

Описание и использование тезаурусов в информационных системах, подходы и реализация

М.Х. Нгуен, А.С. Аджиев

ВЦ РАН, Москва

В статье рассмотрены разные подходы к формализации тезаурусов, а также стандарты ISO, ANSI и ГОСТ. Сделан анализ некоторых возможных платформ для такой формализации, описаны особенности работы с тезаурусами в информационных системах, а также проблемы при этом возникающие, требования к реализации тезауруса в рамках SemanticWeb [12].

Рассмотрены особенности и различия классификаторов ресурсов и обычных терминологических и лингвистических тезаурусов. Дан сравнительный анализ существующих схем данных и подходов к реализации тезаурусов для информационных систем на основе RDF. Рассмотрены также вопросы организации пользовательских интерфейсов для работы с тезаурусами, и использования их при поиске в информационной системе, а также интерфейсы администрирования тезаурусов.

Во второй части статьи на основании проделанного анализа сформулированы требования к описанию тезауруса в ИСИР, и приведена общая универсальная схема данных для представления тезауруса в этой информационной системе, удовлетворяющая перечисленным требованиям, и небольшой пример реализации в ней классификатора MSC.

На основании предложенной общей универсальной схемы и сформулированных требований описана реализация тезауруса в ИСИР.

Тезаурусы в описании информации

Для описания какой-либо предметной области всегда используется определенный набор терминов, каждый из которых обозначает или описывает какое-либо понятие или концепцию из данной предметной области. Совокупность терминов, описывающих данную предметную область, с указанием семантических отношений (связей) между ними является *тезаурусом*. Такие отношения в тезаурусе всегда указывают на наличие смысловой (семантической) связи между терминами.

Основным отношением (связью) между терминами в тезаурусе является связь между *более широкими* (более выразительными) и *более узкими* (более специализированными) понятиями. Часто выделяют 2 подвида этого отношения:

- Один термин обозначает понятие, являющееся частью понятия, обозначаемого другим термином (например, «наука» и «математика», «математика» и «теория чисел»)
- Один термин обозначает элемент класса, обозначаемого другим термином («горные районы» и «Кавказ»).

Это отношение на множестве терминов является отношением частичного порядка, то есть множество терминов с такими связями образует ациклический граф, или полииерархическую структуру.

Существуют также и другие связи между терминами. Например, одно понятие или концепция может быть обозначено несколькими терминами, являющимися синонимами. Некоторые термины могут быть антонимами для других. Часто среди терминов, относящихся к одному понятию, выделяют единственный (для каждого языка тезауруса) *наиболее предпочтительный* (наиболее подходящий) термин, который наиболее хорошо характеризует, или обозначает данное понятие. Остальные термины являются *менее предпочтительными* (менее подходящими).

Помимо вышеописанных, между терминами могут существовать также и другие, *ассоциативные связи*, если понятия, обозначаемые этими терминами, как-либо связаны между собою по своему смыслу, за исключением описанных выше иерархических связей.

В многоязычных тезаурусах существуют также *связи эквивалентности* между терминами на разных языках. Выделяют полную (строгую) эквивалентность, и несколько видов частичной (нестрогой) смысловой эквивалентности терминов на разных языках.

Тезаурус часто содержит *комментарии* к терминам, раскрывающие для пользователя смысл термина, а также поясняющие, как следует его использовать.

Тезаурусы применяются, прежде всего, для классификации и поиска информационных ресурсов. При этом каждому ресурсу при классификации могут быть сопоставлены одно или более понятий, описываемых терминами в тезаурусе, а пользователь, осуществляющий поиск, может по тезаурусу найти интересующие его понятия в данной предметной области, а также все характеризующие их термины. То есть на основе связей тезауруса происходит расширение поискового запроса (расширение слов запроса синонимичными, более общими или более частными по смыслу терминами). Навигация по связям тезауруса помогает четче сформулировать сам запрос.

Существует ряд стандартов разного уровня значимости и проработанности на формат представления тезаурусов. Эти стандарты представляют тезаурус в виде

набора объектов нескольких типов, между которыми может быть несколько типов связей. Некоторые стандарты (например, стандарт ANSI/NISO Z39.19-1993) регламентируют также формат представления тезауруса в линейаризованном (текстовом) виде, пригодном для восприятия, как машиной, так и человеком.

Стандарты ISO и ANSI/NISO Z39.19-1993

Основными документами, регламентирующим формат представления тезауруса, являются стандарты ISO 2788-1986 для описания одноязычных тезаурусов, и ISO 5964-1985 для многоязычных.

Стандарт ISO 2788-1986 определяет тезаурус, как набор терминов, связанных между собою соответствующими связями (отношениями).

Термины могут иметь следующие атрибуты:

- **SN** - *Scope Note*. Комментарий к термину. Например, представляет вербальное пояснение термина, или правила его использования.
- **TT** - *Top Term*. Признак, Выделяющий термины на самом верхнем уровне иерархии (термины наиболее общих понятий в данной иерархии понятий).

Связи между терминами могут быть следующими:

- **USE** - Связывает термин с наиболее предпочтительным (на том же языке) термином для данного понятия. **AUSEB** = термин *B* является наиболее предпочтительным для понятия, обозначаемого термином *A*.
- **UF** - *Used For*. Обращение связи **USE**. Связывает наиболее подходящий термин с синонимами и квазисинонимами (менее подходящими терминами). **AUFBóBUSEA**.
- **BT** - *Broader Term*. Связь термина с термином более общего понятия. **ABTB** = термин *B* обозначает более общее понятие по сравнению с понятием, обозначаемым термином *A*.
- **BTG** - *Broader Term Generic*. Вариант связи **BT** в случае, когда термин характеризует разновидность понятия, определяемого более общим термином. Например, «Попугаи» и «птицы». Наличие связи **BTG** подразумевает наличие связи **BT**.
A BTG B ó A BT B.
- **BTP** - *Broader Term Partitive*. Вариант связи **BT** в случае, когда термин характеризует часть понятия, определяемого более общим термином. Например, «математика» и «теория чисел». Наличие связи **BTP** подразумевает наличие связи **BT**.
A BTP B ó A BT B.
- **NT, NTG, NTP** - *Narrower Term, Narrower Term Generic, Narrower Term Partitive*. Обращение связей **BT**, **BTG** и **BTP** соответственно.
A NT B ó B BT A; A NTG B ó B BTG A; A NTP B ó B BTP A.
- **RT** - *Related Term*. Ассоциативная связь. Связывает семантически связанные между собою термины, не находящиеся при этом в одной иерархии, и не являющиеся синонимами или квазисинонимами. Эта связь проставляется в тех случаях, когда пользователю тезауруса может быть полезно

осуществлять поиск или индексацию не только по данному термину, но и по связанному с ним. Связь должна быть двунаправленной (симметричной):
A RTB ó B RT A.

Структура многоязычных тезаурусов регламентируется стандартом ISO 5964-1985. В нем, помимо всех вышеперечисленных связей и требований к ним, есть также связи между эквивалентными терминами на разных языках. Существуют следующие типы таких связей:

- *Полная эквивалентность*
- *Неполная эквивалентность* (значения терминов не совпадают, но пересекаются)
- *Частичная эквивалентность* (значение одного термина шире, чем значение другого)
- *Эквивалентность один ко многим* (значение одного термина соответствует совокупности значений нескольких терминов).

Американский стандарт ANSI/NISO Z39.19-1993 расширяет и уточняет стандарт ISO 2788-1986 для одноязычных тезаурусов, а также накладывает ряд дополнительных ограничений на структуру тезауруса. Основные его отличия следующие:

Добавлены новые связи между терминами:

- **BTI** - *Broader Term Instance*. Вариант связи **BT** в случае, когда термин характеризует элемент класса, или частный случай понятия, определяемого более общим термином. Например, «Кавказ» и «горные районы». Наличие связи **BTI** подразумевает наличие связи **BT**.
A BTI B ó A BT B.
- **NTI** - *Narrower Term Instance*. Обращение связи **BTI**.
A NTI B ó B BTI A.
- **GS** - *Generic Structure*. Это иерархическая связь, используемая для визуального представления тезауруса. Она может не соответствовать структуре связей **BT/NT**. Эта связь используется потому, что визуальное представление полииерархической структуры, образуемой связями **BT/NT** затруднительно и ненаглядно.
- **USE+** - *Use ... and...* Связь один ко многим. Используется, когда для данного термина более предпочтительными является совокупность нескольких терминов. Например, «Угольные шахты» **USE+** «Уголь» and «Шахты».
- **UF+** - Обращение связи **USE+**.

Добавлены также атрибуты термина:

- **ID** - *Identifier*. Уникальный идентификатор термина.
- **HN** - *History Note*. История модификации связей и атрибутов данного термина.

В стандарте указаны следующие ограничения на структуру тезауруса:

- Из термина, не являющегося наиболее подходящим для какой-либо концепции, могут исходить только связи **USE** и **USE+**, а входить только связи **UF** и **UF+**. Никаких других связей этот термин иметь не может.
- Термин не может иметь связи с самим собою.
- Одна пара терминов не может иметь 2 или более связи (за исключением случаев, когда одна связь следует из другой по правилам стандарта).

Стандарт ANSI/NISO Z39.19-1993 помимо структуры регламентирует также и другие аспекты создания, представления и поддержки тезаурусов. Однако это выходит за рамки рассмотрения данной статьи.

Стандарт ГОСТ 7.25-2001 и Стандарт ГОСТ 7.24-90

Эти стандарты были созданы на базе вышеописанных стандартов ISO и ANSI, и, фактически, мало от них отличаются.

ГОСТ 7.25-2001 – Тезаурус информационно-поисковый одноязычный.

Этот стандарт устанавливает правила разработки, структуру, состав и форму представления информационных тезаурусов, ориентированных на использование лексики русского языка и разрабатываемых в рамках автоматизированных информационных систем и сетей научно-технической информации. ГОСТ 7.25-2001 также как и ANSI/NISO Z39.19-1993, расширяет и уточняет стандарт ISO 2788-1986 для одноязычных тезаурусов. В этом стандарте определены два типа терминов: Дескриптор и Аскриптор.

Таблица - Типы и значение отношения между Дескрипторами и Аскрипторами.

Тип ссылки	Обозначение на русском языке	Символьное обозначение	Значение ссылки	Аналог на английском языке
2 Ссылка от дескриптора к эквивалентному аскриптору	C	=	Синоним	UF (used for)
3 Ссылка от аскриптора к нескольким альтернативно заменяющим его дескриптам	Иа	=:	Используй альтернативно	Нет аналога

орам

4 Ссылка от Ик аскриптора к заменяющей его комбинации дескрипторо в	=+	Используй комбинацию	USE+ (ANI/NISO Z39.19-1993)
5 Ссылка от В дескриптора к вышестоящ ему дескриптору	<	Выше	BT (broader term)
6 Ссылка от Вр дескриптора к вышестоящ ему родовому дескриптору	:<	Выше-род	BTG (broader term generic)
7 Ссылка от Вц дескриптора к вышестоящ ему дескриптору, обозначающ ему целое	-<	Выше-целое	BTP (broader term partitive)
8 Ссылка от Н дескриптора к нижестоящ ему дескриптору	>	Ниже	NT (narrower term)
9 Ссылка от Нв дескриптора к нижестоящ ему видовому дескриптору	>:	Ниже-вид	NTG (narrower term generic)

10 Ссылка от Нч дескриптора к нижестоящему дескриптору, обозначающему часть	>-	Ниже-часть	NTP (narrower term partitive)
11 Ссылка от А дескриптора к ассоциативно связанному дескриптору	—	Ассоциация	RT (related term)
12 Ссылка от Са дескриптора к дескриптору, который заменяется данным дескриптором при альтернативном выборе (обратная ссылка к «иа»)	:=	Сравни альтернативный выбор	Нет аналога
13 Ссылка от Ск дескриптора к дескриптору, который заменяется комбинацией, включающей данный дескриптор (обратная ссылка к «ик»)	+=	Сравни комбинацию	UF+ (ANI/NISO Z39.19-1993)

14	Ср	:	Сравни	Нет аналога
Техническая обратная ссылка				
15	лп	/.../	Лексическое SN (scope примечание note)	
Уточнение значения и области применения				

Стандарт ГОСТ 7.24-90

Этот стандарт распространяется на многоязычные информационно-поисковые тезаурусы (МИПТ) и устанавливает состав, структуру и основные требования к построению МИПТ, применяемых в информационно-поисковых системах.

МИПТ – согласованная совокупность одноязычных информационно-поисковых тезаурусов, содержащая эквивалентные дескрипторы на языках - *компонентах МИПТ*, необходимые и достаточные для межъязыкового обмена, и включающая средства для указания их эквивалентности. *Одноязычной версией МИПТ* называют каждый из одноязычных тезаурусов, входящих в состав МИПТ. *Дескриптором МИПТ* называют совокупность эквивалентных дескрипторов одноязычных версий, связанных связями эквивалентности.

Существуют следующие виды эквивалентности терминов:

1. полная;
2. неполная (понятия, выражаемые терминами, пересекаются);
3. частичная (понятие, выражаемое одним термином, является часть понятия, выражаемого другим);

Допускается также использование вышеперечисленных связей для выражения эквивалентности вида “один ко многим”. В этом случае дескриптор на одном языке может быть связан с несколькими дескрипторами на другом языке.

При наличии в языках-компонентах полностью эквивалентных терминов они считаются представителями одного дескриптора МИПТ. При отсутствии в языках-компонентах полных эквивалентов для выражения одного и того же понятия в качестве дескриптора МИПТ в одноязычных версиях используют неполные и частичные эквивалентные дескрипторы. При этом к связям эквивалентности приписывают *реляторы* или *комментарии*, описывающие степень эквивалентности. Рекомендуются также и некоторые другие способы решения этой проблемы, малоприменимые для машинной реализации.

Особенности применения тезаурусов в информационных системах.

Модель данных

Описанные выше стандарты были разработаны для представления тезаурусов в виде, удобном для ручной индексации информационных ресурсов. Такая модель с ограничениями может быть также использована для машинной индексации с целью осуществления последующего поиска по ключевым словам.

Однако существует ряд тезаурусов, основная задача которых не индексация ресурсов, а их классификация. В этом случае основными объектами таких тезаурусов (*классификаторов*) выступают не термины, а понятия (*рубрики*), и, часто, идентифицирующие их уникальные идентификаторы (коды классификации). Отношения в таком тезаурусе – не семантические связи между терминами, а характеризующие логику описываемой предметной области отношения между понятиями (рубриками). Примерами таких тезаурусов могут служить тематические классификаторы в разных отраслях науки, например, MSC [13], PACS [14], DDC [15].

Структура классификатора соответствует структуре обычного тезауруса, поскольку связи между его рубриками по смыслу те же, что и между терминами тезауруса, и классификатор является его частным случаем. Однако при классификации в соответствие ресурсам ставятся не термины, а обозначаемые ими понятия. Потому в схеме данных информационной системы понятия тезауруса должны быть выделены в самостоятельные объекты. Это означает, что такая схема должна иметь структуру, отличную от вышеописанных стандартов, в которых понятия не выступают отдельными объектами, а есть лишь термины и связи между ними. В то же время, схема должна позволять работать с тезаурусами, описанными в соответствии с этими стандартами, т.е. быть совместима с ними.

Среди связей между терминами в вышеописанных стандартах следует различать связи, которые по смыслу характеризуют фактически соотношения не между терминами, а между термином, и обозначаемым им понятием. К таковым относятся связи **Use, UsedFor** в ISO и ANSI, и связи **Смотри (Use), Синоним (UF), Используй альтернативно, Используй комбинацию (Use+), Сравни альтернативный выбор, Сравни комбинацию(UF+)** в ГОСТ 7.25-2001. В схеме данных для информационной системы стоит ставить такие связи между понятиями и терминами, которые их обозначают.

Аналогично, иерархические и ассоциативные связи по смыслу являются связями между понятиями. Признак *TopTerm* также является признаком понятия, находящегося на вершине иерархии понятий.

Таким образом, получается следующее отображение связей между терминами в стандартах ISO и ANSI для одноязычных тезаурусов на отношения в схеме данных информационной системы: Те связи, которые допустимы между наиболее предпочтительными терминами (дескрипторами) для каких либо понятий, в схеме данных информационной системы становятся отношениями между понятиями. Те связи, которые были допустимы между наиболее предпочтительным термином (дескриптором) и другими терминами (аскрипторами) данного понятия, становятся отношениями между понятием и термином.

Как указывалось выше, в многоязычных тезаурусах термины имеют атрибут *язык*, на котором данный термин обозначает данное понятие. Кроме того, стандартами ISO 5964-1985 и ГОСТ 7.24-90 предусматривается ряд отношений эквивалентности между терминами на разных языках, допускающие, помимо строгой эквивалентности, несколько видов неполной эквивалентности терминов. По смыслу атрибут *язык* – свойство термина, а не понятия. В то же время термины на разных языках, между которыми есть только частичная эквивалентность, строго говоря, соответствуют разным, пусть и близким, понятиям.

Таким образом, более естественной в схеме данных тезауруса для информационных систем будет привязка языка к терминам, а не к понятиям. Более того, такой подход является единственно возможным для классификаторов, в которых именно независимые от языка понятия классифицируют другие ресурсы. Обычно такие классификаторы изначально создаются как одноязычные, и лишь потом для них делаются переводы на другие языки. В этом случае между терминами на разных языках имеет место только строгая эквивалентность, поскольку при переводе для каждого термина дается его строгий эквивалент (который является эквивалентом по определению, в контексте данного классификатора, даже если фактически перевод не совсем точен). Привязка языка к понятию означала бы необходимость делать отдельную копию одного и того же понятия для каждого языка, и делать отдельную связь между каждой копией понятия и классифицируемым им ресурсом. Привязка языка к термину привязать все эквивалентные термины на разных языках к одному и тому же понятию.

Однако в тезаурусах, где много отношений неполной эквивалентности между разноязычными терминами, а также имеются разные иерархии для терминов на разных языках, даже полностью эквивалентные термины могут оказаться в разных иерархиях, а значит, не могут быть привязаны к одному понятию. Все это означает, что для поддержки многоязычных тезаурусов схема данных должна предусматривать описанные в стандартах ISO и ГОСТ соотношения эквивалентности между терминами на разных языках, как отношения между понятиями. При этом для каждого тезауруса, в зависимости от его специфики, необходимо сделать выбор, каким образом реализовывать отношение полной эквивалентности между разными терминами:

1. Приписывать термины к разным понятиям, и ставить между понятиями отношение полной эквивалентности.
2. Приписывать термины к одному и тому же понятию.

Очевидно, для классификаторов необходимо использовать второй подход, а для многоязычных тезаурусов, имеющих разные иерархии на разных языках – первый. Следует заметить, что тезаурус, в котором есть отношение неполной эквивалентности, по смыслу уже подразумевает наличие разных иерархий на разных языках, а значит, необходим первый подход при их реализации.

Еще одним важным атрибутом термина в тезаурусе является комментарий к нему (*ScopeNote*). В тезаурусах-классификаторах, где, по сути, первично понятие, а не термин, комментарий, как правило, также характеризует понятие. Однако, в других тезаурусах комментарий может относиться именно к термину. Например,

описывать случаи предпочтительного употребления именно этого синонима перед другими. Таким образом, в разных тезаурусах комментарии могут относиться, как к понятиям, так и к терминам. Выбор зависит от конкретного тезауруса. Универсальная схема данных в информационной системе должна допускать оба варианта применения комментариев.

Платформа реализации тезауруса, требования SemanticWeb

Модель данных тезауруса, в том числе и учитывающая все перечисленные выше требования, может быть создана практически на любой платформе представления онтологии. В частности, существуют модели тезаурусов на основе TopicMaps [11], RDF [4, 6], DAML [5].

Однако для того, чтобы реализация тезауруса могла в полной мере соответствовать концепциям проекта SemanticWeb, на нее накладываются следующие требования:

1. *Синтаксическая и семантическая интероперабельность.* Любое приложение, работающее в соответствии с требованиями SemanticWeb должно иметь возможность работать с тезаурусом без предварительного согласования форматов.
2. *Расширяемость тезаурусов.* При необходимости любое приложение должно иметь возможность добавить в открытый тезаурус свои элементы и использовать его в таком расширенном виде для своих нужд.
3. *Расширяемость модели.* Схема данных должна допускать расширения и детализацию. То есть любое приложение должно иметь возможность добавить в модель новые типы ресурсов и связей, в частности детализировать уже существующие, если это, например, необходимо для описания нестандартного тезауруса. В то же время приложения, не знающие о таком расширении, должны иметь возможность корректно работать с этим тезаурусом в рамках прежней модели, имея доступ к той части данных тезауруса, которая в нее вписывается.

Эти требования накладывают ограничения и на платформы реализации тезауруса. Например, платформа TopicMaps [10] в формате XTM в целом удовлетворяет пунктам 1 и 2, но не удовлетворяет пункту 3. Наиболее соответствует перечисленным требованиям платформа RDF, а так же ее расширения (например, DAML+OIL) [25]. Платформа RDF принята также в качестве основной для описания онтологии в SemanticWeb.

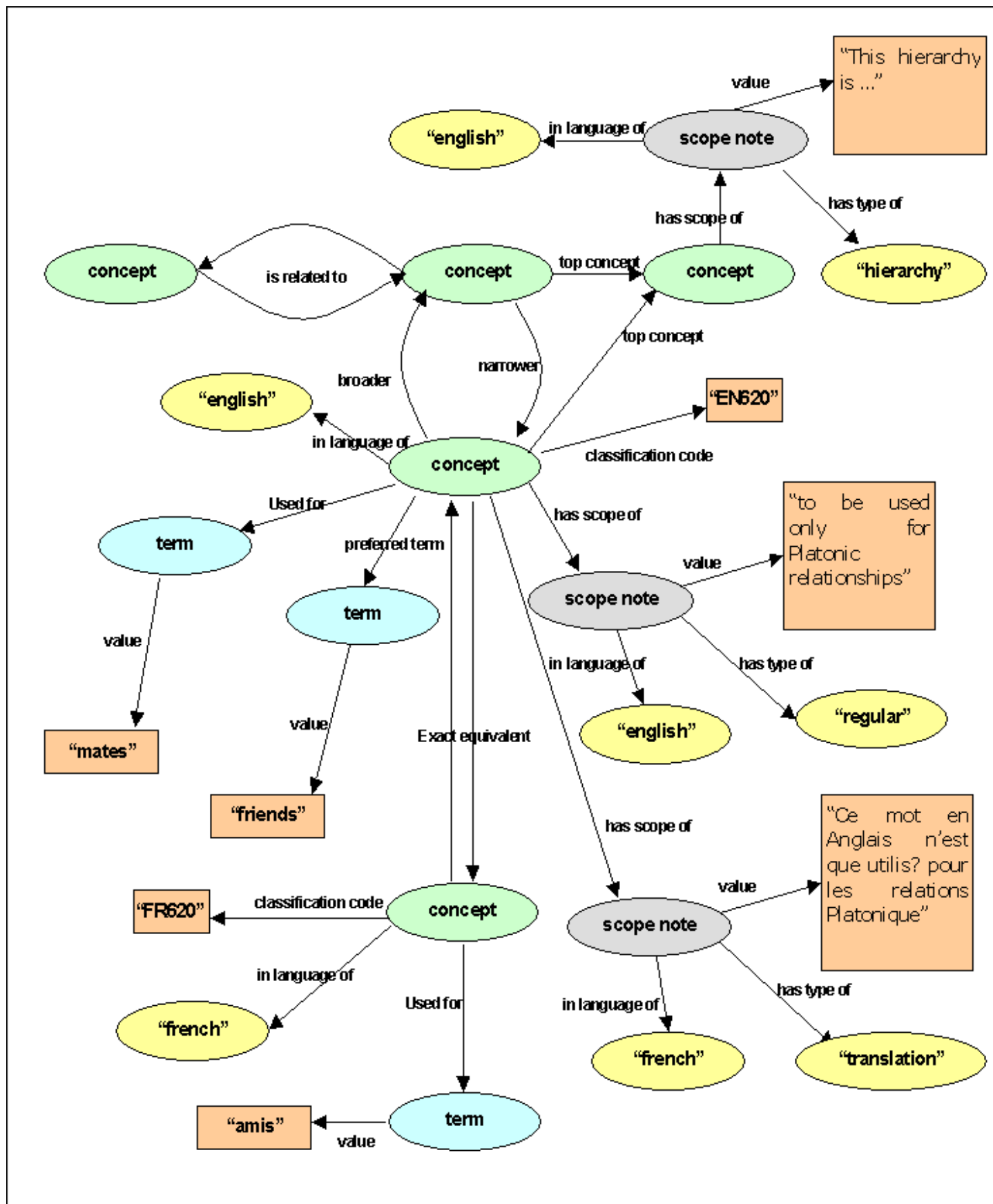
Подходы к описаниям тезаурусов

В этом разделе рассмотрены некоторые существующие схемы данных на основе RDF, предложенные в качестве стандартов для описания тезаурусов в информационных системах.

Формат представления многоязычного тезауруса в RDF, разработанный в рамках проекта LIMBER.

Данный формат изначально разрабатывался для многоязычного тезауруса ELSST

(European Language Social Science Thesaurus) [4]. Однако в настоящий момент LIMBER [3] предлагает данную модель как универсальную, для представления многоязычных тезаурусов.



Пример описания тезауруса в схеме данных LIMBER

Модель имеет следующие основные типы объектов (ресурсов):

- *Понятие (Concept)*
- *Термин (Term)*
- *Комментарий (ScopeNote)*
- *Язык (LanguageCode)*

Существуют следующие свойства понятий:

- *Уникальный идентификатор (ClassificationCode)*
- *Язык (inLanguageOf)*
- *Комментарий к понятию (hasScopeNote)*
- *Наиболее предпочтительный термин (PreferredTerm)*
- *Менее предпочтительный термин (UsedFor)*

Существуют следующие свойства комментариев:

- *Язык (inLanguageOf)*
- *Тип (hasTypeOf)*.

Существуют следующие типы комментариев:

- *General*. Комментарий к понятию на основном языке тезауруса (один из языков тезауруса в модели выделяется как основной или главный).
- *Translation*. Комментарий на неосновных языках.
- *Hierarchy*. Признак понятия, находящегося на вершине иерархии.
- *History*. Пометки об истории изменения этого понятия в предыдущих версиях тезауруса.

Существуют следующие связи между понятиями одного языка:

- *Более широкое понятие (BroaderConcept)*
- *Более узкое понятие (NarrowerConcept)*
- *Связанное понятие (isRelatedTo)*
- *Указатель на корневую концепцию данной иерархии (TopOfHierarchy)*

Существуют следующие связи между понятиями на разных языках:

- *Строгая эквивалентность (ExactEquivalent)*
- *Нестрогая эквивалентность (InexactEquivalent)*
- *Частичная эквивалентность (PartialEquivalent)*
- *Эквивалентность типа «один ко многим» (OneToManyEquivalent)*

Эта модель хорошо подходит для описания многоязычных тезаурусов, в которых существуют разные иерархии терминов на разных языках. Однако здесь язык термина является атрибутом понятия, а не термина. Как было описано выше, такая модель неудобна для описания многоязычных классификаторов ресурсов, в которых понятия семантически не связаны с каким-либо определенным языком.

Схема данных тезауруса ILRT

Эта схема данных строилась в расчете на работу не только с тезаурусами в обычном, «лингвистическом» смысле, но и с классификаторами. Потому язык термина привязан не к понятию, а к самому термину, а термины на разных языках, точно эквивалентные друг другу, привязаны к одному и тому же понятию. Термины на разных языках, не имеющие строгой эквивалентности, должны быть отнесены к разным понятиям.

Модель предполагает 2 уровня детализации описания тезауруса. Первый уровень реализует связи, предусмотренные стандартом ISO 2788-1986 для одноязычных тезаурусов, а также атрибут «язык» для терминов. Второй уровень детализации пока не оформлен в виде RDFS, и предполагает детализацию ряда связей 1 уровня детализации. Например, связь «более общее понятие» распадается на 3 RDF-связи, реализующие 3 описанных выше вида этой связи. Аналогично происходит детализация других связей.

По сути, эта схема предназначена для одноязычных тезаурусов и для тезаурусов-классификаторов, поскольку механизм полной поддержки многоязычных тезаурусов никак не прописан, а обозначено только направление, как это можно сделать в рамках данной модели.

Особенностью данной модели, в сравнении со предыдущей, является отсутствие избыточных связей оптимизирующих скорость исполнения запросов. Например, нет связи «более широкое понятие», поскольку оно является обращением связи «более узкое понятие». Отсутствует также связь понятий с самыми верхними понятиями включающих их иерархий, поскольку она тоже вычисляется из иерархических связей. Это накладывает дополнительные ограничения на техническую реализацию такой модели. В частности, традиционные способы реализации графов не позволяют за один шаг вычислить корневую вершину иерархии для произвольного понятия.

Модель тезауруса DRC

Эта модель наиболее точно соответствует модели одноязычного тезауруса ISO 2788-1986. В частности, в нем отсутствует класс понятий, и все связи существуют только между терминами. Некоторые связи детализированы, в частности выделены разные виды связей менее предпочтительными терминами. Модель реализована на языке DAML [16].

Стоит выделить одну явную ошибку этой модели. Связь *RelatedTerm* является транзитивной, что не соответствует действительности. Например, связанными терминами являются *транспортировка нефти* и *трубы для нефтепроводов*, а также *трубы для нефтепроводов* и *стальной прокат*. Однако прямой связи между понятиями *транспортировка нефти* и *стальной прокат*, очевидно, нет [4].

Поскольку в модели нет понятий, как отдельных объектов, она не удобна для реализации классификаторов.

Интерфейсы работы с тезаурусом в информационных системах

Просмотр тезауруса и поиск ресурсов

В информационной системе тезаурус является не только самостоятельным информационным ресурсом, но и инструментом для классификации или индексации ресурсов. Таким образом, пользователь информационной системы должен иметь возможность:

- Осуществлять просмотр тезауруса.
- Осуществлять поиск ресурсов по ассоциированным с ними терминам или понятиям.

Поиск ресурсов может вестись двумя способами:

- Поиск по ключевым словам, используя тезаурус.
- Навигация по тезаурусу. То есть поиск сначала нужного понятия в тезаурусе с последующим запросом ресурсов, соответствующих этому понятию.

При поиске ресурсов по ключевым словам поисковая система может, используя тезаурус, расширять результаты поиска, выдавая пользователю не только ресурсы, соответствующие введенным пользователем ключевым словам, но и ресурсы, соответствующие связанным с ними терминам, или терминам, обозначающим также более узкие понятия относительно исходного термина. Например, если пользователь ищет ресурсы, соответствующие термину «туннель», в результатах поиска необходимо выдать также все ресурсы, соответствующие термину «тоннель», поскольку оба они являются разными правильными вариантами написания одного и того же слова. Или если ищутся ресурсы, соответствующие понятию *дифференциальные и функциональные уравнения*, имеет смысл включить в результаты поиска также ресурсы, соответствующие рубрике *системы функциональных уравнений и неравенства*. Система поиска может также, используя тезаурус, подсказать пользователю, по каким еще словам ему стоит осуществить поиск (например, квазисинонимы, связанные термины, более широкие термины, и т.д.). Оба этих варианта использования тезауруса широко применяются, например, в поисковых машинах Internet.

Интерфейс просмотра тезауруса должен:

- Показывать все атрибуты данного термина или понятия.
- Показывать, с какими терминами и понятиями связан данный термин или понятие.
- Достаточно наглядно показывать пользователю место термина или понятия в иерархии понятий тезауруса.

Первые 2 пункта выполнимы, если показывать пользователю для каждого понятия тезауруса на отдельном экране (странице) все его атрибуты, все связанные с ним термины (на всех или на определенном языке), и все связанные с ним понятия. Интерфейс должен при этом обеспечивать переход к странице просмотра любого из перечисленных на данной странице понятий. Если в тезаурусе схемой данных

разрешена привязка термина более чем к одному понятию, на той же странице для каждого термина должны быть перечислены также понятия, к которым еще привязан данный термин. Если у понятия есть термины на других языках, не полностью эквивалентные данному понятию, или полностью эквивалентные, но прикрепленные в силу структуры данного тезауруса к другим понятиям, на странице должны присутствовать ссылки на страницы этих понятий.

Наглядно показать пользователю место термина или понятия в тезаурусе достаточно сложно, поскольку достаточно наглядное отображение полииерархической структуры на одной странице, в отличие от иерархии, довольно сложно, как для отображения, так и для восприятия пользователем. В частности, в общем случае невозможно будет обойтись без пересекающихся линий, показывающих иерархические связи между понятиями. Потому имеет смысл показать только часть понятий и связей, которая, с одной стороны, была бы легко отображаемой и воспринимаемой, и в то же время достаточно наглядно показывала бы место понятия в общей иерархии понятий.

Если тезаурус имеет строго древовидную структуру, то представление дерева обычно осуществляется следующими способами:

1. Визуализация пути по дереву от корня к текущему элементу. Например, [\[18\]](#).
2. Визуализация пути по дереву от корня к текущему элементу, а также соседей каждого предка текущего элемента. Например, [\[17\]](#)
3. Визуализация всего дерева целиком. Обычно в таких случаях пользователь может открывать и закрывать отображение на экране потомков любых узлов. Например, программа «Проводник» («Explorer») в операционных системах MicrosoftWindows.

Чтобы обеспечить эффективную выборку (одним запросом) требуемых в первых двух случаях разрезов иерархических структур, представляемых рекурсивной связью между узлами этих структур, соответствующие таблицы БД расширяются вспомогательными столбцами и условиями целостности.

В случае полииерархической структуры первые 2 из вышеописанных способов также могут быть применены. Но в этом случае необходимо задать путь от корня полииерархии к текущей вершине, по которому будет произведена визуализация. При этом известные алгоритмы визуализации дерева одним запросом к реляционной БД неприменимы. Однако максимальное количество запросов к БД в этом случае не будет велико. Оно будет равно максимальной длине пути по полииерархии тезауруса, которая, как правило, сопоставима с логарифмом от общего количества понятий тезауруса. Это вполне приемлемо для информационной системы. Примером реализации такого подхода может служить [\[18\]](#) (см. интерфейс добавления сайта в каталог).

Еще один вариант отображения полииерархии – построение остового дерева, и отображение его вышеописанными способами. В этом случае для каждого элемента тезауруса необходимо выделить из всех его предков одного, связь с которым и станет связью остового дерева (см., например, [\[17\]](#)). В некоторых

тезаурусах заложено решение этой проблемы именно таким путем. Специальная связь *GenericStructure* используется для организации понятий в обычную древовидную иерархию. При этом значимого семантического смысла эта связь не имеет и служит лишь для отображения в интерфейсах и в печатном виде тезауруса в качестве иерархической структуры [2].

Возможно также построение для визуализации полного дерева путей по полииерархии тезауруса. Однако размер такого дерева может оказаться недопустимо большим. Например, в случае полииерархии типа «сетка рабица» (каждый элемент, кроме крайних, имеет ровно по 2 предка и по 2 потомка), количество копий каждого такого элемента будет равно степени двойки, в показателе которой стоит глубина пути к этому элементу от корня (или от предка - крайнего элемента). Это значит, что количество элементов такого дерева, будет расти экспоненциально с ростом количества узлов сетки. Это недопустимо для информационной системы.

Как уже упоминалось. Для первого и второго варианта визуализации окружения текущего элемента полииерархии необходим путь, по которому должна осуществляться визуализация. Если пользователь пришел к данному понятию посредством навигации от корня иерархии, то визуализацию следует осуществлять в соответствии с тем путем, по которому он пришел. Однако если пользователь пришел к просмотру данного элемента другим способом (например, из поискового интерфейса), и путь (или верхняя часть пути) не известен, его (или неизвестную его часть) можно либо вообще не отображать (см., например, [18]), либо вычислять какой-либо путь по умолчанию, например, самый левый, и отображать именно его. Для этого так же можно генерировать и использовать остовое дерево.

Еще один вариант отображения положения элемента в полииерархии, который будет, вероятно, полезен для пользователя – визуализация всех соседей всех его непосредственных предков. Это будет, по сути, двухмерная таблица, легко отображаемая на экране.

Администрирование тезауруса

Интерфейсы администрирования тезауруса должны обеспечивать выполнение следующих операций:

- *Добавить новое понятие к тезаурусу.* При добавлении добавляется так же связь с некоторым другим уже существующем в тезаурусе понятием. Указывается тип этой связи.
- *Добавить связь определенного типа между понятиями.* Должно обеспечиваться ограничение: не более одной связи между двумя понятиями. При добавлении иерархической или ассоциативной связи добавляется так же парная к ней обратная связь (A **BT** B ó B **NT** A ; A **RT** B ó B **RT** A).
- *Изменить тип связи между понятиями.* Должно обеспечиваться ограничение: Связь RT запрещена между понятиями, одно из которых является предком другого.

- *Удалить понятие и все его связи.* При удалении понятия все его потомки, не имеющие других предков, могут либо удаляться вместе с ним, либо выделяться в отдельную иерархию.
- *Удалить связь между понятиями.* При удалении иерархической связи понятие-потомок и все его потомки, не имеющие других предков, могут либо удаляться вместе с ним, либо выделяться в отдельную иерархию.
- *Добавить/изменить наиболее подходящий термин для данного понятия на некотором языке.* Должно обеспечиваться ограничение: Для каждого понятия не более одного наиболее подходящего термина на каждом языке.
- *Добавить/изменить менее подходящий термин для данного понятия на некотором языке.* При добавлении добавляется также связь к этому термину и указывается тип этой связи.
- *Добавить связь определенного типа между термином и понятием.* Должно обеспечиваться ограничение: для каждого термина не более одной связи с одним и тем же понятием.
- *Изменить тип связи между термином и понятием.*
- *Добавить/изменить комментарий к связи между термином и понятием на некотором языке.*
- *Удалить термин и все его связи.*
- *Удалить связь между термином и понятием.* Если термин не имеет других связей, он также удаляется.
- *Изменить код (идентификатор) понятия.*
- *Изменить код (идентификатор) термина.*
- *Добавить/изменить комментарий к понятию.* Должно обеспечиваться ограничение: не более одного комментария к одному понятию на одном языке.
- *Добавить/изменить комментарий к термину.* Должно обеспечиваться ограничение: не более одного комментария к одному термину на одном языке.

Интерфейсы администрирования должны включать и использовать интерфейсы просмотра тезауруса для поиска тех понятий, терминов, комментариев и связей, которые должны быть изменены. Интерфейсы редактирования могут быть также частично интегрированы в интерфейсы просмотра (в виде добавленных органов управления в окна просмотра).

Подход к описанию тезауруса в ИСИР

Формулировка задачи

Для информационной системы ИСИР реализация тезауруса должна удовлетворять следующим свойствам:

1. Позволять хранить любые существующие тезаурусы, в частности, любые классификаторы, имеющие структуру тезауруса в соответствии со стандартами ISO 2788-1986, 5964-1985, ГОСТ-7.25-2001 и ГОСТ-7.24-90. В том числе, реализация должна позволять работать с многоязычными тезаурусами.
2. Позволять, используя тезаурус, индексировать ресурсы терминами данного

тезауруса, а также классифицировать ресурсы понятиями тезаурусов-классификаторов. При этом работа с обоими видами тезаурусов должна осуществляться единообразно.

3. Позволять осуществлять просмотр (навигацию) по тезаурусу, а также поиск ресурсов, проиндексированных или классифицированных тезаурусом. То есть реализация должна обеспечивать эффективное выполнение необходимых для этого запросов, а именно:
 - Получить значение атрибутов понятия.
 - Получить все понятия, связанные с данным понятием, связями заданных видов (для связей в соответствии со стандартами ISO, ГОСТ или их детализаций).
 - Получить самые верхние понятия в иерархии понятий, в которую входит данное понятие.
 - Получить все термины, связанные с данным понятием, связями заданных видов (для связей в соответствии со стандартами ISO, ГОСТ или их детализаций).
 - Получить все термины на данном языке, связанные с данным понятием связями заданных видов (для связей в соответствии со стандартами ISO, ГОСТ или их детализаций).
 - Получить все термины на данном языке, связанные связями заданных видов (для связей в соответствии со стандартами ISO, ГОСТ или их детализаций) с данным понятием, или с понятиями, связанными с данным понятием данными связями эквивалентности терминов на разных языках.
 - Получить значение атрибутов термина.
 - Получить все термины, содержащие данное слово (или ключевое слово).
 - Получить полииерархию понятий тезауруса с отображением некоторого пути, а так же всех соседей всех узлов этого пути в виде дерева.
 - Получить полный список терминов тезауруса.
4. Позволять осуществлять администрирование тезауруса в соответствии с требованиями, описанными в предыдущей части статьи.
5. Быть расширяемой. То есть допускать детализацию при необходимости некоторых связей, а так же добавление новых связей.

Описание схемы данных тезауруса

Эта схема данных основана на платформе RDFS. Приведенный подход к описанию тезауруса опирается на работы [19]. При описании тезауруса мы ориентировались на стандарты ГОСТ, потому что эти стандарты являются, фактически, расширениями стандартов ISO и ANSI/NISO.

Как было описано выше, по ГОСТ 7.25-2001 в тезаурусе существуют два типа терминов: дескриптор и аскриптор (в англоязычных стандартах *preferredandnonpreferredterms*). Между дескрипторами существуют следующие виды отношений: иерархические, ассоциативные, эквивалентные. Отношения дескриптора с аскрипторами бывают: синонимические, «сравни альтернативный выбор», «сравни комбинацию». Отношения аскриптора с дескрипторами могут

быть, соответственно: «смотри», «используй альтернативно», «используй комбинацию».

Исходя из описанных в предыдущей части статьи соображений, схема содержит два основных класса объектов: *ThesaurusConcept* – понятие которое выражается дескриптором, *ThesaurusTerm* – понятие которое выражается аскриптором. Тогда отношения между дескрипторами в ГОСТ 7.25-2001 станут отношениями между *ThesaurusConcept*, отношения между дескрипторами и аскрипторами станут отношениями между *ThesaurusConcept* и *ThesaurusTerm*. В данной схеме реализованы два варианта для установления отношения полной эквивалентности между разными терминами. Приписывать полные разноязычные эквивалентные термины к одному и тому же понятию. Привязка всех полных эквивалентных терминов на разных языках к одному и тому же понятию осуществляется только в том случае, когда все термины на разных языках имеют общую иерархию (это необходимо в классификаторах, а так же возможно в других тезаурусах, где иерархии терминов на разных языках совпадают). В противном случае необходимо привязывать термины к разным понятиям, и ставить между понятиями отношение полной эквивалентности.

Классы модели

- *ThesaurusConcept*. Понятие (дескриптор, *preferredterm*) .

Этот класс имеет следующие атрибуты:

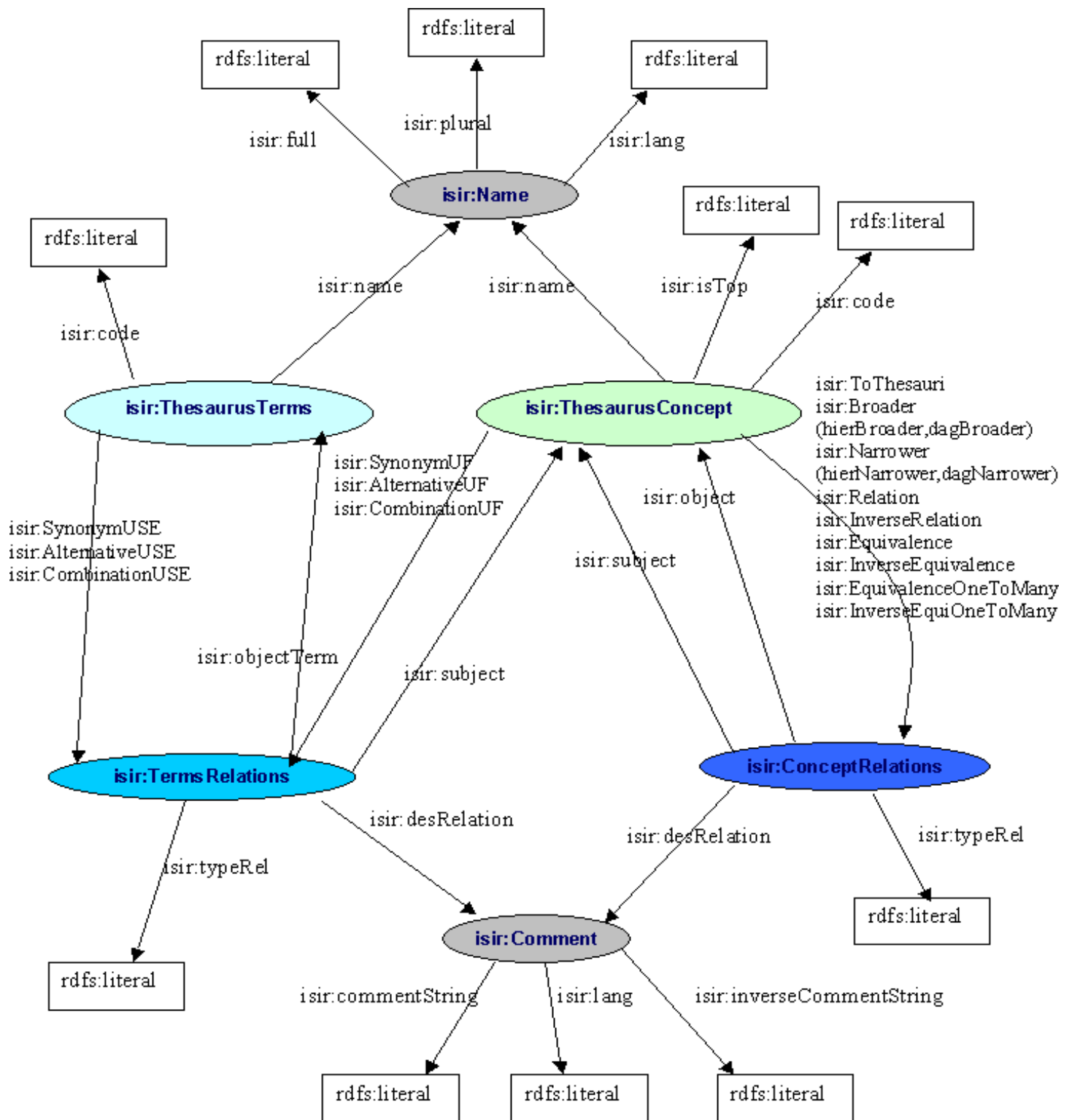
- *code*. Уникальный идентификатор, или код рубрики классификатора. Необязательный атрибут. Этот атрибут присутствует только тогда, когда он имеет смысловую нагрузку в тезаурусе, и не заменяет внутренний системный или технический ID записи в базе данных.
- *isTop*. Признак используется для того, чтобы определить является ли данное понятие самым верхним в иерархии.
- *name*. Термин, обозначающий данное понятие. Значение этого атрибута является экземпляром класса ***Name***.
- *toThesauri*. Связывается с тезаурусом.
- *broader*. Связь с более широким (более общим) понятием. Значение этой связи является экземпляром класса ***ThesaurusConcept***. Эта связь имеет два подвида:
 - *hierbroader*. Для установления иерархических связей между понятиями.
 - *dagbroader*. Для тех случаев, когда одно понятие имеет более одного непосредственного предка. Только один предок связывается с данным понятием через *hierbroader*. Остальные связываются через *dagbroader*.
- *narrower*. Связь с более узким понятием. Значение этой связи является экземпляром класса ***ThesaurusConcept***. Эта связь имеет два подвида, обратных соответствующим подвидам связи *broader*:
 - *hiernarrower*.
 - *dagnarrower*.
- *relation*. Связь с ассоциативным понятием. Значение этой связи является экземпляром класса ***ThesaurusConcept***. (*inverseRel* – обратная связь к этой связи)

- synonymUF. Связь с синонимическим термином. Значение этой связи является экземпляром класса **ThesaurusTerm**. (synonymUSE - обратная связь к этой связи)
 - alternativeUF. Значение этой связи является экземпляром класса **ThesaurusTerm**. (обратная связь к alternativeUSE)
 - combinationUF. Значение этой связи является экземпляром класса **ThesaurusTerm**. (обратная связь к combinationUSE)
 - equivalence. Связь с эквивалентным понятием, у которого есть термины на других языках. Эта связь используется в случае, когда у многоязычного тезауруса на каждом языке есть своя собственная полииерархия терминов. Значение этого атрибута является экземпляром класса **ThesaurusConcept**. (inverseEquivalence - обратная связь к этой связи). Вид эквивалентности указывается в комментарии (*Полная эквивалентность, Неполная эквивалентность (значения терминов не совпадают, но пересекаются), Частичная эквивалентность (значение одного термина шире, чем значение другого)*).
 - equivalenceOneToMany. Эта связь используется для установления отношения один ко многим. Значение этого атрибута является экземпляром класса **ThesaurusConcept**. (inverseEquiOneToMany - обратная связь к этой связи)
- **ThesaurusTerm**. Понятие соответствует аскриптору (*nonpreferredterm*). Этот класс имеет следующие атрибуты:
 - *code*. Уникальный идентификатор, или код рубрики классификатора. Необязательный атрибут. Этот атрибут присутствует только тогда, когда он имеет смысловую нагрузку в тезаурусе, и не заменяет внутренний системный или технический ID записи в базе данных.
 - *name*. Термин, который обозначает данное понятие. Значение этого атрибута является экземпляром класса **Name**.
 - synonymUSE. Обратная связь к synonymUF. Значение этой связи является экземпляром класса **Thesaurusconcept**.
 - alternativeUSE. Связь с альтернативным понятием (дескриптором). Значение этой связи является экземпляром класса **Thesaurusconcept**.
 - combinationUSE. Связь с комбинационным понятием (дескриптором). Значение этой связи является экземпляром класса **Thesaurusconcept**.
- **Name**. Термин. Имеет следующие атрибуты:
 - *full*. Написание (наименование) термина на данном языке.
 - *Plural*. Комментарий к этому термину на данном языке.
 - *Lang*. Язык термина
- **ConceptsRelations**. Реализует отношение (связь) между объектами **ThesaurusConcept**, снабженное комментарием.
 - *desRelation*. Комментарий к отношению. Значение этого атрибута является экземпляром класса **Comment**.
 - *subject*. Субъектсвязи (**ThesaurusConcept**).

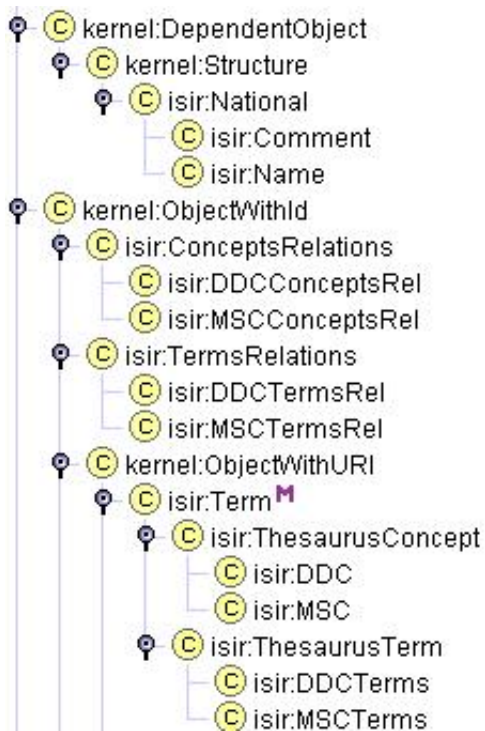
- object. Объектсвязи (***ThesaurusConcept***).
- ***TermsRelations***. Реализует отношение (связь) между ***ThesaurusConcept*** и ***ThesaurusTerm***, снабженное комментарием.
 - *desRelation*. Комментарий к отношению. Значение этого атрибута является экземпляром класса ***Comment***.
 - subject. Субъектсвязи (***ThesaurusConcept***).
 - objectTerm. Объектсвязи (***ThesaurusTerm***)
- ***Comment***. Комментарий к отношению.
 - *commentString*. Значение комментария на данном языке.
 - *InverseCommentString*. Комментарий к обратной связи на данном языке.
 - *Lang*. Язык комментария.

Ряд тезаурусов, например, математический классификатор MSC, имеют отношения между понятиями, которые нельзя отнести строго к какому-либо из определенных в стандарте типов, либо такое отношение требует уточнения (пример см. ниже). Как правило, таких отношений в тезаурусе очень мало, а потому нецелесообразно для них выделять отдельные типы отношений. Средством описания таких отношений может стать приписывание такого отношения к одному из базовых существующих типов, с добавлением к нему комментария, характеризующего его особенности. В данной модели тезауруса любое отношение между двумя понятиями (дескрипторами), или между понятием и термином (аскриптором) может быть снабжено комментарием на любом языке.

Данный подход позволит также минимизировать неминуемое дальнейшее расширение и детализацию наборов связей между терминами или понятиями, которая сейчас наблюдается в различных моделях и национальных стандартах (например, ANSI, ГОСТ), поскольку, как альтернативу детализации, можно использовать комментарии к связям специального вида. На рисунке ниже наглядно показана RDF схема тезауруса



RDF-схема тезауруса в ИСИР



Классы модели в системе Protégé OntoVizplugin

Ниже приведен пример описания одной из рубрик классификатора MSC на официальном сайте MSC <http://www.ams.org/msc>, а так же графически отражено представление этой рубрики и ее связей на двух языках в данной модели. В целях экономии места из всех более узких понятий данной рубрики в реальном MSC здесь оставлено только три.

Term 20Gxx Linear algebraic groups (classical groups)

For arithmetic theory, see 11E57, 11H56

NT 20G05 Representation theory

Здесь связь (**Forarithmetictheory**) между 20Gxx и 11E57, 11H56 является ассоциативной связью (тип связи RT), снабженной комментарием: **Forarithmetictheory**. На Рис 6 показаны RDF данные этого фрагмента MSC в предложенной схеме

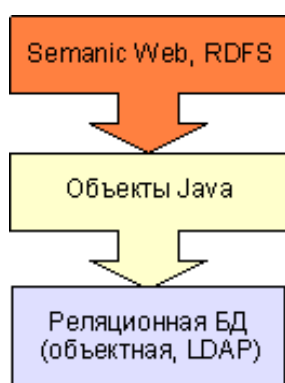
"**A** \$! **T** : **A Preferred Term T** (единственность для каждого языка).
A Related Concept B δ **B** не является ни предком, ни потомком **A** .

В дополнение к этим, модель имеет еще следующие ограничения, вытекающие из предыдущих рассуждений:

" **A** \$ **B** : **A Top Concept B** .
\$ **A** : **B Broader Concept A** δ *IsTop* (**B**) = false.

Реализация тезауруса в ИСИР

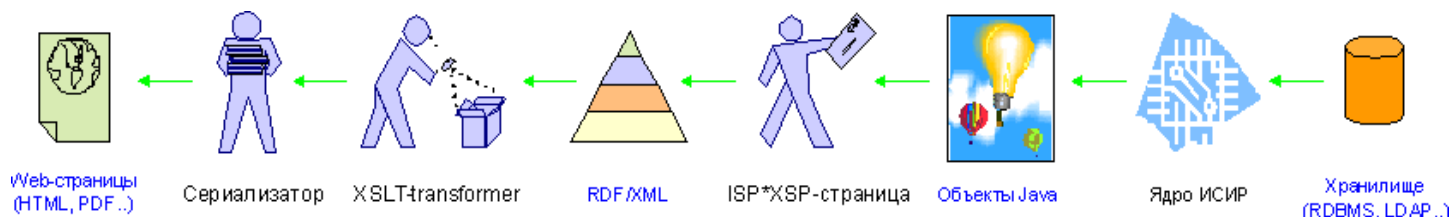
Созданные компоненты реализации тезауруса позволяют просматривать, редактировать и добавлять данные тезауруса в систему через web-формы а также импортировать и экспортировать данные в систему в специальном XML-формате.



Реализация тезауруса выполнена в соответствии с технологией Интегрированной Системы Информационных Ресурсов (ИСИР) [20, 21, 22]. Эта технология имеет следующие особенности архитектуры. Компоненты системы, реализующие бизнес-логику приложения, работают с объектной моделью данных на основе RDF/RDFS. Это достигается с помощью специальных компонентов, осуществляющих отображение объектной модели в реляционную базу данных. При этом хранимые информационные объекты, которыми оперируют компоненты бизнес-логики, реализованы через *JavaBean* классы, отображаемые в соответствующие таблицы СУБД.

Для создания форм редактирования данных тезауруса был использован сервис ИСИР *FormBuilder* [24]. Данный сервис предназначен для автоматизации процесса построения форм редактирования ресурсов, их заполнения начальными данными по полям редактируемых объектов и обработки данных, получаемых из этих форм. Сервис представляет собою набор специальных тегов, с помощью которых создаются JSP-страницы форм редактирования.

Для представления данных тезауруса реализация использует «XML-серверные страницы ИСИР (ISP*XSP)» [23]. ISP*XSP отвечают за динамическую генерацию RDF/XML-документов, то есть за выборку необходимой объектной информации из хранилища при помощи Ядра ИСИР и представление её в RDF/XML-виде. Дальнейшее преобразование информации в необходимое представление (например, в виде HTML для пользователя Internet) осуществляется с помощью XSLT-фильтров.



Язык «XML-серверных страниц» в ИСIP позволяет максимально наглядно и декларативно описать, какие данные следует выбрать. Такая страница похожа на шаблон XML-документа, в который осталось лишь подставить значения из хранилища.

Существует 2 способа функционирования тезауруса в системе:

1. Тезаурус встроен в информационную систему и является ее неотъемлемой частью.
2. Тезаурус является хранимым ресурсом информационной системы.

Во втором случае для просмотра и редактирования всех таких тезаурусов, а также его использования их для классификации и индексации других ресурсов, используются единые интерфейсы (например, предлагаемые в данной реализации). Тезаурусы могут динамически в процессе работы системы добавляться и удаляться из нее.

В первом случае хранение тезауруса в базе данных системы осуществляется так же как и во втором случае, или в отдельных таблицах. Кроме того, для конкретного встроенного тезауруса могут быть созданы отдельные Java-классы – наследники классов, используемых для тезаурусов - хранимых ресурсов. Это позволит при разработке информационной системы для разных встроенных тезаурусов использовать при необходимости разные интерфейсы, например, упрощенные для простых иерархических рубрикаторов.

Краткое описание пользовательских интерфейсов

Главная страница

<i>Thesaurus list</i>					
ru.ccas.isir.core.DDC					
Edit	View	Find	Import	Export	
ru.ccas.isir.core.MSC					
Edit	View	Find	Import	Export	

На этой странице перечисляются все функции реализации: редактирование,

представление, поиск, импорт, экспорт.

Форма редактирования тезауруса.

Resource editing

Data language: english Choose

Mathematics Subject Classification > Mathematical logic and foundations

Code: 03Cxx

Term: Model theory

Description:

Update Reset

Enabled sections

- Common info
- Synonym
- Alternative
- Combination
- BT
- NT
- RT
- Equivalence
- Equivalence one to many
- Top-concepts
- To terms structure
- Find
- View

Меню слева содержит такие пункты:

- *Common info*: для редактирования данного понятия (код, термин, комментарий)
- *Synonym, Alternative, Combination, BT, NT, RT, Equivalence, Equivalence one to many*: для добавления и редактирования понятия, связанного с данным понятием через эти связи.
- *TopConcepts*: для добавления и редактирования самых верхних понятий тезауруса.
- *ToTermsStructure*: Просмотр понятий тезауруса в виде иерархии.
- *Find*: Поиск понятия тезауруса.
- *View*: Просмотр понятий тезауруса

Форма просмотра понятий тезауруса

01-xx
History and biography

BT: (MSC) Mathematics Subject Classification

NT: (01-00) General reference works (handbooks, dictionaries, bibliographies, etc.)

NT: (01-01) Instructional exposition (textbooks, tutorial papers, etc.)

NT: (01-02) Research exposition (monographs, survey articles)

NT: (01-06) Proceedings, conferences, collections, etc.

NT: (01-08) Computational methods

NT: (01Axx) History of mathematics and mathematicians

Enabled sections

- Common info
- Synonym
- Alternative
- Combination
- BT
- NT
- RT
- Equivalence
- Equivalence one to many
- Top-concepts
- To terms structure
- Find
- Edit

Меню слева содержит пункты соответствующие разделу редактирования. При выборе какого-либо пункта показываются соответствующие понятия.

Представление терминов тезауруса в виде иерархии.

Мы можем осуществлять навигацию по тезаурусу и смотреть термины тезауруса. Если в тезаурусе существуют другие связи кроме иерархических, то для открытых и листовых понятий, кроме понятий, расположенных выше, ниже и рядом с данным понятием, показываются все другие понятия, связанные с данным понятием. Например, на показанной ниже странице для открытого понятия «05-xx» показывается, что оно связано с понятием «11Exx» через связь «RT (*Forfinitefields*)».



Форма поиска ресурсов по классификатору

FIND PUBLICATION

Request	Searching result
Publication title:	<input type="text"/>
Publication type:	<input type="text" value="▼"/>
Language:	<input type="text" value="▼"/>
DDC:	<input type="text"/> Browse
MSC:	<input type="text"/> Browse
<input type="button" value="Search"/> <input type="button" value="Clear"/>	

В этой форме поисковый запрос можно ограничивать определенными рубриками классификатора (понятиями тезауруса). Для выбора рубрик необходимо нажать на ссылку Browse, которая соответствует классификатору. Появится новое рорир-окно. В этом окне понятия тезауруса показываются в виде иерархии (или полииерархии) с элементами checkbox, чтобы легко выбирать нужные рубрики.

Ниже показан пример для классификатора MSC.

Mathematics Subject Classification

(MSC) Mathematics Subject Classification

- (00-xx) General
- (01-xx) History and biography
- (03-xx) Mathematical logic and foundations
- (05-xx) Combinatorics
- (06-xx) Order, lattices, ordered algebraic structures
- (08-xx) General algebraic systems
- (11-xx) Number theory
- (12-xx) Field theory and polynomials
- (13-xx) Commutative rings and algebras

...

- (86-xx) Geophysics
 - See also: (76U05) Rotating fluids
 - See also: (76V05) Reaction effects in flows
- (90-xx) Operations research, mathematical programming
- (91-xx) Game theory, economics, social and behavioral sciences
- (92-xx) Biology and other natural sciences
- (93-xx) Systems theory; control
- (94-xx) Information and communication, circuits
- (97-xx) Mathematics education

Choose

Close

Импорт и экспорт тезаурусов.

Для повторного использования и связи с другими системами в системе реализованы средства импорта и экспорта тезаурусов. Данные тезаурусов экспортируются и импортируются в систему через специальный формат XML. В настоящий момент в систему загружены два классификатора: DDC на двух языках (вьетнамском и английском, 1000 терминов), и MSC на английском языке (2000 терминов).

Ниже показан фрагмент тезауруса DDC, записанный в этом формате. Данные DDC представлены на двух языках (английском и вьетнамском).

```
<?xml version="1.0" encoding="utf-8" ?> <rdf:RDF xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance" xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#" xmlns:rdfs="http://www.w3.org/2000/01/rdf-schema#" xmlns:schema="urn:hdl:1016.1/schema/" xmlns:isir="urn:hdl:1016.1/core/" xmlns:kernel="urn:hdl:1016.1/kernel/" xmlns:dc="http://purl.org/dc/elements/1.1/" xmlns:dctype="http://purl.org/dc/d
```



```

cmitype/" xmlns:dcterms="http://purl.org/dc/terms/"> <isir:DDC rdf
:about="urn:hdl:1016.1/core/DDC#DDC"> <isir:name rdf:parseType="Res
ource"> <isir:full>Dewey Decimal Classification</isir:full> <dc:la
nguage>en</dc:language> </isir:name> <isir:name rdf:parseType="Res
ource"> <isir:full>Khung Phân Lo?i Th?p Phân</isir:full> <dc:langu
age>vi</dc:language> </isir:name> <isir:code>DDC</isir:code> <isi
r:isTop>top</isir:isTop> </isir:DDC> <isir:DDC rdf:about="urn:hdl:
1016.1/core/DDC#000"> <isir:name rdf:parseType="Resource"> <isir:f
ull>T?ng Quát</isir:full> <dc:language>vi</dc:language> </isir:nam
e> <isir:name rdf:parseType="Resource"> <isir:full>Generalities</i
sir:full> <dc:language>en</dc:language> </isir:name> <isir:code>0
00</isir:code> <isir:toThesauri rdf:about="urn:hdl:1016.1/core/DDC#
DDC"/> <isir:hiernarrower rdf:parseType="Resource"> <rdf:type rdf:
resource="urn:hdl:1016.1/core/DDCConceptsRel" /> <isir:typeRel>hier
</isir:typeRel> <isir:desRelation rdf:parseType="Resource"> <isir:
commentString>NT</isir:commentString> <isir:inverseCommentString>BT
</isir:inverseCommentString> <dc:language>en</dc:language> </isir:
desRelation> <isir:subject rdf:about="urn:hdl:1016.1/core/DDC#?DDC"
/> <isir:object rdf:about="urn:hdl:1016.1/core/DDC#000"/> </isir:h
iernarrower> </isir:DDC> . . . </rdf:RDF>

```

Литература

[1] Thesaurus Construction

<http://instruct.uwo.ca/gplis/677/thesaur/main00.htm>

[2] Thesaurus Format: Nu search Standard Specification

http://www.excavio.com/pdf/wp_nusearch_thesaurus_spec.pdf

[3] LIMBER (Language Independent Metadata Browsing of European Resources) project:

<http://www.limber.rl.ac.uk/>

[4] A Thesaurus Interchange Format in RDF

http://www.limber.rl.ac.uk/External/SW_conf_thes_paper.htm

[5] Hall, M. (2001) *CALL Thesaurus Ontology in DAML*.

<http://orlando.drc.com/daml/ontology/Thesaurus/CALL/>

[6] RDF Thesaurus Specification

<http://ilrt.org/discovery/2001/01/rdf-thes/>

[7] Web Thesaurus Compendium

<Http://www.darmstadt.gmd.de/~lutes/thesoecd.html>

[8] ISO2788: Guidelines for establishment and development of monolingual thesauri, 2nd ed., Geneva: ISO1986.

[9] ISO5964: Guidelines for establishment and development of multilingual thesauri, 1st

ed., Geneva: ISO1985.

[10] Steve Pepper, The TAO of Topic Maps

<http://www.ontopia.net/topicmaps/materials/tao.html>

[11] Thesaurii, Techquila

<http://www.techquila.com/tmsinia3.html>

[12] Semantic Web project

<http://www.w3.org/2001/sw/>

[13] Mathematical Subject Classification (MSC)

<http://www.ams.org/msc>

[14] Physics and Astronomy Classification Scheme (PACS)

<http://www.aip.org/pacs/>

[15] Dewey Decimal Classification (DDC)

<http://www.oclc.org/dewey/>

[16] DARPA Agent Markup Language (DAML)

<http://www.daml.org/>

[17] Информационная система ИСИР

<http://uis.isir.ras.ru/>

[18] Каталог ресурсов «Кирилл и Мефодий»

<http://search.km.ru/url/index.asp>

[19] Подходы к описанию и использованию тезаурусов в информационных системах. Аджиев Алим Сапарович, Нгуен Мань Хунг, Труды 5-й Всероссийской научной конференции RCDL2003, Санкт-Петербург, Россия, 2003

[20] Бездушный А.А., Сысоев Т.М., Нестеренко А.К., Бездушный А.Н., Серебряков В.А., RDFS как основа среды разработки цифровых библиотек и Web-порталов, Электронные библиотеки, 2003, Том 6, Выпуск 3.

<http://www.elbib.ru/index.phtml?page=elbib/rus/journal/2003/part3/BBNSS>

[21] Бездушный А.А., Сысоев Т.М., Нестеренко А.К., Бездушный А.Н., Серебряков В.А., Архитектура и технологии RDFS-среды разработки цифровых библиотек и Web-порталов, Электронные библиотеки, 2003, Том 6, Выпуск 4.

<http://www.elbib.ru/index.phtml?page=elbib/rus/journal/2003/part4/BNSBS>

[22] Архитектура RDFS-системы. Практика использования открытых стандартов и технологий Semantic Web в системе ИСИР. Бездушный А.А. Бездушный А.Н. Нестеренко А.К. Серебряков В.А. Сысоев Т.М., Труды 5ой Всероссийской научной конференции «Электронные библиотеки: перспективные методы и технологии, электронные коллекции» - RCDL2003, Санкт-Петербург, Россия, 2003

[23] XML-публикация в ИСИР. Бездушный А.А. Препринт ВЦ РАН Москва 2-2004

[24] FormBuilder - средство автоматизации Web -редактирования ресурсов.
Нестеренко А.К. Препринт ВЦ РАН Москва 2-2004

[25] DAML+OIL (March 2001) Reference Description.
<http://www.daml.org/2001/03/daml+oil-index.html>

Об авторах

Нгуен Мань Хунг - аспирант Вычислительного Центра имени А. А. Дородницына РАН. Сфера деятельности: программирование, научные информационные системы, тезаурусы, онтологии, метаданные, Semantic Web, XML, RDFS.

Тел. +7 095 135 54 71

E-mail: nmhungru@yahoo.com

Аджиев Алим Сапарович - младший научный сотрудник Центра научных телекоммуникаций и информационных технологий РАН. Сфера деятельности: программирование, Java, RDFS, Semantic Web, Web-порталы, базы данных, системное программирование, цифровые библиотеки, информационное обеспечение научной деятельности.

Тел. +7 095 938 37 09

E-mail: ajiev@ccas.ru
