

## Markup and meter

Using XML tools to teach a computer to think about versification

David J. Birnbaum and Elise Thorsen  
<http://poetry.obdurodon.org>

Balisage 2015: The markup conference, 2015-08-11

## About this presentation

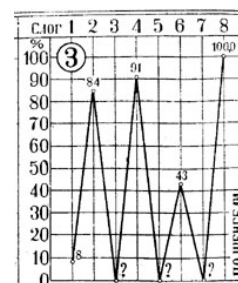
- Poetry scholars care about form
  - Meter, rhyme etc.
  - Can we identify form without massive human effort?
- What else
  - Possibly a different way to think about poetry
  - And because this is Balisage:
    - Possibly a different way to think about overlap

## Outline

- Verse scholarship
  - Research goals and the research context
  - Stress, meter, and rhyme
- Markup
  - Explicit and implicit markup
  - Explicit and implicit overlap
  - Elements without borders
  - Mixed content as a type of overlap

## Research framework

- Input
  - Russian plain-text verse corpus in normal orthography
- Output
  - Meter, rhyme, and other formal features
  - Text- and corpus-level analysis



## Why it's hard intellectually

- Written text is accessible
  - But it lacks information about stress and pronunciation
- Meter and rhyme depend on the place of stress
  - The place of stress is not marked
  - The place of stress is not predictable without linguistic knowledge
  - Metrical and linguistic emphasis may differ
- Rhyme further depends on pronunciation

## Why it's hard technically

- Overlapping hierarchies
- Mixed content

## Assumptions

- Russian quantitative verse studies are worth doing
  - Andrej Belyj, 1910, *Simvolizm*; Jurij Tynjanov, 1924, *Problema stixotvornogo jazyka*; Viktor Žirmunskij 1925, *Vvedenie v metriku. Teorija stixa*; Kiril Taranovski, 1953, *Ruski dvodelni ritmovi*; Boris Ejsenbaum, 1969, *O poëzii*; Mixail Gasparov, 1984, *Očerki istorii russkogo stixa. Metrika, ritmika, rīfma, strofika*
  - Vladimir Nabokov, 1964, *Notes on prosody*; J. Thomas Shaw, 1993, *Pushkin's poetics of the unexpected: The nonrhymed lines in the rhymed poetry and the rhymed lines in the nonrhymed poetry*; Ian K. Lilly, 1995, *The dynamics of Russian verse*
  - Handbooks and textbooks: Boris Unbegaun 1956, Barry Scherr 1986, Michael Wachtel 2004
  - Generative poetics: Morris Halle, Bruce Hayes, Paul Kiparsky
  - Names to watch: James Bailey, Nila Friedberg, Emily Klenin, Barry Scherr, J. Thomas Shaw, Marina Tarlinskaja
- Target corpus is generally regular syllabotonic verse: stanzas, lines, feet, meter, rhyme

## Lexical stress vs metrical ictus

No longer mourn **for me when I** am dead  
 Than you shall hear the surly sullen bell  
 Give warning **to** the world **that I** am fled  
 From this **vile world** with vilest worms to dwell:  
 [Shakespeare, Sonnet 71, iambic pentameter]

○ x | ○ x | ○ ○ | ○ ○ | ○ x  
 ○ x | ○ x | ○ x | ○ x | ○ x  
 ○ x | ○ ○ | ○ x | ○ ○ | ○ x  
 ○ x | x x | ○ x | ○ x | ○ x

## Lexical stress vs metrical ictus

- Pyrrhic (o o), spondee (x x), trochaic (x o) substitutions in iambic (o x) verse
- Metrical variation
  - Preserves meter, while preventing poetry from becoming “sing-song”
  - Establishes associations among words and lines
  - Modulates the tempo
  - Draws attention to important moments
  - Adapts international meter to local linguistic properties (stress system, word length)

## Meter and language: orthography

- In English
  - The relationship between vowel letters and vowel sounds (syllables) is not one to one
- In Russian
  - Every vowel is syllabic
  - No silent vowels (cf. English *Adelaide*)
  - No representation of single vowel sounds by sequences of vowel letters (cf. Eng. *Adelaide*)
- Which means
  - Vowel *letters* in Russian are surrogates for *syllables*

## Meter and language: stress

- English
  - Long words often have secondary stress
- Russian
  - Secondary stress only in compound words: **трѣхэта́жный** trěxetàžnyj ‘three-story’
  - Otherwise Russian words, no matter how long, have only primary stress: **до́стопримеча́тельность** dostoprimečatel’nost’ ‘(tourist) attraction’
  - So what happens with long words in binary meter?

## Implicit meter and actual stress

национальный корпус  
РУССКОГО  
ЯЗЫКА

[перейти на страницу поиска](#) [выбрать корпус](#) [версия без ударений](#) [инструмент](#) [формат KWIC](#)

[Результаты поиска](#)

мой  
на расстоянии 1 от лица  
на расстоянии 1 от имен  
на расстоянии 1 от чисел  
на расстоянии 1 от прилаг

Найдено 1 вхождение

Страница: 1

Л. А. С. Пушкин, Евгений Онегин / Глава первая (1823-1824) [оформление не считая] [Всё примеры \(1\)](#)

«Мой **адам** **самых** **чистых** **прива**,  
 Когда не в шутку **адам**,  
 Он **укажет** **себя** **востан**  
 И **лучше** **вдохнуть** **не** **мо**»

[А. С. Пушкин, Евгений Онегин / Глава первая (1823-1824)] [оформление не считая] [KWIC](#)

## Meter and language: word length

- Average word length in Shakespeare Sonnet 71 is 3.8 letters
  - Lots of short words
- Average word length in first stanza of Pushkin's *Eugene Onegin* in Russian is 9.5 letters
  - Lots of long words
- Neither English nor Russian fits binary meter naturally

## Meter and language: verse convention

- Russian
  - Strong sense of line
  - Strong sense of foot
  - Strong syllabotonic orientation
- English
  - Stronger role for tonic organization

## "The old woman of Berkeley"

Robert Southey (1774–1842; 1799)

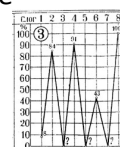
The ra ven croaked   as she sate   at her meal,	2 2 3 3	a
And the Old   Woman knew   what he said;	3 3 3	b
And she   grew pale   at the Ra ven's tale,	2 2 3 2	c
And sick ened, and went   to her bed.	2 3 3	b

Vasilij Andrejevič Žukovskij (1783–1852; 1814/1831)

На кро вле во рон ди ко про кричал —	2 2 2 2	a
Стару шка слы шит и   бледнеет.	2 2 2 2+	B
Понят но ей,   что во рон тот   сказал:	2 2 2 2 2	a
Слегла   в постель,   дрожит,   хладеет.	2 2 2 2+	B

## What quantitative metrics tells us about Russian verse

- Final stress must always be realized
- "Law of regressive accentual dissimilation" (Taranovski)
  - Pre-final foot is weakest
  - Iambic tetrameter: 2 3 1 4
  - Iambic pentameter: 3 2 4 1 5
- Pattern holds over 18<sup>th</sup>, 19<sup>th</sup>, 20<sup>th</sup> centuries (Friedberg), but with changes
- No such regularity in English (Tarlinskaja)



[Image from: N. V. Lapšina, I. K. Romanovič, and B. I. Jarxo, *Metričeskij spravočnik k stixotvorenijam A. S. Puškina*, Moscow: Academia, 1934, p. 134bis. <http://feb-web.ru/feb/pushkin/critics/jar/jar-005-.htm>]

## What the system should tell us

- For individual poems:
  - Identify which syllables are stressed linguistically
  - Identify metrical structures and ambient meter
  - Identify deviations from the ambient meter
  - Identify rhyme schemes
  - (Other formal regularities?)
- Corpus level
  - Historical patterns (authors, periods, movements)
  - Relationships between form and meaning
    - E.g., semantic halo

## TEI

```
<div type="book" n="1"
  met="-+|-+|-+|-+/-" rhyme="aa">
```

- Inline vs standoff
  - Meter and rhyme stand apart from text
- Data vs metadata
  - The human analyzes the line and writes the result
- Tagging the text vs (meta)tagging the poem

## Tagging the text vs the poem

```
<|>
<w>
<orth>Я</orth>
<str>Я</str>
</w>
<w>
<orth>номною,</orth>
<str>но<stress>о</stress>мною</str>
</w>
<w>
<orth>ровопок</orth>
<str>ровоп<stress>о</stress>к</str>
</w>
<w>
<orth>еро</orth>
<str>ер<stress>о</stress></str>
</w>
</>
```

## Sample browsing output

**<oo> → <met> Meter, rhythm, and rhyme**

Maintained by: David J. Hirston ([djhirston@gmail.com](mailto:djhirston@gmail.com))  
Last modified: 2013-05-24T15:00:00Z

**Его речь (Борис Леонидович Пастернак)**

Line	Text	Meter	Rhyme	Stressed Vowels	
1	Я помню говорю <b>его</b>	ox   oo   ox   ox	a	o o o	100
2	Прошли мне искрами загорюс	ox   ox   oo   ox (o)	B	i i i	90
3	Как шорох мильи шаровой	ox   ox   oo   ox	c	o o o	80
4	Все встало с мест глазами <b>втуго</b>	ox   ox   ox   ox (o)	D	A i A u	70
5	Обширная крайний стол	ox   oo   ox   ox	e	A A o	60
6	Как курят он вприс-но <b>пробуго</b>	ox   ox   oo   ox (o)	D	u i u	50
7	И вприс-но раньше чем <b>внуло</b>	ox   ox   ox   ox	e	i A u o	40
8	Он прескромно неусладно	oo   ox   oo   ox (o)	F	u i	30
9	Сквозь строй пролетел и <b>подлет</b>	ox   ox   ox   ox	g	o A o	20
10	Как <b>лет</b> и коннату без <b>дыма</b>	ox   ox   oo   ox (o)	F	i o i	10
11	Грози <b>влет</b> ающий <b>комко</b>	ox   ox   oo   ox	g	i A o	
12	Он <b>был</b> как <b>выпад</b> на <b>раппе</b>	ox   ox   oo   ox (o)	H	i i i	
13	Топнсь за <b>высказанным</b> <b>вска</b>	ox   ox   oo   ox	i	A i i	

## From plain text input to rich output

- Input must be in native Russian orthography
  - Native Russian orthography almost never marks stress
- Meter
  - Meter depends on stress
- Rhyme
  - Rhyme depends on pronunciation
  - Pronunciation can be inferred from orthography only if stress is also known
- But if we can determine stress automatically ...

## Procedure

1. Make stress explicit (dictionary lookup)
2. Metrical valence
3. Strong and weak position
4. Metrical type (binary ~ ternary)
5. Metrical subtype (foot type)
6. Line length
7. Catalexis and hypermetricality
8. (Rhyme)

[All processing is XQuery and XSLT]

## 1. Dictionary lookup

- Input is word in normal orthography
  - Mixed case, punctuation, no stress
  - Dictionary contains word forms with stress and morphological information
    - Morphological information is irrelevant for our purposes
- Eventual output has all vowels tagged
  - Stressed
  - Unstressed
  - Unknown
    - Not in dictionary
    - Dictionary evidence is contradictory

## Dictionary content (eXist)

```
<item>
  <unstressed>капий</unstressed>
  <stressed>к<stress>а</stress>рий</stressed>
  <pos>п</pos>
  <form>
    <categories>
      <category case="N" gender="m" number="sg"/>
      <category case="A" gender="m" animacy="i" number="sg"/>
    </categories>
    <content>к<stress>а</stress>рий</content>
  </form>
  <!-- other <form> elements -->
</item>
```

## Final dictionary output

к  
 <vowel stress="1">а</vowel>  
 п  
 <vowel stress="-1">и</vowel>  
 й

## 2. Metrical valence

- For each vocalic position in the line
  - Ignore text node children of <line>
  - Line as sequence of <vowel> element
- Metrical valence
  - stressed / (stressed + unstressed)
    - Ignore unknowns
  - Varies between 0 and 1
  - Sample output: 0 0.5 0 1 0 0 0 1 0

## 3. Strong and weak position

- Compare valence of each position to preceding and following
  - Assume a 0 value if preceding or following is missing, i.e., at beginning or end of a line
- If target value is higher than both neighbors: strong
- If target value is lower than both neighbors: weak
- Otherwise: weak
  - Provisional; adjacent strong positions do not occur in common Russian meter

## 4. Metrical type (binary ~ ternary)

- Calculate how often the strong ~ weak property of a syllable matches the property two (resp. three) syllables earlier
- Count both strong/strong and weak/weak matches
- The greater number of matches determines the type
- Resolve tagging ambiguities according to positional valence (where possible)
- Retag all vowels with @stress values of 0 or 1
  - Represents strong ~ weak (not necessarily stress)

## 5. Metrical subtype (foot type)

- Having determined metrical type (binary ~ ternary)
- Subtype is based on last foot
  - Last stress is the only obligatory one
  - Iamb ~ trochee
  - Dactyl ~ anapest ~ amphibrach

## 6. Line length (number of feet)

- Number of strong positions = number of feet
- May be global or line-specific

Во всем мне хочется дойти	оx   оx   оx   оx
До самой сути.	оx   оx (o)
В работе, в поисках пути,	оx   оx   оо   оx
В сердечной смуте.	оx   оx (o)

[Pasternak 1956]

## 7 Catalexis and hypermetricality

- Catalexis: Number of syllables
  - Is sufficient for the number of feet
  - Is not sufficient for the number of complete feet
- Hypermetricality
  - Syllables after the final stress are easily identified
  - Hypermetrical caesura: Demarcate feet based on strong position

Как ветер мокрый, ты быешься в ставни,	оx   ox (o)    ox   ox (o)
Как ветер черный, поешь: ты мой!	ox   ox (o)    ox   ox
Я древний хаос, я друг твой давний,	ox   ox (o)    ox   ox (o)
Твой друг единый, отрой, отрой!	ox   ox (o)    ox   ox

[Gippius, Neljubov', 1907]

## Rhyme

- Meter vs rhyme
  - Meter can be identified through vowel *letters*
  - Rhyme requires *sounds*
    - род ~ pot ['rot]
    - ноги [nə'gɪ] ~ ноги ['nogɪ]
  - Rhyme runs from last stressed vowel to end of line
    - Exception: open masculine rhyme: tree ~ see
- Russian orthography can be mapped to coarse phonetics algorithmically, except that
  - Stress must be known
  - e ~ ë must be resolved
  - Both are accessible in dictionary

## Meter ~ rhyme

- Meter: privilege markup
  - Distinguish line/vowel from line/text()
- Rhyme: privilege text
  - Convert mixed content to text
  - String matching
- Imperfect (slant) rhyme: privilege markup
  - Decompose segments into distinctive features
  - Identify rhyme scheme on the basis of exact rhyme
  - Infer imperfect rhyme on the basis of the ambient rhyme (cf. ambient meter)
  - Characterize slant rhyme by neutralized features

## Taking stock

- We can count syllables by counting vowel letters
- If we know the place of stress
  - We get meter
  - We get most pronunciation ...
  - ... and therefore most rhyme
- If we also know e ~ ë
  - We get the rest of pronunciation
  - We also get rhyme

## Overlapping hierarchies: caesura

- Свои | ми кольцами || она, | упорная,
- Iambic tetrameter, dactylic caesura and clausula
- Feet:    o x | o x (o o) || o x | o x (o o)
- Words:  o x o | x o o || o x | o x o o

```
<word>
  СВ
  <vowel stress="0">о</vowel>
  <vowel stress="1">и</vowel>
  М
  <vowel stress="0">и</vowel>
</word>
```

## Implicit and explicit markup

- P5: "A text is not an undifferentiated sequence of words, much less of bytes."
  - Quotation marks delimit quotations
  - Space characters delimit words
  - New line characters delimit lines of poetry
  - Multiple new line characters delimit paragraphs of prose
  - Asterisks or underscores delimit emphasized text,
- Metrical foot and word hierarchies overlap
- Can we use a combination of explicit and implicit markup to represent the *logical* overlap without *syntactic* overlap?
  - Is there pseudo-markup of the metrical hierarchy that we can use?

## Implicit and explicit markup

- Words are tagged explicitly as <word>
- Feet are implicit
  - Represented at level of line/descendant::vowel
  - Ignore everything else

```
let $vowelCount := min(//line/count(descendant::vowel))
let $midPoint := $vowelCount div 2
let $targets := //line/descendant::vowel[position() eq $midPoint]
return
  if ($targets/following-sibling::vowel)
  then
    'no caesura'
  else
    'caesura'
```

## Q: So where are the foot boundaries?

```
<word>сво<stress>и</stress>ми</word>
<word>к<stress>о<stress/>льц<stress/>ам</word>
<word>он<stress>а</stress></word>
<word>ын<stress>о<stress>п<stress>а</word>
```

### • A: *We don't care!*

- If there is a word boundary between feet
  - Word boundary = foot boundary (potential caesura)
- Otherwise
  - We don't need to locate the foot boundary
- Not only do we not need to tag the foot boundaries, but we don't even need to be able to find them

## What if the words weren't tagged?

```
<w>сво<v>о</v><v>и</v>м<v>и</v></w>
<w>к<v>о</v>льц<v>а</v>м<v>и</v></w>
<w><v>о</v>н<v>а</v></w>
<w><v>ы</v>п<v>о</v>п<v>а</v>я</v></w>
```

```
сво<str>и</str>ми
к<str>о</str>льцами
он<str>а</str>
ын<str>о</str>пная
```

## What if the words weren't tagged?

- The data
  - Feet are still implicit
  - Now words are encoded with pseudo-markup white space
- Processing
  - What we want to do is tokenize(line)
    - Oops!
  - Convert markup to text
  - Convert text to markup
    - Replace white space in line/text() with <w/> tags
    - Convert milestones to wrappers (<xsl:for-each-group>)

## What if it isn't just words?

```
<lpa>съ грѣхы. ѿпусти слыш< marginalia>даривъ.
<lb/>ско</ marginalia>рю</lpa>
```

Not generally identified as overlap

- Syntactically it isn't
- Logically it is
- White space as pseudo-markup
- It may raise the same processing challenges as traditional types of overlap

## So what have we learned?

- We can identify an “element” without start or end tags
  - And without knowing where it starts and ends
  - Foot boundaries are not only untagged, but also unknown
- Overlapping hierarchies may hide in plain sight
  - An overlapping hierarchy may be encoded through plain text pseudo-markup (e.g., white space)
  - Avoidance of overlap is only apparent, and may vanish when the implicit hierarchy needs to be processed
- Needing to treat words during processing as though they had been marked up individually is as much overlap as when the markup is overt
  - White space as milestone

## So what have we learned?

- With absent and implicit markup, as with other strategies for avoiding syntactic overlap
  - No well-formedness errors
  - Similar processing challenges as soon as the researcher needs to engage with them explicitly

## Thank you!

Elise Thorsen (enthorsen@gmail.com)  
 David J. Birnbaum (djbpiitt@gmail.com)  
<http://poetry.obdurodon.org>

Assisted by: Sam Depretis, Erin Harrington  
 Thanks to: Elisa Beshero-Bondar, Sibelan Forrester

Birnbaum, David J., and Elise Thorsen. "Markup and meter: Using XML tools to teach a computer to think about versification." Presented at Balisage: The Markup Conference 2015, Washington, DC, August 11–14, 2015. In *Proceedings of Balisage: The Markup Conference 2015*. Balisage Series on markup technologies, vol. 15 (2015). doi:10.4242/BalisageVol15.Birnbaum01. <http://www.balisage.net/Proceedings/vol15/html/Birnbaum01/BalisageVol15-Birnbaum01.html>