

Извлечение научно-технических фактов из отраслевых документов на основе методов семантико-синтаксического и концептуального анализа

Кан А.В.*

ФГБУ «НИЦ «Институт имени Н.Е. Жуковского»
г. Москва, Российская Федерация
ORCID: <https://orcid.org/0000-0001-9410-406X>
e-mail: kanav@nrczh.ru

Козловская Я.Д.**

Московский авиационный институт (МАИ)
г. Москва, Российская Федерация
ORCID: <https://orcid.org/0000-0002-1780-5687>
e-mail: yana_kozlovskaja@mail.ru

Тололова А.А.***

Московский авиационный институт (МАИ)
г. Москва, Российская Федерация
e-mail: tokolovaa@gmail.com

Извлечение научно-технических фактов является трудной задачей с точки зрения корректности получаемой информации. Предлагаемая модель выделения фактов основывается на четких представлениях о смысловой структуре текста, выраженной в виде иерархии синтаксических конструкций единиц смысла, что позволяет выявить межфразовые связи в контактно расположенных предложениях. В качестве смысловых единиц используются отдельные слова, словосочетания, присущие конкретной предметной области, и образующие её понятийный состав. Для обработки исходного текста используются процедуры фразеологического, концептуального и семантико-синтаксического анализа текстов.

Ключевые слова: извлечение фактов, семантико-синтаксический анализ, концептуальный анализ, фразеологический анализ, смысловая структура, единица смысла.

Для цитаты:

Кан А.В., Козловская Я.Д., Тололова А.А. Извлечение научно-технических фактов из отраслевых документов на основе методов семантико-синтаксического и концептуального анализа // Моделирование и анализ данных. 2024. Том 14. № 1. С. 27–40.
DOI: <https://doi.org/10.17759/mda.2024140102>



***Кан Анна Владимировна**, кандидат технических наук, доцент МАИ, начальник аналитического отдела ФГБУ «НИЦ «Институт имени Н.Е. Жуковского», ORCID: <https://orcid.org/0000-0001-9410-406X>, e-mail: kanav@nrczh.ru

****Козловская Яна Дмитриевна**, студент магистратуры, ФГБОУ ВО Московский авиационный институт (МАИ), ORCID: <https://orcid.org/0000-0002-1780-5687>, e-mail: yana_kozlovskai@mail.ru

*****Токолова Алина Александровна**, студент магистратуры, институт «Компьютерные науки и прикладная математика» Московский авиационный института (МАИ), г. Москва, Российская Федерация, e-mail: tokolovaa@gmail.com

1. ВВЕДЕНИЕ

Научные публикации уже давно стали основным источником и способом распространения научно-технических открытий. Согласно отчету [1], с 2018 года ежегодное количество научных публикаций, в том числе опросов, тезисов и обзоров, выросло на 22,4% и за 2022 год достигло значения свыше 5,4 миллионов. Такой объем информации качественно затрудняет ручной поиск релевантных сведений и фактов для их дальнейшей обработки.

Научно-технические факты представляют особый интерес для пользователя и под этим термином понимается событие, понятие или явление. Для корректного выделения фактов из отраслевой документации используются методы извлечения фактологической информации (ИФИ), среди которых семантико-синтаксический, фразеологический и концептуальный анализ текстов. Выделение наименований понятий (НП) основывается на вычислении меры смысловой значимости, а смысловая структура текста представляется в виде предикатно-актантной структуры (ПАС), получаемой методами семантико-синтаксического и концептуального анализа.

2. ОСНОВНЫЕ ПОЛОЖЕНИЯ ФРАЗЕОЛОГИЧЕСКОГО КОНЦЕПТУАЛЬНОГО АНАЛИЗА ТЕКСТОВ

Фразеологический концептуальный анализ [2] используется для установления единиц смысла предложения, которые и формируют смысловое содержание текста.

Основные положения концепции фразеологического концептуального анализа текстов:

- Смысловое содержание текстов выражается с помощью единиц смысла.
- Понятие – самая устойчивая единица смысла.
- Объекты предложения обладают особыми признаками, выражающимися через предикатно-актантную структуру (ПАС) и набором отношений с другими объектами.
- Сверхфразовые единства формируются из предложений и представляются в виде последовательностей предложений (связного текста).

Предикатно-актантная структура (ПАС) – представляет предложения в виде понятий-предикатов, устанавливающих связи между объектами, и в виде понятий-актантов, содержащих признаки объектов. ПАС обеспечивает автоматический перевод текстов с естественных языков на формализованные и наоборот.

3. ОПРЕДЕЛЕНИЕ ЗНАЧИМЫХ НАИМЕНОВАНИЙ ПОНЯТИЙ В ТЕКСТЕ

Значимые наименования понятий [3,4] выделяются методами статистического, синтаксического и концептуального анализа текстов. Сложность этого процесса заключается в фиксировании границ наименований понятий, несущих ключевую смысловую нагрузку.

Семантические методы позволяют с помощью эталонных словарей идентифицировать значимые слова и словосочетания определенной предметной области.

Статистические методы позволяют определить значимые понятия путем присвоения объектам текста весовых коэффициентов на основе частотного анализа встречаемости в конкретном документе и во всем корпусе документов.

Синтаксические методы позволяют определить, какую синтаксическую роль играют значимые слова и словосочетания, определяя их синтаксическую роль в предложении, как субъект, предикат или объект.

Смысловая значимость слов в контексте документа оценивается статистической мерой $TF \cdot IDF$ [5].

$$TF(\text{concept}, \text{Text}) = \frac{l_{\text{concept}}}{\sum_k l_k} - \text{частота слова в документе,}$$

где l_{concept} – число вхождений наименования понятия *concept* в документ;
 $\sum_k l_k$ – общее число наименований понятий в данном документе.

$$IDF(\text{concept}, M) = \log \frac{|M|}{|\{Text_n \in M \mid \text{concept} \in Text_n\}|}$$

$IDF(\text{concept}, M)$ – обратная частота документа, т.е. инверсия частоты, с которой наименование понятия встречается в корпусе текстов, каждому слово соответствует одно значение IDF , где $|M|$ – число документов в корпусе;
 $|\{Text_n \in M \mid \text{concept} \in Text_n\}|$ – число документов из массива документов M , в которых встречается понятие *concept* (когда $l_{\text{concept}} \neq 0$).

Эта статистическая мера увеличивает значение слова пропорционально частоте его появления в тексте, однако, уменьшает это же значение при появлении слова в большом количестве документов.



4. МОДЕЛЬ ВЫЯВЛЕНИЯ ФАКТОВ В КОЛЛЕКЦИИ ТЕКСТОВ

Факт *Fact* состоит из контактно расположенных предложений текста и содержится во фрагменте текста. Таким образом, факт можно представить в виде некоторой последовательности предложений:

$$Fact = (Sent_u, Sent_{u+1}, \dots, Sent_t),$$

где $Sent_u, Sent_{u+1}, \dots, Sent_t$ – последовательность предложений факта;

u – первое предложение факта;

t – последнее предложение факта.

Предложение факта *Sent* представляет собой кортеж элементов факта:

$$Sent = el_1, el_2, \dots, el_i,$$

где el_i – элемент факта,

i – количество элементов факта.

Элементом текста *el* является слово или знак препинания:

$$GPS\ el = \{w, pm\},$$

где w – слово;

pm – знак препинания.

Каждый элемент текста *el* обладает грамматическими и синтаксическими признаками:

$$GSP(el) = (GP, SP),$$

где *GP* – грамматические признаки;

SP – синтаксические признаки.

Термином «понятие» определяется некоторый социально значимый мыслительный образ, представляемый в виде отдельного слова или в виде устойчивого фразеологического словосочетания. Понятие *concept* можно представить математически в виде одного элемента факта или в виде совокупности нескольких элементов факта:

$$concept = \sum \left\{ el_j \right\}_{j=1}^{V\text{-счетное}}.$$

Каждое понятие *concept* имеет унифицированную и нормализованную формы, получаемые функциями унификации и нормализации:

$Unif : concept \rightarrow concept^{un}$ – функция унификации,

$Norm : concept \rightarrow concept^{nf}$ – функция нормализации,

где $concept^{un}$ – унифицированная форма понятия;

$concept^{nf}$ – нормализованная форма понятия.

Смысл предложения факта *MF* выражается через его предикатно-актантную структуру (систему элементарных смысловых триад):

$$MF = \{PSO_1, PSO_2, \dots, PSO_e\},$$

где e – количество элементарных смысловых триад в предложении.

Элементарная смысловая триада представляется в виде кортежа субъекта S , предиката P и объекта O :

$$PSO = S, P, O,$$

где P – смысловая связь между субъектом и объектом (предикат);

S – главное слово или словосочетание элементарного высказывания (субъект);

O – зависимое слово или словосочетание элементарного высказывания (объект).

Полное описание факта [5] представляется в виде совокупности предикатно-актантных структур предложений факта:

$$FullFact = \sum \{MF_v\}_{v=1}^{V\text{-счетное}}, FullFact \in Fact.$$

5. ПРОЦЕСС ИЗВЛЕЧЕНИЯ НАУЧНО-ТЕХНИЧЕСКИХ ФАКТОВ НА ОСНОВЕ СИНТАКСИЧЕСКОЙ МОДЕЛИ ТЕКСТА

Реализованная синтаксическая модель текста на основе системы обобщенных синтагм, которые выражают форму единицы смысла и состоят из сочетания символов обобщенных грамматических классов слов, входящих в состав словосочетаний эталонного словаря, позволяет анализировать его структуру, извлекать понятийный состав и определять смысловые отношения. Эта модель также помогает извлекать ключевые синтаксические конструкции предложения и классифицировать их элементы. Значимые именованные сущности регистрируются и обрабатываются в соответствии с их синтаксической ролью. Затем эти элементы унифицируются в соответствии с требованиями системы.

Алгоритм автоматического выявления и формализации фактов в текстах:

Шаг 1. Разделить исходный текст на предложения и выполнить морфологический анализ предложений.

Шаг 2. Определить именные и глагольные словосочетания и установить их синтаксическую роль в предложении при помощи упрощенного семантико-синтаксического анализа.

Шаг 3. Составить частотный словарь слов и словосочетаний, определить какие из этих слов и словосочетаний являются значимыми для данного текста.

Шаг 4. Присвоить уникальные идентификаторы каждому словосочетанию согласно словарю унифицированных формализованных представлений наименований понятий (УФПНП) и сопоставить их исходные формы с унифицированными формами представления, указав все их местоположения в тексте (номера предложений).

Шаг 5. Определить номера предложений, содержащих значимые слова и словосочетания, используя словарь указателей связей предложений (УСП).



Шаг 6. Определить связи между предложениями, содержащими ключевые понятия, и их окружением. Установить границы описания фактов в тексте, используя разметку текста по указателям смысловых связей и обобщенным наименованиям понятий.

Шаг 7. Присвоить каждому текстовому описанию фактов идентификационный номер, содержащий порядковый номер события, код и тип документа.

Шаг 8. Определить главные и второстепенные члены предложения, границы словосочетаний, построить дерево зависимостей предложения, построить ПАС и сформировать «скелет» предложения, используя информацию, полученную на шаге 2.

Шаг 9. Определить обобщенную синтагму и построить формализованное представление для каждого словосочетания.

Шаг 10. Соотнести полученные методом концептуального анализа словосочетания со словосочетаниями, полученными путем синтаксического анализа.

Шаг 11. Произвести нормализацию.

Шаг 12. Расчленив описание каждого факта на составные элементы – формализованное представление элементов ПАС, «скелет» предложения с указанием номеров словосочетаний в эталонном концептуальном словаре (ЭКС).

Шаг 13. Выполнить генерацию формализованных представлений предложений факта в обобщенное формализованное представление его смысловой структуры.

Шаг 14. Произвести преобразование обобщенного формализованного представления в его машинную форму.

Основным отличием предлагаемой модели от существующих заключается в том, что в рамках этой модели смысловое представление текста выражается в виде иерархии синтаксических конструкций единиц смысла.

Описания конкретных фактов в тексте могут отображаться его контактно расположенной последовательностью предложений, связанных межфразовыми связями [6]. Рассмотрим в качестве такого примера текст, представленный в таблице 1.

Таблица 1

Фрагмент текста

Стратегический бомбардировщик “В -52” более полувека стоит на вооружении ВВС США. В настоящее время продажа за пределами США боинга “В-52” запрещена, так как его летные характеристики многократно превосходят зарубежные аналоги. Самолеты бомбардировщики “Б-52” были введены в строй с 1955 года. Всего было построено 728 бомбардировщиков. По проекту каждый из этих самолетов несет на борту до 51 единицы боеприпасов.

Разобьем исходный текст на предложения и присвоим каждому идентификационный индекс, состоящий из порядкового номера события, номера предложения в нем, источника информации и даты публикации (см. табл. 2).

Таблица 2

Перенумерованные предложения с идентификационными индексами

1. *Стратегический бомбардировщик “В-52” более полувека стоит на вооружении ВВС США. – 458496_1*
2. *В настоящее время продажа за пределами США боинга “В-52” запрещена, так как его летные характеристики многократно превосходят зарубежные аналоги. – 458496_2*
3. *Самолеты бомбардировщики “Б-52” были введены в строй с 1955 года. – 458496_3*
4. *Всего было построено 728 бомбардировщиков. – 458496_4*
5. *По проекту каждый из этих самолетов несет на борту до 51 единицы боеприпасов. – 458496_5*

В приведенном фрагменте текста межфразовая связь между предложением № 4 и предложением № 5 обусловлена местоименной анафорой «**их этих**», связь между предложениями № 1, № 2 № 3 и № 4 обусловлена синонимичными конструкциями «**Стратегический бомбардировщик – самолет бомбардировщик**». Наличие межфразовых связей является необходимым условием выделения описания конкретного факта в тексте [7]. Дополнительным условием, предоставленным в таблице 3, является смысловая связь (родовидовые и синонимичные отношения) между частью наименований понятий этого фрагмента, находящимися в разных предложениях (в скобках указаны номера предложений).

Таблица 3

Смысловая связь между частью наименований понятий фрагмента, находящихся в разных предложениях

стратегический бомбардировщик “В-52” (№ 1) – самолеты бомбардировщики “Б-52” (№ 3) – бомбардировщиков (№ 4)

стратегический бомбардировщик “В-52” (№ 1) – боинга “В-52” (№ 2) – самолеты бомбардировщики “Б-52” (№ 3)

В рамках разработанной синтаксической модели возможно осуществить автоматическое преобразование исходного текста в его формализованную синтаксическую структуру, аналогичное приведенной выше структуре. В процессе такого преобразования для получения унифицированного представления элементов ПАС необходимо дополнительно выполнить нормализацию форм слов наименований понятий и их унификацию по словарю УФПНП. Результат такой операции приведен в таблице 4.

Таблица 4

Унифицированные формы представления наименований понятий

| Наименование понятия | Унифицированная форма представления |
|-------------------------------|-------------------------------------|
| стратегический бомбардировщик | самолет |
| боинг | самолет |



| Наименование понятия | Унифицированная форма представления |
|----------------------|-------------------------------------|
| ВВС США | войско |
| боеприпасы | оружие |
| вооружение | оружие |
| США | страна |

В таблице 5 указывается идентификатор факта и приводится его исходная текстовая форма, указывается состав ПАС (в мнемонике грамматических классов слов), унифицированный «скелет», идентификатор факта, номер предложения; приводится формализованное представление ПАС предложения в сокращенной форме в виде символов обобщенных синтагм и в виде унифицированных форм главных слов ПАС; приводится формализованное представления «скелета» (SkIRus) предложений факта (в сокращенном виде – только главные слова, в полном виде эти элементы представлены номерами элементов словаря ЭКС); приводится формализованное представление предложений (SenRus) в виде последовательности нормальных слов предложения [7].

Таблица 5

Результаты семантико-синтаксического и концептуального анализа предложений описания факта, приведенного в таблице 1

| Исходное предложение № 1 | |
|--|--|
| <i>Стратегический бомбардировщик “В -52” более полувека стоит на вооружении ВВС США.</i> | |
| Формализованное описание элементов ПАС предложения $PSO=\{V, N, N\}=\{\text{стоять; бомбардировщик; вооружение}\}$ | |
| Формализованное представление элементов ПАС предложений | |
| <i>Predicate (P)</i> | $лА = \text{стоять} \#лА = \text{стоять}$ |
| <i>Subject (S)</i> | $FA = \text{бомбардировщик} \#FA = \text{бомбардировщик}$ |
| <i>Object (O)</i> | $ЁК = \text{вооружение} \#ЁК = \text{вооружение}$ |
| Формализованное представления ПАС предложений (PSORus) | |
| $PSO = лАFAЁК = \text{стоять; бомбардировщик; вооружение} \# лАFAЁК$ | |
| Формализованное представления «скелета» предложений (SkIRus) | |
| $N \# N \# NVFN = FA \# \# @A9A \# \# YAlA \# A \# ЁК \# A \# A = \text{бомбардировщик “В -52” полувек стоять на вооружение ВВС США}$ | |
| Формализованное представления предложений (SenRus) | |
| $PiFA \# \# @A9A \# \# cAYAlA \# A \# ЁК \# A \# A = \text{стратегический бомбардировщик “В -52” более полувек стоять на вооружение ВВС США}$ | |
| Исходное предложение № 2 | |
| <i>В настоящее время продажа за пределами США боинга “В-52” запрещена, так как его летные характеристики многократно превосходят зарубежные аналоги.</i> | |
| Формализованное описание элементов ПАС предложения $PSO=\{N, N, V, N, V, N\}=\{\text{продажа; боинг; запрещен; превосходит; аналог}\}$ | |
| Формализованное представление элементов ПАС предложений | |
| <i>Predicate (P)</i> | $aA = \text{запрещен} \# aA = \text{запрещен}$ $Юр = \text{превосходит} \# Юр = \text{превосходит}$ |



| | |
|---|--|
| Subject (S) | <i>uB = продажа #uB =продажа</i> |
| | <i>vB = характеристика # vB = характеристика</i> |
| Object (O) | <i>FA = боинг # FA =боинг</i> |
| | <i>FA = аналог # FA =аналог</i> |
| Формализованное представления ПАС предложений (PSORus) | |
| <i>PSO_1= uBFaA = продажа; боинг; запрещен # uBFaA</i> | |
| <i>PSO_2= vBЮpFA =характеристика; превосходит; аналог# vBЮpFA</i> | |
| Формализованное представления «скелета» предложений (SkIRus) | |
| <i>FNNFN»N»K, NVN =ыАГтvHRWbTKхаИ,, АТежА= в время продажа за предел боинг запрещен, характеристика превосходить аналог</i> | |
| Формализованное представления предложений (SenRus) | |
| <i>ыАПiИмиВъААСъАFA»»9А»»аА,, сАсАяАНТvВсАлАНТFA = в настоящий время продажа за пределам США боинг »В-2» запрещен, так как его летный характеристика многократно превосходить зарубежный аналог</i> | |
| Исходное предложение № 3 | |
| <i>Самолеты бомбардировщики "Б-52" были введены в строй с 1955 года.</i> | |
| Формализованное описание элементов ПАС предложения | |
| <i>PSO={V, N, N}={введен; самолет; строй}</i> | |
| Формализованное представление элементов ПАС предложений | |
| Predicate (P) | <i>aA = введен # aA =введен</i> |
| Subject (S) | <i>AA = самолет # AA =самолет</i> |
| Object (O) | <i>Dx = строй #Dx=строй</i> |
| Формализованное представления ПАС предложений (PSORus) | |
| <i>PSO= aAAADx= введен; самолет; строй # aAAADx</i> | |
| Формализованное представления «скелета» предложений (SkIRus) | |
| <i>N»N»LKFNNFN = eW»»@A9A»»ЯАаАыADхцA9= самолет "Б-52" был введен в строй с 1955 год</i> | |
| Формализованное представления предложений (SenRus) | |
| <i>AAFA»»9А»»ЯАаАыADхцA9АНА=самолет_бомбардировщик_ "Б-52" _был_введен_в_строй_с_1955_год</i> | |
| Исходное предложение № 4 | |
| <i>Всего было построено 728 бомбардировщиков.</i> | |
| Формализованное описание элементов ПАС предложения | |
| <i>PO={V, N}={построен; бомбардировщик}</i> | |
| Формализованное представление элементов ПАС предложений | |
| Predicate (P) | <i>aA = построен #aA=построен</i> |
| Subject (S) | <i>FA = бомбардировщик # FA =бомбардировщик</i> |
| Формализованное представления ПАС предложений (PSORus) | |
| <i>PO=aAFA = построен; самолет # aAFA</i> | |
| Формализованное представления «скелета» предложений (SkIRus) | |
| <i>LKN = был_построен_бомбардировщик</i> | |
| Формализованное представления предложений (SenRus) | |
| <i>ЧМЯАaA9AFA = всего был_построен_728_бомбардировщик</i> | |
| Исходное предложение № 5 | |
| <i>По проекту каждый из этих самолетов несет на борту до 51 единицы боеприпасов.</i> | |



| | |
|---|--|
| Формализованное описание элементов ПАС предложения $PSO=\{V, N, N\}=\{\text{несет}; \text{самолет}; \text{боеприпас}\}$ | |
| Формализованное представление элементов ПАС предложений | |
| Predicate (P) | $\text{БУ} = \text{несет} \# \text{БУ} = \text{несет}$ |
| Subject (S) | $\text{AA} = \text{самолет} \# \text{AA} = \text{самолет}$ |
| Object (O) | $\text{AA} = \text{боеприпас} \# \text{AA} = \text{боеприпас}$ |
| Формализованное представления ПАС предложений (PSORus) | |
| $PSO = \text{БУАААА} = \text{несет}; \text{самолет}; \text{боеприпас} \# \text{БУАААА}$ | |
| Формализованное представления «скелета» предложений (SkIRus) | |
| $\text{FN FN VFNFN} = \text{по_проект_из_этот_нести_на_борт_до_51_единиц}$ | |
| Формализованное представления предложений (SenRus) | |
| $\text{xAAAHNTфAΦAAAЪУыAHфA9ABBA} = \text{по_проект_каждый_из_этот_самолет}$ $\text{несет_на_борт_до_51_единица_боеприпас}$ | |

Результатом работы алгоритма является формализованное представление предложений, обобщающее смысловое содержания фактологической информации. Для генерации формализованного представления обобщенного смыслового содержания факта выделялись унифицированные представления элементов ПАС предложений. Результаты такого представления приведены в таблице 6 в виде формализованного и индексного представления элементов ПАС факта. Каждый элемент ПАС сопровождается его весовым коэффициентом, индексом синтаксической роли, номером словосочетания в словаре ЭКС и уникальным идентификатором факта.

Таблица 6. Формализованное и индексное представление элементов ПАС факта текста № 387473. В этих представлениях номер события был заменен на его уникальный идентификационный номер (ID=008154)

| |
|---|
| Формализованное представление элементов ПАС факта |
| $PSO=\{V, N, N\} \rightarrow$ $P=\{\text{стоять}_{064_065785}, \text{запрещен}_{049_087631}, \text{превосходит}_{057_075324},$ $\text{введен}_{047_087631}, \text{построен}_{049_098713}, \text{несет}_{058_075692}\};$ $S=\{\text{бомбардировщик}_{025_034758}, \text{продажа}_{039_062584}, \text{характеристика}_{039_07234},$ $\text{самолет}_{029_046273}\}; \text{O}=\{\text{вооружение}_{038_086759}, \text{аналог}_{020_015638},$ $\text{строй}_{029_063348}\}$ |
| Индексное представление элементов ПАС факта |
| $\text{стоять}_{\{064_065785_P_008154\}}, \text{запрещен}_{\{049_087631_P_008154\}},$ $\text{превосходит}_{\{057_075324_P_008154\}}, \text{введен}_{\{047_087631_P_008154\}},$ $\text{построен}_{\{049_098713_P_008154\}}, \text{несет}_{\{058_075692_P_008154\}},$ $\text{бомбардировщик}_{\{025_034758_S_008154\}}, \text{продажа}_{\{039_062584_S_008154\}},$ $\text{характеристика}_{\{039_07234_S_008154\}}, \text{самолет}_{\{029_046273_S_008154\}},$ $\text{вооружение}_{\{038_086759_O_008154\}}, \text{аналог}_{\{020_015638_O_008154\}}, \text{строй}_{\{029_063348_P_008154\}}$ |

Таким образом, как видно из таблицы 5 и таблицы 6, основная задача формализации смыслового представления описания факта заключается в генерации совокупности унифицированных представлений ПАС. Полученная совокупность формализованных представлений ПАС, «скелета» и предложения позволяет не только обеспечить возможность поиска по любому элементу их формализованного описания, но и обеспечивает возможность последовательного перехода от каждого нижестоящего элемента формализованного описания к вышестоящему, а также возможность перехода в обратном порядке. Такое представление смысла предложений обеспечивает реализацию всего спектра семантических операций над смысловым содержанием факта.

Полученное формализованное представление обобщенного смыслового содержания факта позволяет производить поиск и сопоставление идентичных или близких по смысловому содержанию фактов, а также обеспечивает возможность их классификации по различным основаниям: по их содержанию, по именам субъектов или объектов фактов или по тональности отношений между ними и др.

6. ЗАКЛЮЧЕНИЕ

В статье описано решение задачи автоматического извлечения фактов из научно-технологических документов на основе методов их семантико-синтаксического и концептуального анализа.

Представленная модель алгоритма автоматического извлечения фактов решает проблему формализации и унификации смыслового содержания фактов. Результатами работы этого алгоритма является представление смысловой структуры фактов, полученных из исходных корпусов текста, в виде совокупности формализованных представлений ПАС, обеспечивающих возможность поиска по любому элементу их формализованного описания. Такое представление смысла предложений обеспечивает реализацию всего спектра семантических операций над смысловым содержанием факта, в том числе классификации и сопоставления, что дает возможность выполнять проверку корректности извлеченной информации.

Литература

1. *Curcic D.* Number of Academic Papers Published Per Year // Wordsrated URL: <https://wordsrated.com/number-of-academic-papers-published-per-year/#:~:text=As%20of%202022%2C%20over%205.14,5.03%20million%20papers%20were%20published.> (дата обращения: 09.10.2023).
2. *Белоногов Г.Г., Калинин Ю.П., Хорошилов А.А.* Компьютерная лингвистика и перспективные информационные технологии. Теория и практика построения систем автоматической обработки текстовой информации. – М.: Изд-во Русский мир, 2004.
3. *Хорошилов Ал-др А., Мусабаев Р.Р., Козловская Я.Д., Никитин Ю.В., Хорошилов А.А.* Автоматическое выявление и классификация информационных событий в текстах СМИ // Научно-техническая информация. Серия 2: Информационные процессы и системы. 2020. № 7. С. 27–38. DOI: 10.36535/0548-0027-2020-07-4
4. *Хорошилов Ал-др. А., Никитин Ю.В., Хорошилов Ал-ей. А., Будзко В.И.* Автоматическое создание формализованного представления смыслового содержания неструктурированных



текстовых сообщений СМИ и социальных сетей // Системы высокой доступности, № 3, том.10, 2014, С. 36–51.

5. Кан А.В., Козловская Я.Д., Кадушкин Н.А., Хорошилов Ал-р А. Автоматическая кластеризация документов СМИ на основе анализа их смыслового содержания // Моделирование и анализ данных. 2020. Том 10. № 3. С. 24–38. DOI: <https://doi.org/10.17759/mda.2020100302>
6. Богатырев М.Ю. Извлечение фактов из текстов естественного языка с применением концептуальных графовых моделей // Известия ТулГУ. Технические науки. – 2016. – № 7. – Ч. 1.
7. Хорошилов Ал-др А., Козловская Я.Д., Мусабаев Р.Р., Красовицкий А.М., Хорошилов Ал-ей А. Определение тональности сообщений СМИ методом их концептуального анализа // Моделирование и анализ данных. 2019. № 4. DOI: [10.17759/mda.2019090405](https://doi.org/10.17759/mda.2019090405)



Extracting Scientific and Technical Facts from Industry Documents Based on Methods of Their Semantic-syntactic and Conceptual Analysis

Anna V. Kan*

FSBI “National Research Center” Institute named after N.E. Zhukovsky”

Moscow, Russian Federation

ORCID: <https://orcid.org/0000-0001-9410-406X>

e-mail: kanav@nrczh.ru

Yana D. Kozlovskaya**

Moscow Aviation Institute (MAI), Moscow, Russian Federation

ORCID: <https://orcid.org/0000-0002-1780-5687>

e-mail: yana_kozlovskaja@mail.ru

Alina A. Tokolova***

Moscow Aviation Institute (MAI), Moscow, Russian Federation

e-mail: tokolovaa@gmail.com

Extraction of scientific and technical facts is a difficult task in terms of correctness of the obtained information. The proposed fact extraction model is based on clear ideas about the semantic structure of the text, expressed as a hierarchy of syntactic constructions of meaning units, which allows identifying interphrase relations in contacted sentences. Individual words, word combinations inherent to a particular subject area and forming its conceptual composition are used as meaning units. The procedures of phraseological, conceptual and semantic-syntactic analysis of texts are used to process the source text.

Keywords: fact extraction, semantic-syntactic analysis, semantic-syntactic analysis, conceptual analysis, semantic triad.

For citation:

Kan A.V., Kozlovskaya Y.D., Tokolova A.A. Extracting Scientific and Technical Facts from Industry Documents Based on Methods of Their Semantic-syntactic and Conceptual Analysis. *Modelirovanie i analiz dannykh = Modelling and Data Analysis*, 2024. Vol. 14, no. 1, pp. 27–40. DOI: <https://doi.org/10.17759/mda.2024140102> (In Russ., abstr. in Engl.).

***Anna V. Kan**, Candidate of Technical Sciences, Associate Professor, Moscow Aviation Institute (MAI), Head of the Analytical Department, Federal State Budgetary Institution “National Research Center” Institute named after N.E. Zhukovsky”, e-mail: kanav@nrczh.ru

****Yana D. Kozlovskaya**, master’s student, Moscow Aviation Institute (National Research University) (MAI), Moscow, Russia, ORCID: <https://orcid.org/0000-0002-1780-5687>, e-mail: yana_kozlovskaja@mail.ru

*****Alina A. Tokolova**, master’s student, Institute of Computer Science and Applied Mathematics Moscow Aviation Institute (National Research University) (MAI), Moscow, Russia, e-mail: tokolovaa@gmail.com



References

1. Curcic D. Number of Academic Papers Published Per Year // Wordsrated URL: <https://wordrated.com/number-of-academic-papers-published-per-year/#:~:text=As%20of%202022%2C%20over%205.14,5.03%20million%20papers%20were%20published> (Accessed: 09.10.2023).
2. Belonogov G.G., Kalinin Y.P., Khoroshilov A.A. Computer linguistics and perspective information technologies. Theory and practice of building systems of automatic processing of text information. – Moscow: Izd-vo Russky Mir, 2004.
3. Khoroshilov AI-Dr. A., Musabaev R.R., Kozlovskaya Y.D., Nikitin Y.V., Khoroshilov A.A. Automatic detection and classification of information events in mass media texts// Scientific and Technical Information. Series 2: Information processes and systems. 2020. № 7. С. 27–38. DOI: 10.36535/0548-0027-2020-07-4
4. Khoroshilov AI-Dr. A., Nikitin Y.V., Khoroshilov AI-ey. A., Budzko V.I. Automatic creation of formalized representation of semantic content of unstructured text messages of mass media and social networks // Systems of High Availability, No.3, Vol.10, 2014, pp. 36–51.
5. Kan A.V., Kozlovskaya Y.D., Kadushkin N.A., Khoroshilov AI-r A. Automatic clustering of media documents based on the analysis of their semantic content // Modeling and Data Analysis. 2020. Vol. 10. No. 3. С. 24–38. DOI: <https://doi.org/10.17759/mda.2020100302>
6. Bogatyrev, M. Yu. Fact extraction from natural language texts using conceptual graph models // Izvestiya TulSU. Technical Sciences. –2016. – № 7. – Vol. 1.
7. Khoroshilov AI-Dr. A., Kozlovskaya Ya.D., Musabaev R.R., Krasovitsky A.M., Khoroshilov AI-ey A. Determination of the tone of media messages by their conceptual analysis method // Modeling and Data Analysis. 2019. № 4. DOI: 10.17759/mda.2019090405

Получена 04.03.2024

Принята в печать 18.03.2024

Received 04.03.2024

Accepted 18.03.2024