



Improving the understanding of spoken referring expressions through syntactic-semantic and contextual-phonetic error-correction[☆]

Ingrid Zukerman^{*}, Andisheh Partovi

Faculty of Information Technology, Monash University, Clayton, Victoria 3800, Australia

Received 15 April 2016; received in revised form 28 February 2017; accepted 12 May 2017

Abstract

Despite recent advances in automatic speech recognition, one of the main stumbling blocks to the widespread adoption of Spoken Dialogue Systems is the lack of reliability of automatic speech recognizers. In this paper, we offer a two-tier error-correction process that harnesses syntactic, semantic and pragmatic information to improve the understanding of spoken referring expressions, specifically descriptions of objects in physical spaces. A syntactic-semantic tier offers generic corrections to perceived ASR errors on the basis of syntactic expectations of a semantic model, and passes the corrected texts to a language understanding system. The output of this system, which consists of pragmatic interpretations, is then refined by a contextual-phonetic tier, which prefers interpretations that are phonetically similar to the mis-heard words. Our results, obtained on a corpus of 341 referring expressions, show that syntactic-semantic error correction significantly improves interpretation performance, and contextual-phonetic refinements yield further improvements.

© 2017 Elsevier Ltd. All rights reserved.

Keywords: Spoken language understanding; Referring expressions; Error correction; Pragmatic interpretation; Physical spaces; Syntactic-semantic model; Contextual-phonetic model

1. Introduction

In recent times, there have been significant improvements in Automatic Speech Recognition (ASR) (Pellegrini and Trancoso, 2010; Chorowski et al., 2015). Nonetheless, ASR errors and inconsistent performance across domains still impede the widespread adoption of Spoken Dialogue Systems (SDSs). For example, a research prototype of a spoken slot-filling dialogue system reported a Word Error Rate (WER) of 13.8% when using “a generic dictation ASR system” (Mesnil et al., 2015), and Google reported an 8% WER for its ASR API,¹ but this API had a WER of 54.6% when applied to the Let’s Go corpus (Lange and Suendermann-Oeft, 2014). The accuracy of the commercial ASR employed in this research (Microsoft Speech SDK 6.1) falls between these numbers, with a WER of 30% for

[☆] This paper has been recommended for acceptance by Prof. R. K. Moore.

^{*} Corresponding author.

E-mail address: Ingrid.Zukerman@monash.edu (I. Zukerman), andi.partovi@monash.edu (A. Partovi).

¹ venturebeat.com/2015/05/28/google-says-its-speech-recognition-technology-now-has-only-an-8-word-error-rate.

9 descriptions of household objects. On one hand, these descriptions are more open-ended than the utterances
10 employed in slot-filling applications, but on the other hand, they do not contain proper nouns, such as names of cities
11 and foreign entities, which are rather error prone (Bulyko et al., 2005).

12 ASR errors not only produce mis-heard (wrongly recognized) entities or actions, but may also yield ungrammati-
13 cal utterances that cannot be processed by subsequent interpretation modules of a Spoken Language Understanding
14 (SLU) system (e.g., “the plate inside the microwave” being mis-heard as “of *plating sight* the microwave”), or yield
15 incorrect results when processed by these modules (e.g., hesitations, often accompanied by fillers such as “hmm” or
16 “ah”, being mis-heard as “and” or “on”) – all of which happened in our trials. Thus, further improvements are
17 required in order to enable the widespread adoption of SDSs.

18 Two approaches for achieving such improvements are (1) enhancing ASR performance, e.g., through the use of
19 deep neural networks (Hinton et al., 2012; Bahdanau et al., 2016) or by reducing the noise in the input signal (Maas
20 et al., 2012); and (2) performing post-ASR error correction, e.g., (Alumäe and Kurimo, 2010; Béchet et al., 2014;
21 Fusayasu et al., 2015). Clearly, perfect ASR would obviate the need for the latter approach, but we are not there yet,
22 even with recent advances in ASRs based on deep neural networks (Tam et al., 2014). Further, Ringger and Allen
23 (1996) argue that “even if the SR engine’s language model can be updated with new domain-specific data, the post-
24 processor trained on the *same* new data can provide additional improvements in accuracy”. Hence, at present, post-
25 ASR error correction is an active area of research (Section 9), which offers a practical way of de-coupling innovation
26 cycles within target applications from ASR innovation cycles (Feld et al., 2012), and increases system
27 portability (Ringger and Allen, 1996).

28 In this paper, we offer a mechanism for improving the understanding of spoken referring expressions, specifically
29 descriptions of objects in physical spaces, by harnessing syntactic, semantic and pragmatic information after obtain-
30 ing output from the ASR. Our mechanism consists of two main stages: syntactic-semantic error correction (Kim
31 et al., 2013) and contextual-phonetic error-correction (Zukerman et al., 2015b).

32 The *syntactic-semantic* tier receives as input textual alternatives returned by the ASR. It first invokes a classi-
33 fier that postulates wrong words in the ASR output (Zavareh et al., 2013) (Section 4) and a shallow semantic
34 parser that breaks up these texts into semantically labelled segments (Kim et al., 2013) (Section 5). It then acti-
35 vates a syntactic-semantic error-correction model that proposes modifications for selected words in each text on
36 the basis of the information provided by the word-error classifier and syntactic expectations of the semantic seg-
37 ments. The following modifications are considered during this stage: removal of noise (which is often due to filled
38 pauses), insertion of missing prepositions, and replacement of mis-heard words. We consider two main types of
39 replacements: (1) words or phrases expected to be closed class are replaced with phonetically similar closed-class
40 words or phrases, e.g., “for there a wave” in prepositional phrase position is replaced with “further away”; and
41 (2) words or phrases expected to be open class incur *generic* replacements, e.g., given a text such as “the *played*
42 inside the microwave” (mis-heard by our ASR from the description “the plate inside the microwave” in the context
43 of Fig. 1b²), “played” is replaced with the generic noun “thing”, yielding “the thing inside the microwave”. The
44 rationale for this replacement is that “played” would lead a parser astray, while the proposed replacement allows
45 an SLU system to proceed.

46 The resultant, possibly modified, texts are then given as input to the *Scusi?* SLU system (Zukerman et al., 2015a),
47 which is a component of an SDS for human-robot interactions. *Scusi?* generates a ranked list of pragmatic interpreta-
48 tions for these texts, where an interpretation comprises candidate things in a physical space and the spatial relations
49 between them, e.g., `plateD-location_in-microwave1`, and the ranking of an interpretation corresponds to the
50 extent to which it matches the description (Section 2). In our example, the interpretation comprising Plate D is
51 ranked first, as it is the only object inside the microwave. When the ASR made the same error for the description
52 “the plate on the table” (which ambiguously refers to Plate E in Fig. 1b), yielding “the *played* on the table”, the syn-
53 tactic-semantic tier generated “the thing on the table”. At this stage, this change offers no benefit in the context of
54 Fig. 1b, as all the items are on the table. However, in the context of Fig. 1a, the three objects on the table are ranked
55 equal first, ahead of the other objects in the scene.

56 The *contextual-phonetic* tier receives as input the top-*N* pragmatic interpretations produced by *Scusi?*, and re-
57 ranks them according to the phonetic similarity between the mis-heard head nouns in the texts that led to each

² The ASR mis-heard “plate” as “played” or “play” in 404 out of the 1697 texts produced for referring expressions describing plates: items D and E in Fig. 1b, and items G, H and I in Fig. 1c.

متن کامل مقاله

دریافت فوری ←

ISIArticles

مرجع مقالات تخصصی ایران

- ✓ امکان دانلود نسخه تمام متن مقالات انگلیسی
- ✓ امکان دانلود نسخه ترجمه شده مقالات
- ✓ پذیرش سفارش ترجمه تخصصی
- ✓ امکان جستجو در آرشیو جامعی از صدها موضوع و هزاران مقاله
- ✓ امکان دانلود رایگان ۲ صفحه اول هر مقاله
- ✓ امکان پرداخت اینترنتی با کلیه کارت های عضو شتاب
- ✓ دانلود فوری مقاله پس از پرداخت آنلاین
- ✓ پشتیبانی کامل خرید با بهره مندی از سیستم هوشمند رهگیری سفارشات