



Citus Con: An Event for Postgres

Hosted by



# Localization on PostgreSQL

**Keisuke Takahashi**

Cloud Solution Architect (Data & Analytics),  
Microsoft

#CITUSCON | April 12-13, 2022



# Agenda

- **Introduction (5min)**
- **A number of different character sets (5min)**
  - Character Set Support
- **Locale-specific collation order (15min)**
  - Locale Support
  - Collation Support



# Agenda

- **Introduction**
- A number of different character sets
  - Character Set Support
- **Locale-specific collation order**
  - Locale Support
  - Collation Support



WILLKOMMEN

欢迎

स्वागत

BIENVENIDA

WELCOME

BIENVENUE ようこそ

환영하다 добро пожаловать

تَرْحِيبٌ

BEM-VINDO

歡迎

Chào mừng

# Who I am



- Keisuke Takahashi
- Cloud Solution Architect (Data&AI) at Microsoft Japan
- Open Source Software enthusiast
- English as a Second Language
- Learned Chinese for two years
- Have some Korean friends

# This presentation is for whom ...

- Already knows what character encoding is.
- Uses PostgreSQL for data that may include non-ASCII characters.
- Be interested in both PostgreSQL and foreign languages.

## Today's Goal

- Learn the basics of Character Sets, Locale and Collation in PostgreSQL
- Learn which locale / collation you should choose

# Azure Database for PostgreSQL deployment options



## Single Server

Fully-managed, single-node PostgreSQL database service with built-in HA

### Example use cases

- Transactional and operational analytics workloads
- Apps requiring JSON, geospatial support, or full-text search
- Cloud-native apps built with modern frameworks

## Flexible Server NEW

Maximum control for your database with a simplified developer experience

### Example use cases

- Support for a variety of workloads with a new simplified architecture
- High-performance apps utilizing zone co-location for low latency

## Hyperscale (Citus)

Worry-free PostgreSQL in the cloud with an architecture built to scale out

### Example use cases

- Scaling PostgreSQL multi-tenant, SaaS applications
- Real-time operational analytics
- Building high throughput transactional apps

# The benefits of Azure Database for PostgreSQL



## Innovate with open-source tools and extensions

Stay productive with full compatibility with community PostgreSQL and support for your favorite PostgreSQL extensions.



## Maximize performance with a fully managed Azure service

Focus on your application innovation, not database management. Enjoy AI-powered performance optimization and advanced security.

Single Server

Flexible Server NEW

Hyperscale (Citus)



## Ultimate control and flexibility for databases

Enjoy maximum control and flexibility with custom maintenance windows, zone redundant high availability and additional configuration parameters for fine grained database tuning.



## Build massively scalable PostgreSQL applications

Scale with ease to hundreds of nodes, with no app rewrites. Save time by running transactions and analytics in one database and avoid the costs of manual sharding.

# Agenda

- Introduction
- A number of different character sets
  - Character Set Support
- Locale-specific collation order
  - Locale Support
  - Collation Support



# Default settings on Azure PostgreSQL

```
SELECT name, setting FROM pg_settings WHERE name IN  
('server_encoding', 'dynamic_shared_memory_type')
```

## Single Server

```
postgres=> SELECT name, setting FROM pg_settings WHERE name IN ('server_encoding', 'dynamic_shared_memory_type');  
          name           | setting  
-----+-----  
dynamic_shared_memory_type | windows  
server_encoding            | UTF8
```

## Flexible Server

```
postgres=> SELECT name, setting FROM pg_settings WHERE name IN ('server_encoding', 'dynamic_shared_memory_type');  
          name           | setting  
-----+-----  
dynamic_shared_memory_type | posix  
server_encoding            | UTF8
```

## Hyperscale (Citus)

```
citus=> SELECT name, setting FROM pg_settings WHERE name IN ('server_encoding', 'dynamic_shared_memory_type');  
          name           | setting  
-----+-----  
dynamic_shared_memory_type | posix  
server_encoding            | UTF8
```

# Default settings on Azure PostgreSQL

```
SELECT datname, datcollate, datctype FROM pg_database  
WHERE datname IN ('template1','template0')
```

## Single Server

```
postgres=> SELECT datname, datcollate, datctype FROM pg_database WHERE datname IN ('template1','template0');  
      datname |      datcollate |      datctype  
-----+-----+-----  
template1 | English_United States.1252 | English_United States.1252  
template0 | English_United States.1252 | English_United States.1252
```

## Flexible Server

```
postgres=> SELECT datname, datcollate, datctype FROM pg_database WHERE datname IN ('template1','template0');  
      datname |      datcollate |      datctype  
-----+-----+-----  
template1 | en_US.utf8 | en_US.utf8  
template0 | en_US.utf8 | en_US.utf8
```

## Hyperscale (Citus)

```
citus=> SELECT datname, datcollate, datctype FROM pg_database WHERE datname IN ('template1','template0');  
      datname |      datcollate |      datctype  
-----+-----+-----  
template1 | en_US.UTF-8 | en_US.UTF-8  
template0 | en_US.UTF-8 | en_US.UTF-8
```

# Supported Character Sets (Server-side) #1-10

Name	Language	ICU?	Bytes/Char
UTF8	<i>all</i>	Yes	1–4
SQL_ASCII	<i>any</i>	No	1
WIN1256	Arabic	Yes	1
LATIN7	Baltic	Yes	1
WIN1257	Baltic	Yes	1
LATIN8	Celtic	Yes	1
LATIN2	Central European	Yes	1
WIN1250	Central European	Yes	1
WIN866	Cyrillic	Yes	1
WIN1251	Cyrillic	Yes	1

# Supported Character Sets (Server-side) #11-20

Name	Language	ICU?	Bytes/Char
KOI8R	Cyrillic (Russian)	Yes	1
KOI8U	Cyrillic (Ukrainian)	Yes	1
WIN1253	Greek	Yes	1
WIN1255	Hebrew	Yes	1
EUC_JP	Japanese	Yes	1–3
EUC_JIS_2004	Japanese	No	1–3
EUC_KR	Korean	Yes	1–3
ISO_8859_6	Latin/Arabic	Yes	1
ISO_8859_5	Latin/Cyrillic	Yes	1
ISO_8859_7	Latin/Greek	Yes	1

# Supported Character Sets (Server-side) #21-30

Name	Language	ICU?	Bytes/Char
ISO_8859_8	Latin/Hebrew	Yes	1
LATIN9	LATIN1 with Euro and accents	Yes	1
MULE_INTERNAL	Multilingual Emacs	No	1–4
LATIN6	Nordic	Yes	1
LATIN4	North European	Yes	1
LATIN10	Romanian	No	1
EUC_CN	Simplified Chinese	Yes	1–3
LATIN3	South European	Yes	1
WIN874	Thai	No	1
EUC_TW	Traditional Chinese, Taiwanese	Yes	1–3

# Supported Character Sets (Server-side) #31-35

Name	Language	ICU?	Bytes/Char
LATIN5	Turkish	Yes	1
WIN1254	Turkish	Yes	1
WIN1258	Vietnamese	Yes	1
LATIN1	Western European	Yes	1
WIN1252	Western European	Yes	1

# Setting the Character Set

\* This operation is not available on Azure PostgreSQL Hyperscale(Citus) because Citus is a DB level extension.

```
CREATE DATABASE new_database TEMPLATE template0 ENCODING  
'encoding' LC_CTYPE 'collation.encoding'
```

Example:

```
CREATE DATABASE my_db TEMPLATE template0 ENCODING 'UTF-8'  
LC_CTYPE 'ja_JP.UTF-8';
```

Result:

```
postgres=> SELECT datname, datcollate, datctype FROM pg_database WHERE datname = 'my_db';  
datname | datcollate | datctype  
-----+-----+-----  
my_db | en_US.utf8 | ja_JP.UTF-8
```

# Setting the Character Set

\* This operation is not available on Azure PostgreSQL Hyperscale(Citus) because Citus is a DB level extension.

```
CREATE DATABASE new_database TEMPLATE template0 ENCODING  
'encoding' LC_CTYPE 'collation.encoding'
```

Example:

```
CREATE DATABASE my_db TEMPLATE template0 ENCODING 'UTF-8'  
LC_CTYPE 'ja_JP.UTF-8';
```

Alternative:

```
$ createdb my_db --locale=japanese --template=template0  
$ initdb --encoding=UTF=8
```

# Agenda

- Introduction
- A number of different character sets
  - Character Set Support
- **Locale-specific collation order**
  - Locale Support
  - Collation Support



# How to check your locale settings

```
SELECT name, setting, context FROM pg_settings WHERE name LIKE 'lc%'
```

Result on Azure Database for PostgreSQL Hyperscale (Citus):

```
citus=> SELECT name, setting, context FROM pg_settings WHERE name LIKE 'lc%';
      name      |   setting   |    context
-----+-----+-----
 lc_collate | en_US.UTF-8 | internal
 lc_ctype   | en_US.UTF-8 | internal
 lc_messages | en_US.utf8  | superuser
 lc_monetary | en_US.utf8  | user
 lc_numeric  | en_US.utf8  | user
 lc_time     | en_US.utf8  | user
```

# Parameters

Name	Description	Note
LC_COLLATE	String sort order	
LC_CTYPE	Character classification (What is a letter? Its upper-case equivalent?)	
LC_MESSAGES	Language of messages	
LC_MONETARY	Formatting of currency amounts	This can also be configured with Azure Portal
LC_NUMERIC	Formatting of numbers	This can also be configured with Azure Portal
LC_TIME	Formatting of dates and times	

# Default settings on Azure PostgreSQL

Single Server

name	setting	context
lc_collate	English_United States.1252	internal
lc_ctype	English_United States.1252	internal
lc_messages	English_United States.1252	superuser
lc_monetary	English_United States.1252	user
lc_numeric	English_United States.1252	user
lc_time	English_United States.1252	user

Flexible Server

name	setting	context
lc_collate	en_US.utf8	internal
lc_ctype	en_US.utf8	internal
lc_messages	en_US.utf8	superuser
lc_monetary	en_US.utf-8	user
lc_numeric	en_US.utf-8	user
lc_time	en_US.utf8	user

Hyperscale  
(Citus)

name	setting	context
lc_collate	en_US.UTF-8	internal
lc_ctype	en_US.UTF-8	internal
lc_messages	en_US.utf8	superuser
lc_monetary	en_US.utf8	user
lc_numeric	en_US.utf8	user
lc_time	en_US.utf8	user

# How to set your locale

\* This operation is not available on Azure PostgreSQL Hyperscale(Citus) because Citus is a DB level extension.

```
CREATE DATABASE new_database TEMPLATE template0 ENCODING  
'encoding' LC_COLLATE 'collation.encoding' LC_CTYPE  
'collation.encoding'
```

Example:

```
CREATE DATABASE my_db TEMPLATE template0 ENCODING 'UTF-8'  
LC_COLLATE 'ja_JP.UTF-8' LC_CTYPE 'ja_JP.UTF-8';
```

Result:

```
postgres=> SELECT datname, datcollate, datctype FROM pg_database WHERE datname = 'my_db';  
datname | datcollate | datctype  
-----+-----+-----  
my_db | ja_JP.UTF-8 | ja_JP.UTF-8
```

# How to set your locale

\* This operation is not available on Azure PostgreSQL Hyperscale(Citus) because Citus is a DB level extension.

```
CREATE DATABASE new_database TEMPLATE template0 ENCODING  
'encoding' LC_COLLATE 'collation.encoding' LC_CTYPE  
'collation.encoding'
```

Example:

```
CREATE DATABASE my_db TEMPLATE template0 ENCODING 'UTF-8'  
LC_COLLATE 'ja_JP.UTF-8' LC_CTYPE 'ja_JP.UTF-8';
```

Alternative:

```
$ createdb ja --locale=japanese --template=template0  
$ initdb --encoding=UTF=8 --locale=ja_JP.UTF=8
```

# How to get Predefined Collations

```
SELECT * FROM pg_collation
```

Example:

```
SELECT collname, collprovider FROM pg_collation WHERE  
collname LIKE 'ja%';
```

Alternative:

```
\dos+
```

# Predefined "ja" Collations on Azure PostgreSQL

```
SELECT collname, collprovider FROM pg_collation WHERE collname LIKE 'ja%';
```

## Single Server

```
postgres=> SELECT collname, collprovider FROM pg_collation WHERE collname LIKE 'ja%';
      collname | collprovider
-----+-----
    ja          | c
  ja.utf8      | c
   ja_JP        | c
 ja_JP.utf8    | c
(4 rows)
```

## Flexible Server

```
postgres=> SELECT collname, collprovider FROM pg_collation WHERE collname LIKE 'ja%';
      collname | collprovider
-----+-----
 ja_JP.eucjp  | c
 ja_JP.utf8   | c
  ja_JP        | c
   ja_JP       | c
  ja-x-icu    | i
 ja-JP-x-icu  | i
(6 rows)
```

## Hyperscale (Citus)

```
citus=> SELECT collname, collprovider FROM pg_collation WHERE collname LIKE 'ja%';
      collname | collprovider
-----+-----
   ja_JP     | c
 ja_JP.eucjp | c
  ja_JP.ujis | c
  ja_JP.utf8 | c
  japanese   | c
japanese.euc | c
   ja_JP     | c
  ja-x-icu   | i
ja-JP-x-icu  | i
(9 rows)
```

# Predefined "zh" Collations on Azure PostgreSQL

## Single Server

```
postgres=> SELECT collname, collprovider
  FROM pg_collation WHERE collname LIKE 'zh%';
   collname | collprovider
-----+-----
    zh          | c
  zh.utf8      | c
  zh_CHS       | c
zh_CHS.utf8    | c
  zh_CHT       | c
zh_CHT.utf8    | c
  zh_CN        | c
zh_CN.utf8     | c
  zh_HK        | c
zh_HK.utf8     | c
  zh_Hans      | c
zh_Hans.utf8   | c
  zh_Hans_HK   | c
zh_Hans_HK.utf8| c
  zh_Hans_M0   | c
zh_Hans_M0.utf8| c
  zh_Hant      | c
zh_Hant.utf8   | c
  zh_M0        | c
zh_M0.utf8     | c
  zh_SG        | c
zh_SG.utf8     | c
  zh_TW        | c
zh_TW.utf8     | c
(24 rows)
```

## Flexible Server

```
postgres=> SELECT collname, collprovider
  FROM pg_collation WHERE collname LIKE 'zh%';
   collname | collprovider
-----+-----
  zh_CN      | c
zh_CN.utf8  | c
  zh_HK.utf8 | c
  zh_SG      | c
zh_SG.utf8  | c
  zh_TW.euctw| c
zh_TW.utf8  | c
  zh_CN      | c
  zh_HK      | c
  zh_SG      | c
  zh_TW      | c
  zh_TW      | c
  zh_x_icu   | i
zh-Hans-x-icu| i
  zh-Hans-CN-x-icu| i
zh-Hans-HK-x-icu| i
zh-Hans-M0-x-icu| i
zh-Hans-SG-x-icu| i
  zh_Hant-x-icu| i
zh-Hant-HK-x-icu| i
zh-Hant-M0-x-icu| i
zh-Hant-TW-x-icu| i
(22 rows)
```

## Hyperscale (Citus)

```
citus=> SELECT collname, collprovider
  FROM pg_collation WHERE collname LIKE 'zh%';
   collname | collprovider
-----+-----
  zh_CN      | c
zh_CN.gb2312| c
  zh_CN.utf8 | c
  zh_HK.utf8 | c
  zh_SG      | c
zh_SG.gb2312| c
  zh_SG.utf8 | c
  zh_TW.euctw| c
zh_TW.utf8  | c
  zh_TW      | c
  zh_TW      | c
  zh_CN      | c
  zh_HK      | c
  zh_SG      | c
  zh_TW      | c
  zh_TW      | c
  zh_x_icu   | i
zh-Hans-x-icu| i
  zh-Hans-CN-x-icu| i
zh-Hans-HK-x-icu| i
zh-Hans-M0-x-icu| i
zh-Hans-SG-x-icu| i
  zh_Hant-x-icu| i
zh-Hant-HK-x-icu| i
zh-Hant-M0-x-icu| i
zh-Hant-TW-x-icu| i
(24 rows)
```

## libc vs. ICU

- **libc**
  - Old-school locale/collation provider of PostgreSQL
- **ICU**
  - International Components for Unicode
  - Have been supported since PostgreSQL 10
  - User-defined locale/collation support

## Create ICU Database

\* This operation is not available on Azure PostgreSQL Hyperscale(Citus) because Citus is a DB level extension.

```
CREATE DATABASE new_database ENCODING 'UTF-8' TEMPLATE  
template0 LC_COLLATE 'collation' LC_CTYPE 'collation'
```

Example:

```
CREATE DATABASE my_icu_db ENCODING 'UTF-8' LC_COLLATE  
'ja-x-icu' LC_CTYPE 'ja-x-icu';
```

## Create ICU Template DB

\* This operation is not available on Azure PostgreSQL Hyperscale(Citus) because Citus is a DB level extension.

```
CREATE DATABASE new_template ENCODING 'UTF-8' TEMPLATE  
template0 LC_COLLATE 'collation' LC_CTYPE 'collation'
```

```
CREATE DATABASE new_database TEMPLATE new_template
```

Example:

```
CREATE DATABASE template_ja_x_icu ENCODING 'UTF-8'  
TEMPLATE template0 LC_COLLATE 'ja-x-icu' LC_CTYPE 'ja-x-  
icu';  
CREATE DATABASE my_db TEMPLATE template_ja_x_icu;
```

# Set Collation for the Column

```
CREATE TABLE tablename (
    colname1 type COLLATE "collation",
    colname2 type COLLATE "collation",
    ...
)
```

Example:

```
CREATE TABLE users (
    id SERIAL PRIMARY KEY,
    name VARCHAR(255) COLLATE "ja-x-icu");
```

# Specify Collation on SELECT

```
SELECT * FROM table ORDER BY column COLLATE "collation"
```

Example:

```
SELECT * FROM mytable ORDER BY mycol COLLATE "ja-x-icu";
```

# Sorting Experiments

- Environment
  - Azure Database for PostgreSQL Hyperscale (Citus)
    - PostgreSQL 14
    - Citus 10.2
    - Basic tier
    - West US 2

# Sorting Experiments

- 10 collations
  - English
    - en\_US.utf8, en-US-x-icu
  - Japanese
    - ja\_JP.utf8, ja-x-icu, ja-JP-x-icu
  - Chinese
    - zh\_CN.utf8, zh-x-icu, zh-Hans-x-icu, zh-Hant-x-icu

# Sorting Experiments

- 165 unicode characters / combinations
  - Latin
  - Japanese (Hiragana, Katakana, Kanji, Simbol, etc.)
  - Chinese (Traditional, Simplified)
  - Hangul
  - Chữ Nôm
  - Arabic
  - Cyrillic

# Tested Characters

Character	UTF-8 (hex)	Description
a	0x61	Latin Small Letter a
z	0x7A	Latin Small Letter z
A	0x41	Latin Capital Letter A
Z	0x5A	Latin Capital Letter Z
a	0xEF 0xBD 0x81	Fullwidth Latin Capital Letter a
z	0xEF 0xBD 0x9A	Fullwidth Latin Capital Letter z
A	0xEF 0xBC 0xA1	Fullwidth Latin Capital Letter A
Z	0xEF 0xBC 0xBA	Fullwidth Latin Capital Letter Z

# Tested Characters

Character	UTF-8 (hex)	Description
0	0x30	ASCII Digit Zero
9	0x39	ASCII Digit Nine
0	0xEF 0xBC 0x90	Fullwidth Digit Zero
9	0xEF 0xBC 0x99	Fullwidth Digit Nine
⓪	0xE2 0x93 0xAA	Circled Digit Zero
⑤⓪	0xE3 0x8A 0xBF	Circled Number Fifty
⓪	0xE3 0x89 0x8C	Circled Number Fifty on Black Square

# Tested Characters

Character	UTF-8 (hex)	Description
①	0xE2 0x91 0xA0	Circled Digit One
①	0xE2 0x9D 0xB6	Dingbat Negative Circled Digit One
I	0xE2 0x85 0xA0	Roman Numeral One
i	0xE2 0x85 0xB0	Small Roman Numeral One
(1)	0xE2 0x91 0xB4	Parenthesized Digit One
(1)	0x28 0x31 0x29	Digit One in Parenthesisses
—	0xE4 0xB8 0x80	CJK Ideograph, First
壹	0xE5 0xA3 0xB1	Number One

# Tested Characters

Character	UTF-8 (hex)	Description
⑨	0xE2 0x91 0xA8	Circled Digit Nine
❾	0xE2 0x9D 0xBE	Dingbat Negative Circled Digit Nine
IX	0xE2 0x85 0xA8	Roman Numeral Nine
ix	0xE2 0x85 0xB8	Small Roman Numeral Nine
(9)	0xE2 0x91 0xBC	Parenthesized Digit Nine
(9)	0x28 0x39 0x29	Digit Nine in Parenthesisses
九	0xE4 0xB9 0x9D	Nine
9.	0xE2 0x92 0x90	Digit Nine Full Stop
9.	0x91 0x2E	Digit Nine and Full Stop

# Tested Characters

Character	UTF-8 (hex)	Description
VII	0xE2 0x85 0xA6	Roman Numeral Seven
vii	0xE2 0x85 0xB6	Small Roman Numeral Seven
VII	0x56 0x49 0x49	Latin Capital Letter V, I, I
vii	0x76 0x69 0x69	Latin Small Letter v, i, i
XII	0xE2 0x85 0xAB	Roman Numeral Twelve
xii	0xE2 0x85 0xBB	Small Roman Numeral Twelve
XII	0x58 0x49 0x49	Latin Capital Letter X, I, I
xii	0x78 0x69 0x69	Latin Small Letter x, i, i

# Tested Characters

Character	UTF-8 (hex)	Description
α	0xCE 0xB1	Greek Small Letter Alpha
ω	0xCF 0x89	Greek Small Letter Omega
Α	0xCE 0x91	Greek Capital Letter Alpha
Ω	0xCE 0xA9	Greek Capital Letter Omega

# Tested Characters

Character	UTF-8 (hex)	Description
(a)	0x28 0x61 0x29	"a" in parenthesisses
(z)	0x28 0x7A 0x29	"z" in parenthesisses
{a}	0x7B 0x61 0x7D	"a" in Curly Brackets
[a]	0x5B 0x61 0x5D	"a" in Square Brackets
(a)	0xEF 0xBC 0x88 0x61 0xEF 0xBC 0x89	"a" in Fullwidth parenthesisses
(z)	0xEF 0xBC 0x88 0x7A 0xEF 0xBC 0x89	"z" in Fullwidth parenthesisses

# Tested Characters

Character	UTF-8 (hex)	Description
{a}	0xEF 0xBD 0x9B 0x61 0xEF 0xBD 0x9D	"a" in Fullwidth Curly Brackets
[a]	0xEF 0xBC 0xBB 0x61 0xEF 0xBC 0xBD	"a" in Fullwidth Square Brackets
【a】	0xE3 0x80 0x90 0x61 0xE3 0x80 0x91	"a" in Black Lenticular Brackets

# Tested Characters

Character	UTF-8 (hex)	Description
a	0x20 0x61	"a" after a space
a	0xE3 0x80 0x80 0x61	"a" after an Ideographic space
_a	0x5F 0x61	"a" after an underscore
__a	0xEF 0xBC 0xBF 0x61	"a" after a Fullwidth Low Line
/a	0x2F 0x61	"a" after a Solidus
／a	0xEF 0xBC 0x8F 0x61	"a" after a Fullwidth Solidus

# Tested Characters

Character	UTF-8 (hex)	Description
あ	0xE3 0x81 0x82	Hiragana Letter A
ア	0xE3 0x82 0xA2	Katakana Letter A
ア	0xEF 0xBD 0xB1	Halfwidth Katakana Letter A
ん	0xE3 0x82 0x93	Hiragana Letter N
ン	0xE3 0x83 0xB3	Katakana Letter N
ン	0xEF 0xBE 0x9D	Halfwidth Katakana Letter N

# Tested Characters

Character	UTF-8 (hex)	Description
は	0xE3 0x81 0xAF	Hiragana Letter Ha
ば	0xE3 0x81 0xB0	Hiragana Letter Ba
は゛	0xE3 0x81 0xB0 0xE3 0x82 0x9B	Hiragana Letter Ha and Voiced Sound Mark
ぱ	0xE3 0x81 0xB1	Hiragana Letter Pa
は゜	0xE3 0x81 0xB0 0xE3 0x82 0x9C	Hiragana Letter Ha and Semi-voiced Sound Mark

# Tested Characters

Character	UTF-8 (hex)	Description
ハ	0xE3 0x83 0x8F	Katakana Letter Ha
バ	0xE3 0x83 0x90	Katakana Letter Ba
ハ゛	0xE3 0x83 0x8F 0xE3 0x82 0x9B	Katakana Letter Ha and Voiced Sound Mark
パ	0xE3 0x83 0x91	Katakana Letter Pa
ハ	0xEF 0xBE 0x8A	Halfwidth Katakana Letter Ha
バ	0xEF 0xBE 0x8A 0xEF 0xBE 0x9E	Halfwidth Katakana Letter Ha and Halfwidth Katakana Voiced Sound Mark
ハ゛	0xEF 0xBE 0x8A 0xEF 0xBE 0x9F	Halfwidth Katakana Letter Ha and Halfwidth Katakana Semi-voiced Sound Mark

# Tested Characters

Character	UTF-8 (hex)	Description
あつ	0xE3 0x81 0x82 0xE3 0x81 0xA4	Hiragana Letter A and Tu
あつ	0xE3 0x81 0x82 0xE3 0x81 0xA3	Hiragana Letter A and Small Tu
ああ	0xE3 0x81 0x82 0xE3 0x81 0x82	Hiragana Letter A and A
ああ	0xE3 0x81 0x82 0xE3 0x81 0x81	Hiragana Letter A and Small A
あゝ	0xE3 0x81 0x82 0xE3 0x82 0x9D	Hiragana Letter A and Hiragana Iteration Mark

# Tested Characters

Character	UTF-8 (hex)	Description
あー	0xE3 0x81 0x82 0xE3 0x83 0xBC	Hiragana Letter A and Prolonged Sound Mark
あ-	0xE3 0x81 0x82 0x2D	Hiragana Letter A and Hyphen/Minus
あ～	0xE3 0x81 0x82 0xE3 0x80 0x9C	Hiragana Letter A and Wave Dash (macOS)
あ～	0xE3 0x81 0x82 0xEF 0xBD 0x9E	Hiragana Letter A and Fullwidth Tilde (Windows)
あ~	0xE3 0x81 0x82 0x7E	Hiragana Letter A and Tilde

# Tested Characters

Character	UTF-8 (hex)	Description
あ...	0xE3 0x81 0x82	Hiragana Letter A and Horizontal Ellipsis
	0xE2 0x80 0xA6	
あ`	0xE3 0x81 0x82	Hiragana Letter A and Voiced Sound Mark
	0xE3 0x82 0x9B	

# Tested Characters

Character	UTF-8 (hex)	Description
キロ	0xE3 0x82 0xAD	Katakana Letter Ki and Ro
	0xE3 0x83 0xAD	
ヰロ	0xEF 0xBD 0xB7	Halfwidth Katakana Letter Ki and Ro
	0xEF 0xBE 0x9B	
ヰロ	0xE3 0x8C 0x94	Square Kiro
km	0x6B 0x6D	Latin Small Letter k and m
k m	0xEF 0xBD 0x8B	Fullwidth Latin Small Letter k and m
	0xEF 0xBD 0x8D	
km	0xE3 0x8E 0x9E	Square km

# Tested Characters

Character	UTF-8 (hex)	Description
キロメートル	0xE3 0x8C 0x96	Square Kiromeetoru
糸	0xE7 0xB2 0x81	Kilometre (Japanese)
cm	0x63 0x6D	Latin Small Letter c and m
c m	0xEF 0xBD 0x83 0xEF 0xBD 0x8D	Fullwidth Latin Small Letter c and m
cm	0xE3 0x8E 0x9D	Square cm
センチ	0xE3 0x8C 0xA2	Square Senti
喱	0xE7 0xB3 0x8E	Centimetre (Japanese)

# Tested Characters

Character	UTF-8 (hex)	Description
kg	0x6B 0x67	Latin Small Letter k and g
k g	0xEF 0xBD 0x8B 0xEF 0xBD 0x87	Fullwidth Latin Small Letter k and g
kg	0xE3 0x8E 0x8F	Square kg
㌘	0xE3 0x8C 0x95	Square Kiroguramu
㌧	0xE7 0x93 0xA9	Kilogram (Japanese)

# Tested Characters

Character	UTF-8 (hex)	Description
(株)	0x28 0xE6 0xA0 0xAA 0x29	numerary adjunct for trees; root, in Parenthesisses
(株)	0xEF 0xBC 0x88 0xE6 0xA0 0xAA 0xEF 0xBC 0x89	numerary adjunct for trees; root, in Fullwidth parentheses
(株)	0xE3 0x88 0xB1	Parenthesized Ideograph Stock

# Tested Characters

Character	UTF-8 (hex)	Description
畠	0xE7 0x95 0x91	dry (as opposed to rice) field; used in Japanese names
畠	0xE7 0x95 0xA0	garden, field, farm, plantation
働	0xE5 0x83 0x8D	labor; work
匂	0xE5 0x8C 0x82	fragrance, smell
♡	0xE2 0x99 0xA1	White Heart Suit
❤	N/A	

# Tested Characters

Character	UTF-8 (hex)	Description
!	0x21	Exclamation Mark
!	0xEF 0xBC 0x81	Fullwidth Exclamation Mark
@	0x40	Commercial At
@	0xEF 0xBC 0xA0	Fullwidth Commercial At
#	0x23	Number Sign
#	0xEF 0xBC 0x83	Fullwidth Number Sign
\$	0x24	Dollar Sign
\$	0xEF 0xBC 0x84	Fullwidth Dollar Sign
?	0x3F	Question Mark

# Tested Characters

Character	UTF-8 (hex)	Description
?	0xEF 0xBC 0x9F	Fullwidth Question Mark
,	0x2C	Comma
,	0xEF 0xBC 0x8C	Fullwidth Comma
.	0x2E	Full Stop
.	0xEF 0xBC 0x8E	Fullwidth Full Stop
、	0xE3 0x80 0x81	Ideographic Comma
。	0xE3 0x80 0x82	Ideographic Full Stop
、	0xEF 0xBD 0xA4	Halfwidth Ideographic Comma
。	0xEF 0xBD 0xA1	Halfwidth Ideographic Full Stop

# Tested Characters

Character	UTF-8 (hex)	Description
高	0xE9 0xAB 0x98	high, tall; lofty elevated
高	0xE9 0xAB 0x99	Variant of 高 U+9ADB, high, tall; lofty elevated
変	0xE5 0xA4 0x89	change, transform, alter (Japanese)
變	0xE8 0xAE 0x8A	change, transform, alter (Traditional Chinese)
变	0xE5 0x8F 0x98	change, transform, alter (Simplified Chinese)
總	0xE7 0xB7 0x8F	collect; overall, altogether (Japanese)
總	0xE7 0xB8 0xBD	collect; overall, altogether (Traditional Chinese)
总	0xE6 0x80 0xBB	collect; overall, altogether (Simplified Chinese)

# Tested Characters

Character	UTF-8 (hex)	Description
가	0xEA 0xB0 0x80	Hangul ga
히	0xED 0x9E 0x88	Hangul hi
까	0xEA 0xB9 0x8C	Hangul kka
찌	0xEC 0xB0 0x8C	Hangul jji
𡇃	N/A	Vietnamese Chữ Nôm
喃	0xE5 0x96 0x83	Vietnamese Chữ Nôm
𦵃	N/A	Vietnamese Chữ Nôm
𩶻	N/A	Vietnamese Chữ Nôm

# Tested Characters

Character	UTF-8 (hex)	Description
Ѐ	0xD0 0x80	Cyrillic Capital Letter Ie with Grave
Ӿ	0xD0 0x81	Cyrillic Capital Letter Io
Ӯ	0xD3 0xBE	Cyrillic Capital Letter Ha with Stroke
ӯ	0xD3 0xBF	Cyrillic Small Letter Ha with Stroke
Ӓ	0xC3 0x84	Latin Capital Letter A with Diaresis
ӓ	0xC3 0xA4	Latin Small Letter a with Diaresis
Ӯ	0xC3 0x9C	Latin Capital Letter U with Diaresis
ӯ	0xC3 0xBC	Latin Small Letter u with Diaresis

# Tested Characters

Character	UTF-8 (hex)	Description
¡	0xC2 0xA1	Inverted Exclamation Mark
¿	0xC2 0xBF	Inverted Question Mark
݂	0xD9 0xB9	Arabic Letter Tteh
݃	0xD9 0xBA	Arabic Letter Tteheh
݄	0xDB 0x92	Arabic Letter Barree
݅	0xDB 0x93	Arabic Letter Barree with Hamza Above

# Setup (on Citus Database)

```
CREATE TABLE t_en_US_utf8 (t text COLLATE "en_US.utf8");
CREATE TABLE t_en_US_x_icu (t text COLLATE "en-US-x-icu");
CREATE TABLE t_ja_JP_utf8 (t text COLLATE "ja_JP.utf8");
CREATE TABLE t_ja_x_icu (t text COLLATE "ja-x-icu");
CREATE TABLE t_ja_JP_x_icu (t text COLLATE "ja-JP-x-icu");
CREATE TABLE t_zh_CN_utf8 (t text COLLATE "zh_CN.utf8");
CREATE TABLE t_zh_x_icu (t text COLLATE "zh-x-icu");
CREATE TABLE t_zh_Hans_x_icu (t text COLLATE "zh-Hans-x-icu");
CREATE TABLE t_zh_Hant_x_icu (t text COLLATE "zh-Hant-x-icu");
```

**SELECT t FROM t\_en\_US\_utf8 ORDER BY t;**

(en\_US.utf8)

#001 -008	#009 -016	#017 -024	#025 -032	#033 -040	#041 -048	#049 -056	#057 -064	#065 -072	#073 -080	#081 -088
,	#	i	9.	は	ン	kg	#	@	ハ	あ…
!	*	vii	①	ば	(株)	cm	\$	A	ン	あ～
i	*	ix	♡	ぱ	50	km	,	Z	あ～	ああ
?	ゞ	xii	①	ん	50	가	.	a	あ-	ああ
ゞ		①	⑨	ア	キロ	까	0	z	♥	あっ
.	VII	⑨	、	ハ	キロ グラム	띠	1	。	字	あつ
@	IX	(1)	。	バ	キロメ ートル	히	9	、	辭	あゝ
\$	XII	(9)	あ	パ	セン チ	!	?	ア	韻	あー

(1/2)

**SELECT t FROM t\_en\_US\_utf8 ORDER BY t;**

(en\_US.utf8)

#089 -096	#097 -104	#105 -112	#113 -120	#121 -128	#129 -136	#137 -144	#145 -152	#153 -160	#161 -168	#169 -176
あ～	k m	9.	a	ä	VII	ꝑ	—	总	總	
あ`	ヰ	a	/a	Ä	xii	ꝑ	九	(株)	總	
ば`	バ`	a	_a	cm	XII	A	働	(株)	變	
は`。	バ`。	_a	(a)	kg	z	α	匁	廷	高	
ヰ	0	/a	(a)	km	(z)	Ω	变	畠	高	
バ`	1	(a)	[a]	ü	(z)	ω	喃	畠		
c m	9	[a]	{a}	Ü	Z	È	壱	糀		
k g	(9)	{a}	A	vii	ꝑ	Ë	麥	糉		

(2/2)

**SELECT t FROM t\_en\_US\_x\_icu ORDER BY t;**

(en-US-x-icu)

#001 -008	#009 -016	#017 -024	#025 -032	#033 -040	#041 -048	#049 -056	#057 -064	#065 -072	#073 -080	#081 -088
a	!	。	(z)	【a】	💓	①	9.	cm	k g	vii
a	!	。	(株)	@	\$	❶	9.	c m	kg	VII
_a	i	(1)	(株)	@	\$	❷	a	cm	km	VII
__a	?	(9)	(株)	/a	0	❸	a	i	k m	xii
,	?	(9)	[a]	/a	0	9	A	l	km	xii
,	¿	(a)	[a]	#	⓪	9	A	ix	ü	XII
`	.	(a)	{a}	#	1	❹	ä	IX	Ü	XII
`	.	(z)	{a}	♡	1	❺	Ä	kg	vii	z

(1/2)

**SELECT t FROM t\_en\_US\_x\_icu ORDER BY t;**

(en-US-x-icu)

#089 -096	#097 -104	#105 -112	#113 -120	#121 -128	#129 -136	#137 -144	#145 -152	#153 -160	#161 -168	#169 -176
z	Ё	까	あ～	あつ	は	パ	九	廷	高	
Z	*	舛	あ…	あつ	ハ	パ <sup>。</sup>	働	畠	高	
Z	*	𠂊	あ～	キ口	ハ	ば <sup>”</sup>	匁	畠	悰	
α	አ	あ	あ～	ヰ	ば	ハ <sup>”</sup>	变	秆	辭	
A	አ	ア	あゝ	ヰ	バ	ん	喃	糧	韙	
ω	አ	ア	あー	ヰ <sup>グム</sup>	バ <sup>”</sup>	ン	壻	總		
Ω	አ	あ <sup>”</sup>	ああ	ヰ <sup>トル</sup>	ぱ <sup>。</sup>	ン	変	總		
Ѐ	가	あ-	ああ	センチ	は <sup>。</sup>	—	总	變		

(2/2)

**SELECT t FROM t\_ja\_JP\_utf8 ORDER BY t;**

(ja\_JP.utf8)

#001 -008	#009 -016	#017 -024	#025 -032	#033 -040	#041 -048	#049 -056	#057 -064	#065 -072	#073 -080	#081 -088
È	VII	⑨	(株)	cm	₩	\$	0	XII	vii	ハ
*	IX	(1)	50	km	♡	(9)	1	Z	xii	バ
*	XII	(9)	50	变	竫	(a)	9	[a]	z	バ°
𩫔	i	9.	キロ	总	辭	(z)	9.	_a	{a}	ン
ߵ	vii	①	キロ グラム	高	齧	(株)	?	a	。	a
ݚ	ix	♥	キロ -トル	가	a	,	@	cm	、	、
ݚ	xii	①	セン チ	까	!	.	A	kg	ア	。
	①	⑨	kg	찌	#	/a	VII	km	キロ	,

(1/2)

**SELECT t FROM t\_ja\_JP\_utf8 ORDER BY t;**

(ja\_JP.utf8)

#089 -096	#097 -104	#105 -112	#113 -120	#121 -128	#129 -136	#137 -144	#145 -152	#153 -160	#161 -168	#169 -176
.	[a]	9	あ	ああ	ば	ぱ	壱	畠	ڦ	
?	{a}	A	あ-	ああ	ぱ	ン	杆	畠	ڦ	
!	【a】	Z	あ~	あつ	ん	A	九	麥	Ü	
_a	\$	a	あ`	あつ	ア	Ω	高	喃	ä	
/a	#	c m	あゝ	あ~	キ口	α	糰	變	ü	
(a)	@	k g	あー	は	ハ	ω	総	廷		
(z)	0	k m	あ~	は°	ハ"	Ё	働	總		
(株)	1	z	あ…	ば"	バ	—	匂	i		

(2/2)

**SELECT t FROM t\_ja\_x\_icu ORDER BY t;**

(ja-x-icu)

#001 -008	#009 -016	#017 -024	#025 -032	#033 -040	#041 -048	#049 -056	#057 -064	#065 -072	#073 -080	#081 -088
a	!	。	(z)	【a】	💓	①	9.	cm	k g	vii
a	!	。	(株)	@	\$	❶	9.	c m	kg	VII
_a	i	(1)	(株)	@	\$	❷	a	cm	km	VII
__a	?	(9)	(株)	/a	0	❸	a	i	k m	xii
,	?	(9)	[a]	/a	0	9	A	I	km	xii
,	¿	(a)	[a]	#	⓪	9	A	ix	ü	XII
`	.	(a)	{a}	#	1	⓫	ä	IX	Ü	XII
`	.	(z)	{a}	♡	1	❹	Ä	kg	vii	z

(1/2)

**SELECT t FROM t\_ja\_x\_icu ORDER BY t;**

(ja-x-icu)

#089 -096	#097 -104	#105 -112	#113 -120	#121 -128	#129 -136	#137 -144	#145 -152	#153 -160	#161 -168	#169 -176
z	Ё	까	あ～	あっ	は	ぱ°	壱	畠	总	
Z	*	舛	あ…	あつ	ハ	パ°	糸	畠	高	
Z	*	𠂊	あ～	キ口	ハ	ば°	九	麥	梓	
α	ㅏ	あ	あ～	ヰ	ば	ハ°	高	喃	辭	
A	ㅓ	ア	あー	ヰ	バ	ん	糰	變	韻	
ω	ㅡ	ア	ああ	ヰグム	バ	ン	総	廷		
Ω	ㅡ	あ°	あゝ	ヰトル	ぱ°	ン	働	總		
Ѐ	가	あ-	ああ	センチ	パ	—	匂	变		

(2/2)

**SELECT t FROM t\_ja\_JP\_x\_icu ORDER BY t;**

(ja-JP-x-icu)

#001 -008	#009 -016	#017 -024	#025 -032	#033 -040	#041 -048	#049 -056	#057 -064	#065 -072	#073 -080	#081 -088
a	!	。	(z)	【a】	💓	①	9.	cm	k g	vii
a	!	。	(株)	@	\$	❶	9.	c m	kg	VII
_a	i	(1)	(株)	@	\$	❷	a	cm	km	VII
__a	?	(9)	(株)	/a	0	❸	a	i	k m	xii
,	?	(9)	[a]	/a	0	9	A	I	km	xii
,	¿	(a)	[a]	#	⓪	9	A	ix	ü	XII
`	.	(a)	{a}	#	1	⓫	ä	IX	Ü	XII
`	.	(z)	{a}	♡	1	❾	Ä	kg	vii	z

(The same as ja-x-icu)

(1/2)

**SELECT t FROM t\_ja\_JP\_x\_icu ORDER BY t;**

(ja-JP-x-icu)

#089 -096	#097 -104	#105 -112	#113 -120	#121 -128	#129 -136	#137 -144	#145 -152	#153 -160	#161 -168	#169 -176
z	Ё	ヵ	あ～	あつ	は	ぱ°	壱	畠	总	
Z	*	刈	あ…	あつ	ハ	パ°	糸	畠	高	
Z	*	𠂊	あ～	キ口	ハ	ば°	九	麥	梓	
α	়	あ	あ～	ヰ	ば	ハ°	高	喃	離	
A	়	ア	あー	ヰ	バ	ん	糰	變	韻	
ω	়	ア	ああ	ヰグム	バ°	ン	総	廷		
Ω	়	あ°	あゝ	ヰトৰ	ぱ°	ン	働	總		
Ѐ	가	あ-	ああ	ヰチ	パ°	—	匂	变		

(The same as ja-x-icu)

(2/2)

**SELECT t FROM t\_zh\_CN\_utf8 ORDER BY t;**

(zh\_CN.utf8)

#001 -008	#009 -016	#017 -024	#025 -032	#033 -040	#041 -048	#049 -056	#057 -064	#065 -072	#073 -080	#081 -088
,	#	i	9.	は	ン	kg	ひ	9	、	辭
!	*	vii	①	ば	(株)	cm	!	?	ア	麌
i	*	ix	♡	ぱ	50	km	#	@	ハ	あ…
?	ゞ	xii	①	ん	50	匁	\$	A	ン	あ～
ゞ	I	①	⑨	ア	キロ	畠	,	Z	あ～	ああ
.	VII	⑨	、	ハ	キロ グラム	가	.	a	あ-	ああ
@	IX	(1)	。	バ	キロメ ートル	까	0	z	♥	あつ
\$	XII	(9)	あ	パ	セン チ	刈	1	。	字	あつ

(Similar to en\_US.utf8)

(1/2)

**SELECT t FROM t\_zh\_CN\_utf8 ORDER BY t;**

(zh\_CN.utf8)

#089 -096	#097 -104	#105 -112	#113 -120	#121 -128	#129 -136	#137 -144	#145 -152	#153 -160	#161 -168	#169 -176
あゝ	c m	9	[a]	{a}	Ü	Z	È	高	壹	
あー	k g	(9)	{a}	A	vii	ঢ	ঃ	九	(株)	
あ～	k m	9.	a	ä	VII	ঢ	变	糧	(株)	
あ`	ヰ	a	/a	Ä	xii	ৢ	麥	喃	总	
ば`	バ	a	_a	cm	XII	A	變	延	總	
ば°	バ°	_a	[a]	kg	z	α	總	杆		
ヰ口	0	/a	(a)	km	(z)	Ω	働	畠		
バ`	1	(a)	[a]	ü	(z)	ω	高	—		

(Similar to en\_US.utf8)

(2/2)

**SELECT t FROM t\_zh\_x\_icu ORDER BY t;**

(zh-x-icu)

#001 -008	#009 -016	#017 -024	#025 -032	#033 -040	#041 -048	#049 -056	#057 -064	#065 -072	#073 -080	#081 -088
a	!	。	(z)	【a】	❤	①	9.	cm	k g	vii
a	!	。	(株)	@	\$	1	9.	c m	kg	VII
_a	i	(1)	(株)	@	\$	50	a	cm	km	VII
__a	?	(9)	(株)	/a	0	50	a	i	k m	xii
,	?	(9)	[a]	/a	0	9	A	l	km	xii
,	¿	(a)	[a]	#	①	9	A	ix	ü	XII
,	.	(a)	{a}	#	1	⑨	ä	IX	Ü	XII
,	.	(z)	{a}	♡	1	9	Ä	kg	vii	z

(Similar to ja-x-icu)

(1/2)

**SELECT t FROM t\_zh\_x\_icu ORDER BY t;**

(zh-x-icu)

#089 -096	#097 -104	#105 -112	#113 -120	#121 -128	#129 -136	#137 -144	#145 -152	#153 -160	#161 -168	#169 -176
z	Ё	까	あ～	あつ	は	パ	变	喃	總	
Z	*	舛	あ…	あつ	ハ	パ <sup>。</sup>	変	廷	總	
Z	*	苟	あ～	キ口	ハ	ば <sup>”</sup>	變	秆	字	
α	ㅏ	あ	あ～	ヰ口	ば	ハ <sup>”</sup>	働	畠	辭	
A	ㅓ	ア	あゝ	ヰ口	バ	ん	高	畠	韻	
ω	ㅡ	ア	あー	ヰ口 ゲム	バ	ン	高	一		
Ω	ㅡ	あ <sup>”</sup>	ああ	ヰ口 メトル	ぱ <sup>。</sup>	ン	九	壹		
Ѐ	가	あ-	ああ	センチ	は <sup>。</sup>	匂	糧	总		

(2/2)

**SELECT t FROM t\_zh\_Hans\_x\_icu ORDER BY t;**

(zh-Hans-x-icu)

#001 -008	#009 -016	#017 -024	#025 -032	#033 -040	#041 -048	#049 -056	#057 -064	#065 -072	#073 -080	#081 -088
a	!	。	(z)	【a】	💓	①	9.	cm	k g	vii
a	!	。	(株)	@	\$	❶	9.	c m	kg	VII
_a	i	(1)	(株)	@	\$	❷	a	cm	km	VII
__a	?	(9)	(株)	/a	0	❸	a	i	k m	xii
,	?	(9)	[a]	/a	0	9	A	I	km	xii
,	¿	(a)	[a]	#	⓪	9	A	ix	ü	XII
`	.	(a)	{a}	#	1	❹	ä	IX	Ü	XII
`	.	(z)	{a}	♡	1	❺	Ä	kg	vii	z

(1/2)

**SELECT t FROM t\_zh\_Hans\_x\_icu ORDER BY t;**

(zh-Hans-x-icu)

#089 -096	#097 -104	#105 -112	#113 -120	#121 -128	#129 -136	#137 -144	#145 -152	#153 -160	#161 -168	#169 -176
z	Ё	까	あ～	あつ	は	パ	变	喃	總	
Z	*	舛	あ…	あつ	ハ	パ <sup>。</sup>	変	廷	總	
Z	*	苟	あ～	キ口	ハ	ば <sup>”</sup>	變	秆	字	
α	ㅏ	あ	あ～	ヰ	ば	ハ <sup>”</sup>	働	畠	辭	
A	ㅓ	ア	あゝ	キロ	バ	ん	高	畠	韻	
ω	ㅡ	ア	あー	キロ ゲム	バ	ン	高	一		
Ω	ㅡ	あ <sup>”</sup>	ああ	キロ ートル	ぱ <sup>。</sup>	ン	九	壹		
Ѐ	가	あ-	ああ	センチ	は <sup>。</sup>	匂	糧	总		

(2/2)

**SELECT t FROM t\_zh\_Hant\_x\_icu ORDER BY t;**

(zh-Hant-x-icu)

#001 -008	#009 -016	#017 -024	#025 -032	#033 -040	#041 -048	#049 -056	#057 -064	#065 -072	#073 -080	#081 -088
a	!	。	(z)	【a】	💓	①	9.	cm	k g	vii
a	!	。	(株)	@	\$	❶	9.	c m	kg	VII
_a	i	(1)	(株)	@	\$	❷	a	cm	km	VII
__a	?	(9)	(株)	/a	0	❸	a	i	k m	xii
,	?	(9)	[a]	/a	0	9	A	I	km	xii
,	¿	(a)	[a]	#	⓪	9	A	ix	ü	XII
`	.	(a)	{a}	#	1	❹	ä	IX	Ü	XII
`	.	(z)	{a}	♡	1	❺	Ä	kg	vii	z

(Similar to zh-Hans-x-icu)

(1/2)

**SELECT t FROM t\_zh\_Hant\_x\_icu ORDER BY t;**

(zh-Hant-x-icu)

#089 -096	#097 -104	#105 -112	#113 -120	#121 -128	#129 -136	#137 -144	#145 -152	#153 -160	#161 -168	#169 -176	
z	Ё	𠮩	あ～	あつ	は	パ	九	杆	總		
Z	*	𠮩	あ…	あつ	ハ	パ <sup>。</sup>	匂	畠	變		
Z	*	𠮩	あ～	キ口	ハ	ば <sup>”</sup>	壱	高	字		
α	়	ା	ାବୁ	କିର୍ତ୍ତ	ବା	ବା <sup>”</sup>	ବିପାନ୍ତ	କାମିକାରୀ	ବିଜ୍ଞାନ		
A	়	ା	ାବୁ	କିର୍ତ୍ତ	ବା	ବା	ବିପାନ୍ତ	କାମିକାରୀ	ବିଜ୍ଞାନ		
ω	়	ା	ାବୁ	କିର୍ତ୍ତ	ବା	ବା	ବିପାନ୍ତ	କାମିକାରୀ	ବିଜ୍ଞାନ		
Ω	়	ା	ାବୁ	କିର୍ତ୍ତ	ବା	ବା	ବିପାନ୍ତ	କାମିକାରୀ	ବିଜ୍ଞାନ		
Ѐ	ା	ାବୁ	ାବୁ	କିର୍ତ୍ତ	ବା	ବା	ବିପାନ୍ତ	କାମିକାରୀ	ବିଜ୍ଞାନ		

(Similar to zh-Hans-x-icu)

(2/2)

Cost Comparison > By EXPLAIN >

## libc vs. ICU (English)

### libc (en\_US.utf8)

```
citus=> EXPLAIN SELECT t FROM t_en_us_utf8 ORDER BY t;
          QUERY PLAN
-----
Sort  (cost=8.73..9.14 rows=165 width=4)
  Sort Key: t COLLATE "en_US.utf8"
  ->  Seq Scan on t_en_us_utf8  (cost=0.00..2.65 rows=165 width=4)
```

### ICU (en-US-x-icu)

```
citus=> EXPLAIN SELECT t FROM t_en_us_x_icu ORDER BY t;
          QUERY PLAN
-----
Sort  (cost=8.73..9.14 rows=165 width=4)
  Sort Key: t COLLATE "en-US-x-icu"
  ->  Seq Scan on t_en_us_x_icu  (cost=0.00..2.65 rows=165 width=4)
```

## Cost Comparison > By EXPLAIN > libc vs. ICU (Japanese)

### libc (ja\_JP.utf8)

```
citus=> EXPLAIN SELECT t FROM t_ja_JP_utf8 ORDER BY t;
          QUERY PLAN
-----
Sort  (cost=8.73..9.14 rows=165 width=4)
  Sort Key: t COLLATE "ja_JP.utf8"
  ->  Seq Scan on t_ja_jp_utf8  (cost=0.00..2.65 rows=165 width=4)
```

### ICU (ja-x-icu)

```
citus=> EXPLAIN SELECT t FROM t_ja_x_icu ORDER BY t;
          QUERY PLAN
-----
Sort  (cost=9.73..10.14 rows=165 width=4)
  Sort Key: t COLLATE "ja-x-icu"
  ->  Seq Scan on t_ja_x_icu  (cost=0.00..3.65 rows=165 width=4)
```

Cost Comparison > By EXPLAIN >

## ja-x-icu vs. ja-JP-x-icu

### ja-x-icu

```
citus=> EXPLAIN SELECT t FROM t_ja_x_icu ORDER BY t;
          QUERY PLAN
-----
Sort  (cost=9.73..10.14 rows=165 width=4)
  Sort Key: t COLLATE "ja-x-icu"
  ->  Seq Scan on t_ja_x_icu  (cost=0.00..3.65 rows=165 width=4)
```

### ja-JP-x-icu

```
citus=> EXPLAIN SELECT t FROM t_ja_JP_x_icu ORDER BY t;
          QUERY PLAN
-----
Sort  (cost=8.73..9.14 rows=165 width=4)
  Sort Key: t COLLATE "ja-JP-x-icu"
  ->  Seq Scan on t_ja_jp_x_icu  (cost=0.00..2.65 rows=165 width=4)
```

Cost Comparison > By EXPLAIN >

# en-US-x-icu vs. ja-JP-x-icu

## en-US-x-icu

```
citus=> EXPLAIN SELECT t FROM t_en_us_x_icu ORDER BY t;
          QUERY PLAN
-----
Sort  (cost=8.73...9.14 rows=165 width=4)
  Sort Key: t COLLATE "en-US-x-icu"
  ->  Seq Scan on t_en_us_x_icu  (cost=0.00..2.65 rows=165 width=4)
```

## ja-JP-x-icu

```
citus=> EXPLAIN SELECT t FROM t_ja_JP_x_icu ORDER BY t;
          QUERY PLAN
-----
Sort  (cost=8.73...9.14 rows=165 width=4)
  Sort Key: t COLLATE "ja-JP-x-icu"
  ->  Seq Scan on t_ja_jp_x_icu  (cost=0.00..2.65 rows=165 width=4)
```

# ja-x-icu vs. ja-JP-x-icu with COLLATE clause

## COLLATE "ja-x-icu"

```
citus=> EXPLAIN SELECT t FROM t_ja_JP ORDER BY t COLLATE "ja-x-icu";  
          QUERY PLAN
```

```
-----  
Sort  (cost=8.73..9.14 rows=165 width=36)  
  Sort Key: t COLLATE "ja-x-icu"  
  ->  Seq Scan on t_ja_jp  (cost=0.00..2.65 rows=165 width=36)
```

## COLLATE "ja-JP-x-icu"

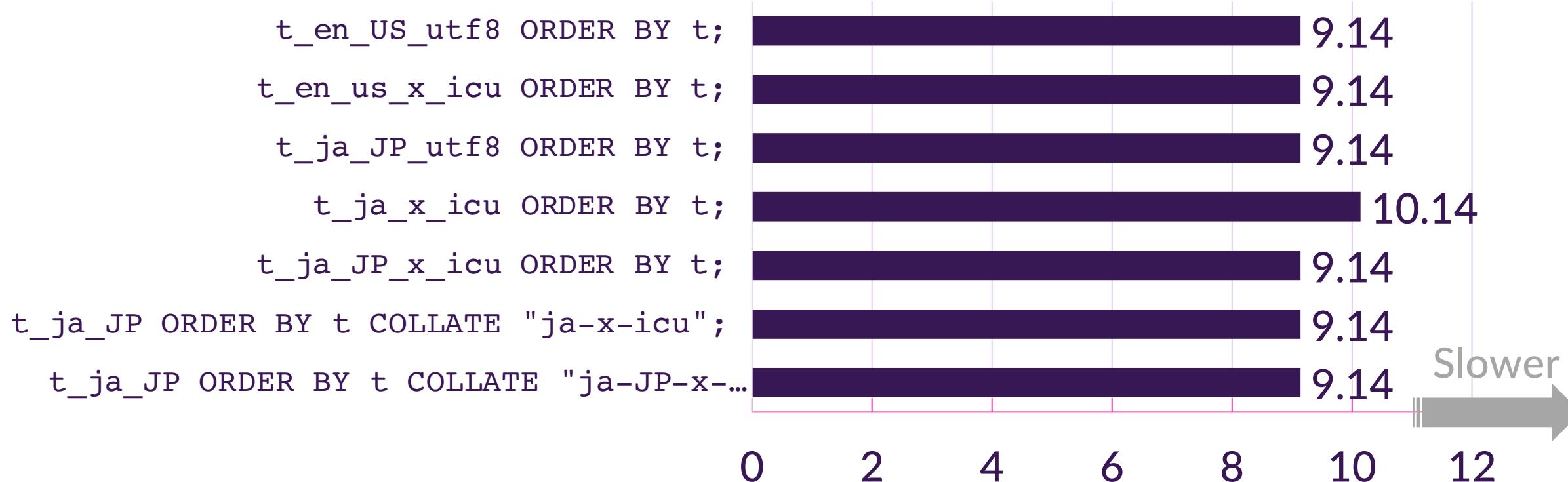
```
citus=> EXPLAIN SELECT t FROM t_ja_JP ORDER BY t COLLATE "ja-JP-x-icu";  
          QUERY PLAN
```

```
-----  
Sort  (cost=8.73..9.14 rows=165 width=36)  
  Sort Key: t COLLATE "ja-JP-x-icu"  
  ->  Seq Scan on t_ja_jp  (cost=0.00..2.65 rows=165 width=36)
```

# Summary of Total Cost by EXPLAIN

SELECT t FROM ...

Total Cost



Cost Comparison > Actual time >

# libc vs. ICU (English)

## libc (en\_US.utf8)

```
citus=> EXPLAIN ANALYZE SELECT t FROM t_en_US_utf8 ORDER BY t;
                                         QUERY PLAN
-----
Sort  (cost=8.73..9.14 rows=165 width=4) (actual time=0.990..1.006 rows=165 loops=1)
  Sort Key: t COLLATE "en_US.utf8"
  Sort Method: quicksort  Memory: 32kB
    -> Seq Scan on t_en_us_utf8  (cost=0.00..2.65 rows=165 width=4) (actual time=0.016..0.109 rows=165 loops=1)
Planning Time: 0.153 ms
Execution Time: 1.046 ms
```

## ICU (en-US-x-icu)

```
citus=> EXPLAIN ANALYZE SELECT t FROM t_en_us_x_icu ORDER BY t;
                                         QUERY PLAN
-----
Sort  (cost=8.73..9.14 rows=165 width=4) (actual time=0.535..0.551 rows=165 loops=1)
  Sort Key: t COLLATE "en-US-x-icu"
  Sort Method: quicksort  Memory: 32kB
    -> Seq Scan on t_en_us_x_icu  (cost=0.00..2.65 rows=165 width=4) (actual time=0.016..0.042 rows=165 loops=1)
Planning Time: 0.152 ms
Execution Time: 0.591 ms
```

Cost Comparison > Actual time >

# libc vs. ICU (Japanese)

## libc (ja\_JP.utf8)

```
citus=> EXPLAIN ANALYZE SELECT t FROM t_ja_JP_utf8 ORDER BY t;
                                         QUERY PLAN
-----
Sort  (cost=8.73..9.14 rows=165 width=4) (actual time=2.498..2.514 rows=165 loops=1)
  Sort Key: t COLLATE "ja_JP.utf8"
  Sort Method: quicksort  Memory: 32kB
  -> Seq Scan on t_ja_jp_utf8  (cost=0.00..2.65 rows=165 width=4) (actual time=0.017..0.044 rows=165 loops=1)
Planning Time: 0.205 ms
Execution Time: 2.553 ms
```

## ICU (ja-x-icu)

```
citus=> EXPLAIN ANALYZE SELECT t FROM t_ja_x_icu ORDER BY t;
                                         QUERY PLAN
-----
Sort  (cost=9.73..10.14 rows=165 width=4) (actual time=0.667..0.688 rows=165 loops=1)
  Sort Key: t COLLATE "ja-x-icu"
  Sort Method: quicksort  Memory: 32kB
  -> Seq Scan on t_ja_x_icu  (cost=0.00..3.65 rows=165 width=4) (actual time=0.011..0.044 rows=165 loops=1)
Planning Time: 0.150 ms
Execution Time: 0.742 ms
```

Cost Comparison > Actual time >

# ja-x-icu vs. ja-JP-x-icu

## ja-x-icu

```
citus=> EXPLAIN ANALYZE SELECT t FROM t_ja_x_icu ORDER BY t;
          QUERY PLAN
-----
Sort  (cost=9.73..10.14 rows=165 width=4) (actual time=0.667..0.688 rows=165 loops=1)
  Sort Key: t COLLATE "ja-x-icu"
  Sort Method: quicksort  Memory: 32kB
    -> Seq Scan on t_ja_x_icu  (cost=0.00..3.65 rows=165 width=4) (actual time=0.011..0.044 rows=165 loops=1)
Planning Time: 0.150 ms
Execution Time: 0.742 ms
```

## ja-JP-x-icu

```
citus=> EXPLAIN ANALYZE SELECT t FROM t_ja_JP_x_icu ORDER BY t;
          QUERY PLAN
-----
Sort  (cost=8.73..9.14 rows=165 width=4) (actual time=0.493..0.509 rows=165 loops=1)
  Sort Key: t COLLATE "ja-JP-x-icu"
  Sort Method: quicksort  Memory: 32kB
    -> Seq Scan on t_ja_jp_x_icu  (cost=0.00..2.65 rows=165 width=4) (actual time=0.016..0.044 rows=165 loops=1)
Planning Time: 0.146 ms
Execution Time: 0.547 ms
```

Cost Comparison > Actual time >

# ja-x-icu vs. ja-JP-x-icu with COLLATE clause

## COLLATE "ja-x-icu"

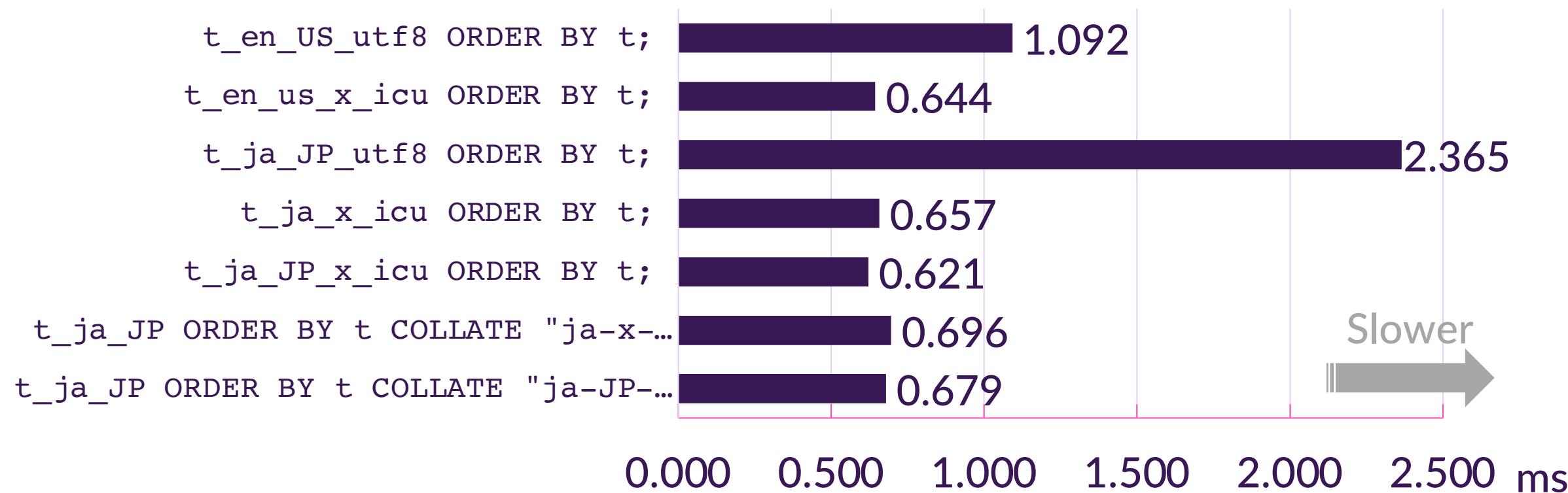
```
citus=> EXPLAIN ANALYZE SELECT t FROM t_ja_JP ORDER BY t COLLATE "ja-x-icu";
                                         QUERY PLAN
-----
Sort  (cost=8.73..9.14 rows=165 width=36) (actual time=0.472..0.488 rows=165 loops=1)
  Sort Key: t COLLATE "ja-x-icu"
  Sort Method: quicksort  Memory: 32kB
  -> Seq Scan on t_ja_jp  (cost=0.00..2.65 rows=165 width=36) (actual time=0.017..0.049 rows=165 loops=1)
Planning Time: 0.171 ms
Execution Time: 0.536 ms
```

## COLLATE "ja-JP-x-icu"

```
citus=> EXPLAIN ANALYZE SELECT t FROM t_ja_JP ORDER BY t COLLATE "ja-JP-x-icu";
                                         QUERY PLAN
-----
Sort  (cost=8.73..9.14 rows=165 width=36) (actual time=0.551..0.567 rows=165 loops=1)
  Sort Key: t COLLATE "ja-JP-x-icu"
  Sort Method: quicksort  Memory: 32kB
  -> Seq Scan on t_ja_jp  (cost=0.00..2.65 rows=165 width=36) (actual time=0.014..0.050 rows=165 loops=1)
Planning Time: 0.106 ms
Execution Time: 0.614 ms
```

# Summary of Actual Time by ANALYZE

SELECT t FROM ... Planning Time + Execution Time



\* Each value is the median of 10-time ANALYZE execution results.

# Conclusion

- Using Azure Database for PostgreSQL Hyperscale (Citus), you can specify locale and collation for each table (not database).
- ICU collation is better than libc not only for the sort order but also the performance.
- Setting up ICU collation for each table is better than using COLLATE clause with ICU collation for libc tables.

# References

- PostgreSQL: Documentation
  - <https://www.postgresql.org/docs/>
- ロケール(国際化と地域化) | Let's POSTGRES
  - <https://lets.postgresql.jp/documents/technical/text-processing/2>
- PostgreSQL 10のICUコレーションを使うと日本語を普通にソートでき、更に文字順序までカスタマイズできる - yohgaki's blog
  - <https://blog.ohgaki.net/postgresql-10-icu-locale-collation-enables-natural-japanese-sorting>
- PostgreSQL Collations | Vertabelo Database Modeler
  - <https://vertabelo.com/blog/collations-in-postgresql/>

# References

- Azure Database for PostgreSQLでCollation(照合順序)を指定してCREATE DATABASEしたい - 浅草橋青空市場
  - <https://asazure.hatenablog.jp/entry/2018/04/13/203114>
- Don't let collation versions corrupt your PostgreSQL indexes - Microsoft Tech Community
  - <https://techcommunity.microsoft.com/t5/azure-database-for-postgresql/don-t-let-collation-versions-corrupt-your-postgresql-indexes/ba-p/1978394>

# Agenda

- Introduction
- A number of different character sets
  - Character Set Support
- Locale-specific collation order
  - Locale Support
  - Collation Support



# THANK YOU AND HAPPY HACKING!

Q&A and discussion  
in Microsoft Open Source Discord



#cituscon



Keep in touch with the speaker?



@kske\_t



@k14i

