

Information Synchronization across Multilingual Semi-structured Tables

Siddharth Khincha¹, Chelsi Jain², Vivek Gupta^{3†}, Tushar Kataria^{3†}, Shuo Zhang⁴

¹IIT Guwahati, ²CTAE, Udaipur, ³University of Utah², ⁴Bloomberg,



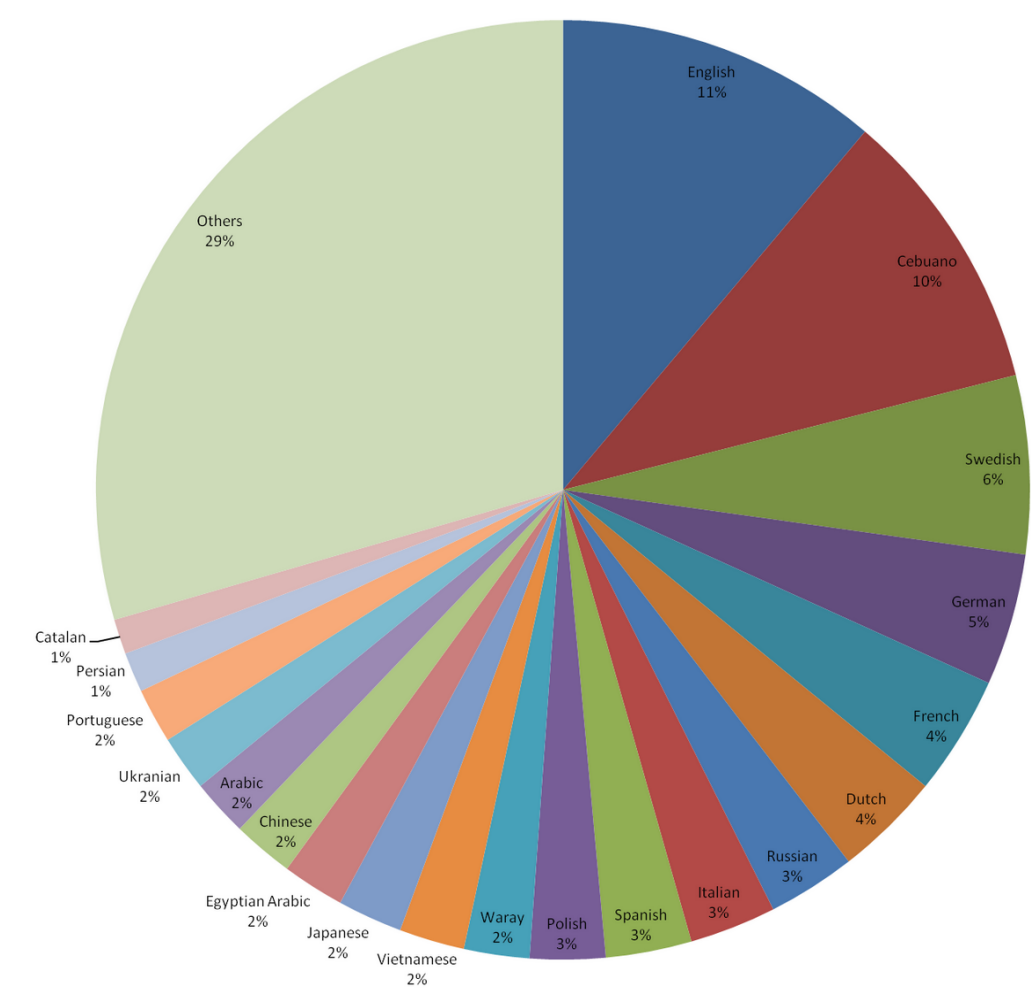
1. Information Mismatch in Tables Across Languages

English Table	Hindi Table
Janaki Ammal Born 4 November 1897 Tellicherry, Madras Presidency, British India Died 7 February 1984 (aged 86) Madras, Tamil Nadu Nationality Indian Alma mater University of Michigan Awards Padma Shri 1977 Scientific career Fields Botany, Cytology Institutions Madras University, John Innes Centre Thesis Chromosome Studies in <i>Nicandra physaloides</i>	जानकी अम्मल जन्म ४ नवम्बर १८९७ तेलिचरी, केरल मृत्यु फरवरी १९८४ (८७ वर्ष की आयु में) आवास भारत राष्ट्रीयता भारतीय क्षेत्र वनस्पति विज्ञान, कोशिका विज्ञान संस्थान यूनिवर्सिटी ऑफ़ मीचिगन, मद्रास <div> <div></div> Unmapped Rows <div></div> Matching Information <div></div> Cultural Context Missing <div></div> Value Mismatch </div>

- 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

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 - Update** update missing/outdated rows for aligned tables across multilingual entity centric tables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

- 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)
- Entities**
 - Each Entity selected contains an Infobox in at least 5 languages
- 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	Corpus-based : Align rows based on keys using their cosine similarity across a category using majority voting.
Key Only	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	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	Key value unidirectional : This module aligns rows with key+value similarity greater than a threshold. They do not have to be mutually most similar.
Multi-Key	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.

6. Method: Rule-Based Update

P.R.	Rule Name	Logical Rule $\forall (R_{Tx}, R_{Ty}) \text{ L} \mapsto \text{R}$	Update Type
1	Row Transfer	$\forall (R_{Tx}, R_{Ty}) \text{ Al}_{Tx}^{Ty}(R_{Tx}; R_{Ty}) = 0$ $\mapsto T_y \cup tr_y^x(R_{Tx}) \wedge \text{Al}_{Ty}^{Tx}(R_{Tx}; tr_y^x(R_{Tx})) = 1$	Row Addition
2	Multi-Match	$\forall (R_{Tx}, R_{Ty}) (\sum_{R_{Ty}} \text{Al}_{Tx}^{Ty}(R_{Tx}; R_{Ty})) > 1$ $\mapsto \{T_y \setminus \cup (\forall R_{Ty} \text{ Al}_{Tx}^{Ty}(R_{Tx}; R_{Ty})=1) R_{Ty}\} \cup tr_y^x(R_{Tx}) \wedge \text{Al}_{Tx}^{Ty}(R_{Tx}; tr_y^x(R_{Tx})) = 1$	Row Delete
3	Time-based	$\forall (R_{Tx}, R_{Ty}) \text{ Al}_{Tx}^{Ty}(R_{Tx}; R_{Ty}) = 1 \wedge (\text{isTime}(R_{Tx}, R_{Ty}) = 1)$ $\wedge (\text{exTime}(R_{Tx}) > \text{exTime}(R_{Ty})) \mapsto R_{Tx} \leftarrow tr_y^x(R_{Tx})$	Value Substitute
4	Positive Trend or Negative Trend	$\forall (R_{Tx}, R_{Ty}, \text{PosTrend}) \text{ Al}_{Tx}^{Ty}(R_{Tx}; R_{Ty}) = 1 \wedge \text{exKey}(R_{Tx}) \in \text{PosTrend}$ $\wedge R_{Tx} > R_{Ty} \mapsto R_{Tx} \leftarrow R_{Ty}$ $\forall (R_{Tx}, R_{Ty}, \text{NegTrend}) \text{ Al}_{Tx}^{Ty}(R_{Tx}; R_{Ty}) = 1 \wedge \text{exKey}(R_{Tx}) \in \text{NegTrend}$ $\wedge R_{Tx} < R_{Ty} \mapsto R_{Tx} \leftarrow R_{Ty}$	Value Substitute
5	Append Value	$R_{Tx} = V \wedge \forall (R_{Tx}, R_{Ty}) \text{ Al}_{Tx}^{Ty}(R_{Tx}; R_{Ty}) = 1 \wedge R_{Tx}[k] > R_{Ty}[k] $ $\mapsto \forall (v \in R_{Tx}[k] \wedge v \notin tr_y^x(R_{Tx}[k])) R_{Ty} \leftarrow R_{Tx} \cup tr_y^x(v)$	Value Addition
6	HR to LR	$(T_x, T_y) \in (H, R, L, R) \wedge \forall (R_{Tx}, R_{Ty}) \text{ Al}_{Tx}^{Ty}(R_{Tx}; R_{Ty}) = 1$ $\wedge tr_x^y(R_{Tx}) \neq tr_y^x(R_{Ty}) \mapsto R_{Ty} \leftarrow tr_x^y(R_{Tx})$	Value Substitute
7	# Rows	$ T_x >> T_y \wedge \forall (R_{Tx}, R_{Ty}) \text{ Al}_{Tx}^{Ty}(R_{Tx}; R_{Ty}) = 1 \wedge tr_x^y(R_{Tx}) \neq tr_y^x(R_{Ty})$ $\mapsto R_{Ty} \leftarrow tr_x^y(R_{Tx})$	Value Substitute
8	Rare Keys	$\forall (R_{Tx}, R_{Ty}, \text{RareKey}) \text{ Al}_{Tx}^{Ty}(R_{Tx}; R_{Ty}) = 1 \wedge tr_x^y(R_{Tx}) \neq tr_y^x(R_{Ty})$ $\wedge \forall (R_{Tx}, R_{Ty}) [\text{exKey}(R_{Tx}) \in \text{RareKey}] > [\text{exKey}(R_{Ty}) \in \text{RareKey}] \mapsto R_{Ty} \leftarrow R_{Tx}$	Value Substitute

7. Result: Alignment

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

Method	Match					UnMatch				
	$T_{en} \leftrightarrow T_x$	$T_x \leftrightarrow T_y$	$T_{en} \leftrightarrow T_{hi}$	$T_{en} \leftrightarrow T_{zh}$	$T_{en} \leftrightarrow T_x$	$T_x \leftrightarrow T_y$	$T_{en} \leftrightarrow T_{hi}$	$T_{en} \leftrightarrow T_{zh}$	$T_{en} \leftrightarrow T_x$	$T_x \leftrightarrow T_y$
SimCSE	75.78	68.46	77.93	80.47	79.11	76.3	73.31	74.91		
LaBSE	85.25	78.44	88.98	89.1	87.03	81.7	88.98	85.06		
mBERT-mp	80.98	73.74	82.9	86.73	82.68	80.22	76.73	81.85		
XML-R	83.38	75.02	86.85	88.08	85.42	80.65	83.14	83.1		
MPNet	82.85	78.63	86.08	87.58	84.2	83.45	83.14	83.76		
distill mBERT	84.55	77.45	87.64	88.7	86.3	82.28	83.14	84.3		
Our Approach										
Corpus-based	61.86	56.74	57.34	69.33	70.51	71.73	54.01	63.11		
+ Key Only	70.41	62.14	73.4	74.67	73.85	73.52	62.49	66.23		
+ Key-Val-Bi	87.71	84.2	90.07	93.04	89.51	85.52	85.06	89.2		
+ Key-Val-Uni	87.89	84.33	90.34	93.12	89.52	85.42	85.16	88.62		
+ Multi-Key	87.91	84.36	90.14	92.8	89.3	85.46	84.98	88.15		

8. Result: Update

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

Rules	Gold			Predicted	
	$T_{en} \rightarrow T_x$	$T_x \rightarrow T_y$	Live Set	$T_{en} \rightarrow T_x$	$T_x \rightarrow T_y$
R1	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
All	14967	9715	2851	14864	10657

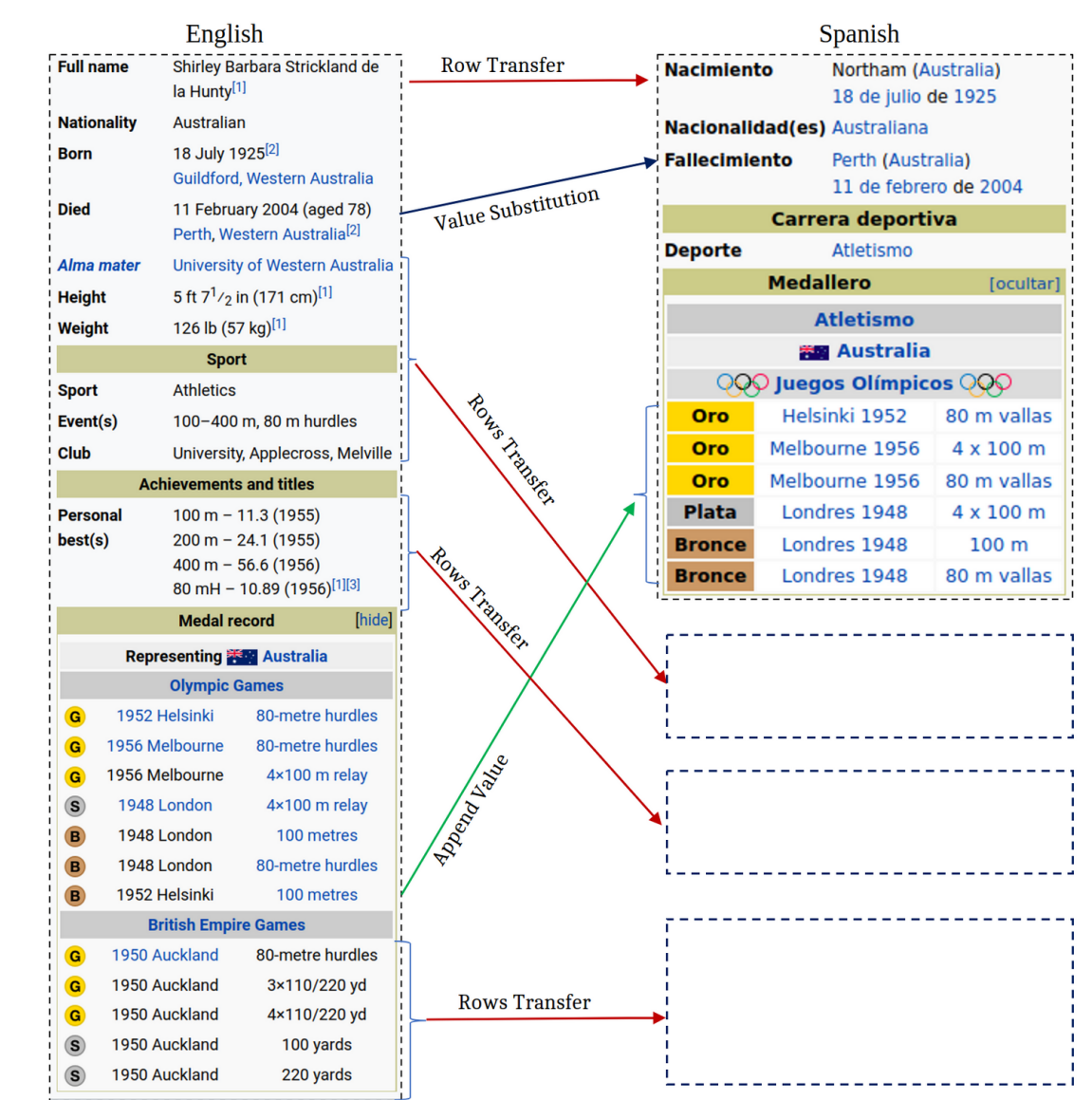
9. Human Assisted Wikipedia Updates

Human evaluator update the Wikipedia Infobox with our mehtod recommendation.

Type	Total	Accept	Reject
Row Transfer	461	368(79.82%)	93(20.17%)
Value Substitution	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	Total	Accept	Reject
$T_{en} \rightarrow T_x$	204	161(78.92%)	43(21.07%)
$T_x \rightarrow T_y$	216	169(78.25%)	47(21.75%)
$T_x \rightarrow T_{en}$	183	136(74.31%)	47(25.68%)
Total	603	466(77.28%)	137(22.71%)

10. Example



11. Key Takeaways

- Multilingual Tabular Information Synchronization is challenging problem.
- Taking Wikipedia Infoboxes as our case study, we created INFOSYNC
- A two-step sequential approach (a.) Alignment and (b.) Updation
 - Alignment method outperforms baseline with an F1-score > 0.85
 - The rule-based method received a 77.28 % approval rate on Wikipedia updates.