# Information Synchronization across Multilingual Semi-structured Tables

Siddharth Khincha<sup>1</sup>, Chelsi Jain<sup>2</sup>, Vivek Gupta<sup>3†</sup>, Tushar Kataria<sup>3†</sup>, Shuo Zhang<sup>4</sup> <sup>1</sup>IIT Guwahati, <sup>2</sup>CTAE, Udaipur, <sup>3†</sup>University of Utah, <sup>4</sup>Bloomberg,

## . Information Mismatch in Tables Across Languages



### English Table

Hindi Table

- Janaki Ammal Infoboxes: English (right) vs. Hindi (left). Hindi lacks "British Rule of India" context.
- Value mismatches: (a) Hindi table doesn't state Died key's state. (b) Institution values differ -Hindi mentions "residence," English doesn't.
- Missing keys in Hindi table: "Thesis," "Awards," and "Alma Mater." Neither mentions parents, early education, or honors.

### 2. Problem Magnitude

- Articles in More than 300 languages.
- English has the most significant Wikipedia covering 23% (11%) of total pages (articles).



• Most users' edits (76%) are also done in English Wikipedia.

## **3. Our Contributions**

### **1** INFOSYNCDataset

- 100K entity-centric wikipedia Infoboxes table across 14 languages
- Approximately 3.5K human annotated table alignment pairs
- Proposed a two-step approach as a solution, include Information
- Alignment to mapped similar rows

• **Update** update missing/outdated rows for aligned tables across multilingual entity centric tables

## 4. Dataset Details :- Language and category selection

### 1 Languages

- Languages are selected to cover all the continents.
- 4 low resource Hindi(hi), Cebuano(ceb), , Turkish(tr), and Afrikaans(ak)
- 7 medium resource German(de), Korean(ko), Russian(ru), Arabic(ar), Chinese(zh), Swedish(sv), Dutch(nl)
- 3 high resource English(en), French(fr), Spanish(es)

### **2** Entities

- Each Entity selected contains an Infobox in at least 5 languages
- 3 Categories Selection
- 21 simple, diverse, and popular topics: Airport, Album, Animal, Athlete, Book, City, College, Company, Country, Food, Monument, Movie Musician, Nobel, Painting, Person, Planet, Shows, and Stadiums.

# 5. Method: Alignment



**Corpus-based** : Align rows based on keys using their cosine similarity across a category using majority voting.

**Key-only** : This module aligns rows with *key* similarity score greater than a threshold value, only if they are mutually most similar keys

Key value bidirectional : This module aligns rows with key+value similarity score greater than a threshold value, only if they are mutually most similar rows.

Key value unidirectional: This module aligns rows with key+value similarity greater than a threshold. They **do not** have to be mutually most

Multi-key : This module considers the case where one row from table needs to be mapped to multiple rows in the second table. It is valid multi-key alignment when the merge value-combination similarity score exceeds that of the most similar key.

Metho SimC LaBS XLM

Corpu + Key + Key + Key + Mu

Our rule-based method efficiently updates a large number of rows, with the highest number of updates being in row transfers.

6. Method: Rule-Based Update

| P.R. | Rule Name      | Logical Rule $\forall_{(\mathbf{R}_{T_x},\mathbf{R}_{T_y})} \mathbf{L} \mapsto \mathbf{R}$   | Update Type      |
|------|----------------|--|------------------|
| 1    | Row Transfer   | $\forall_{(R_{T_x},R_{T_y})}Al_{T_x}^{T_y}(R_{T_x};R_{T_y})=0$   | Row Addition     |
|      |                | $\mapsto T_y \cup tr_x^y(\mathbf{R}_{T_x}) \bigwedge \operatorname{Al}_{T_x}^{T_y}(\mathbf{R}_{T_x}; tr_x^y(\mathbf{R}_{T_x})) = 1$  |                  |
| 2    | Multi-Match    | $\forall_{(\mathbf{R}_{T_x},\mathbf{R}_{T_y})}(\sum_{\mathbf{R}_{T_x}} \operatorname{Al}_{T_x}^{T_y}(\mathbf{R}_{T_x};\mathbf{R}_{T_y})) > 1$  | Row Delete       |
|      |                | $\mapsto \{T_y \setminus \bigcup_{(\forall_{\mathbf{R}_{T_y}} \operatorname{Al}_{T_x}^{T_y}(\mathbf{R}_{T_x};\mathbf{R}_{T_y})=1)} \mathbf{R}_{T_y}\} \bigcup tr_x^y(\mathbf{R}_{T_x}) \bigwedge \operatorname{Al}_{T_x}^{T_y}(\mathbf{R}_{T_x};tr_x^y(\mathbf{R}_{T_x})) = 1$ |                  |
| 3    | Time-based     | $\forall_{(\mathbf{R}_{T_x},\mathbf{R}_{T_y})} \operatorname{Al}_{T_x}^{T_y}(\mathbf{R}_{T_x};\mathbf{R}_{T_y}) = 1 \bigwedge (\operatorname{isTime}(\mathbf{R}_{T_x},\mathbf{R}_{T_y}) = 1)$  | Value Substitute |
|      |                | $\bigwedge (\operatorname{exTime}(R_{T_x}) > \operatorname{exTime}(R_{T_y})) \mapsto R_{T_y} \leftarrow tr_x^y(R_{T_x})$   |                  |
| 4    | Positive Trend | $\forall_{(R_{T_x},R_{T_y},PosTrend)}Al_{T_x}^{T_y}(R_{T_x};R_{T_y}) = 1 \bigwedge exKey(R_{T_x}) \in PosTrend$  | Value Substitute |
|      | or             | $\bigwedge \mathbf{R}_{T_x} > \mathbf{R}_{T_y} \mapsto \mathbf{R}_{T_y} \leftarrow \mathbf{R}_{T_x}$   |                  |
|      | Negative Trend | $\forall_{(R_{T_x},R_{T_y},NegTrend)}Al_{T_x}^{T_y}(R_{T_x};R_{T_y}) = 1 \bigwedge exKey(R_{T_x}) \in NegTrend$  | Value Substitute |
|      |                | $\bigwedge R_{T_x} < R_{T_y} \mapsto R_{T_y} \leftarrow R_{T_x}$   |                  |
| 5    | Append Value   | $\mathbf{R}_{T_x} = \mathbf{V} \bigwedge \forall_{(\mathbf{R}_{T_x}, \mathbf{R}_{T_y})} \mathbf{Al}_{T_x}^{T_y}(\mathbf{R}_{T_x}; \mathbf{R}_{T_y}) = 1 \bigwedge  \mathbf{R}_{T_x}[k]  >  \mathbf{R}_{T_y}[k] $   | Value Addition   |
|      |                | $\mapsto \forall_{(v \in \mathbf{R}_{T_r}[k] \land \notin tr_x^y(\mathbf{R}_{T_r}[k]))} \mathbf{R}_{T_y} \leftarrow \mathbf{R}_{T_y} \cup tr_x^y(v)$   |                  |
| 6    | HR to LR       | $(T_x, T_y) \in (HR, LR) \bigwedge \forall_{(\mathbf{R}_{T_x}, \mathbf{R}_{T_y})} \operatorname{Al}_{T_x}^{T_y}(\mathbf{R}_{T_x}; \mathbf{R}_{T_y}) = 1$   | Value Substitute |
|      |                | $\bigwedge tr_x^{en}(\mathbf{R}_{T_x}) \neq tr_y^{en}(\mathbf{R}_{T_y}) \mapsto \mathbf{R}_{T_y} \leftarrow tr_x^y(\mathbf{R}_{T_x})$  |                  |
| 7    | # Rows         | $ T_x  >>  T_y  \bigwedge \forall_{(R_{T_x},R_{T_y})} \operatorname{Al}_{T_x}^{T_y}(R_{T_x};R_{T_y}) = 1 \bigwedge tr_x^{en}(R_{T_x}) \neq tr_y^{en}(R_{T_y})$   | Value Substitute |
|      |                | $\mapsto \mathbf{\tilde{R}}_{T_y} \leftarrow tr_x^y(\mathbf{R}_{T_x})$   |                  |
| 8    | Rare Keys      | $\forall_{(R_{T_x},R_{T_y},RarKeys)}Al_{T_x}^{T_y}(R_{T_x};R_{T_y}) = 1 \bigwedge tr_x^{en}(R_{t_x}) \neq tr_y^{en}(R_{t_y})$  | Value Substitute |
|      |                | $\bigwedge \forall_{(R_{T_x},R_{T_y}}   exKey(R_{T_x}) \in RarKey  >   exKey(R_{T_y}) \in RarKey  \mapsto R_{T_y} \leftarrow R_{T_x}$  |                  |

## 7. Result: Alignment

Proposed similarity-based alignment method outperforms different multi-lingual baseline.

| Match                        |  |  |   | UnMatch   |  |  |  |
|------------------------------|--|--|---|---|--|--|--|
| $T_{en} \leftrightarrow T_x$ | $T_x \leftrightarrow T_y$  | $T_{en} \stackrel{*}{\leftrightarrow} T_{hi}$  | $T_{en} \stackrel{*}{\leftrightarrow} T_{zh}$   | $T_{en} \leftrightarrow T_x$  | $T_x \leftrightarrow T_y$                              | $T_{en} \stackrel{*}{\leftrightarrow} T_{hi}$  | $T_{en} \stackrel{*}{\leftrightarrow} T_{zh}$  |
| 75.78                        | 68.46  | 77.93  | 80.47   | 79.11   | 76.3   | 73.31  | 74.91  |
| 85.25                        | 78.44  | 88.98  | 89.1  | 87.03   | 81.7   | 88.98  | 85.06  |
| 80.98                        | 73.74  | 82.9   | 86.73   | 82.68   | 80.22  | 76.73  | 81.85  |
| 83.38                        | 75.02  | 86.85  | 88.08   | 85.42   | 80.65  | 83.14  | 83.1   |
| 82.85                        | 78.63  | 86.08  | 87.58   | 84.2  | 83.45  | 83.14  | 83.76  |
| 84.55                        | 77.45  | 87.64  | 88.7  | 86.3  | 82.28  | 83.14  | 84.3   |
|                              |  |  |   |   |  |  |  |
| 61.86                        | 56.74  | 57.34  | 69.33   | 70.51   | 71.73  | 54.01  | 63.11  |
| 70.41                        | 62.14  | 73.4   | 74.67   | 73.85   | 73.52  | 62.49  | 66.23  |
| 87.71                        | 84.2   | 90.07  | 93.04   | 89.51   | 85.52  | 85.06  | 89.2   |
| 87.89                        | 84.33  | 90.34  | 93.12   | 89.52   | 85.42  | 85.16  | 88.62  |
| 87.91                        | 84.36  | 90.14  | 92.8  | 89.3  | 85.46  | 84.98  | 88.15  |
|                              | $\begin{array}{c} T_{en} \leftrightarrow T_x \\ 75.78 \\ 85.25 \\ 80.98 \\ 83.38 \\ 82.85 \\ 84.55 \\ \hline 61.86 \\ 70.41 \\ 87.71 \\ 87.89 \\ 87.91 \\ 87.91 \end{array}$ | $T_{en} \leftrightarrow T_x \ T_x \leftrightarrow T_y$ 75.7868.4685.2578.4480.9873.7483.3875.0282.8578.6384.5577.4561.8656.7470.4162.1487.7184.287.8984.3387.9184.36 | Match $T_{en} \leftrightarrow T_x \ T_x \leftrightarrow T_y \ T_{en} \stackrel{*}{\leftrightarrow} T_{hi}$ 75.7868.4677.9385.2578.4488.9880.9873.7482.983.3875.0286.8582.8578.6386.0884.5577.4587.6461.8656.7457.3470.4162.1473.487.7184.290.0787.8984.3390.34 <b>87.9184.3690.14</b> | Match $T_{en} \leftrightarrow T_x \ T_x \leftrightarrow T_y \ T_{en} \stackrel{*}{\leftrightarrow} T_{hi} \ T_{en} \stackrel{*}{\leftrightarrow} T_{zh}$ 75.7868.4677.9380.4785.2578.4488.9889.180.9873.7482.986.7383.3875.0286.8588.0882.8578.6386.0887.5884.5577.4587.6488.761.8656.7457.3469.3370.4162.1473.474.6787.7184.290.0793.0487.8984.3390.3493.12 <b>87.9184.3690.1492.8</b> | $\begin{array}{c c c c c c c c c c c c c c c c c c c $ | MatchUnl $T_{en} \leftrightarrow T_x \ T_x \leftrightarrow T_y \ T_{en} \stackrel{*}{\leftrightarrow} T_{hi} \ T_{en} \stackrel{*}{\leftrightarrow} T_{zh} \ T_{en} \leftrightarrow T_x \ T_x \leftrightarrow T_y$ 75.7868.4677.9380.4779.1176.385.2578.4488.9889.187.0381.780.9873.7482.986.7382.6880.2283.3875.0286.8588.0885.4280.6582.8578.6386.0887.5884.283.4584.5577.4587.6488.786.382.2861.8656.7457.3469.3370.5171.7370.4162.1473.474.6773.8573.5287.7184.290.0793.0489.5185.5287.8984.3390.3493.1289.5285.42 <b>87.9184.3690.1492.889.385.46</b> | MatchUnMatch $T_{en} \leftrightarrow T_x \ T_x \leftrightarrow T_y \ T_{en} \stackrel{*}{\leftrightarrow} T_{hi} \ T_{en} \stackrel{*}{\leftrightarrow} T_{zh} \ T_{en} \leftrightarrow T_x \ T_x \leftrightarrow T_y \ T_{en} \stackrel{*}{\leftrightarrow} T_{hi}$ 75.7868.4677.9380.4779.1176.373.3185.2578.4488.9889.187.0381.788.9880.9873.7482.986.7382.6880.2276.7383.3875.0286.8588.0885.4280.6583.1482.8578.6386.0887.5884.283.4583.1484.5577.4587.6488.786.382.2883.1461.8656.7457.3469.3370.5171.7354.0170.4162.1473.474.6773.8573.5262.4987.7184.290.0793.0489.5185.5285.0687.8984.3390.3493.1289.5285.4285.16 <b>87.9184.3690.1492.8<b>89.385.4684.98</b></b> |

### 8. Result: Update

| Rules      | $ _{T_{en}} \rightarrow T_x$ | $\begin{array}{c} \textbf{Gold} \\ T_x \to T_y \end{array}$ | Live Set | $\begin{array}{c} \mathbf{Pred} \\ T_{en} \to T_x \end{array}$ | icted $T_x \to T_y$ |
|------------|------------------------------|---|----------|--|---------------------|
| <b>R</b> 1 | 20320                        | 18055   | 4213     | 21246  | 17675               |
| R2         | 648                          | 502   | 207      | 1395   | 1852                |
| R3         | 546                          | 399   | 75       | 443  | 347                 |
| R4         | 142                          | 151   | 4        | 120  | 147                 |
| R5         | 3507                         | 2116  | 784      | 3193   | 1960                |
| R6         | 5237                         | 3047  | 332      | 5062   | 2891                |
| R7         | 2748                         | 1899  | 990      | 2732   | 1855                |
| R8         | 25                           | 77  | 5        | 29   | 82                  |
| Al         | 14967                        | 9715  | 2851     | 14864  | 10657               |

• Multilingual Tabular Information Synchronization is challenging problem.

- (b.) Updation
- 0.85





| Туре             |       | Total | Accept       | Reject      |
|------------------|-------|-------|--------------|-------------|
| Row Transfer     |       | 461   | 368(79.82%)  | 93(20.17%)  |
| Value Substitu   | ition | 70    | 52(74.28%)   | 18(25.72%)  |
| Append Value     | ;     | 72    | 46(63.88%)   | 26(36.12%)  |
| Total            |       | 603   | 466 (77.28%) | 136(22.72%) |
| Ln Pairs         | Tot   | al    | Accept       | Reject      |
| $T_{en} \to T_x$ | 204   | 4 10  | 61(78.92%)   | 43(21.07%)  |
| $T_x \to T_y$    | 210   | 6 10  | 69(78.25%)   | 47(21.75%)  |
| $T_x \to T_{en}$ | 18.   | 3 1.  | 36(74.31%)   | 47(25.68%)  |
| Total            | 60.   | 3 4   | 66(77.28%)   | 137(22.71%) |

# Updates

# 9. Human Assisted Wikipedia





Engineering

Human evaluator update the Wikipedia Infobox with our mehtod recommendation.

### 10. Example



### 11. Key Takeaways

<sup>2</sup> Taking Wikipedia Infoboxes as our case study, we created INFOSYNC

3 A two-step sequential approach (a.) Alignment and

• Alignment method outperforms baseline with an F1-score >

• The rule-based method received a 77.28 % approval rate on Wikipedia updates.