# Translated Wikipedia Biographies

English -> Spanish (516 KB)

English-> German (517 KB)

The Translated Wikipedia Biographies dataset has been designed to evaluate gender accuracy in long text translations (multiple sentences or passages). The set has been designed to analyze common gender errors in machine translation like incorrect gender choices in anaphora resolutions, possessives and gender agreement.

**PUBLISHER(S)**

Google LLC

**FUNDING**

Google LLC

**INDUSTRY TYPE**

Corporate - Tech

**FUNDING TYPE**

Private Funding

**DATASET AUTHORS**

Anja Austermann, Google Michelle Linch, Google Romina Stella, Google Kellie Webster, Google

**DATASET CONTACT**

translate-gender-challenge-sets@google.com

**DATASET PURPOSE(S)**

Testing

**KEY APPLICATION(S)**

Machine Translation, Gender Accuracy

**ACCESS COST**

Open Access

**PRIMARY MOTIVATION(S)**

Study gender accuracy in translations beyond the sentence in demographic and occupations diversity for fairness research.

**INTENDED AND/OR SUITABLE USE CASE(S)**

To evaluate gender accuracy on translations beyond the sentence (multiple sentences or passages). The set is focused on the presence of this specific linguistic

phenomena to evaluate the most common contextual error:

- Pro-drop (Spanish → English)
- Neutral to gender-specific possessives (Spanish → English)
- Gender agreement (English → Spanish, German)

| PRIMARY DATA TYPE(S) | DATASET SNAPSHOT | | DESCRIPTION OF CONTENT |
|---|---|---|---|
| Non-Sensitive Public Data about people | Total Instances | 138 | This dataset is based on publicly available data on public and/or historical figures (Wikipedia articles) at a given snapshot in time. |
| | Masculine biographies (entities) | 63 | |
| | Masculine biographies (countries) | 51 | The dataset has 138 instances and each instance contains the first 8 to 15 sentences from a Wikipedia article. Articles are written in native English and have been professionally translated to Spanish and German. 126 of these instances represent a person with an associated stated gender and 12 are related with rock bands or sport teams (considered genderless). |
| | Feminine biographies (entities) | 63 | |
| | Feminine biographies (countries) | 57 | |
| | Rock bands & sport teams (entities) | 12 | |
| | Rock bands & sport teams (countries) | 12 | |

SOURCE DATASET(S)

- Source text: English Wikipedia
- Target text: Professional translations

| PRIMARY DATA MODALITIE(S) | EXAMPLE OF ACTUAL DATA POINT* | HOW TO INTERPRET A DATAPOINT |
|---|---|---|
| Textual data | - **Source language:** en <br> - **Target Language:** de <br> - **Source text:** Kaisa-Leena Mäkäräinen (born 11 January 1983) is a Finnish former world-champion and 3-time world-cup-winning biathlete, who currently competes for Kontiolahden Urheilijat. <br> - **Translated text:** Kaisa-Leena Mäkäräinen (geboren am 11. Januar 1983) ist eine ehemalige Weltmeisterin und 3-malige Weltcup-Siegerin im Biathlon aus Finnland, die derzeit für Kontiolahden Urheilijat antritt. <br> - **Document ID:** 1 <br> - **Sentence ID:** 1-1 <br> - **Perceived gender associated with the entity:** Female <br> - **Name of the Entity:** Kaisa Mäkäräinen <br> - **Link to the original Wikipedia article:** https://en.wikipedia.org/wiki/Kaisa_M%C3%A4k%C3%A4r%C3%A4inen | **Each datapoint** refers to a central entity that can be a person (stated as feminine or masculine), a rock band or a sport team (considered genderless). <br><br> **Each entity** is represented by a long text translation (multiple connected sentences or continuous passage referring to that main entity). <br><br> Each datapoint will have the following data: <br><br> - **Source language:** Language of the original text <br> - **Target Language:** Language of the translation <br> - **Source text:** Text from Wikipedia in source language (special characters and quotes removed) <br> - **Translated text:** Translation of the Wikipedia source text into the target text <br> - **Document ID:** ID generated to identify all the sentences belonging to the same passage. <br> - **Sentence ID:** Composed by the Document ID and Sentence number in the passage. <br> - **Perceived gender associated with the entity:** identified as Female, Male, Neutral <br> - **Name of the Entity:** Name of the main entity according Wikipedia <br> - **Link to the original Wikipedia article:** Link to the Wikipedia article at the time of extraction. Please consider that content in Wikipedia articles can be modified so differences may be found if the article has been re-edited. |

| LICENSE TYPE(S) | LICENSE BREAKDOWN | LICENSE PERMISSIONS (*CC-BY-SA 3.0*) |
|---|---|---|
| CC-BY-SA 3.0 | Source text has been extracted from English Wikipedia articles, which is made available under the CC-BY-SA 3.0 Unported license. All the rest is synthetic data. | <ul><li>Share — copy and redistribute the material in any medium or format.</li><li>Adapt — remix, transform, and build upon the material for any purpose, even commercially.</li><li>Attribution — You must give appropriate credit, provide a link to the license, and indicate if changes were made</li></ul> |
| **STATUS** <br><br> Limited Maintenance | **DATASET VERSION** <br><br> Version       1.0 <br> Last Updated    06/2021 <br> First Release     06/2021 | **FIRST EDITION** <br><br> Original data was collected late in 2020 and translated at the beginning of 2021. <br><br> **MAINTENANCE PLAN** <br><br> No refreshes planned, but the dataset may be updated to incorporate feedback.. |
| **DATA COLLECTION METHOD(S)** <br><br> Scraped <br><br> Independent Paid Professional(s) | **DATA SOURCES BY COLLECTION METHOD(S)** <br><br> <ul><li>**Scraped:** English Wikipedia (source text)</li><li>**Translation:** Independent paid professional human translations (target text)</li><li>**Annotated:** Human added labels and metadata</li></ul> | **DATA SELECTION CRITERIA - SCRAPING** <br><br> <ul><li>Grouped people from Wikipedia according to their occupation, profession, job and/or activity.</li><li>Divided all these instances based on geographical diversity (optimizing for diversity at the country level), to mitigate the skew to Western-individuals (using regions from census.gov as a proxy of geographical diversity).</li><li>Focused on having equal representation of feminine and masculine entities. Labeled each instance as "Female" or "Male" based on these biographies containing gender-indicative terms to refer to the person (like she, he, woman, son, father etc.)</li></ul> |

**Note:** The set doesn't include non-binary individuals as we couldn't find enough instances to accurately reflect the community.

**EXCLUDED DATA**

- Quotes numbers from Wikipedia sentences were removed.
- Titles from the Wikipedia articles were excluded.
- Images were not considered. The dataset is just text.

**SUMMARIES OF DATA COLLECTION METHOD**

- **Scraping:** Sentences extracted from Wikipedia documents. (Source text)
- **Translations:** Source text has been professionally translated into the target language. For Spanish translations, guidance to focus on pronoun-drop sentences. (Target text)
- **Human added labels and metadata**: Added source and target languages, ids, entity names, links and perceived gender labels.

| SAMPLING METHOD(S) | SAMPLING BREAKDOWN | | SAMPLING CRITERIA |
|---|---|---|---|
| Stratified Sampling | Total data sampled | 2000 entities | • **Country diversity:** Entities that belong to countries that had at least 3 entities were discarded |
| | Sample size | 138 | • **Minimum text length:** 8 - 10 sentences |
| | | | • **Occupational Activity:** subjects played an active role in the field of their occupation, and the |

| | | | | |
|---|---|---|---|---|
| | | | wikipedia article pertains directly to their occupation | |

- **Perceived gender:** inferred based on gender-indicative words in descriptions provided within the article
- **Budgets:** within limits of budget available to project

| HUMAN ATTRIBUTE(S) | PERCEIVED GENDER DISTRIBUTION | | | |
|---|---|---|---|---|
| Perceived gender<br><br>Geography | **Perceived Masculine Biographies**<br><br>Individual Instances   63<br>Country Coverage   51 | **Perceived Feminine Biographies**<br><br>Individual Instances   63<br>Country Coverage   57 | **Genderless Articles (Rock bands & sports teams)**<br><br>Individual Instances   12<br>Country Coverage   12 | |

| GEOGRAPHIC DISTRIBUTION<br><br>**Biographies \***<br><br><br><br><br><br><br><br><br><br><br><br><br><br>\*organized by region for readability. | **Africa**<br>0.79% Cameroon<br>0.79% Central African Republic<br>0.79% Ethiopia<br>1.59% Ghana<br>1.59% Kenya<br>0.79% Liberia<br>0.79% Mauritania<br>0.79% Mauritius<br>0.79% Namibia<br>1.59% Nigeria<br>1.59% Senegal<br>0.79% South Africa<br>0.79% Tunisia<br>1.59% Uganda | 0.79% Philippines<br>0.79% Singapore<br>0.79% South Korea<br>0.79% Sri Lanka<br>0.79% Thailand<br>1.59% Taiwan<br><br>**Europe**<br>0.79% Armenia<br>0.79% Austria<br>0.79% Denmark<br>2.38% England<br>1.59% Finland<br>0.79% France<br>0.79% Georgia<br>1.59% Germany | **Latin America and the Caribbean**<br>0.79% Antigua and Barbuda<br>1.59% Argentina<br>0.79% Barbados<br>1.59% Brazil<br>0.79% Cayman Islands<br>1.59% Chile<br>0.79% Colombia<br>0.79% Cuba<br>0.79% Curaçao<br>0.79% Dominica<br>0.79% Dominican Republic | 0.79% Iraq<br>2.38% Israel<br>0.79% Jordan<br>1.59% Lebanon<br>0.79% Morocco<br>1.59% Pakistan<br>1.59% Turkey<br><br>**North America**<br>0.79% Bahamas<br>0.79% Belize<br>2.38% Canada<br>1.59% Jamaica<br>2.38% United States |

| | | | | |
|---|---|---|---|---|
| | 0.79%  Zambia<br>0.79%  Zimbabwe<br><br>**Asia**<br>1.59%  China<br>0.79%  Hong Kong<br>2.38%  India<br>0.79%  Indonesia<br>2.38%  Japan<br>0.79%  Malaysia<br>0.79%  Mongolia<br>0.79%  Nepal | 0.79%  Hungary<br>0.79%  Iceland<br>0.79%  Ireland<br>0.79%  Italy<br>0.79%  Lithuania<br>0.79% Netherlands<br>0.79%  Norway<br>1.59%  Russia<br>0.79%  Scotland<br>0.79%  Spain<br>0.79%  Sweden<br>0.79%  Ukraine<br>0.79%  Wales | 0.79%  Guatemala<br>0.79%  Mexico<br>0.79%  Paraguay<br>0.79%  Trinidad &<br>Tobago<br>0.79%  Uruguay<br>0.79%  Venezuela<br><br>**Near East**<br>0.79%  Algeria<br>0.79%  Egypt<br>2.38%  Iran | **Oceania**<br>2.38%  Australia<br>0.79%  Fiji<br>0.79%  Micronesia<br>2.38%  New<br>Zealand<br>0.79%  Palau<br>0.79%  Papua New<br>Guinea<br>0.79%  Tonga<br>0.79%  Tuvalu |
| **GEOGRAPHIC DISTRIBUTION**<br><br>**Bands and sport teams** | **Africa**<br>8.33%  Kenya<br>8.33%  Nigeria<br>8.33%  South Africa | **Asia**<br>8.33%  Japan<br>8.33%  South Korea<br>8.33%  India | **Europe**<br>8.33%  Spain<br>8.33%  Sweden<br>8.33%  Russia | **Latin America and<br>the Caribbean**<br>8.33%  Argentina<br>8.33%  Brazil<br><br>**Oceania**<br>8.33%  Australia |

| **LABELING METHOD(S)** | **LABEL TYPE(S)** | **LABELING PROCEDURE** |
|---|---|---|
| Human labels<br><br>Algorithmic labels | **Human labels**<br><ul><li>Perceived gender: Annotated by raters based on gender-indicative words on the source text.</li></ul>**Algorithmic labels**<br><ul><li>Document ID: generated by Google internal system.</li></ul> | **Human labels**<br><ul><li>Perceived gender labels are based on the presence of gender-indicative terms in the article as explained in the "Data Selection Criteria" section. The label "neutral" was used for rock bands and sports teams.</li></ul><br>**Algorithmic labels** |

| | |
|---|---|
| • Sentence ID: sequential number based on the location of the sentence in the paragraph<br>• URL: extracted from Wikipedia<br>• Entity name: extracted from Wikipedia. | • Entity Name was extracted from the title of the Wikipedia article. The URL redirects to the article version when the dataset was created.<br>• Document IDs were assigned based on document ordering. Sentence IDs are based on the location of the sentence in the document. |