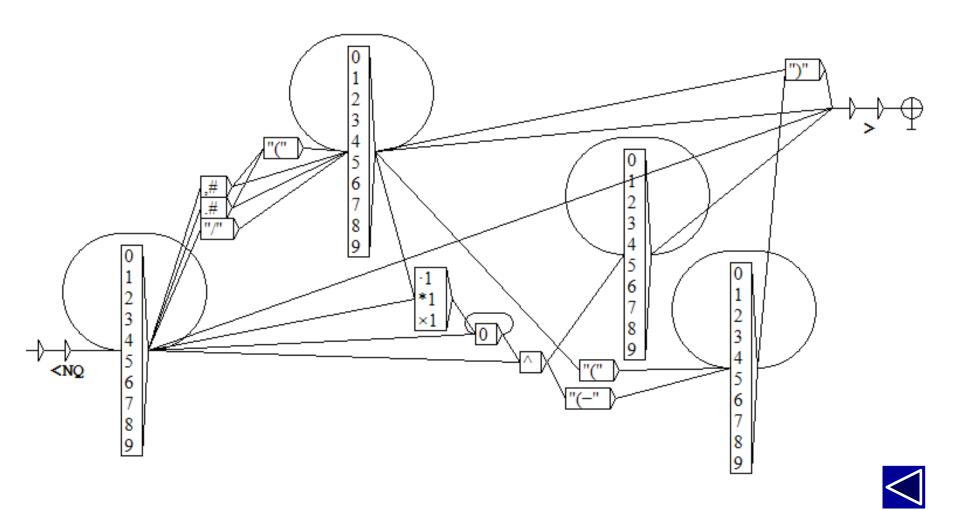
BUT!

All the morphological <u>data are preserved</u> by the algorithm (according to the grammatical characteristics of the basic stem).

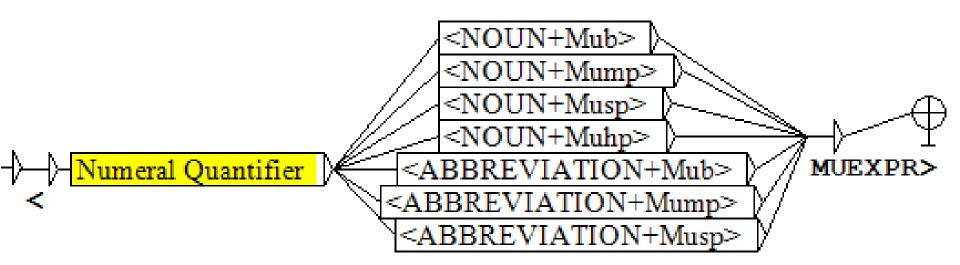
2. Identification & Classification of QE with MU according to Word Formation



S1 for <u>Numeral Quantifiers</u> is the same.



S2 <u>collects</u> M1-M4 data and <u>identifies</u> QE with MU



Excerpts of the results obtained with the S2

<MUEXPR> EN BE

Before Seq. масы - грам (0,001 кг). йна вар'іруецца зблізку на апору з сілай нанафарадамі (пішуць асць шара з радыусам Mbit(альбо проста Mb). святло ў вакууме за (ıазон - 40 кэв-3 МЭВ, 2ы тэлевізара - парадку ьмічных прамянёў - ад

31 мкТл 2.4 м3в 9.81 H 60 000 пф 1 сантыметр 1 мегабіт 1 / 299 792 458) секунды 200 M3B 20 кілаэлектронвольт

After (3,1×10⁽⁻⁵⁾ Тл) - напружанасць у год. 1 Н ёсць Прыбліжэнне, што 1 кг адпавяд а не 60 нф; 2 000 мкф змешчанага ў вакуум. 1 сантым = 1000^2 6it = 10^6 6it = 100000 Метр быў упершыню ўведзень 2-30 кэв, 0,1 Гц-300 кгц, 0-50 кг

Энэргіі касьмічных прамянёў -

да 1000 тэралектронвольтаў.

31 µT 2,4 mZv 9.81 N **60 000 pF** 1 sm 1 megabit (1/299 792 458) seconds **200 MeV** 20 kiloelectronvolts 1 megaelectronvolt

RU

Before Seq. организм ток не превышал могут сказать «файл в тивление величиной от 1 до омегафарад пикотеравольт до 64 Мбит/с) и энширский изумруд» массой ремя жизни мюонов - около а которой оказалась равной то они оказались равными: высота 670 км - наклонение

1 мА

1 мегаэлектронвольта

100 килобайт

100 MОм

13 йоттайоктограммов

137,4 МГц

1383,95 каратов

2.2 мкс

22 фемтограммам

8.1.10^21 Дж

98,00 град

After

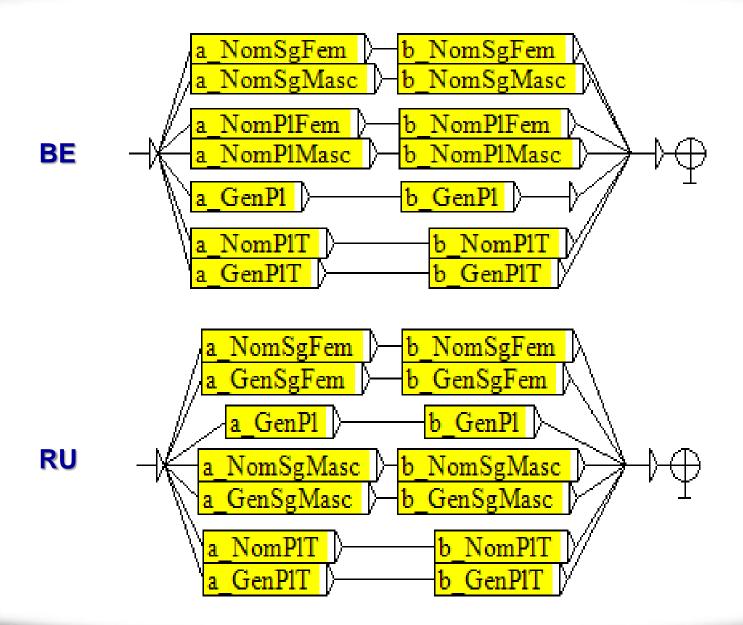
. На человека токи статиче »). При обозначении скоро , чтобы протекающий чере Каждая строка содержит (метровый диапазон, фор Изумруды выращивают і - осложняет задачу созда (1 фг = 1•10^(-15) г). . Мю (уменьшение массы ледні Срок активного существо

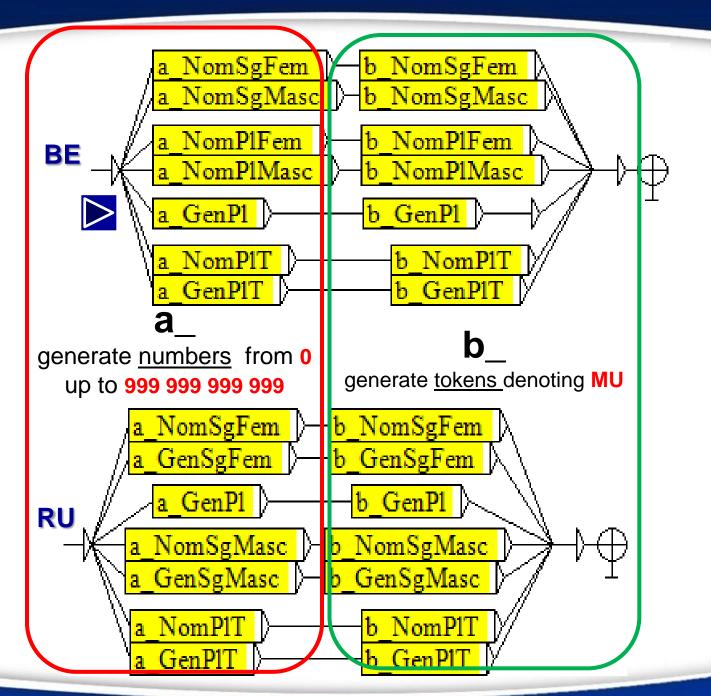
1mA 100 kilobytes **100 Mohm** 13 yottayoktograms 137,4 MHz 1383,95 carats $2.2 \mu s$ 22 femtograms 8.1·10²¹ J 98,00 deg

3. Generating QE with MU into orthographical words

тем тысяч метраў ВЕ 7000 м семь тысяч метров RU

7000 m ----> seven thousand metres EN





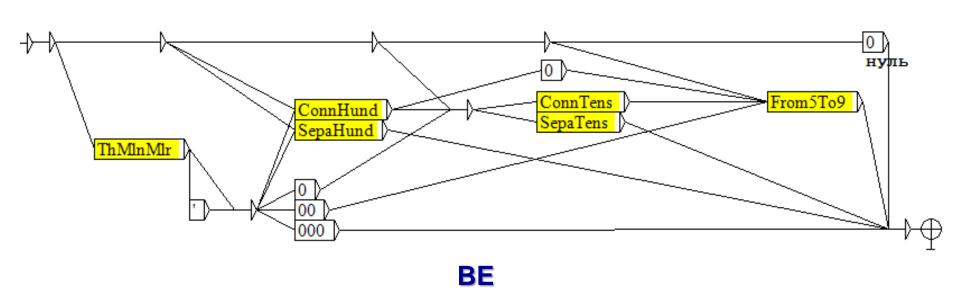


Peculiarities of the inflection of nouns after numerals

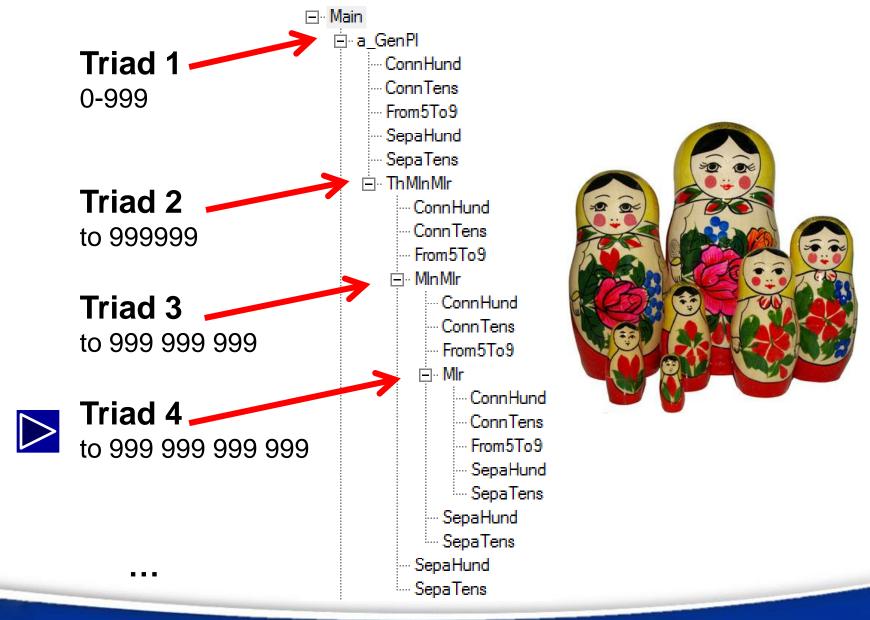
- 1. After number $\mathbf{1}$ (or with $\mathbf{1}$ as a final digit) nouns take endings of the Nominative Singular (NomSg). QE will proceed to branches $\mathbf{1}$ or $\mathbf{2}$, depending on the gender of nouns, in particular Masculine (Masc) or Feminine (Fem).
- 2. After numbers **2**, **3**, **4** (or with 2, 3 or 4 as a final digit) nouns take the Nominative plural (*NomPl*) in Belarusian, whereas in the Russian these numbers require nouns in the Genetive singular (*GenSg*). Depending on the gender, QE will move to branches **3** or **4**.
- 3. Numbers **from 5 to 19** and **round** numbers (or with them as final digits) require nouns in the Genitive plural (*GenPl*) in both languages. QEs will follow the **5**th branch. Branches **6** and **7** are for pluralia tantum nouns.
- 4. Special inflection is demanded by **pluralia tantum nouns**: *cymкu* and *cymкi* (Eng. *twenty-four hours*). They take branches **6** or **7**.

a_GenPl

from 0 to 999 999 999 999 Genitive plural

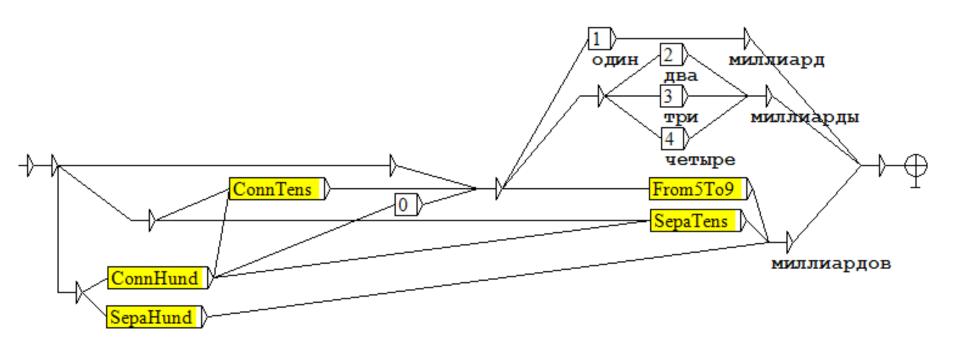


Scheme of the BE ramified algorithmic complexes



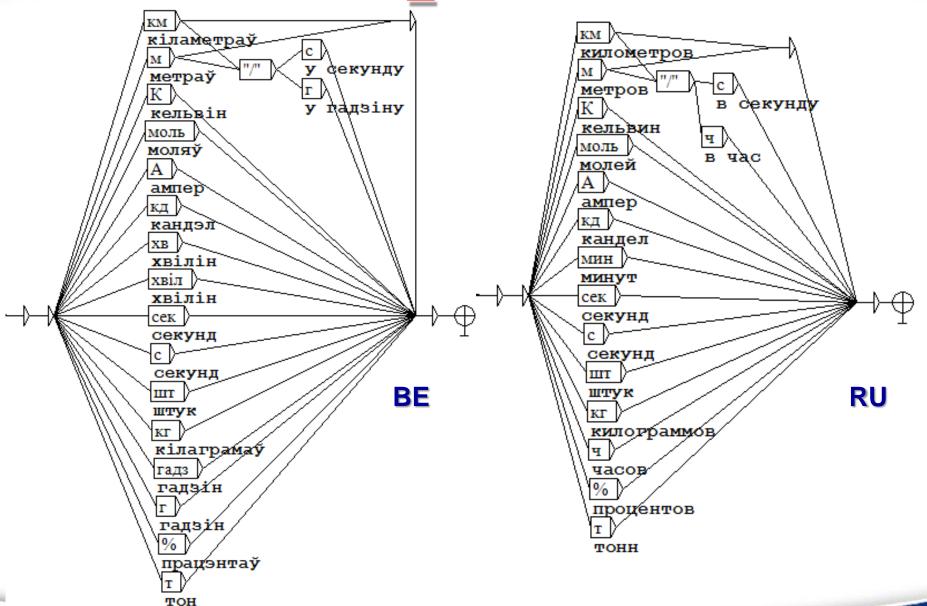


It generates any whole number from the class of **billions**.



RU

b_GenPl



Excerpts of <u>results</u> of <u>generating</u> QE with MU BE into orthographical words

Seq.

15 т/пятнаццаць тонаў

123 кг/сто дваццаць тры кілаграмы

5678904 К/пяць мільёнаў шэсцьсот семдзесят восем тысяч дзевяцьсот чатыры кельвіны 123 А/сто дваццаць тры амперы

6785672 м/шэсць мільёнаў семсот восемдзесят пяць тысяч шэсцьсот семдзесят два метры 787879 сут/семсот восемдзесят сем тысяч восемсот семдзесят дзевяць сутак

6761 сут/шэсць тысяч семсот шэсцьдзесят адны суткі

99999994 моль/дзевяцьсот дзевяноста дзевяць мільёнаў дзевяцьсот дзевяноста дзевяць тысяч дзевяцьсот дзевяноста чатыры молі

RU

6661с/шесть тысяч шестьсот шестьдесят одна секунда 77700 т/семьдесят семь тысяч семьсот тонн 800009кд/восемьсот тысяч девять кандел 120202 мин/сто двадцать тысяч двести две минуты 8600км/ч/восемь тысяч шестьсот километров в час 903 м/с/девятьсот три метра в секунду

The regular expression for a search for QE

```
(or <NB> go <NB>)(or <NB> go <NB>,(NB>,(NB> go <NB>,(NB>,(NB> go <NB>,(NB> go <NB))((NB> go <NB>,(NB> go <NB>,(NB> go <NB))((NB> go <NB) go <NB)((NB> go <NB)((NB> go <NB) go <NB)((NB> go <NB)((NB> go <NB) go <NB)((NB> go <NB> g
```

(<WF><NB>)

Excerpts: variety of QE,

revealed with the help of the obtained formula

Formula's constituents	ormula's constituents Examples of QE, found with the help of a certain constituen								
(от <nb> до <nb>)</nb></nb>	находится в пределах от 40 до 347 г/л. Несмотря лежит в диапазоне от 2000 до 3000 А, и следовательно удаленной на расстояние от 160 до 640 км от восточного мере удаления "Пионера-10" от 40 до 60 а.е. величина								
(от <nb> до <nb><nb>)</nb></nb></nb>	находятся в диапазоне от 4000 до 14 000 МГц и даже вариации) с периодами от 10 до 10 000 лет, сосредоточенными в								
(от <nb><nb> до <nb>)</nb></nb></nb>	см (частотный интервал от 600 000 до 1000 МГц) относится к								
(от <nb>,<nb> до <nb>,<nb>)</nb></nb></nb></nb>	семь камней массой от 94,45 до 4,39 кар каждый. Кроме между рельсами пути от 1,52 до 1,68 м) типичны для кристалла со скоростями от 0,01 до 0,001 скорости света; их								
(от <nb>,<nb> до <nb>)</nb></nb></nb>	видов токсичность составляет от 0,01 до 1 %. Практически это означает								
(от <nb>.<nb> до <nb>.<nb>)</nb></nb></nb></nb>	диапазоне длин волн от 1.3 до 1.7 мкм. На подложке								
(от <nb> до <nb>,<nb>)</nb></nb></nb>	отношениях ионных радиусов от 1 до 0,732 (рис. 4,a). При С-диапазоне и от 12 до 12,7 ГГц в Q								
(от <nb>-<nb> до <nb>- <nb>)</nb></nb></nb></nb>	выемчатые. Их длина от 1-2 до 30-40 см. Самые длинные длиной волны I от 10-3 до 10-8 м. Этот диапазон								
(от <nb>×<nb>-<nb> до <nb>×<nb>-<nb>)</nb></nb></nb></nb></nb></nb>	с удельным сопротивлением от 5×10-8 до 8×10-5 Ом хм. Композиционные								
(от <nb>×<nb>-<nb> до <nb>×<nb>-<nb>)</nb></nb></nb></nb></nb></nb>	в разных материалах: от 3×10-6 до 2×10-5 см. Магнитный поток								
(от <nb> до почти <nb>)</nb></nb>	током (при этом от 50 до почти 100 % его энергии превращается								

Distribution of various ways of expressing QE in thematically distinct texts

Text	Physics	Space Travel Science	Geography	Military Equipment	Mineralogy	Botany	Transportation Communications	History
Number of variations, a	51	23	22	19	18	16	14	9
Number of numeral expressions, b	2841	2245	2765	9961	3668	1407	2066	4198

Intonation (Prosodic) marking in sentences with QE/MU

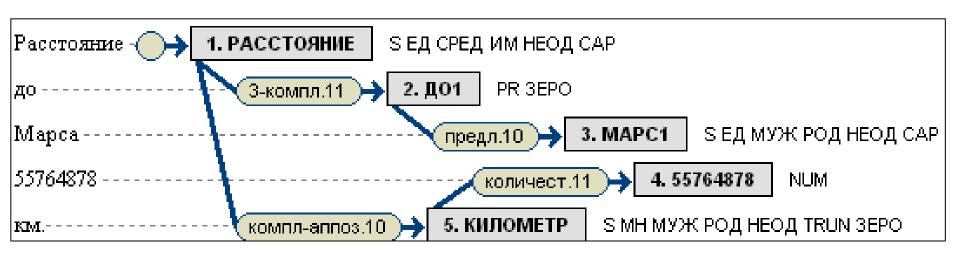
- dividing into syntagmas
- marking emphatically highlighted words
- indicating syntagmas with accent units
- creating a melodic contour of each syntagma

BUT!

First texts are reduced to a normalized orthographic form.

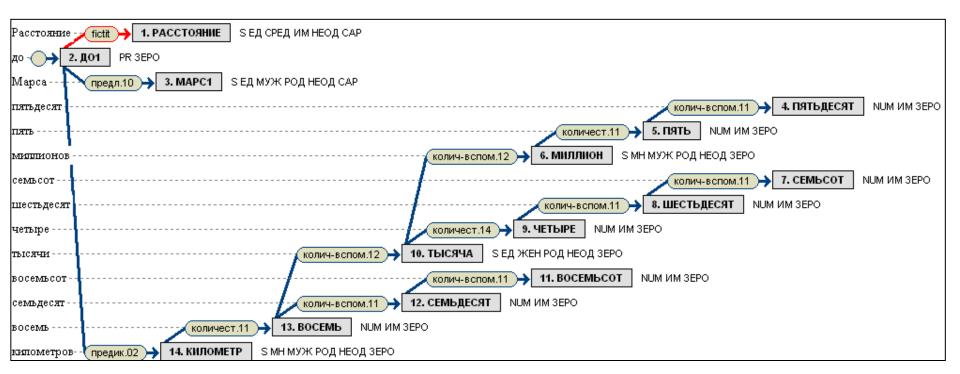
Расстояние до Марса 55764878 км

'The distance to Mars is 55764878 km'



Расстояние до Марса пятьдесят пять миллионов семьсот шестьдесят четыре тысячи восемьсот семьдесят восемь километров

'The distance to Mars is fifty-five million seven hundred sixty-four thousand eight hundred seventy-eight kilometers'

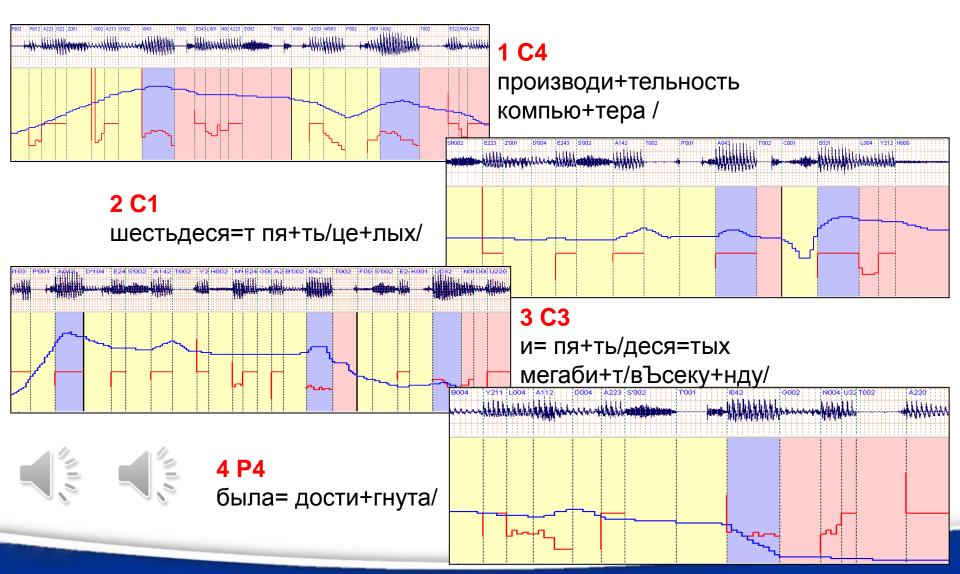






Производительность компьютера 65.5 Мбит/с была достигнута

'The computer performance 65.5 Mbit/s was achieved'



Conclusion

- Three sets of NooJ grammars for BE and RU have been obtained:
- Identification & Classification of QE with MU according to the SI



2. Identification & Classification of QE with MU according to word formation



3. **Generating** QE with MU into **orthographical** words



QE with MUs can get the correct intonation marking only after they are properly generated, i. e. expanded into orthographical words.

Future work includes:

 generating numeral quantifiers, expressed by fractional and decimal numbers

- **disambiguation**, e. g., in such cases when algorithms "confuse" some units (the same initial letter г for год 'year', грам 'gram' на 'hour'
- developing algorithms that will identify numeral quantifiers expressed not only by numbers (mathematical objects), but also by numerals (parts of speech)

THANK YOU FOR ATTENTION!

E-mail contacts:

Yury.Hetsevich@gmail.com skelena777@gmail.com lobanov@newman.bas-net.by