

# Target word prediction and paraphasia classification in spoken discourse

Anonymous NAACL submission

## Abstract

We present a system for automatically detecting and classifying phonologically anomalous productions in the speech of individuals with aphasia. Working from transcribed discourse samples, our system identifies neologisms, and uses a combination of string alignment and language models to produce a lattice of plausible words that the speaker may have intended to produce. We then score this lattice according to various features, and attempt to determine whether the anomalous production represented a phonemic error or a genuine neologism. This approach has the potential to be expanded to consider other types of paraphasic errors, and could be applied to a wide variety of screening and therapeutic applications.

## 1 Introduction

Aphasia is a neuropsychological condition in which an individual's ability to produce or comprehend language is compromised. It can be caused by a number of different underlying pathologies, but can generally be traced back to physical damage to the individual's brain: tissue damage following a stroke or other ischemic event, lesions caused by a traumatic brain injury or infection, etc. It can also be associated with various neurodegenerative diseases, as in the case of Primary Progressive Aphasia. According to the National Institute of Neurological Disorders and Stroke, approximately 1,000,000 people in the United States suffer from aphasia (National Institute of Neurological Disorders and Stroke, 2014), and aphasia is a common consequence of strokes

(prevalence estimates for aphasia among stroke patients vary, but are typically in the neighborhood of 30% (Engelter et al., 2006)).

*Anomia* is the inability to access and retrieve words during language production, and is a common manifestation of aphasia [cite]. An anomic individual will experience difficulty recalling words and naming items, which can cause substantial difficulties in day-to-day communication. Additionally, long-term communication difficulties associated with aphasia in general have been shown to affect the psychological well-being of people with aphasia as well as their families (Cruice et al., 2003; Gaete and Bogousslavsky, 2008; van Dijk et al., 2015).

The process of screening for, diagnosing, and assessing anomia is typically manual in nature, and requires substantial time, labor, and expertise. Compared to other neuropsychological assessment instruments, aphasia-related assessments are particularly difficult to computerize, as they typically depend on subtle and complex linguistic judgments about the phonological and semantic similarity of words, and also require the examiner to interpret phonologically disordered speech. Furthermore, the most commonly used assessments focus for practical reasons on relatively constrained tasks such as picture naming, which may lack ecological validity (Mayer and Murray, 2003).

In this work, we describe an approach to automatically detecting and analyzing certain categories of word production errors characteristic of anomia in connected speech. Our approach is a first step towards an automated anomia assessment tool that

could be used cost effectively in both clinical and research settings,<sup>1</sup> and could also be applied to other disorders of speech production. The method we propose uses statistical language models to identify possible errors, and employs a phonologically-informed edit distance model to determine phonological similarity between the subject’s utterance and a set of plausible “intended words.” We then apply machine learning techniques to determine which of several categories a given erroneous production may fall into. We show XXXX

### 1.1 Anomia and Paraphasias

Anomia can take several different forms, but in this work we are concerned with *paraphasias*, which are unintended errors in word production.<sup>2</sup> There are several categories of paraphasic error. *Semantic errors* arise when an individual unintentionally produces a semantically-related word to their original, intended word (their “target word”). A classic semantic error would be saying “cat” when one intended to say “dog.”

*Phonemic* (sometimes called “formal”) errors occur when the speaker produces an unrelated word that is *phonemically related* to their target: “mat” for “cat”, for example. It is also possible for an erroneous production to be *mixed*, that is both semantically and phonemically related to the target word: “rat” for “cat.” Individuals with anomia also produce *perseverations*, which are words that are neither semantically or phonemically related to their intended target word: for example, producing “skis” instead of “zipper.”

Each of these categories shares the commonality that the word produced by the individual is a “real” word. There is another family of anomic errors, *neologisms*, in which the individual produces *non-word* productions. A neologistic production may be phonemically related to the target, but containing phonological errors: “[d̪aɪnoʊsɔɪ]” for “dinosaur.” These are often referred to as *phono-*

*logical* paraphasias. Alternatively, the individual may produce *abstruse neologisms*, in which the produced phonemes bear no discernable similarity to any “real” lexical item (“[æpməl]” for “comb”<sup>3</sup>).

A full discussion of the theoretical basis for this typology of paraphasias is beyond the scope of this paper. That said, it is worth noting that the standard model explaining these sorts of anomic errors is Dell’s two-step word production model (Dell et al., 1997; Dell, 1986). In Dell’s model, language production occurs in two primary phases. First, the speaker forms some semantic representation of what they wish to say. Next, that semantic representation leads to the activation of some number of appropriate lexical items. Finally, those selected lexical items are translated into spoken language.

Under this model, then, there are two primary ways that paraphasic productions can occur: the “wrong” lexical item may be selected (as in the case of a semantic error or a primary perseveration), and/or the translation process from lexical item to produced speech may go awry, resulting in formal (or other such phonemic) errors. When problems occur in both steps of the process, we see a mixed error (i.e., an error containing both phonological and semantic components).

The present work focuses exclusively on neologisms, both of the phonological variety as well as the abstruse variety. However, our fundamental approach can be extended to include other forms, as described in section 5.1.

Typical methods of diagnosing, staging, and otherwise characterizing anomia involve determining the number and kinds of paraphasias produced by an individual while undergoing some structured language elicitation process, for example a confrontation naming test (see (Kendall et al., 2013) and (Brookshire et al., 2014) for examples of such a study). As alluded to previously, producing these counts and classifications is a complex and laborious process. Furthermore, it is also often an inherently subjective

<sup>1</sup>As in the computer-administered (but manually-scored) assessments developed by Fergadiotis and colleagues (Fergadiotis et al., 2015; Hula et al., 2015).

<sup>2</sup>Note that individuals *without* any sort of language disorder do occasionally produce errors in their speech; this fact has led to a truly shocking amount of study by linguists. Frisch & Wright (2002) provide a reasonable overview of the background and phonology of the phenomenon.

<sup>3</sup>This example was taken from a corpus of responses to a confrontation naming test (Mirman et al., 2010), in which the subject is shown a picture and required to name its contents. As such, in the case of this specific error, we have *a priori* knowledge of what the target word “should” have been. Obviously, in a more naturalistic task or setting, we would not have this advantage.

process: are “carrot” and “banana” semantically related? What about “hose” and “rope”?

Reliability estimates of expert human performance at paraphasia classification in confrontation naming scenarios reflect the difficulty in this task. One recent study reported a kappa-equivalent score of 0.76 — a score that is certainly acceptable, but that leaves much room for disagreement on the status of specific erroneous productions (Minkina et al., 2015). Other reported scores fall in a similar range (Kristensson et al., 2015), including when the productions are from neurotypical individuals (Nicholas et al., 1989). Automating this aspect of the task would not only improve efficiency, but would also decrease scoring variability.

Having a reliable, automated method to analyze paraphasic errors would also expand the scope of what is currently possible in terms of assessment methodologies. Confrontation naming tests are often used (both in the clinic as well as in research settings) not because they are thought to optimally characterize the speech capabilities of their subjects, but rather out of concerns regarding feasibility of scoring. Naturalistic language samples may be more ecologically valid, but such data present a very complex and challenging scoring scenario (for example, see (Nicholas and Brookshire, 1993; Berndt et al., 2000; Rochon et al., 2000)), and so practitioners often eschew them in favor of simpler, more structured assessments.

Notably, the approach we outline in this paper is explicitly designed to work on samples of natural, connected speech— though it did grow out of work that the authors are currently doing on automated confrontation naming test scoring. It is our hope that, by enabling automated calculation of error frequencies and types on narrative speech, we might make using such material far easier in practice than it is today.

## 1.2 Related work

The general problem of applying NLP to language assessment has seen a phenomenal amount of work in recent years.

Computational analysis of aphasic speech: Using aphasiabank (MacWhinney et al., 2011) automatic sub-typing from narrative transcripts (Fraser et al., 2014b) statistical parsing to detect (Fraser et al., 2014a) Syntactic analysis of same: (Goodglass et al.,

1994) Segmentation of aphasic speech: (Fraser et al., 2015)

Guessing the next word: (Shannon, 1951)

## 2 Methods

Our overall approach was as follows. Beginning with transcribed narrative samples from people with aphasia, we used a corpus-driven method to identify possible neologisms in the subjects’ speech. Once we have identified candidate neologisms, we must Next, for each neologism, we used an n-gram language model to build a weighted lattice of plausible “target words” (see figure 2.4 for an unweighted example of such a lattice). For each sentence, we then compute a metric of phonological similarity between the erroneous utterance produced by the subject and the candidate target words. We then attempt to classify this production as either a phonological paraphasia or an abstruse neologism. We will describe each step of the process in detail below, beginning with our the data set.

### 2.1 Dataset

For the work described in this paper, we relied on the AphasiaBank project (MacWhinney et al., 2011), which has assembled a large database of transcribed interactions between examiners and people with aphasia, nearly all of whom have suffered a stroke. Notably, AphasiaBank also includes some number of transcribed sessions with neurotypical controls. Each interaction follows a common protocol and script, and is transcribed in great detail using a standardized set of annotation guidelines. The transcripts include word-level error codes, according to a spectacularly detailed taxonomy of errors and associated annotations. Erroneous productions are not simply flagged as erroneous. In the case of semantic, formal, and phonemic errors, the word-level annotations include a “best guess” on the part of the transcriber as to what the speaker’s intended production may have been. In the case of non-lexical productions (phonemic errors, neologisms, etc.), the annotations include an IPA transcription of the subject’s precise utterance. In some cases, the transcribers include information about gestures and other non-verbal communication that the subjects may have produced. The transcripts are stored in

the CLAN (Computerized Language Analysis) format (MacWhinney, 2000), and are therefore highly amenable to automated analysis.

Each transcribed session consists of a prescribed sequence of language elicitation activities, including a set of personal narratives (e.g., “Tell me a story about an important event that happened to you”, “Do you remember when you had your stroke? Please tell me about it.”), standardized picture description tasks, a story retelling task (involving the story of *Cinderella*), and a procedural discourse task (in which the subjects are asked to describe for the examiner the process of making a peanut-butter and jelly sandwich).

In addition to the narrative sessions, each subject’s data also includes the results of a battery of standardized assessments, including a confrontation naming test, the Aphasia Quotient sub-test of the Western Aphasia Battery, and so forth.

We obtained an up-to-date copy of the AphasiaBank database, and applied a series of minor normalizations as a first step in our analytical pipeline. First, we harmonized the names by which examiners referred to the various tasks, as this varied slightly from study site to study site (e.g., one site referred to the “Important event” task, while another referred to “Important\_Event”), and dealt with several other such minor orthographic irregularities across study sites.

Less trivially, we collapsed certain word-level error codes. As an example, the AphasiaBank protocol includes a special annotation used to indicate that a given neologism error represents an utterance that recurs frequently within a particular subject, as occasionally happens with individuals with aphasia (e.g., an individual might repeatedly utter the phonemes XXX). Since the present analysis was not concerned with this aspect of the transcripts, we collapsed instances of this and several other similarly specific error codes into their more general forms.

Certain non-standard productions that are nevertheless typical of a dialect (“gotta”, etc.) are occasionally labeled by the AphasiaBank transcribers as such, along with their “canonical” form (e.g. “got to” for “gotta”). In these cases, we replaced the “true” production with its canonical form. This was motivated by the fact that we would be using language models that were not explicitly trained on conversa-

tional speech, and we did not want incidental dialect usage of this sort to complicate matters.

Finally, we transformed the human-generated IPA representation of non-word productions into a relatively impoverished graphemic representation. Our purpose in doing this was to simulate a scenario where we did *not* have high-quality human-produced phonetic transcriptions. Obviously, one major obstacle to using narrative samples in any sort of clinical assessment is the need for transcription. Two of the most plausible solutions to this problem are automatic speech recognition (ASR)<sup>4</sup> and/or the use of *non-expert* transcribers to produce “quick-and-dirty” transcripts (perhaps via Amazon Mechanical Turk, or some other such crowd-sourced platform).

In the case of ASR-derived transcripts, we would presumably have some sort of phonemic representation of what was said, though this representation would almost certainly be incomplete and error-prone (modern ASR systems are notorious for struggling with neologisms). In the case of non-expert transcribers, we would of course not have the professional-quality IPA transcriptions found in the AphasiaBank database. We would instead be working with a “best effort” attempt at rendering what the transcriber heard the subject produce.

We therefore anticipate that our system will need to operate in a world with imperfect phonemic information. As such, we XXX

We chose two of the AphasiaBank personal narrative tasks— describing the individual’s stroke, and telling the story of an important event that happened sometime during the individual’s life— to work with, and excerpted the intervals containing those sections from each transcript. This resulted in XXX sentences from YYY individuals. We then identified sentences containing instances of our errors of interest: phonological paraphasia (codes XXX, YYY, etc.) or abstruse neologism Note that the distribution of errors within sentences was relatively Zipfian, in that the majority of error-containing sentences contained a single error, followed somewhat distantly by sentences containing two errors, with a relatively steep dropoff thereafter. For the present study, we restricted our analysis to sentences that contained ei-

<sup>4</sup>See Fraser et al. (2015) for a detailed discussion of the implications of using ASR on the speech of individuals with aphasia

ther one or two errors.

Our reasoning for this restriction was that we do not presently have a theoretically-informed model of what, if any, relationship there may be between multiple errors within a sentence, but that it seems likely that the errors occurring in a sentence containing (for instance) five paraphasic errors might be somehow related to one another. We anticipate exploring this phenomenon in the future (see section 5.1).

Our final data set consisted of XXX sentences from YYY individuals, containing errors from ZZZ distinct categories.

## 2.2 Identification of Neologisms

We next turned to the question of identifying neologisms in our sentences. Simply using a standard dictionary to determine lexicality could result in numerous “false positives,” driven largely by proper names of people, brands, etc. To avoid this, we used the Subtlex-US corpus (Brysbaert and New, 2009) to identify neologisms. Subtlex-US was built using subtitles from English-language television shows and movies, and Brysbaert and New have demonstrated that it correlates with a number of psycholinguistic behavior measures (very notably including naming latencies) better than several commonly used corpora such as the Brown corpus and CELEX2.

Note that, in the present analysis, this step in our pipeline was something of a contrived exercise, as (thanks to the detailed annotations present in the AphasiaBank transcripts) we already “knew” which tokens represented neologisms. However, in a “real-world” scenario, when we did *not* know *a priori* which tokens represented non-word productions, this step would be of particular importance, and we wished to simulate it with as much fidelity as possible. Furthermore, recall that while the present analysis only concerns itself with non-word productions, there are a number of paraphasia types in which a valid word is produced. Determining which productions represent a semantic or formal error is much more complex than simply performing a corpus lookup. We expect this stage of our pipeline to grow in complexity in the future.

Upon identifying a possible non-word production, recall that our next goal is to determine whether it represents a *phonemic* error (substituting “[dʌmɒʊsɔɪ]” for “dinosaur”) or an *abstruse neolo-*

*gism* (a completely novel sequence of phonemes that does not correspond to an actual word). To help accomplish this, we use a language model to identify plausible words that *could* fit in the slot occupied by the erroneous production, and produce a lattice of these candidate target words (i.e., words that the subject may have been intending to produce, given what we know about the context in which they were speaking).

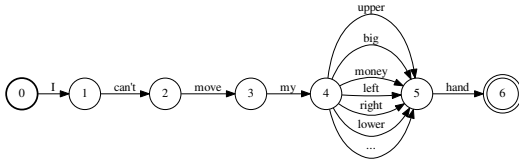
## 2.3 Language Model Construction

Our language models for this study were built using the New York Times section of Gigaword, 5th edition (LDC2011T07). We tokenized using the standard Penn Treebank tokenizer, left stopwords intact, and case-folded all sentences to upper-case. Cardinal numbers were collapsed into a category token, as were ordinal numbers and dates (each category was given its own token). We used the the OpenGrm-Ngram language modeling toolkit (Roark et al., 2012) to build the language models themselves, using an n-gram order of 4, with XXX smoothing.

We investigated two different language model approaches. In the first approach, we trained our models on the totality of the New York Times data. However, given that many of the AphasiaBank narrative tasks consist of fairly topic-constrained language (e.g., the Cinderella retelling, the personal narrative about the subject’s stroke, etc.), we hypothesized that we would get better results (i.e., higher-quality word predictions) if we could train our models on a more focused set of data.

To accomplish this, we used the Gensim topic modeling package (Řehůřek and Sojka, 2010) to train a Latent Dirichlet Allocation topic model (Blei et al., 2003) on the entire New York Times data set. For the present analysis, we instructed the model to use 20 topics. We next projected the text of each of the narrative samples into the topic space described by the model, and calculated the centroids for each of the narrative task. This gave us, for example, the estimated topic distribution most representative of the Cinderella retellings, the PB&J instructions, and so on.

Then, we calculated the Euclidean distance between each article in the New York Times corpus and each narrative task’s centroid. This allowed us to determine the “most similar” New York Times articles



**Figure 1:** An example candidate word lattice for the production “I can’t move my [var@u] hand.”

for each narrative task, which in turn enabled us to produce “task-specific” subsets of the larger corpus.

The end result was that we were able to produce a more topically homogeneous collection of New York Times articles to use in training per-narrative-task language models. We ended up identifying two of the narrative tasks—“tell me about your stroke” and “tell me about an important event”—as fitting particularly well into the topic space described by the New York Times data, and focused on those two tasks for the remainder of the analysis.

In this study, we evaluated the performance implications of using either the omnibus model or a task-specific model. In future work (see section 5.1), we anticipate experimentation with interpolating between the two.

## 2.4 Lattice Construction & Scoring

## 2.5 Phonological Similarity

Phonetisaurus (Novak et al., 2011)

## 2.6 Ranking & Scoring

## 3 Results

## 4 Discussion

## 5 Conclusion & Future Work

Limitations: - We are only using sentences with 1 error- excluding sentences with >1 error (N = 1,866) - more generally, we don’t have a clean idea of how/whether sentences with multiple errors are different from mono-error sentences - open question for future work: are paraphasic errors within a sentence related to one another in some way? - our general

finite-state approach can be generalized to sentences with additional errors, and we will explore such possibilities in future work!

## 5.1 Future Work

Future work: - better LM, better phonology, etc. etc. - interpolation between omnibus and task-specific LMs - using PoS/syntax to help word prediction - subject-level adaptation: using characteristics of the subject and their language to help identify erroneous productions - use of surprisal

## References

- R Berndt, S Wayland, E Rochon, E Saffran, and M Schwartz. 2000. *Quantitative production analysis: A training manual for the analysis of aphasic sentence production*. Psychology Press.
- David M Blei, Andrew Y Ng, and Michael I Jordan. 2003. Latent Dirichlet Allocation. *J. Mach. Learn. Res.*, 3:993–1022, March.
- C Elizabeth Brookshire, Tim Conway, Rebecca Hunting Pompon, Megan Oelke, and Diane L Kendall. 2014. Effects of Intensive Phonomotor Treatment on Reading in Eight Individuals With Aphasia and Phonological Alexia. *American Journal of Speech-Language Pathology*, 23(2):S300–S311, May.
- Marc Brysbaert and Boris New. 2009. Moving beyond Kucera and Francis: a critical evaluation of current word frequency norms and the introduction of a new and improved word frequency measure for American English. *Behavior research methods*, 41(4):977–990, November.
- Madeline Cruice, Linda Worrall, Louise Hickson, and Robert Murison. 2003. Finding a focus for quality of life with aphasia: Social and emotional health, and psychological well-being. *Aphasiology*, 17(4):333–353.
- G S Dell, M F Schwartz, N Martin, E M Saffran, and D A Gagnon. 1997. Lexical access in aphasic and nonaphasic speakers. *Psychological review*, 104(4):801–838, October.
- G S Dell. 1986. A spreading-activation theory of retrieval in sentence production. *Psychological review*, 93(3):283–321, July.
- Stefan T Engelter, Michal Gostynski, Susanna Papa, Maya Frei, Claudia Born, Vladeta Ajdacic-Gross, Felix Gutzwiller, and Phillipe A Lyrer. 2006. Epidemiology of aphasia attributable to first ischemic stroke: incidence, severity, fluency, etiology, and thrombolysis. *Stroke; a journal of cerebral circulation*, 37(6):1379–1384, June.

- Gerasimos Fergadiotis, Stacey Kellough, and William D Hula. 2015. Item Response Theory Modeling of the Philadelphia Naming Test. *Journal of Speech, Language, and Hearing Research*, 58(3):865–877, June.
- Kathleen C Fraser, Graeme Hirst, Jed A Meltzer, Jennifer E Mack, and Cynthia K Thompson. 2014a. Using statistical parsing to detect agrammatic aphasia. In *Proceedings of BioNLP 2014*, pages 134–142, Baltimore, Maryland, June. Association for Computational Linguistics.
- Kathleen C Fraser, Jed A Meltzer, Naida L Graham, Carol Leonard, Graeme Hirst, Sandra E Black, and Elizabeth Rochon. 2014b. Automated classification of primary progressive aphasia subtypes from narrative speech transcripts. *Cortex; a journal devoted to the study of the nervous system and behavior*, 55:43–60, June.
- Kathleen C Fraser, Naama Ben-David, Graeme Hirst, Naida Graham, and Elizabeth Rochon. 2015. Sentence segmentation of aphasic speech. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 862–871, Denver, Colorado, May. Association for Computational Linguistics.
- Stefan A Frisch and Richard Wright. 2002. The phonetics of phonological speech errors: An acoustic analysis of slips of the tongue. *Journal of Phonetics*, 30(2):139–162.
- Jorge Moncayo Gaete and Julien Bogousslavsky. 2008. Post-stroke depression. *Expert review of neurotherapeutics*, 8(1):75–92, January.
- Harold Goodglass, Julie Ann Christiansen, and Roberta E Gallagher. 1994. Syntactic constructions used by agrammatic speakers: Comparison with conduction aphasics and normals. *Neuropsychology*, 8(4):598–613.
- William D Hula, Stacey Kellough, and Gerasimos Fergadiotis. 2015. Development and Simulation Testing of a Computerized Adaptive Version of the Philadelphia Naming Test. *Journal of Speech, Language, and Hearing Research*, 58(3):878–890, June.
- Diane L Kendall, Rebecca Hunting Pompon, C Elizabeth Brookshire, Irene Minkina, and Lauren Bislick. 2013. An Analysis of Aphasic Naming Errors as an Indicator of Improved Linguistic Processing Following Phonomotor Treatment. *American Journal of Speech-Language Pathology*, 22(2):S240–S249, May.
- Joana Kristensson, Ingrid Behrns, and Charlotta Saldert. 2015. Effects on communication from intensive treatment with semantic feature analysis in aphasia. *Aphasiology*, 29(4):466–487.
- Brian MacWhinney, Davida Fromm, Margaret Forbes, and Audrey Holland. 2011. AphasiaBank: Methods for Studying Discourse. *Aphasiology*, 25(11):1286–1307.
- Brian MacWhinney. 2000. *The CHILDES Project: Tools for Analyzing Talk*. Lawrence Erlbaum Associates, Mahwah, NJ, 3rd edition.
- Jamie Mayer and Laura Murray. 2003. Functional measures of naming in aphasia: Word retrieval in confrontation naming versus connected speech. *Aphasiology*, 17(5):481–497, January.
- Irene Minkina, Megan Oelke, Lauren P Bislick, C Elizabeth Brookshire, Rebecca Hunting Pompon, Joann P Silkes, and Diane L Kendall. 2015. An investigation of aphasic naming error evolution following phonomotor treatment. *Aphasiology*, pages 1–19, August.
- Daniel Mirman, Ted J Strauss, Adelyn Brecher, Grant M Walker, Paula Sobel, Gary S Dell, and Myrna F Schwartz. 2010. A large, searchable, web-based database of aphasic performance on picture naming and other tests of cognitive function. *Cognitive neuropsychology*, 27(6):495–504, September.
- National Institute of Neurological Disorders and Stroke. 2014. NINDS aphasia information page. Technical report.
- Linda E Nicholas and Robert H Brookshire. 1993. A system for quantifying the informativeness and efficiency of the connected speech of adults with aphasia. *Journal of speech and hearing research*, 36(2):338–350, April.
- Linda E Nicholas, Robert H Brookshire, Donald L MacLennan, James G Schumacher, and Shirley A Porrazzo. 1989. Revised administration and scoring procedures for the Boston Naming test and norms for non-brain-damaged adults. *Aphasiology*, 3(6):569–580.
- Josef Novak, Dong Yang, Nobuaki Minematsu, and Kei-kichi Hirose. 2011. Initial and evaluations of an open source wfst-based phoneticizer. Technical report, The University of Tokyo.
- Radim Řehůřek and Petr Sojka. 2010. Software Framework for Topic Modelling with Large Corpora. In *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*, pages 45–50, Valletta, Malta, May. ELRA. <http://is.muni.cz/publication/884893/en>.
- Brian Roark, Richard Sproat, Cyril Allauzen, Michael Riley, Jeffrey Sorensen, and Terry Tai. 2012. The OpenGrm open-source finite-state grammar software libraries. In *Proceedings of the ACL 2012 System Demonstrations*, ACL ’12, pages 61–66, Stroudsburg, PA, USA. Association for Computational Linguistics.
- E Rochon, E M Saffran, R S Berndt, and M F Schwartz. 2000. Quantitative analysis of aphasic sentence production: further development and new data. *Brain and language*, 72(3):193–218, May.
- C Shannon. 1951. Prediction and entropy of printed English. *Bell System Technical Journal*, 30:51–64.

672	Mariska J van Dijk, Janneke M de Man-van Ginkel,	720
673	Thóra B Hafsteinsdóttir, and Marieke J Schuurmans.	721
674	2015. Identifying depression post-stroke in patients	722
675	with aphasia: A systematic review of the reliability, va-	723
676	lidity and feasibility of available instruments. <i>Clinical</i>	724
677	<i>rehabilitation</i> , page 0269215515599665, August.	725
678		726
679		727
680		728
681		729
682		730
683		731
684		732
685		733
686		734
687		735
688		736
689		737
690		738
691		739
692		740
693		741
694		742
695		743
696		744
697		745
698		746
699		747
700		748
701		749
702		750
703		751
704		752
705		753
706		754
707		755
708		756
709		757
710		758
711		759
712		760
713		761
714		762
715		763
716		764
717		765
718		766
719		767