

### Finite-State Guesser

- ♦ Xerox's fsm software (Beesley & Karttunen 2003)
- ♦ 92 LEXICON units (continuation classes)
- ♦ regexes for orthographic normalization and phonological rules (18)
- ♦ small lexicon of stems and frequent function words
- ♦ If no stem in the lexicon fits, "guessed stems" inserted, following very basic phonotactic rules
- ♦ guesser vastly overgenerates, e.g.

начЯТИ	+Guess+Part+Prt+Pass+Kf+Masc+Nom+Pl
начЯ	+Guess+Verb+Inf
начЯТ	+Guess+Noun+Fem+Anim+Hum+Loc+Sg
начЯТ	+Guess+Noun+Fem+Anim+Hum+Nom+Du
начЯТ	+Guess+Noun+Fem+Anim+Hum+Dat+Sg

### Pure Guesser

Otecъ	Verb:Sg:Aor:2P Verb:Sg:Aor:3P Noun:Nom:Sg:Masc:Anim:Hum Noun:Nom:Sg:Masc:Anim:NonHum PPoss:Nom:Sg:Masc:IP Verb:Sg:Masc:Prt:PI:NonRefl Punc Con Part:St:Acc:Sg:Masc:Prt:Act Noun:Acc:Sg:Masc:NonAnim Verb:Sg:Pres:3P Noun:Acc:Sg:Fem:Anim:NonHum Noun:Nom:Sg:Fem:NonAnim Noun:Acc:Sg:Fem:NonAnim Noun:Acc:Sg:Fem:Anim:Hum Adj:St:Nom:Sg:Masc Noun:Acc:Sg:Masc:Anim:Hum Noun:Acc:Sg:Masc:NonAnim Noun:Nom:Sg:Masc:NonAnim Noun:Nom:Sg:Masc:Anim:Hum Adj:St:Acc:Sg:Masc Noun:Nom:Sg:Masc:Anim:NonHum Noun:Acc:Sg:Masc:Anim:NonHum Noun:Nom:Sg:Fem:Anim:NonHum Noun:Nom:Sg:Fem:Anim:Hum Name:Loc:Sg:Masc:NonAnim	father  my died  and S. sits
moj		
umeritъ		
a		
Svjatopolkъ		
séditъ		
Kyevé		(in) Kiev

### Guesser with lexicon (best case)

Otecъ	Verb:Sg:Aor:2P Verb:Sg:Aor:3P Noun:Nom:Sg:Masc:Anim:Hum Noun:Nom:Sg:Masc:Anim:NonHum PPoss:Nom:Sg:Masc:IP Verb:Sg:Masc:Prt:PI:NonRefl Punc Con Part:St:Acc:Sg:Masc:Prt:Act Noun:Acc:Sg:Masc:NonAnim Noun:Nom:Sg:Masc:NonAnim Verb:Sg:Pres:3P Noun:Nom:Sg:Fem:NonAnim Noun:Acc:Sg:Fem:NonAnim Noun:Acc:Sg:Fem:Anim:Hum Adj:St:Nom:Sg:Masc Noun:Acc:Sg:Masc:Anim:Hum Noun:Acc:Sg:Masc:NonAnim Noun:Nom:Sg:Masc:NonAnim Noun:Nom:Sg:Masc:Anim:Hum Adj:St:Acc:Sg:Masc Noun:Nom:Sg:Masc:Anim:NonHum Noun:Acc:Sg:Masc:Anim:NonHum Noun:Nom:Sg:Fem:Anim:NonHum Noun:Nom:Sg:Fem:Anim:Hum Name:Loc:Sg:Masc:NonAnim	father
moj		
umeritъ		
a		
Svjatopolkъ		
séditъ		
Kyevé		

### Word Alignment (NRu—ORu normalized)

- ♦ statistical word alignment
- ♦ general-purpose tools ignore the close relation between the aligned languages
- ♦ Uplug (GIZA++) tested
- ♦ potentially useful for aligning the orthographically non-normalized original and the normalized version — sentence boundaries?
- ♦ alternative approach: edit distance
- ♦ for each ORu word, look for a perfect match in the NRu parallel segment
- ♦ where there is no perfect match, select the modern word in the same segment which is minimally orthographically distinct from the ORu form
- ♦ Levenshtein edit distance, adapted by
- ♦ disregarding predictable orthographic variation
- ♦ relativized to word length
- ♦ Introducing an upper threshold value for possible cognates

### Simple edit distance with obligatory matches

1	8	начати	+Guess+Part+Pr+Pass+Kf+Masc+Nom+Pl ...
1	9	начать	+Guess+Part+Pr+Pass+Kf+Masc+Nom+Pl ...
1	10	старыми	+Guess+Adj+Lf+Neut+Ins+Pl ...
1	11	словеса	+Noun+Neut+Nom+Anim+Ins+Pl ...
1	12	трудныхъ	+Guess+Adj+Lf+Neut+Loc+Pl ...
1	13	повѣстїи	+? ...
1	14	о	+Guess+Verb+Aor+2P+Sg ...
1	15	пылку	+Guess+Noun+Fem+Anim+Hum+Gen+Du ...
1	16	Игоревъ	+Guess+Verb+Pres+1P+Du ...
1	17	Игоря	+? ...
1	18	Святъспавлича	+Guess+Verb+Aor+3P+Sg ...
1	19	?	+? ...

старыми  
словеса  
братья  
повести  
о  
ли  
Игоревом  
Игоря  
Святославица  
?2

### Improved edit distance

- "Yers" (ь, ъ) are disregarded
- **ѣ** / **и** and **ѣ** / **е** alternations do not count
- threshold relative to word length

Отець	Отец	father
мој	мой	my
умерѣ	umer	died
а	а	and
Свјатополкъ	Svjatopolk	S
сѣдѣ	sedit	sit
Кыевѣ	Kieve	(in) Kiev
избиваја		beating-up
brateju	bratej	brothers
svoju	svoich	his

### Tag projection (1)

- analyze w's and v's tags and transform them into a common tagset (properties of a tag object)
- simple transfer of NR tag where *no* OR tag could be guessed
- common case: use *part* of the NRu tag to select from guessed OR tags
- unification of part-of-speech and grammatical categories
- beware of ...
  - gender (syncretisms in declensions not yet related to animacy)
  - number (OR had a dual)
  - tense (OR had many past tenses - aorist, imperfect, pluperfect ...)
- tag projection removes implausible tags from the guessed ORu tags
- (guessed tags are preserved independently)

### Tag projection (2)

- The NRu corpus is not tagged perfectly -- potential errors due to NRu syncretisms should be avoided.
- *Nominal gender and animacy* may be projected from NRu to ORu.
- *Number* may be projected with some proviso, i.e., Sg (NRu) => Sg (ORu), but Pl (NRu) => Pl or Du (ORu)
- *Unambiguous noun tags* in NRu show no systematic errors of the tagger; thus they may transfer case and number (Sg vs. non-Sg).
- **Currently**,
  - full projection into ORu for identical word forms without a ORu tag
  - some variability for names vs. nouns
  - rules hard-wired into program code, no separate rule language

### Output with projection rules

Otecъ	Noun:Nom:Sg:Masc:Anim:Hum
moj	PPos:Nom:Sg:Masc:IP
umeriъ	Verb:Sg:Masc:Prt:Pf:NonRefl
,	Punc
a	Con
Svjatopolkъ	Noun:Acc:Sg:Masc:NonAnim
sěditъ	Noun:Nom:Sg:Masc:NonAnim
Kyevě	Verb:Sg:Pres:3P
,	Name:Loc:Sg:Masc:NonAnim
izbivaja	Punc
brateju	Noun:Ins:Sg:Neut:NonAnim
svoju	PPos:Gen:Pl:X
».	Punc

### Many problems for pure guessing + projection ...

I	Con	
sъbra	NumCard:Acc:X	[should be Verb in Aorist]
jaroslavъ	Name:Nom:Sg:Masc:Anim:Hum	
varjagъ	Name:Nom:Sg:Masc:Anim:NonHum	
tyjsačju	Noun:Gen:Pl:Masc:Anim:NonHum	
,	Noun:Gen:Pl:Masc:Anim:Hum	
i	Noun:Acc:Sg:Neut:NonAnim	
,	Punc	
poide	Con	
na	Verb:Sg:Masc:Prt:Pf:NonRefl	[should be Verb in Aorist]
Svjatopolka	Prp:Acc	
,	Noun:Nom:Sg:Neut:NonAnim	
Ne	Picl	
ja	PPER:Nom:Sg:IP	
počachъ	PPos:Loc:Pl:X:IP	[should be Verb in Aorist]
izbivati	Verb:Inf	Verb:Sg:Aor:3P
bratsju	Noun:Gen:Pl:Masc:Anim	

### Syncretism in NRu noun paradigms

SG	NOM	Masc:Anim	Masc:NonAnim	i-Class	Fem:Anim	Fem:NonAnim	Neut
	GEN	volk	leš	put'	žena	zemlja	okno
	DAT	volka	leša		ženj	zemlji	okna
	ACC	volku	lešu		žene	zemle	oknu
	LOC	volka	leš	put'	ženu	zemlju	okno
	NOM	volke	lešu (-e)		žene	zemle	okne
PL	GEN	volki	leš		ženj	zemlj	okna
	DAT	volkam	lešov	putěj	žen	zemel'	okon
	ACC	volčim	lešj	putj	žen	zemljam	okn

- ♦ Anim => Pl:Acc = Pl:Gen; NonAnim => Pl:Acc = Pl:Nom
- ♦ Masc:Anim => Sg:Acc = Sg:Gen; Masc:NonAnim => Sg:Acc = Sg:Nom
- ♦ i-Class => Sg:Gen = Sg:Dat = Sg:Loc = Pl:Nom
- ♦ Fem => Sg:Dat = Sg:Loc; Fem => Sg:Gen = Pl:Nom

### Syncretism in ORu nominal paradigms

SG	NOM	vъlkъ	konь	synъ	gostъ	kamy	žena	zemlja
GEN	vъlka	konja	synovi	gosti	kamene	ženě	ženě	zemli
DAT	vъlky	konj	synovj	gostj	kamene	ženě	ženě	zemli
ACC	vъlkъ	konь	synъ	gostъ	kamene	ženě	ženě	zemli
LOC	vъlcě	konj			Ø	ženo	ženě	zemle
VOC	vъlče	konju				ženo	ženě	zemli
DU	NOM=ACC	vъlka	konja	synj	gostj			
	GEN=LOC	vъlčj	konj	synovj	gostj	kamenu		
PL	NOM	vъlci	konj	synove	gostje	kamene	ženj	zemlj
	GEN	vъlkъ	konь	synovъ	gostj	kamene	ženъ	zemlъ
	ACC	vъlky	koně	synj	gostj			
	INS	vъlky	konj	synjmi	gostjami		ženami	zemljami

### Unambiguous items in NRu

- Masc:Sg:Loc (excluding the "second locative");  
Masc:Anim:Sg:Nom, Masc:Anim:Sg:Dat, Masc:Anim:Pl:Nom;  
Masc:NonAnim:Sg:Gen, Masc:NonAnim:Pl:Gen
- Fem:Sg:Acc; Fem:NonAnim:Pl:Gen
- Neut:Sg:Dat, Neut:Sg:Loc, Neut:Pl:Gen
- Masc:Sg:Ins, Neut:Sg:Ins
- Pl:Dat, Pl:Ins, Pl:Loc
- respective rule for each of the above tags:
  - if ORu guessed set contains a tag compatible with this tag, remove all tags containing a different case or number
  - (irrespective of ORu syncretisms)

21

### Informative syncretisms in NRu

- partial overlap in potential tags, e.g.
  - ORu: {{Sg:Gen}}... & NRu: {{Sg:Gen}}{{Anim:Sg:Acc}}  
=> ORu: {Sg:Gen}
  - ORu: {{Pl:Gen}}... & NRu: {{Pl:Gen}}{{Anim:Pl:Acc}}  
=> ORu: {Pl:Gen}
- NB: *We do not find Gen/Acc syncretisms in the Sg or Pl in early ORu*
- ORu: {{Sg:Nom}}{Sg:Acc}}{Pl:Gen}}...  
& NRu: {{Acc}}{Nom}} => ORu: {{Sg:Nom}}{Sg:Acc}}
- NB: *The NRu tag is never more than two-way ambiguous.*

22

### Perspectives

- **Implemented**
  - sentence alignment Nru/ORu-normalized; word alignment (cf. references)
  - guesser and small ORu lexicon
  - edit distance
  - tag projection rules
- **Fine-tuning**
  - word alignment: adding safe bets; hand-correction?
- **Further gains**
  - Bootstrapping an enlarged lexicon via word alignment
  - Extracting orthographic variants

23

### References

- Gale, W. & Church, K. 1994. "A Program for Aligning Sentences in Bilingual Corpora", in: Armstrong, S.(ed.). *Using Large Corpora*. MIT Press, 91-102.
- Meyer, R. 2009. "Old wine in new wineskins? — Tagging Old Russian via projection from modern translations", submitted to *Russian Linguistics*.
- Tiedemann, J. 2006. "ISA & ICA – Two Web Interfaces for Interactive Alignment of Bitexts", in: *Proceedings of LREC 2006*. Genova.

Thank you!

roland.meyer@sprachlit.uni-regensburg.de  
<http://www-slavistik.uni-r.de/subjekte>

24