

# Исследование графа категорий английской версии Wikipedia

© А. В. Шкотин

Государственный Геологический Музей РАН

Москва

ashkotin@acm.org

## Аннотация

Wikipedia является выдающимся проектом по накоплению знаний, как общего пользования, так и различных областей специализации. Проверка качества этих знаний, особенно автоматическая, чрезвычайно важна. В работе представлены результаты изучения строения английской версии ГKB (орграф категорных статей Википедии) в целом. Являясь по идее системой тем он поддерживает систематизацию знаний и мы интересуемся из чего эта систематизация состоит и как она устроена. Показано, что в графе есть неприемлемые логические нарушения и обсуждаются организационные и технические методы их устранения.

## 1 Введение

ГKB [1] есть подграф графа в котором статьи Википедии приписаны категорным статьям. Выделение ГKB из этого полного графа есть первая техническая задача. Важно, что далее изучается дампы ГKB на некоторый момент времени и в нём есть незавершённая "строящаяся" часть. Поэтому выводы надо делать с осторожностью. Естественно ввести термин "точка роста", когда мы натываемся "в дампе" на часть, которая ещё не завершена. Дамп полного графа получен из ИСП РАН и соответствует 16 сентября 2010г. Дамп состоит из двух текстовых файлов: файла отображения номера страницы Википедии в номер категорной страницы, что приписывает страницу категории; а также файла в котором номеру страницы Википедии приписано её наименование. Математически ГKB есть орграф каждый узел которого взаимно-однозначно соответствует категорной странице и помечен её номером. Дуга из узла N1 в узел N2 идёт тогда и только тогда, когда страница с номером N1 есть под-категория страницы с номером N2. Всего таких стрелок 1221133.

В статье исторически вместо термина «дуга» употребляется термин «стрелка».

Множество узлов ГKB (593796 штук), как и

Труды 14-й Всероссийской научной конференции «Электронные библиотеки: перспективные методы и технологии, электронные коллекции» — RCDL-2012, Переславль-Залесский, Россия, 15-18 октября 2012 г.

любого произвольного графа, распадается на два подмножества: изолированные узлы (26272) и узлы связанные стрелками (567524 узлов). Изолированная категория это скорее всего "точка роста": на момент снятия дампа она уже есть, но в ГKB ещё не включена.

Далее анализируется только "граф стрелок", т.е. все характеристики даны без учёта изолированных узлов. Состав изолированных узлов можно посмотреть в отчёте [4] (далее - отчёт) в таблице указанной во введении. Состав и характеристики узлов со стрелками можно посмотреть в таблице указанной там же, равно как и граф стрелок. Важный вопрос - количество связанных компонент графа, т.к. в дальнейшем их строение можно изучать отдельно. Таких компонент оказалось 1987. Изолированные узлы при этом учитываются отдельно. Алгоритм разбиения описан в отчёте [5]. Впрочем проще воспользоваться программой, например, Rajek [2] умеющий разбивать узлы графа на слабо связанные компоненты.

Первые 10 самых больших компонент:

cn	count	cn	count
1	561636	6680	20
21727	210	19212	19
14332	36	20868	19
2863	29	13325	17
20842	27	13287	16

Здесь cn - уникальный номер компоненты, присвоенный при разбиении. Конечно в случае с Wikipedia малые компоненты это точки роста. Петель (N1 → N1) в графе нет.

Источников (узлов в которые нет входящих стрелок) - 345597. Это категории нижнего уровня. Стоков (узлов из которых нет исходящих стрелок) - 11767. Это категории верхнего уровня дампа и скорее всего "точки роста". Промежуточных узлов, соответственно - 210160.

Максимальное количество исходящих из одного узла стрелок - 85. Это промежуточный узел № 690451, а заголовок "Category:World War II", т.е. эта категория приписана 85 над-категориям. Максимальное количество входящих стрелок (12625)



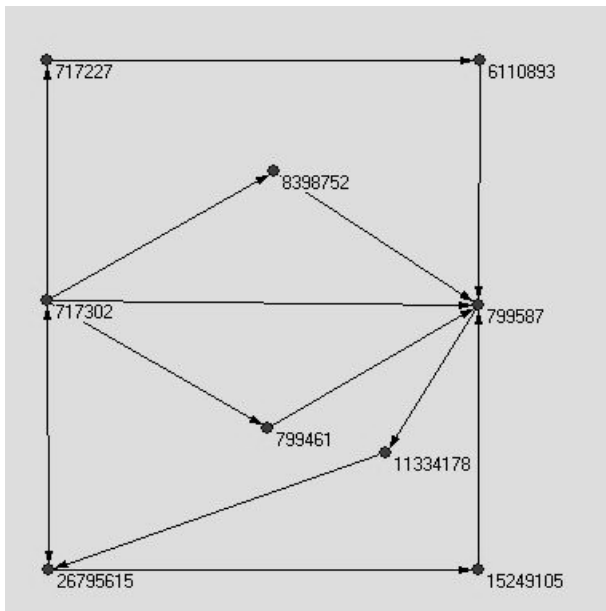


Рис. 1 Рисунок графа компоненты ССК 41

## 5.2 Сильно связанные компоненты ядра

В ядре нас интересует зацикливание отношения под-категория - над-категория. Тут есть два подхода:

- общий - применить алгоритм поиска сильно связанных компонент (ССК);

- частный - найти так называемые "линзы" - два узла ссылающиеся друг на друга (как под-категория - над-категория).

Второй путь вполне приемлем для ГКВ, т.к. по идее в нём вообще не должно быть циклов. Впрочем как для линз так и для циклов большей длины следует заметить, что они математически утверждают эквивалентность соответствующих терминов, т.е. синонимию, что в принципе возможно. Но конкретно в Википедия может быть реализовано через `gedirect`. Интуитивно же в большинстве случаев мы обнаружим ошибку, т.е. какие-то стрелки цикла ошибочны.

Чтобы получить состав сильно связанных компонент ядра была использована программа `Rajek` [2]. Заметим, что петель в ГКВ нет, а поэтому узлы ядра не попавшие в ССК это узлы на путях между циклами (см. выше).

ССК оказалось 457. Узлов не входящих в ССК, так сказать связующих ядра — 7646. Есть одна гигантская по сравнению с остальными ССК - в ней 3967 узлов.

В отчёте в разделе «Сильно связанные компоненты ядра» приведена таблица самых больших ССК. Рассмотрим для примера компоненту №41 у которой всего 9 узлов (см. рис. 1).

Если номер накладывается на стрелку то под ним конечника (треугольничка) нет. Это важно т.к. `Rajek` рисует "линзы" ( $U1 \rightarrow U2 \rightarrow U1$ ) как одну стрелку с конечниками на обоих окончаниях. В данной ССК линза одна - слева внизу вертикально.

Заголовки узлов:

ni	title
717227	Category:Orthodox rabbis
717302	Category:Talmud rabbis
799461	Category:Mishnah
799587	Category:Talmud
6110893	Category:Talmudists
8398752	Category:Talmud people
11334178	Category:Rabbinic literature
15249105	Category:Talmud concepts and terminology
26795615	Category:Chazal

## 5.3 Линзы

Линза это два узла такие что:  $U1 \rightarrow U2$  и  $U2 \rightarrow U1$ . Она может быть отдельной ССК, а может входить в ССК как часть.

В ядре оказалось 1269 линз. Из них 1260 имеют заголовки для обоих узлов. Их можно посмотреть в таблице указанной в разделе «Линзы» отчёта.

## 6 Мантия - ациклическая часть ГКВ

Чтобы получить мантию мы удаляем из ГКВ ядро. При этом оказывается, что часть источников и стоков станет изолированными. В первом случае все из них исходящие стрелки попали в ядро, во втором - все входящие в них стрелки шли из ядра. Изолировавшихся источников - 14421, а стоков - 60.

Кроме того в мантии появляются ложные вершины (пики). Это те её узлы, которые стали стоками после удаления ядра, а вообще-то имели исходящие стрелки, которые все попадали в ядро. Таких вершин 18157. Причём максимальная высота - 28. Для сравнения, стоков ГКВ получивших уровень, т.е. не изолированных - 11707, максимальная высота - 24.

Ложная вершина - рекордсмен (высоты 28) имеет №15715670, а заголовок "Category:Creation myths".

Замечание. Конечно ГКВ можно представить и в виде "галстука-бабочки" как в работе [6], где оргграф был использован для представления схемы связей между транснациональными корпорациями. Но в данном случае сравнение с горами нагляднее - вверх к более обширным темам; горами в которых есть ядро из 20-ти связанных компонент. Одна из которых большая, а 19 - линзы.

Количество узлов на уровнях показано ниже в табличке и оправдывает сравнение с горами:

level	count	level	count	level	count
NULL	14481	NULL	14481	NULL	14481
28	1	18	50	9	1915
27	2	17	57	8	3103
26	3	16	71	7	4858
25	3	15	100	6	7754
24	5	14	149	5	13019
23	7	13	226	4	23302
22	12	12	425	3	45323
21	16	11	697	2	105958
20	20	10	1187	1	1331205
19	30			0	13545

В таблице в строке с level = NULL - количество изолированных узлов мантии, а у 0 - количество узлов в ядре.

## 7 Связующие стрелки

Между мантией и ядром есть стрелки-связующие. Стрелок из ядра в мантию - 591. Стрелок из мантии в ядро — 210514.

## 8 Другие способы исследования

Можно напрямую изучать <http://dbpedia.org> через точку входа для SPARQL: <http://dbpedia.org/sparql>. Привязка к категории идёт через свойство <http://purl.org/dc/terms/subject>.

Вот пример запроса, который начинает выдавать полный граф связи страниц и категорий:

```
select ?x ?z where {?x dct:subject ?z}
```

Надо только поставить timeout, например, 1000.

Запрос

```
select ?x ?z where {?x skos:broader ?z}
```

выдаёт отношение "x is a sub-category of z". см. с. p.5 "Categories." [3]

А вот запрос

```
select ?x ?z where
{?x skos:broader ?z. ?z skos:broader ?x.}
```

выдаёт "линзы".

Вот узлы первой:

```
http://dbpedia.org/resource/Category:Political_philosophers
и
http://dbpedia.org/resource/Category:Political_theorists
```

Она действительно есть в Wikipedia(en).

А всего запрос выдаёт 2000 линз, что наверно не предел.

## 9 Заключение

Естественно считать, что ГКВ должен быть ациклическим графом. Таким образом исследование показало, что аномалии значительны.

Можно создать средства, которые обнаруживая аномалию, например линзу, будут размещать на соответствующих страницах в Discussion уведомление о логическом противоречии.

Основных вопросов два:

1. Как к такому подходу отнесутся авторы страниц категорий? Это можно проверить экспериментально.

2. Как к логическим противоречиям относятся идеологи Википедии? Те кто задаёт правила классификации. Судя по всему индифферентно.

Общая рекомендация. Многие отношения между категориями попавшие в sub-category of следует перенести в See also.

Оценить предстоящую работу можно так: для начала надо разобраться с 1269 линзами. Они сильно убавят размер ССК.

Только если это нужно википедистам можно было бы продолжить и:

- Исследовать длинные пути.

- Попытаться представить архитектуру графа в целом. Например применить 3D визуализацию.

- Проанализировать состав и логику связи заголовков (особенно ССК).

Особняком стоит задача получить и проанализировать русский ГКВ. В проекте dbpedia можно получить дампы русской версии, надо только перекодировать с rdf-кодов букв (например, \u0432) в UTF-8.

## Литература

- Anton Korshunov, Denis Turdakov, Jingu Jeong, Minh Lee, Changsung Moon. A Category-Driven Approach to Deriving Domain Specific Subset of Wikipedia. Proceedings of SYRCoDIS'11: The Seventh Spring Researchers Colloquium on Databases and Information Systems, 2011, pp. 43-53.
- Batagelj V., Mrvar A. Pajek reference manual. Ljubljana, April 16, 2012.
- Christian Bizer, Jens Lehmann, Georgi Kobilarov, Soren Auer, Christian Becker, Richard Cyganiak, Sebastian Hellmann. DBpedia - A Crystallization Point for the Web of Data. May 25 2009.
- Шкотин А., Исследование графа категорий английской версии Wikipedia, Сообщение о результатах первого этапа, Интернет, 2011. <https://sites.google.com/site/alex0shkotin/grafy/wikipedia-category-graph>
- Шкотин А., Разбиение графа на связанные компоненты, Алгоритм и программа, Интернет, 2011. <https://sites.google.com/site/alex0shkotin/grafy/svaznye-komponenty>
- Stefania Vitali, James B. Glattfelder, Stefano Battiston. The network of global corporate control. ArXiv.org, 2011 <http://arxiv.org/abs/1107.5728>
- OWL 2 Web Ontology Language: Structural Specification and Functional-Style Syntax. Boris Motik, Peter F. Patel-Schneider, Bijan Parsia, eds.

W3C Recommendation, 27 October 2009.  
<http://www.w3.org/TR/owl2-syntax/>

## **Investigation of the English version of the Wikipedia categories graph**

Alexander Shkotin

Wikipedia is the outstanding project of knowledge accumulation. The knowledge is both of the general use,

as well of various specialization domains. Quality check of this knowledge, especially automatic, is very important. In this paper results of studying of a structure of the English version of WCG (Wikipedia Categories Graph) as a whole are presented. WCG is a system that supports structure of knowledge and we are interested in WCG content and its arrangement. It is shown that in graph there are unacceptable logical violations; organizational and technical methods for their elimination are discussed.