



JRC-Names-Retrieval

A Standardized Benchmark for Name Search



BABEL STREET

©2024 Babel Street, Inc. All Rights Reserved.

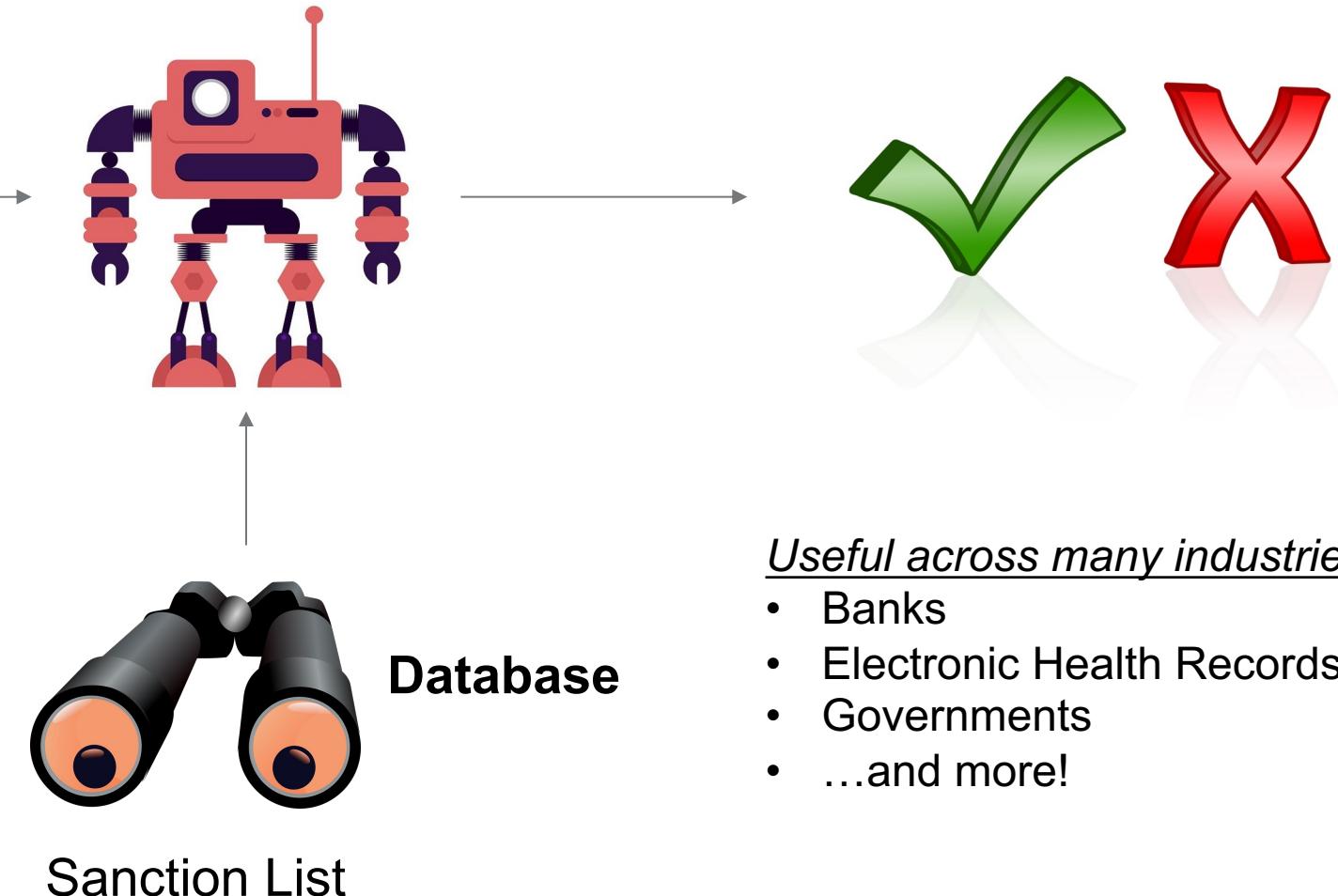
Philip Blair & Kfir Bar
22-24 May 2024 | Torino, IT

Name Search is a Common Problem



Query

Match?

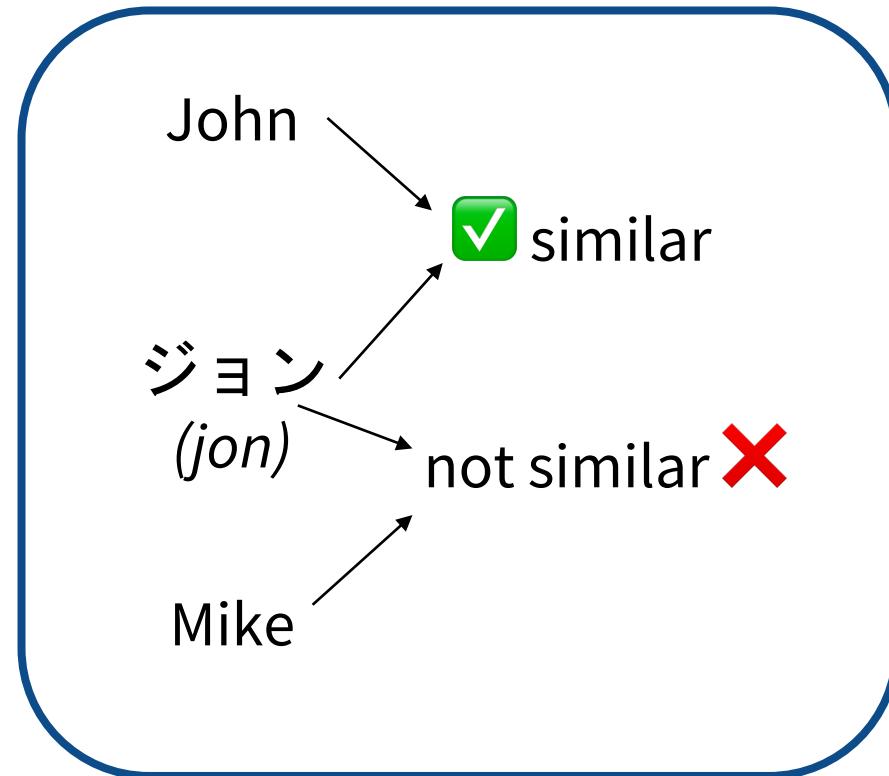


Useful across many industries:

- Banks
 - Electronic Health Records
 - Governments
 - ...and more!



Cross-Lingual Name Search is Hard!



What types of variations are we interested in?

- **Phonetically/Lexically Similar Names**, e.g. “John Smith” and “Jon Smith”, or “Phillippe Jones” and “Felipe Jones”
- **Initials**, e.g. “George Walker Bush” and “George W. Bush”
- **Missing Components**, e.g. “John A. Macdonald” and “John Macdonald”
- **Similar Names/Nicknames**, e.g. “William Clinton” and “Bill Clinton”, or “Michael Scott” and “Prison Mike”
- **Out-of-Order Components**, e.g. “François Mitterand” and “MITTERAND François”
- **Titles**, e.g. “Sir Tony Blair” and “Tony Blair”
- **Different Writing Scripts**, e.g. “Hu Jintao” and “胡锦涛”



Particularities for Organizations

- **Semantically Similar Components**, e.g. “Raven Train Company” and “Raven Locomotive, Inc.”.
- **Semantically Similar Components in Different Languages**, e.g. “株式会社京都アニメーション” (*Kabushiki-gaisha Kyōto Animēshon*) and “Kyoto Animation Co. Ltd.”.



Problem: No Benchmark!



Why a Standard Benchmark is Useful

1. Consistency

- All systems measure the same query/result pairs
- Future works able to compare performance without needing to run existing systems

2. Multilinguality

- We can design multiple (balanced) data splits to look at different types of multilingual support
 - Four writing script combinations published:
 - Latin, Cyrillic, and Arabic
 - Latin and Hani
 - Hangul and Hebrew
 - Devanagari and Kana (Hiragana+Katakana)
-
- ```
graph TD; A[Latin, Cyrillic, and Arabic] --> C[Realistic]; B[Latin and Hani] --> C; D[Hangul and Hebrew] --> E[Stress Test]; E[Devanagari and Kana (Hiragana+Katakana)] --> E;
```



# Background: JRC-Names



|         |   |   |                            |
|---------|---|---|----------------------------|
| 2453512 | P | u | George Tabayan             |
| 2453513 | P | u | Richard Lyal               |
| 2453510 | P | u | Francisco Pereira de Souza |
| 2450124 | P | u | Ana Tajadura-Jiménez       |
| 2450124 | P | u | Ana Tajadura-Jimenez       |
| 2450124 | P | u | Ana Tajadura Jimenez       |
| 2450124 | P | u | Ana Tajadura Jiménez       |
| 2453511 | P | u | Přemysl Kovář              |
| 2453511 | P | u | Přemysl Kovař              |
| ...     |   |   |                            |



# JRC-Names-Retrieval: Basic Idea

|         |   |   |                            |
|---------|---|---|----------------------------|
| 2453512 | P | u | George Tabayan             |
| 2453513 | P | u | Richard Lyal               |
| 2453510 | P | u | Francisco Pereira de Souza |
| 2450124 | P | u | Ana Tajadura-Jiménez       |
| 2450124 | P | u | Ana Tajadura-Jimenez       |
| 2450124 | P | u | Ana Tajadura Jimenez       |
| 2450124 | P | u | Ana Tajadura Jiménez       |
| 2453511 | P | u | Přemysl Kovář              |
| 2453511 | P | u | Přemysl Kovař              |
| ...     |   |   |                            |

**JRC-Names**

**Query:** Ana Tajadura-Jimenez

**Expected Results:**

- Ana Tajadura-Jiménez
- Ana Tajadura Jimenez
- Ana Tajadura Jiménez

**Query:** Přemysl Kovář

**Expected Results:**

- Přemysl Kovař

**JRC-Names-Retrieval**



# Data Statistics

| Type | Script         | # Names   | # Clusters |
|------|----------------|-----------|------------|
| PER  | Latn+Cyrl+Arab | 1,178,357 | 737,361    |
| PER  | Latn+CJK       | 97,077    | 38,268     |
| PER  | Hang+Hebr      | 1,331     | 640        |
| PER  | Deva+Kana      | 1,792     | 454        |
| ORG  | Latn+Cyrl+Arab | 30,113    | 2,942      |
| ORG  | Latn+CJK       | 27,702    | 2,857      |
| ORG  | Hang+Hebr      | 494       | 224        |
| ORG  | Deva+Kana      | 380       | 122        |

Table 1: JRC-Names-Retrieval training subsplit sizes for each entity type/script combination split.

| Entity Type | Script         | Database Size | # Queries | Avg. # Results |
|-------------|----------------|---------------|-----------|----------------|
| PER         | Latn+Cyrl+Arab | 5602          | 1885      | 2.97           |
| PER         | Latn+CJK       | 1457          | 268       | 5.43           |
| PER         | Hang+Hebr      | 316           | 266       | 1.18           |
| PER         | Deva+Kana      | 550           | 160       | 3.43           |
| ORG         | Latn+Cyrl+Arab | 841           | 152       | 5.53           |
| ORG         | Latn+CJK       | 738           | 135       | 5.51           |
| ORG         | Hang+Hebr      | 119           | 97        | 1.23           |
| ORG         | Deva+Kana      | 90            | 38        | 2.36           |



Table 2: JRC-Names-Retrieval evaluation subsplit sizes for each entity type/script combination split.

# Baselines



# Preexisting Baselines

## String-Based

- Double Metaphone
  - Approximates pronunciation of a string (European language-centric)
  - Index names by taking bigrams
- Lucene Fuzzy-Query
  - Damerau-Levenshtein edit distance

## Neural Network-Based

- Pretrained ByT5
- SimCSE
- Sentence-T5



# Additional Baseline: Metric Learning

- Architecture: ByT5 encoder, followed by linear layer and mean pool (across all characters) into single vector

$$\mathcal{L}_{contrast}(n_a, n_+) = \frac{\delta(n_a, n_+)^2}{2}, \quad \mathcal{L}_{contrast}(n_a, n_-) = \frac{[m - \delta(n_a, n_-)^2]^+}{2}$$
$$\mathcal{L}_{triplet}(n_a, n_+, n_-) = [\delta(n_a, n_+)^2 - \delta(n_a, n_-)^2 + m]^+$$

Figure 1: Loss functions used in our metric learning baseline.

- Negative mining based on NSEEN (Fakhraei et al., 2019)
  - After each epoch, create an approximate nearest neighbor index with each name
  - Look up each query in the ANN index and take the closest negative(s)



# Results (Person, Realistic)

*Latin, Cyrillic, and Arabic*

| Algorithm                            | MAP                               | Recall@1                           | Recall@5                          | Recall@10                         | Recall@50                         | Recall@100                        |
|--------------------------------------|-----------------------------------|------------------------------------|-----------------------------------|-----------------------------------|-----------------------------------|-----------------------------------|
| Double Metaphone <sup>†</sup>        | 0.16427                           | 0.08546                            | 0.16556                           | 0.16968                           | 0.17626                           | 0.18214                           |
| + pre-transliteration                | <b>0.91355</b>                    | <b>0.50809</b>                     | <b>0.89106</b>                    | <b>0.94381</b>                    | <b>0.97900</b>                    | <b>0.98530</b>                    |
| Sentence-T5 <sup>†</sup>             | 0.29894                           | 0.17942                            | 0.30980                           | 0.32471                           | 0.35936                           | 0.37559                           |
| + pre-transliteration                | 0.84549                           | 0.48011                            | 0.83323                           | 0.88937                           | 0.94981                           | 0.96674                           |
| ByT5                                 | 0.06841                           | 0.04018                            | 0.07354                           | 0.09702                           | 0.15778                           | 0.19513                           |
| + pre-transliteration                | 0.23765                           | 0.17140                            | 0.24862                           | 0.28274                           | 0.38930                           | 0.44801                           |
| SimCSE <sup>†</sup>                  | 0.22508                           | 0.11528                            | 0.24852                           | 0.29057                           | 0.36891                           | 0.39210                           |
| + pre-transliteration                | 0.74304                           | 0.42829                            | 0.74217                           | 0.80665                           | 0.88724                           | 0.91324                           |
| Lucene FuzzyQuery <sup>†</sup>       | 0.39133                           | 0.24901                            | 0.39832                           | 0.41042                           | 0.42456                           | 0.428090                          |
| + pre-transliteration                | 0.76186                           | 0.44020                            | 0.75026                           | 0.80762                           | 0.88085                           | 0.89819                           |
| <b>Ours (Latin-only)<sup>†</sup></b> | $0.50 \pm 0.01$                   | $0.289 \pm 0.00$                   | $0.47 \pm 0.02$                   | $0.50 \pm 0.03$                   | $0.62 \pm 0.01$                   | $0.66 \pm 0.01$                   |
| + pre-transliteration                | $0.93 \pm 0.01$                   | $0.513 \pm 0.00$                   | $0.91 \pm 0.01$                   | $0.95 \pm 0.00$                   | $0.97 \pm 0.00$                   | $0.98 \pm 0.00$                   |
| <b>Ours (La+Ar+Cy)</b>               | <b><math>0.95 \pm 0.01</math></b> | <b><math>0.521 \pm 0.01</math></b> | <b><math>0.93 \pm 0.00</math></b> | <b><math>0.97 \pm 0.00</math></b> | <b><math>0.99 \pm 0.00</math></b> | <b><math>0.99 \pm 0.00</math></b> |
| + pre-transliteration                | $0.94 \pm 0.01$                   | $0.520 \pm 0.00$                   | $0.92 \pm 0.01$                   | $0.96 \pm 0.01$                   | $0.98 \pm 0.00$                   | $0.98 \pm 0.00$                   |

*Latin and Hanzi*

| Algorithm                            | MAP                               | Recall@1                           | Recall@5                          | Recall@10                         | Recall@50                         | Recall@100                        |
|--------------------------------------|-----------------------------------|------------------------------------|-----------------------------------|-----------------------------------|-----------------------------------|-----------------------------------|
| Double Metaphone <sup>†</sup>        | 0.31134                           | 0.10988                            | 0.30951                           | 0.31853                           | 0.33370                           | 0.33619                           |
| + pre-transliteration                | <b>0.74852</b>                    | <b>0.17531</b>                     | <b>0.65205</b>                    | <b>0.77245</b>                    | <b>0.87325</b>                    | <b>0.91044</b>                    |
| Sentence-T5 <sup>†</sup>             | 0.31130                           | 0.11126                            | 0.30871                           | 0.31692                           | 0.33371                           | 0.34055                           |
| + pre-transliteration                | 0.68389                           | 0.17071                            | 0.60647                           | 0.69882                           | 0.79372                           | 0.83053                           |
| ByT5                                 | 0.07486                           | 0.03837                            | 0.07718                           | 0.10056                           | 0.18022                           | 0.24148                           |
| + pre-transliteration                | 0.15977                           | 0.06965                            | 0.14502                           | 0.16835                           | 0.27705                           | 0.37376                           |
| SimCSE <sup>†</sup>                  | 0.42119                           | 0.17357                            | 0.41617                           | 0.43420                           | 0.45417                           | 0.46238                           |
| + pre-transliteration                | 0.63187                           | 0.16978                            | 0.55939                           | 0.64179                           | 0.73881                           | 0.78750                           |
| Lucene FuzzyQuery <sup>†</sup>       | 0.26278                           | 0.09944                            | 0.26380                           | 0.27437                           | 0.28613                           | 0.29670                           |
| + pre-transliteration                | 0.58183                           | 0.15951                            | 0.52562                           | 0.60317                           | 0.65018                           | 0.66013                           |
| <b>Ours (Latin-only)<sup>†</sup></b> | $0.45 \pm 0.00$                   | $0.17 \pm 0.00$                    | $0.44 \pm 0.00$                   | $0.45 \pm 0.00$                   | $0.46 \pm 0.00$                   | $0.48 \pm 0.00$                   |
| + pre-transliteration                | $0.73 \pm 0.03$                   | $0.17 \pm 0.00$                    | $0.64 \pm 0.03$                   | $0.73 \pm 0.03$                   | $0.83 \pm 0.03$                   | $0.86 \pm 0.03$                   |
| <b>Ours (Latin+CJK)</b>              | <b><math>0.88 \pm 0.01</math></b> | <b><math>0.183 \pm 0.00</math></b> | <b><math>0.76 \pm 0.00</math></b> | <b><math>0.88 \pm 0.01</math></b> | <b><math>0.94 \pm 0.01</math></b> | <b><math>0.95 \pm 0.00</math></b> |
| + pre-transliteration                | $0.82 \pm 0.00$                   | $0.180 \pm 0.00$                   | $0.72 \pm 0.00$                   | $0.83 \pm 0.00$                   | $0.90 \pm 0.00$                   | $0.93 \pm 0.00$                   |

## Takeaway:

Fine-tuned ByT5 neural network performs comparably to systems which rely on transliteration.



# Results (Organization, Realistic)

*Latin, Cyrillic, and Arabic*

| Algorithm                                                     | MAP                               | Recall@1                          | Recall@5                          | Recall@10                         | Recall@50                         | Recall@100                        |
|---------------------------------------------------------------|-----------------------------------|-----------------------------------|-----------------------------------|-----------------------------------|-----------------------------------|-----------------------------------|
| Double Metaphone <sup>†</sup><br>+ pre-transliteration        | 0.35200<br>0.61014                | 0.11590<br>0.15976                | 0.34803<br>0.53487                | 0.37544<br>0.63575                | 0.44287<br>0.78958                | 0.47522<br>0.83476                |
| Sentence-T5 <sup>†</sup><br>+ pre-transliteration             | 0.40391<br><b>0.61739</b>         | 0.13180<br><b>0.16086</b>         | 0.38311<br><b>0.53640</b>         | 0.43399<br><b>0.64846</b>         | 0.51118<br><b>0.79956</b>         | 0.54079<br><b>0.85362</b>         |
| ByT5<br>+ pre-transliteration                                 | 0.08393<br>0.10012                | 0.03191<br>0.03355                | 0.06985<br>0.08542                | 0.10614<br>0.11469                | 0.24485<br>0.22763                | 0.33048<br>0.34737                |
| SimCSE <sup>†</sup><br>+ pre-transliteration                  | 0.37026<br>0.53519                | 0.13180<br>0.15099                | 0.35461<br>0.47588                | 0.40450<br>0.55888                | 0.47401<br>0.71107                | 0.50581<br>0.76173                |
| Lucene FuzzyQuery <sup>†</sup><br>+ pre-transliteration       | 0.26533<br>0.38923                | 0.10186<br>0.11754                | 0.26447<br>0.35471                | 0.28202<br>0.42708                | 0.31272<br>0.48684                | 0.32149<br>0.50164                |
| <b>Ours (Latin-only)<sup>†</sup></b><br>+ pre-transliteration | <b>0.41 ± 0.02</b><br>0.57 ± 0.03 | <b>0.13 ± 0.00</b><br>0.15 ± 0.00 | <b>0.37 ± 0.02</b><br>0.51 ± 0.03 | <b>0.42 ± 0.02</b><br>0.59 ± 0.03 | <b>0.54 ± 0.03</b><br>0.73 ± 0.04 | <b>0.61 ± 0.04</b><br>0.79 ± 0.04 |
| <b>Ours (La+Ar+Cy)</b><br>+ pre-transliteration               | <b>0.49 ± 0.04</b><br>0.49 ± 0.01 | <b>0.14 ± 0.00</b><br>0.14 ± 0.00 | <b>0.44 ± 0.04</b><br>0.43 ± 0.01 | <b>0.52 ± 0.05</b><br>0.51 ± 0.02 | <b>0.67 ± 0.03</b><br>0.65 ± 0.01 | <b>0.75 ± 0.04</b><br>0.73 ± 0.02 |

*Latin and Hanzi*

| Algorithm                                                     | MAP                               | Recall@1                          | Recall@5                          | Recall@10                         | Recall@50                         | Recall@100                        |
|---------------------------------------------------------------|-----------------------------------|-----------------------------------|-----------------------------------|-----------------------------------|-----------------------------------|-----------------------------------|
| Double Metaphone <sup>†</sup><br>+ pre-transliteration        | 0.43352<br>0.47043                | 0.13000<br>0.13580                | 0.41173<br>0.43247                | 0.46605<br>0.49914                | 0.55247<br>0.64037                | 0.59346<br>0.69901                |
| Sentence-T5 <sup>†</sup><br>+ pre-transliteration             | 0.51244<br><b>0.54248</b>         | 0.14235<br><b>0.14963</b>         | 0.47728<br><b>0.49852</b>         | 0.52914<br><b>0.55840</b>         | 0.59580<br><b>0.65346</b>         | 0.60963<br><b>0.69815</b>         |
| ByT5<br>+ pre-transliteration                                 | 0.08231<br>0.08535                | 0.02790<br>0.02605                | 0.06728<br>0.07432                | 0.09321<br>0.09802                | 0.22617<br>0.22506                | 0.34309<br>0.35148                |
| SimCSE <sup>†</sup><br>+ pre-transliteration                  | 0.47097<br>0.48958                | 0.14136<br>0.14593                | 0.43519<br>0.44420                | 0.49074<br>0.50593                | 0.57617<br>0.63889                | 0.60457<br>0.67284                |
| Lucene FuzzyQuery <sup>†</sup><br>+ pre-transliteration       | 0.34896<br>0.38029                | 0.10309<br>0.10889                | 0.34235<br>0.37691                | 0.37568<br>0.40901                | 0.40037<br>0.44728                | 0.40901<br>0.45593                |
| <b>Ours (Latin-only)<sup>†</sup></b><br>+ pre-transliteration | <b>0.48 ± 0.02</b><br>0.50 ± 0.03 | <b>0.14 ± 0.00</b><br>0.14 ± 0.00 | <b>0.45 ± 0.03</b><br>0.46 ± 0.03 | <b>0.49 ± 0.03</b><br>0.51 ± 0.04 | <b>0.58 ± 0.01</b><br>0.62 ± 0.02 | <b>0.64 ± 0.01</b><br>0.68 ± 0.01 |
| <b>Ours (Latin+Hanzi)</b><br>+ pre-transliteration            | <b>0.18 ± 0.00</b><br>0.18 ± 0.00 | <b>0.07 ± 0.00</b><br>0.07 ± 0.00 | <b>0.16 ± 0.00</b><br>0.17 ± 0.00 | <b>0.20 ± 0.00</b><br>0.21 ± 0.00 | <b>0.33 ± 0.00</b><br>0.33 ± 0.00 | <b>0.42 ± 0.00</b><br>0.42 ± 0.01 |

## Takeaway:

Sentence-T5 is a hard baseline to beat when transliteration is available, but fine-tuning ByT5 is better when it is not.



# Thank You!



BABEL STREET

©2024 Babel Street, Inc. All Rights Reserved.