

# Benchmarking Google Translate

Philosophy of Computation Lab IV

Henry Blanchette

April 5, 2019

## Abstract

TODO

## 1 Introduction

## 2 Transcript Setup

I selected transcript sections that reflect a variety of writing styles, including modern English, English, technical writing, storytelling, and English translated from other languages.

Transcripts:

T1. The Bible, Genesis

T2. Melville's Moby Dick, Chapter 1

T3. Marriam-Webster English Dictionary, definition of Abdicate

T4. Bedau's patentsample.txt

T5. Shakespeare's Henry IV, Part 1

## 3 Translation task

Start with a sequence of languages  $L_0, \dots, L_n$  and a transcript in  $L_0$ , called the *original transcript*. GT translates the  $L_0$ -transcript to  $L_1$ , and then translates the resulting  $L_1$ -transcript to  $L_2$ , and so on until the transcript has been translated to  $L_n$ . At the end, there is left a  $L_n$ -transcript. Then, GT translates this  $L_n$ -transcript back to  $L_0$  - the result is the *processed transcript*. The differences between the original and processed transcripts are measured to rate GT's success at this task. The goal of this scoring is to rate GT according to how well it preserve the meaning and grammatic structure of the original transcript.

## 4 Translation Success Measure

I rated GT's success at the task by how close the processed transcript was to the original transcript in terms of meaning and grammar. For each of these dimensions, I categorized an ranking of ordered performance classes.

**Grammar** is how properly-constructed the processed transcript is according to the rules of  $L_0$  and the grammatical structure of the original transcript. The following are the classes of grammar performance I partitioned in order of increasing success. They are meant to be “evenly spaced” in the space of deviations from the exact grammar and structure of the original transcript.

Grammar Class	Description
$G1$	completely confused
$G2$	mostly confused
$G3$	often confused
$G4$	sparsely confused
$G5$	effectively perfect

**Meaning** is how close the processed transcript is to the original transcript in meaning. The following are the classes of meaning performance I partitioned in order of increasing success. They are meant to be “evenly spaced” in the space of deviations from reflecting the exact meaning of the original transcript.

Meaning Class	Description
$M1$	irrelevant
$M2$	sparsely relevant
$M3$	often relevant
$M4$	mostly accurate
$M5$	effectively perfect

## 5 Experiment 1: Translation Ring with Well-Documented Languages

### 5.1 Language Setup

I selected from the top 5 languages (without English) by native speaker count . I hypothesized that this would correlate with the amount of effort that Google has put into training translations to and from these languages, which should yield more coherent and thus easier-to-score processed texts from this task.

The following are the languages used in this experiment in order of decreasing native speakers count:

- L1. Chinese (simplified)
- L2. Spanish
- L3. Hindi
- L4. Arabic
- L5. Portuguese

### 5.2 Experimental Design

I ran each transcript through the following trials, where the selected languages and their order was chose randomly:

Trial 1: Chinese  $\rightarrow$  Arabic  $\rightarrow$  Spanish  $\rightarrow$  Portuguese  $\rightarrow$  Hindi

Trial 2: Hindi  $\rightarrow$  Chinese  $\rightarrow$  Portuguese  $\rightarrow$  Arabic  $\rightarrow$  Spanish

Trial 3: Hindi  $\rightarrow$  Spanish  $\rightarrow$  Arabic  $\rightarrow$  Chinese  $\rightarrow$  Portuguese

Trial 4: Chinese  $\rightarrow$  Arabic  $\rightarrow$  Spanish  $\rightarrow$  Portuguese  $\rightarrow$  Hindi

Trial 5: Arabic  $\rightarrow$  Chinese  $\rightarrow$  Portuguese  $\rightarrow$  Spanish  $\rightarrow$  Hindi

### 5.3 Predictions

I predicted that GT will perform better on translating these languages than it w

### 5.4 Results

### 5.5 Analysis

## 6 Experiment 2: Under-Documented Language Translation Ring

### 6.1 Language Setup

I selected from the bottom 5 languages by native speakers that Google Translate supports .

Languages:

L1. Nepali

L2. Sinhala

L3. Greek

L4. Hungarian

L5. Zulu

### 6.2 Experimental Design

I ran each transcript through the following trials, where the selected languages and their order was chose randomly:

Trial 1: Zulu  $\rightarrow$  Hungarian  $\rightarrow$  Nepali  $\rightarrow$  Sinhala  $\rightarrow$  Greek

Trial 2: Nepali  $\rightarrow$  Hungarian  $\rightarrow$  Greek  $\rightarrow$  Sinhala  $\rightarrow$  Zulu

Trial 3: Nepali  $\rightarrow$  Sinhala  $\rightarrow$  Zulu  $\rightarrow$  Greek  $\rightarrow$  Hungarian

Trial 4: Sinhala  $\rightarrow$  Greek  $\rightarrow$  Nepali  $\rightarrow$  Zulu  $\rightarrow$  Hungarian

Trial 5: Hungarian  $\rightarrow$  Greek  $\rightarrow$  Sinhala  $\rightarrow$  Zulu  $\rightarrow$  Nepali

**6.3 Predictions**

**6.4 Results**

**6.5 Analysis**

**7 Conclusion**

**Bibliography**