# Digraph Analysis of Dictionary Preposition Definitions

Kenneth C. Litkowski
CL Research
9208 Gue Road
Damascus, MD 20872
ken@clres.com

## Abstract

We develop a model of preposition definitions in a machine-readable dictionary using the theory of labeled directed graphs and analyze the resulting digraphs to determine a primitive set of preposition senses. We characterize these primitives and show how they can be used to develop an inheritance hierarchy for prepositions, representing the definitions by a type and slots for its arguments. By analyzing the definitions, we develop criteria for disambiguating among the highly polysemous primitives. We show how these criteria can be used in developing the inheritance hierarchy and how they may be used in assigning theta roles to the objects of transitive verbs. Finally, we describe the use of the disambiguation criteria to parse and represent the meaning of the prepositions as used in encyclopedia articles.

## 1    Introduction

Prepositions have generally been viewed as function words to be discarded in many natural language processing applications. However, prepositions have considerable importance as identifiers of semantic relations tying various elements of a sentence together. Since many prepositions are highly polysemous, it is necessary to develop a finer-grained analyses of their meanings so that the semantic relations can be more accurately identified.

We have modeled preposition definitions in a dictionary using directed labeled graphs (digraphs) to identify primitive preposition senses. We have used the definitions from a machine-readable version of a comprehensive English dictionary.

## 2    Modeling Preposition Definitions

A preposition is "a word governing, and usually preceding, a noun or pronoun and expressing a relation to another word or element in the clause." The definition of a preposition takes two principal forms: (1) a usage expression characterizing the relation or (2) an expression that can be substituted for the preposition. A substituting preposition definition usually consists of a prepositional phrase (including both a preposition and a noun phrase) and a terminating preposition (e.g., for *around*, one definition is "on every side of").

### 2.1    Headwords as Digraph Nodes

A digraph consists of nodes and directed arcs between the nodes. In general, an arc should correspond to a transitive relation. Modeling a dictionary with a digraph entails assigning an interpretation to the nodes and arcs. For our initial model, we subsume all the definitions of a preposition as one node in the digraph, labeled by the preposition. An arc is drawn from one node (e.g., *of*) to another (e.g., *around*) if the preposition represented by the first node **contributes a typed meaning component with an open slot to** the preposition represented by the second node, e.g., "**part-of of around**" would arise from the definition of *around* ("on every side of").

Loosely, for our purposes, the terminating preposition acts as a genus term in an ISA hierarchy and makes it possible to use the results from digraph theory to analyze the relationships between definitions. In particular, digraph analysis identifies definitional cycles and "primitives" and arranges the nodes into an inheritance hierarchy. When a dictionary is modeled like this, digraph theory (Harary, et al. 1965) indicates that there is a "basis set" of nodes, which may be viewed as a

set of primitives.[1]

Many prepositions are not used as the final preposition of other preposition definitions (specifically, their nodes have an outdegree of 0). These are the leaves of the inheritance hierarchy. When these are removed from the dictionary, other prepositions will come to have outdegree 0, and may in turn be removed. After all such iterations, the remaining nodes are "strongly connected", that is, for every node, there is a path to each other node; a strong component is an equivalence class and corresponds to a definitional cycle.

Each strong component may now be viewed as a node. Some of these nodes also have the property that they have outdegree 0; these strong component may also be removed from the dictionary. This may introduce a new round where individual nodes or strong components have outdegree 0 and hence may be removed from the dictionary.

After all removals, what is left is a set of one or more strong components, each of which is unreachable from the other. This final set is viewed as the set of primitives. What this means is that we have converted the preposition dictionary into an inheritance hierarchy. If we can characterize the meanings of the primitives, we can then inherit these meanings in all the words and definitions that have been previously been removed.

## 2.2    Definitions as Digraph Nodes

This model of prepositions is very coarse, lumping all senses into one node. Having reduced the set of prepositions with this model, we can initiate a new round of digraph analysis by disambiguating the final preposition. In this new model, each node represents a single sense and the arc between two nodes indicates that one specific sense is used to define one specific sense of another word (i.e., "contributes a typed meaning component with an open slot to").

---

[1]The determination of the "basis set" of a digraph is NP-complete However, as pointed out in (Litkowski, 1988), this process will not involve millions of nodes. In our implementation of the algorithm for finding strong components (Even 1980), the digraph analysis of prepositions takes less than two seconds.

With this new model, we can enter into a further round of digraph analysis. In this round, which proceeds as above, instead of a set of primitive prepositions, the outcome will be a set of primitive preposition definitions. However, as mentioned above, preposition definitions come in two flavors. The usage expressions are lumped into the digraph analysis when a node corresponded to all definitions, but they do not do so in the definition digraph analysis.

## 3    NODE Prepositions

As the data for the digraph analysis, we began with the 155 prepositions identified in a machine-readable dictionary (The New Oxford Dictionary of English, 1998) (NODE). Additional prepositions are found as unmarked phrases under noun or adjective headwords, but not so labeled, e.g., *in spite of* under the headword *spite*. To find these prepositions, we developed a more rigorous specification of a **preposition signature**. A preposition definition is either (1) a preposition; (2) a prepositional phrase + a preposition; (3) (an optional leading string) + a transitive present participle; or (4) a leading string + an infinitive of a transitive verb. This led to the addition of 218 phrasal prepositions, for a total of 373 entries, with 847 senses, shown in the Appendix.

We may have missed other subsenses that have a preposition signature. In all likelihood, these patterns would enter the digraph analysis as nodes with outdegree 0 and hence would be eliminated in the first stage of the primitive analysis.

## 3.1    Substitutable Definitions

Most preposition definitions are in a form that can be substituted for the preposition. For a sense of *against* ("as protection from"), with an example "he turned up his collar against the wind", the definition can be fully substituted to obtain "he turned up his collar as protection from the wind."

The preposition definitions were parsed, putting them into a generic sentence frame, usually "Something is [prepdef] something." For example, the definition of *ahead of* ("in store for") would be parsed as "Something is in store for something."

For definitions with a selectional restriction on the preposition's object (identifiable by a parenthesized expression in the definition), the parentheses were removed in the sentence frame, e.g., *above* ("higher than (a specified amount, rate, or norm)") would be parsed as "Something is higher than a specified amount, rate, or norm."

The parse tree would then be analyzed to obtain the final preposition, treated as the hypernym. For definitions containing a verb at the end, e.g., another sense of *above* ("overlooking", parsed as "Something is overlooking something") would yield "overlooking" as the hypernym.

## 3.2 Usage Note Definitions

Many preposition definitions are not substitutable, but rather characterize how the preposition is used syntactically and semantically. One sense of *of* ("expressing the relationship between a part and a whole") characterizes the semantic relationship (in this case, the partitive). One of its subsenses ("with the word denoting the part functioning as the head of the phrase") indicates syntactic characteristics when this sense is used. These definitions are not parsed and do not lead to the identification of hypernyms. As shown below, these definitions will emerge as the primitives.

## 3.3 Definition Modifications

The automatic generation of preposition hypernyms was less than perfect. We examined each definition and made various hand modifications. Our editing process included hand entry of hypernyms: adding or modifying automatically generated hypernyms, making hypernymic links for "non-standard" entries (e.g., making *upon* the hypernym of *'pon*), and creating hypernymic links from a subsense to a supersense

## 4 Digraph Analysis Results

The digraph analysis described above eliminated 309 of the 373 entries. The remaining 64 entries were grouped into 25 equivalence classes, as shown in Table 1 and portrayed in Figure 1 in the appendix. Figure 1 shows how these strong components are related to one another. The strong components highlighted in the table are primitives. Seven of the primitive strong components (**in**, **of**, **than**, **as**, **from**, **as far as**, and **including**) have paths into strong component 12. Strong components 14 to 18 arise essentially from the primitive strong component **of**. The eighth strong component (23) and other entries defined by words in this class exist somewhat independently.

It would seem that the largest strong component (12, with 33 entries) should be broken down into smaller classes; this would occur in the sense-specific digraph analysis. Specialized senses of **with**, **by**, **to**, **for**, and **before** give rise to definitional cycles within this strong component.

In addition to the strong components shown above, 62 non-prepositional primitives have been identified. The first 42 of these primitives were used in defining entries that were removed in the first phase of the digraph analysis. The 20 beginning with **affect** were used in defining entries in the primitive strong components.

There are 155 preposition senses (out of 847) that are defined solely with usage notes. Of these, 71 are subsenses, leaving 74 senses in 26 entries (as shown in Table 3) that can be considered the most primitive senses and deserving initial focus in attempting to lay out the meanings of all preposition senses.

## 5 Interpretation of Results

The digraph analysis of prepositions provides additional perspectives in understanding their meanings and their use. To begin with, the analysis enables us to identify definitional cycles and move toward the creation of an inheritance hierarchy. The large number of senses that have verb hypernymic roots indicates a close kinship between prepositions and verbs, suggesting that a verb hierarchy may provide an organizing principle for prepositions (discussed further below). The large number of senses rooted in usage notes, which essentially characterize how these senses function, encapsulates the role of prepositions as "function words;" however, as described below, these functions are not simply syntactic in nature, but also capture semantic roles.

| Table 1 Strong Components | |
|---|---|
| **Entries** | |
| 1 | over, above |
| 2 | against |
| 3 | but |
| 4 | along |
| 5 | on |
| 6 | via, by way of |
| 7 | through |
| 8 | touching |
| 9 | until, up to |
| 10 | below, underneath |
| 11 | inside, within |
| 12 | in favour of, along with, with respect to, in proportion to, in relation to, in connection with, with reference to, in respect of, as regards, concerning, about, with, in place of, instead of, in support of, except, other than, apart from, in addition to, behind, beside, next to, following, past, beyond, after, to, before, in front of, ahead of, for, by, according to |
| 13 | in |
| 14 | across |
| 15 | by means of |
| 16 | in the course of |
| 17 | during |
| 18 | on behalf of |
| 19 | of |
| 20 | than |
| 21 | as |
| 22 | from |
| 23 | by reason of, because of, on account of |
| 24 | as far as |
| 25 | including |

| Table 2 |
|---|
| **Non-Prepositional Primitives** |
| embrace, incur, lose, injure, called, taking into consideration, taking account of, help, guide, interest, impress, providing, exceeding, requiring, needing, losing, injuring, restrain, see, attaining, support, defend, award, subtracting, nearly, cover, exclude, involving, undergoing, do, encircle, separating, taking into account, concerns, lacking, encircling, hit, achieving, using, involve, **affect**, overlooking, awaiting, having, being, reach, preceding, constituting, affecting, representing, facing, promote, obtain, containing, approaching, almost, taking, complete, reaching, concern, possessing, wearing |

The frequency with which the various prepositions are used as hypernyms in defining other prepositions reveals something about their relative importance. The most frequent hypernyms are **of** (175), **to** (74), **than** (45), **with** (44), **by** (39), **from** (30), **for** (22), **as** (20), and **in** (12). These prepositions correspond to the primitives identified

in Table 1, as well as those with the largest number of usage notes shown in Table 3.

| Table 3 |
|---|
| **Usage-Note Primitives** |
| about (2), as (1), as from (1), as of (1), at (6), between (1), but (1), by (7), for (6), from (11), in (7), in relation to (1), into (8), like (1), of (9), on (1), on the part of (1), out of (1), over (1), than (2), this side of (1), to (7), towards (1), under (1), up to (1), with (4) |

On the other hand, the relative frequencies may not correspond well with our intuitions about a semantic classification of prepositions. (Quirk, et al. 1985) give the greatest prominence to spatial and temporal meanings, followed by the cause/purpose spectrum, the means/agentive spectrum, accompaniment, and support and opposition, and finally, several miscellaneous categories. In the semantic relations hierarchy of the Unified Medical Language System (UMLS) (Unified Medical Language System 2002), five general types of associations are identified: physical, spatial, functional (causal), temporal, and conceptual. The leaves of the UMLS hierarchy are realized as verbs, but have a strong correspondence to the classification in (Quirk, et al. 1985).

In our identification of primitives, including the usage notes, spatial and temporal senses are conspicuously reduced in significance, while a comparative term (**than**) seems to have a much greater presence. The explanation for these two observations is that (1) many of the basic spatial and temporal prepositions were located in the largest strong component (12 in Table 1) or were derived from it and (2) many of the senses of these spatial and temporal prepositions have "than" as hypernym. This suggests that a considerable amount of the meaning of such prepositions lie principally in describing relative position in a spatio-temporal continuum.

# 6 Developing an Inheritance Hierarchy

As suggested earlier, the next stage of digraph analysis involves disambiguating the hypernymic preposition, so that individual nodes of the digraph represent senses or concepts. As suggested in (Litkowski, 1978), these nodes will consist of a

gloss and the various lexicalizations the concept, much like the synsets in WordNet (Fellbaum 1998). A prototypical case would be strong component 23 which may be lexicalized as {**by reason of, because of, on account of**}; our analysis suggests that, in this case, some further characterization of the usage of this concept by the lexicographers would be desirable, since otherwise we have only a vicious definitional cycle.

The creation of the hierarchy would involve assigning a label or type to the individual concepts and then characterizing the information that is to be inherited. The typology can be developed from the bottom up, rather than developing some *a priori* structure. In other words, since the digraph analysis has identified primitive senses, these provide an appropriate starting point. Each sense can be examined on its own merits with an initial assignment of a type and later examination of the full set of primitives for organization into a data-driven set of types and subtypes.

As to what gets inherited, we begin with the fact that in general, each preposition has two arguments, **arg1** (the object of the preposition) and **arg2** (the attachment point, or head, of the prepositional phrase). We may take these as the two slots associated with each representation and we may give the slots names according to the type (or just implicitly understand that a type has particular types of arguments). When considering the general structure of a non-primitive preposition definition (a prepositional phrase with an ending preposition), the NP of the prepositional phrase is the value of **arg2**. This value will be useful in disambiguating the hypernymic preposition (as described in the next section). In considering the slots for prepositions whose hypernym is a verb (as identified in Table 2), **arg1** will be the object of the verb.

## 7 Definition Use

To describe the process by which preposition senses will be disambiguated and also how the representations of their meaning will be used in processing text, Table 4 shows the definitions for "of", the most frequently used hypernym and perhaps the second most frequent word in the English language. In the table, we have assigned a type to each of nine main senses. In the definition column, the main sense is given first, with any subsenses given in parentheses, separated by semicolons if there is more than one subsense.

First, we consider the disambiguation of hypernyms in preposition definitions, that is, those whose final word is "of". One sense of "after" is "in imitation of" (e.g., "a mystery story after Poe"); examining the table suggests that this is a **deverbal** use of "of", where the object of "after" would be the object of the underlying verb of "imitation", so that when "after" is used in this sense, its **arg1** is the object of the verb "imitate". A sense of "on behalf of" is "as a representative of"; this is the **partitive** sense, so that **arg1** of "on behalf of" is a "whole". Finally, one sense of "like" is "characteristic of"; this is the **predicative deverbal**. Carrying out this process throughout the preposition definitions will thus enable us not only to disambiguate them, but also to identify characteristics of their arguments when the prepositions they define are used in some text.

In addition, prepositions very often appear at the end of the definitions of transitive verbs. For example, one sense of "accommodate" is "provide lodging or sufficient space for", where the sense of "for" is "to the benefit of", where "of" is used in the **genitive** sense (i.e., "someone's or something's benefit). With this interpretation, we can say that the object of "accommodate" is a benefactive and that a benefactive role has been lexicalized into the meaning of "accommodate". With disambiguation of the final preposition in such definitions, we will be able to characterize the objects of these verbs with some theta role.

The ultimate objective of this analysis of prepositions is to be able to characterize their occurrences in processing text. Specifically, we would like to disambiguate a preposition, so that we can assign each instance a type and characterize its arguments. In this way, processing a text would identify the semantic relations present in the text. We have performed some initial investigations into the viability of this goal.

We have begun implementing a discourse analysis of encyclopedia articles. At the base of this analysis, we are identifying and characterizing discourse entities, essentially the noun phrases. Our

| Table 4. Definitions of "of" | |
|---|---|
| **Type** | **Definition (Subsense(s))** |
| 1. Partitive | relationship between a part and a whole (part functioning as head; after a number, quantifier, or partitive noun, with the word denoting the whole functioning as the head of the phrase) |
| 2. Scale-Value | relationship between a scale or measure and a value (an age) |
| 3. Genitive | association between two entities, typically one of belonging (relationship between an author, artist, or composer and their works collectively) |
| 4. Direction | relationship between a direction and a point of reference |
| 5. Hypernym | relationship between a general category and the thing being specified which belongs to such a category (governed by a noun expressing the fact that a category is vague) |
| 6. Deverbal | relationship between an abstract concept having a verb-like meaning and (a noun denoting the subject of the underlying verb; the second noun denotes the object of the underlying verb; head of the phrase is a predicative adjective) |
| 7. Indirect Object | relationship between a verb and an indirect object (a verb expressing a mental state; expressing a cause) |
| 8. Substance | the material or substance constituting something |
| 9. Time | time in relation to the following hour |

analysis includes identification of the syntactic role and semantic type of the noun phrases, along with attributes such as number and gender. The analysis also includes resolution of anaphora, coreferences, and definite noun phrases. The modules analyzing the discourse entities come after a full parse of each sentence. We have now introduced a module to examine prepositions and build semantic relations. The results of these analyses generate an XML representation of discourse segments, discourse entities, and semantic relations, each with an accompanying set of attributes.

Our implementation of the semantic relation module has identified several issues of interest. First, the characterization of the semantic relation needs to come after the object of the prepositional phrase has been analyzed for its discourse entity properties. For example, if the object is an anaphor, the antecedent needs to be established. Second, the attachment points of the prepositional phrase need to be identified; our parser establishes a stack of possible attachment points (index positions in the sentence), with the most likely at the top of the stack. (Attachment tests could be implemented at this point, although we have not yet done so.) The attachment point is necessary to identify the arguments to be analyzed.

Having identified the arguments, the information subject to analysis includes the literal arguments (both the full phrase and their roots), the parts of speech of the arguments, any semantic characterizations of the arguments that are available (such as the WordNet file number), and access to the dictionary definitions of the root heads. The analysis for the semantic relation is specific to the preposition. We are encoding a semantic relation type and one or more tests with each sense. Some of these tests are simple, such as string matches, and others are complex, involving function calls to examine semantic relationships between the arguments.

In the case of "of", the first test was whether **arg2** is an adjective, in which case we assigned a **type** of "predicative". Next, if **arg2** was a vague general category ("form", "type", or "kind"), we set the type to "hypernymic". If neither of these conditions was satisfied, we looked up the root of **arg2** in WordNet to determine if the word had a "part-of" relation (resulting in a "partitive" type) or "member-of" relation (resulting in a "hypernymic" type). If a type had not been established by this point, we used the WordNet file number to establish an intermediate type. Thus, for example, if **arg2** was an "action" or "process" word, we set the type for the semantic relation to "deverbal"; for a "quantity", we set the type to "partitive". Finally, we can make use of the definition for **arg1** (parsed to identify its hypernym) to determine if **arg2** is the hypernym of **arg1**. When these criteria are not sufficient, we label the type "undetermined".

In our encyclopedia project, we parse and process the articles to generate XML files. We then apply an XSL transformation to extract all the semantic relations that were identified, including the preposition, the type assigned, and the values of
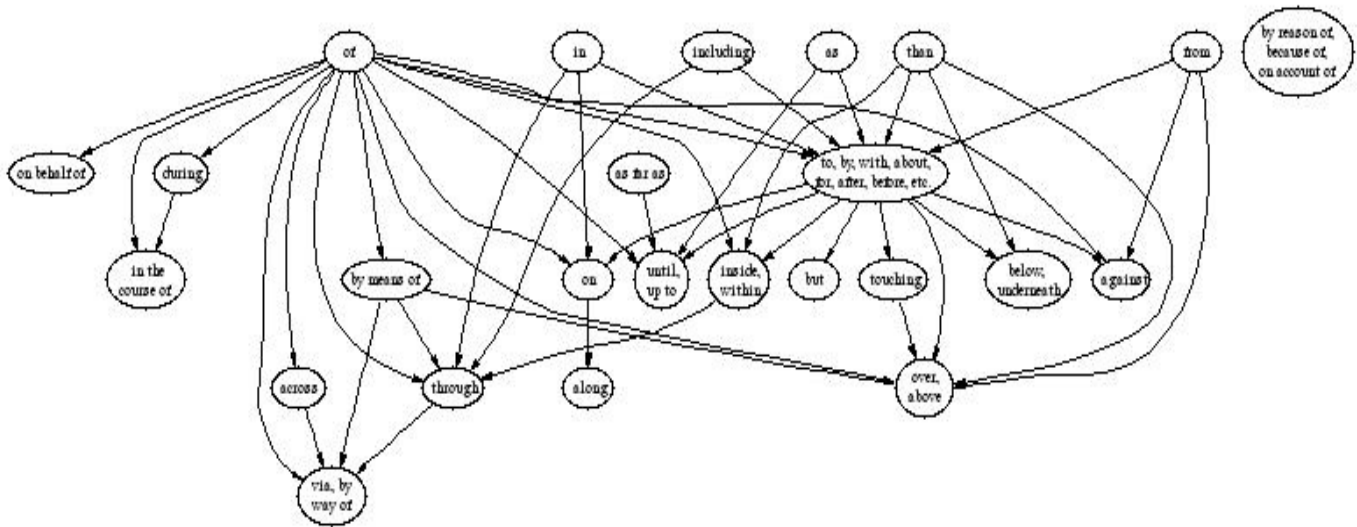
**Figure 1**. Basis Digraph of NODE Prepositions

**arg1** and **arg2**. We can sort on these fields to facilitate analysis of our success and to identify situations in need of further work.

After the initial implementation, we were able to assign semantic relations to 50 percent of the instances of "of", although many of these were given incorrect assignments. However, the method is useful for identifying instances for which improved analysis is necessary. For example, we can identify where improved characterization of discourse entities is needed, or where additional lexical information might be desirable (such as how to identify a partitive noun).

## 8   Conclusions and Further Work

We have shown that a digraph analysis of preposition definitions provides a useful organizing principle for analyzing and understanding the meanings of prepositions. The definitions themselves provide sufficient information for developing an inheritance hierarchy within a typed-feature structure arrangement and also provide a rich set of criteria for disambiguating among the many senses. By incorporating these criteria in a text processing system, it is possible to develop semantic triples that characterize intrasentential relationships among discourse entities. Further, the characterization of meanings may prove useful in identifying theta roles implied by the ending prepositions of transitive verb definitions

Much work remains to be done to develop the full set of information for all prepositions. We believe we have established a suitable framework for carrying out this work.

## References

Even, S. (1980). *Graph Algorithms.* Rockville, MD: Computer Science Press.

Fellbaum, C. (1998). (Ed.), *WordNet: An Electronic Lexical Database* (pp. 69-104). Cambridge, Massachusetts: The MIT Press.

Harary, F., Norman, R. Z., & Cartwright, D. (1965). *Structural models: An introduction to the theory of directed graphs.* New York: John Wiley and Sons, Inc.

Litkowski, K. C. (1988). On the search for semantic primitives. *Computational Linguistics, 14*(1), 52.

Litkowski, K. C. (1978). Models of the semantic structure of dictionaries. *American Journal of Computational Linguistics, Mf.81,* 25-74.

*The New Oxford Dictionary of English.* (1998) (J. Pearsall, Ed.). Oxford: Clarendon Press.

Quirk, R., Greenbaum, S., Leech, G., & Svartik, J. (1985). *A comprehensive grammar of the English language.* London: Longman.

*Unified Medical Language System Knowledge Sources.* (13th ed.). (2002). Bethesda, MD: National Library of Medicine.

| | | | | |
|---|---|---|---|---|
| #à la | beyond | in connection with | modulo | saving |
| #'cept | but | in consideration of | more like | short for |
| #'gainst | but for | in contravention of | near | short of |
| #'mongst | by | in consequence of | near to | shot through with |
| #'pon | by courtesy of | in default of | neath | sick and tired of |
| a cut above | by dint of | in despite of | next | since |
| abaft | by force of | in excess of | next door to | strong on |
| abaht | by means of | in face of | next to | subsequent to |
| aboard | by reason of | in favor of | nigh | than |
| about | by the hand of | in favour of | none the worse for | thanks to |
| above | by the hands of | in front of | not a patch on | the better part of |
| absent | by the name of | in honor of | not someone's idea of | this side of |
| according to | by the side of | in honour of | nothing short of | thro' |
| across | by virtue of | in keeping with | notwithstanding | through |
| afore | by way of | in lieu of | o' | throughout |
| after | care of | in light of | o'er | thru |
| after the fashion of | chez | in line with | of | thwart |
| against | circa | in memoriam | of the name of | till |
| agin | come | in need of | of the order of | to |
| ahead of | complete with | in obedience to | off | to the accompaniment of |
| all for | con | in peril of | offa | to the exclusion of |
| all of | concerning | in place of | on | to the tune of |
| all over | considering | in proportion to | on a level with | to windward of |
| along | contrary to | in re | on a par with | together with |
| along of | counting | in reference to | on account of | touching |
| along with | courtesy of | in regard to | on behalf of | toward |
| alongside | cum | in relation to | on pain of | towards |
| amid | dan | in respect of | on the order of | uh |
| amidst | dehors | in restraint of | on the part of | under |
| among | depending on | in sight of | on the point of | under pain of |
| amongst | despite | in spite of | on the right side of | under cover of |
| an apology for | despite of | in succession to | on the score of | under sentence of |
| anent | down | in support of | on the strength of | under the auspices of |
| anti | due to | in terms of | on the stroke of | under the banner of |
| anything like | during | in the act of | on the wrong side of | under the baton of |
| anywhere near | ere | in the cause of | on top of | under the heel of |
| apart from | even as | in the course of | onto | underneath |
| apropos | every bit as | in the face of | opposite | unknown to |
| around | ex | in the fashion of | other than | unlike |
| as | except | in the gift of | out | until |
| as far as | except for | in the grip of | out for | unto |
| as for | excepting | in the heat of | out of | up |
| as from | excluding | in the interest of | out of keeping with | up against |
| as of | exclusive of | in the interests of | out of line with | up and down |
| as regards | failing | in the light of | outa | up before |
| as to | following | in the matter of | outboard of | up for |
| aside from | for | in the midst of | outside | up on |
| aslant | for all | in the name of | outside of | up to |
| astraddle | for the benefit of | in the nature of | outta | up to one's elbows in |
| astride | for the love of | in the pay of | outwith | up to one's neck in |
| at | forby | in the person of | over | upon |
| at a range of | forbye | in the shape of | over against | upward of |
| at peril of | fore | in the teeth of | over and above | upwards of |
| at right angles to | fornenst | in the throes of | overtop | v |
| at the expense of | fornent | in the way of | owing to | v. |
| at the hand of | frae | in token of | pace | versus |
| at the hands of | from | in view of | past | via |
| at the heels of | give or take | in virtue of | pending | vice |
| at the instance of | given | in with | per | vis-à-vis |
| at the mercy of | gone | including | plus | vs |
| athwart | good for | inclusive of | preparatory to | while |
| atop | having regard to | inshore of | previous to | with |
| back of | head and shoulders above | inside | prior to | with a view to |
| bar | in | inside of | pro | with one eye on |
| bare of | in accord with | instead of | pursuant to | with reference to |
| barring | in addition to | into | qua | with regard to |
| because of | in advance of | into the arms of | re | with respect to |
| before | in aid of | irrespective of | regarding | with the exception of |
| behind | in answer to | less | regardless of | withal |
| below | in back of | like | relative to | within |
| beneath | in bed with | little short of | respecting | within a measurable distance of |
| beside | in behalf of | mid | round | within sight of |
| besides | in case of | midst | round about | without |
| between | in common with | minus | sans | |
| betwixt | in company with | mod | save | |

**Table A-2 Prepositions in the New Oxford Dictionary of English**