

Computational Linguistics and Intellectual Technologies:  
Proceedings of the International Conference “Dialogue 2019”

Moscow, May 29—June 1, 2019

## **SIMULATION OF BACKGROUND KNOWLEDGE AND BRIDGING IN RUSSIAN**

**Dikonov V. G.** (dikonov@iitp.ru)

IITP RAS, Moscow, Russia

This paper introduces a knowledge-based semantic approach towards bridging annotation of Russian texts. Our method simulates human background knowledge by using compact domain descriptions based on an extended version of SUMO ontology and lexical-semantic data from the “Universal Dictionary of Concepts”. Our approach supports a wide and extensible range of bridging relations. The tagger that implements it can build complex bridges with multiple arcs, supports making assumptions and can be adapted to annotate other languages supported by the underlying dictionary of concepts.

**Keywords:** computational linguistics, bridging, reference, anaphora, linguistic ontology

## **МОДЕЛИРОВАНИЕ ФОНОВЫХ ЗНАНИЙ И ПОИСК АССОЦИАТИВНЫХ АНАФОР В ТЕКСТАХ НА РУССКОМ ЯЗЫКЕ**

**Диконов В. Г.** (dikonov@iitp.ru)

ИППИ РАН, Москва, Россия

## 1. Bridging

The notion of bridging, also known as indirect or associative anaphora, introduced by Clark [2] captures an essential mechanism of text interpretation inside a human mind. Every comprehensible piece of text contains some identifiable entities and statements about them. As we read or listen, we encounter new entities (new information) and refer them with previously mentioned ones (given information in Clark's terms). This is a way to construct a mental model of the reality described in the text. Pairs of related entities are viewed as a kind of anaphora where the previously given entity becomes the antecedent of the new. Unlike the common definition of anaphora, the relation between the entities in bridging is not limited to identity. It can be of any semantic type meaningful to the Listener. Building references requires deep background knowledge of the relevant domain, especially when the Speaker skips "redundant" details to improve the speed of communication. When the Listener fails to find a directly related antecedent of some new entity, he is forced to insert a suitable intermediate concept. This leads to creation of "bridges" with multiple "arcs".

- (1) На станции метро «Владыкино» в Москве найдено взрывное устройство.  
Найденный предмет обследовали с использованием служебных **собак**.  
(An *explosive device* was found at the Moscow underground station Vladykino.  
The object was *examined* with service **dogs**)

Here the Listener assumes the existence of policemen, who were not mentioned directly, and constructs the following possible bridge: *explosive device* <sup>isObjectOf</sup> *examine* <sup>hasAgent</sup> policemen <sup>isUserOf</sup> *dog*. This assumption is based on common knowledge about police work and terrorism, implanted by television newscasts. It is possible to build arbitrarily long bridges by adding new assumptions.

Introduction of new concepts associated with the given information from background knowledge is also a productive mechanism of creativity. Its proper modeling combined with good plausibility filters might give AI the ability to invent.

## 2. Overview of the approach

Existing works in the field of bridging fall into two groups: semantic approaches and syntactic ones. Syntactic approaches choose particular syntactic patterns, usually definite NPs, and treat the ability of certain words to fill such patterns as a criterion of a non-typed bridging relation. Later research by Hou [5] departs from a single pattern restriction, but still lacks the ability to explicitly represent the meaning of the detected bridging relations. The first published paper on bridging in Russian [9] follows the same path and uses Russian genitive NPs as the clue pattern.

A semantic approach always ascribes a semantic type to the discovered relations. The authors of such approaches often impose restrictions on the types of bridges they detect in order to accommodate to their resources and relation search methods. Papers by Poesio [7], Lassale [6] concentrate only on part-whole relations. Recasens [8], Zikánová [10], etc. add set-subset, cohyponymy, predicate-argument and symptom relations. Many studies rely on Princeton Wordnet as the source of lexical data and

a knowledge base to estimate semantic relatedness of words. Unfortunately, English Wordnet provides only part of the information needed to simulate the mental mechanism of bridging. It offers good lexical coverage, usable (though poorly organized) taxonomy, but is very limited in the field of semantic relations other than part-whole. In particular, it lacks cause-result and predicate-role relations. Roitberg et al. [9] wrote that absence of a (large scale) Russian Wordnet prevents the use of semantic methods on Russian material. We would answer that there are alternative resources for Russian and they have some advantages over the English Wordnet. One of them is briefly described in [section 2.1](#).

We take a semantic approach based on a rich background knowledge base (KB). The target language is Russian, but our KB is a language-neutral semantic resource. As a result, our bridging tool can be adapted to work with other languages supported by the underlying semantic dictionary UNLDC [3] (English, Hindi, French etc). Our project bears resemblance with the work by Fan [4], yet it is different in some key points. Both projects use a knowledge base encoded as a semantic graph and support simple taxonomy based inference. However, the structure and contents of the KBs are different. The set of relations in our study is wider. Our tool supports making assumptions and builds complex chain relations with intermediate concepts like the relation between the bomb and dogs in [example 1](#).

## 2.1. Resources

Our relation search engine operates with ontology concepts instead of words. We use a modified version of SUMO ontology with greatly extended taxonomy (extended ontology). This extension exists in the framework of developing the “Universal Dictionary of Concepts” (UNLDC) [3]. The extended ontology is an experimental resource and is different from the internal ontology of the linguistic processor ETAP<sup>1</sup> (ETAP ontology), which is also based on SUMO and mentioned further in this paper.

UNLDC translates the concepts of the extended ontology into Russian and several other languages. The Russian lexicon used in this project contains 42,973 Russian words and multi-word expressions with 66,896 senses total. These senses are linked with 48,883 concepts of the extended ontology (both original SUMO concepts and the added ones). UNLDC also has a growing semantic network that includes many relation types not available in the Wordnet, including the cause-result and argument ones. The types of semantic relations supported in this project are described in [section 3.2](#). UNLDC is an open public resource. Its core parts are available for download at GitHub2. The extended ontology is a supplement to UNLDC.

---

<sup>1</sup> ETAP is a multipurpose linguistic processor developed by the laboratory of computer linguistics at the Institute of Information Transmission Problems (IITP) in Moscow. It supports robust syntactic parsing, English ↔ Russian machine translation, paraphrasing, semantic parsing using two different frameworks, question answering and more.

<sup>2</sup> <https://github.com/dikonov/Universal-Dictionary-of-Concepts>

### 3. Knowledge base

Modeling the mechanism of human association reference requires an imitation of human knowledge about the subject domain of the text, which consists of:

- a) set of concepts relevant to the domain,
- b) semantic relations that hold between such concepts.

It also needs imitation of the relevant subset of human linguistic ability sufficient for transition from an NL text to a set of concepts. The latter includes at least chunking and morphology engines to identify sentences and lemmatize words, a semantic lexicon linking the words with concepts and some kind of lexical disambiguation.

#### 3.1. Concept inventory

Using ontology concepts to abstract away from lexical variation and peculiarities of different natural languages always poses the problem of choosing the right degree of abstraction or “semantic grain” for the task. Consider the following example:

- (2) Во Владимирской области произошло столкновение товарного поезда с застрявшим на переезде *грузовиком*. **Водитель** успел выскочить из кабины. **Машинист** получил травмы.

(A freight train hit a *truck* stuck at a crossing in the Vladimir region. The **driver** managed to jump out of the cabin. The **train driver** was injured.)

Bridging is expected to establish relations of association between *машинист* (*train driver*) and *поезд* (*train*), *водитель* (*driver*) and *грузовик* (*truck*) based on the fact that each type of driver controls a particular type of vehicle. This information is embedded in definitions of Russian words.

Initially we had three different sets of concepts offered by SUMO, ETAP Ontology and UNLDC to choose from. Straight ontology rendering of this example would use the same class label “SocialRole” (SUMO) / “DriverRole” (Etap Ontology) for the truck and the train drivers. This would not allow the bridging process to see the difference between the two driver entities and link them with the Train and Automobile concepts correctly.

On the other hand, UNLDC uses a very fine-grained set of concepts, that correspond to word senses from several natural languages. In particular, it includes most of the English Wordnet senses. UNLDC concepts can reflect even stylistic distinctions between members of the same Wordnet synset. Semantic classes roughly parallel to NL POS categories are imposed on top. This level of detail is an overkill for most text processing tasks except translation.

The extended ontology offers a fourth option—an optimized set of concepts, more general than lexical senses and more specific than most SUMO/Etap Ontology concepts. It is produced by an automatic procedure. We a) merge into one concept all synonymous senses regardless of the POS class of the source words, e.g. *катанье* (*act of rolling as a ball*) gets merged with *катить* (*cause to move by turning like a ball*) and all their synonyms b) merge pairs of predicates like *катить* (*cause to move by turning*) and *катиться* (*move by turning*), which differ only by the regular

transformation of their argument frames. Each new concept receives a unique OWL-compatible name and a link to an upper SUMO class or another new concept. The new concepts inherit semantic relations from UNLDC semantic network, including *is\_a* and *instance\_of*, which create subtrees of new concepts within SUMO classes, and other types translated into the bridging relation set, e.g. *катание (roll—act of rolling as a ball)* <sup>subProcess</sup> *боулинг (bowling game)* = “rolling (balls) is part of playing bowling”.

### 3.2. Relations

The relation types supported by our bridging tool are listed in **Table 1**. This set of relations can be extended through editing of the knowledge base. All relations are directed and have corresponding reverse relation types. Type labels are taken from the Etap Ontology or follow the same style.

**Table 1:** Bridging relation types

Group	Relation / Reverse relation	Examples (X—Y)	Comment
Function	hasFunction / isFunctionOf	restaurant—serve meals baker—to bake	Y is what X does or is for.
	hasRoleAt / isRoleAt	company—accountant tourists—guide cathedral—priest	Y is a function in respect to the group or object X. There may be multiple persons/objects with the same function.
	hasChief / isChiefOf	team—trainer company—director country—president	The leader of a group
Part ↔ whole	hasPart / isPartOf	room—wall	Parts that are always present
	hasOptionalPart / isOptionalPartOf	room—chandelier	Parts that may be absent
	hasDetachablePart / isDetachablePartOf	lock—key violin—bow	Required accessories that are not physically attached
	hasMember / isMemberOf	parliament—MP government—minister	All members of the group X are Y-s.
	hasSubEvent / isSubEventOf	eat—swallow	

Group	Relation / Reverse relation	Examples (X—Y)	Comment
Object ↔ matter	hasIngredient / isIngredientOf	tea—water water—oxygen	Y is one of the raw materials used and irrevocably changed or chemically bound in making X.
	hasSubstance / isSubstanceOf	table—wood ocean—water	X is a mass of pure Y. There may be parts made of other substances.
Event ↔ role	hasAgent / isAgentOf	buy—buyer fly—airplane	
	hasAgent2 / isAgent2Of	buy—seller	
	hasObject / isObjectOf	write—letter	
	hasInstrument / isInstrumentOf	eat—spoon	
	hasLocation / isLocationOf	study—school	
	hasStartingPlacePoint / isStartingPlacePoint	delivery—warehouse (as an order in a webshop)	
	hasTerminalPlacePoint / isTerminalPlacePoint	delivery—home (as an order in a webshop)	
	hasRecipient / isRecipientOf	delivery—customer (as an order in a webshop)	
	hasBeneficiary / isBeneficiaryOf	sing—audience	X has an object or message delivered to Y
	hasSource / isSourceOf	passport—Russia	
Cause ↔ result	hasResult / isResultOf	murder—death	
	newstatus-agent	compete—winner compete—loser	Y is a new social role of the agent of the event X
	newstatus-object	matriculation—student	Y is a new social role of the object of the event X

Group	Relation / Reverse relation	Examples (X—Y)	Comment
Temporal	before / after	grab—arrest—jail	Relative position at the timeline. Used in describing typical sequences of events concurrence.
	concurrence		Events that occur at the same time but neither is a subEvent of the other.
	hasTime / isTimeOf	breakfast—morning	Customary period
Misc. association	hasResident / isResidentOf	Berlin—Berliner	Resident of a place
	hasBeliever / isBeliefOf	Pope—Christianity socialist—socialism	Supporter and teaching supported
	hasAuthor / isAuthorOf	writer—book	Y is an object designed by X.
	hasMaker / isMakerOf	blacksmith—horseshoe	Y is one of many manufactured objects
	hasFrame / isFrameOf	clash—public protest study—university study—seminar	A typical scene (event, institution, proposition) associated with event X and forming its background.
	isUserOf / isUsedBy	woodcutter—ax pilot—airplane	Y is a default instrument of X e.g an attribute of profession
	hasOwner / isOwnerOf	cop—uniform	X typically possesses Y
	hasAttribute / isAttributeOf	exam—passing grade	
Cohyponyms	cohyponym	hands—legs mother—son	Only usable at low taxonomy levels.
Equivalence	SameAs	projector—apparatus shopper—client	Y is another name for X in the given domain

### 3.3. Domain descriptions

Relations and concepts are used to make semantic graphs containing generalized descriptions of different subject domains. Together such domain descriptions and the extended ontology constitute our knowledge base for bridging.

The graphs consist of triplets, where the relation labels take the place of predicates. A domain description can be saved as an RDF document. Each triplet has an additional annotation field, containing a list of domain names. Domain annotation is used to limit the scope of statements applicable only to certain parts of actual reality.

This kind of data can be imported from various domain ontologies that a) cover the domains relevant to the text to be processed, b) provide non-taxonomic relations used to construct bridges, c) have their concepts linked with the dictionary used to lemmatize/disambiguate the text. The extra fourth field (domain annotation) can be filled with the ontology's declared domain.

Our goal is to interpret texts dealing with everyday life and typical news topics: shopping, medical care, education, traffic, crime and police, sport, banking, politics etc. We model a very basic level of common background knowledge of Russian people, essential to understand contemporary Russian texts, reflecting the reality of Russia and late USSR. We did not have a suitable ontology to fill the knowledge base. The domain descriptions used in our experiments are written manually and later augmented with data from the UNLDC semantic network. We found that the amount of labor needed to describe a single domain is agreeable.

We start by enumerating a few key concepts of the domain (not including any individual persons and institutions). At the next step we link them to each other using the relations from Table 1. Later we enumerate key events concerning the domain and corresponding predicates, e.g. *matriculation*, *studying*, *reading*, *writing*, *answering*, *evaluation*, *passing exams*, *graduation*, etc. in the educational domain, list their default argument slot fillers, e.g. *student*, *professor*, *textbook*, etc. and specify typical temporal and causal relations between the events, e.g. *studying* <sup>before</sup> *passing exams*. Everything is done ad-hoc to replicate human background knowledge.

The main reason to do it is to capture typical domain-bound sequences of events (scripts) that people follow in their life and work. Scripts are presented as chains of predicate concepts placed along the abstract timeline and connected by the temporal and/or causal relations. A typical scripted activity is fishing where an angler has to *attach (a fly to the hook)* <sup>before</sup> *throw (the line into the river)* <sup>before</sup> *wait* <sup>concurrency</sup> *watch (the cork)* <sup>before</sup> *strike (fish)* <sup>hasResult</sup> *pull (the line)*, etc. This information is not present in Wordnet or general purpose ontologies, but it turned out to be very useful for bridging. It can explain the relations between participants of the events e.g. *fish* and *cork* by connecting the events they take part in *fish* <sup>isAgentOf</sup> *strike* <sup>hasResult</sup> *bob* <sup>hasAgent</sup> *cork*. UNLDC does provide some cause-result links, but they are limited to universal connections between concepts, embedded in their definitions, e.g. *to grow (vegetables)* <sup>hasResult</sup> *growth (of the plants)*.

Another reason is that manually formulated domain descriptions help to identify most important keywords making up the lexical footprint of the domain. We take pre-made concepts from UNLDC, which already have associated Russian words. Consequently, the domain description graphs are accompanied by a cloud of keywords that help to identify domain texts.

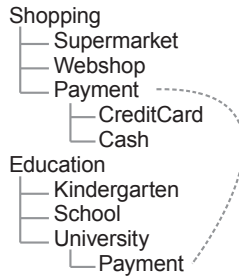


The resulting sketch description is immediately useful and can be tested with our bridging tool. We make a test run and check, if there are any important keywords/concepts missing from the domain description. This final step can be repeated many times to improve the recall of bridging relations. The typical size of a domain description is 100–1000 triplets. A fragment is shown in **Figure 1**.

Domain	Subdomains			Triplet	
EducationalProcess	SchoolEducationalInstitution		SchoolEducationalInstitution	hasChief	Headmaster
EducationalProcess	SchoolEducationalInstitution	Matriculation	OrderRequest	hasAgent	ParentGenitor
EducationalProcess	SchoolEducationalInstitution	Matriculation	OrderRequest	hasRecipient	SchoolEducationalInstitution
EducationalProcess	SchoolEducationalInstitution	Matriculation	OrderRequest	hasRecipient	Headmaster
EducationalProcess	SchoolEducationalInstitution	Matriculation	OrderRequest	hasTopic	ChildJuvenile
EducationalProcess	SchoolEducationalInstitution	Matriculation	OrderRequest	hasResult	Matriculation
EducationalProcess	SchoolEducationalInstitution	Matriculation	Matriculation	hasAgent	ChildJuvenile
EducationalProcess	SchoolEducationalInstitution	Matriculation	Matriculation	hasTerminalPoint	SchoolEducationalInstitution
EducationalProcess	SchoolEducationalInstitution	Matriculation	Matriculation	hasResult	EnrollRegister
EducationalProcess	SchoolEducationalInstitution	Matriculation	EnrollRegister	hasAgent	SchoolEducationalInstitution
EducationalProcess	SchoolEducationalInstitution	Matriculation	EnrollRegister	hasObject	ChildJuvenile
EducationalProcess	SchoolEducationalInstitution	Matriculation	EnrollRegister	Result-newstatus	Schoolchild

**Figure 1:** A few lines of a domain description showing the process of enrolling a child in a school. “The parents make an application to the school. The child gets enrolled and becomes a pupil”

The domains have their own taxonomy. Statements made in the general domains, such as *Education* and *Shopping* apply together with all statements from more specific domains, such as *University* and *Supermarket*.



**Figure 2:** A fragment of the taxonomy of domains

The KB describes a default general state of affairs. The statements in the domain descriptions are just “usually true”. No claim for universal truth can be made here. Actual truth in the real world or a fictional reality described in some concrete text has to be determined during understanding of the text or a situation in the real world. For example, the *TerminalPlacePoint* argument slot of the concept *Carrying* is always filled by some *Region*. The statement *Carrying* *hasTerminalPlacePoint* *Region* is universally true. However, in the domain of supermarkets shoppers usually carry goods to the checkout counter. Therefore, the description of the supermarket domain contains the statement *Carrying* *hasTerminalPlacePoint* *Checkout*, which is expected to be true in the domain. It makes the content of the domain descriptions unfit for a general purpose ontology, where all statements must be universally true. Instead each sub-domain section of a domain description could be viewed as a small domain ontology.

## 4. Bridging annotation

Our bridging annotation tool has two major functions:

- 1) search through a corpus and detect fragments of text that match known domains,
- 2) generate a set of potential bridging relations for the fragments found.

We use a corpus of newspaper texts as a source of examples. It consists of automatically parsed news feeds and full articles in the ETAP TGT format. The current version of the program uses only lemmatization tags. Syntax relations are used only to detect multiword expressions. It can use ETAP combinatorial dictionary entry tags for disambiguation and falls back to lemmas if they are not available. A simple TF-IDF ranked keyword search is used to extract fragments that contain higher than average density of keywords linked with available domain descriptions. The length of the fragments is not set and usually falls between 3 and 15 sentences. Each fragment gets tagged with the applicable domains with weight numbers. There is a weigh threshold which can be adjusted to tune the output between better domain detection and quantity of examples.

The bridging annotation option works as follows: an example text is scanned for any nouns, verbs and multi-word expressions (MWEs), e.g. *банк России* (*bank of Russia*), *барная стойка* (*bar stand*), present in UNLDC. The words/expressions whose lexical senses match the domains ascribed to the example text form a set of possible reference words and antecedents. The set contains all words suitable for bridging in the whole text. MWEs are represented by their head words that carry the lexical meaning of the corresponding expression.

The words are taken one by one in the linear order of the text and paired with every preceding word of the same set within a rolling window of configurable number of sentences. This creates candidate pairs of words which are turned into two sets of concepts, associated with different senses of both words/MWEs. The concepts linked with the possible reference word and mentioned in the background knowledge base are paired with all concepts linked with the possible antecedent.

Resulting pairs of concepts are fed to a search function which returns all possible bridges between them, if any. The bridges may consist of either a single semantic relation or a chain of 1–2 intermediate concepts with relations between them. Since the extended ontology has taxonomic relations between its extra concepts within SUMO/Etap Ontology classes, the search function can use `subclass_superclass_sibling` criterion [4] to improve recall and relate antecedent concepts not mentioned in the domain description.

The resulting bridges are filtered by applying such criteria as number of intermediate concepts (“bridge arcs”), number of assumed extra concepts, distance between the reference and antecedent, saliency, etc. Each confirmed antecedent word receives a list of discovered reference words and clusters of bridge links are formed.

An interesting feature of our tool is building of a possible associations list. Intermediate concepts, which are not linked with any words in the text but occur in complex bridge relation, e.g. *Policeman* in **Example 1**, are remembered. Most frequent associations are returned together with the list of discovered bridging pairs.

## 5. Problems

Like every other ontology based system, our approach falls prey to the expert knowledge input bottleneck. The amount of background knowledge provided by domain descriptions is never enough (just like with us humans) but extending them manually is a labor intensive process that gets harder with more elaborate descriptions.

The tool demonstrates domain bias. Lack of a relevant domain description provokes our tool to switch to other domains which have partially similar lexical footprint. As a result, news reports about politics and wars, for instance, get interpreted in terms of crimes and terrorism. Sport events can get mixed with theater performances because both actors and athletes play and win contests and those domains share a certain amount of keywords. This problem can be mitigated by making brief descriptions of interfering domains that cover problematic keywords.

Use of very general ontology classes in domain descriptions creates spurious assumptions, yet it is hard to avoid. For example, the domain of police work includes the concept of arresting some *Human*. It makes the system assume that every entity of a *Policeman* arrests every entity of a *Human* mentioned in the text. It is impossible to enumerate all possible objects of arresting. A text-wide resolution of identity anaphora and semantic parsing is needed to filter out bad bridges and select correct ones.

The system can make bridges that are irrelevant or redundant from a human point of view. For example, it can link words *зачетка* (*student's grade book*) and *дверь* (*door*): *grade book* <sup>hasOwner</sup> *student* <sup>isAgentOf</sup> *opening* <sup>hasObject</sup> *door*. It is hard to make a filter that would prevent such cases. Such filter must introduce the notion of the reader's intention, i.e. what we want to learn from the text.

There is no good stopping rule in assumption generation, except to ban all concepts not explicitly mentioned in the text. In a story about hijacking of a car that results in a chase, crash and explosion, the computer will happily (mis-)assume an existence of a bomb and some terrorists, because the bag of concepts (*Automobile, Impacting, Explosion, Policeman, Criminal*) has enough similarity with the domain of terrorism. This problem can also affect humans when there is no sufficient context to rule out wrong assumptions.

## 6. Evaluation

We used two different methods to assess the performance of the bridging tagger. The standard approach, which relies on precision/recall measurement against a manually tagged test corpus, hit its limits and proved to be impractical for our project. It happened because the very nature of the modeled process implies high variability and individual bias.

The second evaluation, described in [section 6.2](#), was based on manual expert assessment of the tagger output without a reference corpus. This method is better suited for evaluation of highly variable results, such as translation, which also records a particular interpretation of a text.

## 6.1. Standard approach

A pilot sample of a test corpus was made and tagged by 6 annotators. The sample consists of two short texts, 2,627 words in total, from the domains of shopping and cinema. We tried to follow formal rules similar to the ones implemented in the software, but inter-annotator agreement was so bad that we rejected the idea of making a larger corpus following the same procedure. Every annotator seemed to have a different set of associations. Out of 197 unique pairs of reference+antecedent words in the test material only 1 pair was universally accepted by all annotators and 148 pairs (75.1%) were chosen by only one person. **Table 2** provides an overview.

**Table 2:** Percentage of detected bridges vs number of annotators sharing them

Annotators	1	2	3	4	5	6
<b>Pairs %</b>	75.1%	13.2%	7.1%	3%	1%	0.5%

In most cases when several annotators selected the same bridging pair with a semantically complex relation, they interpreted it differently. For example, three annotators expressed the relation between *покупатель* (*buyer*) and *магазин* (*store*) in the following three different ways: 1) *buyer* <sup>hasLocation</sup> *store*, 2) *buyer* <sup>isAgentOf</sup> *buying* <sup>hasFrame</sup> *store*, 3) *buyer* <sup>isRecipientOf</sup> *retailing* <sup>isFunctionOf</sup> *store*. All three variants are correct and acceptable.

Identification of the words that represented referring and antecedent entities in the texts was much more uniform. 59% of the words were chosen by at least 3 annotators and 44% were chosen by more than 3 people.

This situation is well aligned with the theory of bridging explained in section 1. The Listener produces associations based on his unique background knowledge, prior information and current goals. Every instance of understanding, even by the same person, may follow a different path of associations. It is unrealistic to expect that several people will produce identical sets of bridging links.

### 6.1.1. Tagger performance

Given the same test data our bridging tool produced 532 candidate bridges. Comparison between the collective of human annotators and the program shows that the computer was able to tag 78 (39.5%) out of 197 referent+antecedent word pairs tagged by at least one human annotator. It compares favorably against the numbers of bridges found by each single human. **Table 3** shows individual recall of human annotators and the computer.

**Table 3:** Number of detected bridges per annotator out of the total pool of 197 relevant pairs

Annotators	A	B	C	D	E	F	Computer
<b>Bridges</b>	22	23	35	46	72	84	<b>78</b>
<b>%</b>	11.1%	11.6%	17.7%	23.3%	36.5%	42.6%	<b>39.5%</b>

Only one of the annotators managed to find more bridges than the system. All semantic types ascribed by the computer to the 78 bridging relations it detected were correct.

The remaining 454 links that were not confirmed by human annotators still contained some valid bridges that were overlooked by all six annotators and a lot of plausible but wrong assumptions, i.e. relations that are “usually true” but in the given context they became false.

Since people could not collectively exhaust all possible ways to interpret the test texts using an open set of bridging relations and build a gold standard corpus, we decided to analyze output of the tagger as is.

## 6.2. Expert evaluation

The second evaluation used a new set of 10 short texts, covering the domains of banking, crime/terrorism and education, 3,649 words together. They were processed automatically with different settings of the tagger and the output was manually assessed by an expert in two rounds.

Round 1 was used to evaluate the plausibility of associations produced by computer. Each bridging link was marked “good” or “bad” without reading the text itself and considering only a pair of words/expressions linked and semantic type of the proposed bridging relation. The “bad” mark was given to bridges that contradicted the expert’s knowledge of the world. This procedure evaluates the system’s ability to make good assumptions and uncovers eventual defects of the knowledge base. Here are some common cases:

- Overly general classification of some lexical senses makes some relations look improbable. For example, the text about evacuation of Russian specialists from Iraq yielded the following bad assumption: *граждане SubjectPerson isAgentOf DepartureByAircraft hasInstrument Airplane hasOwner Human некарь (baker)*. It is correct to assume that a person may own an airplane, but it is highly unlikely that a common baker would be that person. Such situations are caused by the problem of general class labels in the domain descriptions, as mentioned in section 5. On the other hand, it is hard to justify existence of e.g. a special class of “people that are rich enough to own a plane” in a general purpose ontology.
- Lack of validation with complex reasoning while building multi-arc bridges. The bridge *президент (president) Human isRecipientOf Payment hasTerminalPlacePoint BudgetFund бюджет* is wrong, even though both its parts are feasible. The knowledge that a head of state does not personally receive payments to the state budget is not available. This is mitigated by the existence of another possible relation between the same words: *президент President isChiefOf Nation isOwnerOf BudgetFund бюджет*.

Despite such problems, the system demonstrated overall high quality of generated assumptions, no less than 85% and reaching 96%, depending on the settings of the bridging tagger (see 6.2.1).

Round 2 evaluated the same bridging relations again, but this time they were put in context. The expert had to carefully read the source text and decide, whether the assumptions were correct or contradicted the text contents. Each bridging link received a second “good” or “bad” mark. False bridges were further tagged according to the nature of the error. There are two important problems:

- Lexical disambiguation errors. One or both bridged words may be labeled with wrong lexical senses, which results in false assumptions. For example, the word “life” in

*Виктор Д. перепробовал в своей жизни много профессий.*  
*Victor D. tried many occupations in his life.*

got wrongly interpreted as *Human* (as in “saved many lives”). This caused a false bridge *\*Human\**<sup>isOwnerOf</sup> *PrivilegeAdvantage* *льготами* (*social benefits*). The correct sense label here is *Life* and it does not take the relation of ownership.

- Reference errors. Many texts contain several entities of the same semantic type and the tagger can falsely relate what is said about one of them with another. For example,

*В сельской школе N4 Аксайского района ... ребята осваивают компьютер просто играючи. В селе Покровском Неклиновского района школьники уже и сами разрабатывают учебные программы.*  
*Children easily acquire skills while working with computers ... in the village school N.4 of Aksaisky district. In the settlement Pokrovskoe of Neklinovsky district pupils started to develop their own educational software.*

- It is evident from the text that the computers used by the pupils of the two schools are different. Therefore, the bridging relation between pupils from Pokrovskoe and the computers from the school N.4 is false.
- Another possible case is a bridge between two members of the same coreference group, e.g. *Директор*<sup>isChiefOf</sup> *предприятие*<sup>hasMember</sup> *руководитель* (*Director*<sup>isChiefOf</sup> *enterprise*<sup>hasMember</sup> *manager*), where director and manager are the same person.
- Reference errors constitute 15–20% of all bridges deemed to be good assumptions at Round 1. A recent paper [11] by Pagel and Rösiger applies a partially similar rule-based approach to German and confirms positive effect of using coreference resolution. They report a 3.3% improvement of the F1 measure.

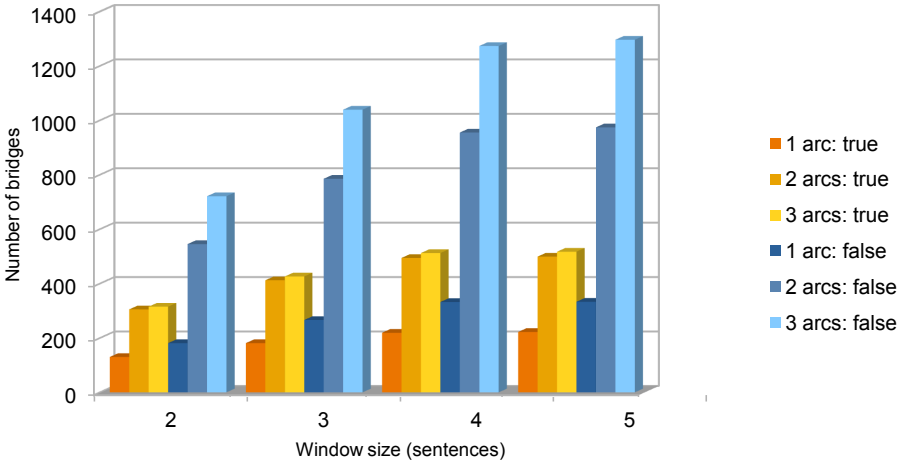
The first evaluation demonstrated that the recall of perfectly true bridges produced by the computer was on par with humans or better and the problem lies in its precision. Full evaluation of the tagger’s output during round 2 ensures quality sufficient for a “gold standard” annotation. The resulting set of texts with labeled bridging relations following our approach and manual coreference annotation can be released as a small corpus, if there is some public interest.

### 6.2.1. Tuning parameters

The tagger has several tuning parameters that influence its performance. It is possible to adjust the length of the look-back window, which tells how many preceding sentences are taken into consideration while searching for antecedents. Another parameter is the maximum number of arcs, i.e. number of intermediate concepts used to explain the semantic relation between the reference word and its antecedent.

The tested window size range is 1–5. Narrowing the window improves the ratio of perfectly good bridges to failed hypotheses and lowers the number of found antecedents. **Figure 3** shows that windows of 4 and 5 sentences applied to our test set show almost equal number of detected bridges. Smaller windows of 3 and 2 sentences caused steep decline of that number.

The bridge length could be set to 1–3 arcs. “One” means that all referent-antecedent pairs must be connected by a single semantic relation. “Two” allows a single intermediate concept. “Three” provided two intermediate concepts, one of which may be an associated entity not actually mentioned in the text. Greater bridge length brings more freedom in constructing complex bridging relations, e.g. *actor* <sup>*isObjectOf*</sup> *makeup* <sup>*hasLocation*</sup> *dressing room* <sup>*isLocationOf*</sup> *mirror*, but increases the number of false bridges.



**Figure 3:** Number of true and false bridges in relation to 1) number of arcs and 2) antecedent search window size

The optimal balanced combination for a regular text seems to be window size 4 with 2 arcs. Three arcs may help when no bridges are found in the regular way. This can happen with a text that tries to say more with fewer words and leaves out more information than average.

### 6.2.2. Results

There are several combinations of settings that provide best results in one or another aspect. [Table 4](#) sums them up. The F1 score calculation bears a special note because, as explained in [section 6.1](#), we do not know the total number of all conceivable true bridges (relevant samples). The expert evaluation procedure does not consider any bridges except those generated by the tagger itself. Therefore the number of relevant samples used to determine the recall is always equal to the total amount of true bridges marked during round 2 (see [section 6.2](#)).

Comparing our results with other studies cannot be straightforward because of differences in methods and bridge identification criteria. In particular, paper [9] by Roitberg and Khachko reports F1 measure 0.65 for Russian with a completely different approach based on syntactic criteria that does not explain the semantics of the relation between reference words and antecedents and covers only nouns, while we also include verbs. Pagel and Rösiger apply a closer approach to German and report F1 measure of 11.1% (no coreference) to 14.1% (with coreference).

**Table 4:** Tagger performance using different settings

	Arcs	Window	% plausible bridges (round 1)	% true bridges (round 2)	F1	# true bridges	# false bridges
<b>Highest recall</b>	3	5	85.4	22.04	0.361	<b>801</b>	2561
<b>Balanced</b>	2	4	90.47	30.5	<b>0.44</b>	633	1390
<b>Best precision</b>	1	2	96.57	<b>37.32</b>	0.255	156	258

## 7. Conclusion

We develop an extensible semantic knowledge base geared towards bridging resolution. It opens up a possibility to explore semantic approaches in Russian and use richer background information than previous studies. All established bridging relations receive a semantic interpretation, which is not limited by a fixed set of pre-defined labels. Flexibility granted by combining multiple relations and intermediate concepts in “multi-arc” bridges allows to represent complex associations but creates problems mentioned in sections 5 and 6. Support for complex semantic relations is an important feature of our bridging tagger.

Most authors in the field narrow down the problem by imposing artificial constraints on the types of bridges they consider. For example, papers [6], [7] limit the relation types to part-whole. Paper [9] ignores semantic types but imposes a syntactic limitation. It greatly simplifies formal evaluation but results in ignoring most of the possible bridges in any text. Such works fail to cover the full scope of the studied phenomenon. We prefer to look at the problem in a more general and holistic way and build a model which covers wider range of possible implicit relations than previous studies following the semantic approach.

The mental process explained in [section 1](#) is always subjective and implies great variability. Different people see different relations because 1) they have different background knowledge (it includes education, prior experience, cultural bias, etc.), 2) different intentions and 3) the space of possible implicit relations is so vast that no one can exhaust it. In our study we cap variability which stems from factors 1 and 2. We look for relations that are based on an explicitly formulated knowledge base (KB) and follow explicitly defined rules. However, experiments showed that even this controlled space of possible relations is bigger than six expert annotators could collectively cover during the first evaluation. It is very difficult to make a “gold standard” corpus with a rich set of semantic relation types and complex “muti-arc” bridges, because annotators naturally produce different interpretations of the same text. There is no way to be certain that the reference tagging is complete and any other bridges in the same corpus will be wrong.

Output of our bridging tool is hard to rate using the traditional method which requires comparison with a reference corpus because bridging belongs to the class of problems that allow many alternative solutions, just like translation. This is why the alternative method of expert evaluation is more feasible.



We can confirm that the list of most useful features for bridging in paper [11] is true. It includes 1) semantic connectivity and 2) distance between reference word and its antecedent. It is also worth to explore the possibilities of populating the domain descriptions and ranking plausibility of different alternative antecedents by ML.

## 8. Acknowledgments

This study was supported by the Russian Science Foundation (grant No. 16-18-10422)

## References

1. Boguslavsky I. M., Dikonov V. G., Frolova T. I., Iomdin L. L., Lazurski F. V., Rygaev I. P., Timoshenko S. P. (2016) Plausible Expectations-Based Inference for Semantic Analysis. Proceedings of the 2016 International Conference on Artificial Intelligence (ICAI'2016). USA: CSREA Press, 2016. pp. 477–483. ISBN: 1-60132-438-3.
2. Clark, H. H. (1975) Bridging. In Proceedings of the 1975 workshop on Theoretical issues in natural language processing, Association for Computational Linguistics, pp. 169–174.
3. Dikonov V. G. (2013) Development of lexical basis for the Universal Dictionary of UNL Concepts; Proceedings of the International Conference "Dialogue". Issue 12(19), Moscow, RGGU Publishers. P 212–221.
4. Fan, J., Barker, K., & Porter, B. (2005) Indirect anaphora resolution as semantic path search. In Proceedings of the 3rd international conference on Knowledge capture ACM, pp. 153–160.
5. Hou, Y., Markert, K., & Strube, M. (2013) Global Inference for Bridging Anaphora Resolution. In HLT-NAACL pp. 907–917.
6. Lassalle, E., & Denis, P. (2011) Leveraging different meronymy discovery methods for bridging resolution in French. In Discourse Anaphora and Anaphora Resolution Colloquium, Springer Berlin Heidelberg, pp. 35–46.
7. Poesio, M., Mehta, R., Maroudas, A., & Hitzeman, J. (2004-B) Learning to resolve bridging references. In Proceedings of the 42nd Annual Meeting on Association for Computational Linguistics, Association for Computational Linguistics, p. 143.
8. Recasens M., Martí M. A., Taulé M. (2007) Text as scene: Discourse deixis and bridging relations. *Procesamiento del lenguaje natural*, 39:205–212.
9. Roitberg A. M., Khachko D. V. (2017) Bridging Anaphora Resolution for the Russian Language. Proceedings of the International Conference "Dialogue 2017", Moscow, May 31 — June 3, 2017.
10. Zikánová Š., Hajičová E., Hladká B., Jínová P., Mírovský J., Nedoluzhko A., Poláková L., Rysová K., Rysová M., Václ J. (2015) Discourse and Coherence. From the Sentence Structure to Relations in Text, volume 14 of Studies in Computational and Theoretical Linguistics. Charles University in Prague, Praha, Czechia.
11. Pagel J., Rösiger I. (2018) Towards Bridging Resolution in German: Data Analysis and Rule-based Experiments. In Proceedings of the Workshop on Computational Models of Reference, Anaphora, and Coreference (CRAC), NAACL, New Orleans, US.