

# How can Research Help in Reducing the Language Gap on Wikipedia

**WIKIMANIA**  
**STOCKHOLM**



**Hady Elsahar**  
wd:Q50290890  
[[User:Hadyelsahar]]  
@hadyelsahar



**Lucie-Aimée Kaffee**  
wd:Q37860261  
[[User:Frimelle]]  
@frimelle

- **Information Poverty & Language Gap** on Wikipedia

- Research Topics

- Scribe Project



**Key Points**



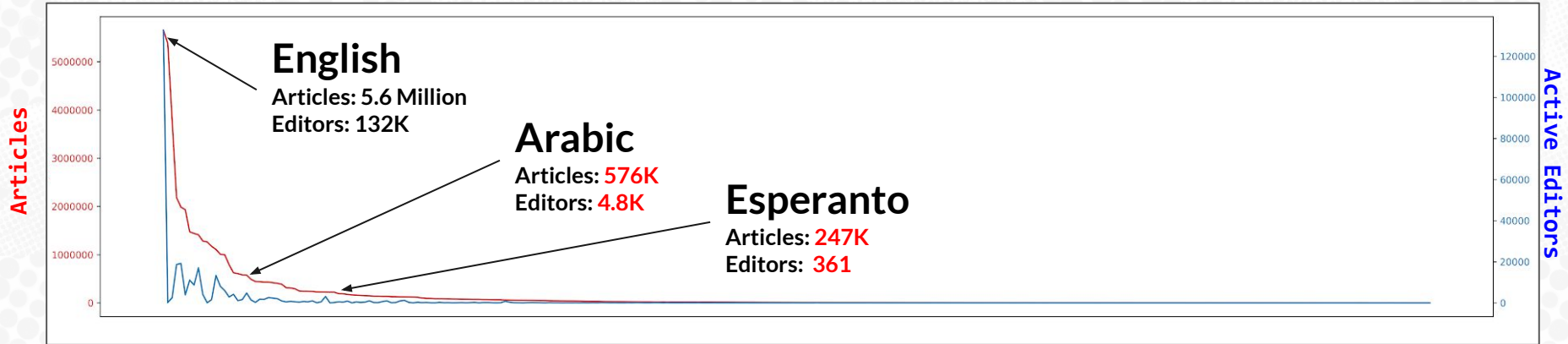
**WIKIPEDIA**  
The Free Encyclopedia

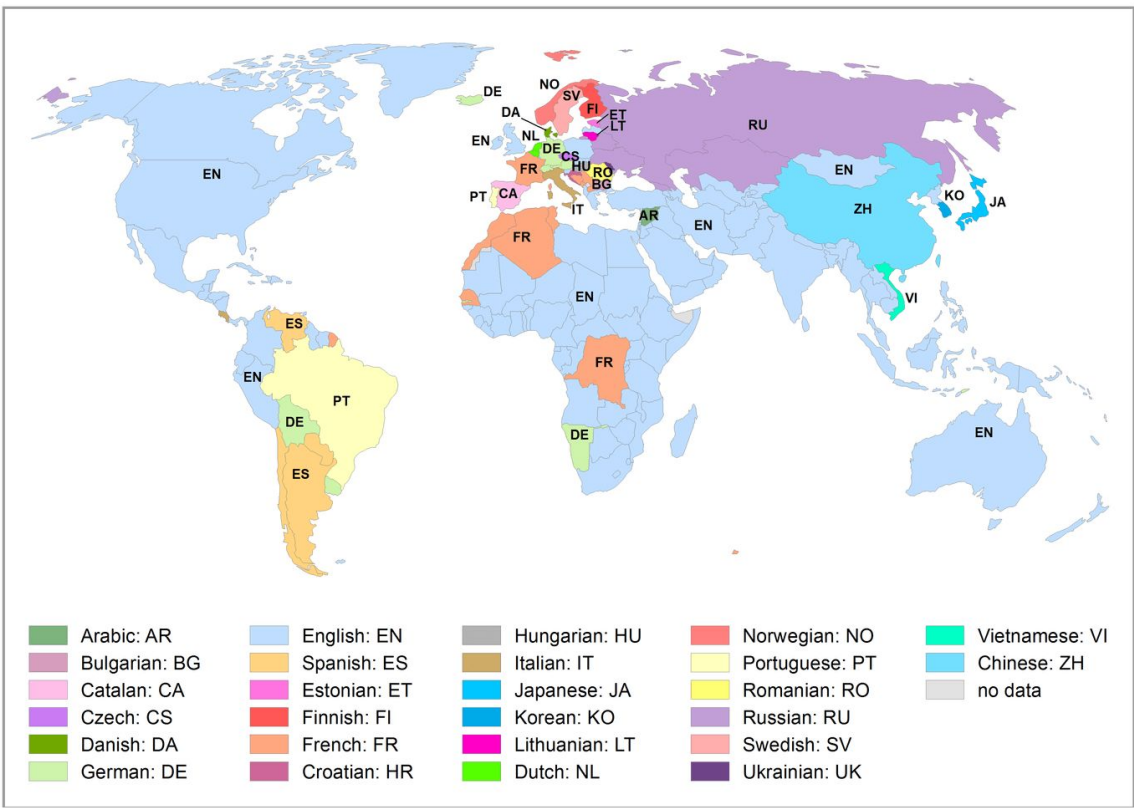
**=**

**Σ All Human Knowledge**

**?**

# Language Distribution of Articles and Active Users on Wikipedia





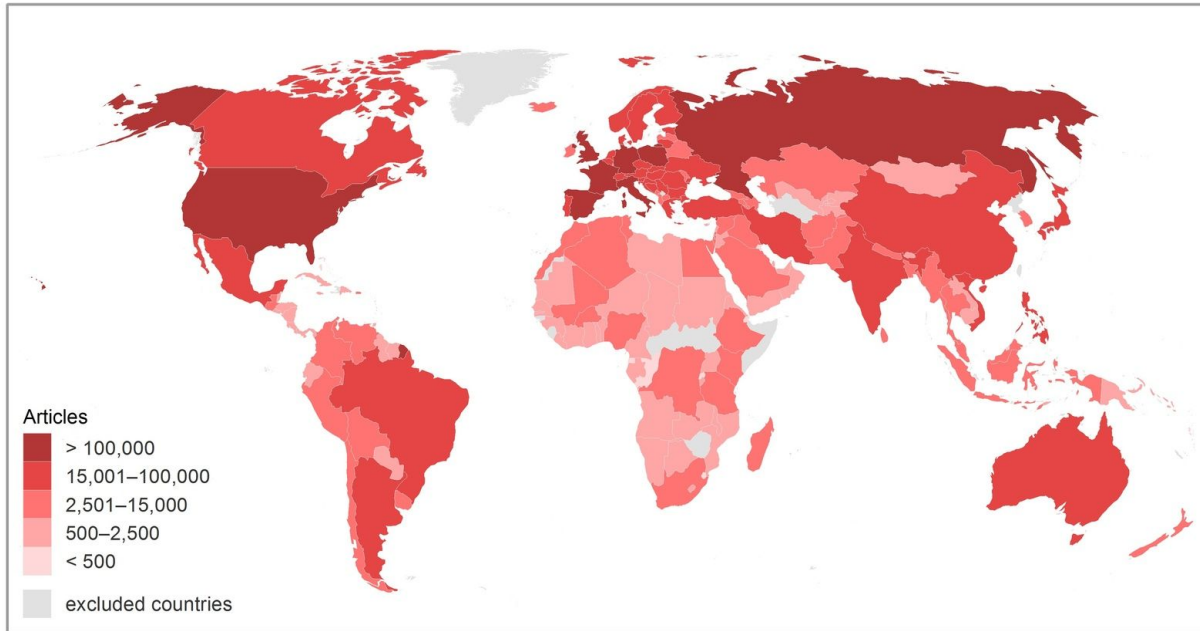
**Language Gap**

or

**Content Gap**

Few articles are  
being written by locals

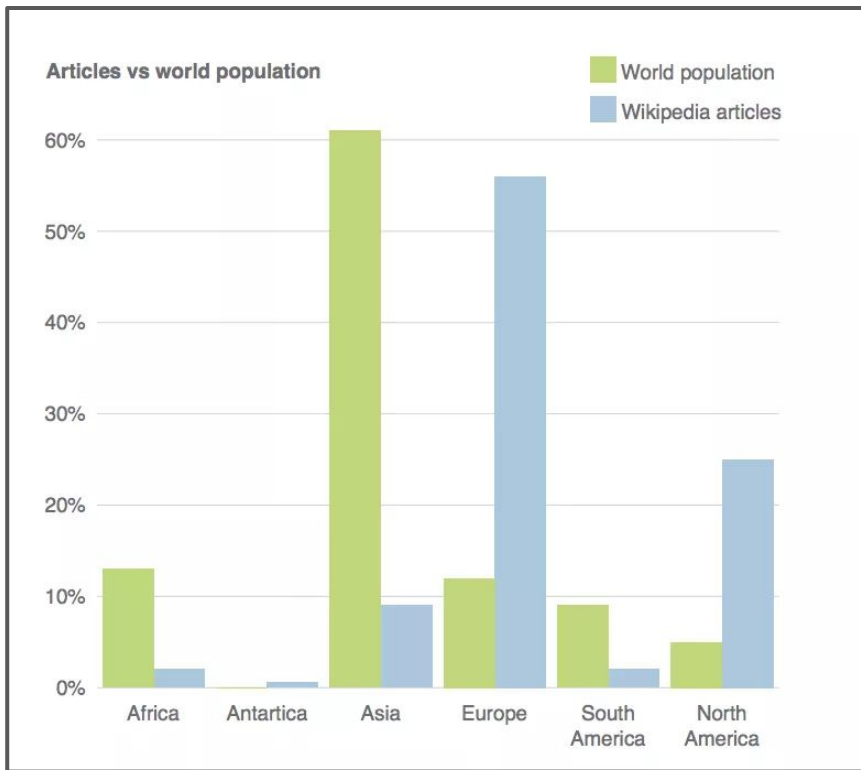
Dominant language of Wikipedia articles (geo tagged by country)  
(Graham et. al. 2014)



Number of geotagged Wikipedia articles per country (Graham et. al. 2014)

# Information Poverty

“There are more articles about Japan than about the whole middle east combined.”



Articles vs world population (Graham et. al. 2011)

# Information Poverty

The World relevant  
topics are not well  
represented on  
Wikipedia



**Solutions** to  
**Information poverty?**



## Current

- Automatic Machine Translation

## Alternatives

- Natural Language Generation
- Summarization
- Cross Lingual Word Embeddings



ويكيبيديا Translate page

< All translations Saved 2 minutes ago Publish

## Linguistic prescription Ⓞ

English [view page](#) العربية

6 categories categories 6

**Linguistic prescription**, or **prescriptive grammar**, is the attempt to lay down rules defining preferred or "correct" use of [language](#).<sup>[1][2]</sup> These rules may address such [linguistic](#) aspects as [spelling](#), [pronunciation](#), [vocabulary](#), [syntax](#), and [semantics](#). Sometimes informed by [linguistic purism](#),<sup>[3]</sup> such normative practices may suggest that some usages are incorrect, illogical, lack communicative effect, or are of low aesthetic value.<sup>[4][5]</sup> They may also include judgments on [socially proper](#) and [politically correct](#) language use.<sup>[6]</sup> Ⓞ

Linguistic prescriptivism may aim to establish a [standard language](#), teach what a particular society perceives as a correct form, or advise



# Content Translation Tool

Relies on automatic Machine Translation

- Google Translate
- Yandex Translate

ويكيبيديا Translate page

< All translations Saved just now Publish

Linguistic prescription معيارية لغوية

English العربية

6 categories categories 6

Linguistic prescription, or prescriptive grammar, is the attempt to lay down rules defining preferred or "correct" use of language.<sup>[1][2]</sup> These rules may address such linguistic aspects as spelling, pronunciation, vocabulary, syntax, and semantics. Sometimes informed by linguistic purism,<sup>[3]</sup> such normative practices may suggest that some usages are incorrect, illogical, lack communicative effect, or are of low aesthetic value.<sup>[4][5]</sup> They may also include judgments on socially proper and politically correct language use.<sup>[6]</sup>

Linguistic prescriptivism may aim to establish a standard language, teach what a particular society perceives as a correct form, or advise



# Content Translation Tool

Relies on automatic Machine Translation

- Google Translate
- Yandex Translate

ويكيبيديا Translate page

< All translations Saved just now Publish

## Linguistic prescription معيارية لغوية

English [view page](#) العربية categories 6

6 categories

Linguistic prescription, or prescriptive grammar, is the attempt to lay down rules defining preferred or "correct" use of language.<sup>[1][2]</sup> These rules may address such linguistic aspects as spelling, pronunciation, vocabulary, syntax, and semantics. Sometimes informed by linguistic purism,<sup>[3]</sup> such normative practices may suggest that some usages are incorrect, illogical, lack communicative effect, or are of low aesthetic value.<sup>[4][5]</sup> They may also include judgments on socially proper and politically correct language use.<sup>[6]</sup>

+ Add translation

Linguistic prescriptivism may aim to establish a standard language, teach what a particular society perceives as a correct form, or advise



# Content Translation Tool

Relies on automatic Machine Translation

- Google Translate
- Yandex Translate

## Linguistic prescription

## معيارية لغوية

English

[view page](#)

العربية

6 categories

**Linguistic prescription**, or **prescriptive grammar**, is the attempt to lay down rules defining preferred or "correct" use of language.<sup>[1][2]</sup> These rules may address such linguistic aspects as **spelling**, **pronunciation**, **vocabulary**, **syntax**, and **semantics**. Sometimes informed by **linguistic purism**,<sup>[3]</sup> such normative practices may suggest that some usages are incorrect, illogical, lack communicative effect, or are of low aesthetic value.<sup>[4][5]</sup> They may also include judgments on **socially proper** and **politically correct** language use.<sup>[6]</sup>

Linguistic prescriptivism may aim to establish a **standard language**, teach what a particular society perceives as a correct form, or advise

categories 6

**الوصفة اللغوية** ، أو **القواعد الإرشادية** ، هي محاولة لوضع قواعد تحدد الاستخدام المفضل أو "الصحيح" للغة .<sup>[1]</sup> <sup>[2]</sup> قد تناول هذه القواعد جوانب لغوية مثل الإملاء والنطق والمفردات وبناء الجملة والدلالات . في بعض الأحيان قد تكون هذه الممارسات المعيارية مطلعة في بعض الأحيان عن طريق التطهير اللغوي ،<sup>[3]</sup> قد تشير إلى أن بعض الاستخدامات غير صحيحة أو غير منطقية أو تنفقر إلى التأثير التواصلي أو تكون ذات قيمة جمالية منخفضة.<sup>[4]</sup> <sup>[5]</sup> وقد تتضمن أيضًا أحكامًا بشأن استخدام اللغة بشكل مناسب اجتماعيًا وصحيًا من الناحية السياسية .<sup>[6]</sup>

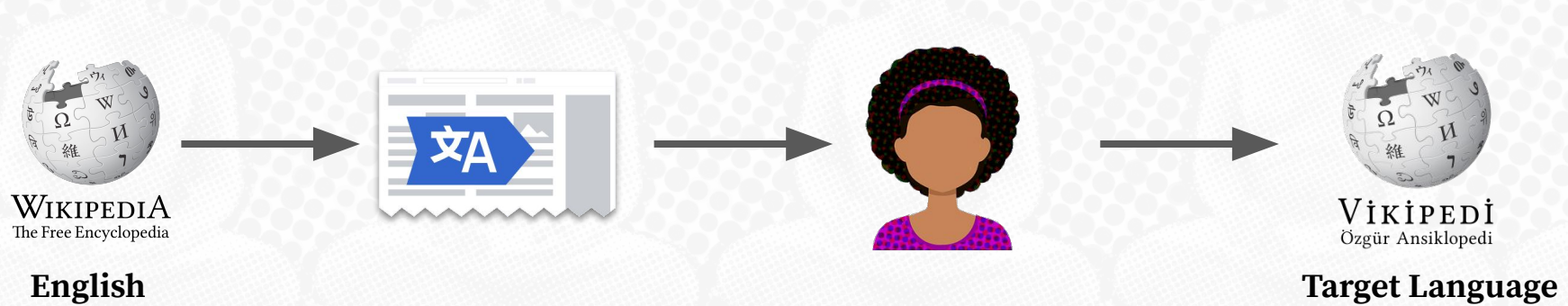


# Content Translation Tool

Relies on automatic Machine Translation

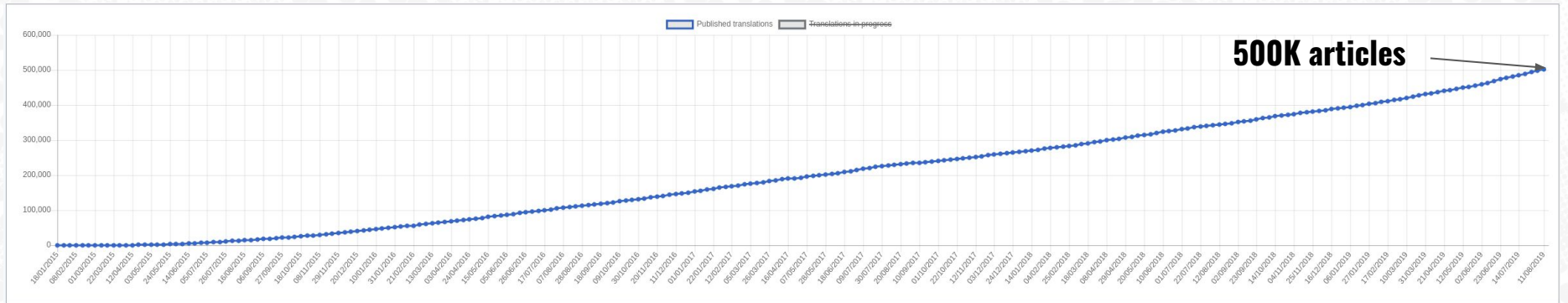
- Google Translate
- Yandex Translate

# Machine Translation with human in the loop.



# Machine Translation with human in the loop

Successful in filling part of the language Gap



# Translations to all languages

<https://en.wikipedia.org/wiki/Special:ContentTranslationStats>

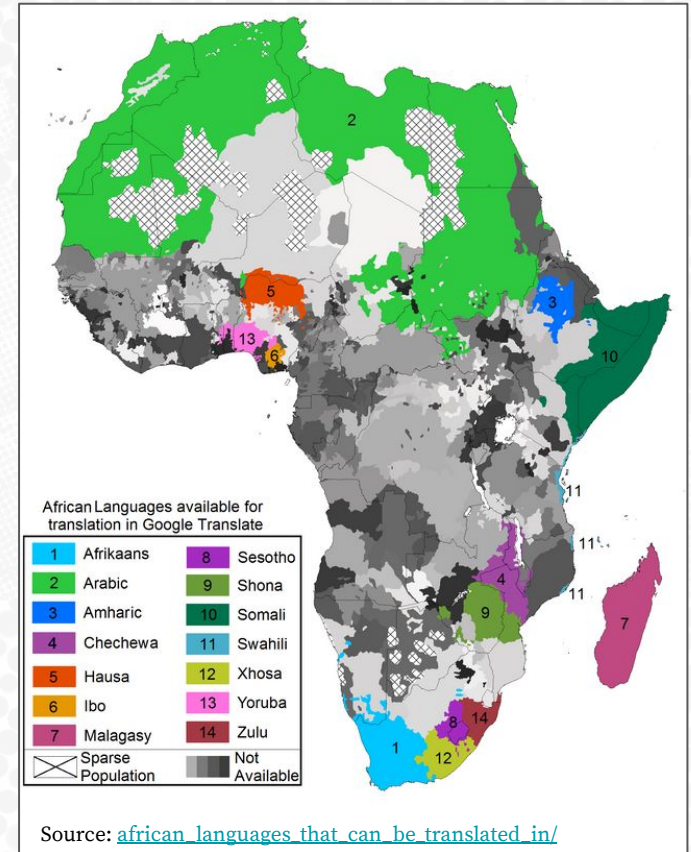
# Machine Translation with human in the loop

**Fails** in representing all languages

**301** on Wikipedia

**103** on Google Translate

**2000-6000** in the world





# Machine Translation with human in the loop

## Quality? For small languages

- No enough parallel data
- **Poor Quality** makes Translation is not useful to editors



**THE VERGE** TECH ▾ REVIEWS ▾ SCIENCE ▾ CREATORS ▾ ENTERTAINMENT ▾ VIDEO FEATURES MORE ▾

POLICY \ REPORT \ US & WORLD \

## Wikipedia has a Google Translate problem

*Smaller editions badly need a machine translation tool — but it isn't good enough to use on its own*

By [Kyle Wilson](#) | May 8, 2019, 10:15am EDT

Source: <https://www.theverge.com/2019/5/8/18526739/wikipedia-translation-tool-machine-learning-ai-english>

# Machine Translation with human in the loop

## Translation of Labels?

**Goodwill Zwelithini**



**King of the Zulus**

**Reign** 17 September 1968 – present

**Coronation** 3 December 1971

**Predecessor** Cyprian Bhekuzulu kaSolomon

---

**Born** July 14, 1948 (age 71)  
Nongoma, Union of South Africa

**Wives** Mantombi Dlamini (*great wife*)  
Sibongile Winifred Dlamini  
Buthele MaMathe  
Thandekile Jane Ndlovu  
Nompumelelo Mchiza  
Zola Zekusiwwe MaFu

**Issue** 27 including: [\[show\]](#)

**Full name**  
Goodwill Zwelithini kaBhekuzulu

**House** House of Zulu

**Father** Cyprian Bhekuzulu kaSolomon

Text Documents

DETECT LANGUAGE XHOSA ENGLISH ARABIC FRENCH SPANISH ARABIC

Goodwill Zwelithini kaBhekuzulu Bonne volonté Zwelithini kaBhekuzulu

DETECT LANGUAGE XHOSA ENGLISH ARABIC FRENCH SPANISH ARABIC

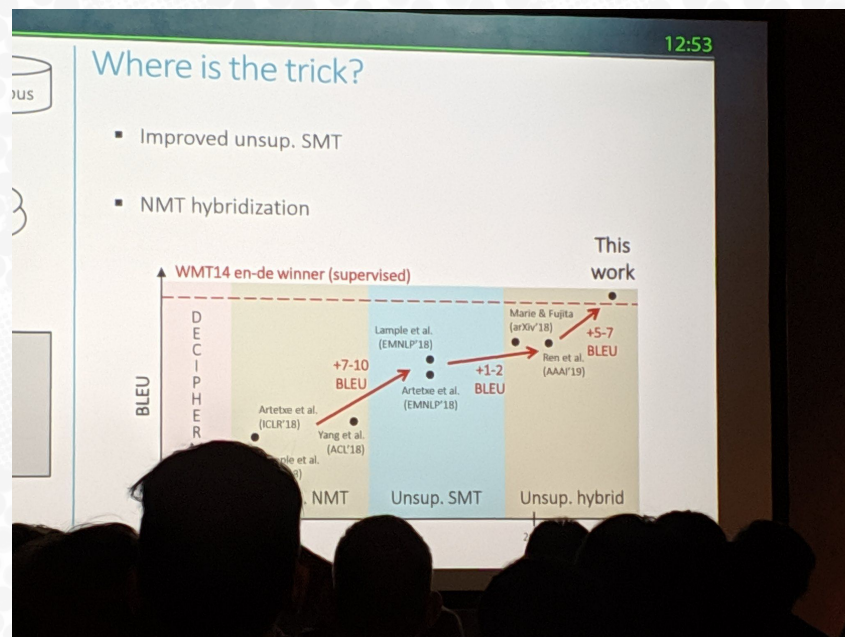
Goodwill Zwelithini kaBhekuzulu Zwelithini kaBhekuzulu النوايا الحسنة

al nawaya alhasanat Zwelithini kaBhekuzulu

# Unsupervised Machine Translation

A **solution** to Quality and Coverage

- No need for parallel corpora
- Very suitable for small languages
- Getting better rapidly



(Artetxe et al. 2019) An Effective Approach to Unsupervised Machine Translation. ACL2019

**Solutions**

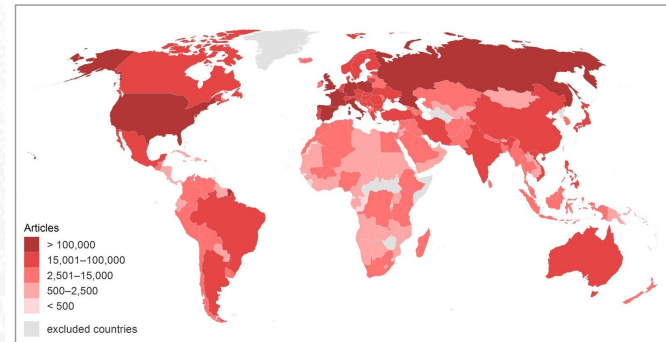
**Label Translation**

# Translation from English

**will not solve  
Information Poverty**

Large language wikipeidias  
get to decide what is notable.

Small Language Wikipeidias  
will enforce that bias through translations.



Number of geotagged Wikipedia articles per country (Graham et. al. 2014)

# Translation of English References

## Harms The Neutral Point of View

**English:** Alaa Salah is a Sudanese student and anti-government protester. The image of Salah has been dubbed as "Woman in White" or "Lady Liberty" of Sudan.

- [CNN](#) (en)
- [2x Guardian](#) (en)
- [Le Monde](#) (fr)
- [BuzzFeed](#) (en)
- [ABC](#) (en)
- [smh](#) (en)
- [2x NYT](#) (en)
- [Nzherald](#) (en)
- [France24](#)(ar)

## آلاء صلاح [edit|edit source]

Wikidata: آلاء صلاح - Q63126345: "ناشطة سودانية" - Aliases: None



آلاء صلاح في الصورة التي أصبحت مشهورة كرمز للاحتجاجات السودانية 2018-19.

آلاء صلاح هي طالبة سودانية ومنظاهرة مناضفة للحكومة. اكتسبت شهرتها من صورة التقطتها (لانا هارون)، وانتشرت في شهر أبريل 2019. وأصبحت آلاء ابقونة ورمز للتمرد في الثورة السودانية والتي تعرف باسم "كنداكة" أو "الملكة النوبية" أو "تمثال الحرية السودانية".<sup>[1][2]9]</sup>

### Contents [hide]

- 1 الحياة المبكرة والتعليم
- 2 احتجاجات السودان
- 3 الصورة
- 4 انظر أيضا
- 5 المراجع

### الحياة المبكرة والتعليم [edit|edit source]

ولدت آلاء صلاح في عام 1996 أو 1997. وكانت والدتها تعمل مصممة أزياء ووالدها يعمل في البناء والتشييد.<sup>[4]</sup> تدرس الهندسة والعمارة في جامعة السودان العالمية في الخرطوم.<sup>[5]</sup>

### المراجع [edit|edit source]

- ↑ CNN, Gianluca Mezzofiore. "This woman has come to symbolize Sudan's protests". *CNN*. 2019. اطلع عليه بتاريخ 11 أبريل 2019.
- ↑ *The Sydney Morning Herald* (باللغة الإنجليزية). 2019-04-11. مؤرشف من الأصل في 29 أبريل 2019. اطلع عليه بتاريخ 11 أبريل 2019.
- ↑ جنوبي كنداكة.. امرأة سودانية تحول لافونة "الثورة" - فرانس 24. 2019-04-11 / 24. مؤرشف من الأصل في 12 أبريل 2019. اطلع عليه بتاريخ 19 أبريل 2019.
- ↑ "Sudan's singing protester speaks out: 'I was raised to love our home'". *The Guardian*. ISSN 0261-0261. Salihi, Zeinab Mohammed (2019-04-10). "I was raised to love our home": Sudan's singing protester speaks out. مؤرشف من الأصل في 9 مايو 2019. اطلع عليه بتاريخ 11 أبريل 2019.
- ↑ "This Woman Stood On Top Of A Car And Became An Icon Of Sudan's Historic Protests". *BuzzFeed News*. مؤرشف من الأصل في 29 مايو 2019.
- ↑ "Sudanese police fire on protesters demanding president step down". *The Guardian*. 27 أبريل 2019. اطلع عليه بتاريخ 11 أبريل 2019.
- ↑ EXCLUSIVE: Sudanese spy chief 'met head of Mossad to discuss Bashir succession plan'. *Middle East Eye*. 7 مايو 2019. مؤرشف من الأصل في 19 مايو 2019.
- ↑ "Le mouvement de protestation s'embrase au Soudan". *The New York Times*. 08-04-2019 (باللغة الفرنسية). مؤرشف من الأصل في 29 مايو 2019. اطلع عليه بتاريخ 11 أبريل 2019.
- ↑ "Sudan's Military to Make Announcement Amid Protests Against Omar Hassan al-Bashir". *The New York Times*. ISSN 0362-4331. Mullany, Gerry (2019-04-11). "Sudan's Military to Make Announcement Amid Protests Against Omar Hassan al-Bashir". مؤرشف من الأصل في 11 أبريل 2019.
- ↑ "Mood in Sudan shifts to anger as the army prepares to seize power". *The Guardian*. 10-04-2019 (باللغة الإنجليزية). Maclean, Ruth (2019-04-11). "Mood in Sudan shifts to anger as the army prepares to seize power". مؤرشف من الأصل في 11 أبريل 2019.
- ↑ Poetic photo of Sudan's 'Lady Liberty' sheds light on anti-government protests". *ABC News*. 11 أبريل 2019. اطلع عليه بتاريخ 11 أبريل 2019.
- ↑ "It's Going to Be the Image of the Revolution". *The New York Times*. ISSN 0362-4331. Friedman, Vanessa (2019-04-10). "It's Going to Be the Image of the Revolution". مؤرشف من الأصل في 28 مايو 2019. اطلع عليه بتاريخ 11 أبريل 2019.
- ↑ "The woman in white: Why a photo from the Sudan protests has gone viral". *ABC News*. 13 مايو 2019. اطلع عليه بتاريخ 10 أبريل 2019.

Source: [https://en.wikipedia.org/wiki/Alaa\\_Salah](https://en.wikipedia.org/wiki/Alaa_Salah) (en, ar)

## Current

- Automatic Machine Translation

## Alternatives

- **Natural Language Generation**
- Summarization
- Cross Lingual Transfer Learning

**Research  
For  
Solutions**

# Natural Language Generation

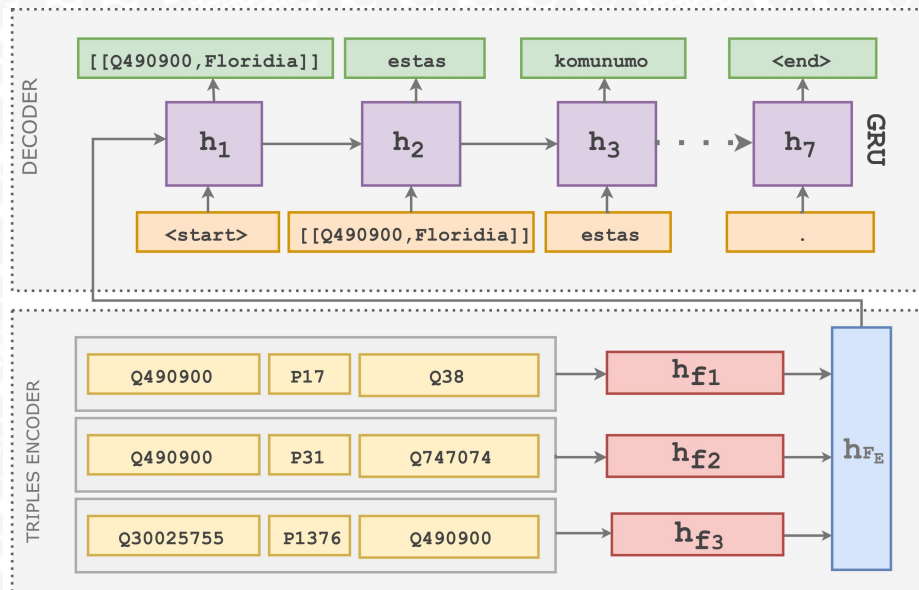
Generation of Wikipedia Articles from **Wikidata**



WIKIPEDIA  
The Free Encyclopedia

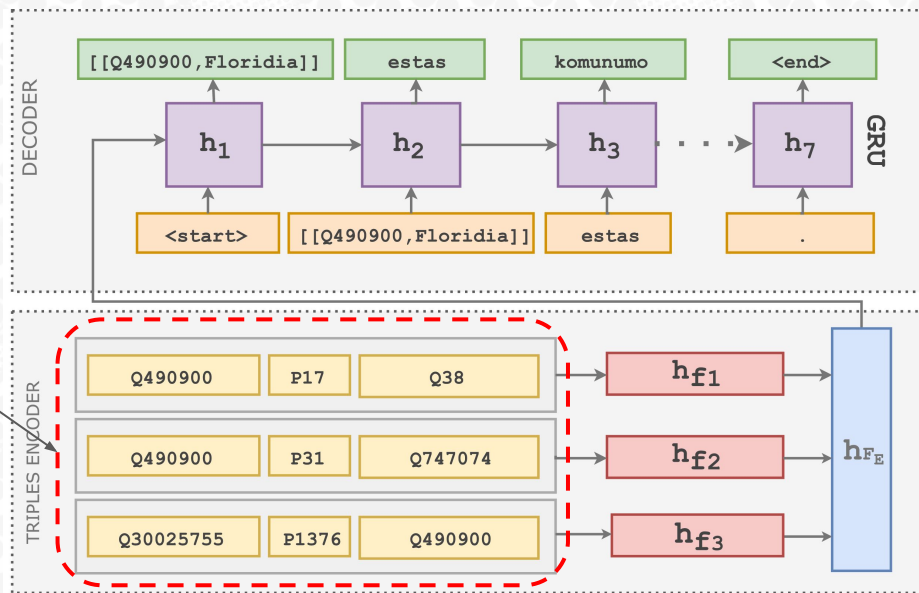


# Using same type of Encoder-Decoder Models of MT

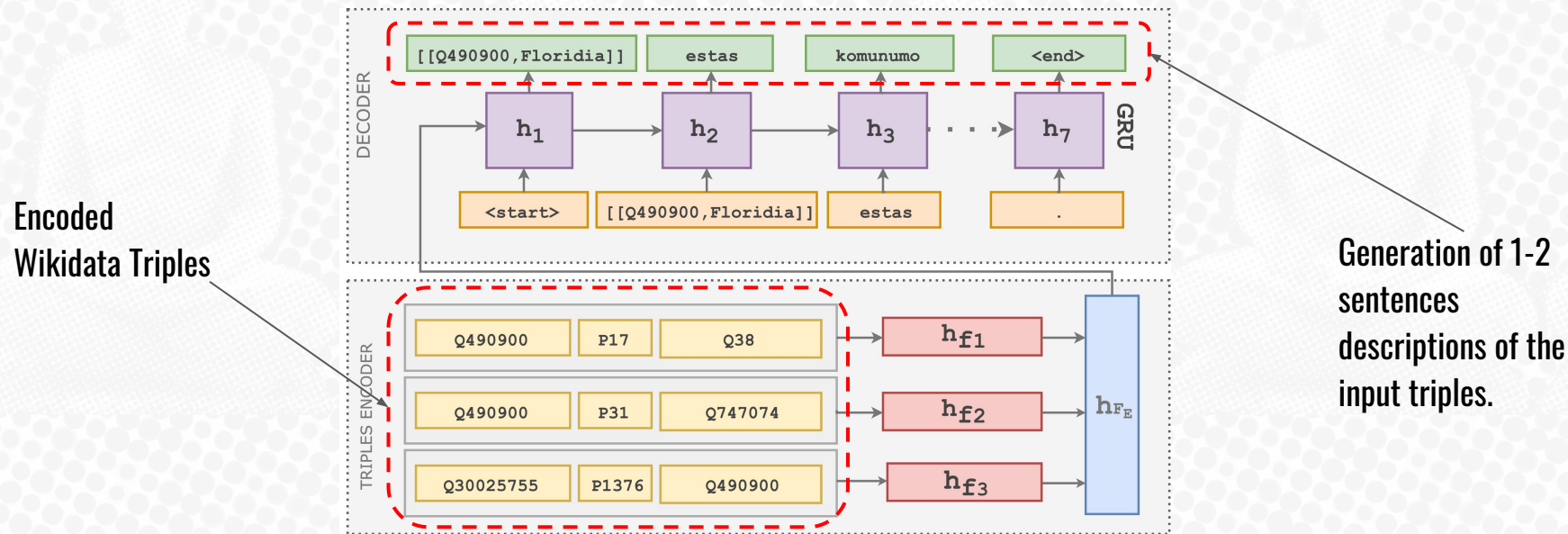


# Using same type of Encoder-Decoder Models of MT

Encoded Wikidata Triples



# Using same type of Encoder-Decoder Models of MT



# Annotated Data For Free

Automatic alignment between knowledge base triples and Free text

Nigragorĝa pigogarolo (Q1586267)

|                      |              |
|----------------------|--------------|
| instance of          | taksono      |
| taksonomia nomo      | Calocitta    |
| supera taksono       | Pigogarolo   |
| original combination | Pica colliei |



WIKIPEDIA  
The Free Encyclopedia

La **Nigragorĝa pigogarolo** (**Calocitta colliei**) estas rimarkinda longvosta pigogarolo de la familio de Korvedoj kaj ordo de Paseroformaj kiu loĝas en nordokcidenta Meksiko.

*The black-throated magpie-jay (Calocitta colliei) is a strikingly long-tailed magpie-jay of northwestern Mexico.*

## Good

- Integrating Wikidata knowledge
- Better than MT for small languages
- Translation of Names (using Wikidata)

## Bad

- Generation of More than 1 Paragraph
- factual mistakes

# Natural Language Generation

Generation of Wikipedia Articles from **Wikidata**



WIKIDATA



WIKIPEDIA  
The Free Encyclopedia

## Current

- Automatic Machine Translation

## Alternatives

- Natural Language Generation
- Summarization
- Cross-Lingual Transfer Learning



## Building sentence representations using:

- Frequency
- Length
- Similarity
- Location
- Centrality

## Optimizing:

- Min Redundancy
- Max Relevance
- Max Informativeness

A Simple Theoretical Model of Importance for Summarization  
(Peyrard ACL2019)

# Multi-document Unsupervised Extractive Summarization

## Current

- Automatic Machine Translation

## Alternatives

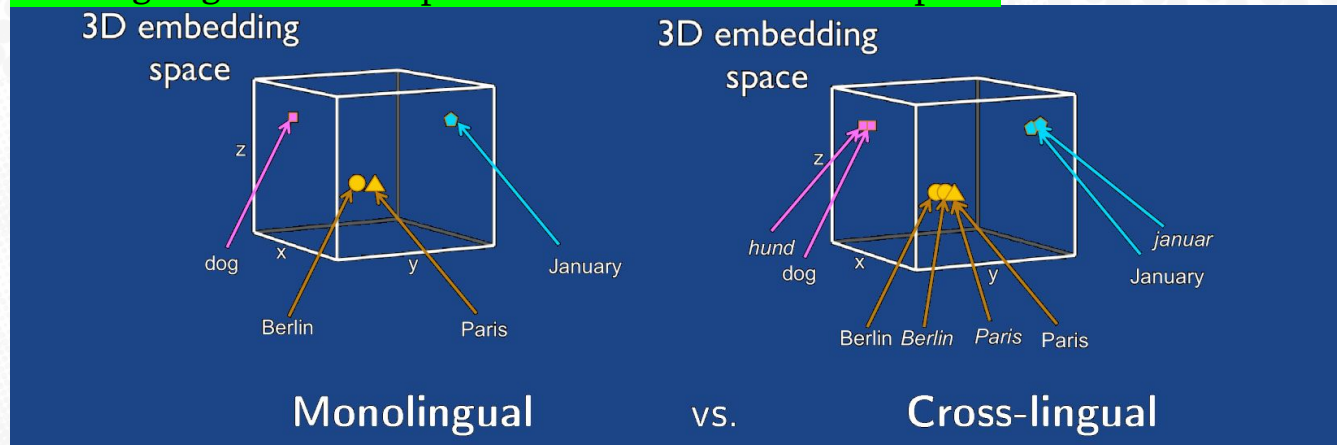
- Natural Language Generation
- Summarization
- Cross-Lingual Transfer Learning





# Unsupervised Cross-Lingual Word Embeddings

Building aligned word representations in the feature space.



**“Unsupervised” = Requires very minimal datasets**

- Small bilingual dictionaries
- Weak supervision (shared vocabulary)

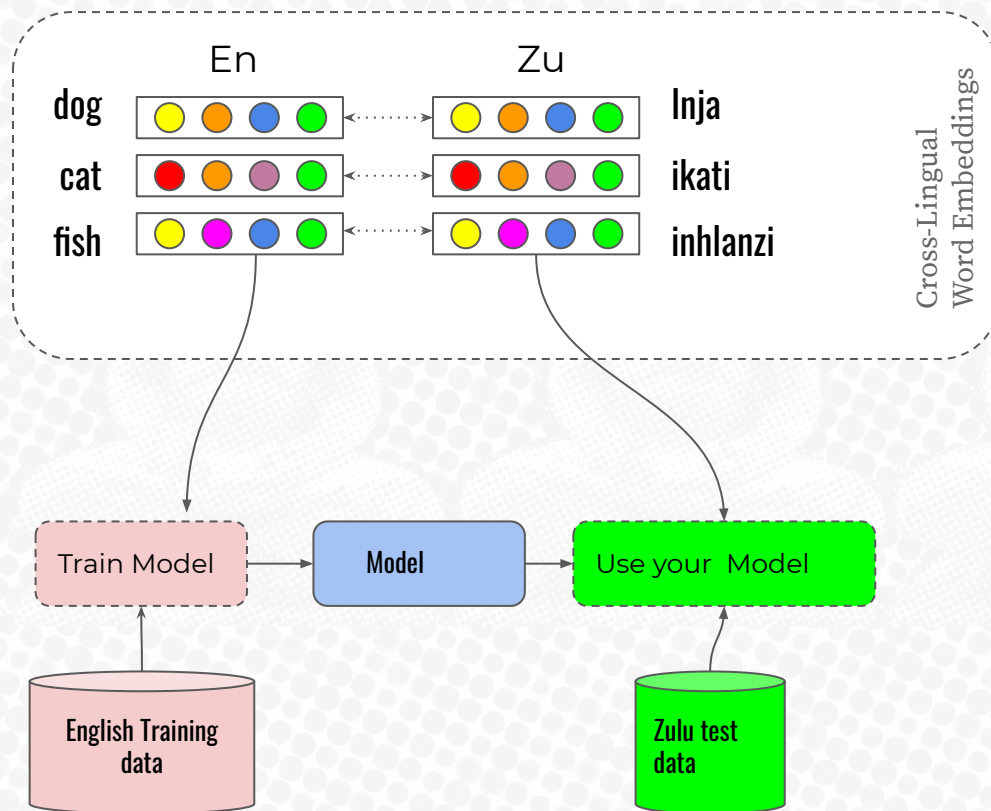
# Unsupervised Cross-Lingual Word Embeddings

**Useful for**

“Cross lingual Transfer Learning”

Acquiring good results for tasks on languages with limited/no annotated datasets (e.g. Zulu)

Using large annotated datasets available in other languages (e.g. English)



# Unsupervised Cross-Lingual Word Embeddings

## Applications: Reference Retrieval

Retrieval using Entities (Entity Linking)

### Cross-lingual Wikification Using Multilingual Embeddings

**Chen-Tse Tsai** and **Dan Roth**

University of Illinois at Urbana-Champaign  
201 N. Goodwin, Urbana, Illinois, 61801  
{ctsai12, danr}@illinois.edu

# Unsupervised Cross-Lingual Word Embeddings

## Applications: Reference Retrieval

### Retrieval using Facts (OpenIE)

#### Multilingual Open Relation Extraction Using Cross-lingual Projection

**Manaal Faruqi**  
Carnegie Mellon University  
Pittsburgh, PA 15213  
mfaruqi@cs.cmu.edu

**Shankar Kumar**  
Google Inc.  
New York, NY 10011  
shankarkumar@google.com

#### MT/IE: Cross-lingual Open Information Extraction with Neural Sequence-to-Sequence Models

**Sheng Zhang** and **Kevin Duh** and **Benjamin Van Durme**  
Johns Hopkins University  
{zsheng2, kevinduh, vandurme}@cs.jhu.edu

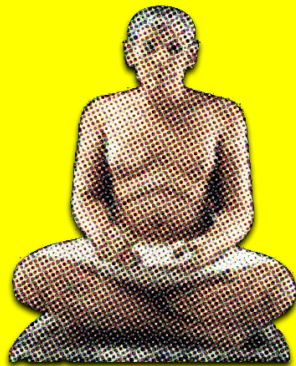
- Multilingual Open Relation Extraction Using Cross-lingual Projection (Faruqi et al. NAACL 2015)
- MT/IE: Cross-lingual Open Information Extraction with Neural Sequence-to-Sequence Models (Zhang et al. EACL2017)

# Research solutions

vs

# problems

|                                | Language Gap | Information Poverty | Reference Bias | Quality | Label Translation | Long articles Generation |
|--------------------------------|--------------|---------------------|----------------|---------|-------------------|--------------------------|
| Machine Translation            | ✓            | □                   | □              | □       | (□)               | ✓                        |
| NLG (data2text)                | ✓            | ✓                   | □              | ✓       | (✓)               | □                        |
| Summarization                  | ✓            | ✓                   | ✓              | ✓       | (✓)               | ✓                        |
| Crosslingual Transfer Learning | —            | —                   | ✓              | —       | —                 | —                        |



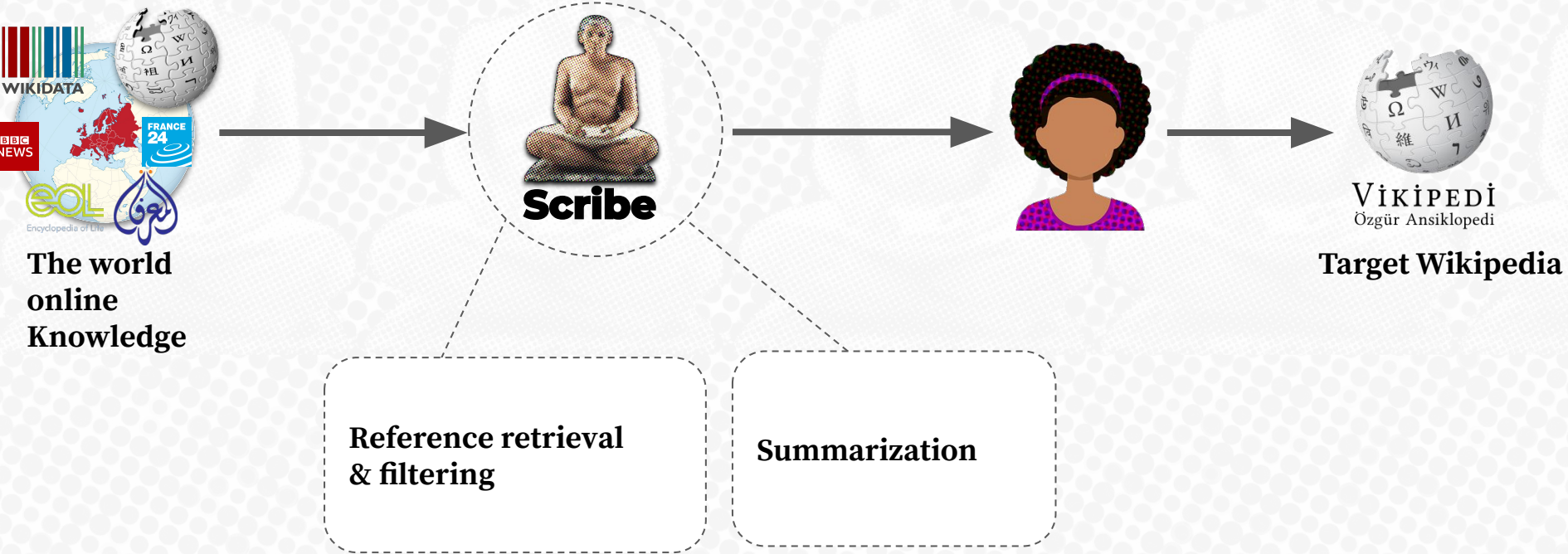
# Scribe

**Wikimedia Funded Project (2019/2020)**

Helping editors of under-resourced languages to  
create new high-quality Wikipedia articles

<https://meta.wikimedia.org/wiki/Scribe>

# Unsupervised Summarization human in the loop.




مقالة نقاش أنشئ أنشئ المصدر ☆ ابحث في ويكيبيديا

نشر الصفحة... نشر الصفحة... إدرج Ω استشهد

# فلافل

سكرايب هي أداة تساعدك على كتابة المقالات بشكل أكثر فعالية



- 1 يشتهر الفلافل في العالم العربي
- 2 مصنوعة من الفول المطحون أو من الحمص
- 3 يعتقد بأن أول من حضر الفلافل هم المصريون
- 4 بدأت الفلافل بالانتشار في الغرب
- 5 تعد الوجبة الأمثل للنباتيين

## المعلومات الغذائية

1 5 4 السعرات الحرارية 288, الدهون 7, الدهون المشبعة 1

Collected references from online resources

Collected images from Wikidata

Suggested Section planning

Key points from external references for each section from online resources

صفحة الرئيسية  
تدات الجارية  
ت التغييرات  
ت التغييرات الأساسية  
فج  
واضع  
دي  
بات  
له عشوائية  
فج بدون إنترنت  
ماركة  
ل بالموسوعة  
باعدة  
بدان

سياسة الخصوصية حول ويكيبيديا إخلاء مسؤولية المطورون بيان تعريف الارتباطات



# THANK YOU!



## Interested in Scribe?

- **Know More:**  
<https://meta.wikimedia.org/wiki/Scribe>
- **You also care about filling the language Gap?**  
**Collaborate with us:**



**Hady Elshar**  
wd:Q50290890  
[[User:Hadyelsahar]]  
@hadyelsahar  
hadyelsahar@gmail.com



**Lucie-Aimée Kaffee**  
Q37860261  
[[User:Frimelle]]  
@frimelle  
lucie.kaffee@gmail.com