# Evaluating Performance Metrics in Bias Mitigation for Generated Patient Vignettes Using Large Language Models

Shaurya Kumar, Chinmay Agrawal

## Abstract

Recent advancements in Large Language Models (LLMs) have introduced new possibilities for automating medical vignette generation, a process that could significantly aid in patient diagnosis. However, concerns about bias in the outputs of these models necessitate a comprehensive evaluation of their performance and fairness. Our study leverages a sophisticated LLM pipeline to generate medical vignettes from a dataset curated from clinical guidelines, biomedical literature, and clinical trials. The LLM pipeline involves multiple stages: a context retriever that loads and embeds clinical guidelines, biomedical literature, and clinical trials; a prompt generator that creates specific prompts based on user queries; a strong LLM (e.g., GPT-4) that generates vignettes; and a post-processing stage that refines the outputs for benchmarking. This workflow ensures that the generated vignettes are contextually relevant and tailored to the specified medical conditions.To explain the motivations for using these specific metrics, we delve into the statistical and model-based scoring methods. ROUGE and BLEU scores measure n-gram overlaps, while GPTScore and Semantic Entropy assess semantic similarity and fluency. Our dataset comprises clinical guidelines and biomedical literature, and we detail the sources and examples used in our analysis. By performing rigorous tests, we present the numerical outcomes of these metrics, demonstrating their significance in evaluating the quality of generated patient vignettes. In conclusion, our research emphasizes the importance of selecting appropriate evaluation metrics to ensure the accuracy and reliability of patient vignettes generated by LLMs. These findings contribute to improving the effectiveness of LLMs in clinical settings, ultimately aiding in unbiased and accurate patient diagnosis.

## Introduction

New developments in Large Language Models (LLMs) have introduced transformative possibilities for automating the generation of medical vignettes, which can significantly enhance the accuracy and efficiency of patient diagnosis. The project aims to develop a robust tool that integrates these advanced models with comprehensive datasets, ensuring that the vignettes generated are both contextually accurate and clinically relevant.
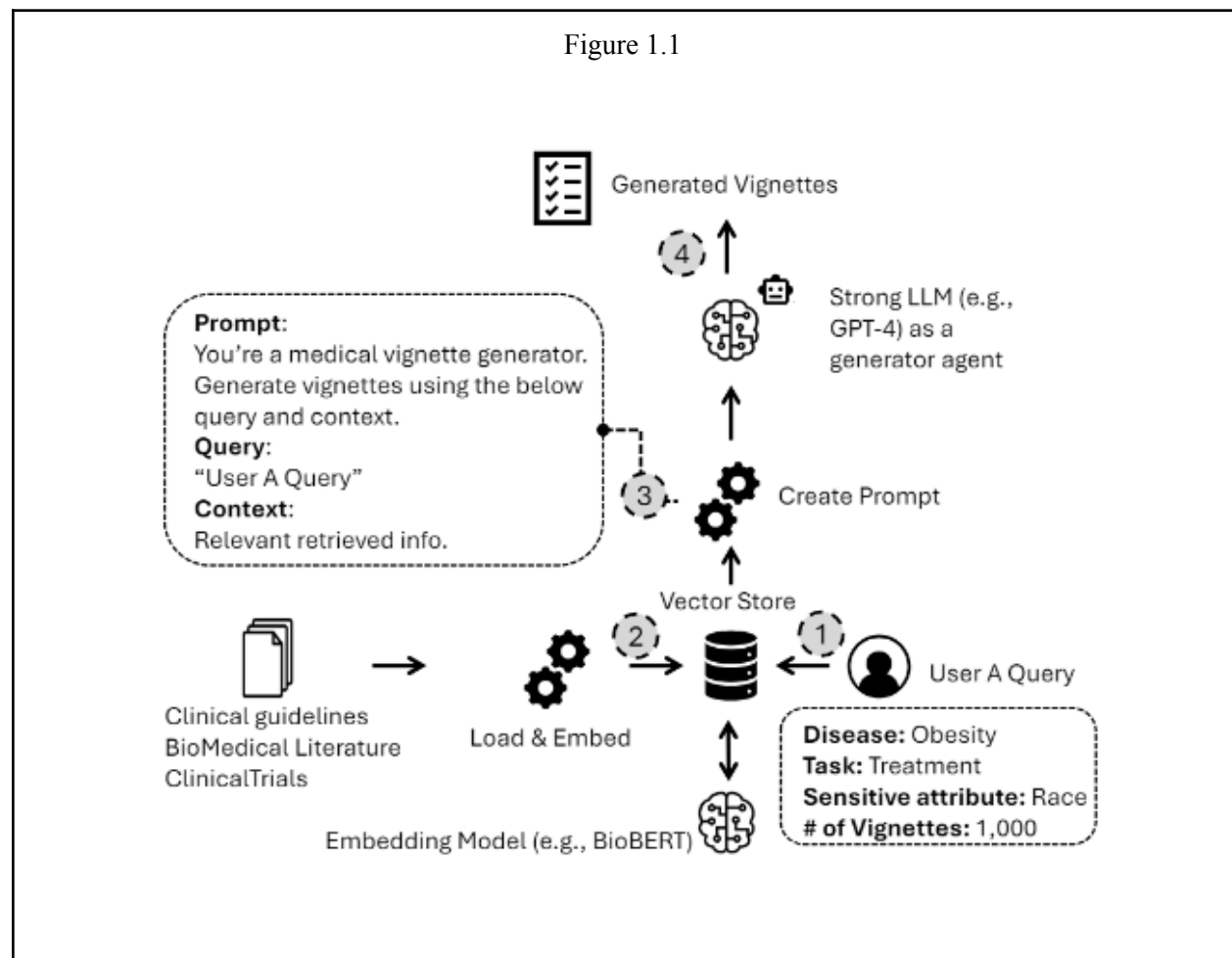
Our tool utilizes a multi-stage LLM pipeline designed to process and generate medical vignettes. The pipeline begins with a context retriever that loads and embeds data from clinical guidelines, biomedical literature, and clinical trials. This embedded information is then used by a prompt generator to create specific queries tailored to the user's needs. The generated prompts are processed by a powerful LLM, such as GPT-4, which produces detailed medical vignettes. Finally, a post-processing stage refines these vignettes to ensure they meet high standards of relevance and accuracy.

A critical aspect of our project is the evaluation of the summarization quality of the generated vignettes. To assess this, we employ a variety of evaluation metrics, including ROUGE and BLEU scores for measuring n-gram overlaps, and GPTScore and Semantic Entropy for evaluating semantic similarity and

fluency. These metrics provide a comprehensive view of the performance and reliability of the vignettes, highlighting areas for improvement and ensuring that the outputs are unbiased and clinically useful.

By rigorously testing and presenting the outcomes of these metrics, we demonstrate their importance in the evaluation process. Our findings underscore the necessity of selecting appropriate evaluation metrics to guarantee the accuracy and dependability of LLM-generated patient vignettes. This research contributes to the ongoing efforts to refine LLM applications in clinical settings, ultimately aiding in the development of unbiased and precise diagnostic tools.

Figure 1.1



Generated Vignettes

Strong LLM (e.g., GPT-4) as a generator agent

**Prompt:**
You're a medical vignette generator. Generate vignettes using the below query and context.
**Query:**
"User A Query"
**Context:**
Relevant retrieved info.

Create Prompt

Vector Store

Clinical guidelines
BioMedical Literature
ClinicalTrials

Load & Embed

User A Query

Embedding Model (e.g., BioBERT)

**Disease:** Obesity
**Task:** Treatment
**Sensitive attribute:** Race
**# of Vignettes:** 1,000

As seen in Figure 1.2 , the high-level workflow diagram of an AI-driven system details the interactions between components to process user inputs and generate refined outputs. The steps below detail the process:

1. **User Inputs**: The process begins with the user providing input data, which could be a query, a request for information, or any other type of data needing processing.
2. **Context Retriever**: This component is responsible for accessing and retrieving relevant information from the knowledge base. It ensures that the system has all the necessary context to understand and process the user input effectively.
3. **Knowledge Base**: A structured repository where all relevant information and data are stored. This database is essential for providing contextual information to the system.
4. **Prompt Generator**: Using the retrieved context, this component generates prompts or initial responses. These prompts serve as the foundation for further processing by the system.
5. **LLM (Large Language Model)**: The core processing unit of the system, the LLM takes the generated prompts and performs sophisticated language modeling tasks. It generates an initial response based on the given prompts and context.
6. **Refinement**: This step involves refining the initial response generated by the LLM. The refinement process may include additional processing, adjustments, or improvements to ensure the output is accurate and meets the user's needs.
7. **Post Processing**: After refinement, the output undergoes post-processing to enhance readability, format the response appropriately, or apply any final adjustments. This ensures the output is polished and ready for delivery.
8. **Benchmark**: The final output is compared against benchmark criteria to evaluate its quality, accuracy, and relevance. This step ensures the system's performance is consistently meeting the desired standards.
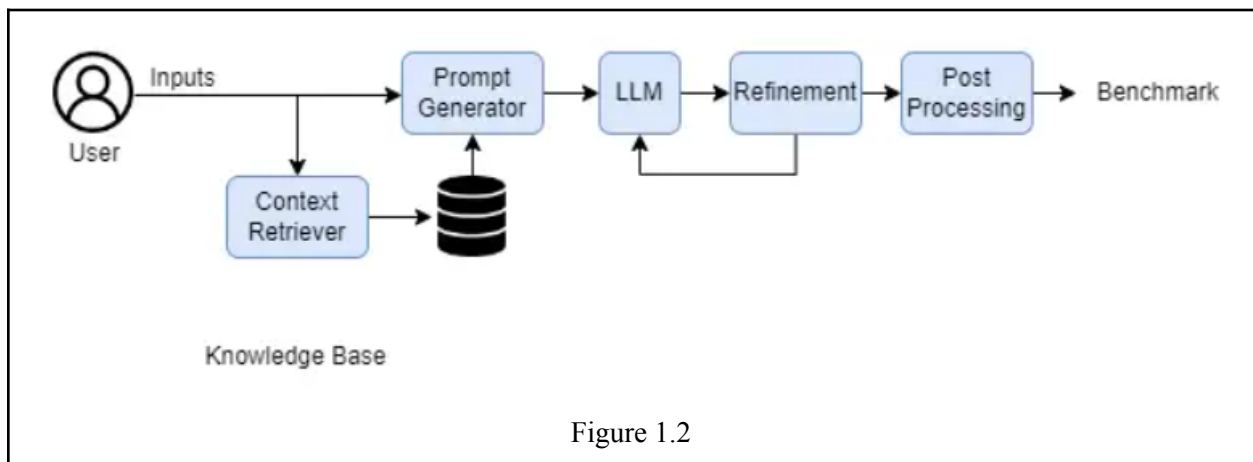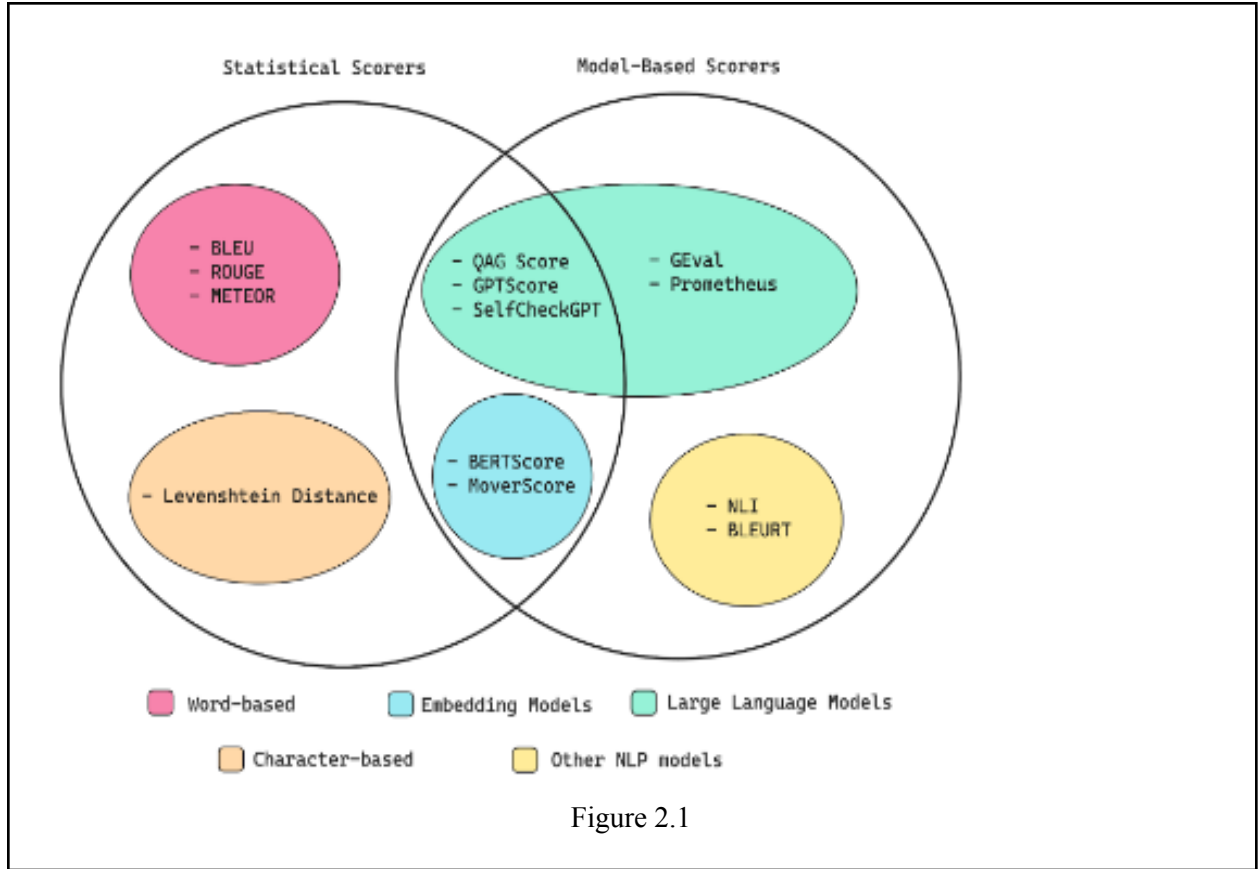


Figure 1.2

# Methods

Figure 2.1

As seen in figure 2.1 , various Statistical and Model - Based scores are utilized as benchmark performance measures for Large Language Models. The use of every one would be redundant hence this project aims to explore the most efficient ones. Below outlines the 3 scores used in our Project along with the motivations behind their choosing and implementation and benefits:

1. **BLEU (Bilingual Evaluation Understudy)**
   - *Motivation*: BLEU is primarily used to evaluate the quality of machine-translated text against a set of high-quality reference translations. The core motivation behind BLEU is to measure the correspondence of the candidate translation to the references at the word or phrase level, focusing on precision—the proportion of words in the machine-generated output that appear in the reference translations.
   - *Utility*: BLEU is computationally inexpensive and correlates moderately well with human judgment, making it a standard baseline in machine translation evaluation. It's particularly useful for assessing consistency in vocabulary usage and basic grammatical structures in translated texts.

$$\text{BLEU} = \text{BP} \cdot \exp\left(\sum_{n=1}^{N} w_n \log p_n\right).$$

Key:
BP - brevity penalty, which penalizes short generated texts.
$W_n$ - weight assigned to n-gram precision
$P_n$ - precision of n-grams

## 2. ROUGE (Recall-Oriented Understudy for Gisting Evaluation)

- *Motivation*: ROUGE is designed to evaluate text summarization and other tasks where recall is as important as precision. It measures the overlap of n-grams, word sequences, and sometimes longer units like sentences between the generated text and reference texts.
- *Utility*: Unlike BLEU, which emphasizes precision, ROUGE captures how much of the reference content appears in the generated text, making it valuable for tasks where covering all relevant information is crucial, such as summarization.

ROUGE-N

$$= \frac{\sum_{S \in \{ReferemceSummaries\}} \sum_{gram_n \in S} Count_{match}(gram_n)}{\sum_{S \in \{ReferenceSummaries\}} \sum_{gram_n \in S} Count(gram_n)} \quad (1)$$

Key:

$Count_{match}(gram_n)$ - number of n-grams in the reference text that match the n-grams in the generated text.

$Count(gram_n)$ - total number of n-grams in the reference text.

## 3. Semantic Entropy

- *Motivation*: While not explicitly shown in your diagram under the name "Semantic Entropy," if we consider metrics like METEOR or newer embeddings-based approaches, these are often motivated by a desire to capture semantic similarity more effectively than simple n-gram matching. Metrics like METEOR, which incorporate synonyms and stemming, or embeddings-based metrics, aim to understand the underlying meanings rather than just the surface string similarity.
- *Utility*: These metrics are useful when the precise words used are less important than the conveyed meaning, such as in dialogue systems, question answering, and more nuanced translations.

$$SE(x) = -\sum_c p(c \mid x) \log p(c \mid x) = -\sum_c \left( \left( \sum_{s \in c} p(\mathbf{s} \mid x) \right) \log \left[ \sum_{s \in c} p(\mathbf{s} \mid x) \right] \right).$$

```python
rouge_score.py > save_rouge_scores
1    import pandas as pd
2    from rouge import Rouge
3
4    # access to .txt files
5    def open_file(file_path):
6        with open(file_path, 'r') as file:
7            return file.read()
8
9
10   # score calculation using ROUGE technique
11   def sum_measurement_calculation(context, vignettes):
12       rouge = Rouge()
13       scores = rouge.get_scores(context, vignettes)
14       return scores[0]
15
16   # put results into new CSV
17   '''
18       explanation of each of the categories:
19           metric: ROUGE used (rouge-1 unigrams, rouge-2 bigrams, rouge-l longest common subsequence)
20           precision: comparison of generated text to reference text
21           recall: comparison of reference text to generated text
22           F1 score: harmonic mean of precision and recall scores
23
24       ROUGE score results are in decimal places, but correspond to a percentage score (i.e .3591 = 35.91%)
25   '''
26   def save_rouge_scores(scores, output_csv):
27       data = {
28           'Metric': [],
29           'Precision': [],
30           'Recall': [],
31           'F1':[]
32       }
33
34       for metric, score in scores.items():
35           data['Metric'].append(metric)
36           data['Precision'].append(score['p'])
37           data['Recall'].append(score['r'])
38           data['F1'].append(score['f'])
39
40       df = pd.DataFrame(data)
41       df.to_csv(output_csv, index=False)
42
43   # main function to run all at consecutively
44   def main(context_file, vignettes_file, output_csv):
45       context = open_file(context_file)
46       vignettes = open_file(vignettes_file)
47       scores = sum_measurement_calculation(context, vignettes)
48       save_rouge_scores(scores, output_csv)
49
50   # test case
51   context = 'context_hyperthyroidismabstract.txt'
52   vignettes = 'vignettes_hyperthyroidism.txt'
53   output_csv = 'result_scores.csv'
54
55   main(context, vignettes, output_csv)
```

```python
bleu_score.py > ...
  1    import pandas as pd
  2    import nltk
  3    from nltk.translate.bleu_score import sentence_bleu, SmoothingFunction
  4
  5    # Download the necessary NLTK resources
  6    nltk.download('punkt')
  7
  8    # Function to read the content of a file
  9    def read_file(file_path):
 10        with open(file_path, 'r') as file:
 11            return file.read()
 12
 13    # Function to calculate BLEU score
 14    def calculate_bleu(context, vignette):
 15        context_tokens = [nltk.word_tokenize(context)]
 16        vignette_tokens = nltk.word_tokenize(vignette)
 17
 18        # Calculate BLEU score with smoothing
 19        smoothie = SmoothingFunction().method4
 20        bleu_score = sentence_bleu(context_tokens, vignette_tokens, smoothing_function=smoothie)
 21        return bleu_score
 22
 23    # Function to save scores to a CSV file
 24    def save_scores_to_csv(score, output_csv):
 25        data = {
 26            'Metric': ['BLEU'],
 27            'Score': [score]
 28        }
 29
 30        df = pd.DataFrame(data)
 31        df.to_csv(output_csv, index=False)
 32
 33    # Main function to execute the whole process
 34    def main(context_file, vignettes_file, output_csv):
 35        context = read_file(context_file)
 36        vignettes = read_file(vignettes_file)
 37        score = calculate_bleu(context, vignettes)
 38        save_scores_to_csv(score, output_csv)
 39
 40    # Example usage
 41    context_file = 'context_hyperthyroidismabstract.txt'
 42    vignettes_file = 'vignettes_hyperthyroidism.txt'
 43    output_csv = 'result_scores.csv'
 44
 45    main(context_file, vignettes_file, output_csv)
```

```python
1    import numpy as np
2    import torch
3    from transformers import GPT2Tokenizer, GPT2LMHeadModel, GPT2Model
4    import nltk
5    import math
6
7    model_name = "gpt2"
8    tokenizer = GPT2Tokenizer.from_pretrained(model_name)
9    model_1 = model = GPT2Model.from_pretrained(model_name)
10   model_2 = GPT2LMHeadModel.from_pretrained(model_name)
11   model.eval()
12
13   # Word entropy functions
14
15 ∨ def word_entropy(text):
16       words = nltk.word_tokenize(text)
17       embeddings = []
18       for word in words:
19           inputs = tokenizer(word, return_tensors="pt")
20           with torch.no_grad():
21               outputs = model_1(**inputs)
22           embeddings.append(outputs.last_hidden_state.mean(dim=1).squeeze().numpy())
23
24       total_words = len(embeddings)
25       entropy = 0.0
26
27       sim_matrix = np.zeros((total_words, total_words))
28       for i in range(total_words):
29           for j in range(i, total_words):
30               if i == j:
31                   sim_matrix[i, j] = 1.0
32               else:
33                   cos_sim = np.dot(embeddings[i], embeddings[j]) / (np.linalg.norm(embeddings[i]) * np.linalg.norm(embeddings[j]))
34                   sim_matrix[i, j] = cos_sim
35                   sim_matrix[j, i]= cos_sim
36
37       probabilities = sim_matrix.mean(axis=1)
38       probabilities = probabilities / probabilities.sum()
39
40       for prob in probabilities:
41           entropy -= prob *math.log2(prob)
42
43       return entropy
44
45 ∨ def compute_word_entropy(vignettes):
46       scores = {}
47       count = 1
48       for vignette in vignettes:
49           scores.update({str(count) : word_entropy(vignette)})
50           count += 1
51       return scores
52
53   # Semantic entropy functions
54
55 ∨ def semantic_entropy(prob_dist):
56       epsilon = 1e-9
57       prob_dist = np.clip(prob_dist, epsilon, 1.0)
58       entropy = -np.sum(prob_dist * np.log2(prob_dist))
59       return entropy
60
61 ∨ def get_probs(text):
62       inputs = tokenizer(text, return_tensors='pt')
63       with torch.no_grad():
64           outputs = model_2(**inputs, labels=inputs['input_ids'])
65           logits = outputs.logits
66       probs = torch.softmax(logits, dim=-1).numpy()
67       return probs
68
69 ∨ def compute_semantic_entropy(vignettes):
70       scores = {}
71       count = 1
72       for vignette in vignettes:
73           probs = get_probs(vignette)
74           entropies = [semantic_entropy(probs[0, i]) for i in range(probs.shape[1])]
75           sem_entropy = np.mean(entropies)
76           scores.update({str(count) : sem_entropy})
77           count += 1
78       return scores
```

# Results

When running tests, an example set of patient contexts and their corresponding vignettes were used for BLEU and ROUGE, each with 50 patients with a variety of differences in symptoms and medical and physical history.

**Rouge CSV Context/Output:**

```
result_scores.csv
1    Metric,Precision,Recall,F1
2    rouge-1,0.35940409683426444,0.5230352303523035,0.42604856029333515
3    rouge-2,0.11431143114311432,0.14126807563959956,0.12636815425960266
4    rouge-l,0.3538175046554935,0.5149051490514905,0.41942604373704373
```

**metric: ROUGE used (rouge-1 unigrams, rouge-2 bigrams, rouge-l longest common subsequence)**
**precision: comparison of generated text to reference text**
**recall: comparison of reference text to generated text**
**F1 score: harmonic mean of precision and recall scores**
**ROUGE score results are in decimal places, but correspond to a percentage score (i.e .3591 = 35.91%)**

**BLEU CSV Context/Output:**

```
result_scores.csv
1    Metric,Score
2    BLEU,0.10271064682349844
```

**score: ranges from 0 to 1, higher score indicates better quality of summary**

**Semantic Entropy Context/Example Output:**
**Unable to conduct tests with semantic entropy, here is an example output file:**

| (a) Scenario 1: No semantic equivalence | | | (b) Scenario 2: Some semantic equivalence | | |
|---|---|---|---|---|---|
| Answer $s$ | Likelihood $p(s \mid x)$ | Semantic likelihood $\sum_{s \in c} p(s \mid x)$ | Answer $s$ | Likelihood $p(s \mid x)$ | Semantic likelihood $\sum_{s \in c} p(s \mid x)$ |
| Paris | 0.5 | 0.5 | **Paris** | 0.5 | 0.9 |
| Rome | 0.4 | 0.4 | **It's Paris** | 0.4 | |
| London | 0.1 | 0.1 | London | 0.1 | 0.1 |
| Entropy | 0.94 | 0.94 | Entropy | 0.94 | 0.33 |

**Output: A numerical value representing the entropy of the text, typically ranging from 0 (low entropy) to higher positive values (high entropy).**

**Meaning: This value quantifies the level of unpredictability or diversity in the semantic content of the generated text.**

# Discussion

In the provided scenarios comparing "No Semantic Equivalence" and "Some Semantic Equivalence," the Semantic Entropy metric effectively quantifies the difference in semantic consistency of potential answers. In the first scenario, Semantic Entropy does not change, reflecting no consideration for semantic similarity among answers like Paris, Rome, and London. However, in the second scenario, acknowledging semantic equivalents like "It's Paris" alongside "Paris" significantly reduces entropy, demonstrating this metric's utility in measuring the nuanced understanding required in contexts like medical vignette generation, where capturing semantic depths is crucial for accuracy.

The results from the ROUGE and BLEU evaluation metrics as depicted in the uploaded images reveal insights into the quality of machine-generated medical vignettes compared to reference texts. ROUGE scores, represented in decimal form but interpreted as percentages, show a moderate level of n-gram overlap between the generated texts and the references, indicating a reasonable but not perfect capture of necessary information. For instance, the highest ROUGE-L score suggests that the longest common subsequences align fairly well, yet the overall scores suggest there is room for improvement in capturing finer details and nuances of the reference texts. On the other hand, the BLEU score, which is on a scale from 0 to 1, reflects a more critical evaluation, focusing on the precision of n-gram matching, and indicates that the exact wording and phrasing have lower correspondence. This suggests that while the gist of the medical content might be captured, the specific phrasing and terminology usage requires closer alignment with standard medical language to enhance accuracy and reliability in clinical settings.

# References