

Методы повышения эффективности процесса коллективного построения лексических ресурсов

Д. А. Усталов

ИММ УрО РАН

31 марта 2016 г.

Ресурсы: словари, тезаурусы, корпуса текстов, и т. д.

Определение

Тезаурус — словарь, в котором слова и словосочетания с близкими значениями сгруппированы в единицы, называемые понятиями, и в котором явно указываются семантические отношения между этими понятиями.

Применение электронных тезаурусов:

- снятие семантической неоднозначности;
- расширение поисковых запросов;
- анализ вопросов в системах общения;
- и др.

- Сложность и длительность процесса создания тезауруса.
- Высокие требования к квалификации лексикографов.
- Доступность и лицензирование существующих ресурсов.

Предмет, цели и задачи исследования

Предмет исследования

Процесс построения лексических ресурсов.

Цель исследования

Разработать эффективные методы построения лексических ресурсов при помощи краудсорсинга.

Задачи исследования

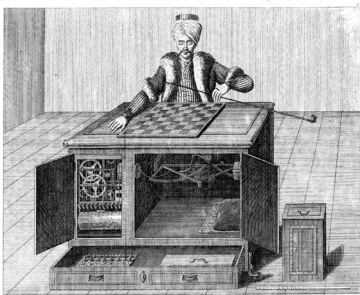
- 1 Разработка методики построения тезауруса при помощи краудсорсинга.
- 2 Разработка вычислительной модели для выполнения процедур разметки.
- 3 Разработка комплекса программ.

- **Традиционный** подход — построение ресурсов узким коллективом экспертов-лексикографов. Процесс оказывается длительным и дорогостоящим.
Пример: WordNet.
- **Автоматический** подход — построение производного ресурса на основе существующих путём автоматического сопоставления. Необходимы качественные данные.
Пример: BabelNet.
- **Автоматизированный** подход — построение языковых ресурсов при помощи краудсорсинга с применением методов обеспечения качества.
Пример: OpenCorpora.

- **RussNet**: 30 тыс. синсетов, 45 тыс. отношений. Данные не представлены в едином формате и не опубликованы.
- **PyТез**: 55 тыс. понятий, 210 тыс. отношений. Успешный ресурс с ограничительной лицензией.
- Автоматический перевод WordNet: **WordNet.ru** (31 тыс. синсетов) и **Russian Wordnet** (недоступен).
- **Викисловарь** (русский): 190 тыс. слов, 70 тыс. отношений. Отношения между словами, оценка качества не проводилась.
- **UNL**: 62 тыс. универсальных слов, 90 тыс. связей между ними. Оценка качества не проводилась.
- **BabelNet** (русский): 2,5 млн. понятий, 380 млн. отношений (всего). Ограничительная лицензия, оценка качества не проводилась.

Определение

Краудсорсинг — коллективный процесс решения задачи, поставленной заказчиком перед толпой участников на специализированной человеко-машинной платформе.



*Il se représente
Un modèle d'un jeu de
P. G. Pezay, le
Par lequel on peut résoudre les problèmes de l'Arithmétique qui se trouvent dans le jeu de l'Arithmétique.*

Основное внимание в данной работе посвящено краудсорсингу микрозадачами.

Определение

Биржа краудсорсинга — человеко-машинная платформа для размещения и выполнения микрозадач.

Биржа	Цена	Кол-во задач	Заказчик?	Участник?
<i>MTurk</i>	\$0.60	3050	Нет*	Да
<i>CrowdFlower</i>	\$0.05	47	Да	Да
<i>microWorkers</i>	\$0.65	93	Да	Да
<i>Яндекс.Толока</i>	\$0.01	3	Да*	Да
<i>TurboText</i>	\$0.06	289	Да	Да

Данные от 1 октября 2015 г.

- Методики решения задач.
- Методы обеспечения качества.
- Программное обеспечение.
- Подходы к оценке качества языковых ресурсов.

Исследователи разрабатывают специализированные методики для решения задач при помощи краудсорсинга.

- **ESP Game** (2006) — игрофицированная разметка изображений.
- **Find-Fix-Verify** (2010) — вычитка, перефразирование и улучшение документов Microsoft Word.
- **PlateMate** (2011) — оценка калорийности блюд по фотографиям.
- **CrowdER** (2012) — обнаружение дубликатов и связывание записей о товарах.
- **TWSI** (2013) — построение языковых ресурсов на основе явления лексической замещаемости.

Предпринимаются попытки адаптации популярных моделей программирования к области краудсорсинга микрозадачами.

- **CrowdForge** (2011) — адаптация модели вычислений MapReduce.
- **CrowdDB** (2011) — оператор CROWD для размещения задач при выполнении SQL-запроса.
- **Qurk** (2011) — SQL-подобный язык для формулирования и выполнения задач.
- **Turkomatic** (2011) — формулирование цепочек заданий на естественном языке.
- **CrowdWeaver** (2012) — потоковая модель вычислений для краудсорсинга.

Оплата труда не гарантирует высокого качества результата.
Создаются специализированные методы вероятностного вывода.

- Метод **Давина-Скина** (1979) — построение матриц ошибок при помощи EM-алгоритма.
- **GLAD** (2009) — вывод сложности заданий, квалификации участников и правильных ответов на основе EM-алгоритма.
- Алгоритм **Каргера-Оха-Шаха** (2011) — оптимальное по порядку назначение заданий и вывод правильных ответов.
- **ZenCrowd** (2012) — вывод квалификации участников и правильных ответов на основе фактор-графов.
- **iCrowd** (2015) — онлайн-алгоритм назначения заданий, оценки участников и вывода ответов.

Создаётся программное обеспечение для упрощения запуска микрозадач и обработки результатов работы.

- **TurKit** (2009) — инструменты для упрощения размещения и выполнения заданий на MTurk.
- **SQUARE** (2013) — средства оценки качества методов обеспечения качества.
- **WebAnno** (2013) — среда для разметки текстов при помощи краудсорсинга.
- **ActiveCrowdToolkit** (2015) — средства оценки качества методов обеспечения качества.
- **CEKA** (2015) — средства анализа результатов выполнения микрозадач.
- **psiTurk** (2015) — сервис организации воспроизводимого процесса разметки.

- **Золотой стандарт** — автоматическое сравнение ресурса с аналогичным ресурсом известного качества.
Пример: WordNet.
- **Экспертная оценка** — ручной анализ ресурса с привлечением эксперта.
Пример: OpenCorpora.
- **Разделённая задача** — соревнование («дорожка») по применению ресурса для решения поставленной прикладной задачи.
Пример: RUSSE.

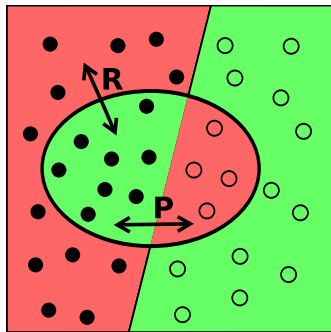
Подходы к оценке качества языковых ресурсов II

Ответы:

- верные положительные (TP),
- верные отрицательные (TN),
- ложные положительные (FP),
- ложные отрицательные (FN).

Метрики:

- Точность $P = \frac{TP}{TP+FP}$.
- Полнота $R = \frac{TP}{TP+FN}$.
- F_1 -мера $F_1 = 2 \frac{P \cdot R}{P+R}$.



Информационно-поисковый
подход.

Разработать методы повышения эффективности процесса построения лексических ресурсов путём обработки существующих словарей и тезаурусов неизвестного качества при помощи краудсорсинга.

- Определены понятия языкового ресурса и краудсорсинга.
- Рассмотрено их практическое применение.
- Сформулирована исследуемая задача.

Постановка задачи

- **Дано:** словарь предметной области \mathbb{D} , множество синсетов \mathbb{S} , множество родовидовых пар слов \mathbb{R} неизвестного качества.
- **Необходимо:** провести очистку данных и построить тезаурус предметной области.

- 1 Извлечь из \mathbb{S} все понятия, содержащие слова из \mathbb{D} .
- 2 Извлечь из \mathbb{R} все пары, слова которых представлены в синсетах.
- 3 Построить неоднозначные родовидовые отношения при помощи краудсорсинга.
- 4 Выполнить уточнение лексикализации понятий при помощи краудсорсинга.
- 5 Объединить дубликаты-когипонимы автоматически.

Краудсорсинг микрозадачами предполагает обработку некоторого набора исходных данных в один или несколько связанных друг с другом этапов.

Определение

Коллективные потоковые вычисления — потоковая вычислительная модель, объединяющая человеко-машинные этапы обработки реляционных данных.

- Данные выражаются расширенной реляционной моделью с операциями упаковывания и распаковывания.
- Процедура решения задачи записывания в виде *схемы коллективных вычислений*.

Определение

Схема коллективных вычислений — слабо связный ориентированный ациклический граф W со множеством рёбер E , множество вершин которого образуется объединением множества этапов разметки S , множества этапов синхронизации Y и множества источников данных D .

$$W = (S \cup Y \cup D, E) \quad (1)$$

Условные обозначения:

- множество элементов схемы: $V = S \cup Y \cup D$;
- реляционное отношение $v \in V$: $(H(v), B(v))$;
- первичный ключ элемента v : $PK(v) \subseteq H(v)$;
- множество входящих вершин в вершину v : $In(v) \subset V$.

Определение

Этап разметки $s \in S$ — это реляционное отношение, тело которого получено путём преобразования толпой участников кортежей единственного входящего отношения:

$$In(s) \subset V \wedge |In(s)| = 1.$$

Определение

Этап синхронизации $y \in Y$ — это реляционное отношение, тело которого получено путём автоматической обработки двух и более входящих отношений: $In(y) \subset V \wedge |In(y)| > 1$.

Определение

Источник данных $d \in D$ — это реляционное отношение, тело которого получено заранее и не зависит от других элементов схемы коллективных вычислений: $In(d) = \emptyset$.

Этапы синхронизации и средства предварительной обработки данных реализуются в виде программного обеспечения различной сложности, что вносит дополнительные требования к детерминированности описания методики краудсорсинга.

Определение

Согласованная схема коллективных вычислений — схема, каждый кортеж каждого элемента которой однозначно идентифицирует породившие его кортежи.

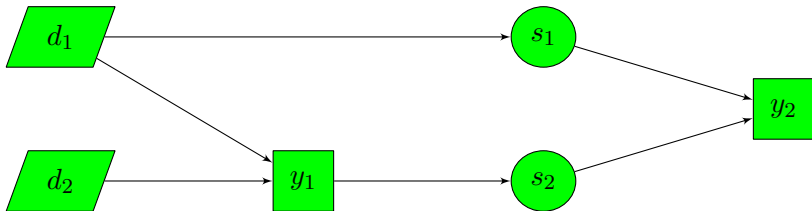
$$\forall v \in V \left(\forall v' \in In(v) (H(v) \cap H(v') \supseteq PK(v')) \right) \quad (2)$$

Алгоритм 1 Синхронный алгоритм выполнения

Require: $V = S \cup Y \cup D, |D| > 0$

- 1: **for all** $v \in V$ **do**
- 2: $M_v \leftarrow (v \in D)$
- 3: **end for**
- 4: **parallel for all** $v \in V \setminus D$ **do**
- 5: $\text{Wait}(\forall v' \in \text{In}(v)(M_{v'} = \text{true}))$
- 6: $B(v) \leftarrow \text{Run}(v)$
- 7: $M_v \leftarrow \text{true}$
- 8: **end for**

Ensure: $(\forall v \in V)(M_v = \text{true})$



Определение

Синсет (синонимический ряд) — множество квазисинонимов, выражающих понятие.

Пример

{автомобиль, машина, авто́, ...}

Синсеты могут быть не идеальны.

- Недостающие слова: {мундир, униформа} \leftarrow {**форма**}.
- Посторонние слова: {знать, понимать, **аристократия**}.

Постановка задачи

- **Дано:** множество синсетов \mathbb{S} и множество слов-кандидатов для включения в них \mathbb{W} .
- **Необходимо:** добавить в синсеты из \mathbb{S} недостающие слова из \mathbb{W} и удалить из их посторонние слова.

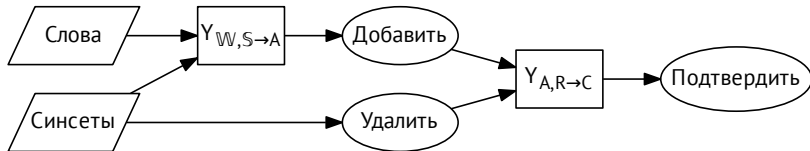
Процедура «добавить–удалить–подтвердить» I

Процедура *ARC*: «добавить–удалить–подтвердить»

Добавить: участник выбирает слова-кандидаты на включение в синсет.

Удалить: участник выбирает посторонние слова для удаления из синсета.

Подтвердить: участник выбирает между оригинальным синсетом и модифицированным.



Элемент	Определение
W	$H(W) = \{S.id, (words, TEXT[])\}$
S	$H(S) = \{(\underline{id}, INT), (words, TEXT[])\}$
$Y_{W,S \rightarrow A}$	$\pi_{S.id, words=S.words, (\underline{W} \bowtie S)}^{candidates=W.words}$
A	$H(A) = \{(\underline{id}, INT), S.id, (added, TEXT[])\}$
R	$H(R) = \{(\underline{id}, INT), S.id, (removed, TEXT[])\}$
$Y_{A,R \rightarrow C}$	$\sigma_{words \neq words'} \left(\pi_{S.id, A.id, R.id, words=S.words, (\underline{S} \bowtie A \bowtie R)}^{words'=S.words \cup A.added \setminus R.removed} \right)$
C	$H(C) = \{S.id, A.id, R.id, Y_{A,R \rightarrow C}.words', (b, BOOL)\}$

Определение

$$S_{ARC} = \{A, R, C\} \quad (3.1)$$

$$Y_{ARC} = \{Y_{\mathbb{W}, \mathbb{S} \rightarrow A}, Y_{A, R \rightarrow C}\} \quad (3.2)$$

$$D_{ARC} = \{\mathbb{W}, \mathbb{S}\} \quad (3.3)$$

$$E_{ARC} = \{(\mathbb{S}, Y_{\mathbb{W}, \mathbb{S} \rightarrow A}), (\mathbb{W}, Y_{\mathbb{W}, \mathbb{S} \rightarrow A}), (\mathbb{S}, R), \\ (Y_{\mathbb{W}, \mathbb{S} \rightarrow A}, A), (A, Y_{A, R \rightarrow C}), \\ (R, Y_{A, R \rightarrow C}), (Y_{A, R \rightarrow C}, C)\} \quad (3.4)$$

$$ARC = (S_{ARC} \cup Y_{ARC} \cup D_{ARC}, E_{ARC}) \quad (3.5)$$

Построение родовидовых отношений I

Определение

Родовидовое (гипо-гиперонимическое) **отношение** — семантическое отношение между парой понятий, при котором одно понятие является разновидностью другого.

Пример

$\{\text{автомобиль}, \dots\} \xrightarrow{is-a} \{\text{транспортное средство}, \dots\}$

Формирование родовидовых отношений между синсетами на основе пар слов приводит к неоднозначности.

- *Отношение:* (ткань, джинса).
- *Синсеты:* {джинса, реклама}, {джинсовая ткань, джинса}.

Постановка задачи

- **Дано:** множество синсетов \mathbb{S} и множество родовидовых пар слов \mathbb{R} .
- **Необходимо:** построить отношения между синсетами в \mathbb{S} на основе \mathbb{R} .

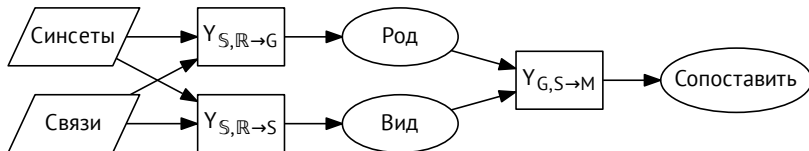
Процедура «род-вид-сопоставить» I

Процедура *GSM*: «род-вид-сопоставить»

Род: участник устанавливает синсет-род для пары *СЛОВ*.

Вид: участник устанавливает синсет-вид для пары *СЛОВ*.

Сопоставить: участник подтверждает осмысленность связи пары *синсетов*.



Процедура «род–вид–сопоставить» II

Элемент	Определение
\mathbb{S}	$H(\mathbb{S}) = \{(\underline{id}, INT), (words, TEXT[])\}$
\mathbb{R}	$H(\mathbb{R}) = \{(\underline{hypernym}, TEXT), (\underline{hyponym}, TEXT)\}$
$Y_{\mathbb{S}, \mathbb{R} \rightarrow G}$	$\pi_{id, words, hypernym, hyponym} (\mathbb{R} \bowtie \pi_{id, words, hypernym = \mu(words)} (\mathbb{S}))$
$Y_{\mathbb{S}, \mathbb{R} \rightarrow S}$	$\pi_{id, words, hypernym, hyponym} (\mathbb{R} \bowtie \pi_{id, words, hyponym = \mu(words)} (\mathbb{S}))$
G	$H(G) = \{(\underline{id}, INT), \mathbb{S}.id, \mathbb{R}.hyponym, (b, BOOL)\}$
S	$H(S) = \{(\underline{id}, INT), \mathbb{S}.id, \mathbb{R}.hypernym, (b, BOOL)\}$
$Y_{G, S \rightarrow M}$	$\sigma_{G.b \wedge S.b = 1} (\pi_{G.id, S.id, G.b, S.b} (\pi_{G.S.id=s_1, S.S.id=s_2} (G \bowtie S)))$
M	$H(M) = \{G.id, S.id, Y_{G, S \rightarrow M}.s_1, Y_{G, S \rightarrow M}.s_2, (b, BOOL)\}$

Определение

$$S_{GSM} = \{G, S, M\} \quad (4.1)$$

$$Y_{GSM} = \{Y_{\mathbb{S}, \mathbb{R} \rightarrow G}, Y_{\mathbb{S}, \mathbb{R} \rightarrow S}, Y_{G, S \rightarrow M}\} \quad (4.2)$$

$$D_{GSM} = \{\mathbb{S}, \mathbb{R}\} \quad (4.3)$$

$$\begin{aligned} E_{GSM} = \{ & (\mathbb{S}, Y_{\mathbb{S}, \mathbb{R} \rightarrow G}), (\mathbb{S}, Y_{\mathbb{S}, \mathbb{R} \rightarrow S}), \\ & (\mathbb{R}, Y_{\mathbb{S}, \mathbb{R} \rightarrow G}), (\mathbb{R}, Y_{\mathbb{S}, \mathbb{R} \rightarrow S}), \\ & (Y_{\mathbb{S}, \mathbb{R} \rightarrow G}, G), (Y_{\mathbb{S}, \mathbb{R} \rightarrow S}, S), \\ & (S, Y_{G, S \rightarrow M}), (G, Y_{G, S \rightarrow M}), \\ & (Y_{G, S \rightarrow M}, M) \} \end{aligned} \quad (4.4)$$

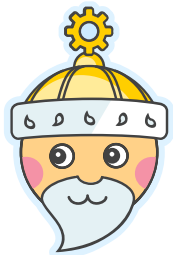
$$GSM = (S_{GSM} \cup Y_{GSM} \cup D_{GSM}, E_{GSM}) \quad (4.5)$$

- Предложена методика коллективного построения тезауруса предметной области.
- Предложена модель коллективных потоковых вычислений.
- Предложена процедура «добавить–удалить–подтвердить» для уточнения лексикализации понятий.
- Предложена процедура «род–вид–сопоставить» для построения родовидовых отношений между понятиями.

Сервис управления процессом краудсорсинга

Сервис управления процессом краудсорсинга реализован в виде веб-сервиса на основе архитектуры REST.

- **Платформа:** Java 8.
- **Программный каркас:** Dropwizard (JAX-RS + Java EE).
- **Хранилище данных:** PostgreSQL.
- **Среда развёртывания:** Docker.



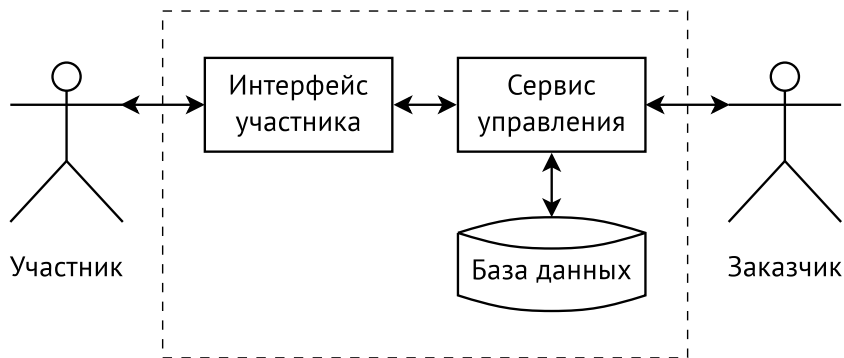
🌐 <http://mtsar.nlpub.org/>

🔗 <https://github.com/mtsar/mtsar/>

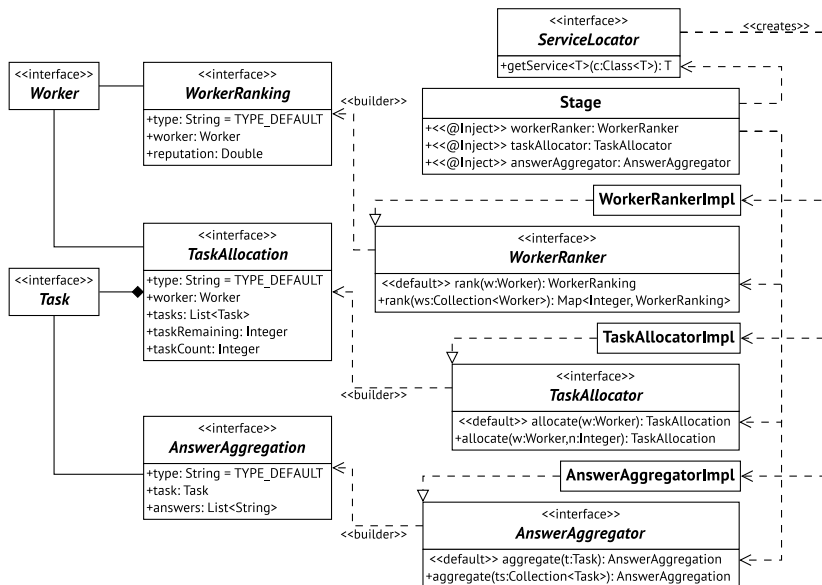
- Назначение заданий участникам.
- Приём ответов от участников.
- Агрегация ответов на задания.
- Оценка согласованности ответов.
- Ввод и вывод данных заказчика.

- Запрос регистрации участника.
- Запрос назначения задания.
- Запрос приёма ответа.

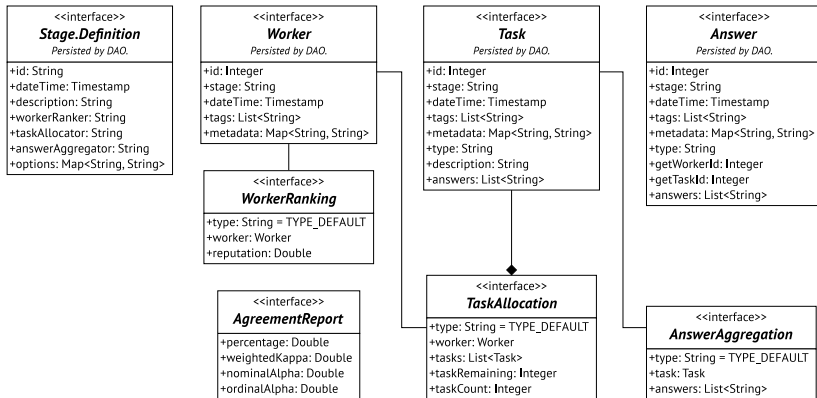
Архитектура системы



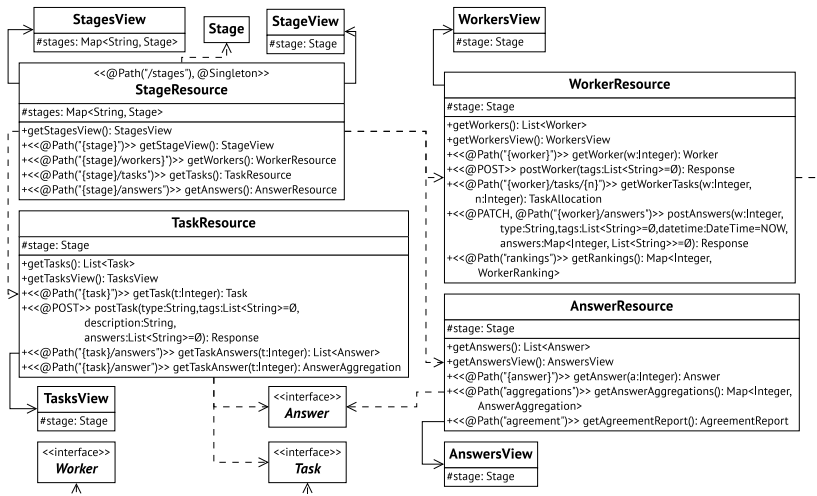
UML-диаграмма классов: процессоры



UML-диаграмма классов: сущности



UML-диаграмма классов: ресурсы



Dashboard

Key	Value
Version	0.0.1-SNAPSHOT
Java	1.8.0_45-internal-b14
Processes	3
Workers	73
Tasks	300
Answers	1313

 Dashboard

 Processes

 GitHub

[Mechanical Tsar](#)

Processes

ID	Description	Workers	Tasks	Answers
arc-add	Add-Remove-Confirm: Add.	24	100	501
arc-rm	Add-Remove-Confirm: Remove.	29	100	512
arc-confirm	Add-Remove-Confirm: Confirm.	20	100	300

 Dashboard

 Processes

 GitHub

[Mechanical Tsar](#)

Process "arc-add"

Key	Value	Action
description	Add-Remove-Confirm: Add.	
workerCount	24	Details
workerRanker	mtsar.processors.worker.ZeroRanker	
taskCount	100	Details
taskAllocator	mtsar.processors.task.FixedNumberAllocator	
answerCount	501	Details
answerAggregator	mtsar.processors.answer.EmptyAggregator	

Additional Options

Key	Value
answersPerTask	5

Genus-Species-Match

Stage “Species”

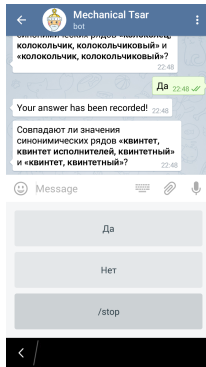
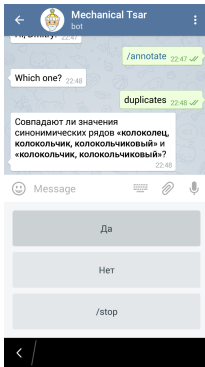
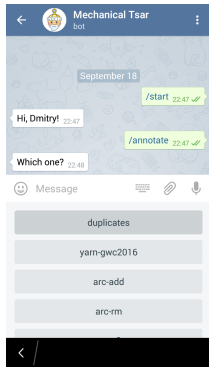
Task	Your Answer
Правда ли, что « лицо , клиент , контрагент » — это частный случай понятия сторона ?	<input type="radio"/> no <input checked="" type="radio"/> yes

Your ID is **653**. There are 833 tasks left.

Submit

You may [return](#) to annotation processes at any moment.

Telegram-интерфейс участника



- Разработан комплекс программ для управления процессом краудсорсинга.
- Комплекс программ поддерживает разметку с разных устройств и платформ.

- Исследование применимости процедуры «добавить–удалить–подтвердить» (*ARC*).
- Исследование применимости процедуры «род–вид–сопоставить» (*GSM*).
- Построение тезауруса предметной области.

Условия эксперимента

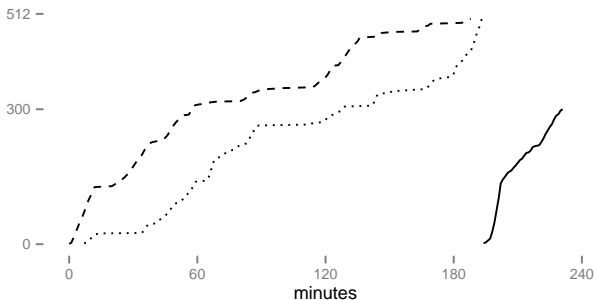
- **Данные:** 100 синсетов YARN с наибольшим количеством дубликатов по эвристике «два общих слова».

$$\exists s_1 \in S, s_2 \in S : s_1 \neq s_2 \wedge |s_1 \cap s_2| \geq 2 \quad (5)$$

- **Участники:** открытый вызов в VK, Facebook, Twitter.
- **Агрегация:** голос большинства.

Применимость *ARC II*

Этап	Участников	Заданий	Ответов	Длит-ть
Добавить	24	100	501	188
Удалить	29	100	512	194
Подтвердить	4	100	300	37
Итого	36	300	1313	231



- Привлечены два эксперта: ставилась оценка 1, если обработанный синсет улучшился; если нет — 0.
- Всего изменилось 84 синсета, из них улучшилось 70.

Участники		Эксперты	
Изменилось	84	Улучшилось	70
Не изменилось	16	Не улучшилось	14
Всего	100	Всего	84

Оценки экспертов согласуются: различается лишь 20 оценок из 84, индекс Жаккара равен $1 - \frac{20}{84} = 74\%$.

Классы ошибок

- Неоднозначный синсет остался неоднозначным после обработки.
- Добавлен гипероним, гипоним или когипоним вместо синонима.
- Лишнее слово не удалено несмотря на добавление недостающих.
- Общее значение синсета изменилось.

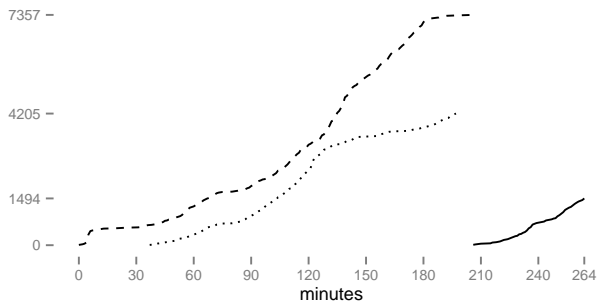
Результаты подтверждают применимость процедуры *ARC* для уточнения лексикализации понятий.

Условия эксперимента

- **Данные:** 2271 синсет YARN по тематике «безопасность жизнедеятельности», 383 кандидата-отношения.
- **Участники:** пользователи биржи TurboText.
- **Агрегация:** голос большинства, KOS, ZenCrowd.

Применимость *GSM II*

Этап	Участников	Заданий	Ответов	Длит-ть
Род	39	1438	7357	206
Вид	18	833	4205	197
Сопоставить	18	287	1494	58
Итого	47	2558	13056	264



Результаты проанализированы экспертом и сопоставлены с эвристическим методом:

$$|\{s' : \exists(g, s') \in \mathbb{R}\} \cap \{g' : \exists(g', s) \in \mathbb{R}\}| > 1. \quad (6)$$

Метод	TP	TN	FP	FN	P	R	F ₁
Эвристика	40	102	17	128	0,70	0,24	0,36
MV	129	57	62	39	0,68	0,77	0,72
KOS	142	63	56	26	0,72	0,84	0,78
ZenCrowd	146	69	50	22	0,74	0,87	0,80

Метод ZenCrowd продемонстрировал лучшие результаты с точки зрения точности, полноты и F_1 -меры.

Классы ошибок

- Некорректное понимание лексических значений в заданиях.
- Ошибки в синсетах: чрезмерная общность или узость понятий.
- Ошибки в данных: недоверенные зашумлённые источники.

Результаты подтверждают применимость процедуры *GSM* для построения родовидовых отношений между понятиями.

Условия эксперимента

- **Данные:** словарь предметной области, синсеты, кандидаты-отношения.
- **Участники:** пользователи биржи TurboText и CrowdFlower.
- **Агрегация:** голос большинства, KOS, ZenCrowd.

Эксперимент в процессе

Исследование выполняется в настоящее время.

- Подтверждена применимость процедуры «добавить–удалить–подтвердить».
- Подтверждена применимость процедуры «род–вид–сопоставить».
- Построен русскоязычный тезаурус предметной области.

Заключение

Цели достигнуты, задачи выполнены.

- **Интеграция с медицинскими технологиями:** использование данных ЭЭГ и других датчиков для отправки ответов.
- **Развитие модели вычислений:** асинхронное выполнение, автоматическое бюджетирование, и т. д.
- **Снижение входных барьеров:** априорная оценка сложности заданий, профилирование участников.
- Построение тезаурусов других предметных областей.

Работ в библиографической базе **Web of Science**: 2.

- *Ustalov D. Enhancing Russian Wordnets Using the Force of the Crowd // Analysis of Images, Social Networks and Texts. — Springer International Publishing, 2014. — Vol. 436 of Communications in Computer and Information Science. — P. 257–264.*
- *Ustalov D. Towards Crowdsourcing and Cooperation in Linguistic Resources // Information Retrieval. — Springer International Publishing, 2015. — Vol. 505 of Communications in Computer and Information Science. — P. 348–358.*

Работ в библиографической базе **Scopus**: 4.

- *Ustalov D., Kiselev Y. Add-Remove-Confirm: Crowdsourcing Synset Cleansing* // Application of Information and Communication Technologies (AICT), 2015 IEEE 9th International Conference on. — IEEE, 2015. — P. 143–147.
- *Ustalov D. Crowdsourcing Synset Relations with Genus-Species-Match* // Proceedings of the AINL-ISMW FRUCT. — 2015. — P. 118–124.
- *Kiselev Y., Ustalov D., Porshnev S. Eliminating Fuzzy Duplicates in Crowdsourced Lexical Resources* // Proceedings of the Eighth Global Wordnet Conference. — 2016. — P. 161–167.
- *YARN: Spinning-in-Progress* / P. Braslavski, D. Ustalov, M. Mukhin, Y. Kiselev // Proceedings of the Eighth Global Wordnet Conference. — 2016. — P. 58–65.

Работ в журналах перечня **ВАК**: 2 (ожидается).

- Усталов Д. Инструментарий краудсорсинга для механизированного труда // *Труды Института системного программирования РАН*. — 2015. — Т. 27, № 3. — С. 351–364.
- Усталов Д. Коллективные потоковые вычисления: реляционные модели и алгоритмы.

Прочие работы и публикации: 3.

- *Ustalov D. Teleboyarin—Mechanized Labor for Telegram* // Proceedings of the AINL-ISMW FRUCT. — 2015. — P. 195–197.
- *Усталов Д.* Свидетельство Роспатента о государственной регистрации программы для ЭВМ «Адаптивная система управления процессом краудсорсинга» № 2015662640 от 30.11.2015.
- *Усталов Д.* Свидетельство Роспатента о государственной регистрации программы для ЭВМ «Система автоматизации процесса коллективного построения баз данных» № 2015662780 от 01.12.2015.

- 14-я конференция европейского отделения Ассоциации по компьютерной лингвистике **EACL 2014** (г. Гётеборг, Швеция).
- Международная суперкомпьютерная конференция «**Научный сервис в сети Интернет**: многообразие суперкомпьютерных миров» (г. Новороссийск, 2014 г.)
- 14-я национальная конференция по искусственному интеллекту **КИИ-2014** (г. Казань).
- 5-я международная конференция по инженерии знаний и Семантической паутине **KESW 2014** (г. Казань).
- 16-я всероссийская научная конференция **RCDL 2014** (г. Дубна).

- 3-я и 4-я международная конференция по анализу изображений, социальных сетей и текстов **АИСТ'2014** и **АИСТ'2015** (г. Екатеринбург).
- 9-й весенне-летний коллоквиум молодых ученых по программной инженерии **SYRCoSE 2015** (г. Самара).
- 21-я международная конференция по компьютерной лингвистике «**Диалог 2015**» (г. Москва).
- 9-я международная конференция по использованию информационно-коммуникационных технологий **АИСТ2015** (г. Ростов-на-Дону).
- Международная конференция **AINL-ISMW FRUCT** (г. Санкт-Петербург, 2015 г.)

- 8-я глобальная конференция по ворднетам **GWC 2016** (г. Бухарест, Румыния).
- Международная молодёжная школа-конференция «**Современные проблемы математики и её приложений**» (г. Екатеринбург, 2016 г.)

- Исследование выполнено при финансовой поддержке РГНФ: проект «Новый открытый электронный тезаурус русского языка» № 13-04-12020 и проект «Интеграция тезаурусов RussNet и YARN» № 16-04-12019;
- Поддержка данного проекта осуществлена в рамках благотворительной деятельности, на средства, предоставленные Фондом Михаила Прохорова.
- Работа выполнена при финансовой поддержке стипендии Президента Российской Федерации молодым учёным и аспирантам № СП-773.2015.5.
- Исследование выполнено при финансовой поддержке РФФИ в рамках научного проекта № 16-37-00354 мол_а «Методы автоматизации процесса коллективного построения лингвистических ресурсов».

Спасибо за внимание!

Дмитрий Усталов

 <https://linkedin.com/in/ustalov>

 <https://ustalov.name/>

 dau@imm.uran.ru