



Computer Processing of Written Cultural Heritage Sources

*Jiří Bubník, Jiří Čelák, Vojtěch Janota,
Alexandr Kára, Václav Novák,
Kiril Ribarov, Tomáš Vondra*



We will focus on:

Lexical and corpus processing of written cultural sources

We place the written cultural sources in an electronic contextual (e-context) field with the following connecting elements:

- source image along with language based contextual structure of the word mass present in the sources
- connections (inner and outer links) among various types of written cultural sources within wider cultural environment.

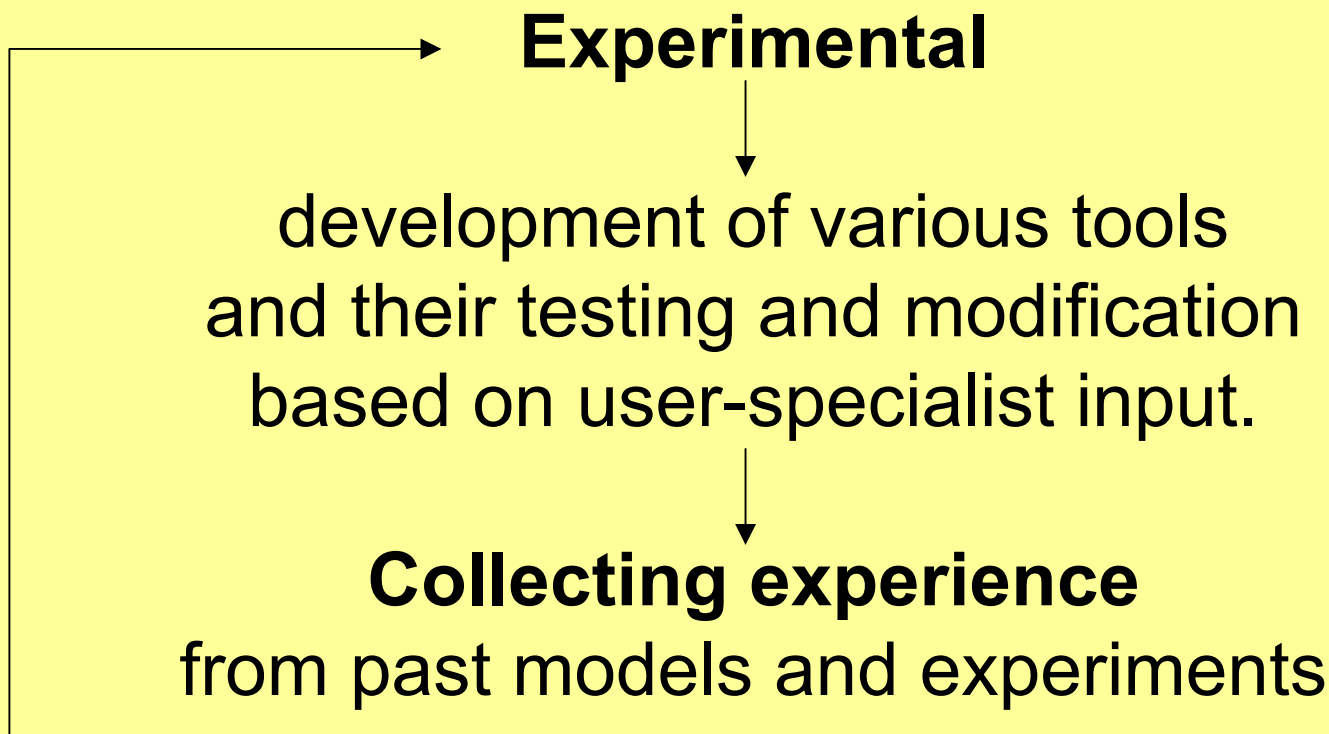


We aim at:

Development of technologies and tools
for
large-scale activities
aimed towards
multi-aspectual presentation
of
written cultural heritage
in a
highly distributed manner



Our approach is:



*Developed as individual work and recently also as student project at the
Center for Computational Linguistics, Charles University.*



Corpus Processing

- Developed for contemporary languages
- Existence of various corpus managers

but

None of them covers specialties present in processing of old sources (dead languages).

**We need special approaches,
but whenever possible,
we prefer the use existent modules.**

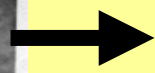
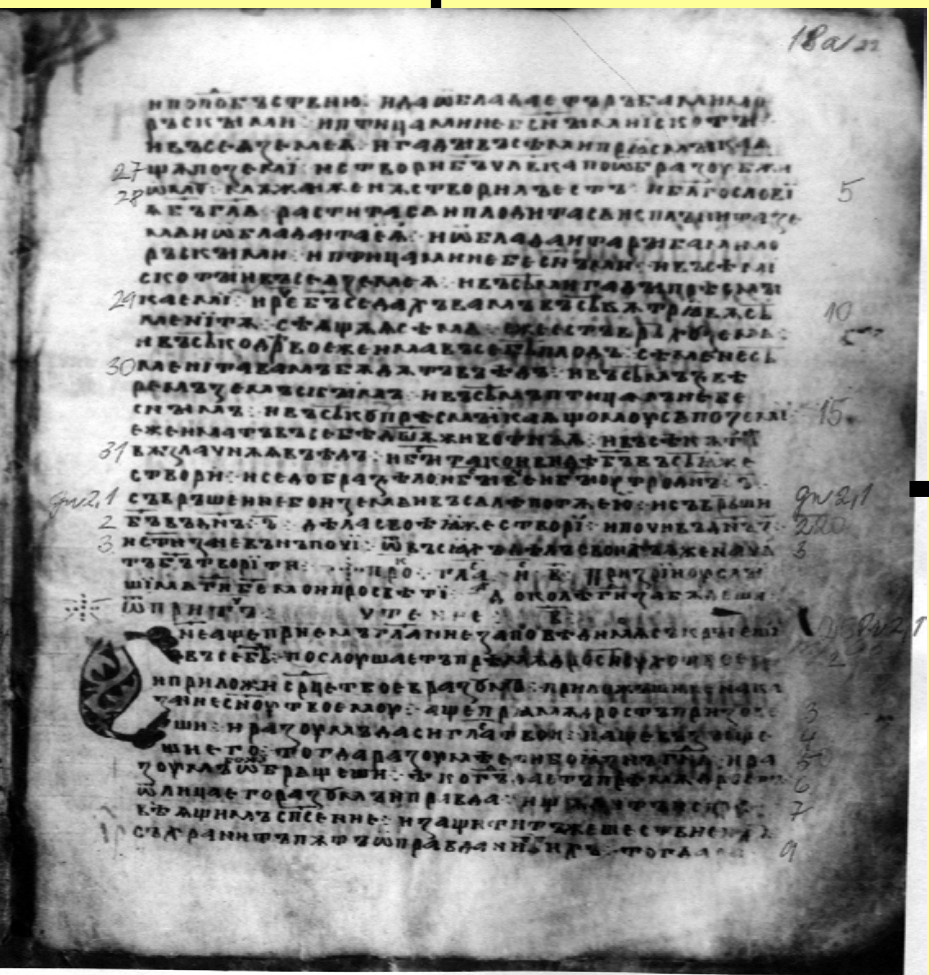


ACT

The input

Gn 1, 26-32; Gn 2, 1-3; Pr 2, 1-9

18



ι πο π6β7σπαι ♥;;ι δα ρβλοδοετ ρ-βαμι μο !
 ρ7σκ-μι;;ι πιχιαμι νεβ9ν-μι Ψσκοτ-;;
 ι τωσε3 ζεμε3;;ι γοδ-τ6τμι πρλ-αμ-κοβ
]4 πο ζεμΨ;≡ει σπαι ρ7 =λπκα πο ρβροζου βαι
 Ωμ8;;μ4εαι εεν4 σπαι ρ7 εοτ7;;≡* ι βλγσολσπΨ %
 3 β7 γλ2;;ροσπιτα εε ι πλοδιτα εε ι σπλ7νιτα ζε
 μ2 ι ρβλοδιτα ε3;;ι ρβλοδιτα ρ-βαμι μο
 ρ7σκ-μι;;ι πιχιαμι νεβεσν7μι;;ι τ6τμΨ
 σκοτ@ ι τωσε3 ζεμε3;;ι τ6τμι γοδ-πρδ-αμ-
 κοεμΨ;≡(ι ρββ7 σε δαη7 σπαι τ6τκ4 τρλ-α τ6
 μενΨ4;;;ε3] 43 τ6μ2;;εε εοτ7 σπ7η8 ζεμ2;;
 ι τ6τκο δρδσ εε ιμα τ6 σεβδ πλδ7;;;τ6μνε τ6
 μενΨα σπαι β4δ4τ7 τ6 5δ7;;#) ι τ6τμ7 / τ6
 ρεμ7 ζεμ7σκ@μ7;;ι τ6τμ7 πιχιαμ7 νεβε
 σν-μ7;;ι τ6τκ8 πρδ-αμ@κοβ] ομου εε πο ζεμΨ;
 εε ιμα τ7 τ6 σεβδ §[4 εαι σπαι ν43;;ι τ6τκ4 τρδ
 α ζλα=ν43 τ6 5δ7;;ι βα-τακο #! ι παδ5 β7 τ6τδ ↓ε
 σπαι;;ι σε δοβρα / 5λα;;#≡ ι βα-αβι βα-αυτρο δν7;;Ζ;;
 ! σπαι τ7 ει νεβι ι ζεμ2 ι τ6τδ λ5ποτ4 ε' ;;≡ ι σπαι τ7 ι
 β7 τ6 δν7;;Ζ;;δ5λα σπαι ↓ε σπαι Ψ;ι πο=ι τ6 δν7 ζ;;
 # ι σπαι ζανε τ6 ν7 πο=Ψ;;. τ6τδ η7 δ5λ7 σπαι η7 3ε να=2
 τ7 β7 τπαι Ψι;;;πΨ0;;;γΔα;;;Δ;;;π;;;πρι ζρΨι αυθ-
 [Ψμ2 Ψι βε μοι προσπαι Ψ;| Δοκολ5 Ψι ζαβ4δε[ι;;
 ;; . πρι @7;;=τ ενι ε;;;π;;;
Cνε α] ε πρι εμ7 γλανε ζσποτδ ι μ3 σπαι κρ-ε[Ψ ≡%
 τ6 σεβδ;;≡ ποσπαι[αετ7 πρδ μ4δρσπαι αυηρ τπαι;;
 ι πρι λωι σπαι ε τπαι σπαι ραζ8μ7;;πρι λω7[ι εε ναικα
 ζανι ε θναι τπαι εμου;;# α] ε πρι μ4δρσπαι πρι ζσπαι
 [ι;;ι ραζουμ7 δαα γλα τπαι;;;ε ι α] ε τ6τ-] ε
 [ι εγα;;%τογδα ραζουμ5ε[ι βολ ζν7 γΨ2;;ι ρα #)
 ζουμ7 ρβρ2] ε[ι;;;L5κο γ7 δαετ7 πρδ μ4δρσπαι;;
 . λιχα εγο ραζ8μ7 ι προσπαι;;&ι] 4δ ι τ7 ι σπαι



ACT

Gn 1, 26-32; Gn 2, 1-3; Pr 2, 1-9

Original forms:

ι πο π̄β7σπαι ♥;;ι δα Ωβλοδοετ7 ρ-βαμι μο
 ρ7σκ-μι;;ι πιτιχομι νεβ̄ον-μι Ψακοτ-;;
 ι τ̄σε3 ζεμε3;;ι γαδ-τ̄σμι πρ̄αμ-κοβ
] 4 πο ζεμΨ;≡&ι σπ̄ρι β7 =λ̄κα πο Ωβροζου β̄ι
 Ωμ8;;μ4εα ι εεν4 σπ̄ρι λ7 εστ7;;≡* ι β̄ιγολοσΨ
 3 β7 γ̄λ2;;ροσπιτα ε2 ι πλοδιτα ε2 ι σπ̄λ7νιτα ζε
 μ2 ι Ωβλοδαιτα ε3;;ι Ωβλοδαιτα ρ-βαμι μο
 ρ7σκ-μι;;ι πιτιχομι νεβ̄ον7μι;;ι τ̄σμΨ
 σκοτ© ι τ̄σε3 ζεμε3;;ι τ̄σμι γαδ-πρ̄αμ-
 κοεμΨ;≡(ι ρ̄β7 σε δαη7 π̄μ7 τ̄σκ4 τρ̄ τ̄ σ̄
 μενΨ4;;σ3] 43 τ̄μ2;;εε εστ7 π̄7η8 ζεμ2;;
 ι τ̄σκο δρ̄τ̄ εε ιμα τ̄ σεβ5 πλοδ7;;σ̄μενε σ̄
 μενΨα π̄μ7 β̄δ4τ7 τ̄ 5δ7;;#) ι τ̄σμ7 / τ̄
 ρεμ7 ζεμ7σκ©μ7;;ι τ̄σμ7 πιτιχομ7 νεβ̄ε
 σν-μ7;;ι τ̄σκ8 πρ̄αμ©κοβ] ομου ε2 πο ζεμΨ;
 εε ιμα7 τ̄ σεβ5 δ̄[4 ειπ̄τν43;;ι τ̄σκ4 τρ̄
 τ̄ ζλα=ν43 τ̄ 5δ7;;ι β̄α-τακο #! ι π̄δ5 β7 τ̄σ̄ ↓εε
 σπ̄ρι;;ι σε δοβρα / 5λο;;#≡ ι β̄α-σ̄ι β̄α-ουτρο δ̄ν7;;Ζ;;
 ! σ̄π̄7[ει νεβο ι ζεμ2 ι τ̄σ̄ λ5ποτ4 ε' ;;≡ ι σ̄π̄7[ι
 β7 τ̄ δ̄ν7;;Ζ;;δ5λα σ̄σ̄ ↓εε σπ̄ρΨ;ι πο=ι τ̄ δ̄ν7 ζ;;
 # ι φ̄ι ζανε τ̄ ν7 πο=Ψ;.: τ̄σ̄ η7 δ5λ7 σ̄σ̄ι η7 3εε να=2
 τ7 β7 τ̄σπ̄Ψι;;;;π̄βο;;γ̄Δα;;Δ;;π̄;π̄ρι ζρ̄Ψι ουσ̄-
 [Ψμ2 φ̄ι βε μοι προσ̄τ̄Ψ;| Δοκολ5 φ̄ι ζαβ̄δε[ι;;
 ;; ∴ πρ̄ι ⊗7;;=τ̄ενι ε;;τ̄;;

! ,!*,
 % ΕGn, ≡Δ ι πο π̄β7σπαι ♥;;ι δα Ωβλοδοετ7 ρ-βαμι μο?
 ρ7σκ-μι;;ι πιτιχομι νεβ̄ον-μι Ψακοτ-;;
 ι τ̄σε3 ζεμε3;;ι γαδ-τ̄σμι πρ̄αμ-κοβ?
] 4 πο ζεμΨ;≡&ι σπ̄ρι β7 =λ̄κα πο Ωβροζου β̄ι
 1) Ωμ8;;μ4εα ι εεν4 σπ̄ρι λ7 εστ7;;≡* ι β̄ιγολοσΨ
 3 β7 γ̄λ2;;ροσπιτα ε2 ι πλοδιτα ε2 ι σπ̄λ7νιτα ζε?
 μ2 ι Ωβλοδαιτα ε3;;ι Ωβλοδαιτα ρ-βαμι μο?
 ρ7σκ-μι;;ι πιτιχομι νεβ̄ον7μι;;ι τ̄σμΨ
 !% σκοτ© ι τ̄σε3 ζεμε3;;ι τ̄σμι γαδ-πρ̄αμ-?
 κοεμΨ;≡(ι ρ̄β7 σε δαη7 π̄μ7 τ̄σκ4 τρ̄ τ̄ σ̄
 τ̄

.....

νε α] ε πριεμ7 γ̄λ̄ανιε ζαποσ̄δι μ3 σ̄κρ-ε[Ψ ≡%
 τ̄ σεβ5;;≡ ποσ̄λου[εστ7 πρ̄μ4δροσπι ουηρ τ̄σπε;;
 ι πριλοσι φ̄χε τ̄σπε π̄ραζ̄μ7;;π̄ριλοσ7[ι εε νακα
 ζανιε σ̄γου τ̄σπεμου;;# α] ε πρ̄μ4δροσ7 πριζσπε
 [ι;;ι ραζ̄ουμ7 δαα γ̄Δα τ̄σπαι;;∃ ι α] ε τ̄ζ-] ε
 [ι εγα;;%τογδα ραζ̄ουμ5ε[ι βολ̄ ζν7 γ̄φ2;;ι ρα
 ζουμ7 ρ̄β̄2] ε[ι;;;L5κο φ̄7 δαετ7 πρ̄μ4δροσ7;;
 ∴ λιχα εγο ραζ̄μ7 ι προσ̄τ̄α;;&ι] 4διτ7 ι σπ̄ρα

#) INFORUM 2004



,!*

ΕΓη, ≅1 ι πο π⁶β7σπ α ♥;;; ι δα Ωβλοδοετ7 ρ-βομι μο?
 ρ7ακ-μι;;; ι πιχομι νεβ⁶ον-μι Ψακοτ-;;;
 ι α7σε3 ζεμεε3;;; ι γοδ-α7ομι πρ⁶αμ-κοβ?
] 4 πο ζεμΨ; ≅& ι σπωρι β7 = λ⁶κα πο Ωβροζ ου β⁶ι
 Ωμδ;;; μ4 ε α ι εεν4 σπωρι λ7 εστ7;;; * ι β⁶λοσλοσΨ
 3 β7 γλ2;;; ρασιτα α ι πλοδιτα α ι σπλ7νιτα ζε?
 μ2 ι Ωβλοδοατα ε3;;; ι Ωβλοδοατα ρ-βομι μο?
 ρ7ακ-μι;;; ι πιχομι νεβ⁶ον7μι;;; ι α7ομΨ
 ακοτ© ι α7σε3 ζεμεε3;;; ι α7ομι γοδ-πρ⁶αμ-?
 κοεμΨ; ≅ ι ρ⁶β7 σε δαη7 πωμ7 α7οκ4 τρ⁶α α7 ο

18r1,3:

πο<δο>β7σπ α ι φο

18r4,5: β<ογ>7

18r5,7:

βλ<α>γοσλοπι?3

Original forms:



18r10,2: οε<=ε>
Rendered forms:



Multiple meanings, Multiple Understanding, Multiple.....

ι||εγοαε||πιδι[ι||πλοδα||σε||σ7τπορι||ω7||μν5
(and the fruit you see created in me),

2 variants!

the string σ7τπορι||ω7||μν5

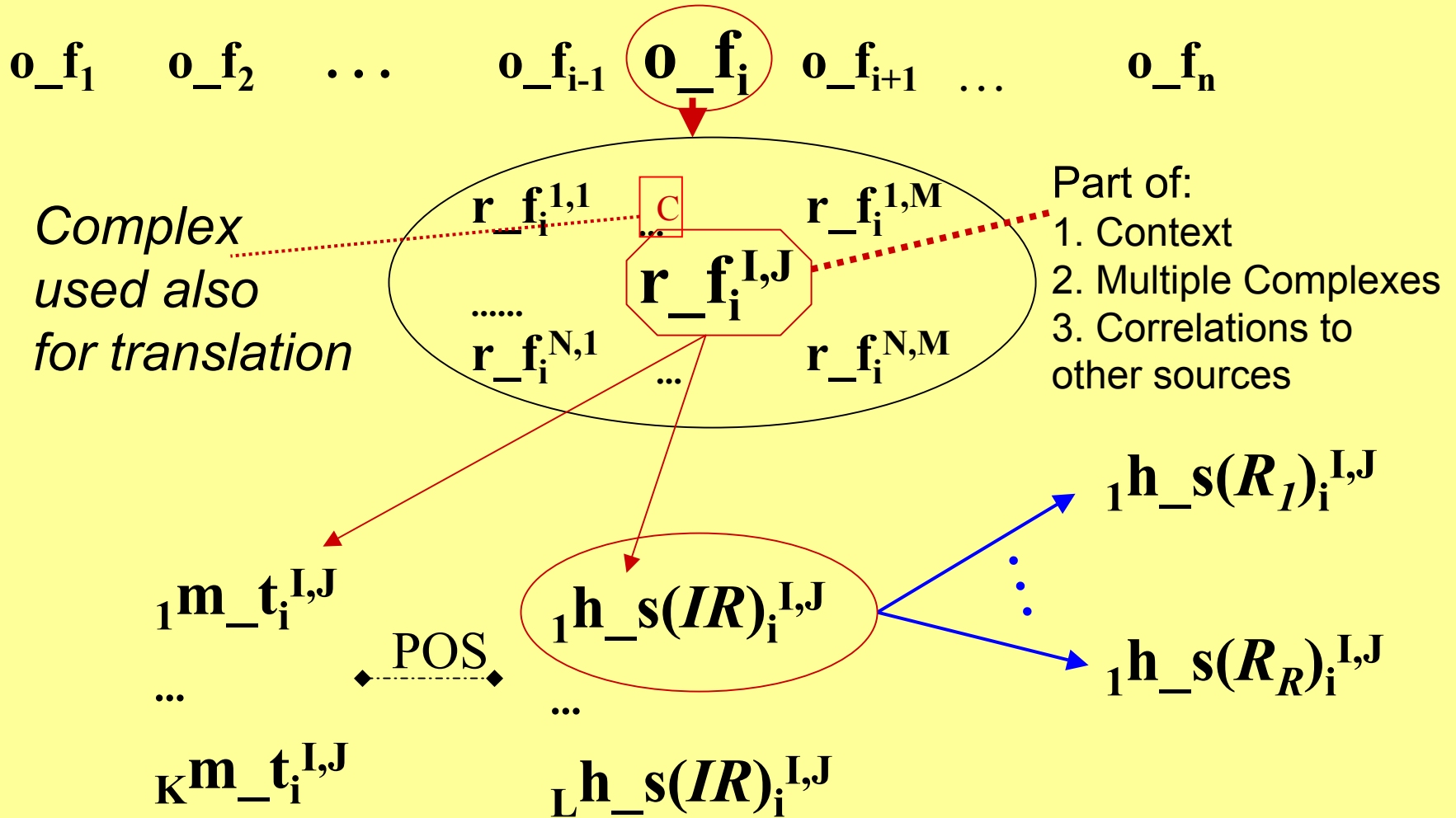
could also be divided as σ7τποριω7||μν5

(where σ7τποριω7 is the past participle - active mood of *create*)

both are grammatically correct

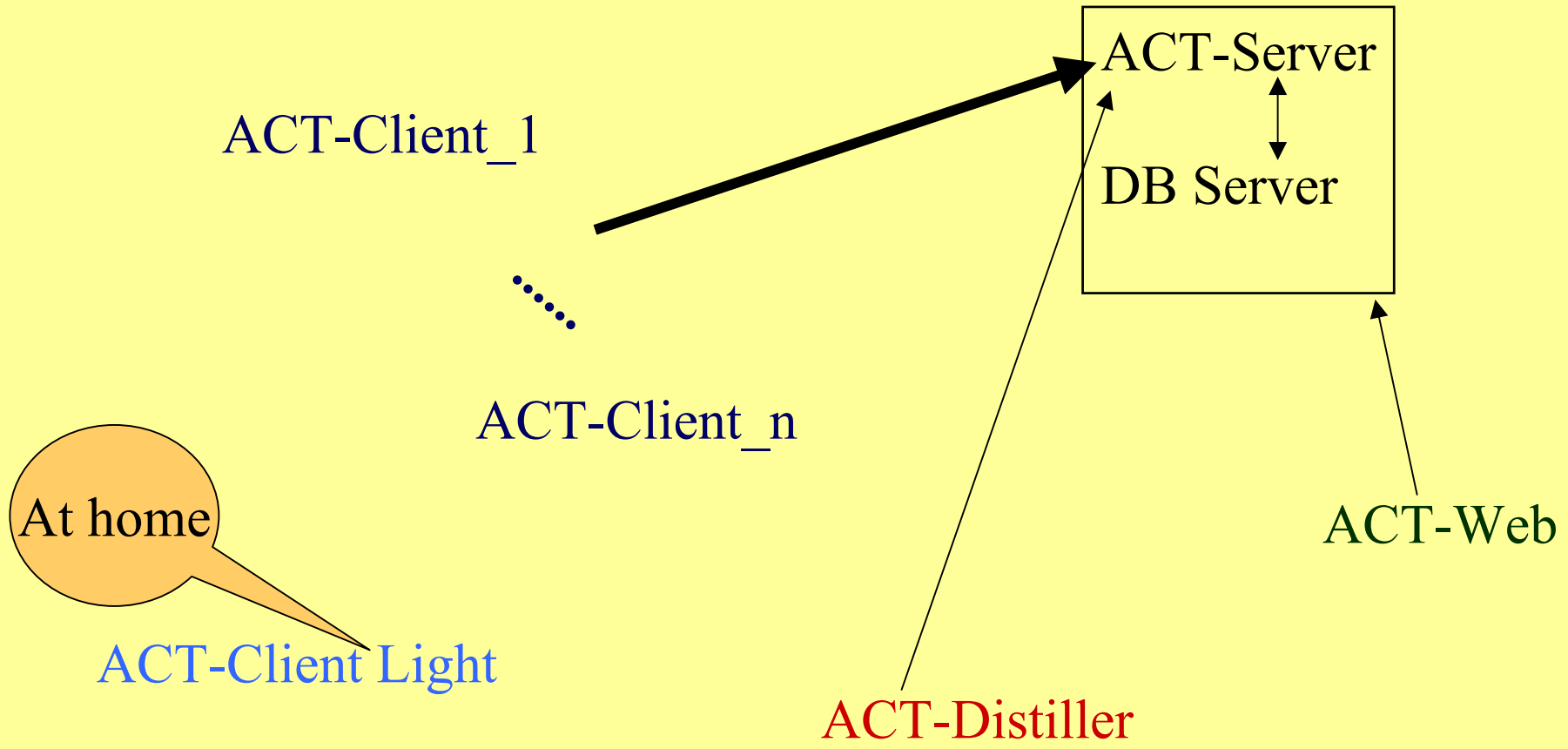


...multiple...multiple...multiple...





ACT Structure



ДА БЖДЕТЪ НА ТЕРЪДИ ПО СРЪАДЪ КОДЪИ И ДА БЖ/ДЕТЪ
 ДА [БЖДЕТЪ [НА [ТЕРЪДИ [ПО [СРЪАДЪ [КОДЪИ [И [ДА [БЖ/ДЕТЪ []
 РАЗЛЖУЪАЩИ МЕЖДОУ КОДОА С БО(ДО)/А
 РАЗЛЖУЪАЩИ [МЕЖДОУ [КОДОА [С [БО(ДО)/А []
 И БЪИСЪТЪ ТАКО И СТВОРИ БОГЪЪ ТЕРЪДЪ И РАЗ(ЛЖ)/УИ
 И [БЪИСЪТЪ [ТАКО [И [СТВОРИ [БОГЪЪ [ТЕРЪДЪ [И [РАЗ(ЛЖ)/УИ []
 БОГЪЪ МЕЖДОУ КОДОА ЪКОЖЕ БЪ НАДЪ (ТЕРЪ)/ДИА
 БОГЪЪ [МЕЖДОУ [КОДОА [ЪКОЖЕ [БЪ [НАДЪ [(ТЕРЪ)/ДИА []
 МЕЖДОУ КОДОА ЪЖЕ БЪ ПОДЪ ТЕРЪ/ДИА
 МЕЖДОУ [КОДОА [ЪЖЕ [БЪ [ПОДЪ [ТЕРЪ/ДИА []
 И БЪИСЪТЪ ТАКО НАРЕУЪЕ БОГЪЪ ТЕРЪДЪ НЕБО
 И [БЪИСЪТЪ [ТАКО [НАРЕУЪЕ [БОГЪЪ [ТЕРЪДЪ [НЕБО []
 И ВИДЪ БОГЪЪ ЪКО ДОБРО И БЪИСЪТЪ ВЕУЪЕРЪ И БЪИСЪТЪ ОУТРО
 И [ВИДЪ [БОГЪЪ [ЪКО [ДОБРО [И [БЪИСЪТЪ [ВЕУЪЕРЪ [И [БЪИСЪТЪ [ОУТРО []

-Rendered form Rendered form full-
 НА НА Previous Next Next ambiguous

All variants

Keyword	Tag	Add
		Delete

Current morphological variant

Part of speech: - Keyword: - Edit assistant settings

Morphological tag



Inputs and Outputs

- **Input:** RTF, XML, TXT
- **Output** (for any subset of documents) HTML, XML, RTF, TXT:
 - various search outputs
 - complete index (index verborum)
 - concordance index
 - word frequencies
 - bigrams
 - retrograde indexes
 - etc...



Documents in ACT

Catalogue: language, sorting, morphology

Text: *Document*₁ *Document*_n

Pictures: *Document*₁ *Document*_n

Pictures: *Card File*₁ *Card File*_{15 000 000}

Location+ Oform --> guess context



ACT

A single card

ГЪВНЯТИ *ipf.*

Ap Slǽpč. 2Cor 4, 3.
Fol. 44v. 2.

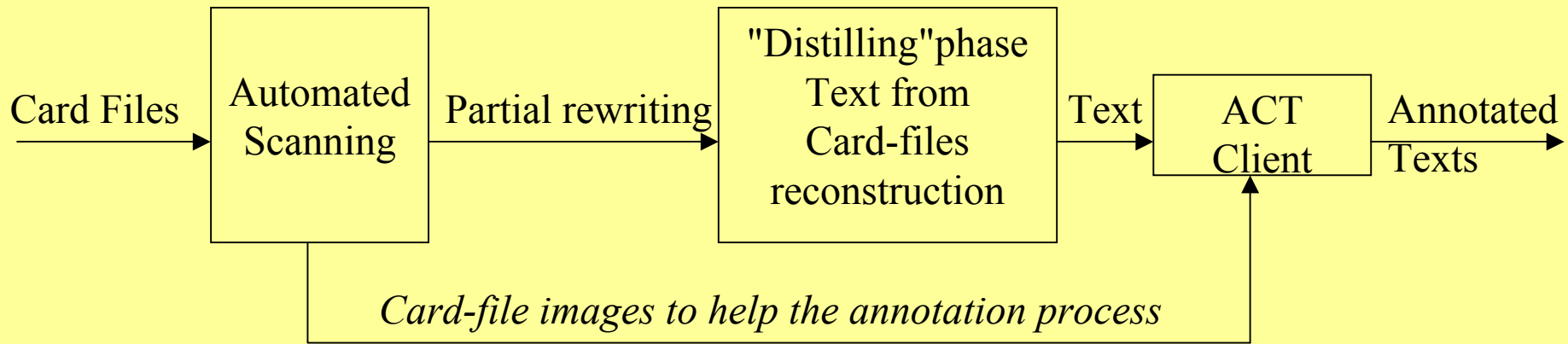
ἀπόλλυθαι, perire, giti & žihubě,

ГЪВНЯЩИХЪ *loc. pl. m. ptc. pass. act.*

ἄψΕ ΑἸ ΖΕ ΕΣΤὺ ΠΟΚΡΩΒΕΝΟ ΕΒΓΛΙΕ ΝΑΩΕ:
Βὺ ΓΒΝΙΑΧΙΧ ΕΣΤὺ ΠΟΚΡΩΒΕΝΟ:

εὺ δὲ καὶ ἔστιν κεκαλυμμένον τὸ εὐαγγέλιον
ἡμῶν, ἐν τοῖς ἀπολλυμένοις ἔστιν κεκαλυμ-
μένον,

= Ḥiḥ. (гыбноуцихъ)





Search criteria

Rendered form (like)

Use equivalences ([show all](#))

Lemma

Part of speech

Index

Morphology (_ is a wildcard)
[morphology help](#)

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15

Order by

	order	priority
Rendered form	ascending	first
First occurrence	no	none
Lemma	no	none
Index	no	none
Part of speech	no	none
Morphology	no	none

Reset Search

Search is restricted to the following documents: Tr, Orb, VK, Mkd, Lesn, Grig, GreekTrans, Hlud, BT, ST, KR, Bon

page: 1 | 2 3 4 5 6 7 8 9 10 11 | [Next page](#) go to page: items per page: font: size:

Rendered form	First occurrence	Lemma	Index	Part of speech	Morphology	#	Detail	Context
БАВНЛА	ST (1577)	БАВНЛА	o	S	S-m-----	1	Detail	Context
БАВНЛО	ST (1652)	БАВНЛА	o	S	S-m-----	1	Detail	Context
БАВНЛЫ	ST (1566)	БАВНЛА	o	S	S-m-----	3	Detail	Context
БАВИЛОНЪ	ST (27992)	БАВИЛОНЪ	o	S	S-m-----	3	Detail	Context
(БАВИЛОНА	ST (30806)	БАВИЛОНЪ	o	S	S-m-----	1	Detail	Context
БАВИЛОНОР	ST (30791)	БАВИЛОНЪ	o	S	S-m-----	1	Detail	Context
БАВИ/ЛОНѢ	ST (30781)	БАВИЛОНЪ	o	S	S-m-----	1	Detail	Context
БАВИ/ЛА	Orbelsky triod (8486)	БАВУЛА		S	S-m-----	1	Detail	Context
БАВИЛОНѢ	BT (9138)	БАВУЛОНЪ		S	S-m-----	4	Detail	Context



Search is restricted to the following documents: Tr, Orb, VK, Mkd, Lesn, Grig, GreekTrans, Hlud, BT, ST, KR, Bon

page: 1 | 2 | [Next page](#)

go to page: items per page: font: size:

left part	form	right part
<p>⟨-а-⟩ ВЪЗВЕЛИЧИМЪ ВСИ УЛОВѢЖ⟨⟩ЛЮБИЕ ТИ ХРИСТЕ БОЖЕ НА ШЬ СЛАВА ТВОИХЪ</p>	рабь	<p>И ВѢНЕЦЪ ВѢРНЫМЪ ВЪЗВЕЛИЧАВЫ ЖТРОБЖ РОЖДЪШЖ ТА РА ШТ⟨ВРЬСТЬ БЫС⟨ТЬ⟩</p>
<p>ВЪДѢЖШЕ УНСТОЖ М⟨⟩Л⟨⟩ДѢТВИ СЛАЩЕ ВПИЕМЪ ХРИСТЪ ПРИ/МИ М⟨⟩Л⟨⟩БЫ СП⟨⟩С⟨⟩С⟨⟩Е ТВОИХЪ</p>	рабь	<p>И М⟨⟩Л⟨⟩ЕНИИА АЗЪ ТЕБѢ СА М⟨⟩Л⟨⟩А М⟨⟩ТИ ХРИСТОУВАГА НЕ ПРѢСТАИ</p>
<p>⟨П⟩ОНТЕ И ПРѢВЪЗНОСИТЕ НЖЕ ПОДВИГОМЪ ДОБРЫМЪ ПОДВИГЪШЕ СА СВ⟨⟩АТИ МЖУ⟨⟩⟨Е⟩НИ⟨⟩ЦИ</p>	рабь	<p>НЕДОСТОИИЪ КО ХРИСТЪУ ХОДАТАИ НЕ ЗАБЖАѢТЕ ДА ДОСТОИИ ТЕЧЕНИЕ ПОСТНОЕ</p>
<p>РАКО/ЖЕ ГКОСПОД⟨⟩БЪНЬ АВРАМІЕ ШТ⟨ВѢ⟩ЩАВ ЖЕ КАНРОСНИ/ЦИ РЕКО ШЖ ЕМОУ ТЪИ ЕС⟨ТЬ⟩</p>	рабь	<p>БОЖІИ И УРЪНЦЪ СВЪРШЕНЬ ЕПН⟨⟩С⟨⟩КОПЬ РЕУ⟨⟩Е КЪ НИМЪ ЕГО</p>
<p>⟨ЗЛА НЫИ БОГЫ ВЪСА РАСКО/ПАЛЪЕС⟨ТЬ⟩ И НЕМОГОШЖ ЕГО ШЗЛОБИТИ ВЪИСТИ/НЖ</p>	рабь	<p>ЕС⟨ТЬ⟩ БОГА ЖИВАГО И ВСѢ РЕЧЕНѢИДА ШТ (НЕГО И/СТИИНА</p>
<p>ДЪНѢР ГКОСПОДНѢ ПОТЫНИМСА ОУБВО НННѢ И ПОДВИГНѢМСА ДА САПЕШБРА/ЩЕМЪ РАКО</p>	рабь	<p>ШНЬ ЕГОЖЕ ПРИШЪД (ГКОСПОДЪ ЕГО ШВРАЩЕТЪ ПИТАЩАСА И НЕРАДАЩА ТѢ/МЖЕО</p>
<p>РАКО НИ ІШСИФΟΥ ОУПОД⟨⟩БИСА НИ БЛАЖЕННѢИ СОУСА/НѢ АЩ</p>	рабь	<p>ПРОДАНЪБЫС⟨ТЬ⟩ ЕГУПТѢНИИЖ РАДИ НЖ НЕПОКОРИСА ЛЫЩЕ/НИ</p>



Current State

- 8 manuscripts, 700 000 forms
 - with keywords (lemma, POS, hom), correlation to other sources, translations
- Annotation tool - ACT
 - distributed annotation
 - direct web exporting
- Created conditions for application of NLP tools
- Possibility to process enormous card files