
bistring

Release 0.4.0

Tavian Barnes

Apr 30, 2020

CONTENTS:

- 1 Introduction** **3**

- 2 Frequently Asked Questions** **5**
 - 2.1 What is a bistring, anyway? 5
 - 2.2 Why am I getting more text than I expect when slicing? 5
 - 2.3 What if I don't know the alignment? 6
 - 2.4 How do I get the actual indices, rather than just substrings? 6
 - 2.5 How do I perform case-insensitive operations? 6
 - 2.6 Forget case insensitivity, how do I make sure that identical looking strings compare equal? 7
 - 2.7 What about similar-looking strings, that aren't necessarily identical? 7
 - 2.8 How do I ensure I get the same results on every machine? 7
 - 2.9 Tokenization 8

- 3 Python** **11**
 - 3.1 bistr 11
 - 3.2 BistrBuilder 17
 - 3.3 Alignment 20
 - 3.4 Tokenization 24
 - 3.5 Tokenizer 26

- 4 JavaScript** **29**
 - 4.1 BiString 29
 - 4.2 BiStringBuilder 34
 - 4.3 Alignment 37
 - 4.4 Tokenization 41
 - 4.5 Tokenizer 44

- 5 Indices and tables** **47**

- Index** **49**

The bistring library provides non-destructive versions of common string processing operations like normalization, case folding, and find/replace. Each bistring remembers the original string, and how its substrings map to substrings of the modified version.

INTRODUCTION

Many operations commonly performed on text strings are destructive; that is, they lose some information about the original string. Systems that deal with text will commonly perform many of these operations on their input, whether it's changing case, performing unicode normalization, collapsing whitespace, stripping punctuation, etc. This helps systems behave in a more uniform manner regarding the many different ways you or I might express the same thing. But the consequence is that when handling parts of this processed text, it may be hard to know what exactly the user originally wrote. Sometimes those details can be very important to the user.

Consider an AI personal assistant, for example, that is helping a user send a text message to a friend. The user writes,

send jane a text that says, "Hey! How are you? Haven't seen you in a while, what's up "

The system may perform some normalization on that text, such that it ends up looking like this, with casing and punctuation gone:

send jane a text that says hey how are you havent seen you in a while whats up emoji

The AI may then identify that the body of the message should be:

hey how are you havent seen you in a while whats up emoji

However, that message wouldn't make much sense as-is. If the assistant uses *bistring* though, it's easy for it to match that with the original text the user intended:

```
>>> from bistring import bistr
>>> query = bistr(
...     'send jane a text that says, '
...     '"Hey! How are you? Haven\'t seen you in a while, what\'s up "'
... )
```

```
>>> # Get rid of upper-/lower-case distinctions
>>> query = query.casefold()
>>> print(query.modified)
send jane a text that says, "hey! how are you? haven't seen you in a while, what's up
↵"
```

```
>>> import regex
>>> # Remove all punctuation
>>> query = query.sub(regex.compile(r'\pP'), '')
>>> # Replace all symbols with 'emoji'
>>> query = query.sub(regex.compile(r'\pS'), 'emoji')
>>> print(query.modified)
send jane a text that says hey how are you havent seen you in a while whats up emoji
```

```
>>> # Extract the substring we care about, the message body
>>> message = query[27:84]
>>> print(message.modified)
hey how are you havent seen you in a while whats up emoji
>>> print(message.original)
Hey! How are you? Haven't seen you in a while, what's up
```

Every *bistr* keeps track of the original string it started with, and maintains a sequence alignment between the original and the modified strings. This alignment means that it knows exactly what substring of the original text is associated with every chunk of the modified text. So when you slice a *bistr*, you get the matching slice of original text automatically!

FREQUENTLY ASKED QUESTIONS

2.1 What is a bistring, anyway?

Simply put, a *bistring* is a pair of strings, an original string and a modified one, along with information about how they align with each other. The *bistring.bistr* class has an API very similar to the built-in *str*, but all its operations keep track of the original string and the alignment for you.

```
>>> from bistring import bistr
>>> s = bistr('HELLO WORLD')
>>> print(s)
'HELLO WORLD'
>>> s = s.lower()
>>> print(s)
('HELLO WORLD' 'hello world')
>>> print(s[6:])
('WORLD' 'world')
```

2.2 Why am I getting more text than I expect when slicing?

When a bistring doesn't have precise enough alignment information to slice exactly, it will give you back the smallest string it knows for certain contains a match for the region you requested. In the worst case, that may be the entire string! This happens, for example, when you use the two-argument *bistr* constructor, which makes no effort to infer a granular alignment between the strings:

```
>>> s = bistr('color', 'colour')
>>> print(s[3:5])
('color' 'ou')
```

Instead, you should start from your original string as a *bistr*, and then transform it how you want:

```
>>> s = bistr('color')
>>> s = s.sub(r'(?<=col)o(?:=r)', 'ou')
>>> print(s)
('color' 'colour')
>>> print(s[3:5])
('o' 'ou')
```

Alternatively, you can piece many smaller bistrings together to achieve the alignment you want manually:

```
>>> s = bistr('col') + bistr('o', 'ou') + bistr('r')
>>> print(s)
('color' 'colour')
>>> print(s[3:5])
('o' 'ou')
```

2.3 What if I don't know the alignment?

If at all possible, you should use *bistring* all the way through your text processing code, which will ensure an accurate alignment is tracked for you. If you don't control that code, or there are other reasons it won't work with *bistring*, you can still have us guess an alignment for you in simple cases with *bistring.bistr.infer()*.

```
>>> s = bistr.infer('color', 'colour')
>>> print(s[0:3])
'col'
>>> print(s[3:5])
('o' 'ou')
>>> print(s[5:6])
'r'
```

infer() is an expensive operation ($O(N \times M)$ in the length of the strings), so if you absolutely need it, try to use it only for short strings.

2.4 How do I get the actual indices, rather than just substrings?

Use *bistring.bistr.alignment*:

```
>>> s = bistr('The quick, brown ')
>>> s = s.replace(', ', '')
>>> s = s.replace(' ', 'fox')
>>> print(s[16:19])
(' ' 'fox')
>>> s.alignment.original_bounds(16, 19)
(17, 18)
>>> s.alignment.modified_bounds(11, 16)
(10, 15)
>>> print(s[10:15])
'brown'
```

See *bistring.Alignment* for more details.

2.5 How do I perform case-insensitive operations?

Use *bistring.bistr.casefold()*. Do not use *lower()*, *upper()*, or any other method, as you will get wrong results for many non-English languages.

To check case-insensitive equality, you don't even need *bistring*:

```
>>> 'HELLO WORLD!'.casefold() == 'HeLlO wOrLd!'.casefold()
True
```

To search for a substring case-insensitively:

```
>>> s = bistr('Bundesstraße').casefold()
>>> s.find_bounds('STRASSE'.casefold())
(6, 13)
>>> print(s[6:13])
('straße' 'strasse')
```

2.6 Forget case insensitivity, how do I make sure that identical looking strings compare equal?

This is a hard problem with Unicode strings. To start with, you should at least perform some kind of [Unicode normalization](#). That ensures that different ways of writing the semantically identical thing (e.g. with precomposed accented characters vs. combining accents) become actually identical:

```
>>> a = bistr('\u00EAtre') # 'être' with a single character for the ê
>>> b = bistr('e\u0302tre') # 'etre' with an 'e' and a combining '^'
>>> a.normalize('NFC').modified == b.normalize('NFC').modified
True
>>> a.normalize('NFD').modified == b.normalize('NFD').modified
True
```

Normalization form NFC tries to keep precomposed characters together whenever possible, while NFD always decomposes them. In general, NFC is more convenient for people to work with, but NFD can be useful for things like removing accents and other combining marks from text.

2.7 What about similar-looking strings, that aren't necessarily identical?

Unicode contains things like ligatures, alternative scripts, and other oddities than can result in similar-looking strings that are represented very differently. Here is where the “compatibility” normalization forms, NFKC and NFKD, can help:

```
>>> s = bistr(' ')
>>> s = s.normalize('NFKC')
>>> print(s)
(' ' 'Hello world')
>>> print(s[6:])
(' ' 'world')
```

2.8 How do I ensure I get the same results on every machine?

Always pass an explicit locale to any *bistr* method that takes one. Many of Python's string APIs implicitly use the system's default locale, which may be quite different than the one you developed with. While this may be the right behaviour if you're displaying strings to the current user, it's rarely the right behaviour if you're dealing with text that originated or will be displayed elsewhere, e.g. for cloud software. *bistr* always accepts a locale parameter in these APIs, to ensure reproducible and sensible results:

```
>>> # s will be 'I' in most locales, but 'İ' in Turkish locales!
>>> s = bistr('i').upper()
>>> # An English locale guarantees a dotless capital I
>>> print(bistr('i').upper('en_US'))
('i' 'I')
>>> # A Turkish locale gives a dotted capital İ
>>> print(bistr('i').upper('tr_TR'))
('i' 'İ')
```

2.9 Tokenization

2.9.1 How do I tokenize text in a reversible way?

bistring provides some convenient tokenization APIs that track string indices. To use Unicode word boundary rules, for example:

```
>>> from bistring import WordTokenizer
>>> tokenizer = WordTokenizer('en_US')
>>> tokens = tokenizer.tokenize('The quick, brown fox jumps over the lazy dog')
>>> print(tokens[1])
[4:9]='quick'
```

2.9.2 How do I find the whole substring of text for some tokens?

bistring.Tokenization.substring() gives the substring itself. *bistring.Tokenization.text_bounds()* gives the bounds of that substring.

```
>>> print(tokens.substring(1, 3))
'quick, brown'
>>> tokens.text_bounds(1, 3)
(4, 16)
```

2.9.3 How do I find the tokens for a substring of text?

bistring.Tokenization.bounds_for_text()

```
>>> tokens.bounds_for_text(4, 16)
(1, 3)
>>> print(tokens.substring(1, 3))
'quick, brown'
```

2.9.4 How to I snap a substring of text to the nearest token boundaries?

bistring.Tokenization.snap_text_bounds()

```
>>> print(tokens.text[6:14])
'ick, bro'
>>> tokens.snap_text_bounds(6, 14)
(4, 16)
>>> print(tokens.text[4:16])
'quick, brown'
```

2.9.5 What if I don't know the token positions?

If at all possible, you should use a *bistring.Tokenizer* or some other method that tokenizes with position information. If you can't, you can use *bistring.Tokenization.infer()* to guess the alignment for you:

```
>>> from bistring import Tokenization
>>> tokens = Tokenization.infer('hello, world!', ['hello', 'world'])
>>> print(tokens[0])
[0:5]='hello'
>>> print(tokens[1])
[7:12]='world'
```


3.1 bistr

class `bistring.bistr`

Bases: `object`

A bidirectionally transformed string.

A *bistr* can be constructed from only a single string, which will give it identical original and modified strings and an identity alignment:

```
>>> s = bistr('test')
>>> s.original
'test'
>>> s.modified
'test'
>>> s.alignment
Alignment.identity(4)
```

You can also explicitly specify both the original and modified string. The inferred alignment will be as course as possible:

```
>>> s = bistr('TEST', 'test')
>>> s.original
'TEST'
>>> s.modified
'test'
>>> s.alignment
Alignment([(0, 0), (4, 4)])
```

Finally, you can specify the alignment explicitly too, if you know it:

```
>>> s = bistr('TEST', 'test', Alignment.identity(4))
>>> s[1:3]
bistr('ES', 'es', Alignment.identity(2))
```

original

The original string, before any modifications.

modified

The current value of the string, after all modifications.

alignment

The sequence alignment between *original* and *modified*.

classmethod infer (*original, modified, cost_fn=None*)

Create a *bistr*, automatically inferring an alignment between the *original* and *modified* strings.

This method can be useful if the modified string was produced by some method out of your control. If at all possible, you should start with `bistr(original)` and perform non-destructive operations to get to the modified string instead.

```
>>> s = bistr.infer('color', 'colour')
>>> print(s[0:3])
'col'
>>> print(s[3:5])
('o' 'ou')
>>> print(s[5:6])
'r'
```

infer() tries to be intelligent about certain aspects of Unicode, which enables it to guess good alignments between strings that have been case-mapped, normalized, etc.:

```
>>> s = bistr.infer(
...     ' ',
...     'the quick brown fox jumps over the lazy dog',
... )
>>> print(s[0:3])
(' ' 'the')
>>> print(s[4:9])
(' ' 'quick')
>>> print(s[10:15])
(' ' 'brown')
>>> print(s[16:19])
(' ' 'fox')
```

Warning: this operation has time complexity $O(N*M)$, where N and M are the lengths of the original and modified strings, and so should only be used for relatively short strings.

Parameters

- **original** (*str*) – The original string
- **modified** (*str*) – The modified string.
- **cost_fn** (Optional[Callable[[Optional[str], Optional[str]], Union[int, float]]]) – A function returning the cost of performing an edit (see `Alignment.infer()`).

Return type *bistr*

Returns A *bistr* with the inferred alignment.

__getitem__ (*index*)

Indexing a *bistr* returns the *n*th character of the modified string:

```
>>> s = bistr('TEST').lower()
>>> s[1]
'e'
```

Slicing a *bistr* extracts a substring, complete with the matching part of the original string:

```
>>> s = bistr('TEST').lower()
>>> s[1:3]
bistr('ES', 'es', Alignment.identity(2))
```


Return type `Union[str, bistr]`

inverse()

Return type `bistr`

Returns

The inverse of this string, swapping the original and modified strings.

```
>>> s = bistr('HELLO WORLD').lower()
>>> s
bistr('HELLO WORLD', 'hello world', Alignment.identity(11))
>>> s.inverse()
bistr('hello world', 'HELLO WORLD', Alignment.identity(11))
```

chunks()

Return type `Iterable[bistr]`

Returns All the chunks of associated text in this string.

count(sub, start=None, end=None)

Like `str.count()`, counts the occurrences of `sub` in the string.

Return type `int`

find(sub, start=None, end=None)

Like `str.find()`, finds the position of `sub` in the string.

Return type `int`

find_bounds(sub, start=None, end=None)

Like `find()`, but returns both the start and end bounds for convenience.

Return type `Tuple[int, int]`

Returns The first `i, j` within `[start, end)` such that `self[i:j] == sub`, or `(-1, -1)` if not found.

rfind(sub, start=None, end=None)

Like `str.rfind()`, finds the position of `sub` in the string backwards.

Return type `int`

rfind_bounds(sub, start=None, end=None)

Like `rfind()`, but returns both the start and end bounds for convenience.

Return type `Tuple[int, int]`

Returns The last `i, j` within `[start, end)` such that `self[i:j] == sub`, or `(-1, -1)` if not found.

index(sub, start=None, end=None)

Like `str.index()`, finds the first position of `sub` in the string, otherwise raising a `ValueError`.

Return type `int`

index_bounds(sub, start=None, end=None)

Like `index()`, but returns both the start and end bounds for convenience. If the substring is not found, a `ValueError` is raised.

Return type `Tuple[int, int]`

Returns The first `i, j` within `[start, end)` such that `self[i:j] == sub`.

Raises `ValueError` if the substring is not found.

`rindex` (*sub*, *start=None*, *end=None*)

Like `str.index()`, finds the last position of *sub* in the string, otherwise raising a `ValueError`.

Return type `int`

`rindex_bounds` (*sub*, *start=None*, *end=None*)

Like `rindex()`, but returns both the start and end bounds for convenience. If the substring is not found, a `ValueError` is raised.

Return type `Tuple[int, int]`

Returns The last *i, j* within [*start*, *end*) such that `self[i:j] == sub`.

Raises `ValueError` if the substring is not found.

`startswith` (*prefix*, *start=None*, *end=None*)

Like `str.startswith()`, checks if the string starts with the given *prefix*.

Return type `bool`

`endswith` (*suffix*, *start=None*, *end=None*)

Like `str.endswith()`, checks if the string starts with the given *suffix*.

Return type `bool`

classmethod `join` (*iterable*)

Like `str.join()`, concatenates many (bi)strings together.

Return type `bistr`

`split` (*sep=None*, *maxsplit=-1*)

Like `str.split()`, splits this string on a separator.

Return type `List[bistr]`

`partition` (*sep*)

Like `str.partition()`, splits this string into three chunks on a separator.

Return type `Tuple[bistr, bistr, bistr]`

`rpartition` (*sep*)

Like `str.rpartition()`, splits this string into three chunks on a separator, searching from the end.

Return type `Tuple[bistr, bistr, bistr]`

`center` (*width*, *fillchar=' '*)

Like `str.center()`, pads the start and end of the string to center it.

Return type `bistr`

`ljust` (*width*, *fillchar=' '*)

Like `str.ljust()`, pads the end of the string to a fixed width.

Return type `bistr`

`rjust` (*width*, *fillchar=' '*)

Like `str.rjust()`, pads the start of the string to a fixed width.

Return type `bistr`

`casefold` ()

Computes the case folded form of this string. Case folding is used for case-insensitive operations, and the result may not be suitable for displaying to a user. For example:

```
>>> s = bistr('straße').casefold()
>>> s.modified
'strasse'
>>> s[4:6]
bistr('ß', 'ss')
```

Return type `bistr`

lower (*locale=None*)

Converts this string to lowercase. Unless you specify the *locale* parameter, the current system locale will be used.

```
>>> bistr('HELLO WORLD').lower()
bistr('HELLO WORLD', 'hello world', Alignment.identity(11))
>>> bistr('I').lower('en_US')
bistr('I', 'i')
>>> bistr('I').lower('tr_TR')
bistr('I', 'ı')
```

Return type `bistr`

upper (*locale=None*)

Converts this string to uppercase. Unless you specify the *locale* parameter, the current system locale will be used.

```
>>> bistr('hello world').upper()
bistr('hello world', 'HELLO WORLD', Alignment.identity(11))
>>> bistr('i').upper('en_US')
bistr('i', 'I')
>>> bistr('i').upper('tr_TR')
bistr('i', 'İ')
```

Return type `bistr`

title (*locale=None*)

Converts this string to title case. Unless you specify the *locale* parameter, the current system locale will be used.

```
>>> bistr('hello world').title()
bistr('hello world', 'Hello World', Alignment.identity(11))
>>> bistr('istanbul').title('en_US')
bistr('istanbul', 'Istanbul', Alignment.identity(8))
>>> bistr('istanbul').title('tr_TR')
bistr('istanbul', 'İstanbul', Alignment.identity(8))
```

Return type `bistr`

capitalize (*locale=None*)

Capitalize the first character of this string, and lowercase the rest. Unless you specify the *locale* parameter, the current system locale will be used.

```
>>> bistr('hello WORLD').capitalize()
bistr('hello WORLD', 'Hello world', Alignment.identity(11))
```

(continues on next page)

(continued from previous page)

```
>>> bistr('').capitalize('el_GR')
bistr('', '', Alignment.identity(2))
```

Return type `bistr`**swapcase** (*locale=None*)

Swap the case of every letter in this string. Unless you specify the *locale* parameter, the current system locale will be used.

```
>>> bistr('hello WORLD').swapcase()
bistr('hello WORLD', 'HELLO world', Alignment.identity(11))
```

Some Unicode characters, such as title-case ligatures and digraphs, don't have a case-swapped equivalent:

```
>>> bistr('epòta').swapcase('hr_HR')
bistr('epòta', 'EPÒTA', Alignment.identity(6))
```

In these cases, compatibility normalization may help:

```
>>> s = bistr('epòta')
>>> s = s.normalize('NFKC')
>>> s = s.swapcase('hr_HR')
>>> print(s)
('epòta' 'lJEPÒTA')
```

Return type `bistr`**expandtabs** (*tabsize=8*)

Like `str.expandtabs()`, replaces tab (`\t`) characters with spaces to align on multiples of *tabsize*.

Return type `bistr`**replace** (*old, new, count=None*)

Like `str.replace()`, replaces occurrences of *old* with *new*.

Return type `bistr`**sub** (*regex, repl*)

Like `re.sub()`, replaces all matches of *regex* with the replacement *repl*.

Parameters

- **regex** (`Union[str, Pattern[str]]`) – The regex to match. Can be a string pattern or a compiled regex.
- **repl** (`Union[str, Callable[[Match[str]], str]]`) – The replacement to use. Can be a string, which is interpreted as in `re.Match.expand()`, or a *callable*, which will receive each match and return the replacement string.

Return type `bistr`**strip** (*chars=None*)

Like `str.strip()`, removes leading and trailing characters (whitespace by default).

Return type `bistr`**rstrip** (*chars=None*)

Like `str.rstrip()`, removes leading characters (whitespace by default).

Return type `bistr`

`rstrip` (*chars=None*)

Like `str.rstrip()`, removes trailing characters (whitespace by default).

Return type `bistr`

`normalize` (*form*)

Like `unicodedata.normalize()`, applies a Unicode [normalization form](#). The choices for *form* are:

- 'NFC': Canonical Composition
- 'NFKC': Compatibility Composition
- 'NFD': Canonical Decomposition
- 'NFKD': Compatibilty Decomposition

Return type `bistr`

3.2 BistrBuilder

class `bistring.BistrBuilder` (*original*)

Bases: `object`

Bidirectionally transformed string builder.

A *BistrBuilder* builds a transformed version of a source string iteratively. Each builder has an immutable original string, a current string, and the in-progress modified string, with alignments between each. For example:

```
original: |The| |quick,| |brown| || |jumps| |over| |the| |lazy| ||
          | | | | | | | | \ \ \ \ \ \ \ \ \ \ \
current:  |The| |quick,| |brown| |fox| |jumps| |over| |the| |lazy| |dog|
          | | | / / /
modified: |the| |quick| |brown| ...
```

The modified string is built in pieces by calling `replace()` to change *n* characters of the current string into new ones in the modified string. Convenience methods like `skip()`, `insert()`, and `discard()` are implemented on top of this basic primitive.

```
>>> b = BistrBuilder('The quick, brown jumps over the lazy ')
>>> b.skip(17)
>>> b.peek(1)
''
>>> b.replace(1, 'fox')
>>> b.skip(21)
>>> b.peek(1)
''
>>> b.replace(1, 'dog')
>>> b.is_complete
True
>>> b.rewind()
>>> b.peek(3)
'The'
>>> b.replace(3, 'the')
>>> b.skip(1)
>>> b.peek(6)
'quick,'
```

(continues on next page)

```
>>> b.replace(6, 'quick')
>>> b.skip_rest()
>>> s = b.build()
>>> s.modified
'the quick brown fox jumps over the lazy dog'
```

Parameters `original` (`Union[str, bistr]`) – The string to start from.

property `original`

The original string being modified.

Return type `str`

property `current`

The current string before modifications.

Return type `str`

property `modified`

The modified string as built so far.

Return type `str`

property `alignment`

The alignment built so far from `self.current` to `self.modified`.

Return type `Alignment`

property `position`

The position of the builder in `self.current`.

Return type `int`

property `remaining`

The number of characters of the current string left to process.

Return type `int`

property `is_complete`

Whether we've completely processed the string. In other words, whether the modified string aligns with the end of the current string.

Return type `bool`

peek (`n`)

Peek at the next `n` characters of the original string.

Return type `str`

skip (`n`)

Skip the next `n` characters, copying them unchanged.

Return type `None`

skip_rest ()

Skip the rest of the string, copying it unchanged.

Return type `None`

insert (`string`)

Insert a substring into the string.

Return type `None`

discard (*n*)

Discard a portion of the original string.

Return type `None`

discard_rest ()

Discard the rest of the original string.

Return type `None`

replace (*n*, *repl*)

Replace the next *n* characters with a new string.

Return type `None`

append (*bs*)

Append a bistr. The original value of the bistr must match the current string being processed.

Return type `None`

skip_match (*regex*)

Skip a substring matching a regex, copying it unchanged.

Parameters **regex** (`Union[str, Pattern[str]]`) – The (possibly compiled) regular expression to match.

Return type `bool`

Returns Whether a match was found.

discard_match (*regex*)

Discard a substring that matches a regex.

Parameters **regex** (`Union[str, Pattern[str]]`) – The (possibly compiled) regular expression to match.

Return type `bool`

Returns Whether a match was found.

replace_match (*regex*, *repl*)

Replace a substring that matches a regex.

Parameters

- **regex** (`Union[str, Pattern[str]]`) – The (possibly compiled) regular expression to match.
- **repl** (`Union[str, Callable[[Match[str]], str]]`) – The replacement to use. Can be a string, which is interpreted as in `re.Match.expand()`, or a *callable*, which will receive each match and return the replacement string.

Return type `bool`

Returns Whether a match was found.

replace_next (*regex*, *repl*)

Replace the next occurrence of a regex.

Parameters

- **regex** (`Union[str, Pattern[str]]`) – The (possibly compiled) regular expression to match.
- **repl** (`Union[str, Callable[[Match[str]], str]]`) – The replacement to use.

Return type `bool`

Returns Whether a match was found.

replace_all (*regex*, *repl*)

Replace all occurrences of a regex.

Parameters

- **regex** (`Union[str, Pattern[str]]`) – The (possibly compiled) regular expression to match.
- **repl** (`Union[str, Callable[[Match[str]], str]]`) – The replacement to use.

Return type `None`

build ()

Build the *bistr*.

Return type `bistr`

Returns A *bistr* from the original string to the new modified string.

Raises `ValueError` if the modified string is not completely built yet.

rewind ()

Reset this builder to apply another transformation.

Raises `ValueError` if the modified string is not completely built yet.

Return type `None`

3.3 Alignment

class `bistring.Alignment` (*values*)

Bases: `object`

An alignment between two related sequences.

Consider this alignments between two strings:

```
|it's| |aligned|
|  \ \      |
|it is| |aligned|
```

An alignment stores all the indices that are known to correspond between the original and modified sequences. For the above example, it would be

```
>>> a = Alignment([
...     (0, 0),
...     (4, 5),
...     (5, 6),
...     (13, 13),
... ])
```

Alignments can be used to answer questions like, “what’s the smallest range of the original sequence that is guaranteed to contain this part of the modified sequence?” For example, the range (0, 5) (“it is”) is known to match the range (0, 4) (“it’s”) of the original sequence:

```
>>> a.original_bounds(0, 5)
(0, 4)
```

Results may be imprecise if the alignment is too coarse to match the exact inputs:


```
>>> a.original_bounds(0, 2)
(0, 4)
```

A more granular alignment like this:

```
|i|t|'s| |a|l|i|g|n|e|d|!|
| | | \ \ \ \ \ \ \ \ /
|i|t| i|s| |a|l|i|g|n|e|d|
```

```
>>> a = Alignment([
...     (0, 0), (1, 1), (2, 2), (4, 5), (5, 6), (6, 7), (7, 8),
...     (8, 9), (9, 10), (10, 11), (11, 12), (12, 13), (13, 13),
... ])
```

Can be more precise:

```
>>> a.original_bounds(0, 2)
(0, 2)
```

Parameters **values** (`Iterable[Tuple[int, int]]`) – The sequence of aligned indices. Each element should be a tuple (x, y) , where x is the original sequence position and y is the modified sequence position.

classmethod **identity** (**args*)

Create an identity alignment, which maps all intervals to themselves. You can pass the size of the sequence:

```
>>> Alignment.identity(5)
Alignment.identity(5)
```

or the start and end positions:

```
>>> Alignment.identity(1, 5)
Alignment.identity(1, 5)
```

or a range-like object (`range`, `slice`, or `Tuple[int, int]`):

```
>>> Alignment.identity(range(1, 5))
Alignment.identity(1, 5)
```

Return type `Alignment`

classmethod **infer** (*original, modified, cost_fn=None*)

Infer the alignment between two sequences with the lowest edit distance.

```
>>> Alignment.infer('color', 'color')
Alignment.identity(5)
>>> a = Alignment.infer('color', 'colour')
>>> # 'ou' -> 'o'
>>> a.original_bounds(3, 5)
(3, 4)
```

Warning: this operation has time complexity $O(N \cdot M)$, where N and M are the lengths of the original and modified sequences, and so should only be used for relatively short sequences.

Parameters

- **original** (`Sequence[~T]`) – The original sequence.
- **modified** (`Sequence[~U]`) – The modified sequence.
- **cost_fn** (`Optional[Callable[[Optional[~T], Optional[~U]], Union[int, float]]]`) – A function returning the cost of performing an edit. `cost_fn(a, b)` returns the cost of replacing *a* with *b*. `cost_fn(a, None)` returns the cost of deleting *a*, and `cost_fn(None, b)` returns the cost of inserting *b*. By default, all operations have cost 1 except replacing identical elements, which has cost 0.

Return type `Alignment`

Returns The inferred alignment.

`__getitem__` (*index*)

Indexing an alignment returns the *n*th pair of aligned positions:

```
>>> a = Alignment.identity(5)
>>> a[3]
(3, 3)
```

Slicing an alignment returns a new alignment with a subrange of its values:

```
>>> a[1:5]
Alignment.identity(1, 4)
```

Return type `Union[Tuple[int, int], Alignment]`

`shift` (*delta_o*, *delta_m*)

Shift this alignment.

Parameters

- **delta_o** (`int`) – The distance to shift the original sequence.
- **delta_m** (`int`) – The distance to shift the modified sequence.

Return type `Alignment`

Returns An alignment with all the positions shifted by the given amounts.

`original_bounds` (**args*)

Maps a subrange of the modified sequence to the original sequence. Can be called with either two arguments:

```
>>> a = Alignment.identity(5).shift(1, 0)
>>> a.original_bounds(1, 3)
(2, 4)
```

or with a range-like object:

```
>>> a.original_bounds(range(1, 3))
(2, 4)
```

With no arguments, returns the bounds of the entire original sequence:

```
>>> a.original_bounds()
(1, 6)
```

Return type `Tuple[int, int]`

Returns The corresponding bounds in the original sequence.

original_range (*args)

Like `original_bounds()`, but returns a `range`.

Return type `range`

original_slice (*args)

Like `original_bounds()`, but returns a `slice`.

Return type `slice`

modified_bounds (*args)

Maps a subrange of the original sequence to the modified sequence. Can be called with either two arguments:

```
>>> a = Alignment.identity(5).shift(1, 0)
>>> a.modified_bounds(2, 4)
(1, 3)
```

or with a range-like object:

```
>>> a.modified_bounds(range(2, 4))
(1, 3)
```

With no arguments, returns the bounds of the entire modified sequence:

```
>>> a.modified_bounds()
(0, 5)
```

Return type `Tuple[int, int]`

Returns The corresponding bounds in the modified sequence.

modified_range (*args)

Like `modified_bounds()`, but returns a `range`.

Return type `range`

modified_slice (*args)

Like `modified_bounds()`, but returns a `slice`.

Return type `slice`

slice_by_original (*args)

Slice this alignment by a span of the original sequence.

```
>>> a = Alignment.identity(5).shift(1, 0)
>>> a.slice_by_original(2, 4)
Alignment([(2, 1), (3, 2), (4, 3)])
```

Return type `Alignment`

Returns The slice of this alignment that corresponds with the given span of the original sequence.

slice_by_modified (*args)

Slice this alignment by a span of the modified sequence.

```
>>> a = Alignment.identity(5).shift(1, 0)
>>> a.slice_by_modified(1, 3)
Alignment([(2, 1), (3, 2), (4, 3)])
```

Return type Alignment

Returns The slice of this alignment that corresponds with the given span of the modified sequence.

compose (*other*)

Return type Alignment

Returns A new alignment equivalent to applying this one first, then the *other*.

inverse ()

Return type Alignment

Returns The inverse of this alignment, from the modified to the original sequence.

3.4 Tokenization

class bistring.**Token** (*text, start, end*)

Bases: object

A token extracted from a string.

Parameters

- **text** (Union[str, bistr]) – The text of this token.
- **start** (int) – The starting index of this token.
- **end** (int) – The ending index of this token.

text: bistr = None

The actual text of the token.

start: int = None

The start position of the token.

end: int = None

The end position of the token.

property original

The original value of this token.

Return type str

property modified

The modified value of this token.

Return type str

classmethod slice (*text, start, end*)

Create a Token from a slice of a bistr.

Parameters

- **text** (Union[str, bistr]) – The (bi)string to slice.

- **start** (`int`) – The starting index of the token.
- **end** (`int`) – The ending index of the token.

Return type `Token`

class `bistring.Tokenization` (`text`, `tokens`)

Bases: `object`

A string and its tokenization.

Parameters

- **text** (`Union[str, bistr]`) – The text from which the tokens have been extracted.
- **tokens** (`Iterable[Token]`) – The tokens extracted from the text.

text: `bistr = None`

The text that was tokenized.

alignment: `Alignment = None`

The alignment from text indices to token indices.

classmethod `infer` (`text`, `tokens`)

Infer a `Tokenization` from a sequence of tokens.

```
>>> tokens = Tokenization.infer('hello, world!', ['hello', 'world'])
>>> tokens[0]
Token(bistr('hello'), start=0, end=5)
>>> tokens[1]
Token(bistr('world'), start=7, end=12)
```

Due to the possibility of ambiguity, it is much better to use a `Tokenizer` or some other method of producing `Tokens` with their positions explicitly set.

Return type `Tokenization`

Returns The inferred tokenization, with token positions found by simple forward search.

Raises `ValueError` if the tokens can't be found in the source string.

__getitem__ (`index`)

Indexing a `Tokenization` returns the `n`th token:

```
>>> tokens = Tokenization.infer(
...     "The quick, brown fox",
...     ["The", "quick", "brown", "fox"],
... )
>>> tokens[0]
Token(bistr('The'), start=0, end=3)
```

Slicing a `Tokenization` returns a new one with the requested slice of tokens:

```
>>> tokens = tokens[1:-1]
>>> tokens[0]
Token(bistr('quick'), start=4, end=9)
```

Return type `Union[Token, Tokenization]`

substring (`*args`)

Map a span of tokens to the corresponding substring. With no arguments, returns the substring from the first to the last token.

Return type `bistr`

text_bounds (*args)

Map a span of tokens to the bounds of the corresponding text. With no arguments, returns the bounds from the first to the last token.

Return type `Tuple[int, int]`

original_bounds (*args)

Map a span of tokens to the bounds of the corresponding original text. With no arguments, returns the bounds from the first to the last token.

Return type `Tuple[int, int]`

bounds_for_text (*args)

Map a span of text to the bounds of the corresponding span of tokens.

Return type `Tuple[int, int]`

bounds_for_original (*args)

Map a span of original text to the bounds of the corresponding span of tokens.

Return type `Tuple[int, int]`

slice_by_text (*args)

Map a span of text to the corresponding span of tokens.

Return type `Tokenization`

slice_by_original (*args)

Map a span of the original text to the corresponding span of tokens.

Return type `Tokenization`

snap_text_bounds (*args)

Expand a span of text to align it with token boundaries.

Return type `Tuple[int, int]`

snap_original_bounds (*args)

Expand a span of original text to align it with token boundaries.

Return type `Tuple[int, int]`

3.5 Tokenizer

class `bistring.Tokenizer`

Bases: `abc.ABC`

Abstract base class for tokenizers.

abstract `tokenize` (*text*)

Tokenize some text.

Parameters `text` (`Union[str, bistr]`) – The text to tokenize, as either an *str* or *bistr*. A plain *str* should be converted to a *bistr* before processing.

Return type `Tokenization`

Returns A `Tokenization` holding the text and its tokens.

class `bistring.RegexTokenizer` (*regex*)
 Bases: `bistring._token.Tokenizer`

Breaks text into tokens based on a regex.

```
>>> tokenizer = RegexTokenizer(r'\w+')
>>> tokens = tokenizer.tokenize('the quick brown fox jumps over the lazy dog')
>>> tokens[0]
Token(bistr('the'), start=0, end=3)
>>> tokens[1]
Token(bistr('quick'), start=4, end=9)
```

Parameters `regex` (`Union[str, Pattern[str]]`) – A (possibly compiled) regular expression that matches tokens to extract.

class `bistring.SplittingTokenizer` (*regex*)
 Bases: `bistring._token.Tokenizer`

Splits text into tokens based on a regex.

```
>>> tokenizer = SplittingTokenizer(r'\s+')
>>> tokens = tokenizer.tokenize('the quick brown fox jumps over the lazy dog')
>>> tokens[0]
Token(bistr('the'), start=0, end=3)
>>> tokens[1]
Token(bistr('quick'), start=4, end=9)
```

Parameters `regex` (`Union[str, Pattern[str]]`) – A (possibly compiled) regular expression that matches the regions between tokens.

class `bistring.CharacterTokenizer` (*locale*)
 Bases: `bistring._token._IcuTokenizer`

Splits text into user-perceived characters/grapheme clusters.

```
>>> tokenizer = CharacterTokenizer('th_TH')
>>> tokens = tokenizer.tokenize('')
>>> tokens[0]
Token(bistr(''), start=0, end=2)
>>> tokens[1]
Token(bistr(''), start=2, end=4)
>>> tokens[2]
Token(bistr(''), start=4, end=5)
```

Parameters `locale` (`str`) – The name of the locale to use for computing user-perceived character boundaries.

class `bistring.WordTokenizer` (*locale*)
 Bases: `bistring._token._IcuTokenizer`

Splits text into words based on Unicode rules.

```
>>> tokenizer = WordTokenizer('en_US')
>>> tokens = tokenizer.tokenize('the quick brown fox jumps over the lazy dog')
>>> tokens[0]
Token(bistr('the'), start=0, end=3)
```

(continues on next page)

(continued from previous page)

```
>>> tokens[1]
Token(bistr('quick'), start=4, end=9)
```

Parameters **locale** (*str*) – The name of the locale to use for computing word boundaries.

class **bistring.SentenceTokenizer** (*locale*)

Bases: `bistring._token._IcuTokenizer`

Splits text into sentences based on Unicode rules.

```
>>> tokenizer = SentenceTokenizer('en_US')
>>> tokens = tokenizer.tokenize(
...     'Word, sentence, etc. boundaries are hard. Luckily, Unicode can help.'
... )
>>> tokens[0]
Token(bistr('Word, sentence, etc. boundaries are hard. '), start=0, end=42)
>>> tokens[1]
Token(bistr('Luckily, Unicode can help.'), start=42, end=68)
```

Parameters **locale** (*str*) – The name of the locale to use for computing sentence boundaries.

4.1 BiString

class BiString (*original, modified, alignment*)
exported from bistring

Implements:

- `Iterable()`

A bidirectionally transformed string.

A `BiString` can be constructed from only a single string, which will give it identical original and modified strings and an identity alignment:

```
new BiString("test");
```

You can also explicitly specify both the original and modified strings. The inferred alignment will be as course as possible:

```
new BiString("TEST", "test");
```

Finally, you can specify the alignment explicitly, if you know it:

```
new BiString("TEST", "test", Alignment.identity(4));
```

Arguments

- **original** (*string*) – The original string, before any modifications.
- **modified** (*undefined/string*) – The modified string, after any modifications.
- **alignment** (*alignment.Alignment*) – The alignment between the original and modified strings.

`BiString.[Symbol.iterator]()`
Iterates over the code points in the modified string.

Returns IterableIterator –

`BiString.alignment`
type: `alignment.Alignment`

The sequence alignment between *original* and *modified*.

`BiString.boundsOf` (*searchValue, fromIndex*)
Like `indexOf()`, but returns both the start and end positions for convenience.

Arguments

- **searchValue** (*string*) –
- **fromIndex** (*undefined|number*) –

Returns alignment.Bounds –

BiString.charAt (*pos*)

Like `String.prototype.charAt()`, returns a code unit as a string from the modified string.

Arguments

- **pos** (*number*) –

Returns string –

BiString.charCodeAt (*pos*)

Like `String.prototype.charCodeAt()`, returns a code unit as a number from the modified string.

Arguments

- **pos** (*number*) –

Returns number –

BiString.codePointAt (*pos*)

Like `String.prototype.codePointAt()`, returns a code point from the modified string.

Arguments

- **pos** (*number*) –

Returns number|undefined –

BiString.concat (*others*)

Concatenate this string together with one or more others. The additional strings can be either BiStrings or normal strings.

Arguments

- **others** (*bistring.AnyString[]*) –

Returns bistring.BiString –

BiString.endsWith (*searchString, position*)

Like `String.prototype.endsWith()`, returns whether this string ends with the given prefix.

Arguments

- **searchString** (*string*) –
- **position** (*undefined|number*) –

Returns boolean –

BiString.equals (*other*)

Arguments

- **other** (*bistring.BiString*) –

Returns boolean – Whether this BiString is equal to another.

BiString.from (*str*)

Create a *BiString* from a string-like object.

Arguments

- **str** (*bistring.AnyString*) – Either a *string* or a *BiString*.

Returns `bistring.BiString` – The input coerced to a *BiString*.

`BiString.indexOf (searchValue, fromIndex)`

Like `String.prototype.indexOf ()`, finds the first occurrence of a substring.

Arguments

- **searchValue** (*string*) –
- **fromIndex** (*undefined|number*) –

Returns `number` –

`BiString.infer (original, modified)`

Create a *BiString*, automatically inferring an alignment between the *original* and *modified* strings.

Arguments

- **original** (*string*) – The original string.
- **modified** (*string*) – The modified string.

Returns `bistring.BiString` –

`BiString.inverse ()`

Returns `bistring.BiString` – The inverse of this string, swapping the original and modified strings.

`BiString.join (items)`

Like `Array.prototype.join ()`, joins a sequence together with this *BiString* as the separator.

Arguments

- **items** (*Iterable*) –

Returns `bistring.BiString` –

`BiString.lastBoundsOf (searchValue, fromIndex)`

Like `lastIndexOf ()`, but returns both the start and end positions for convenience.

Arguments

- **searchValue** (*string*) –
- **fromIndex** (*undefined|number*) –

Returns `alignment.Bounds` –

`BiString.lastIndexOf (searchValue, fromIndex)`

Like `String.prototype.lastIndexOf ()`, finds the last occurrence of a substring.

Arguments

- **searchValue** (*string*) –
- **fromIndex** (*undefined|number*) –

Returns `number` –

`BiString.length`

type: `number`

The length of the modified string.

`BiString.match (regexp)`

Like `String.prototype.match ()`, returns the result of a regular expression match.

Arguments

- **regexp** (*RegExp*) –

Returns `RegExpMatchArray|null` –

`BiString.prototype.matchAll` (*regexp*)

Like `String.prototype.matchAll()`, returns an iterator over all regular expression matches.

Arguments

- **regexp** (*RegExp*) –

Returns `IterableIterator` –

`BiString.prototype.modified`

type: `string`

The current value of the string, after all modifications.

`BiString.prototype.normalize` (*form*)

Like `String.prototype.normalize()`, applies a Unicode normalization form.

Arguments

- **form** (`"NFC" | "NFD" | "NFKC" | "NFKD"`) – The normalization form to apply, one of “NFC”, “NFD”, “NFKC”, or “NFKD”.

Returns `bistring.BiString` –

`BiString.prototype.original`

type: `string`

The original string, before any modifications.

`BiString.prototype.padEnd` (*targetLength*, *padString*)

Like `String.prototype.padEnd()`, pads a string at the end to a target length.

Arguments

- **targetLength** (*number*) –
- **padString** (*string*) –

Returns `bistring.BiString` –

`BiString.prototype.padStart` (*targetLength*, *padString*)

Like `String.prototype.padStart()`, pads a string at the beginning to a target length.

Arguments

- **targetLength** (*number*) –
- **padString** (*string*) –

Returns `bistring.BiString` –

`BiString.prototype.replace` (*pattern*, *replacement*)

Like `String.prototype.replace()`, returns a new string with regex or fixed-string matches replaced.

Arguments

- **pattern** (*string|RegExp*) –
- **replacement** (*string|regex.Replacer*) –

Returns `bistring.BiString` –

`BiString.search(regex)`

Like `String.prototype.search()`, finds the position of the first match of a regular expression.

Arguments

- **regex** (*RegExp*) –

Returns `number` –

`BiString.searchBounds(regex)`

Like `search()`, but returns both the start and end positions for convenience.

Arguments

- **regex** (*RegExp*) –

Returns `alignment.Bounds` –

`BiString.slice(start, end)`

Extract a slice of this `BiString`, with similar semantics to `String.prototype.slice()`.

Arguments

- **start** (*number*) –
- **end** (*undefined|number*) –

Returns `bistring.BiString` –

`BiString.split(separator, limit)`

Like `String.prototype.split()`, splits this string into chunks using a separator.

Arguments

- **separator** (*string|RegExp*) –
- **limit** (*undefined|number*) –

Returns `bistring.BiString[]` –

`BiString.startsWith(searchString, position)`

Like `String.prototype.startsWith()`, returns whether this string starts with the given prefix.

Arguments

- **searchString** (*string*) –
- **position** (*undefined|number*) –

Returns `boolean` –

`BiString.substring(start, end)`

Extract a substring of this `BiString`, with similar semantics to `String.prototype.substring()`.

Arguments

- **start** (*number*) –
- **end** (*undefined|number*) –

Returns `bistring.BiString` –

`BiString.toLowerCase()`

Like `String.prototype.toLowerCase()`, converts a string to lowercase.

Returns `bistring.BiString` –

`BiString.toUpperCase()`

Like `String.prototype.toUpperCase()`, converts a string to uppercase.

Returns `bistring.BiString` –

`BiString.trim()`

Like `String.prototype.trim()`, returns a new string with leading and trailing whitespace removed.

Returns `bistring.BiString` –

`BiString.trimEnd()`

Like `String.prototype.trim()`, returns a new string with trailing whitespace removed.

Returns `bistring.BiString` –

`BiString.trimStart()`

Like `String.prototype.trim()`, returns a new string with leading whitespace removed.

Returns `bistring.BiString` –

4.2 BiStringBuilder

class `BiStringBuilder` (*original*)

exported from `builder`

Bidirectionally transformed string builder.

A *BiStringBuilder* builds a transformed version of a source string iteratively. Each builder has an immutable original string, a current string, and the in-progress modified string, with alignments between each. For example:

```
original: |The| |quick,| |brown| || |jumps| |over| |the| |lazy| ||
          | | |         | |         | | \ \         \ \         \ \         \ \         \ \
current:  |The| |quick,| |brown| |fox| |jumps| |over| |the| |lazy| |dog|
          | | |         / /         /
modified: |the| |quick| |brown| ...
```

The modified string is built in pieces by calling `replace()` to change *n* characters of the current string into new ones in the modified string. Convenience methods like `skip()`, `insert()`, and `discard()` are implemented on top of this basic primitive.

Construct a `BiStringBuilder`.

Arguments

- **original** (*bistring.AnyString*) – Either an original string or a `BiString` to start from.

`BiStringBuilder.alignment`

type: `alignment.Alignment`

The alignment as built so far from *this.current* to *this.modified*.

`BiStringBuilder.append(bs)`

Append a `BiString`. The original value of the `BiString` must match the current string being processed.

Arguments

- **bs** (*bistring.BiString*) –

`BiStringBuilder.build()`

Build the *BiString*.

Returns `bistring.BiString` –

`BiStringBuilder.current`

type: string

The current string before modifications.

`BiStringBuilder.discard(n)`

Discard a portion of the original string.

Arguments

- **n** (*number*) –

`BiStringBuilder.discardMatch(pattern)`

Discard a substring that matches a regex.

Arguments

- **pattern** (*RegExp*) – The pattern to match. Must have either the sticky flag, forcing it to match at the current position, or the global flag, finding the next match.

Returns boolean – Whether a match was found.

`BiStringBuilder.discardRest()`

Discard the rest of the original string.

`BiStringBuilder.insert(str)`

Insert a substring into the string.

Arguments

- **str** (*string*) –

`BiStringBuilder.isComplete`

type: boolean

Whether we've completely processed the string. In other words, whether the modified string aligns with the end of the current string.

`BiStringBuilder.modified`

type: string

The modified string as built so far.

`BiStringBuilder.original`

type: string

The original string being modified.

`BiStringBuilder.peek(n)`

Peek at the next few characters.

Arguments

- **n** (*number*) – The number of characters to peek at.

Returns string –

`BiStringBuilder.position`

type: number

The position of the builder in *this.current*.

`BiStringBuilder.remaining`

type: number

The number of characters of the current string left to process.

`BiStringBuilder.replace(n, str)`

Replace the next *n* characters with a new string.

Arguments

- **n** (*number*) –
- **str** (*bistring.AnyString*) –

`BiStringBuilder.replaceAll(pattern, replacement)`

Replace all occurrences of a regex, like `String.prototype.replace()`.

Arguments

- **pattern** (*RegExp*) – The pattern to match. The global flag (*/g*) must be set to get multiple matches.
- **replacement** (*string|regex.Replacer*) – The replacement string or function, as in `String.prototype.replace()`.

`BiStringBuilder.replaceMatch(pattern, replacement)`

Replace a substring that matches a regex.

Arguments

- **pattern** (*RegExp*) – The pattern to match. Must have either the sticky flag, forcing it to match at the current position, or the global flag, finding the next match.
- **replacement** (*string|regex.Replacer*) – The replacement string or function, as in `String.prototype.replace()`.

Returns boolean – Whether a match was found.

`BiStringBuilder.rewind()`

Reset this builder to apply another transformation.

`BiStringBuilder.skip(n)`

Skip the next *n* characters, copying them unchanged.

Arguments

- **n** (*number*) –

`BiStringBuilder.skipMatch(pattern)`

Skip a substring matching a regex, copying it unchanged.

Arguments

- **pattern** (*RegExp*) – The pattern to match. Must have either the sticky flag, forcing it to match at the current position, or the global flag, finding the next match.

Returns boolean – Whether a match was found.

`BiStringBuilder.skipRest()`

Skip the rest of the string, copying it unchanged.

4.3 Alignment

class Alignment (*values*)
exported from alignment

An alignment between two related sequences.

Consider this alignment between two strings:

```
|it's| |aligned!|
|   \ \   |
|it is| |aligned|
```

An alignment stores all the indices that are known to correspond between the original and modified sequences. For the above example, it would be

```
let a = new Alignment([
  [0, 0],
  [4, 5],
  [5, 6],
  [13, 13],
]);
```

Alignments can be used to answer questions like, “what’s the smallest range of the original sequence that is guaranteed to contain this part of the modified sequence?” For example, the range (0, 5) (“it is”) is known to match the range (0, 4) (“it’s”) of the original sequence.

```
console.log(a.original_bounds(0, 5));
// [0, 4]
```

Results may be imprecise if the alignment is too coarse to match the exact inputs:

```
console.log(a.original_bounds(0, 2));
// [0, 4]
```

A more granular alignment like this:

```
|i|t|'s| |a|l|l|i|g|n|e|d|!|
| | | \ \ \ \ \ \ \ \ \ /
|i|t| is| |a|l|l|i|g|n|e|d|
```

```
a = new Alignment([
  [0, 0], [1, 1], [2, 2], [4, 5], [5, 6], [6, 7], [7, 8],
  [8, 9], [9, 10], [10, 11], [11, 12], [12, 13], [13, 13],
]);
```

Can be more precise:

```
console.log(a.original_bounds(0, 2));
// [0, 2]
```

Create a new Alignment.

Arguments

- **values** (*Iterable*) – The pairs of indices to align. Each element should be a pair destructurable as $[x, y]$, where x is the original sequence position and y is the modified sequence position.

`Alignment.compose` (*other*)

Arguments

- **other** (*alignment.Alignment*) –

Returns alignment.Alignment – An alignment equivalent to applying *this* first, then *other*.

`Alignment.concat` (*others*)

Concatenate this alignment together with one or more others.

Arguments

- **others** (*alignment.Alignment []*) –

Returns alignment.Alignment –

`Alignment.equals` (*other*)

Arguments

- **other** (*alignment.Alignment*) –

Returns boolean – Whether this alignment is the same as *other*.

`Alignment.identity` (*size*)

Create an identity alignment of the given size, which maps all intervals to themselves.

Arguments

- **size** (*number*) – The size of the sequences to align.

Returns alignment.Alignment – The identity alignment for the bounds $[0, size]$.

`Alignment.identity` (*start, end*)

Create an identity alignment between the given bounds.

Arguments

- **start** (*number*) – The first index to align.
- **end** (*number*) – The last index to align.

Returns alignment.Alignment – The identity alignment for the bounds $[start, end]$.

`Alignment.identity` (*bounds*)

Create an identity alignment with the given bounds.

Arguments

- **bounds** (*alignment.Bounds*) – The bounds of the alignment (e.g. $[1, 5]$).

Returns alignment.Alignment – The identity alignment for the given bounds.

`Alignment.infer` (*original, modified, costFn*)

Infer the alignment between two sequences with the lowest edit distance.

Warning: this operation has time complexity $O(N * M)$, where N and M are the lengths of the original and modified sequences, and so should only be used for relatively short sequences.

Arguments

- **original** (*Iterable*) – The original sequence.
- **modified** (*Iterable*) – The modified sequence.

- **costFn** (*alignment.CostFn*) – A function returning the cost of performing an edit. *costFn(a, b)* returns the cost of replacing *a* with *b*. *costFn(a, undefined)* returns the cost of deleting *a*, and *costFn(undefined, b)* returns the cost of inserting *b*. By default, all operations have cost 1 except replacing identical elements, which has cost 0.

Returns alignment.Alignment – The inferred alignment.

`Alignment.inverse()`

Returns alignment.Alignment – The inverse of this alignment, from the modified to the original sequence.

`Alignment.length`

type: number

The number of entries in this alignment.

`Alignment.modifiedBounds()`

Returns alignment.Bounds – The bounds of the modified sequence.

`Alignment.modifiedBounds(start, end)`

Maps a subrange of the original sequence to the modified sequence.

Arguments

- **start** (*number*) – The start of the span in the original sequence.
- **end** (*number*) – The end of the span in the original sequence.

Returns alignment.Bounds – The bounds of the corresponding span in the modified sequence.

`Alignment.modifiedBounds(bounds)`

Maps a subrange of the original sequence to the modified sequence.

Arguments

- **bounds** (*alignment.Bounds*) – The bounds of the span in the original sequence.

Returns alignment.Bounds – The bounds of the corresponding span in the modified sequence.

`Alignment.originalBounds()`

Returns alignment.Bounds – The bounds of the original sequence.

`Alignment.originalBounds(start, end)`

Maps a subrange of the modified sequence to the original sequence.

Arguments

- **start** (*number*) – The start of the span in the modified sequence.
- **end** (*number*) – The end of the span in the modified sequence.

Returns alignment.Bounds – The bounds of the corresponding span in the original sequence.

`Alignment.originalBounds(bounds)`

Maps a subrange of the modified sequence to the original sequence.

Arguments

- **bounds** (*alignment.Bounds*) – The bounds of the span in the modified sequence.

Returns alignment.Bounds – The bounds of the corresponding span in the original sequence.

`Alignment.shift(deltaO, deltaM)`

Shift this alignment.

Arguments

- **deltaO** (*number*) – The distance to shift the original sequence.
- **deltaM** (*number*) – The distance to shift the modified sequence.

Returns alignment.Alignment – An alignment with all the positions shifted by the given amounts.

`Alignment.slice` (*start*, *end*)

Extract a slice of this alignment.

Arguments

- **start** (*undefined|number*) – The position to start from.
- **end** (*undefined|number*) – The position to end at.

Returns alignment.Alignment – The requested slice as a new *Alignment*.

`Alignment.sliceByModified` (*start*, *end*)

Slice this alignment by a span of the modified sequence.

Arguments

- **start** (*number*) – The start of the span in the modified sequence.
- **end** (*number*) – The end of the span in the modified sequence.

Returns alignment.Alignment – The requested slice of this alignment.

`Alignment.sliceByModified` (*bounds*)

Slice this alignment by a span of the modified sequence.

Arguments

- **bounds** (*alignment.Bounds*) – The bounds of the span in the modified sequence.

Returns alignment.Alignment – The requested slice of this alignment.

`Alignment.sliceByOriginal` (*start*, *end*)

Slice this alignment by a span of the original sequence.

Arguments

- **start** (*number*) – The start of the span in the original sequence.
- **end** (*number*) – The end of the span in the original sequence.

Returns alignment.Alignment – The requested slice of this alignment.

`Alignment.sliceByOriginal` (*bounds*)

Slice this alignment by a span of the original sequence.

Arguments

- **bounds** (*alignment.Bounds*) – The bounds of the span in the original sequence.

Returns alignment.Alignment – The requested slice of this alignment.

`Alignment.values`

type: `keyof:BiIndex[]`

The pairs of aligned indices in this alignment.

4.4 Tokenization

class Token (*text, start, end*)

exported from token

A token extracted from a string.

Create a token.

Arguments

- **text** (*bistring.AnyString*) – The text of this token.
- **start** (*number*) – The start position of the token.
- **end** (*number*) – The end position of the token.

Token.end

type: number

The end position of the token.

Token.modified

type: string

The modified value of the token.

Token.original

type: string

The original value of the token.

Token.slice (*text, start, end*)

Create a token from a slice of a string.

Arguments

- **text** (*bistring.AnyString*) – The text to slice.
- **start** (*number*) – The start index of the token.
- **end** (*number*) – The end index of the token.

Returns token.Token –

Token.start

type: number

The start position of the token.

Token.text

type: bistring.BiString

The actual text of the token.

class Tokenization (*text, tokens*)

exported from token

A string and its tokenization.

Create a *Tokenization*.

Arguments

- **text** (*bistring.AnyString*) – The text from which the tokens have been extracted.
- **tokens** (*Iterable*) – The tokens extracted from the text.

`Tokenization.alignment`

type: `alignment.Alignment`

The alignment between the text and the tokens.

`Tokenization.boundsForOriginal` (*start*, *end*)

Map a span of original text to the bounds of the corresponding span of tokens.

Arguments

- **start** (*number*) –
- **end** (*number*) –

Returns `alignment.Bounds` –

`Tokenization.boundsForText` (*start*, *end*)

Map a span of text to the bounds of the corresponding span of tokens.

Arguments

- **start** (*number*) –
- **end** (*number*) –

Returns `alignment.Bounds` –

`Tokenization.infer` (*text*, *tokens*)

Infer a *Tokenization* from a sequence of tokens.

Due to the possibility of ambiguity, it is much better to use a *Tokenizer()* or some other method of producing *Token()*s with their positions explicitly set.

Arguments

- **text** (*bistring.AnyString*) – The text that was tokenized.
- **tokens** (*Iterable*) – The extracted tokens.

Returns `token.Tokenization` – The inferred tokenization, with token positions found by simple forward search.

`Tokenization.length`

type: `number`

The number of tokens.

`Tokenization.originalBounds` (*start*, *end*)

Map a span of tokens to the bounds of the corresponding original text.

Arguments

- **start** (*undefined|number*) –
- **end** (*undefined|number*) –

Returns `alignment.Bounds` –

`Tokenization.slice` (*start*, *end*)

Compute a slice of this tokenization.

Arguments

- **start** (*undefined|number*) – The position to start from.
- **end** (*undefined|number*) – The position to end at.

Returns `token.Tokenization` – The requested slice as a new *Tokenization*.

Tokenization.**sliceByOriginal** (*start*, *end*)

Map a span of original text to the corresponding span of tokens.

Arguments

- **start** (*number*) –
- **end** (*number*) –

Returns token.Tokenization –

Tokenization.**sliceByText** (*start*, *end*)

Map a span of text to the corresponding span of tokens.

Arguments

- **start** (*number*) –
- **end** (*number*) –

Returns token.Tokenization –

Tokenization.**snapOriginalBounds** (*start*, *end*)

Expand a span of original text to align it with token boundaries.

Arguments

- **start** (*number*) –
- **end** (*number*) –

Returns alignment.Bounds –

Tokenization.**snapTextBounds** (*start*, *end*)

Expand a span of text to align it with token boundaries.

Arguments

- **start** (*number*) –
- **end** (*number*) –

Returns alignment.Bounds –

Tokenization.**substring** (*start*, *end*)

Map a span of tokens to the corresponding substring.

Arguments

- **start** (*undefined|number*) –
- **end** (*undefined|number*) –

Returns bistring.BiString –

Tokenization.**text**

type: bistring.BiString

The text that was tokenized.

Tokenization.**textBounds** (*start*, *end*)

Map a span of tokens to the bounds of the corresponding text.

Arguments

- **start** (*undefined|number*) –
- **end** (*undefined|number*) –

Returns alignment.Bounds –

`Tokenization.tokens`

type: keyof:Token[]

The tokens extracted from the text.

4.5 Tokenizer

class Tokenizer()

interface, exported from token

A tokenizer that produces *Tokenization()*s.

`Tokenizer.tokenize(text)`

Tokenize a string.

Arguments

- **text** (*bistring.AnyString*) – The text to tokenize, either a string or a *BiString()*.

Returns token.Tokenization – A *Tokenization()* holding the text and its tokens.

class RegExpTokenizer(pattern)

exported from token

Implements:

- `token.Tokenizer()`

Breaks text into tokens based on a `RegExp()`.

Create a *RegExpTokenizer*.

Arguments

- **pattern** (*RegExp*) – The regex that will match tokens.

`RegExpTokenizer.tokenize(text)`

Arguments

- **text** (*bistring.AnyString*) –

Returns token.Tokenization –

class SplittingTokenizer(pattern)

exported from token

Implements:

- `token.Tokenizer()`

Splits text into tokens based on a `RegExp()`.

Create a *SplittingTokenizer*.

Arguments

- **pattern** (*RegExp*) – A regex that matches the regions between tokens.

`SplittingTokenizer.tokenize(text)`

Arguments

- **text** (*bistring.AnyString*) –
Returns **token.Tokenization** –

INDICES AND TABLES

- genindex
- modindex
- search

Symbols

`__getitem__()` (*bistring.Alignment* method), 22
`__getitem__()` (*bistring.Tokenization* method), 25
`__getitem__()` (*bistring.bistr* method), 12

A

`alignment` (*bistring.bistr* attribute), 11
`alignment` (*bistring.Tokenization* attribute), 25
`Alignment` (*class in bistring*), 20
`alignment()` (*bistring.BistrBuilder* property), 18
`Alignment()` (*class*), 37
`Alignment.compose()` (*Alignment* method), 37
`Alignment.concat()` (*Alignment* method), 38
`Alignment.equals()` (*Alignment* method), 38
`Alignment.identity()` (*Alignment* method), 38
`Alignment.infer()` (*Alignment* method), 38
`Alignment.inverse()` (*Alignment* method), 39
`Alignment.length` (*Alignment* attribute), 39
`Alignment.modifiedBounds()` (*Alignment* method), 39
`Alignment.originalBounds()` (*Alignment* method), 39
`Alignment.shift()` (*Alignment* method), 39
`Alignment.slice()` (*Alignment* method), 40
`Alignment.sliceByModified()` (*Alignment* method), 40
`Alignment.sliceByOriginal()` (*Alignment* method), 40
`Alignment.values` (*Alignment* attribute), 40
`append()` (*bistring.BistrBuilder* method), 19

B

`bistr` (*class in bistring*), 11
`BistrBuilder` (*class in bistring*), 17
`BiString()` (*class*), 29
`BiString.alignment` (*BiString* attribute), 29
`BiString.boundsOf()` (*BiString* method), 29
`BiString.charAt()` (*BiString* method), 30
`BiString.charCodeAt()` (*BiString* method), 30
`BiString.codePointAt()` (*BiString* method), 30
`BiString.concat()` (*BiString* method), 30
`BiString.endsWith()` (*BiString* method), 30

`BiString.equals()` (*BiString* method), 30
`BiString.from()` (*BiString* method), 30
`BiString.indexOf()` (*BiString* method), 31
`BiString.infer()` (*BiString* method), 31
`BiString.inverse()` (*BiString* method), 31
`BiString.join()` (*BiString* method), 31
`BiString.lastBoundsOf()` (*BiString* method), 31
`BiString.lastIndexOf()` (*BiString* method), 31
`BiString.length` (*BiString* attribute), 31
`BiString.match()` (*BiString* method), 31
`BiString.matchAll()` (*BiString* method), 32
`BiString.modified` (*BiString* attribute), 32
`BiString.normalize()` (*BiString* method), 32
`BiString.original` (*BiString* attribute), 32
`BiString.padEnd()` (*BiString* method), 32
`BiString.padStart()` (*BiString* method), 32
`BiString.replace()` (*BiString* method), 32
`BiString.search()` (*BiString* method), 32
`BiString.searchBounds()` (*BiString* method), 33
`BiString.slice()` (*BiString* method), 33
`BiString.split()` (*BiString* method), 33
`BiString.startsWith()` (*BiString* method), 33
`BiString.substring()` (*BiString* method), 33
`BiString.toLowerCase()` (*BiString* method), 33
`BiString.toUpperCase()` (*BiString* method), 33
`BiString.trim()` (*BiString* method), 34
`BiString.trimEnd()` (*BiString* method), 34
`BiString.trimStart()` (*BiString* method), 34
`BiString.[Symbol.iterator]()` (*BiString.[Symbol]* method), 29
`BiStringBuilder()` (*class*), 34
`BiStringBuilder.alignment` (*BiStringBuilder* attribute), 34
`BiStringBuilder.append()` (*BiStringBuilder* method), 34
`BiStringBuilder.build()` (*BiStringBuilder* method), 34
`BiStringBuilder.current` (*BiStringBuilder* attribute), 34
`BiStringBuilder.discard()` (*BiStringBuilder* method), 35
`BiStringBuilder.discardMatch()` (*BiString-*

Builder method), 35
 BiStringBuilder.discardRest() (*BiStringBuilder method*), 35
 BiStringBuilder.insert() (*BiStringBuilder method*), 35
 BiStringBuilder.isComplete (*BiStringBuilder attribute*), 35
 BiStringBuilder.modified (*BiStringBuilder attribute*), 35
 BiStringBuilder.original (*BiStringBuilder attribute*), 35
 BiStringBuilder.peek() (*BiStringBuilder method*), 35
 BiStringBuilder.position (*BiStringBuilder attribute*), 35
 BiStringBuilder.remaining (*BiStringBuilder attribute*), 35
 BiStringBuilder.replace() (*BiStringBuilder method*), 35
 BiStringBuilder.replaceAll() (*BiStringBuilder method*), 36
 BiStringBuilder.replaceMatch() (*BiStringBuilder method*), 36
 BiStringBuilder.rewind() (*BiStringBuilder method*), 36
 BiStringBuilder.skip() (*BiStringBuilder method*), 36
 BiStringBuilder.skipMatch() (*BiStringBuilder method*), 36
 BiStringBuilder.skipRest() (*BiStringBuilder method*), 36
 bounds_for_original() (*bistring.Tokenization method*), 26
 bounds_for_text() (*bistring.Tokenization method*), 26
 build() (*bistring.BistrBuilder method*), 20

C

capitalize() (*bistring.bistr method*), 15
 casefold() (*bistring.bistr method*), 14
 center() (*bistring.bistr method*), 14
 CharacterTokenizer (*class in bistring*), 27
 chunks() (*bistring.bistr method*), 13
 compose() (*bistring.Alignment method*), 24
 count() (*bistring.bistr method*), 13
 current() (*bistring.BistrBuilder property*), 18

D

discard() (*bistring.BistrBuilder method*), 18
 discard_match() (*bistring.BistrBuilder method*), 19
 discard_rest() (*bistring.BistrBuilder method*), 19

E

end (*bistring.Token attribute*), 24

endswith() (*bistring.bistr method*), 14
 expandtabs() (*bistring.bistr method*), 16

F

find() (*bistring.bistr method*), 13
 find_bounds() (*bistring.bistr method*), 13

I

identity() (*bistring.Alignment class method*), 21
 index() (*bistring.bistr method*), 13
 index_bounds() (*bistring.bistr method*), 13
 infer() (*bistring.Alignment class method*), 21
 infer() (*bistring.bistr class method*), 11
 infer() (*bistring.Tokenization class method*), 25
 insert() (*bistring.BistrBuilder method*), 18
 inverse() (*bistring.Alignment method*), 24
 inverse() (*bistring.bistr method*), 13
 is_complete() (*bistring.BistrBuilder property*), 18

J

join() (*bistring.bistr class method*), 14

L

ljust() (*bistring.bistr method*), 14
 lower() (*bistring.bistr method*), 15
 lstrip() (*bistring.bistr method*), 16

M

modified (*bistring.bistr attribute*), 11
 modified() (*bistring.BistrBuilder property*), 18
 modified() (*bistring.Token property*), 24
 modified_bounds() (*bistring.Alignment method*), 23
 modified_range() (*bistring.Alignment method*), 23
 modified_slice() (*bistring.Alignment method*), 23

N

normalize() (*bistring.bistr method*), 17

O

original (*bistring.bistr attribute*), 11
 original() (*bistring.BistrBuilder property*), 18
 original() (*bistring.Token property*), 24
 original_bounds() (*bistring.Alignment method*), 22
 original_bounds() (*bistring.Tokenization method*), 26
 original_range() (*bistring.Alignment method*), 23
 original_slice() (*bistring.Alignment method*), 23

P

partition() (*bistring.bistr method*), 14
 peek() (*bistring.BistrBuilder method*), 18

position() (*bistring.BistrBuilder* property), 18

R

RegexTokenizer() (*class*), 44

RegexTokenizer.tokenize() (*RegexTokenizer* method), 44

RegexTokenizer (*class in bistring*), 26

remaining() (*bistring.BistrBuilder* property), 18

replace() (*bistring.bistr* method), 16

replace() (*bistring.BistrBuilder* method), 19

replace_all() (*bistring.BistrBuilder* method), 20

replace_match() (*bistring.BistrBuilder* method), 19

replace_next() (*bistring.BistrBuilder* method), 19

rewind() (*bistring.BistrBuilder* method), 20

rfind() (*bistring.bistr* method), 13

rfind_bounds() (*bistring.bistr* method), 13

rindex() (*bistring.bistr* method), 14

rindex_bounds() (*bistring.bistr* method), 14

rjust() (*bistring.bistr* method), 14

rpartition() (*bistring.bistr* method), 14

rstrip() (*bistring.bistr* method), 17

S

SentenceTokenizer (*class in bistring*), 28

shift() (*bistring.Alignment* method), 22

skip() (*bistring.BistrBuilder* method), 18

skip_match() (*bistring.BistrBuilder* method), 19

skip_rest() (*bistring.BistrBuilder* method), 18

slice() (*bistring.Token* class method), 24

slice_by_modified() (*bistring.Alignment* method), 23

slice_by_original() (*bistring.Alignment* method), 23

slice_by_original() (*bistring.Tokenization* method), 26

slice_by_text() (*bistring.Tokenization* method), 26

snap_original_bounds() (*bistring.Tokenization* method), 26

snap_text_bounds() (*bistring.Tokenization* method), 26

split() (*bistring.bistr* method), 14

SplittingTokenizer (*class in bistring*), 27

SplittingTokenizer() (*class*), 44

SplittingTokenizer.tokenize() (*SplittingTokenizer* method), 44

start (*bistring.Token* attribute), 24

startswith() (*bistring.bistr* method), 14

strip() (*bistring.bistr* method), 16

sub() (*bistring.bistr* method), 16

substring() (*bistring.Tokenization* method), 25

swapcase() (*bistring.bistr* method), 16

T

text (*bistring.Token* attribute), 24

text (*bistring.Tokenization* attribute), 25

text_bounds() (*bistring.Tokenization* method), 26

title() (*bistring.bistr* method), 15

Token (*class in bistring*), 24

Token() (*class*), 41

Token.end (*Token* attribute), 41

Token.modified (*Token* attribute), 41

Token.original (*Token* attribute), 41

Token.slice() (*Token* method), 41

Token.start (*Token* attribute), 41

Token.text (*Token* attribute), 41

Tokenization (*class in bistring*), 25

Tokenization() (*class*), 41

Tokenization.alignment (*Tokenization* attribute), 41

Tokenization.boundsForOriginal() (*Tokenization* method), 42

Tokenization.boundsForText() (*Tokenization* method), 42

Tokenization.infer() (*Tokenization* method), 42

Tokenization.length (*Tokenization* attribute), 42

Tokenization.originalBounds() (*Tokenization* method), 42

Tokenization.slice() (*Tokenization* method), 42

Tokenization.sliceByOriginal() (*Tokenization* method), 42

Tokenization.sliceByText() (*Tokenization* method), 43

Tokenization.snapOriginalBounds() (*Tokenization* method), 43

Tokenization.snapTextBounds() (*Tokenization* method), 43

Tokenization.substring() (*Tokenization* method), 43

Tokenization.text (*Tokenization* attribute), 43

Tokenization.textBounds() (*Tokenization* method), 43

Tokenization.tokens (*Tokenization* attribute), 44

tokenize() (*bistring.Tokenizer* method), 26

Tokenizer (*class in bistring*), 26

Tokenizer() (*class*), 44

Tokenizer.tokenize() (*Tokenizer* method), 44

U

upper() (*bistring.bistr* method), 15

W

WordTokenizer (*class in bistring*), 27