

# ICU4X 2.0: Next Level i18n

Shane Carr, Unicode Technology Workshop 2024



# Intro to ICU4X



**INTERNATIONALIZATION**

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18

**I18N**

# Examples of i18n Operations

Today's date in various locales and calendars:

- en: Oct 23, 2024 (Gregorian)
- de-CH: 23. Okt. 2024 (Gregorian)
- he: 21 5785 בתשרי (Hebrew)
- bn: ২১ রবিউস সানি, ১৪৪৬ যুগ (Islamic)
- zh: 2024年九月21 (Chinese)
- ja: 令和6年10月23日 (Japanese)

Breakpoints in Japanese text (words may be multiple characters wide):

中|ワ|況|写|イ|ノ|ナ|開|億|が  
|や|へ|者|43|俳|寺|式|7|沢|  
暮|ル|材|年|る|あん|移|酔|む  
ぎ|え|す|写|情|主|逃|69|引|  
ぎ|う|ぱ|何|81|昇|フ|ネ|イ|マ|  
本|因|ぱ|不|真|え|断|候|ら|。

# The i18n Stack

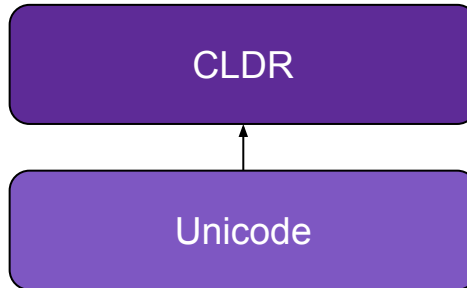
**Unicode:**  
Foundational algorithms, specification,  
character set



Unicode

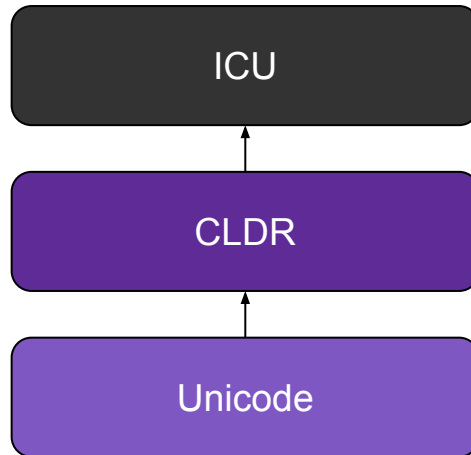
# The i18n Stack

**CLDR:**  
Data for hundreds of locales, specs for  
MessageFormat, person names,  
keyboards, ...



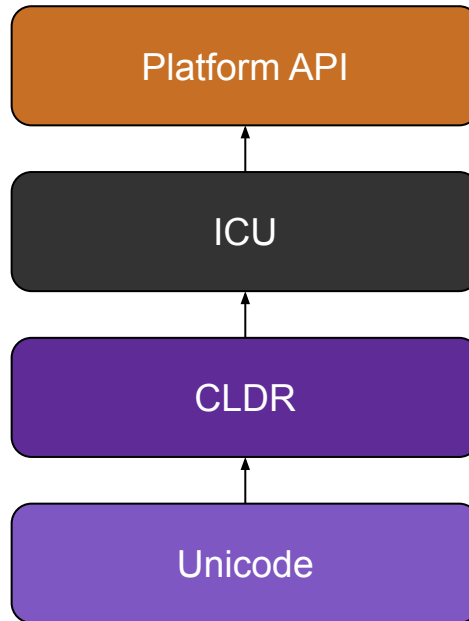
# The i18n Stack

ICU:  
Code that implements the algorithms in  
Unicode with the data in CLDR



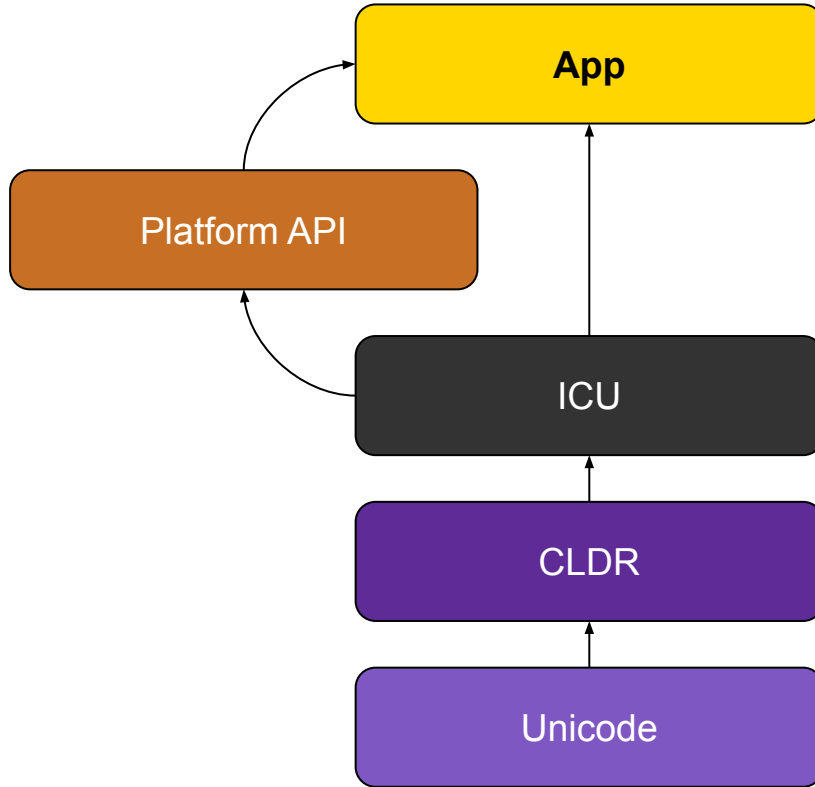
# The i18n Stack

**Platform API:**  
Making i18n easy to use from applications  
(example: ECMA-402, android.icu)

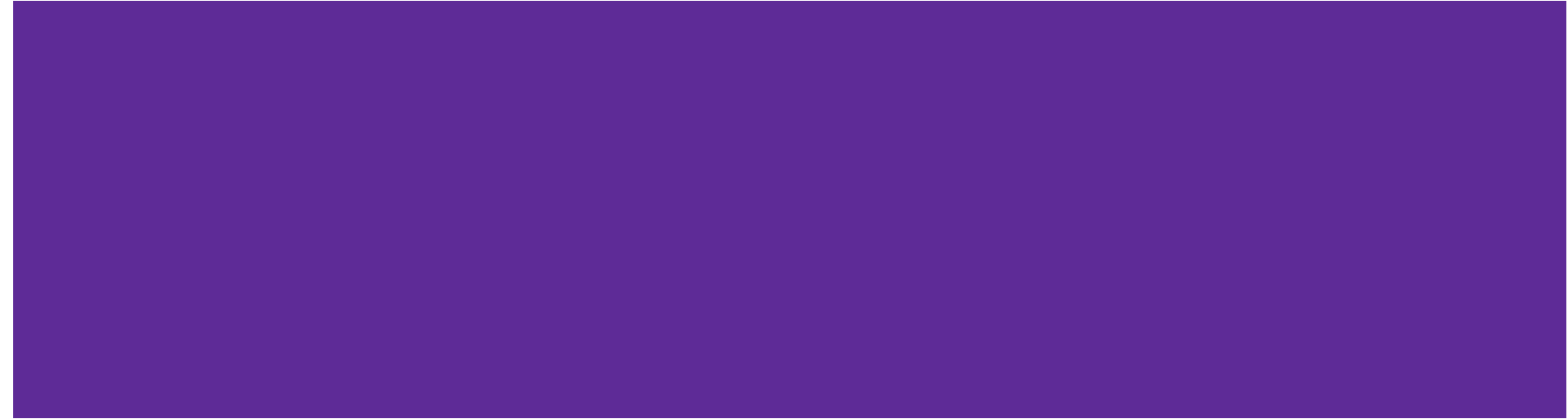




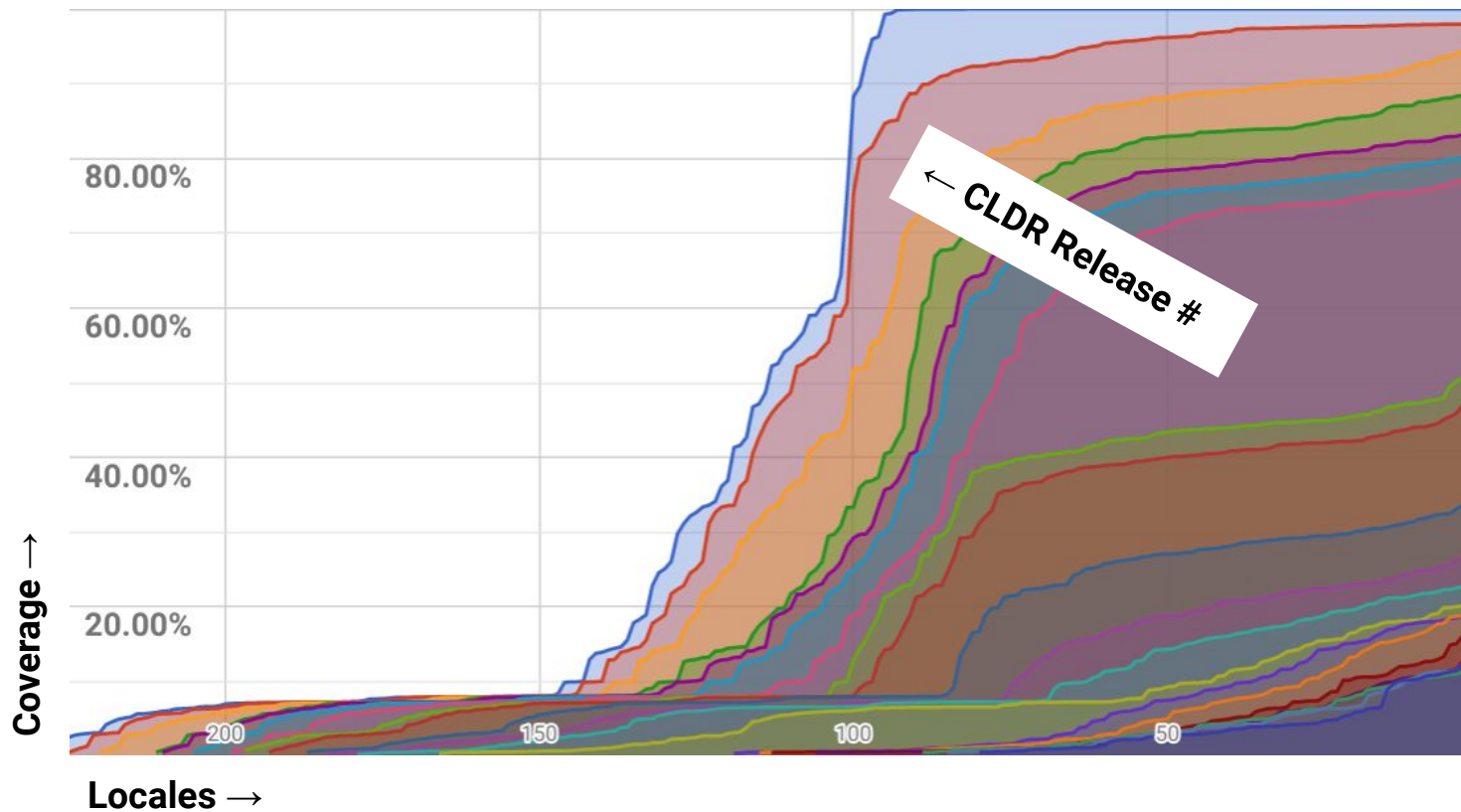
# The i18n Stack



# Challenges with Scaling $i18n$

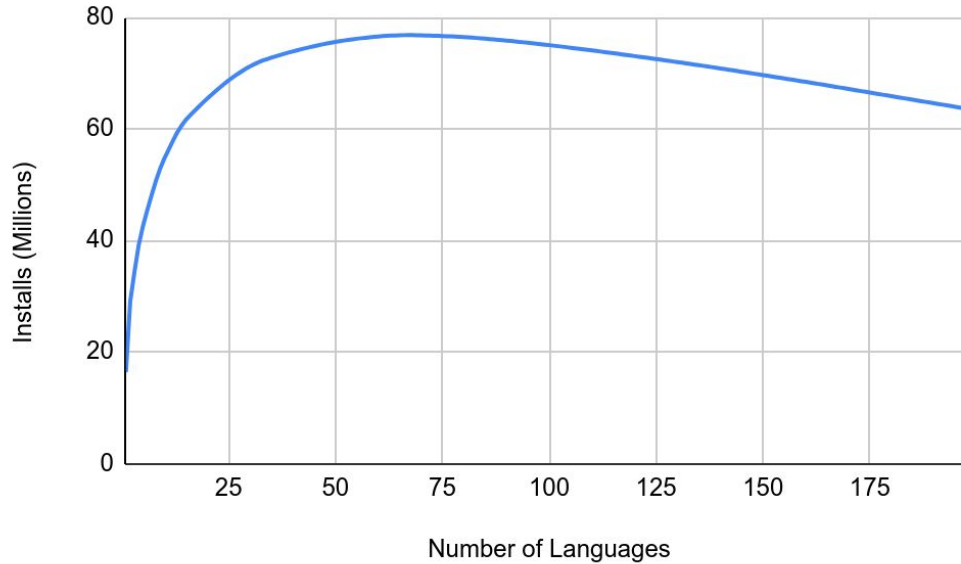


# Quadratic Growth of Data



# Adding more default languages has diminishing returns

Cumulative Installs by # of Default Languages



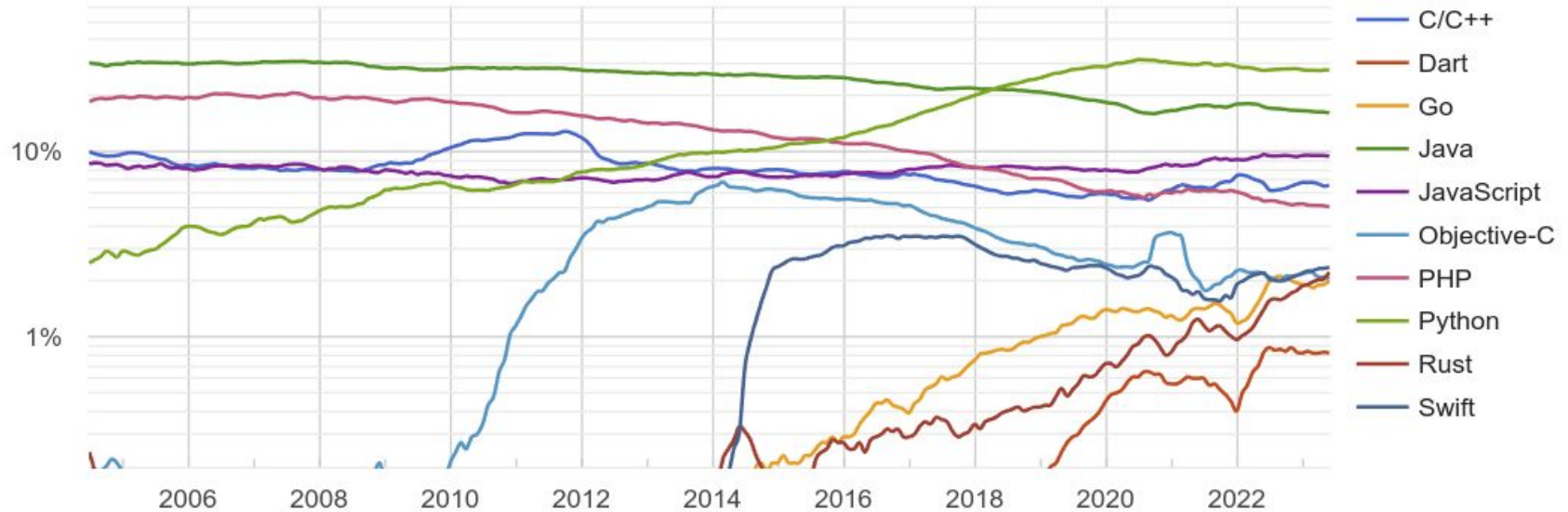
*If adding a language increases app size by 100 kB,  
reducing total downloads by 0.25%*

# Devices have different requirements



# Hot new programming languages every year

PYPL Popularity of Programming Language



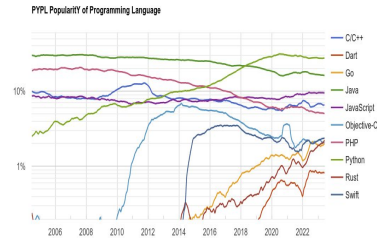
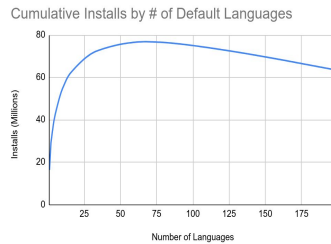
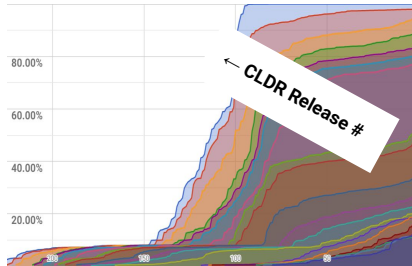
# i18n-as-a-service doesn't make sense

- Low latency requirements
- Data-heavy algorithms
- Privacy implications

中|ワ|況|写|イノナ|開|億|が|や|へ|者|43|  
俳|寺|式|7|沢|暮|ル|材|年|る|あん|移|酔|  
むぎ|え|す|写|情|主|逃|69|引|ぎ|う|ぱ|何|  
|81|昇|フ|ネ|イ|マ|本|因|ぱ|不|真|え|断|候|  
ら|。|吉|群|コ|ク|チ|賀|伝|ツ|エ|別|写|ユ|リ|  
キ|施|見|ら|カ|1|取|目|れ|づ|ぱ|ぞ|然|野|  
テ|倍|展|び|る|て|け|京|答|週|ネ|コ|入|整|  
料|勉|ら|げ|。

# Wish List for a Scalable i18n Library

- Pay for what you use and not what you don't use
- Add extra locales on demand
- Run on all types of devices
- Work in both today's and tomorrow's programming languages
- Run everything on-device



中|ワ|況|写|イ|ノ|ナ|開|徳|が|や|へ|者  
|43|俳|寺|式|7|沢|暮|ル|材|年|る|あ  
ん|移|醉|む|ぎ|え|す|写|情|主|逃|69|  
引|ぎ|う|ば|何|81|昇|フ|ネ|イ|マ|本|因  
ば|不|真|え|断|候|ら|。|吉|群|コ|ク|チ  
賀|伝|ツ|エ|別|写|ユ|リ|キ|施|見|ら|カ  
|1|取|目|れ|づ|ば|ぞ|然|野|テ|倍|展  
び|る|て|け|京|答|週|ネ|コ|入|整|料|  
勉|ら|げ|。



# Enter ICU4X

ICU4X is an internationalization library that is:

1. **Portable by design**
  - Runs on all types of devices with supported wrappers for multiple programming languages
2. **Lightweight**
  - Modular design, great for compile-time dead-code elimination and data slicing
3. **Secure**
  - Written in Rust, a memory-safe language

ICU4X-TC



**moz://a**

**Google**

**amazon**

# ICU4X, CLDR, and ICU

- ICU4C, ICU4J, ICU4X: Same data (CLDR), same spec (UTS 35), different code, similar behavior
- See the [conformance dashboard](#) for details on feature coverage and behavior differences

## ICU Data Driven Test Summary

Read more about the [Unicode Conformance project](#).

### Tests and platforms

Report generated: 2024-10-10 23:59:22

Executors verified: cpp\_71.1, cpp\_72.1, cpp\_73.1, cpp\_74.2, cpp\_75.1, dart\_web\_v20.1.0, icu4j\_73.2, icu4j\_74.2, icu4j\_75.1, node\_v14.18.3, node\_v14.21.3, node\_v18.7.0, node\_v18.14.2, node\_v20.1.0, node\_v21.6.0, node\_v22.9.0, rust\_1.3.2, rust\_1.4.0  
Tests verified: collation\_short, number\_fmt, plural\_rules, datetime\_fmt\_list\_fmt, lang\_names, rdt\_fmt, likely\_subtags, message\_fmt2

### Summary of all tests

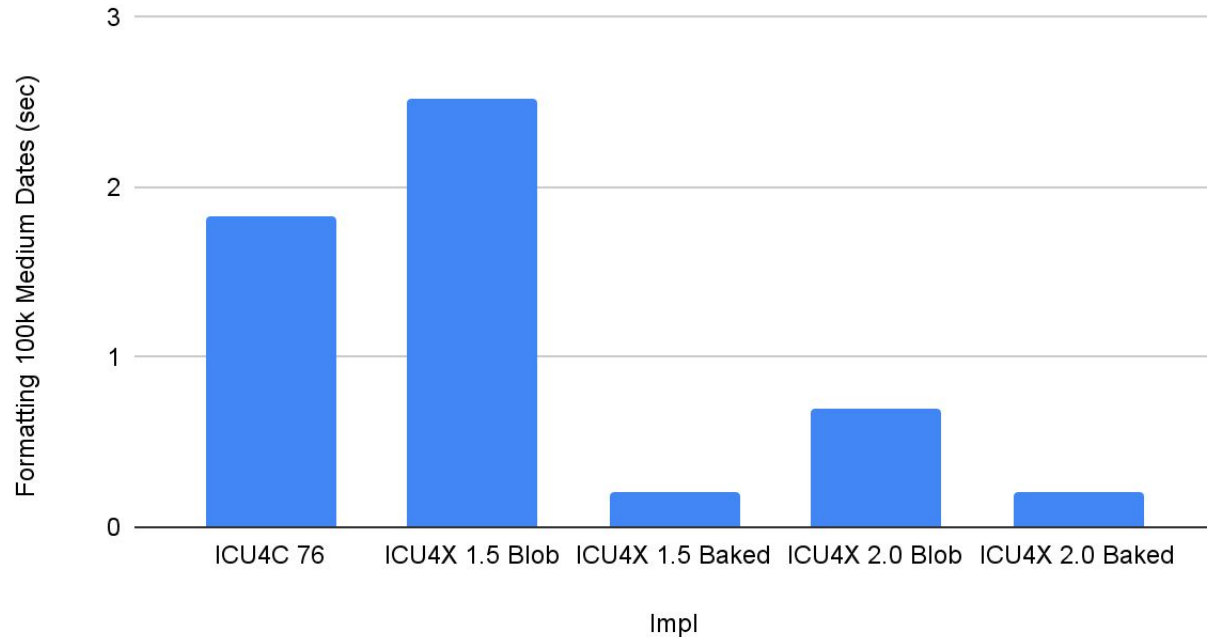


# Benchmarks



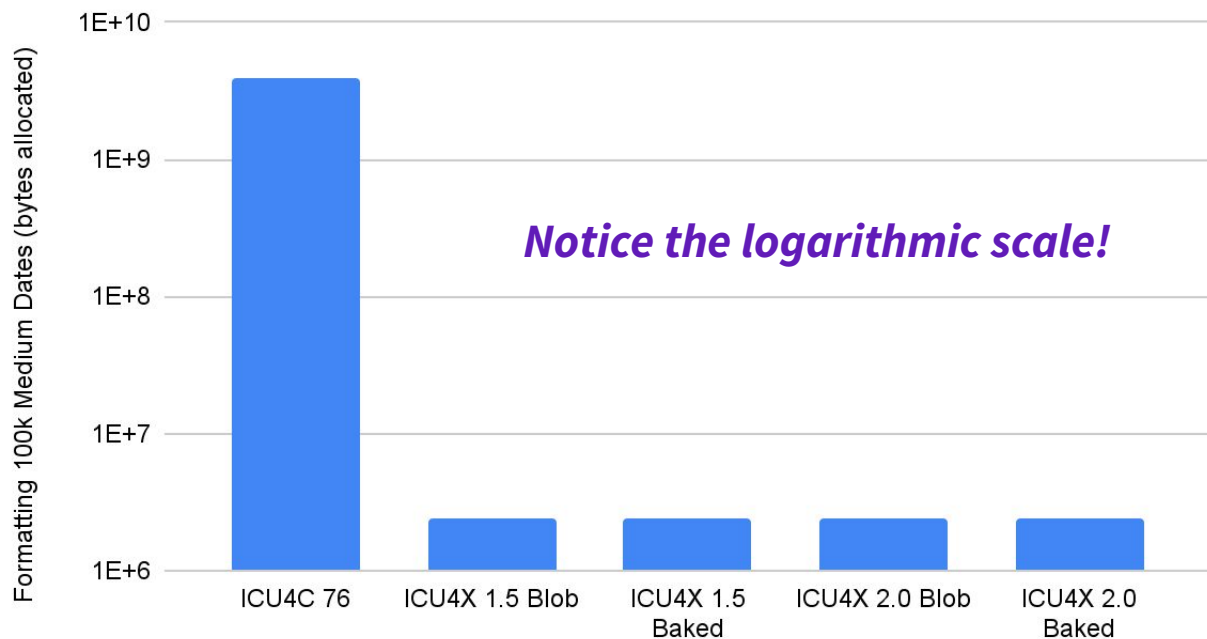
# ICU4C 76 vs ICU4X 1.5 vs ICU4X 2.0: Performance

DateTimeFormatter Performance (lower is better)



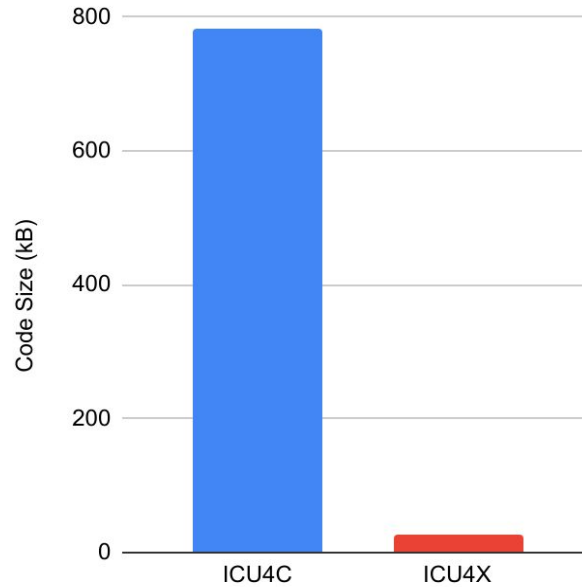
# ICU4C 76 vs ICU4X 1.5 vs ICU4X 2.0: Memory Use

DateTimeFormatter Memory Use (lower is better)



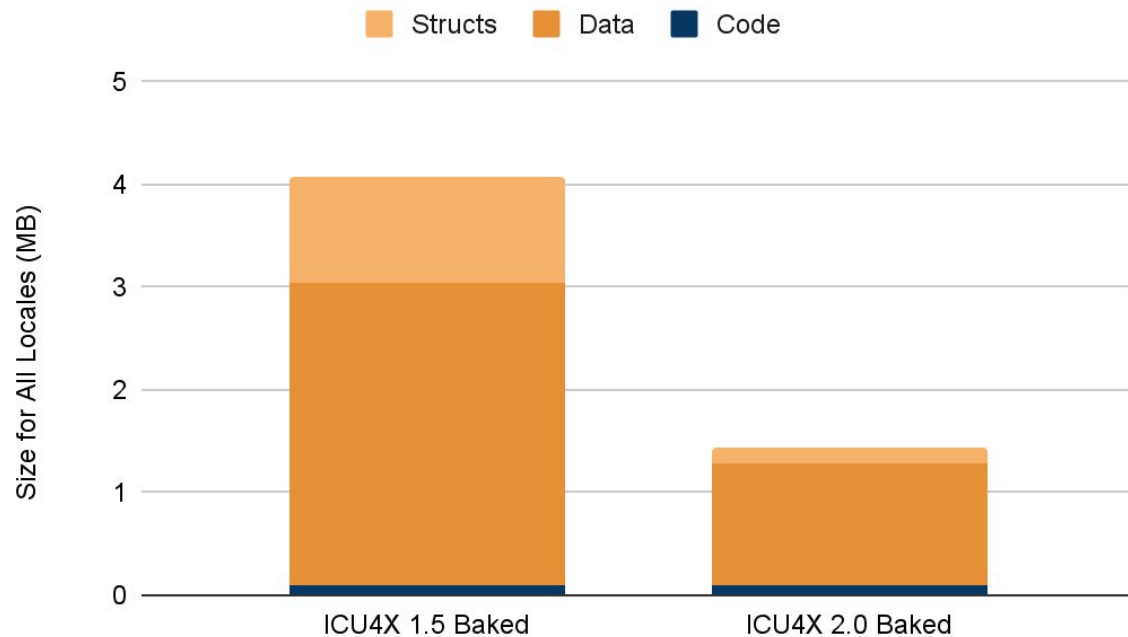
# ICU4C vs ICU4X: Binary size (code, no data)

Basic Number Format: Code Size



# ICU4X 1.5 vs ICU4X 2.0: Binary size

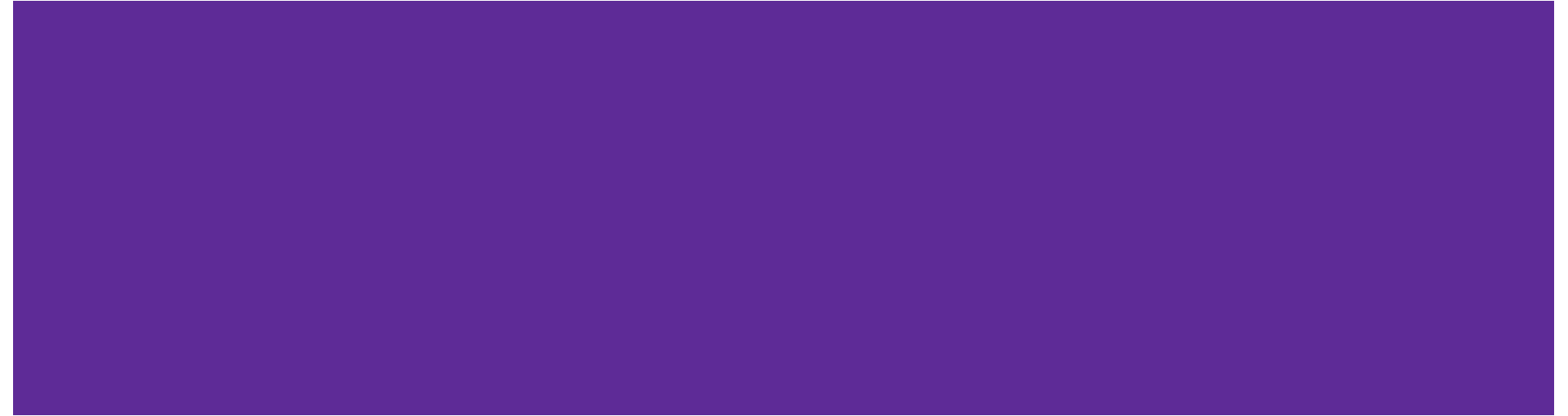
Binary Size: ICU4X 1.5 vs 2.0 (Multi-Component Demo)





**What other  
metrics would you  
like to see?**

# **New in ICU4X 2.0: Semantic Date and Time**



# Why Semantic Skeletons?

Correctness

Smaller and Faster

# Field Sets: Restrict to only valid combinations!

## Date

Day  
Weekday  
Day + Weekday  
Month + Day  
Month + Day + Weekday  
Year + Month + Day  
Year + Month + Day + Weekday

## Calendar Period

Month  
Year  
Year + Month

## Composite

Date + Time  
Date + Time Zone  
Date + Time + Time Zone  
Time + Time Zone

## Time

Hour  
Hour + Minute  
Hour + Minute + Second

## Time Zone

*This is a comprehensive list of CLDR 46 field sets!*

# Options

Option Name	Values	Fields
Length	Long, Medium, Short	Year, Month, Day, Weekday, Hour, Zone
Alignment	Inline, Column	Year, Month, Day, Hour
Year Style	Auto, Full, With Era	Year
Hour Cycle	Auto, H11, H12, ...	Hour
Fractional Second Digits	Auto, 0, 1, ..., 9	Second
Time Zone Style	Specific, Generic, ...	Zone

# Live Demo! (Of ICU4X Docs)

Note: This is 2.0 Alpha, and it could still change. Please send your feedback!

<https://unicode-org.github.io/icu4x/rustdoc/icu/datetime/fieldset/index.html>

# Technical Note: How This Achieves Smaller Data

Your formatter is *specific to your field set!* In many cases, hidden by Rust type inference.

- `DateTimeFormatter<YMD> ~= YearMonthDayFormatter`
- `DateTimeFormatter<HMZ> ~= HourMinuteSpecificTimeZoneFormatter`

The smaller types enable compile-time dead code elimination.

*Also!* The smaller data payloads are more easily deduplicated.

# More About ICU4X 2.0

*A deep dive to give appreciation to everything  
that makes us achieve our goals*



# Top 8 Contributors in ICU4X 1.4, 1.5, and 2.0



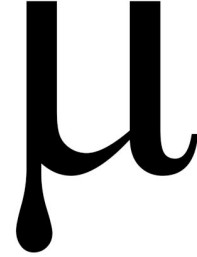
Robert Bastian  
396 Commits



Shane Carr  
319 Commits



Younies Mahmoud  
117 Commits



Manish Goregaokar  
108 Commits



José Julián Espina  
17 Commits



Kartavya  
17 Commits



Makoto Kato  
14 Commits



Henri Sivonen  
11 Commits

# And so many more!

- 10 Zibi Braniecki
- 8 Bruce Mitchener
- 8 Kevin Ness
- 8 Tyler K
- 5 Robin Leroy
- 4 Ashutosh Jha
- 4 Kevin
- 4 Steven R. Loomis
- 3 Blayne Marjama
- 3 Colin Rofls
- 3 Kartavya Vashishtha
- 3 Sean Burke
- 2 Adam Chalmers
- 2 Ben Beasley
- 2 CanadaHonk
- 2 Elango Cheran
- 2 Linus Groh

# New Features: A few highlights beyond datetime

- CLDR 46
- New experimental components: Duration Formatter, Unit Formatter, Currency Formatter, Person Names, Unit Converter
- Graduated casemap component to stable
- All new Diplomat WASM demo
- Add LocaleData parameter for word/sentence segmenter
- Support Unicode 15.1 for line segmenter
- Add PluralElements for algorithmic plural selection
- IXDTF Parsing: Date, Time, DateTime, ZonedDateTime try\_from\_str
- Time zone end-to-end display names without TZDB

# Infrastructure Upgrades: A few highlights

- Revamped data request architecture with locale and attributes
- All string input functions robustly handle ill-formed UTF-8
- Improvements to ICU4X datagen APIs and options (e.g. locale families)
- DryDataProvider for dry-run loading of data (e.g. supported locale queries)
- Add script fallback priority; improve language/region fallback priority
- Introduce borrowed variants of normalizer and collator structs
- Underlying zerovec data format optimized to pack variable-length items better
- Improvements in non-Gregorian calendrical calculations
- Use ZeroTrie in more places, reducing size and improving efficiency
- Fine-grained error enums

# FFI Improvements in ICU4X 2.0

- Dart now supported, on top of C, C++, JS, and TypeScript
- More usage of idiomatic features in languages that provide them, like getters, setters, constructors, and stringifiers
- APIs follow more idiomatic per-language naming conventions
- Namespacing in C++

# Status of ICU4X 2.0

- Done: All major feature work. Today is "alpha" quality.
- To-do before "beta": full API re-review; migration guide
- To-do before "final": FFI, finish time zone API, small function renames
- Stretch: type-safe locales, more fallback/datagen, nice-to-haves

ICU4X-TC ships when things are ready, not necessarily on a regular cadence like ICU-TC.  
(Hint: if you want something shipped sooner, become a contributor!)

**Truly exceptional  
engineering for 4+  
years has brought  
ICU4X from a  
concept to reality.**

# Plans for 2025

1. Compatibility with ICU4C/ICU4J for easier partial migration
2. Reduce duplication of work between ICU-TC and ICU4X-TC
3. ECMA-402 100% coverage, tested
4. Super-well polished for 3P adoption

What do *you* want to see in this list?



# Conclusion and Call to Action

Considering ICU4X? Let's chat!  
Help us improve the offering!

We love mentoring! Lots of  
projects both big and small!

Learn more on our GitHub:  
[github.com/unicode-org/icu4x](https://github.com/unicode-org/icu4x)

Docs: [icu4x.unicode.org](https://icu4x.unicode.org)

Version 2.0 Shipping Soon!

Subscribe to our low traffic  
mailing list! QR Code:

My contact info:

[shane@unicode.org](mailto:shane@unicode.org)

[linkedin.com/in/shanecarr](https://linkedin.com/in/shanecarr)

Twitter: @\_sffc

Mastodon, QR Code:

