# Supplemental Material for Mixing Modalities of 3D Sketching and Speech for Interactive Model Retrieval in Virtual Reality

ANONYMOUS

## 1 INTRODUCTION

In this supplemental material, we describe the motivation of our choice of using a Wizard of Oz (WoZ) approach for the speech interaction (Section 2). Then we provide the sequence of user sessions for our second experiment (Section 4), show the dictionary and chair collection (shape and colors) we generated by colour permutations (Section 3).

## 2 WHY DO WE CHOOSE A WIZARD OF OZ APPROACH?

In this section we introduce the state of the art of Natural Language Process algorithms, we define the entire speech interaction pipeline and discuss the advantages and disadvantages of using automatic processes or the human component.

### 2.1 State of The Art of Natural Language Process

In the last decade, deep learning algorithms have been designed to process text data for classification tasks such as sentiment analysis or text comprehension. Recurrent Neural Networks (RNN) have been effectively utilised for text analysis [5], however they are prone to suffer from vanishing gradients [4], implying that they can forget the context and thus neglect parts of the sentence that are critical for correctly parsing it. More recent models such as Long Short-Term Memory (LSTM) [2] and its variants overcome this memory problem, making it easier to remember data. In addition, recent techniques of mass parallel processing developed by Google Brain and Google Research have helped to improve attention mechanisms, implemented in Transformer model architecture [6]. These NLP models had a considerable improvement in 2018, comparable to the impact of ImageNet [3] in 2012 for Computer Vision. These results make us consider to build a software stack for an automatic pipeline. This pipeline is described in the next section, considering the advantages and disadvantages of each step. Based on this analysis we conclude the automatic pipeline can introduce an error that percolates and/or accumulates in the subsequent steps, we chose to have an experimenter-in-the-loop approach, and so a semi-automatic. We motivate our choice in the next section.

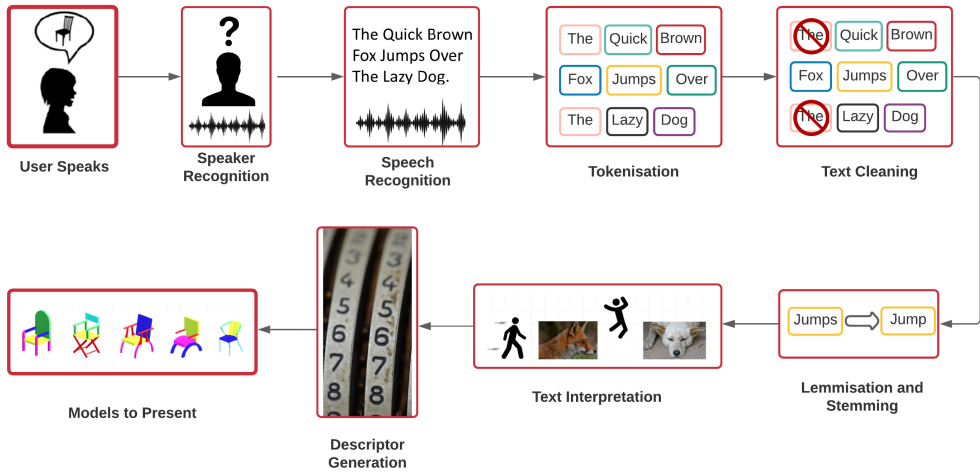**Unpublished working draft. Not for distribution.**

Fig. 1. Our speech pipeline includes: speaker identification, speech recognition, tokenisation, lemmatisation, stemming, text interpretation, descriptor generation, selection of the proposed results.

## 2.2 Speech Interaction Experiment Requirements

In the experiment, where three types of interaction are compared, while the sketch method was based on Giunchi *et al.* [1] work, the speech interaction was to be designed to be easily included in the system. We analysed the phase of speech interaction, and we split the entire process into stages. All these stages are part of a pipeline. This pipeline manipulates the speech query, as shown in Figure 1.

In our experiment, speech interaction involves a sequence of steps that starts with the user verbally describing the model, and finishes with the presentation of the results to the user. Between these two events, we can enumerate stages identified by the specific required task. This pipeline needs to last a maximum of 5 seconds with the following steps:

*Speaker identification.* Speaker identification or recognition is the process that allows identifying who is the active speaker. During the experiment, the user can interact with the experimenter by asking clarifications. Thus, it is essential to understand who is talking to avoid another input channel that injects words in the system.

*Speech recognition.* Speech recognition is the process that converts speech to plain text. This process is challenging, but efficient dictation software can be very accurate when trained with the speaker's voice. We tested three different speech recognition services without training phase, achieving poor results with audio recorded with Oculus equipment. In addition, another potential problem can be the words spoken by the user that are not relevant to the search, such as clarifications asked the experimenter during the test.

*Tokenisation.* Tokenisation or lexical analysis is the process that converts a sequential text in tokens. Each token is a string with its own length and meaning. Tokenising a text is a very simple operation.

99 *Text cleaning.* Text cleaning is the process that removes from the list of tokens all the items that
100 do not contain relevant information for the training of a machine learning model. NLP libraries
101 implement dictionaries that are used to exclude such words from the tokens list.
102
103 *Lemmatisation and Stemming.* Lemmatisation and Stemming are two processes that prepare
104 the text to be used for training a model, grouping words with the same root—while Lemmatisa-
105 tion focuses on the morphological analysis of the word, stemming cut prefixes or suffixes of the
106 considered word.
107
108 *Text interpretation.* Text interpretation is the process that maps the corpus with the meanings that
109 are being tried to convey. In our case, each group of words represents information that describes
110 the object or a part of it. The output of this process is a feature vector that describes the speech
111 query.
112
113 *Selection of the models.* The descriptor generated by the previous step is used to measure the
114 Euclidean distance with the feature vectors of the models in the dataset, to find the closest matches.
115
116 ## 2.3 Motivation for the Wizard Of Oz
117 In this section, we enumerate the problems that each stage of the pipeline can introduce in the
118 system. Because all the steps of the pipeline can be automatic, we analyse them, one by one and we
119 motivate our choice of using a Wizard Of Oz instead of a fully automatic pipeline.
120 Speaker identification is requested because during the experiment can occur that the participant
121 asks the experimenter some clarifications about the experiment itself, and the nature of the questions
122 can be very different. Normally the experimenter answers and his voice is a different input that
123 should not be used for the experiment. One possible alternative is to instruct the user to avoid any
124 kind of interaction with the experimenter. This stage could be considered optional. An automatic
125 step, in this case, needs to have a 100% accuracy since injecting a different and unrelated input can
126 compromise the query. If a human handles, this phase is naive to consider the voice coming only
127 from the participant.
128 Speech recognition is an essential step in the pipeline that converts all the words pronounced by
129 the speaker to plain text. During this process, we noticed that the participants sometimes speak
130 without providing relevant information, instead they are react to some events happening in the
131 virtual environment, e.g., commenting on the results or also make some questions to themselves.
132 These chunks of information need to be disregarded as they can inject incorrect information in
133 to the pipeline. Accuracy of speech to text is the most important metric in this case. If the stage
134 is handled in an automatic way, there are different software that can be used. We tested three
135 different speech-to-text services, not fine-tuned with the speaker's voice, that represents a plausible
136 condition for our experiment. The accuracy for all the three services was not satisfactory, and we
137 present the detailed results in Section 2.4. To solve the problem of excluding uninformative chunk,
138 one possibility is to give the user a command that enables the speech to be interpreted, but this is
139 not an optimal solution as one of our pre-condition is to avoid the use of the hands during speech
140 interaction. An experimenter is able to convert speech to text very quickly, but it isn't easy in
141 our experiment configuration to write text in real-time. An efficient alternative is to provide the
142 experimenter with a GUI and instructing the participant to follow a limited dictionary (always
143 visible in VR). This dictionary contains all the meanings valid for the dataset that we designed. In
144 addition, we describe a speech query in order to generate one search each 10 seconds, this is because
145 a continuous speech is challenging to deal with for both for an automatic and a semi-automatic
146 process.
147

Tokenisation stage is a very simple task to achieve from a corpus both if a computer or a human manages the stage.

Text cleaning stage is another phase of the pipeline that is simple to achieve. If handled by computer NLP libraries contains dictionaries that list all the words that do not give useful information. Also, for a human that is listening to a speaker, the task of dropping words that do not add semantic value to the sentence is easy.

Lemmatisation and stemming stage can be managed by NLP libraries, and also easily by a human with the limitation that we put in our experiment configuration.

Text interpretation and descriptor generation step can be managed automatically by one of the state-of-the-art models such as Transformer model. In this case, anyway, it is necessary to train the model, labelling all the models with many descriptions. A possible way to achieve such meta-information dataset is to use Amazon Mechanical Turk or to hire additional participants to describe a large part of the chairs and complete the colour variational dataset, replacing the colours in the description. However, this activity requires a big effort. On the other hand, the experimenter can be provided with an interface where he can clock buttons that increase the value of specific entries connected to the meanings that the user expressed. In this case, at the end of the query, a feature vector is automatically created and, moreover, can bee synchronised with the current selection in the virtual scene.

Selection of the models is a stage that anticipates the models' presentation and can be handled automatically.

Each stage can introduce errors, independently if a human or a machine manages it. These errors are also difficult to calculate for each stage, but they accumulate over the stages. We notice that speech recognition gave us very low accuracy in some cases when managed by speech-to-text software. The reasons can be found in the different component of the audio profile of each participant (tempo, rhythm, pitch, context), as well as fluency and accents which all can affect the accuracy of the transcript. In addition, we do not use software trained with the voice of each participant, and it is reasonable to experience a worse accuracy if compared with trained dictation software. However, fine-tuning would add an additional training step within the study. Moreover, the microphone we used could inject noise in the system.

On the other side, a human can deal with speech to text conversion easily, and with some limitation related to the domain of the meanings extracted from the dataset, we achieve a reliable semi-automatic system.

## 2.4 Speech to Text services

In this section, we present the results of some speech to text services fed with the audio files coming from some users that tested our apparatus. We tested 10 audio files in the following services:

(1) Watson IBM (https://speech-to-text-demo.ng.bluemix.net/ from heree)
(2) SONIX (https://sonix.ai/)
(3) google speech to text (https://cloud.google.com/speech-to-text)

We determined the accuracy considering groups of keywords in the speech queries. Each group conveys specific information about the style and/or colour of the chair, or the style and/or colour of a part of the chair. For example, if the speaker told "red straight arm" only that sequence of words (or eventually "straight red arm") will trigger a point in the accuracy score. We achieved 37% of accuracy for the Watson IBM, 60% for the SONIX service, and even only 16% with google speech.

## 3 DICTIONARY AND DATASET

We define a dictionary that contains the following characteristics describing a chair. Each feature is associated with a value that represents how much that feature impact on the description of the chair. Our dictionary contains the following concepts: Height-Length, Size, Thick, Decorated, Curvy, Modern, Antique, Slatted, Swiveling, Flexible, Stable, Reclinable, Padded, Slanted, Canvased, Missing. Each concept can be associated to the chair globally, or to each component of the segmented chairs: back, seat, arms, legs.

### 3.1 Chair Shapes and Colours

This section shows the shapes (45) of the chairs we include in the dataset, in Figure 2, and for one chair all the colour permutation (360) in Figure 3 for a total of 16200. This dataset is 5 time larger than the version used in Giunchi *et al.* [1] work.



Fig. 2. The 45 shapes selected for the chairs in the dataset.

Fig. 3. The 360 permutations without repetitions of colour in one type of chair.

## 4 SECOND EXPERIMENT SCHEDULE

This section shows the schedule of the experiment, in Figure 4.



Fig. 4. This table shows the schedule of the experiment with ten users (rows). Each column represents a searching session, where the first number is the searched shape model, while the second number and the colour reveals the interaction type. Each user searched 27 chairs, split into three sessions with 9 chairs with the same type of interaction. Both the models and the interaction are randomised.

## REFERENCES

[1] Daniele Giunchi, Stuart James, and Anthony Steed. 2018. 3D sketching for interactive model retrieval in virtual reality. In *Proceedings of the Joint Symposium on Computational Aesthetics and Sketch-Based Interfaces and Modeling and Non-Photorealistic Animation and Rendering.* ACM, 1.

[2] Klaus Greff, Rupesh Kumar Srivastava, Jan Koutník, Bas R. Steunebrink, and Jürgen Schmidhuber. 2017. LSTM: A Search Space Odyssey. *IEEE Trans. Neural Networks Learn. Syst.* 28, 10 (2017), 2222–2232. https://doi.org/10.1109/TNNLS.2016.2582924

[3] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. 2012. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*. 1097–1105.

[4] Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. 2013. On the difficulty of training recurrent neural networks. In *Proceedings of the 30th ICML (Proceedings of Machine Learning Research, Vol. 28)*, Sanjoy Dasgupta and David McAllester (Eds.). PMLR, Atlanta, Georgia, USA, 1310–1318.

[5] Richard Socher, Cliff C Lin, Chris Manning, and Andrew Y Ng. 2011. Parsing natural scenes and natural language with recursive neural networks. In *Proceedings of the 28th ICML*. 129–136.

[6] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA*, Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (Eds.). 5998–6008. http://papers.nips.cc/paper/7181-attention-is-all-you-need