

# A R A N E A : ΝΑΙΙΑΕΝΟΑΙ ΙΕΕΕΕΑΔΑΙΥΟ ΑΑΑ-ΕΙΘΙΟΝΙΑ

© Αεααει]δ Ααίει (Vladimir Benko). Νετααεεγ, Αδαοεηεααα. Ειηοεοοο  
γυειγίαιεγ ει. Επαταεοα Οοοδα Νετααοειε  
Αεαααιεε ιαοε, ΠΙΑΝΕΙ  
εαοααδα ιιγαιγυ-ιιε ε ιαεεοεουοοδιε ειιοιεεαοεε Οιεααδηεοαοα ει. Βια  
Εταιηεται α Αδαοεηεααα

Εαεοεε ιδααηοααγο ιοταεο  
ηαιαεηοαα ααα-ειθιοηια Aranea αεγ γυειτα εηηιευγίααιυο ε/εεε  
ιδαηιααιυο α ηετααοεεε οιεααδηεοαοαο, ειθιουα ιδααιαγία-αιυ αεγ ιδαηιαααιεγ  
οεεηεταε-αηεεο ε οδαηηεαοηεταε-αηεεο ιδααιαοια, α οαεαα ε αεγ εεηαεηοε-αηεεο  
εηηεααηααιεε.

Ααα-ειθιοη ιδααηοααεγαο ηιαιε ηηιαυε  
αεα εεηαεηοε-αηεηαη ειθιοηα, ειθιουε ηηααη ιοδαη ηηοαηαηιε ααδογεε οαεηοια  
εγ ειοαδταοα ιδε ηηηε ααοηαοεγεδτααιυο ιθιοααοδ, ειθιουα ια εαοο ηηαααεγπο  
γυει ε ειαεδταεο ιααεηιυο ααα ηοδαηεο, οααεγπο οααειιυ, γεαηαου ιααεααοεε,  
ηηυεεε ε δαεεαο (δ.ι. boilerplate),  
ηηουαηοαεγπο οδαηηοηιαοεη ια οαεηο, οεεουοδαοεη, ηηιαεεαοεη ε ααοηεεεαοεη  
ηηεο-αηιυο αηεοιαοια, ειθιουα αοαοη ηηαι ιαδααηοαοη οδααεοειηιυε  
ειηοοηοιαοιαε ειθιοηηιε εεηαεηοεεε (οιεαηεαοεγ, ιεθοηηεηοαεηε-αηεαγ ε  
ηεηοαεηε-αηεαγ αηηοαοεγ) ε αηααδεοη α ηηεηεηαοη ειθιοηηοη ηεηοαηο. Νηααηεα  
ααα-ειθιοηα ια οηεηεη ιαηηηα ααοααεα, ηη ιδαααα αηααη ααη δαγιαδ ηηαο αουοη  
αααα ια ηδγαιε αηεηοα οδααεοειηιυο ειθιοηηια.

Ιδε ηηααηεε αηαο ειθιοηηια αυεα  
ηδεηαηα ιαεηαεηαγ ιαοηεηαεγ ε ιαεηαηα ιααηδ ιδααηαηυο ειηοοηοιαοια: SpiderLing, Onion, Unitok [1] ε TreeTagger [2]. Α εα-α  
ηηεηεηαηε ηεηοαηη εηηηευγοαοηγ NoSketchEngine [3] (open-ource) εεε SketchEngine [4] (ηεαοηαγ). Ο ειθιοηηια ιαααηεγ ια «ιαεοο  
(εαοεηεηηη) γυηεα ιαηηα-αηυαη γυηε ε οηεα δαγιαδ ειθιοηηα, ιαηδεηαδ AraneumAnglicumMinus,  
AraneumRussicumMaius, ε ο. ι. Α ιαηοηηυαα αδαηγ ηαιαεηοαη ηηααδεεο  
18 ειθιοηηια ια 14 γυηεαα α ααοο δαγιαδαο ε αηα ειθιοηηη αηηοοηιυ α ααηηεαοηη  
δααηεα ια ειθιοηηηη ηηοαεα ιοηαεοα [5].

Α ηοεε-εα ηο αη-εηεεοαεηιυο εεηαεηοηα,  
ειθιουα ιαδαααοηααπο ειθιοηηηα ααηηα α ιαεαοηηη δααηεα, ηηοαεηιυα ηηευγίαοαεε  
ειθιοηηια ιαη-ηη αεηοαδαηηααηηη α ηηεηεαο ειηεδαοηυο ηδεηαδηα ηεηα, ηη-αοαηεε ε  
ηεηοαεηε-αηεεο ηοδοεοοδ δαηηα-αοαηηυο α εεαα ειηεηααηηα, -αηοηοηυο ηηεηεηα ε  
ηθιοεεαε ια γεδαηα. Ο-εοηαγ δαγιαδη ηηαδαηαηηυο ειθιοηηηα ηοαηαο ηηηηη, -οη  
γοδαεοεαηηηο ε οαηαηηοη ηηεηεηαηε γαεγαοηγ η-αηη αααηηη οαεοηηη αεγ αηγεηε  
δααηοη η ειθιοηηηη.

Ναηεηαδ ηηηαγυαη ιδαεοεεα δααηοη η  
ηαιαεηοαηη ειθιοηηηη ηηεηεηαηηο ηεηοαηη NoSketch  
Engine [6]ε Sketch Engine [7],  
ηδεηαεεααηεηε ε ηαηηη εο-οεη ιεδα ειηοοηοιαοιαη αεγ δααηοη ηη «ηααδοαηεηοεε»  
ειθιοηηαηε (δαγιαδηη α ααηγοεε ιεεεεαδαηα οιεαηηα). Ιαα ηεηοαηη αηεε ηηααηη α Εααηδαοηδεε ιαδααηοεε αηοαηοααηηαη γυηεα  
Οαεοεηοαοα ειθιουαοεε Οιεααδηεοαοα ει. Ιαηαδεεα α Αδηη, ηδε-αη οοηεοεε  
ααηηεαοηηε (open-source) ηεηοαηηη NoSketch  
Engine γαεγποηγ ηηαιηεαηοαηη οοηεοεε Sketch Engine ε αεεη-απο α ηηεηη  
ιαυαηα δααηοη ηη ηηεηεηαηε ηεηα (Word List)  
ε ειηεηδααηηαδηη, ο. α. ηηεηε ηη ηεηαηοηηα, εαηηα, ηη-αοαηεε ε ηηδοηηεηοαεηε-αηεηε  
ιαδεα α δαγηηο ηηηεηαηοεγθ ια γυηεα CQL  
(Corpus Query Language), ε οηεα

áñ-èñéáíéá éíééíéáòéé íá áàçá ñòàòèñòè-áñéèò íáð ñí-áòàáíñòè (T-score, MI, MI3, log likelihood, min. sensitivity è logDice). Ñèñòáìá ðááíòááò á ðáæèìá ñáðááð/éééáíò, ááá íá ñáðááðá òðáíçòñý áñá ááííúá è ñòóáñòáéççòñý ñèñéíáúá ñáðáòéé è ñéçíááòáéú ðááíòááò ñ éééáíòí ðáðáç ááá-éíòáððáéñ ïðè ñííè ñòáíáàðòííáí áðáòçáðá.

Íèàòíáç ñèñòáìá Sketch Engine ñíááðæèò éðíìá áñáð óóíéòéé NoSketch Engine òðè ñóòáñòááííúð ðáñèððáíçý &ndash; éíééíéá ïðíòééé (ñéáò-è) ñíñòðíáííúá íá áàçá ñéççíááòáéúñéíé ñéáò-áðáìíàòééé, áèñòðéáóòéáíúé òáçàòðòñ è óóíéòèð ñðááíáíçý ñéáò-áé áéç ááóó éáèñè-áñéèð ááéíéò. Áñá çòè óóíéòéé ðááíòáðò ñ ááííúé áñ-èñéáííúé çáðáíáá, ðí ááéááò ñèñòáíó í-áíú áúñòðíé è óáíííé. Ñèñòáìá ïðáíñòááéççòñý á áéáá ñáðáéñá (ñíáíèñéè) íá ñáðááðáò éíííáíéè Lexical Computing, íá éíòíðúð òðáíçòñý éíðíóñú áíéáá ðáí íá 80 ççúéáð, áéèð-àç 15-ìéééèáðáíúé éíðíóñ ðóññéíáí ççúéá.

Éííéíðááíñ òíðíàòá KWIC áéç ñéíáíòíðíú «ííáíñéáéðñééé»

Íðááíñòíðííúá  
éíééíéáòú ñéíáíòíðíú «ííáíñéáéðñééé»

Áèñòðéáóòéáíúé òáçááð áéç éáìíú  
«ííáíñéáéðñééé»

çáñòíóíáç áèñòðéáóòéç éáìíú  
«ííáíñéáéðñééé» ñí TLD

Íðíáðáìá  
éóðñá

«A r a n e a : Ñáíáéñòáí ìéééèáðáíúð  
ááá-éíðíóñíá»

(10 ðáñíá éáéòéííúð çáíçòéé è ïðáéòééóííá,  
16-20 ñíçáðç 2015 á.)

Éáéòéç 1. Ááááíéá. Ééíáéèñòè-áñééé éíðíóñ éáé èñòí-íéé  
éíòíðíàòéé í ççúéá

Íñííáíúá ñíçòéç: áéáú çéáéòðííúð  
éíééáéòéé òáéñòíá

Éñòíðéç ñíçááíçý éíðíóñíá, ááíáðáòéé, ïðéíáíáíçý

Íðíáéòú íáöèííáéúíúò èíðíóñíá

Éíðíóñíáíä èèíááèñòèèá èáè íáòíá èèè íñíááííáíä ááòèá  
ÿçúéíçííáíéÿ

Íðèíáíáíéá íáòíáíá èíðíóñííé èèíááèñòèèè á ñèíððíííí è  
áèáòðíííí èññèááííááíèè ÿçúéá

Íííÿòèá ááá-èíðíóñá, íñíááíííñòè è íðèè-èÿ íò  
òðááèòèííúò èíðíóñíá

Èáèòèÿ 2 Aranea &ndash; Nĩáíáéñòáí íèèèèáðáíúò  
ááá-èíðíóñíá

Íñííáíúá ðáøáíéÿ íðíáéòá Aranea: ÿçúéè, ðàçíáðú è áàðèáíòú, íàçááíéÿ

Éíñòðóíáíòú äéÿ íáðááíòèè: èðáóèèíá, ííðáááèáíéá ÿçúéá, óáàèáíéá  
øááèííá, ááíóíéèèáòèÿ, òíèáíéçáòèÿ, ííðóíñèíòáéñè-áñéáÿ ðàçíáðèá,óíèòèèáòèÿ  
òááñáòíá

Íóáèèèáòèÿ è èñíèüçííááíéá èíðíóñíá: èíðíóñíúá íáíááæáðú

Óíðíàò èíðíóñíá Aranea [8]:  
àðèéáóòú è ñòðóèòóðú

Íðáèòèèóíú:

A. Ðááíòá  
ñ èíðíóñííé íñèñéíáíé ñèñòáííé (No)SketchEngine

Íðáèòèèóí

1. Íñèñé íí ñèíáíòíðíá, èáíá è  
ñèíáííí-áòáíèð

Ðáæèíú èçíáðáæáíéÿ, íáñòðíéèè

Óèèòóðú, èííóáéñò

×áñòíòíúá àèñòðèáóòèè

Óááñáòú (tagsets), íñèñé íí òáááí

## 2. Βζύε ζαίδηηιά CQL

Δαάόεγδύα άύδάαείεγ

Ίδελάίεα CQL:

ηεηε ηελάεηε-αηεεο ηόδóεóóδ

Έηεηεάóεε, ιάδú αηηίεάóεάίηηóε

Άú-εηεάίεά εηεεηεάóεηίίúó εαίεάάóτá

### B. Δαάτò

η έηδύóηίε ηεηεηάίε ηεηόαίε SketchEngine

Ίδελόεέóι 1: Έηδύóηίε ίαίάαεάδ Sketch Engine (https://www.sketchengine.co.uk/)

Έηεηεάóεηίίúά ηδύóεε (ηεάó-ε)

Νεάó--αδύιιáóεε: ηελάεηε-αηεεε ε εηεεηεάóεηίίúε ητáτá

Νεεó--αδύιιáóεε äëγ έηδύóηίá Aranea

Νδάαίáίεά ηεάó-άε (Sketch-diff)

ε äεηóδεάóóεάίúε óαζαάδ

Άάóγζú-ίúά ηεάó-ε

Ίδελόεέóι

2: Δαηόδηú ε είóδύáίóú ηάεóá Sketch Engine

Έηδύóηú ηαίáεηόάα TenTen

Ίδελεεάεúίúά έηδύóηú

Είηόδύáίóú äëγ ηίζαάίεγ ηεúζίáάóáεüηεεó έηδύóηίá: Corpus Architect ε WebBootCaT

Ýεηóδáεóεγ óáδύεηίετáεε

Ά δάιεάó

Ίδελόεέóιηá ηόóάáίóú ίζάεηίγóηγ η δαάίóίε η ίάάεε ηεηόαίáε ε ηεó-άó

ίáτááίε-áίúε áηόóί ε έηδύóηá ίá έηδύóηίηι ηδóάεά ηδύáεά Aranea (http://ucts.uniba.sk/) ε áδάίáίúε 3-ó

ίáηγ-ίúε áηηεάóίúε áεεάóίó äëγ Sketch Engine.

## Webography

<https://savba.academia.edu/VladimirBenko>

[http://ucts.uniba.sk/aranea\\_about/](http://ucts.uniba.sk/aranea_about/)

[1]  
<http://corpus.tools/>

[2]  
<http://www.cis.uni-muenchen.de/~schmid/tools/TreeTagger/>

[3]  
<http://nlp.fi.muni.cz/trac/noske>

[4]  
<https://www.sketchengine.co.uk/>

[5]  
[http://ucts.uniba.sk/aranea\\_about/](http://ucts.uniba.sk/aranea_about/)

[6]<http://nlp.fi.muni.cz/trac/noske>

[7]<https://www.sketchengine.co.uk/>

[8] Nĩáéò ĩđĩáêòà Aranea: [http://ucts.uniba.sk/aranea\\_about/](http://ucts.uniba.sk/aranea_about/)

Êĩđĩóńĩúá ĩđòàèû: Aranea: <http://ucts.uniba.sk/>  
; <http://ella.juls.savba.sk/>

