

# AMBIPUN: Generating Puns with Ambiguous Context

Anonymous ACL submission

## Abstract

Computational humor has garnered interest in the natural language processing community due to its wide applications to real-world scenarios. One way to express humor is via the use of puns. In this paper, we propose a simple yet effective way to generate pun sentences that does not require any training on existing puns. Our approach is inspired by humor theories that ambiguity comes from the context rather than the pun word itself. Given a pair of definitions of a pun word, our model first produces a list of related concepts through a reverse dictionary. We then utilize one-shot GPT3 to generate context words and then generate puns incorporating context words from both concepts. We also investigate how the position of a pun word appearing in the sentence will influence the generation results. Human evaluation shows that our method successfully generates pun 52% of the time, outperforming well crafted baselines and the state-of-the-art models by a large margin.

## 1 Introduction

Humor has the tendency to provoke laughter and provide amusement. By creating an engaging and conducive environment, it is one of the most important forms of human communication (Booth-Butterfield and Wanzer, 2018). Teaching computers to understand and generate humorous texts such as puns pave the way for various practical applications, such as improving creativity in machine-aided writing and making chat-bots more engaging.

In this paper, we tackle the problem of generating homographic puns (Miller et al., 2017): two or more meanings of a **polysemy** for an intended humorous or rhetorical effect. For example, the three punning jokes listed in Figure 1 exploits two contrasting meanings of the word *sentence*: 1) a string of words that are grammatical, and 2) the punishment by a court assigned to a guilty person. Compared with heterographic puns where the ambiguity comes from two near homophones spelled

Sense 1 Definition	a string of words that is complete in itself, typically containing a subject and predicate
Sense 2 Definition	(criminal law) a final judgment of guilty in a criminal case and the punishment that is imposed
Ours 1	The <u>sentence</u> is ungrammatical. The jury didn't hear it.
Ours 2	I'm sorry I said the <u>sentence</u> was too long but punishments are endless.
Human	The Judge has got a <u>stutter</u> . Looks like I am not getting a <u>sentence</u> .

Figure 1: An illustration of homographic puns. The target pun word “*sentence*” and the two sense definitions are given to the model as input, and the desired outputs are many punning sentences. To make the target word interpretable in both senses, we propose to include context words related to both sense definitions. We highlight the context words of each sense in blue and pink.

in a different way, the challenge of processing homographic puns is even bigger: we must differentiate contrasting senses of words that sound and are spelled in the same way.

Due to the lack of sizable training data, existing approaches to generate puns are all heavy-weighted in order to *not* rely on pun sentences for training. For example, (Yu et al., 2018) train a constrained neural language model Mou et al. (2015) from a general text corpus, and then use a joint decoding algorithm to guarantee that both definitions of the target pun word will make sense in the generated sequence. He et al. (2019) propose a local-global surprisal principle, and Luo et al. (2019) leverage the Generative Adversarial Nets (Goodfellow et al., 2014) to encourage ambiguity of the outputs via reinforcement learning. We, on the other hand, propose a simple yet effective way to tackle this problem: *encouraging ambiguity by incorporating context words related to each sense*.

Inspired by humor theories (Lippman and Dunn, 2000), we hypothesize that it is the contextual connections rather than the pun word itself that are

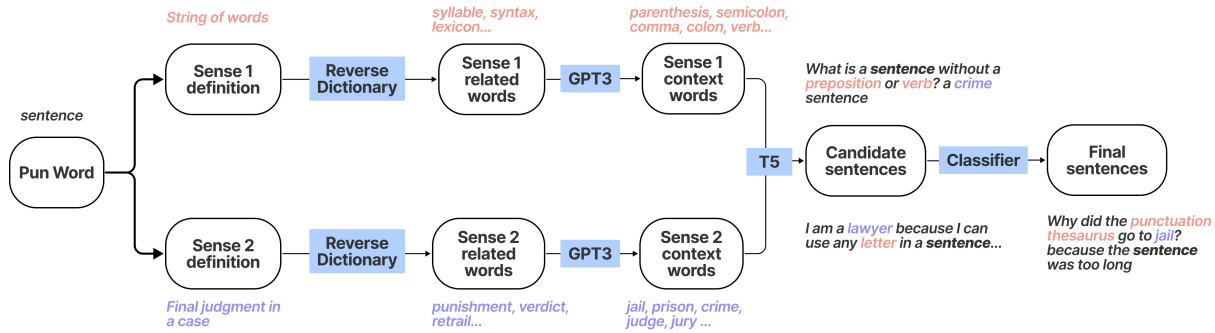


Figure 2: Overview of the approach. Given a pun word and its senses, we convert them to their sense definitions. then we use reverse dictionary to obtain the related words. Using few-shot GPT3, we generate context word for each related word. Using a combination of context words along with the pun word, we generate several candidates sentences using T5. Finally, we use a classifier to choose the most humorous sentences. We also give an example for pun word ‘sentence’ for each part of the approach.

065 crucial for the success of pun generation. For instance, in Figure 1 we observe that context related to both senses (e.g., *ungrammatical* and *jury*) appear in a punning sentence. Such observation is important as the error analysis of the state-of-the-art model (Luo et al., 2019) shows that 46% of the generated sentences fail to be puns due to single word sense, and another 27% fail due to being too general, both of which can be resolved by introducing more context.

066 Specifically, given the two sense definitions of a target pun word, we first use a reverse dictionary<sup>1</sup> to generate **related words** that are monosemous for both senses. This first step helps us circumvent the obstacle of processing pun words with the same written form. However, related words alone are not enough to generate coherent pun sentences because they are clustered and tend to be synonyms. We hence propose to use **context words** (described in Section 2.3) to link the contrasting senses and make our target pun word reasonable when interpreted in both definitions. We explore three different settings: retrieval-based (TF-IDF), similarity-based (Word2Vec), and generative-based (Few-shot GPT3). Finally, we finetune the T5 model (Raffel et al., 2020) on general non-humorous texts to generate coherent sentences given the pun word and contexts words as input.

067 Interestingly, our experimental results show that to retrieve-and-extract context words outperforms the giant few-shot GPT3 model in terms of generating funny pun sentences, although the latter has shown to be much more powerful in many sophisticated tasks (Brown et al., 2020). Our simple pipeline remarkably outperforms all the

068 more heavy-weighted approaches including the constrained language models with special decoding (Yu et al., 2018) and the state-of-the-art Pun-GAN model (Luo et al., 2019).<sup>2</sup>

## 2 Methodology

### 2.1 Overview and Motivation

069 We first give an overview of different steps in our approach. Our input is the target pun word ( $\mathcal{P}$ ) and its two senses ( $S_1, S_2$ ), and the expected output is a *list* of humorous punning sentences where  $\mathcal{P}$  can be interpreted in both senses. We implement the ambiguity principle proposed in (Kao et al., 2016): a pun sentence should contain one or more context words corresponding to each of the two senses of the pun word.

070 The overview of our approach is visualized in Figure 2 and formally written in Algorithm 1. Given two different sense descriptions (*SenseDef*), we first use a reverse dictionary to generate a list of words that semantically match the query descriptions. We call them **related words** ( $RW$ ) and describe full details in Section 2.2. However, those related words are synonyms of each other and only cover a few focused topics, thus failing to compose humorous punning sentences. For example, for the sentence: “The Judge has got a stutter. Looks like I am not getting a sentence.”, The word representing the first sense (i.e. a final judgment of guilty in a criminal case and the punishment that is imposed) is represented by *Judge*. *Judge* could not be generated using the sense definition, but is used frequently along with the sense definition.

<sup>1</sup><https://reversdictionary.org/>

<sup>2</sup>Our code and data will be released upon acceptance.

---

**Algorithm 1** Pun Generation

---

```
1: function GENPUN( $P, S_1, S_2$ )
2: Input: Tuple of pun word and its sense -  $P, S_1, S_2$ 
3: Output: List of final sentences -  $Sent_{final}$ 
4: for  $P, S_1, S_2$  in  $P, S_1, S_2$  do
5:    $SenseDef = \text{get\_sense\_definitions}(S_1, S_2)$ 
6:    $RW = \text{reverse\_dictionary}(SenseDef) \triangleright \text{get related}$ 
   words from sense definitions
7:    $RW_{refined} = \text{refine}(RW)$ 
8:    $CW = \text{get\_context\_words}(RW_{refined}) \triangleright \text{get context}$ 
   words from related words
9:    $CW_{refined} = \text{refine}(CW)$ 
10:   $Sent_{candidates} = \text{generate\_sentences}(CW_{refined})$ 
11:   $Sent_{final} = \text{classify\_sentences}(Sent_{candidates})$ 
   return  $Sent_{final}$ 
```

---

133 Considering such nuances, in Section 2.3 we  
134 propose three different methods to obtain the **con-**  
135 **text words** ( $CW_1, CW_2$ ) of all the related words.  
136 They are TF-IDF (retrieve-and-extract), similarity-  
137 based (word2vec), and generative model (few-shot  
138 GPT3). Finally, in Section 2.4 and Section 2.5,  
139 we introduce a keyword-to-text generator to gener-  
140 ate candidate sentences ( $Sent_{candidates}$ ), and a  
141 humor classifier to rule out some of the non-pun  
142 sentences. Final sentences ( $Sent_{final}$ ) are then  
143 randomly sampled for evaluation.<sup>3</sup> All our training  
144 data is general, non-humorous corpus except for  
145 the humor classifier.

## 146 2.2 Related words

147 We aim at differentiating the two senses of a poly-  
148 semy by taking the related words, so that each  
149 sense will be represented by a set of monosemous  
150 words. To this end, we leverage the Reverse Dic-  
151 tionary (Qi et al., 2020; Zhang et al., 2020) which  
152 takes as input a description and generates multiple  
153 related words whose semantic meaning match the  
154 query description. For each sense definition, we  
155 generate five words.

## 156 2.3 Context words

157 For context words, we compare three different ap-  
158 proaches. Refinement of the context words is men-  
159 tioned in Section 3.2.

160 **Method 1: TF-IDF** For each related word, we re-  
161 trieve sentences from the One Billion Word dataset  
162 that contains that word and then extract a few key-  
163 words. Next, we implement TF-IDF (Ramos, 2003)  
164 to rank them. For a given word and corpora, the

---

<sup>3</sup>Note that all previous works produce only the *best* sen-  
tence during decoding time, while we aim at generating *tens*  
or *hundreds* of sentences for a target pun word, so that our  
task is actually more challenging.

---

**Algorithm 2** TF-IDF for context words

---

```
function GETCW( $KW_1, KW_2$ )
2: Input: List of related words -  $KW_1, KW_2$ 
   Output: List of context words -  $CW_1, CW_2$ 
4: Initialize  $CW_1, CW_2$  to empty
   for  $KW$  in  $\text{zip}(KW_1, KW_2)$  do
6:   for  $w$  in  $KW$  do
      $S = \text{generate\_sentences}(w) \triangleright \text{Retrieve sentences}$ 
     for each word
8:      $kw = \text{generate\_keywords}(S)$ 
      $F_S = \text{get\_freq\_sentences}(kw)$ 
10:     $F_C = \text{get\_freq\_corpora}(kw)$ 
      $FK = \text{get\_tfidf}(F_S, F_C)$ 
12:     $CW_1 = \text{append}(FK) \triangleright \text{Add top TF-IDF words}$ 
     to context words list
   return  $CW_1, CW_2$ 
```

---

TF-IDF value is given in Equation 1, where  $\mathcal{F}_S$   
corresponds to the frequency of that word in the  
retrieved sentences and  $\mathcal{F}_C$  corresponds to the fre-  
quency of the word in the entire corpora. Based on  
this value, we choose the top 10 context words that  
are mostly likely to be used along with the related  
words and therefore the pun word. Detailed steps  
for the process are listed in Algorithm 2.

$$tf(W, C) = \frac{\mathcal{F}_S}{\mathcal{F}_C + 1} \quad (1)$$

174 **Method 2: Word2Vec** Inspired by the idea that  
175 “a word is characterized by the company it keeps”  
176 (Firth, 1957), we propose to get context words  
177 from word2vec (Mikolov et al., 2013), which pro-  
178 vides distributed word representations. Following  
179 a previous work (Ghazvininejad et al., 2016), we  
180 train a continuous-bag-of-words model with win-  
181 dows size 40 and word vector dimension 200, and  
182 then calculate the cosine similarity between words.  
183 Ghazvininejad et al. (2016) have also shown that  
184 the training corpus for word2vec plays a crucial  
185 role on the quality of generated context words.  
186 Hence, we try to train word2vec models on three  
187 different corpus: the largest available humorous  
188 Dataset, rJokes (Weller and Seppi, 2020), the En-  
189 glish Gigaword (Graff et al., 2003) which is an  
190 archive of newswire text data, and the one-billion  
191 Wikipedia corpus<sup>4</sup>. We find that the topics covered  
192 by rJokes is far from what it needs to train a good  
193 word2vec model, and that the word2vec model  
194 trained on Gigaword strongly favors newsy words  
195 than the others. Hence, we train on Wikipedia.

196 **Method 3: GPT3** For the generative version, we  
197 use the powerful language model, one-shot GPT3

---

<sup>4</sup><http://mattmahoney.net/dc/enwik9.zip>

(Brown et al., 2020) to generate context words. We choose not to train another model because the output of one-shot GPT3 is already satisfactory. An example can be seen in Table 1, where we compare the output of context words for the pun word ‘sentence’.

## 2.4 Candidate Sentence generation

After receiving context words for each sense, we generate humorous puns. For this step, we finetune a keyword-to-sentence model using T5 (Raffel et al., 2020), as it is capable of handling text-to-text tasks. To train this we need to create a dataset that replicate the expected behaviour i.e. given a prompt (in a specific format), generate a well formed and humorous sentence. The prompt will contain information about the pun word ( $P$ ), and 2 words from each of the two senses ( $S_{1a}$ ,  $S_{1b}$ ,  $S_{2a}$ ,  $S_{2b}$ ). We expect our generated output to contain the pun word and context words related to each sense. The following prompt is given to the trained model:

**generate sentence:**  $P, S_{1a}, S_{1b}, S_{2a}, S_{2b}$ .

For example for the word ‘sentence’, a possible prompt can be *generate sentence: sentence, judge, trail, noun, comma*. However, we also investigate whether the position of the pun word will affect the quality of generated sentences. We insert the pun word in the start (first place), middle (third place), and end (fifth place) of the prompt and generate candidate sentences using these prompt configurations. We discuss our findings in Section 4.3.

## 2.5 Humor Classification

Finally, we introduce a classification model to assist us in selecting (i.e., ranking) punning sentences. Since we do not have sizable training data for puns, we propose to train our classification model on humorous dataset in a distantly supervised fashion. Specifically, we train BERT-large (Devlin et al., 2019) on the ColBERT dataset that contains 200,000 jokes and non-jokes used for humor detection. We use the probability produced by the classification model to rank our candidate sentences.

Our error analysis shows that our distantly supervised classification model can successfully rule out the bad generations, i.e., non-puns, as puns are humorous by nature. However, the model is not

great at choosing the best samples.<sup>5</sup> Therefore, we use this classifier only to remove the bottom third candidates. We leave this for open future work to accurately pick out high-quality punning sentences instead of funny sentences.

## 3 Experiments

### 3.1 Datasets

**Training dataset:** For the context word generation steps, we use the One Billion word dataset (Chelba et al., 2013) to retrieve sentences for a given word. To calculate TF-IDF, we use this dataset to calculate the frequency of words. This dataset contains roughly 0.8B words and is obtained from WMT 2011 News crawl data.

For training the candidate generation module, we use ColBERT dataset (Khattab and Zaharia, 2020). It contains 100k positives and 100k negative samples collected from various sources like Reddit, news headlines, etc. For each sentence, we extract the keywords using RAKE (Rose et al., 2010). We also use the same data to finetune BERT-large to develop our humor classifier.

**Evaluation dataset:** On lines of other recent pun generation works, we use the SemEval 2017 Task 7 (Miller et al., 2017) for evaluation. The dataset contains 1,163 human written pun sentences with a total of 895 unique pun words. Each sentences has the target pun word, location of the pun word and the WordNet sense keys of the two senses.

### 3.2 Implementation Details

**Experimental Settings** For the word2vec model we train a continuous-bag-of-words model with window size 40 and word vector dimension 200. For the candidate generation module, we train the T5-base model on 10 epochs and select the best performing model based on validation loss. Max sequence length for target and source is set to 30. Batch size is set to 64.

**Data Refinement** The process to generate keywords (i.e., both related and context words) can entail many words that are not ideal. Continuing with these words would further propagate and enlarge the noise. Hence, to minimize this noise, we implement the following data refinement steps to ensure the keywords stick to our standards: we avoid using polysemous words as keywords during

<sup>5</sup>A few samples along with their assigned probabilities can be seen in Table ?? in the appendix



	Sense 1	Sense 2
<b>Definition</b>	a string of words satisfying the grammatical rules of a language	a final judgment of guilty in a criminal case and the punishment that is imposed
<b>Related words</b>	syllable, syntax, lexicon, thesaurus, grammatical	punishment, verdict, sentencing, retrial, penalty
<b>TF-IDF</b>	syllables, words, three, spelling, even, said, describe, typos	cruel, expected, end, court, scheduled, set, spectator, seeking
<b>Word2Vec</b>	syllable, pronounced, words, rhyme, verbs, meaning, hence, example	punished, crimes, offender, torture, moral, guilt, abuse, offender
<b>GPT3</b>	words, letters, punctuation, grammar, synonym, dictionary, meaning, comma	prison, judge, jury, trial, justice, lawyer, court, evidence

Table 1: Comparison of the three different context word generation mechanism. We take this example for the word ‘sentence’. The table lists sense definitions of the two senses. Then list the related words obtained from the sense definitions. For these related words, we obtain context words using three different mechanisms.

intermediate steps because their perceived sense is highly ambiguous. We also disregard any numbers and special characters produced by our systems.

### 3.3 Baselines

There are two existing works on homographic pun generation, the same task as ours. Besides, we also compare our model with the powerful few-shot learner, GPT3 (Brown et al., 2020).

**Neural Pun** Yu et al. (2018) propose the first neural approach to homographic puns based on a conditional neural language model. A constrained beam search algorithm is proposed to jointly decode the two distinct senses of the same word.

**Pun-GAN** The state-of-the-art-model introduced by Luo et al. (2019) that adopts the Generative Adversarial Net (GAN) (Goodfellow et al., 2014) to generate homographic puns. Specifically, a generator is responsible for generating a pun sentence, and a discriminator is trained to tell human-written puns from machine generated puns. Such setting encourages the ambiguity of the generated sentence via reinforcement learning (RL).

**Few-shot GPT3** We also generate puns with a few examples feeding into GPT3 davinci-instruct-beta, the most capable model in the GPT3 family to follow the instructions and generate creative language.<sup>6</sup> We provide the target pun word and its two senses in our prompt along with the instruction.

**Ablations of our own models** We also compare three methods proposed by us to obtain the context words (described in Section 2.3). We call them Ext AMBIPUN, Sim AMBIPUN, and Gen AMBIPUN.

<sup>6</sup><https://beta.openai.com/docs/engines/instruct-series-beta>

### 3.4 Evaluation

**Automatic Evaluation** Previous works use two metrics to automatically evaluate the quality of the generated puns. First, both (He et al., 2019) and Luo et al. (2019) report the the *unusualness* of an  $n$ -length output, which is defined as the normalized log-probability of each token  $x_i$  subtracted by its training probability under a language model

$$U \triangleq -\frac{1}{n} \log \left( p(x_1, \dots, x_n) / \prod_{i=1}^n p(x_i) \right). \quad (2)$$

Although He et al. (2019) further show that unusualness does not correlate well with human ratings of puns, we still follow the same procedure. Besides, (Luo et al., 2019) and (Yu et al., 2018) use distinct unigram and bigrams (Li et al., 2015) to measure the *diversity* of each system on a sentence-level. However, we observe that certain systems tend to generate sentences with fixed patterns. Namely, those generation models lack diversity corpus-wise, but could still gain high distinctiveness score sentence-wise. Hence, we propose to measure the diversity from both levels. We also report the the average sentence length produced.

**Human Evaluation** It is known that currently available automatic evaluation metrics could not reflect the nuances of language, including humor and creativity. Following the procedure of previous works (Yu et al., 2018; He et al., 2019), we randomly shuffle and select 100 sentences for human evaluation. We collected our human ratings on Amazon Mechanical Turk (AMT). For each sentence, three workers are explicitly given the target pun word. We first ask them to judge if a given sentence is a pun sentence on a binary scale. Then, they are asked the questions: “How funny is this

Model	Avg Sequence Length	Sentence-level Diversity		Corpus-level Diversity		Unusualness
		Dist-1	Dist-2	Dist-1	Dist-2	
Few-shot GPT3	12.3	<b>37.1</b>	<b>80.4</b>	94.5	91.5	0.09
Neural Pun	12.6	30.2	73.0	91.3	90.5	0.22
Pun GAN	9.7	<u>34.6</u>	71.9	90.2	87.6	<b>0.47</b>
Sim AMBIPUN	13.4	32.4	77.1	92.9	91.2	0.26
Gen AMBIPUN	<u>13.5</u>	32.8	77.8	93.6	91.2	0.26
Ext AMBIPUN	<b>14.0</b>	31.7	<u>78.7</u>	<b>96.3</b>	<b>92.3</b>	<u>0.28</u>
Human	14.1	36.6	81.9	95.5	92.4	0.35

Table 2: Results of automatic evaluation on average sequence length, sentence-level and corpus-level diversity, and the unusualness scores. Boldface denotes the best performance and underline denotes the second best performance among systems. We compare with three strong baselines: Few-shot GPT3, Neural Pun (Yu et al., 2018), and Pun GAN (Luo et al., 2019), and three variations of our own method: similarity-based context generation (Sim AMBIPUN), generative context generation (Gen AMBIPUN) and extraction-based context generation (Ext AMBIPUN). Note that unusualness has been shown to have weak correlation with human ratings by He et al. (2019).

Model	Success Rate	Fun	Coherence
Few-shot GPT3	13.0%	1.82	<b>3.77</b>
Neural Pun	35.3%	2.17	3.21
Pun GAN	35.8%	2.28	2.97
Sim AMBIPUN	45.5%	2.69	3.38
Gen AMBIPUN	50.5%	2.94	3.53
Ext AMBIPUN	<b>52.2%</b>	<b>3.00</b>	3.48
Human	70.2%	3.43	3.66

Table 3: Human evaluation results on all the pun generation systems. We show the success rates, and average scores of funniness and coherence of each system. Overall, Ext AMBIPUN performs the best.

sentence?” and “How coherent or fluent is this sentence?” on a scale from 1 (not at all) to 5 (extremely). We provide detailed instructions and examples with explanation for each criteria. We also adopt attention questions and qualification types to make our collected results more reliable. For pun judgement (binary), we take the majority vote among three workers, while for funniness and coherence (1 to 5), we take the average ratings. We then use the pairwise kappa coefficient to measure the inter-annotator agreement. The average inter-annotator agreement of all raters for pun success, funniness and coherence are 0.55, 0.48 and 0.40, meaning that our collected results are reliable.

## 4 Results and Analysis

### 4.1 Pun Generation Results

**Automatic Evaluation** Results of the automatic evaluation can be seen in Table 2. We compare three baselines, three variations of our own model AMBIPUN, and the human written puns. First, the average length of our generated sentence are closest

to human written sentences. Although our baseline Pun-GAN has higher distinct ratio at sentence level, we observe that is mainly due to a short sequence length. Moreover, it falls short in corpus-level diversity, meaning that the generated sentences have similar syntax patterns. On the other hand, our Ext AMBIPUN achieves the highest corpus-level diversity. As for unusualness, Pun GAN also obtains unreasonably high score compared with gold. A possible explanation is that the model generate incoherently, which is also verified by our human ratings. Our experimental results resonate with the findings by He et al. (2019) that unusualness does not correlate well with human ratings.

**Human Evaluation** Results from the automatic evaluation can be seen in Table 3. We evaluate the success rate, funniness, and coherence of the generated outputs. The superiority of our models are obvious. All three of our systems outperform the baselines in terms of success rate and funniness. On the other hand, GPT3 could generate even more coherently than humans.

**Analysis between extractive and generative method.** Interestingly, Ext AMBIPUN has higher success rates and is funnier than Gen AMBIPUN, indicating that extracting salient words from human written sentences could introduce more surprising and uncommon words than language models. We posit that those atypical words refresh people’s eyes and thus boost the pun success rate as well as the funniness score. On the other hand, we also tried to equip GPT3 with greater creativity by top-k sampling with a large temperature  $T$ . However, larger  $T$ s also result in arbitrary responses that human may find unreadable. We hope our discovery could

<b>Pun word</b>	<b>Irrational</b>		
Sense 1	Real but not expressible as the quotient of two integers		
Sense 2	Not consistent with or using reason		
<b>Model</b>	<b>Example</b>	<b>Pun</b>	<b>Funny</b>
GPT3	I can't make a decision with all this irrationality going on.	No	1.4
Neural Pun	Note that this means that there is an irrational problem.	Yes	2.4
Pun-GAN	It can be use the irrational system.	No	1.2
Ext AMBIPUN	I have an irrational <u>paranoia</u> about <u>mathematical integers</u> .	Yes	<b>3.8</b>
Gen AMBIPUN	My <u>calculator</u> is unjust and <u>illogic</u> . It's irrational.	Yes	3.4
Human	Old math teachers never die, they just become irrational.	Yes	<b>3.8</b>

<b>Pun word</b>	<b>Drive</b>		
Sense 1	A journey in a vehicle (usually an automobile)		
Sense 2	The trait of being highly motivated		
<b>Model</b>	<b>Example</b>	<b>Pun</b>	<b>Funny</b>
GPT3	I am exhausted, I need a nap before I can drive any more.	No	2.0
Neural Pun	It is that it can be use to drive a variety of function?	No	1.6
Pun-GAN	In he drive to the first three years.	No	1.2
Ext AMBIPUN	What do you call a <u>genius</u> with cunning drive? <u>racecar driver</u> .	Yes	3.6
Gen AMBIPUN	I have the <u>determination</u> to <u>travel</u> to my <u>destination</u> . But i don't have the drive.	Yes	4.0
Human	A boy saving up for a car has a lot of drive.	Yes	<b>4.2</b>

Table 4: We show generated sentences for the word 'Irrational' and 'Drive' in the above table, along with their two senses. For the results of our top performing models Gen AMBIPUN and Ext AMBIPUN, we underline the context words that are related to each sense. All the generations are evaluated by external annotators, not the authors.

	Success Rate
Beginning	46.7%
Middle	52.0%
End	54.7%

Table 5: The pun success rate sentences based on their position annotated by human.

draw the community's attention to those traditional techniques for creative generation.

## 4.2 Case Study

To better understand the advantages of our method from a qualitative perspective, we conduct a case study for the pun word "Irrational" and "Drive" and evaluate the generated samples by our top performing models as well as the baselines. The generated outputs along with human evaluation results can be seen in Table 4. For both the examples pun words, at most one of the baselines successfully generates a punning sentence. As discussed earlier, one possible reason is the absence of both senses. On the other hand, both Ext AMBIPUN and Sim AMBIPUN introduce context words for the two senses and thus are able to generate of high quality puns that almost match the human written puns in terms of the funniness score.

## 4.3 The Position of Pun Words

As is mentioned in Section 2.4, we play with the position of the pun word in the prompt given to the candidate generation model. We try three variants

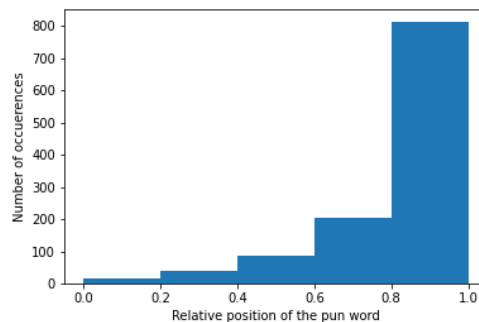


Figure 3: Analysis of the position of pun word in human written puns. The y-axis indicates the number of sentences and the x-axis indicates the position of pun word on a scale from 0 (start) to 1 (end). The analysis is based on 1,163 human written sentences included in the SemEval 2017 Task 7 (Miller et al., 2017).

by putting the target pun word at the start, in the middle, and at the end. For each variant, we then ask Mechanical Turkers to judge if the given sentences are puns. Again, each sentence is rated by three Turkers and we take the majority answer if the workers disagree. Results from this analysis can be seen in Table 5. We observe that people find a sentence more likely to be a pun when the target word appears at the end.

To verify such hypothesis, we also calculate the position of the pun words of 1,163 human written pun sentences and report the distribution in Figure 3. The histogram corroborates with the human

448 annotated samples in that both suggest that keeping  
449 the pun word at the end of the sentence generates  
450 funnier puns. Theory of humor which says that the  
451 "joke" in a funny sentences some towards the end  
452 of the sentence (Shahaf et al., 2015) validates our  
453 analysis.

## 454 5 Related Works

### 455 5.1 Creative Language Generation

456 **Pun generation.** Many of the previous works  
457 on pun generation have focused on phonological  
458 or syntactic pattern rather than semantic pattern  
459 (Miller and Gurevych, 2015; Hong and Ong, 2009;  
460 Petrović and Matthews, 2013; Valitutti et al., 2013)  
461 thus lacking creativity and flexibility. He et al.  
462 (2019) make use of local-global surprisal principle  
463 to generate homophonic puns and Yu et al. (2020)  
464 uses constrained lexical rewriting for the same task.  
465 Hashimoto et al. (2018) use a retrieve and edit ap-  
466 proach to generate homographic puns and Yu et al.  
467 (2018); Luo et al. (2019) propose complex neural  
468 model architecture such as constrained language  
469 model and GAN, and do not put emphasis on the  
470 linguistic structure of puns. We identify their ab-  
471 sence of both the senses as a shortcoming and build  
472 our approach from there.

473 **Figurative language generation.** There have  
474 been several attempts to generate other types  
475 of figurative language such as metaphor, simile  
476 (Chakrabarty et al., 2020b), sarcasm, etc. Yu  
477 and Wan (2019) use metaphorically used verbs  
478 to generate metaphors in an unsupervised fashion.  
479 (Chakrabarty et al., 2021) generates metaphors us-  
480 ing symbolism and discriminative decoding. Stowe  
481 et al. (2021) study diverse metaphor generation  
482 using conceptual mapping. Mishra et al. (2019)  
483 propose a modular architecture for unsupervised  
484 sarcasm generation Chakrabarty et al. (2020a) use  
485 commonsense knowledge for the same task. Tian  
486 et al. (2021) leverage semantic structure and com-  
487 monsense and counterfactual knowledge to gener-  
488 ate hyperbole.

489 As for stories, recent works focus on hierar-  
490 chical story generation that first plans a plot and  
491 then writes stories based on the storyline (Martin  
492 et al., 2018; Yao et al., 2019; Fan et al., 2019).  
493 Goldfarb-Tarrant et al. (2020) incorporates SRL  
494 extracted event representations in storylines with  
495 several event related decoding objectives.

**Humor generation.** With the recent advent of  
496 diverse datasets (Hasan et al. (2019), Mittal et al.  
497 (2021), Yang et al. (2021)), it has become easier  
498 to detect and generate humor. While large pre-  
499 trained model have fairly successful in detection,  
500 humor generation still remains an unsolved prob-  
501 lem. Therefore, humor generation is usually stud-  
502 ied in a specific setting. Petrović and Matthews  
503 (2013) generates joke of the type 'I like my X like  
504 I like my Y, Z'. Garimella et al. (2020) develops a  
505 model to fill blanks in madlibs format to generate  
506 humorous sentences and Yang et al. (2020) edit  
507 headlines to make them funny. More research is  
508 required to generate humorous sentences that are  
509 not constrained by their semantic structure. 510

### 511 5.2 Pun detection

512 Being able to detect puns can be an essential step  
513 to generate them as will be evident in the com-  
514 ing sections. SemEval 2017 Task 7 (Miller et al.,  
515 2017) introduced the challenge of pun detection,  
516 location detection and sense interpretation for ho-  
517 mographic and homophonic puns. It also released  
518 a dataset which becomes the backbone of our and  
519 several other related works. Diao et al. (2019) make  
520 use of Gated Attention network to detection ho-  
521 mophonic puns. Zou and Lu (2019) introduces a  
522 tagging schemes which lets them detect puns as  
523 well as their location. They apply this approach to  
524 both homophonic and homographic puns.

## 525 6 Conclusion

526 We propose a novel approach towards homographic  
527 puns generation. Unlike previous works that are  
528 mathematically heavy, our approach is back-boned  
529 by the humor theory that ambiguity is achieved  
530 by the context. Both automatic and human evalua-  
531 tions show that our model AMBIPUN outperforms  
532 the current state-of-the-art model by a significant  
533 margin. We also analyze why our extraction-based  
534 variation are more humorous than generation-based  
535 variation, and investigate the role of the position  
536 of pun words, which corresponds with human writ-  
537 ten sentences. In future work, we want to make a  
538 step further and explore the part of speech tags by  
539 filtering out the context words based on their POS  
540 tags and make combinations accordingly. Another  
541 interesting direction could be to apply our proposed  
542 approach to set phrases, which also make use of  
543 different senses.



544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593  
594  
595  
596  
597

## References

Melanie Booth-Butterfield and Melissa Wanzer. 2018. [Humor in interpersonal communication](#).

Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*.

Tuhin Chakrabarty, Debanjan Ghosh, Smaranda Muresan, and Nanyun Peng. 2020a.  [\$\mathcal{R}\$ : Reverse, retrieve, and rank for sarcasm generation with commonsense knowledge](#). *arXiv preprint arXiv:2004.13248*.

Tuhin Chakrabarty, Smaranda Muresan, and Nanyun Peng. 2020b. Generating similes effortlessly like a pro: A style transfer approach for simile generation. *arXiv preprint arXiv:2009.08942*.

Tuhin Chakrabarty, Xurui Zhang, Smaranda Muresan, and Nanyun Peng. 2021. [Mermaid: Metaphor generation with symbolism and discriminative decoding](#).

Ciprian Chelba, Tomas Mikolov, Mike Schuster, Qi Ge, Thorsten Brants, Phillipp Koehn, and Tony Robinson. 2013. [One billion word benchmark for measuring progress in statistical language modeling](#).

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [Bert: Pre-training of deep bidirectional transformers for language understanding](#).

Yufeng Diao, Hongfei Lin, Liang Yang, Xiaochao Fan, Di Wu, Dongyu Zhang, and Kan Xu. 2019. [Heterographic pun recognition via pronunciation and spelling understanding gated attention network](#). In *The World Wide Web Conference, WWW '19*, page 363–371, New York, NY, USA. Association for Computing Machinery.

Angela Fan, Mike Lewis, and Yann Dauphin. 2019. Strategies for structuring story generation. In *ACL*.

John R Firth. 1957. A synopsis of linguistic theory, 1930-1955. *Studies in linguistic analysis*.

Aparna Garimella, Carmen Banea, Nabil Hossain, and Rada Mihalcea. 2020. [“judge me by my size \(noun\), do you?” YodaLib: A demographic-aware humor generation framework](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 2814–2825, Barcelona, Spain (Online). International Committee on Computational Linguistics.

Marjan Ghazvininejad, Xing Shi, Yejin Choi, and Kevin Knight. 2016. Generating topical poetry. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1183–1191.

Seraphina Goldfarb-Tarrant, Tuhin Chakrabarty, Ralph Weischedel, and Nanyun Peng. 2020. Content planning for neural story generation with aristotelian rescoring. *arXiv preprint arXiv:2009.09870*.

Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative adversarial nets. *Advances in neural information processing systems*, 27. 598  
599  
600  
601  
602

David Graff, Junbo Kong, Ke Chen, and Kazuaki Maeda. 2003. English gigaword. *Linguistic Data Consortium, Philadelphia*, 4(1):34. 603  
604  
605

Md Kamrul Hasan, Wasifur Rahman, AmirAli Bagher Zadeh, Jianyuan Zhong, Md Iftekhar Tanveer, Louis-Philippe Morency, and Mohammed (Ehsan) Hoque. 2019. [Ur-funny: A multimodal language dataset for understanding humor](#). *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. 606  
607  
608  
609  
610  
611  
612  
613  
614

Tatsunori B. Hashimoto, Kelvin Guu, Yonatan Oren, and Percy Liang. 2018. [A retrieve-and-edit framework for predicting structured outputs](#). 615  
616  
617

He He, Nanyun Peng, and Percy Liang. 2019. Pun generation with surprise. *arXiv preprint arXiv:1904.06828*. 618  
619  
620

Bryan Anthony Hong and Ethel Ong. 2009. [Automatically extracting word relationships as templates for pun generation](#). In *Proceedings of the Workshop on Computational Approaches to Linguistic Creativity*, pages 24–31, Boulder, Colorado. Association for Computational Linguistics. 621  
622  
623  
624  
625  
626

Justine T Kao, Roger Levy, and Noah D Goodman. 2016. A computational model of linguistic humor in puns. *Cognitive science*, 40(5):1270–1285. 627  
628  
629

Omar Khattab and Matei Zaharia. 2020. [Colbert: Efficient and effective passage search via contextualized late interaction over bert](#). 630  
631  
632

Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2015. A diversity-promoting objective function for neural conversation models. *arXiv preprint arXiv:1510.03055*. 633  
634  
635  
636

Louis G Lippman and Mara L Dunn. 2000. Contextual connections within puns: Effects on perceived humor and memory. *The journal of general psychology*, 127(2):185–197. 637  
638  
639  
640

Fuli Luo, Shun Yao Li, Pengcheng Yang, Baobao Chang, Zhifang Sui, Xu Sun, et al. 2019. Pun-gan: Generative adversarial network for pun generation. *arXiv preprint arXiv:1910.10950*. 641  
642  
643  
644

Lara J Martin, Prithviraj Ammanabrolu, Xinyu Wang, William Hancock, Shruti Singh, Brent Harrison, and Mark O Riedl. 2018. Event representations for automated story generation with deep neural nets. In *Thirty-Second AAAI Conference on Artificial Intelligence*. 645  
646  
647  
648  
649  
650

651	Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. In <i>Advances in neural information processing systems</i> , pages 3111–3119.	
652		
653		
654		
655		
656	Tristan Miller and Iryna Gurevych. 2015. <a href="#">Automatic disambiguation of English puns</a> . In <i>Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 719–729, Beijing, China. Association for Computational Linguistics.	
657		
658		
659		
660		
661		
662		
663	Tristan Miller, Christian F Hempelmann, and Iryna Gurevych. 2017. Semeval-2017 task 7: Detection and interpretation of english puns. In <i>Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)</i> , pages 58–68.	
664		
665		
666		
667		
668	Abhijit Mishra, Tarun Tater, and Karthik Sankaranarayanan. 2019. A modular architecture for unsupervised sarcasm generation. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 6144–6154.	
669		
670		
671		
672		
673		
674		
675	Anirudh Mittal, Pranav Jeevan, Prerak Gandhi, Diptesh Kanojia, and Pushpak Bhattacharyya. 2021. "so you think you're funny?": Rating the humour quotient in standup comedy.	
676		
677		
678		
679	Lili Mou, Rui Yan, Ge Li, Lu Zhang, and Zhi Jin. 2015. Backward and forward language modeling for constrained sentence generation. <i>arXiv preprint arXiv:1512.06612</i> .	
680		
681		
682		
683	Saša Petrović and David Matthews. 2013. <a href="#">Unsupervised joke generation from big data</a> . In <i>Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 228–232, Sofia, Bulgaria. Association for Computational Linguistics.	
684		
685		
686		
687		
688		
689	Fanchao Qi, Lei Zhang, Yanhui Yang, Zhiyuan Liu, and Maosong Sun. 2020. Wantwords: An open-source online reverse dictionary system. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 175–181.	
690		
691		
692		
693		
694		
695	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. <a href="#">Exploring the limits of transfer learning with a unified text-to-text transformer</a> .	
696		
697		
698		
699		
700	Juan Ramos. 2003. Using tf-idf to determine word relevance in document queries.	
701		
702	Stuart Rose, Dave Engel, Nick Cramer, and Wendy Cowley. 2010. <a href="#">Automatic Keyword Extraction from Individual Documents</a> , pages 1 – 20.	
703		
704		
	Dafna Shahaf, Eric Horvitz, and Robert Mankoff. 2015. <a href="#">Inside jokes: Identifying humorous cartoon captions</a> . In <i>Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '15</i> , page 1065–1074, New York, NY, USA. Association for Computing Machinery.	705
		706
		707
		708
		709
		710
	Kevin Stowe, Tuhin Chakrabarty, Nanyun Peng, Smaranda Muresan, and Iryna Gurevych. 2021. <a href="#">Metaphor generation with conceptual mappings</a> .	711
		712
		713
	Yufei Tian, Arvind krishna Sridhar, and Nanyun Peng. 2021. <a href="#">HypoGen: Hyperbole generation with commonsense and counterfactual knowledge</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 1583–1593, Punta Cana, Dominican Republic. Association for Computational Linguistics.	714
		715
		716
		717
		718
		719
		720
	Alessandro Valitutti, Hannu Toivonen, Antoine Doucet, and Jukka Toivanen. 2013. "let everything turn well in your wife": Generation of adult humor using lexical constraints. volume 2.	721
		722
		723
		724
	Orion Weller and Kevin Seppi. 2020. The rjokes dataset: a large scale humor collection. In <i>Proceedings of The 12th language resources and evaluation conference</i> , pages 6136–6141.	725
		726
		727
		728
	Ziqing Yang, Yiming Cui, Zhipeng Chen, Wanxiang Che, Ting Liu, Shijin Wang, and Guoping Hu. 2020. <a href="#">TextBrewer: An Open-Source Knowledge Distillation Toolkit for Natural Language Processing</a> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations</i> , pages 9–16, Online. Association for Computational Linguistics.	729
		730
		731
		732
		733
		734
		735
		736
	Zixiaofan Yang, Shayan Hooshmand, and Julia Hirschberg. 2021. <a href="#">CHoRaL: Collecting humor reaction labels from millions of social media users</a> . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 4429–4435, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	737
		738
		739
		740
		741
		742
		743
	Lili Yao, Nanyun Peng, Ralph Weischedel, Kevin Knight, Dongyan Zhao, and Rui Yan. 2019. Plan-and-write: Towards better automatic storytelling. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 33, pages 7378–7385.	744
		745
		746
		747
		748
	Zhiwei Yu, Jiwei Tan, and Xiaojun Wan. 2018. A neural approach to pun generation. In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1650–1660.	749
		750
		751
		752
		753
	Zhiwei Yu and Xiaojun Wan. 2019. <a href="#">How to avoid sentences spelling boring? towards a neural approach to unsupervised metaphor generation</a> . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 861–871, Minneapolis, Minnesota. Association for Computational Linguistics.	754
		755
		756
		757
		758
		759
		760
		761

762 Zhiwei Yu, Hongyu Zang, and Xiaojun Wan. 2020. [Ho-](#)  
763 [mophonic pun generation with lexically constrained](#)  
764 [rewriting](#). In *Proceedings of the 2020 Conference on*  
765 *Empirical Methods in Natural Language Processing*  
766 *(EMNLP)*, pages 2870–2876, Online. Association for  
767 Computational Linguistics.

768 Lei Zhang, Fanchao Qi, Zhiyuan Liu, Yasheng Wang,  
769 Qun Liu, and Maosong Sun. 2020. Multi-channel  
770 reverse dictionary model. In *Proceedings of the AAAI*  
771 *Conference on Artificial Intelligence*, pages 312–319.

772 Yanyan Zou and Wei Lu. 2019. [Joint detection and](#)  
773 [location of English puns](#). In *Proceedings of the 2019*  
774 *Conference of the North American Chapter of the*  
775 *Association for Computational Linguistics: Human*  
776 *Language Technologies, Volume 1 (Long and Short*  
777 *Papers)*, pages 2117–2123, Minneapolis, Minnesota.  
778 Association for Computational Linguistics.

## Appendix

### A Humor Classifier Results for Selecting Puns

To further discuss the accuracy and recall of our humor classifier, we show a representative output in Table 6. The table contains a few selected sentences ranked by the humor classifier. We also label each sentence as *yes*, *no*, and *maybe* to indicate if it is a pun or not. As discussed in the methodology, we train our classifier on humor dataset. As puns are an important part of humor generation, this model can help rule out some options. Basic theories of humor such as incongruity and surprise apply to both of them. As can be seen in the table, our classifier is able to successfully pull aside unfunny or non-coherent sentences. Looking at the examples at the top and the middle, it can be observed that some better examples are classified lower than others. Making this observation across many pun words, we decided to use the classifier only to rule out the bottom third samples. For the rest of the generations, we randomly sample them.

On manual observation, we realised that when we cherry-pick samples, we're able to find many sentences that meet our expectations. Therefore, building a classifier that can accurately find these sentences can increase the accuracy by a large margin. We treat this as an opportunity for future work.

### B More Examples of Generated Puns

We compile more examples generated by AM-BIPUN in Table 7 for the following pun words: *sentence*, *case*, *bugs*, *delivery*. This table further supports our point that on manual selection our model is able to generate human-like sentences.



Sentence	Rank	Pun
What's the interest rate on a home mortgage? No interest.	1	Yes
My bank said I think they're interested in me. I said no.	2	No
My girlfriend said she had an interest in banking so i loan her a quarter	3	Yes
I have no interest in being a guardian. It's free.	4	Maybe
I've never had interest placed on borrowings. It's a waste of time.	5	Yes
Why did the republican attack the bank? Because it was in its interest.	6	Maybe
What is the republican's strategy? The interest rate.	7	No
What is the most dispensable interest in investment?	8	No
If trump had an interest in president he would make it an president-of-interest.	9	No

Table 6: An example of candidate pun sentences ranked by the humor classifier. As can be seen, the model is able to rule out non-pun sentences but fails to pick out high-quality ones.

Target word	<b>sentence</b>
Sense 1	A string of words satisfying the grammatical rules of a language
Sense 2	(Criminal law) a final judgment of guilty in a criminal case and the punishment that is imposed
1	The word jail is a sentence.
2	What's the punishment for using antonyms in a sentence syntax is it a sentence?
3	I'm sorry I said the sentence was too long but punishments are endless.
4	The sentence in the dictionary doesn't sound very guilty.
Target word	<b>case</b>
Sense 1	A portable container for carrying several objects
Sense 2	A statement of facts and reasons used to support an argument
1	What's the most durable luggage for a detective? jury case
2	A jury just found a container of leather there's no reason to argue it's a case
3	What do you call a container used for investigation research? a case study
4	Why did the cardboard get into a trial? because it was an investigation case
Target word	<b>bugs</b>
Sense 1	General term for any insect or similar creeping or crawling invertebrate
Sense 2	A fault or defect in a computer program, system, or machine
1	Why did the garden restart its computer? it had bugs in it.
2	What do you call a pest that's slow programmer? bugs bug
3	Why did the compost crash? it had bugs in it.
4	What do you call a bug that's disgusting? a glitch in the internet
Target word	<b>delivery</b>
Sense 1	the act of delivering or distributing something (as goods or mail)
Sense 2	your characteristic style or manner of expressing yourself orally
1	What did the letter say to the parcel? clear delivery!
2	What do you call a trucking truckdriver with no articulation? delivery driver.
3	The distribution center has a pronunciation dictionary. it's a delivery service
4	What do you call a parcel with no dialogue and an accent? delivery service.

Table 7: More examples generated by Ext AMBIPUN.