

# ICU Done Wrong: How to Build a Polyglot Text Search Index

Sharzy Manaka, <https://github.com/SharzyL>

Apr 24, 2026

# Story 2020

In 2020, I was still new to the Telegram world, finding it hard to search for Telegram messages.

So I built a bot!

Commit `6c5a8a1`

Sharzy Luo committed on Jul 10, 2020

first commit

master · v0.5.0 ··· v0.1.1

4 files changed



# Story 2026

---

In 2026, I decided to rewrite it in   Rust.

- The search library [whoosh](#) has been unmaintained since 2022
- Unsatisfied with the current indexing mechanism
- Yes, we can vibe code in 2026

So I rewrote it in Rust + Tantivy, 3 months ago.

But why still not released?

# Recall: How Full-Text Search (FTS) Works

---

$doc_0$  = 武汉市长江大桥 Mayor of Wuhan, Jiang Daqiao

$doc_1$  = 南京市长江大桥 Nanjing Changjiang Bridge

Jieba( $doc_0$ ) = 武汉市 | 长江大桥 | Mayor | of | Wuhan | , | Jiang | Daqiao  
                  0          1 2          3 4          5 6          7 8 9          10 11          12

Jieba( $doc_1$ ) = 南京市 | 长江大桥 | Nanjing | Changjiang | Bridge  
                  0          1 2          3 4          5 6          7

Inserted into the inverted index:

Term	doc	freq	positions
■	$doc_0$	5	[2, 4, 6, 9, 11]
	$doc_1$	3	[2, 4, 6]
,	$doc_0$	1	[8]
Bridge	$doc_1$	1	[7]
Changjiang	$doc_1$	1	[5]
Daqiao	$doc_0$	1	[12]
.....			

# Recall: How Full-Text Search (FTS) Works (Cont'd)

Term	doc	freq	positions
	doc <sub>0</sub>	5	[2, 4, 6, 9, 11]
	doc <sub>1</sub>	3	[2, 4, 6]
,	doc <sub>0</sub>	1	[8]
Bridge	doc <sub>1</sub>	1	[7]
Changjiang	doc <sub>1</sub>	1	[5]
Daqiao	doc <sub>0</sub>	1	[12]
Jiang	doc <sub>0</sub>	1	[10]
Mayor	doc <sub>0</sub>	1	[3]
Nanjing	doc <sub>1</sub>	1	[3]
Wuhan	doc <sub>0</sub>	1	[7]
of	doc <sub>0</sub>	1	[5]
南京市	doc <sub>1</sub>	1	[0]
武汉市	doc <sub>0</sub>	1	[0]
长江大桥	doc <sub>0</sub>	1	[1]
	doc <sub>1</sub>	1	[1]

Query  $Q = \text{武汉市}$  is parsed as `TermQuery(武汉市)` and hits  $\text{doc}_0$ .

$$\begin{aligned} & \text{Score}_{\text{BM25}}(\text{doc}_0, Q) \\ &= \text{IDF} \cdot \frac{f \cdot (k_1 + 1)}{f + k_1 \cdot \left(1 - b + b \cdot \frac{|D|}{\text{avgdl}}\right)} \\ &= \ln\left(1 + \frac{N - n + 0.5}{n + 0.5}\right) \cdot \frac{1 \cdot 2.2}{1 + 1.2 \cdot \left(0.25 + 0.75 \cdot \frac{13}{10.5}\right)} \\ &= \ln 2 \cdot \frac{2.2}{2.41} \approx 0.63 \end{aligned}$$

- $N = 2$ : total number of documents
- $n = 1$ : documents containing 武汉市
- $f = 1$ : term frequency in  $\text{doc}_0$
- $k_1 = 1.2, b = 0.75$ : constant
- $|D| = 13$ : length of  $\text{doc}_0$  (in tokens)
- $\text{avgdl} = 10.5$ : average document length

# *CJK: From Jieba to Bigram*

# Jieba Fails for Japanese

---

$doc_2 =$  天気がいいから、散歩しましょう

$jieba(doc_2) =$  天 気 が い い か ら 、 散 歩 し ま し ょ う

$jieba(doc_2, HMM) =$  天 気 が い い か ら 、 散 歩 し ま し ょ う

ParseQuery(天気) is parsed as either

- TermQuery(天)  $\vee$  TermQuery(気)
- PhraseQuery(天, 気) (matches doc with 天 and 気 with adjacent positions)

Bad performance for both cases

# Phrase Query Not Working As Expected

---

Jieba has a `cut_for_search` mode that generates overlapping tokens, to optimize recall

JiebaForSearch(武汉市长江大桥) = 武汉<sub>0</sub> 武汉市<sub>0</sub> 长江<sub>1</sub> 长江大桥<sub>1</sub> 大桥<sub>2</sub>

---

ParseQuery(武汉长江) = PhraseQuery(武汉, 长江) **matches** 武汉市长江大桥.

Not necessarily problematic, but semantically incorrect

# CJK Bigram: Simple but Working

---

- Jieba and cang-jie work for Chinese but fail for Japanese
- Lindex works for Chinese, Japanese and Korean, but is weak for Chinese
- No single segmenter works well for all
- Really?

---

$doc_0 =$  武汉市长江大桥 Mayor of Wuhan, Jiang Daqiao

CJKBigram( $doc_0$ ) = 武汉 汉市 市长 长江 江大 大桥 Mayor of Wuhan Jiang Daqiao  
0 1 2 3 4 5 6 7 8 9 10

---

ParseQuery(武汉) = TermQuery(武汉), **matches**  $doc_0$

ParseQuery(武汉市) = PhraseQuery(武汉, 汉市), **matches**  $doc_0$

Drawback: 市长 matches 武汉市长江大桥, not necessarily bad

# Japanese Bigram: More Than Kanji=漢字

---

Japanese is a mix of Kanji=漢字 and Kana=仮名

CJK Bigram in Lucene only generates inner-script bigrams

$\text{CJKBigram}_{\text{Lucene}}(\text{可愛くてごめん}) = \text{可愛} \text{ くて} \text{ てご} \text{ ごめ} \text{ めん}$

Consider the query 可愛く

- $\text{TermQuery}(\text{可愛}) \vee \text{TermQuery}(\text{く})$ : matches too many
- $\text{PhraseQuery}(\text{可愛}, \text{く})$ : **does not match** 可愛くてごめん

# Japanese Bigram: More Than 漢字 Cont'd

---

We choose to generate kanji-kana bigrams

$\text{CJKBigram}_{\text{ours}}(\text{可愛くてごめん}) = \text{可愛 愛く くて てご ごめ めん}$

Consider the query 可愛く

- `PhraseQuery(可愛, 愛く)`: **matches** 可愛くてごめん

Drawback: producing undesired bigrams

# Unigram Desired

---

What if we want to search for a single CJK Unified Ideograph (汉字/漢字/.....)?

$doc_0 =$  ふたりの時間、言い訳の五千円

---

Double-field index:

**bigram** ( $doc_0$ ) = ふた たり りの の時 時間 言い い訳 訳の の五 五千 千円

**unigram** ( $doc_0$ ) = 時 間 言 訳 五 千 円

---

Only route unigram query to **unigram** index

$ParseQuery(\text{ふたり 訳}) = PhraseQuery_{\text{bigram}}(\text{ふた, たり}) \vee TermQuery_{\text{unigram}}(\text{訳})$

# ***ICU: Segmentation, Normalization and More***

# Why ICU?

---

Unicode is full of weird things:

- 17 kinds of spaces (Space\_Separator, Zs)
- (E202+U RIGHT-TO-LEFT OVERRIDE) tes retcarahc lortnoC
- Composable (𐀀-𐀁𐀂𐀃、𐀄𐀅、 ⇒ 𐀆) and decomposable (株式会社 = 株式会社) things
- Same codepoint for different scripts (U+201C “)
- Different codepoints for the same ideograph (神 U+FA19, 神 U+795E, 神 U+795E U+FE00)<sup>1</sup>

ICU is a library suite to handle Unicode.

---

<sup>1</sup>Typst still lacks complete IVS/SVS support

# Normalization

---

UAX #15 defines 4 normalization forms:

- **Canonical vs Compatibility:** canonical decomposition preserves formatting distinctions; compatibility decomposition erases them.
- **Decomposed vs Composed:** whether to recombine sequences back into precomposed codepoints after decomposition.

	Decomposed	Composed
<b>Canonical</b>	<b>NFD</b> — break apart only ü → u + ◌̈; 神 (U+FA19) → 神 (U+795E)	<b>NFC</b> — break apart, then recombine Å (Angstrom) → A + ◌° → Å
<b>Compatibility</b>	<b>NFKD</b> — also flatten formatting 株式会社 → 株式会社, ⑨ → 9, الله →   ج ج ه	<b>NFKC</b> — flatten, then recombine fí (ligature fi + acute) → fi◌' → fí

# CaseFold

---

Unicode defines 4 case folds

- C (Common)  $A \rightarrow a$ ,  $\zeta \rightarrow \sigma$
- S (Simple) German  $\beta \rightarrow \text{ß}$ , Greek  $A_i \rightarrow \alpha$
- F (Full) German  $\beta \rightarrow \text{ss}$ , Greek  $A_i \rightarrow \alpha\iota$
- T (Turkic) Turkic  $\text{I} \rightarrow \text{ı}$ , because uppercase  $\text{i}$  in Turkic is  $\text{İ}$

Possible combinations: C + F, C + S, C + F + T, C + S + T

$\text{CaseFold}_{C+F}(\text{İ}) = ?$

# CaseFold

---

Unicode defines 4 case folds

- C (Common)  $A \rightarrow a, \zeta \rightarrow \sigma$
- S (Simple) German  $\beta \rightarrow \text{ß}$ , Greek  $A_i \rightarrow \alpha$
- F (Full) German  $\beta \rightarrow \text{ss}$ , Greek  $A_i \rightarrow \alpha\iota$
- T (Turkic) Turkic  $I \rightarrow \text{ı}$ , because uppercase  $i$  in Turkic is  $\dot{I}$

Possible combinations: C + F, C + S, C + F + T, C + S + T

$\text{CaseFold}_{C+F}(\dot{I}) = ?$

$\text{CaseFold}_{C+F}(\dot{I}) = \text{i} \text{ ◌}$

# Segmentation (UAX #29)

---

ICU word break (UAX #29): rule-based segmentation using character properties.

和三聯齋の下北沢店で naïve な Thé Noir と phở と كَبَابَ を注文、oṅṅṅ と नमस्ते で先輩に挨拶した。8月10日、二人 幸 終。

NFKC + CaseFold (C + F)

令和 ㄨ 明治株式会社の下北沢店で naïve な thé noir と phở と كَبَابَ を注文、oṅṅṅ と नमस्ते で先輩に挨拶した。8月10日、二人 幸 終。

ICU Segmentation (UAX \#29)

令 和 ㄨ 明 治 株 式 会 社 の 下 北 沢 店 で naïve な thé noir と phở と كَبَابَ を 注 文 oṅṅṅ と नमस्ते で 先  
輩 に 挨 拶 し た 8 月 10 日 二 人 幸 終

# Arabic and Hebrew

---

Arabic and Hebrew use **diacritical marks** for vowels that are often omitted in everyday writing. The same word may appear with or without them:

- كِتَابٌ (with harakat) = كتاب (without) = “kitāb” (book)
- שְׁלוֹמִי (with niqqud) = שלום (without) = “shalom”

Input	Normalized	Rule
الله	الله	Alif wasla → Alif
مدرسة	مدرسه	Ta marbuta → Ha
ي / ك	ي / ك	Farsi → Arabic form
الله	الله / الله	Tatweel (kashida) removed
٢٠٢٦ / ٢٠٢٦	2026	Arabic / Persian digits → ASCII
ت	ت	Don't touch this

Not accurate for Uyghur Arabic alphabet

# More Than Normalization: Diacritic Folding

---

- NFKC Casefold normalizes `A p p l e` → `apple` and `Ω` (U+2126) → `ω` (U+03C9)
- But precomposed accented characters survive: `café` stays `café`.

For search, we want `cafe` to match `café`. Solution: **diacritic folding**.

---

**Algorithm:** NFD decompose → strip foldable combining marks → NFC recompose

`café`  $\xrightarrow{\text{NFD}}$  `cafe◌́`  $\xrightarrow{\text{strip}}$  `cafe`  $\xrightarrow{\text{NFC}}$  `cafe`

Naïve idea: remove all combining marks (ccc > 0)

# Combining Marks (CC) ≠ Diacritic Marks

---

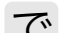
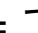
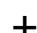

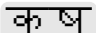
Language	CC Usage	CC in Daily Writing	CC Stripping Acceptable
European Languages (?)	Accents, Variants	Yes	Mostly Yes
Japanese	Consonant Variants	Yes	Yes/No
Arabic, Hebrew	Vowels	No	Yes
Brahmic scripts	Vowels	Yes	No

---

**Foldable ranges** (Latin/Greek/Cyrillic/Vietnamese accents):

- U+0300–036F: Combining Diacritical Marks
- U+1AB0–1AFF: Combining Diacritical Marks Extended
- U+1DC0–1DFF: Combining Diacritical Marks Supplement

**Preserved** (structurally significant):

- Japanese dakuten (U+3099):  =  +   — folding would destroy meaning
- Devanagari virama (U+094D):  — consonant conjuncts
- Arabic harakat, Hebrew niqqud — handled by [SemiticNormalizationFilter](#)

# Vietnamese: Diacritic Folding Unacceptable

---

Unmarked	Grave	Hook	Perispomeni	Acute	Dot
	`(U+0300)	´(U+0309)	˘(U+0303)	´(U+0301)	.(U+0323)
Uppercase letters					
A	À (U+00C0)	Ả (U+1EA2)	Ã (U+00C3)	Á (U+00C1)	Ạ (U+1EA0)
Ă (U+0102)	Ằ (U+1EB0)	Ẳ (U+1EB2)	Ẵ (U+1EB4)	Ẳ (U+1EAE)	Ẵ (U+1EB6)
Â (U+00C2)	Ằ (U+1EA6)	Ẳ (U+1EA8)	Ẵ (U+1EAA)	Ẳ (U+1EA4)	Ẵ (U+1EAC)

1

---

<sup>1</sup>[https://en.wikipedia.org/wiki/Vietnamese\\_alphabet](https://en.wikipedia.org/wiki/Vietnamese_alphabet)

# Solution: Triple-Field Indexing

---

Field	Content	Purpose
<code>folded_bigram</code>	diacritics folded + CJK bigrams	primary recall (broad match)
<code>unigram</code>	isolated Han unigrams	single-char CJK fallback
<code>diacritic</code>	only tokens <b>with</b> foldable diacritics, pre-fold form	sparse, exact-accent matching

---

Query routing:

- `cafe` (no diacritics)  $\rightarrow$   $\text{TermQuery}_{\text{bigram}}(\text{cafe})$  — matches both
- `café` (has diacritics)  $\rightarrow$   $\text{TermQuery}_{\text{bigram}}(\text{cafe}) \vee \text{TermQuery}_{\text{diacritic}}(\text{café})$ 
  - So `café` match is scored higher than `cafe` match

# Solution: Triple-Field Indexing (Cont'd)

令和三年八月十日、三軒茶屋の下北沢店で naive な Thé Noir と phở と كَبَابَ を注文、お礼と नमस्ते で先輩に挨拶した。8月10日、二人幸終。

NormalizingICUTokenizer + SemiticNormalizationFilter

令和三年八月十日、三軒茶屋の下北沢店で naive な the noir と pho と كَبَابَ を注文、お礼と नमस्ते で先輩に挨拶した。8月10日、二人幸終。

DiacriticFoldingFilter  
CJKBigramFilter

DiacriticOnlyFilter

HanOnlyFilter

令和三年八月十日、三軒茶屋の下北沢店で naive な the noir と pho と كَبَابَ を注文、お礼と नमस्ते で先輩に挨拶した。8月10日。

naive thé phở

令和三年八月十日、三軒茶屋の下北沢店で先輩に挨拶した。8月10日。

folded\_bigram

diacritic

unigram

# Scoring: Is BM25 the Silver Bullet?

---

Recall the BM25 formula:

$$\text{score} = \text{IDF} \cdot \frac{f \cdot (k_1 + 1)}{f + k_1 \cdot \left(1 - b + b \cdot \frac{|D|}{\text{avgdl}}\right)}$$

The  $b \cdot \frac{|D|}{\text{avgdl}}$  term penalizes long documents. Makes sense for article search (a 10-page match is less relevant than a short article match).

**But Telegram messages are short (5–50 tokens).**

Matching a long message is **not necessarily less relevant** than matching a short message.

# Scoring: Maybe IDF is Enough

---

**Our approach:** replace BM25 with `ConstScoreQuery`, scored by IDF only

$$\text{IDF}(\text{term}) = \ln\left(\frac{N - n(\text{term}) + 0.5}{n(\text{term}) + 0.5} + 1\right)$$

- IDF can be understood as information entropy (with numerical corrections)
- Rare terms get higher scores

Tradeoff: we lose term frequency signal entirely

# Project Status

---

- [tantivy-analyzer-icu](#) is almost complete, will be a standalone repo later
- [tg-searcher](#) still lacks some features

## Remarks:

- Prefer recall
- Achieve reasonable language-agnostic accuracy
- No language-specific features
  - Synonym
  - Stemmer
  - Stop word filter
  - OpenCC

## 今天信什么主义



**Internationalism(国际主义)**

**Thanks & Questions**