



IndicWiki Internship 2025

On

Indigenous Language Development

Submitted by

Bhagwan Ji Jha

Under the guidance of

Prof. Krupal Kashyap

Program Manager – IndicWiki, IIIT Hyderabad

14 July 2025

Acknowledgment

First and foremost, I would like to express my sincere thanks to my project investigator(Dr. Radhika Mamidi), program manager(Dr. Krupal Kasyap), Assistant Professor at the Language Technologies Research Centre(LTRC)- Dr. Parameswari Krishnamurthy, All mentors, teammates, and everyone who supported me on this journey, who gave me this valuable opportunity to work in such a learning environment through this Internship.

I would like to extend my heartfelt gratitude to my project mentors, Dr. Krupal Kasyap, Dr. Neechal Karan, Mr. Somyadip Ghosh, Mr. Nagaraju Vappala, and the entire IndicWiki team. Your guidance, constant supervision, and the wealth of knowledge and expertise you have imparted have been invaluable throughout this journey. Thank you for your unwavering support and encouragement.

Finally, I am deeply thankful to my parents and teachers who helped and inspired me throughout this Internship.

Sincerely

Bhagwan Ji Jha,

Government Engineering College Ajmer

Congratulations! You're Shortlisted for IndicWiki Indigenous Language Development Internship 2025 [Inbox](#)

◆ Summarise this email

P

PM Indicwiki <pm.indicwiki@iiit.ac.in>
to me

Sat 10 May, 12:30

☆

😊

↶

⋮

Hi Jijha Bhagwan,

We are pleased to inform you that you have been shortlisted for the **Indigenous Language Development Intern** role under the **IndicWiki Summer Internship 2025**. **Congratulations** on this achievement! Ref: <https://indicwiki.iiit.ac.in/internship2025/>

About the Internship

- **Duration:** May 12 to June 30, 2025
- **Mode:** Remote (Only Online)
- **Focus:** Digital preservation and promotion of Telangana's rich indigenous languages

Congratulations on Completing the IndicWiki Summer Internship 2025 Program!

◆ Summarise this email

P

PM Indicwiki
to me

Sun 20 Jul, 14:39

☆

😊

↶

⋮

Hi Bhagwan Ji Jha,

Congratulations on successfully completing the IndicWiki Indigenous Language Development Internship Program held from 12 May to 14 July 2025.

We truly appreciate your dedication and contributions throughout the internship. Your efforts have played a valuable role in promoting and preserving indigenous languages through digital platforms.

Please find your **certificate of completion** attached to this email.

We invite you to continue your journey with Wikipedia, the Wikimedia sister projects, and the broader Wikimedia community. Your expertise and dedication are truly valued, and we're excited to see how your future contributions will shape and enrich the Wiki ecosystem.

Remember, the work you do today helps build a more inclusive, accessible, and informative world for millions of users across the globe.

Wishing you all the best in your future endeavors—and we hope you'll stay actively involved in the vibrant Wikimedia community!

With sincere appreciation,

Cordially

IndicWiki Project Team

International Institute of Information Technology

Professor CR Rao Rd, Gachibowli,

Hyderabad. Telangana 500032

Internship Completion Certificate

INTERNATIONAL INSTITUTE OF
INFORMATION TECHNOLOGY

HYDERABAD

CERTIFICATE OF INTERNSHIP

This certificate is presented to

Bhagwan Ji Jha

for completing IndicWiki Indigenous Language Development

Internship program from 12 May-14 July 2025.


Krupal Kasyp

PROGRAM MANAGER


Dr. Radhika Mamidi

PROJECT INVESTIGATOR

Contents

ACKNOWLEDGMENT	1
CONTENTS	4
PROJECT INTRODUCTION	5
GETTING STARTED WITH WIKIPEDIA	9
WIKIMEDIA COMMONS	14
WIKIDATA	15
WIKTIONARY AND WIKISOURCE	16
INDICWIKI DATA COLLECTION & ANNOTATION PROJECT, LINGUA LIBRE	19
WIKIMEDIA TECHNICAL SPACE	23
REFERENCES AND OUR DETAILS (PROJECTS)	24

Chapter 1

Introduction

The IndicWiki Project, spearheaded by IIIT-Hyderabad, is a transformative initiative dedicated to enriching Wiki content in Indian languages, with a focus on Telugu and other regional languages. By creating comprehensive, freely accessible knowledge repositories, the project bridges the digital divide, promotes linguistic diversity, and empowers non-English-speaking communities to contribute to the global knowledge ecosystem. Through innovative tools, community engagement, and cultural preservation efforts, IndicWiki is redefining digital knowledge dissemination in India.

India's linguistic diversity, with over 1,600 languages and dialects, stands as a vibrant testament to its rich cultural heritage. This diversity is not merely a statistical marvel but a living, breathing embodiment of the country's history, traditions, and identity. However, the digital knowledge landscape has historically been dominated by English, creating a profound disparity in access to information. Millions of non-English speakers, particularly those in rural and semi-urban areas, have been left on the periphery of the digital revolution, struggling to find reliable, relevant, and culturally resonant content in their native languages.

This linguistic gap has far-reaching implications, from limiting educational opportunities to hindering socio-economic development and eroding cultural pride. Recognizing this critical issue, the IndicWiki Project has emerged as a transformative initiative, dedicated to bridging this divide and fostering inclusivity in the digital realm. By creating a collaborative, open-access platform,

IndicWiki enables the generation and dissemination of encyclopaedic content in Indian languages, leveraging cutting-edge language technologies and grassroots community participation.

At its core, the IndicWiki Project is a movement to democratize knowledge. It empowers native speakers to become active contributors and consumers of content in their mother tongues, ensuring that their voices are heard and their cultural narratives preserved. This participatory approach not only revitalizes Indian languages but also strengthens their relevance in the digital age. By building robust online encyclopaedias, IndicWiki is creating a repository of knowledge that reflects the depth and diversity of India's linguistic and cultural tapestry.

Moreover, the project plays a pivotal role in preserving traditional knowledge systems, which are often at risk of being overshadowed by globalized narratives. By

documenting indigenous practices, folklore, and historical insights in local languages, IndicWiki ensures that this invaluable heritage is safeguarded for future generations. This effort is particularly crucial in an era where globalization threatens to homogenize cultural expressions.

The impact of the IndicWiki Project extends beyond language preservation; it is a catalyst for social and economic empowerment. Access to information in one's native language enhances educational outcomes, fosters innovation, and promotes digital literacy among underserved communities. It also encourages a sense of pride and ownership among speakers of lesser-known languages, many of which are endangered due to neglect and lack of representation.

In essence, the IndicWiki Project is not just about creating content; it is about reclaiming identity, fostering equity, and ensuring that India's linguistic heritage thrives in the digital age. By addressing the long-standing disparity in the digital knowledge landscape, IndicWiki is paving the way for a more inclusive, diverse, and culturally vibrant future. It is a testament to the power of collaboration, technology, and community-driven efforts in shaping a knowledge ecosystem that truly belongs to all.

1.1 Core Objectives

The IndicWiki Project is driven by four primary goals:

- Language Preservation and Promotion: Develop comprehensive knowledge repositories in Indian languages to document cultural heritage, literature, history, and traditional knowledge systems.
- Accessibility and Knowledge Sharing: Make high-quality information accessible to non-English speakers, reducing language barriers and enabling knowledge exchange in native languages.
- Community Collaboration: Foster vibrant communities of contributors, including volunteers, language enthusiasts, and subject matter experts, to drive collaborative content creation.
- Free and Open Knowledge: Uphold the principles of open access, ensuring all content is freely available and encouraging a culture of knowledge sharing.

1.2 Strategic Approach

IndicWiki employs a multi-faceted strategy to achieve its objectives, encompassing content creation, community engagement, technological innovation, and knowledge preservation.

- Content Creation Framework IndicWiki adopts a systematic approach to develop high-quality, culturally relevant content:

- Thematic Approach: Prioritizes topics of cultural significance, such as literature, history, science, and local knowledge.
- Expert Validation: Engages subject matter experts to review and verify content for accuracy.
- Quality Standards: Implements guidelines to ensure consistency and reliability across languages.
- Translation-Plus Model: Enhances translations with cultural contextualization to resonate with local audiences.

A notable achievement is the creation of 10 lakh articles across 45 subjects, published through the Wiki Sandbox, an isolated testing environment that allows users to experiment with content rendering.

- Community Engagement Initiatives

Building sustainable contributor communities is central to IndicWiki's success:

- Workshops and Editathons: Regular training events introduce new contributors to Wikipedia editing.
- Mentorship Programs: Pair experienced editors with newcomers to facilitate skill development.
- Recognition Systems: Celebrate contributors' efforts to sustain motivation.
- WikiClubs: Establish student-led communities in colleges to promote collaborative content creation.

High-profile events, such as the Wikimedia Technology Summit (2021, 2023, 2024), WikiConference India, and Wiki Women Camp, alongside hackathons and Code Labs, have bolstered community participation. Internship programs further enhance technical and editorial skills among participants.

- Technological Innovations

IndicWiki develops specialized tools to overcome challenges in Indian language computing:

- Transliteration Tools: Enable seamless script conversion between languages.

- Content Suggestion Systems: Use AI to identify priority articles for creation or improvement.
- Quality Assessment Tools: Automate evaluation of article completeness and quality.
- Specialized Editors: Provide user-friendly interfaces optimized for Indian languages.

The Micro Content Development interface simplifies the process of updating existing content, while advancements in structured data, information extraction, language translation systems, and natural language generation (NLG) enhance content creation efficiency.

- Knowledge Preservation Efforts Beyond encyclopedic content, IndicWiki digitizes and preserves traditional knowledge:
- Oral History Documentation: Records and transcribes oral narratives and traditions.
- Indigenous Knowledge Systems: Documents traditional sciences, arts, and practices to ensure their longevity.

1.3 Impact and Outcomes

While The IndicWiki Project has delivered significant outcomes:

- Content Creation: Produced high-quality, reliable Wiki articles in Indian languages, with 10 lakh pages published across diverse subjects.
- Digital Inclusion: Expanded access to knowledge for non-English speakers, particularly in rural and semi-urban areas, fostering greater participation in the digital ecosystem.
- Preservation of Culture: Documented India's rich cultural, historical, and linguistic heritage, safeguarding traditional knowledge for future generations.
- Capacity Building: Trained volunteers, students, and researchers in Wikipedia editing, empowering native speakers to sustain Indic-language content growth.

2.3 Create Wikipedia Account

The screenshot shows a user profile page with the following sections:

- Suggested edits:** 10 topics, 3 of 120 suggestions. A preview box shows the article "Individual psychology".
- Your Impact:** 142 Total edits, 0 Thanks received. Last edited 5 months ago, longest streak 4 days.
- Your recent activity (last 60 days):** 0 Edits, Mar 15 to May 13.
- Most viewed (since your edit):** Assiniboine Park Zoo (1,324 views), Data philanthropy (738 views), Currie Media.
- Appearance:** General settings for text size (Small, Standard, Large), width (Standard, Wide), color (beta), and font (Automatic).

- Easy / Beginner Friendly
- Spellings
- Grammar
- Adding References

2.4 Wikipedia Page Structure

The screenshot shows the Wikipedia page for "Earth" with the following structure:

- Title:** The title "Earth" is highlighted with a blue arrow.
- Body:** The main content of the page is highlighted with a large blue arrow.
- InfoBox:** An info box on the right side contains the following information:
 - Image:** A photograph of Earth from space, labeled "The Blue Marble, Apollo 17, December 1972".
 - Designations:** The world, the globe, Sol III, Terra, Tellus, Gaia, Mother Earth.
 - Alternative names:** Earthly, terrestrial, terran, tellurian.
 - Symbol:** ♦ and ♀.

WIKIPEDIA The Free Encyclopedia

Article Talk **Discussions**

Coronavirus

From Wikipedia, the free encyclopedia

For the ongoing outbreak, see [2019–20 Wuhan coronavirus outbreak](#). For the specific virus causing the outbreak, see [Novel coronavirus \(2019-nCoV\)](#).

Coronaviruses are a group of viruses that cause diseases in mammals and birds. In humans, they can cause common colds, SARS, MERS, and COVID-19. In cattle, they can cause respiratory infections which are typically mild including the common cold but rarer forms like SARS and MERS can be lethal. In cows and pigs they may cause diarrhea, while in chickens they can cause an upper respiratory disease. There are no vaccines or antiviral drugs that are approved for prevention or treatment.

Hyper Links

Contents [hide]

- 1 Discovery
- 2 Name and morphology
- 3 Replication
- 4 Taxonomy
- 5 Evolution
- 6 Human coronaviruses
 - 6.1 Severe acute respiratory syndrome (SARS)
 - 6.2 Middle East respiratory syndrome
 - 6.3 Novel coronavirus (2019-nCoV)
- 7 Other animals
- 7.1 Diseases caused

History

Citations/References

Orthocoronavirinae

Electron micrograph of infectious bronchitis virus virions

Render of 2019 nCoV virion

References

References

1. ^{a b} Simon, J. L.; et al. (February 1994). "Numerical expressions for precession formulae and mean elements for the Moon and planets". *Astronomy and Astrophysics*. **282** (2): 663–683. Bibcode:1994A&A...282..663S.
2. ^{a b c d e} Staff (13 March 2021). "Useful Constants". International Earth Rotation and Reference Systems Service. Archived from the original on 29 October 2012. Retrieved 8 June 2022.
141. ^a Pollack, Henry N.; Hurter, Suzanne J.; Johnson, Jeffrey R. (August 1993). "Heat flow from the Earth's interior: Analysis of the global data set". *Reviews of Geophysics*. **31** (3): 267–280. Bibcode:1993RvGeo..31..267P. doi:10.1029/93RG01249.
142. ^a Richards, M. A.; Duncan, R. A.; Courtillot, V. E. (1989). "Flood Basalts and Hot-Spot Tracks: Plume Heads and Tails". *Science*. **246** (4926): 103–107. Bibcode:1989Sci...246..103R. doi:10.1126/science.246.4926.103. PMID 17837760.

V · T · E

Earth [show]

V · T · E **Structure of Earth** [show]

V · T · E **Earth science** [show]

V · T · E **Solar System** [show]

V · T · E **Geology** [show]

Other articles related to Earth [show]

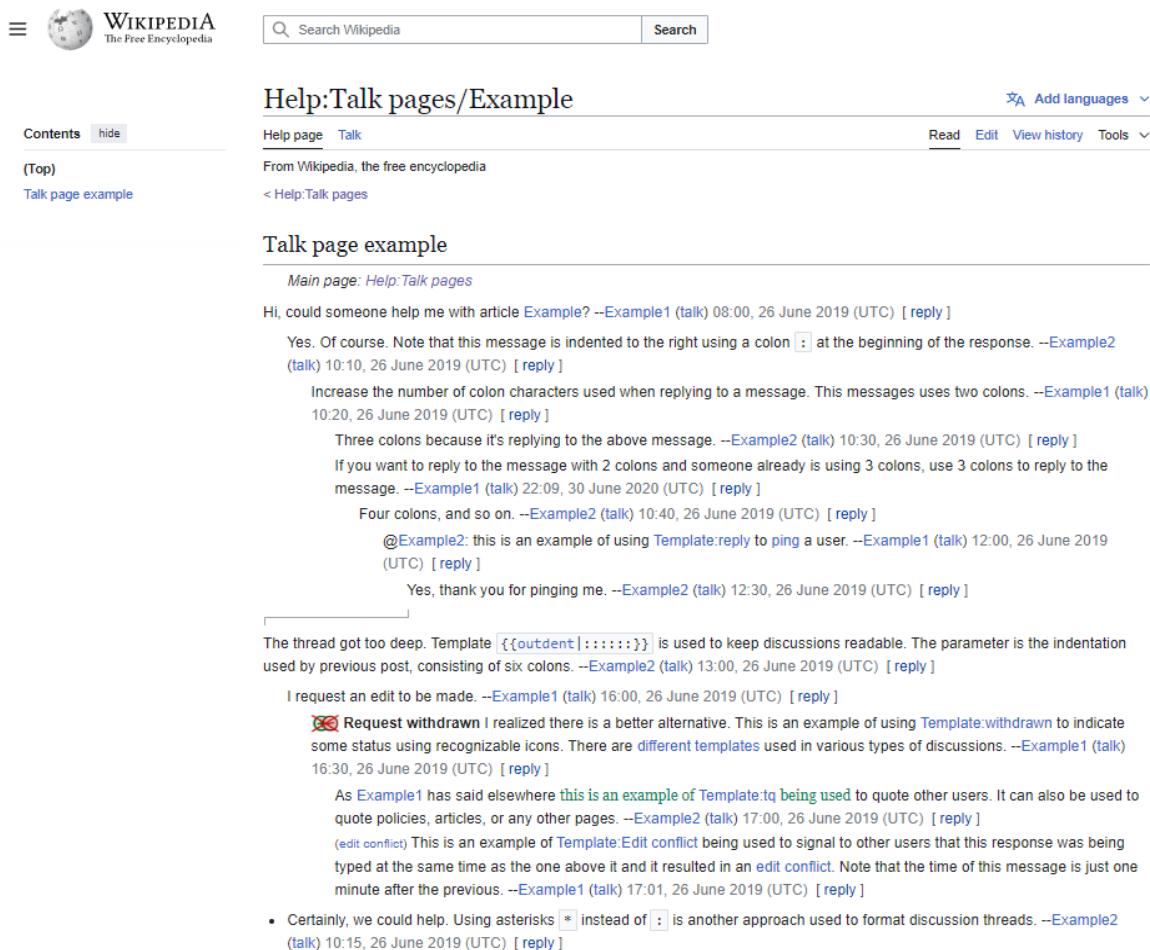
Portals: [Biology](#) [Earth sciences](#) [Ecology](#) [Geography](#) [Volcanoes](#) [Solar system](#) [Outer space](#)
[Weather](#) [World](#)

Earth at Wikipedia's sister projects: [Definitions](#) from Wiktionary [Media](#) from Commons [News](#) from Wikinews [Quotations](#) from Wikiquote [Texts](#) from Wikisource [Textbooks](#) from Wikibooks [Resources](#) from Wikiversity

Categories

Categories: [Solar System](#) | [Earth](#) | [Astronomical objects known since antiquity](#) | [Global natural environment](#) | [Nature](#)
[Planets of the Solar System](#) | [Terrestrial planets](#)

2.5 Talk Page



The screenshot shows a Wikipedia talk page titled "Help:Talk pages/Example". The page content is as follows:

Hi, could someone help me with article Example? --Example1 (talk) 08:00, 26 June 2019 (UTC) [reply]

Yes. Of course. Note that this message is indented to the right using a colon : at the beginning of the response. --Example2 (talk) 10:10, 26 June 2019 (UTC) [reply]

Increase the number of colon characters used when replying to a message. This messages uses two colons. --Example1 (talk) 10:20, 26 June 2019 (UTC) [reply]

Three colons because it's replying to the above message. --Example2 (talk) 10:30, 26 June 2019 (UTC) [reply]

If you want to reply to the message with 2 colons and someone already is using 3 colons, use 3 colons to reply to the message. --Example1 (talk) 22:09, 30 June 2020 (UTC) [reply]

Four colons, and so on. --Example2 (talk) 10:40, 26 June 2019 (UTC) [reply]

@Example2: this is an example of using Template:reply to ping a user. --Example1 (talk) 12:00, 26 June 2019 (UTC) [reply]

Yes, thank you for ping me. --Example2 (talk) 12:30, 26 June 2019 (UTC) [reply]

The thread got too deep. Template {{outdent}} is used to keep discussions readable. The parameter is the indentation used by previous post, consisting of six colons. --Example2 (talk) 13:00, 26 June 2019 (UTC) [reply]

I request an edit to be made. --Example1 (talk) 16:00, 26 June 2019 (UTC) [reply]

☒ Request withdrawn I realized there is a better alternative. This is an example of using Template:withdrawn to indicate some status using recognizable icons. There are different templates used in various types of discussions. --Example1 (talk) 16:30, 26 June 2019 (UTC) [reply]

As Example1 has said elsewhere this is an example of Template:tg being used to quote other users. It can also be used to quote policies, articles, or any other pages. --Example2 (talk) 17:00, 26 June 2019 (UTC) [reply]

(edit conflict) This is an example of Template:Edit conflict being used to signal to other users that this response was being typed at the same time as the one above it and it resulted in an edit conflict. Note that the time of this message is just one minute after the previous. --Example1 (talk) 17:01, 26 June 2019 (UTC) [reply]

- Certainly, we could help. Using asterisks * instead of : is another approach used to format discussion threads. --Example2 (talk) 10:15, 26 June 2019 (UTC) [reply]

This component is responsible for displaying a set of banner images in the top bar of the website.

Code Overview:

- Displays a set of banner images at the top of the website.
- Initially hardcoded data for the banner images, later updated to use dynamic Markdown files for easier content management.

Features

- Responsive design ensures it looks good on different screen sizes.
- Includes links wrapped around banner images, directing users to specified URLs.

2.2.2 Navbar

This component handles the main navigation menu of the site, including a responsive hamburger menu for mobile views.

Code Overview:

- Manages the main navigation menu, including a responsive design that adapts to mobile screens with a hamburger menu.

2.6 Why References are needed?

- Verifiability
- Neutral Point of View
- Notability

2.7 What Makes a Good Reference?

- Independent & reliable sources
- News articles
- Books
- Academic Journals
- Government Websites

2.8 Creating User page

- Login to Wikipedia
- Go to User page (Click on your Wikipedia Username > top right corner)
- Click Edit
- Start Editing (Write about yourself)
- Name
- Studies
- Passion
- Add images
- Preview & Publish

Chapter 3

Wikimedia Commons

3.1 What is Wikimedia Commons?

Wikimedia Commons is a **free online media repository** that provides access to **images, audio, videos, and other media files** that can be freely used by anyone.

It is a project of the **Wikimedia Foundation**, the same organization that operates **Wikipedia**, and is used to support all Wikimedia projects by providing a central place to store freely licensed media files.

3.2 Features of Wikimedia Commons

- Free To Use
- Multilingual Support
- Support Multiple formats (JPG, PNG, SVG, MP3, OGG, WebM, Pdf etc..)
- Used by Wikipedia and Other Projects
- Community Driven

3.3 How to Use Wikimedia Commons Files ?

1. Visit the Website (<https://commons.wikimedia.org>)
2. Login or Signup
3. Search for Media
4. Choose the File
5. Check the License
6. Download the File
7. Use in Your Work

Chapter 4

Wikidata

4.1 Introduction

Wikidata is a free and open knowledge base that stores structured data to support Wikipedia, Wikimedia Commons, and other Wikimedia projects. It was launched on October 29, 2012, by the Wikimedia Foundation and is collaboratively edited by volunteers from around the world. Unlike Wikipedia, which stores textual content for human readers, Wikidata is designed to store machine-readable data, making it highly useful for applications, websites, data analysis, and AI systems.

4.2 How Wikidata Works

Wikidata stores information in the form of items and statements. Each item represents a concept or object — like a person, place, book, or idea — and has a unique Q-number (e.g., Q42 for Douglas Adams). Each item contains statements, which are structured as property-value pairs. For example, the item for "Albert Einstein" (Q937) has properties like "date of birth", "occupation", and "country of citizenship". These statements can also include references and qualifiers, making Wikidata not just a database of facts, but a context-aware knowledge system.

4.3 Uses and Applications

Wikidata is widely used in Wikipedia infoboxes, where it supplies real-time, updatable data like population, coordinates, or political leadership. It is also a valuable resource for researchers, developers, and data scientists, who use it to build tools, visualizations, AI systems, and more. For instance, voice assistants like Siri or Google Assistant can access Wikidata indirectly to answer factual questions. It is also used in academic projects, educational platforms, and linked open data ecosystems.

Chapter 5

Wiktionary and Wikisource

5.1 Wiktionary

Wiktionary is a free, online, multilingual dictionary that aims to document all words in all languages. It is one of the sister projects of Wikipedia, and it is also maintained by the Wikimedia Foundation. Launched on December 12, 2002, Wiktionary is entirely community-driven, meaning that anyone with internet access can contribute to or edit its content. Unlike traditional printed dictionaries, Wiktionary is not limited by space or language boundaries — its mission is to include as many words, meanings, and translations as possible.

The platform provides a wide range of linguistic information for each word, including definitions, etymology (word origin), pronunciation (with audio), grammatical usage, inflections, example sentences, and even synonyms and antonyms. Many entries also include translations of the word into other languages, making it a valuable resource for language learners, translators, and linguists. Its multilingual and cross-lingual nature allows users to find information about a word in one language while learning how it is used in many others.

Wiktionary entries are structured in a very detailed and organized manner. For each word, it distinguishes different meanings based on parts of speech — for example, the word "run" can be listed separately as a noun and as a verb, each with its own definitions, usage notes, and grammatical forms. This structure allows users to not only understand the meaning but also how the word behaves in sentences. Audio pronunciations and phonetic notations (like IPA – International Phonetic Alphabet) also enhance learning and pronunciation accuracy.

Being freely licensed under Creative Commons Attribution-ShareAlike, the content on Wiktionary can be copied, reused, modified, and redistributed even for commercial purposes — as long as proper credit is given. This makes it an open and powerful linguistic resource for developers, educators, and students alike. In summary, Wiktionary is more than just a dictionary; it is a living, expanding language database built by the people, for the people.

5.2 Wikisource

Wikisource is a free online digital library of source texts — documents, books, historical writings, legal texts, and other published works that are either in the public domain or available under free licenses. It is a project run by the Wikimedia Foundation, the same organization behind Wikipedia, Wiktionary, and other open knowledge initiatives. Launched in 2003, Wikisource's main goal is to provide a platform where users can access, read, and contribute to authentic textual content from around the world.

Unlike Wikipedia, which provides articles summarizing knowledge, Wikisource hosts original full texts of written works. This includes literary classics, religious texts, speeches, constitutions, letters, autobiographies, historical records, scientific papers, and translations. For example, you can read the entire Indian Constitution, Shakespeare's plays, or Tagore's poetry on Wikisource — all free of cost and freely shareable. One of the core principles of Wikisource is accuracy, so each document is proofread and verified by contributors, often using scanned images of original books hosted on Internet Archive or Wikimedia Commons.

A unique feature of Wikisource is the side-by-side view of scanned pages and their transcriptions, allowing readers to verify that the typed version matches the original scan. This makes it an excellent platform for academic referencing, historical research, and language learning. It is also multilingual, with language-specific editions (e.g., Hindi Wikisource, Tamil Wikisource, English Wikisource), allowing native readers to access literature and cultural heritage in their own language.

Wikisource is entirely community-driven and freely editable, meaning anyone can contribute by uploading public domain works, proofreading scanned documents, or translating texts into other languages. All content is licensed under the Creative Commons Attribution-ShareAlike license or is in the public domain, so it can be reused freely with proper attribution. Overall, Wikisource is a powerful tool for preserving, sharing, and democratizing access to the world's literary and historical documents.

5.2.1 Types of Content on Wikisource

- Historical texts (e.g., speeches, treaties, letters)
- Literary works (e.g., novels, plays, poetry)
- Legal documents (e.g., constitutions, laws, court decisions)
- Religious texts (e.g., scriptures, commentaries)
- Translations of classic texts

- Government records and official reports

5.2.2 Unique Features of Wikisource

- Side by Side View
- Proofreading Interface
- Transclusion
- Searchable Text

5.2.3 Accessing Wikisource

- Just visit the main page:
 <https://wikisource.org>
- Or language-specific versions like:
 <https://hi.wikisource.org> for Hindi

Chapter 6

IndicWiki Data Collection & Annotation Project, Lingua Libre

6.1 NER TAGs

- NER is a task in **Natural Language Processing (NLP)** that involves identifying and classifying named entities in text into predefined categories.
- Example: In the sentence "**Narendra Modi visited New York in June**", NER identifies:
 - "Bhagwan Jha" → **Person**
 - "New York" → **Location**
 - "June" → **Date**

6.2 Common NER Tags

1. PERSON (PER)

- Meaning: Identifies names of individual people or fictional characters.
- Examples:
 - Mahatma Gandhi, Rani Lakshmi Bai, Virat Kohli
- Use Case:
 - To extract names from articles, resumes, social media, etc.

2. LOCATION (LOC)

- Meaning: Identifies names of natural or geographical locations (not geopolitical).

- Examples:
 - Mount Everest, Thar Desert, Ganga River
- Use Case:
 - Used in travel applications, geography-based search, etc.

3. ORGANIZATION (ORG)

- Meaning: Names of companies, schools, government agencies, institutions.
- Examples:
 - ISRO, Indian Railways, Google, IIT Delhi
- Use Case:
 - Detect organizations in contracts, news, and social media.

4. GPE (Geo-Political Entity)

- Meaning: Political entities like countries, states, or cities with governance.
- Examples:
 - India, Rajasthan, New York, Bengaluru
- Difference from LOC:
 - LOC is geographical; GPE is political. E.g., "Himalayas" is LOC, but "India" is GPE.

5. DATE

- Meaning: Any mention of calendar dates.
- Examples:
 - 15 August 1947, January 26, next Monday
- Use Case:
 - Timeline generation, appointment scheduling, history applications.

6. TIME

- Meaning: Mentions of specific times or durations.
- Examples:
 - 5:30 PM, midnight, 2 hours later
- Use Case:
 - Alarms, reminders, event planning.

7. NRP (Nationalities, Religious & Political Groups)

- Meaning: Nationalities, religions, political parties.
- Examples:
 - *Indian, Hindu, Congress, BJP, Muslim*

6.2 Annotation Process

Step 1: Prepare for Annotation

- Prepare for Annotation
- Locate your raw files

Step 2: Using the Annotation Tool

- Access the tool
- Upload and annotate
- Tag Guidelines

Step 3: Upload Annotated Files

- Save annotated files
- Commit changes

6.3 Lingua Libre

Lingua Libre is a collaborative project created by the **Wikimedia France** community. Its main goal is to **record and share pronunciations** of words, phrases, and sentences in as many languages as possible — especially **underrepresented**,

regional, minority, and endangered languages. It is a part of the **Wikimedia movement**, meaning all its content is available for free and can be reused on platforms like **Wiktionary**, **Wikipedia**, and **Wikimedia Commons**.

Lingua Libre helps bridge the gap between **written text and spoken language**, making it easier for language learners, researchers, and digital systems (like speech recognition tools) to understand how words are actually pronounced by native speakers. All recordings are released under **Creative Commons licenses**, making them accessible and reusable for educational and technological purposes.

The recordings are made using a **web-based tool** that allows users to quickly and efficiently record multiple words in one sitting. These are automatically uploaded to **Wikimedia Commons** and can be linked to corresponding entries in Wiktionary and other projects. Lingua Libre supports **hundreds of languages**, including many that are not well represented in digital form.

6.3.1 Key Features

- Global and Inclusive
- Easy Recording Tool
- Reusable Content
- Integration with Wikimedia Projects
- Open and Volunteer-Based
- Impact and Reach

Chapter 7

Wikimedia Technical Space

7.1 Version Control System

- [Gerrit](#) - Primary version control system
- [Get started with Gerrit](#)
- [GitHub](#) - Used by some projects
- Gerrit repos are mirrored
- [GitLab](#) - Self hosted

7.2 Project management

- [Phabricator](#)
 - Task management
- GitHub
 - Some projects use GitHub

Developer Resources :

<https://developer.wikimedia.org/>

<https://codesearch.wmcloud.org/search/>

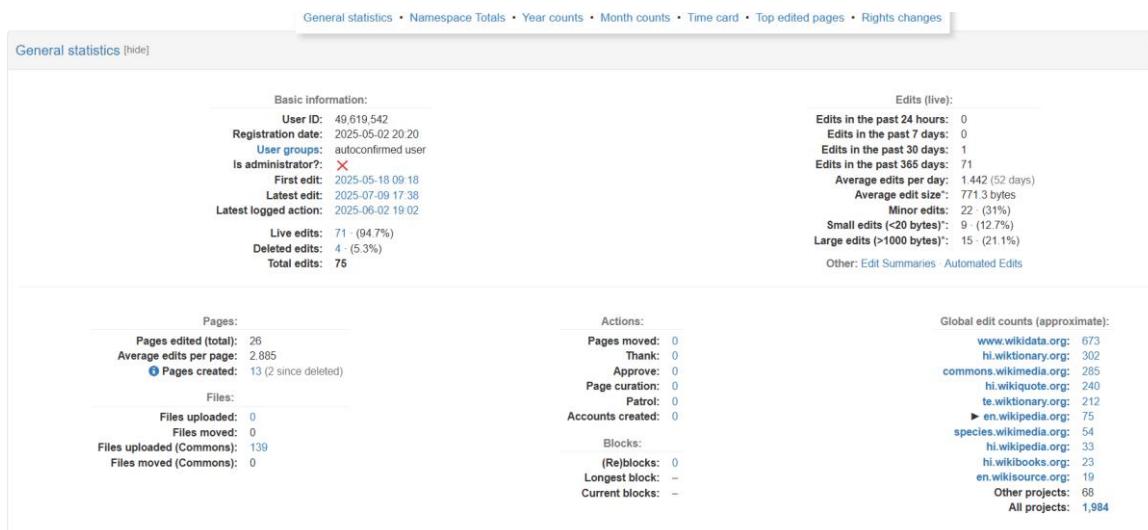
7.3 Developer setup

- Login / Sign up to MediaWiki.org
- Login to Phabricator using MediaWiki.org account
- Sign up for a Wikimedia developer account
- Setup Git / Gerrit
- Create a PR / patch for the sandbox Gerrit repo

Chapter 8

References and Our Details (Projects)

- GITHUB Indic repository : <https://github.com/Soumyadip0806/indicWikiData/tree/indicWiki>
- Full Contributions : <https://xtools.wmcloud.org/ec/en.wikipedia.org/Baap8969>
- User_Name on Wikipedia and its sister platform : Baap8969
- Projects Contributions :



- Outreach Dashboard : [INTERNSHIP POSITION](#)
- GitHub Page: [Bhagwanjha85/IIITH Indicwiki Internship](#)
- Annotation Tool: <https://plural.iiit.ac.in/ner-annotator/>
- Worked/Working on : wikipedia, wikidata, Wikisource, Wikimedia commons, Wiktionary, Wikivoyage, Wikinews, Wikispecies, Wikitech, Wikiquotes, Metawiki, Wikibook, Phabricator, GIT, NER- Data Annotations, Lingua Libre
- Contact : 9014120442 , indicwiki@iiit.ac.in