
Graph Transliterator Documentation

Release 1.2.4

A. Sean Pue

Oct 16, 2023

CONTENTS

1	Transliteration... What? Why?	3
2	Features	5
3	Sample Code and Graph	7
3.1	Get It Now	7
3.2	Citation	7
3.3	Indices and tables	86
	Python Module Index	87
	Index	89


```
[ pypi v1.2.4 ]  
[ docs failing ]
```

```
[ python 3.9 | 3.10 | 3.11 ]
```

```
[ DOI 10.5281/zenodo.3558385 ]
```

```
[ JOSS 10.21105/joss.01717 ]
```

A graph-based transliteration tool that lets you convert the symbols of one language or script to those of another using rules that you define.

- Free software: MIT license
- Documentation: <https://graphtransliterator.readthedocs.io>
- Repository: <https://github.com/seanpue/graphtransliterator>

TRANSLITERATION... WHAT? WHY?

Moving text or data from one script or encoding to another is a common problem:

- Many languages are written in multiple scripts, and many people can only read one of them. Moving between them can be a complex but necessary task in order to make texts accessible.
- The identification of names and locations, as well as machine translation, benefit from transliteration.
- Library systems often require metadata be in particular forms of romanization in addition to the original script.
- Linguists need to move between different methods of phonetic transcription.
- Documents in legacy fonts must now be converted to contemporary Unicode ones.
- Complex-script languages are frequently approached in natural language processing and in digital humanities research through transliteration, as it provides disambiguating information about pronunciation, morphological boundaries, and unwritten elements not present in the original script.

Graph Transliterator abstracts transliteration, offering an “easy reading” method for developing transliterators that does not require writing a complex program. It also contains bundled transliterators that are rigorously tested. These can be expanded to handle many transliteration tasks.

Contributions are very welcome!

FEATURES

- Provides a transliteration tool that can be configured to convert the tokens of an input string into an output string using:
 - user-defined types of input **tokens** and **token classes**
 - **transliteration rules** based on:
 - * a sequence of input tokens
 - * specific input tokens that precede or follow the token sequence
 - * classes of input tokens preceding or following specified tokens
 - **“on match” rules** for output to be inserted between transliteration rules involving particular token classes
 - defined rules for **whitespace**, including its optional consolidation
- Can be setup using:
 - an **“easy reading” YAML** format that lets you quickly craft settings for the transliteration tool
 - a **JSON** dump of a transliterator (quicker!)
 - **“direct”** settings, perhaps passed programmatically, using a dictionary
- **Automatically orders rules** by the number of tokens in a transliteration rule
- **Checks for ambiguity** in transliteration rules
- Can provide **details** about each transliteration rule match
- Allows **optional matching of all possible rules** in a particular location
- Permits **pruning of rules** with certain productions
- **Validates**, as well as **serializes** to and **deserializes** from JSON and Python data types, using accessible **marshmallow** schemas
- Provides **full support for Unicode**, including Unicode **character names** in the “easy reading” YAML format
- Constructs and uses a **directed tree** and performs a **best-first search** to find the most specific transliteration rule in a given context
- Includes **bundled transliterators** that *you* can add to that check for full test coverage of the nodes and edges of the internal graph and any “on match” rules
- Includes a command-line interface to perform transliteration and other tasks

SAMPLE CODE AND GRAPH

```
from graphtransliterator import GraphTransliterator
GraphTransliterator.from_yaml("""
  tokens:
    h: [consonant]
    i: [vowel]
    " ": [whitespace]
  rules:
    h: \N{LATIN SMALL LETTER TURNED I}
    i: \N{LATIN SMALL LETTER TURNED H}
    <whitespace> i: \N{LATIN CAPITAL LETTER TURNED H}
    (<whitespace> h) i: \N{LATIN SMALL LETTER TURNED H}!
  onmatch_rules:
    - <whitespace> + <consonant>: ;
  whitespace:
    default: " "
    consolidate: true
    token_class: whitespace
  metadata:
    title: "Upside Down Greeting Transliterator"
    version: "1.0.0"
""").transliterate("hi")
```

```
'iᴉʞ!'
```

3.1 Get It Now

```
$ pip install -U graphtransliterator
```

3.2 Citation

To cite Graph Transliterator, please use:

Pue, A. Sean (2019). Graph Transliterator: A graph-based transliteration tool. Journal of Open Source Software, 4(44), 1717, <https://doi.org/10.21105/joss.01717>

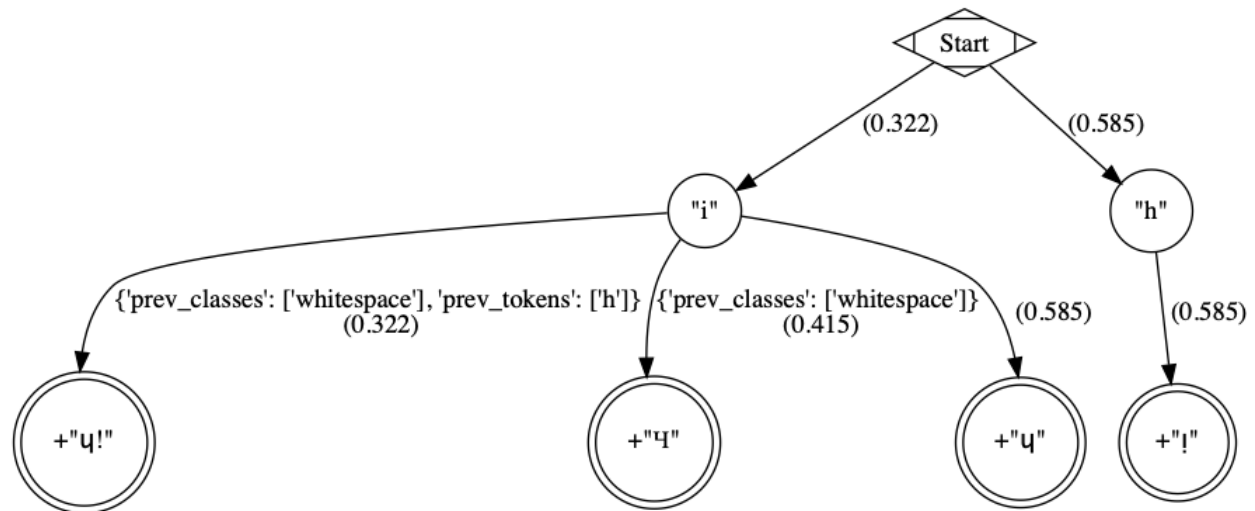


Fig. 1: Sample directed tree created by Graph Transliterators. The *rule* nodes are in double circles, and *token* nodes are single circles. The numbers are the cost of the particular edge, and less costly edges are searched first. Previous token classes and previous tokens that must be present are found as constraints on the edges incident to the terminal leaf *rule* nodes.

3.2.1 Installation

Stable release

To install Graph Transliterators, run this command in your terminal:

```
$ pip install graphtransliterators
```

This is the preferred method to install Graph Transliterators, as it will always install the most recent stable release.

If you don't have [pip](#) installed, this [Python installation guide](#) can guide you through the process.

From sources

The sources for Graph Transliterators can be downloaded from the [Github repo](#).

You can either clone the public repository:

```
$ git clone git://github.com/seanpue/graphtransliterators
```

Or download the [tarball](#):

```
$ curl -OL https://github.com/seanpue/graphtransliterators/tarball/master
```

Once you have a copy of the source, you can install it with:

```
$ python setup.py install
```

Required modules

Graph Transliterator requires three Python modules, `click`, `marshmallow` and `pyyaml`. These modules will be installed automatically using the methods described above.

3.2.2 Usage

To use Graph Transliterator in a project:

```
1 from graphtransliterator import GraphTransliterator
```

Overview

Graph Transliterator requires that you first configure a `GraphTransliterator`. Then you can transliterate an input string using `transliterate()`. There are a few additional methods that can be used to extract information for specific use cases, such as details about which rules were matched.

Configuration

Graph Transliterator takes the following parameters:

1. The acceptable types of **tokens** in the input string as well as any associated **token classes**.
2. The **transliteration rules** for the transformation of the input string.
3. Rules for dealing with **whitespace**.
4. “**On match**” **rules** for strings to be inserted in particular contexts right before a transliteration rule’s output is added (optional).
5. **Metadata** settings for the transliterator (optional).

Initialization

Defining the rules for transliteration can be difficult, especially when dealing with complex scripts. That is why Graph Transliterator uses an “easy reading” format that allows you to enter the transliteration rules in the popular **YAML** format, either from a string (using `from_yaml()`) or by reading from a file or stream (`GraphTransliterator.from_yaml_file()`). You can also initialize from the loaded contents of **YAML** (`GraphTransliterator.from_easyreading_dict()`).

Here is a quick sample that parameterizes `GraphTransliterator` using an easy reading **YAML** string (with comments):

```
2 yaml_ = """
3     tokens:
4         a: [vowel]                # type of token ("a") and its class (vowel)
5         bb: [consonant, b_class]  # type of token ("bb") and its classes (consonant, b_
6         ' ': [wb]                # type of token (" ") and its class ("wb", for wordbreak)
7     rules:
8         a: A                      # transliterate "a" to "A"
9         bb: B                    # transliterate "bb" to "B"
10        a a: <2AS>                # transliterate ("a", "a") to "<2AS>"
```

(continues on next page)

(continued from previous page)

```

11     ' ': ' ' # transliterate ' ' to ' '
12     whitespace:
13         default: " " # default whitespace token
14         consolidate: false # whitespace should not be consolidated
15         token_class: wb # whitespace token class
16     """
17     gt_one = GraphTransliterator.from_yaml(yaml_)
18     gt_one.transliterate('a')

```

```
'A'
```

```
19 gt_one.transliterate('bb')
```

```
'B'
```

```
20 gt_one.transliterate('aabb')
```

```
'<2AS>B'
```

The example above shows a very simple transliterator that replaces the input token “a” with “A”, “bb” with “B”, “ ” with “ ”, and two “a” in a row with “<2AS>”. It does not consolidate whitespace, and treats “ ” as its default whitespace token. Tokens contain strings of one or more characters.

Input Tokens and Token Class Settings

During transliteration, Graph Transliterator first attempts to convert the input string into a list of tokens. This is done internally using `GraphTransliterator.tokenize()`:

```
21 gt_one.tokenize('abba')
```

```
[' ', 'a', 'bb', 'a', ' ']
```

Note that the default whitespace token is added to the start and end of the input tokens.

Tokens can be more than one character, and longer tokens are matched first:

```

22 yaml_ = """
23     tokens:
24         a: [] # "a" token with no classes
25         aa: [] # "aa" token with no classes
26         ' ': [wb] # " " token and its class ("wb", for wordbreak)
27     rules:
28         aa: <DOUBLE_A> # transliterate "aa" to "<DOUBLE_A>"
29         a: <SINGLE_A> # transliterate "a" to "<SINGLE_A>"
30     whitespace:
31         default: " " # default whitespace token
32         consolidate: false # whitespace should not be consolidated
33         token_class: wb # whitespace token class
34     """
35     gt_two = GraphTransliterator.from_yaml(yaml_)
36     gt_two.transliterate('a')

```

```
'<SINGLE_A>'
```

```
37 gt_two.transliterate('aa')
```

```
'<DOUBLE_A>'
```

```
38 gt_two.transliterate('aaa')
```

```
'<DOUBLE_A><SINGLE_A>'
```

Here the input “aaa” is transliterated as “<DOUBLE_A><SINGLE_A>”, as the longer token “aa” is matched before “a”.

Tokens can be assigned zero or more classes. Each class is a string of your choice. These classes are used in transliteration rules. In YAML they are defined as a dictionary, but internally the rules are stored as a dictionary of token strings keyed to a set of token classes. They can be accessed using `GraphTransliterator.tokens`:

```
39 gt_two.tokens
```

```
{'a': set(), 'aa': set(), ' ': {'wb'}}
```

Transliteration Rules

Graph Transliterator can handle a variety of transliteration tasks. To do so, it uses transliteration rules that contain **match settings** for particular tokens in specific contexts and also a resulting **production**, or string to be appended to the output string.

Match Settings

Transliteration rules contain the following parameters (ordered by where they would appear in a list of tokens):

- **previous token classes** : a list of token classes (optional)
- **previous tokens** : a list of tokens (optional)
- **tokens** : a list of tokens
- **next tokens** : a list of tokens (optional)
- **next token classes** : a list of token classes (optional)

One or more (**tokens**) must be matched in a particular location. However, specific tokens can be required before (**previous tokens**) or behind (**next tokens**) those tokens. Additionally, particular token classes can be required before (**previous token classes**) and behind (**next token classes**) all of the specific tokens required (previous tokens, tokens, next tokens).

Depending on their complexity, these match conditions can be entered using the “easy reading” format in the following ways.

If there are no required lookahead or lookbehind tokens, the rule can be as follows:

```
rules:
  a a: aa # two tokens (a,a), with production "production_aa"
```

If, in an addition to tokens, there are specific previous or following tokens that must be matched, the rule can be entered as:

```
tokens:
  a: []
  b: []
  c: []
  d: []
rules:
  a (b): a_before_b # matches token 'a' with the next token 'b'
  (c) a: a_after_c # matches token 'a' when the previous token is 'c'
  a (b c): a_before_b_and_c # matches token 'a' when next tokens are 'b' then 'c'
  (d) a (b c): a_after_d_and_before_b,c # matches the token 'a' after 'd' and
↳before 'b' and 'c'
```

Token class names are indicated between angular brackets (“<classname>”). If preceding and following tokens are not required but classes are, these can be entered as follows:

```
tokens:
  a: []
  b: [class_b]
  c: []
  ' ': [wb]
rules:
  c <class_b>: c_after_class_b # match token 'c' before a token of class 'class_b'
  <class_b> a: a_before_class_b # match token 'a' after a token of class 'class_b'
  <class_b> a <class_b>: a_between_class_b # match token 'a' between tokens of class
↳'class_b'
```

If token classes must precede or follow specific tokens, these can be entered as:

```
tokens:
  a: []
  b: []
  c: [class_c]
  d: [class_d]
  ' ': [wb]
rules:
  d (b <class_c>): a_before_b_and_class_c # match token 'd' before 'b' and a token of
↳class 'class_c'
  (<class_c> b) a: a_after_b_and_class_c # match token 'a' after 'b' and a token of
↳class 'class_c'
  (<class_c> d) a (b <class_c> <class_d>): x # match 'a' after token of 'class_c' and
↳'d' and before a token of 'class_c' and of 'class_d'
whitespace:
  default: ' '
  token_class: wb
  consolidate: false
```

Automatic Ordering of Transliteration Rules

Graph Transliterator automatically orders the transliteration rules based on the number of tokens required by the rule. It picks the rule requiring the longest match in a given context. It does so by assigning a cost to each transliteration rule that decreases depending on the number of tokens required by the rule. More tokens decreases the cost of a rule causing it to be matched first:

```
40 yam1_ = """
41 tokens:
```

(continues on next page)

(continued from previous page)

```

42     a: []
43     b: []
44     c: [class_of_c]
45     ' ': [wb]
46 rules:
47     a: <<A>>
48     a b: <<AB>>
49     b: <<B>>
50     c: <<C>>
51     ' ': _
52     <class_of_c> a b: <<AB_after_C>>
53 whitespace:
54     default: " "
55     consolidate: false
56     token_class: wb
57 """
58 gt_three = GraphTransliterator.from_yaml(yaml_)
59 gt_three.transliterate("ab") # should match rule "a b"

```

```
'<<AB>>'
```

```
60 gt_three.transliterate("cab") # should match rules: "c", and "<class_of_c> a b"
```

```
'<<C>><<AB_after_C>>'
```

Internally, Graph Transliterator uses a special TransliterationRule class. These can be accessed using GraphTransliterator.rules. Rules are sorted by cost, lowest to highest:

```
61 gt_three.rules
```

```

[TransliterationRule(production='<<AB_after_C>>', prev_classes=['class_of_c'], prev_
→tokens=None, tokens=['a', 'b'], next_tokens=None, next_classes=None, cost=0.
→32192809488736235),
  TransliterationRule(production='<<AB>>', prev_classes=None, prev_tokens=None,
→tokens=['a', 'b'], next_tokens=None, next_classes=None, cost=0.41503749927884376),
  TransliterationRule(production='<<A>>', prev_classes=None, prev_tokens=None, tokens=[
→'a'], next_tokens=None, next_classes=None, cost=0.5849625007211562),
  TransliterationRule(production='<<B>>', prev_classes=None, prev_tokens=None, tokens=[
→'b'], next_tokens=None, next_classes=None, cost=0.5849625007211562),
  TransliterationRule(production='<<C>>', prev_classes=None, prev_tokens=None, tokens=[
→'c'], next_tokens=None, next_classes=None, cost=0.5849625007211562),
  TransliterationRule(production='_', prev_classes=None, prev_tokens=None, tokens=['
→'], next_tokens=None, next_classes=None, cost=0.5849625007211562)]

```

Whitespace Settings

Whitespace is often very important in transliteration tasks, as the form of many letters may change at the start or end of words, as in the right-to-left Perso-Arabic and left-to-right Indic scripts. Therefore, Graph Transliterator requires the following **whitespace settings**:

- the **default** whitespace token
- the whitespace **token class**
- whether or not to **consolidate** whitespace

A *whitespace token and token class must be defined for any Graph Transliterator*. A whitespace character is added temporarily to the start and end of the input tokens during the transliteration process.

The `consolidate` option may be useful in particular transliteration tasks. It replaces any sequential whitespace tokens in the input string with the default whitespace character. At the start and end of input, it removes any whitespace:

```
62 yaml_ = """
63     tokens:
64         a: []
65         ' ': [wb]
66     rules:
67         <wb> a: _A
68         a <wb>: A_
69         <wb> a <wb>: _A_
70         a: a
71         ' ': ' '
72     whitespace:
73         default: " "      # default whitespace token
74         consolidate: true  # whitespace should be consolidated
75         token_class: wb    # whitespace token class
76 """
77 gt = GraphTransliterator.from_yaml(yaml_)
78 gt.transliterate('a')    # whitespace present at start of string
```

```
'_A_'
```

```
79 gt.transliterate('aa')  # whitespace present at start and end of string
```

```
'_AA_'
```

```
80 gt.transliterate(' a')  # consolidate removes whitespace at start of string
```

```
'_A_'
```

```
81 gt.transliterate('a ')  # consolidate removes whitespace at end of string
```

```
'_A_'
```

Whitespace settings are stored internally as `WhitespaceRules` and can be accessed using `GraphTransliterator.whitespace`:

```
82 gt.whitespace
```

```
WhitespaceRules(default=' ', token_class='wb', consolidate=True)
```

On Match Rules

Graph Transliterator allows strings to be inserted right before the productions of transliteration rules. These take as parameters:

- a list of **previous token classes**, preceding the location of the transliteration rule match
- a list of **next token classes**, from the index of the transliteration rule match
- a **production** string to insert

In the easy reading YAML format, the `onmatch_rules` are a list of dictionaries. The key consists of the token class names in angular brackets (“<classname>”), and the previous classes to match are separated from the following classes by a “+”. The production is the value of the dictionary:

```
83 yaml_ = """
84     tokens:
85         a: [vowel]
86         ' ': [wb]
87     rules:
88         a: A
89         ' ': ' '
90     whitespace:
91         default: " "
92         consolidate: false
93         token_class: wb
94     onmatch_rules:
95         - <vowel> + <vowel>: ', ' # add a comma between vowels
96     """
97 gt = GraphTransliterator.from_yaml(yaml_)
98 gt.transliterate('aa')
```

```
'A,A'
```

On Match rules are stored internally as a `OnMatchRule` and can be accessed using `GraphTransliterator.onmatch_rules`:

```
99 gt.onmatch_rules
```

```
[OnMatchRule(prev_classes=['vowel'], next_classes=['vowel'], production=', ')]
```

Metadata

Graph Transliterator allows for the storage of metadata as another input parameter, `metadata`. It is a dictionary, and fields can be added to it:

```
100 yaml_ = """
101     tokens:
102         a: []
103         ' ': [wb]
104     rules:
105         a: A
```

(continues on next page)

(continued from previous page)

```

106     ' ': ' '
107     whitespace:
108         default: " "
109         consolidate: false
110         token_class: wb
111     metadata:
112         author: Author McAuthorson
113         version: 0.1.1
114         description: A sample Graph Transliterator
115     """
116 gt = GraphTransliterator.from_yaml(yaml_)
117 gt.metadata

```

```

{'author': 'Author McAuthorson',
 'version': '0.1.1',
 'description': 'A sample Graph Transliterator'}

```

Unicode Support

Graph Transliterator allows Unicode characters to be specified by name, including in YAML files, using the format “\N{UNICODE CHARACTER NAME}” or “\u{#####}” (where ##### is the hexadecimal character code):

```

118 yaml_ = """
119     tokens:
120         b: []
121         c: []
122         ' ': [wb]
123     rules:
124         b: \N{LATIN CAPITAL LETTER B}
125         c: \u0043      # hexadecimal Unicode character code for 'C'
126         ' ': ' '
127     whitespace:
128         default: " "
129         consolidate: false
130         token_class: wb
131     """
132 gt = GraphTransliterator.from_yaml(yaml_)
133 gt.transliterate('b')

```

```
'B'
```

```
134 gt.transliterate('c')
```

```
'C'
```

Configuring Directly

In addition to using `GraphTransliterator.from_yaml()` and `GraphTransliterator.from_yaml_file()`, Graph Transliterator can also be configured and initialized directly using basic Python types passed as dictionary to `GraphTransliterator.from_dict()`

```

135 settings = {
136     'tokens': {'a': ['vowel'],
137               ' ': ['wb']},
138     'rules': [
139         {'production': 'A', 'tokens': ['a']},
140         {'production': ' ', 'tokens': [' ']}],
141     'onmatch_rules': [
142         {'prev_classes': ['vowel'],
143          'next_classes': ['vowel'],
144          'production': ', '}],
145     'whitespace': {
146         'default': ' ',
147         'consolidate': False,
148         'token_class': 'wb'},
149     'metadata': {
150         'author': 'Author McAuthorson'}
151 }
152 gt = GraphTransliterator.from_dict(settings)
153 gt.transliterate('a')

```

```
'A'
```

This feature can be useful if generating a Graph Transliterator using code as opposed to a configuration file.

Ambiguity Checking

Graph Transliterator, by default, will check for ambiguity in its transliteration rules. If two rules of the same cost would match the same string(s) and those strings would not be matched by a less costly rule, an `AmbiguousTransliterationRulesException` occurs. Details of all exceptions will be reported as a `logging.warning()`:

```

155 yaml_ = """
156 tokens:
157     a: [class1, class2]
158     b: []
159     ' ': [wb]
160 rules:
161     <class1> a: A
162     <class2> a: AA # ambiguous rule
163     <class1> b: BB
164     b <class2>: BB # also ambiguous
165 whitespace:
166     default: ' '
167     consolidate: True
168     token_class: wb
169 """
170 gt = GraphTransliterator.from_yaml(yaml_)

```

```

WARNING:root:The pattern [{'a'}, {'a'}, {' ', 'b', 'a'}] can be matched by_
↳both:
    <class1> a

```

```
<class2> a
```

```
WARNING:root:The pattern [{'a'}, {'b'}, {'a'}] can be matched by both:
<class1> b
b <class2>
```

AmbiguousTransliterationRulesException

The warning shows the set of possible previous tokens, matched tokens, and next tokens as three sets.

Ambiguity checking is only necessary when using an untested Graph Transliterator. It can be turned off during initialization. To do so, set the initialization parameter `check_ambiguity` to *False*.

Ambiguity checking can also be done on demand using `check_for_ambiguity()`.

Ambiguity checking is not performed if loading from a serialized GraphTransliterator using `GraphTransliterator.load()` or `GraphTransliterator.loads()`.

Setup Validation

Graph Transliterator validates both the “easy reading” configuration and the direct configuration using the `marshmallow` library.

Transliteration and Its Exceptions

The main method of Graph Transliterator is `GraphTransliterator.transliterate()`. It will return a string:

```
171 GraphTransliterator.from_yaml(  
172     '''  
173     tokens:  
174         a: []  
175         ' ': [wb]  
176     rules:  
177         a: A  
178         ' ': '_'  
179     whitespace:  
180         default: ' '  
181         consolidate: True  
182         token_class: wb  
183     ''').transliterate("a a")
```

'A_A'

Details of transliteration error exceptions will be logged using `logging.warning()`.

Unrecognizable Input Token

Unless the `GraphTransliterator` is initialized with or has the property `ignore_errors` set as `True`, `GraphTransliterator.transliterate()` will raise `UnrecognizableInputTokenException` when character(s) in the input string do not correspond to any defined types of input tokens. In both cases, there will be a `logging.warning()`:

```

184 from graphtransliterator import GraphTransliterator
185 yaml_ = """
186     tokens:
187         a: []
188         ' ': [wb]
189     rules:
190         a: A
191         ' ': ' '
192     whitespace:
193         default: " "
194         consolidate: true
195         token_class: wb
196 """
197 GraphTransliterator.from_yaml(yaml_).transliterate("a!a") # ignore_errors=False

```

WARNING:graphtransliterator:Unrecognizable token ! at pos 1 of a!a

UnrecognizableInputTokenException

```

198 GraphTransliterator.from_yaml(yaml_, ignore_errors=True).transliterate("a!a") #_
    ↪ ignore_errors=True

```

WARNING:graphtransliterator:Unrecognizable token ! at pos 1 of a!a

'AA'

No Matching Transliteration Rule

Another possible error occurs when no transliteration rule can be identified at a particular index in the index string. In that case, there will be a `logging.warning()`. If the parameter `ignore_errors` is set to `True`, the token index will be advanced. Otherwise, there will be a `NoMatchingTransliterationRuleException`:

```

199 yaml_='''
200     tokens:
201         a: []
202         b: []
203         ' ': [wb]
204     rules:
205         a: A
206         b (a): B
207     whitespace:
208         default: ' '
209         token_class: wb
210         consolidate: False
211 '''

```

(continues on next page)

(continued from previous page)

```

212 gt = GraphTransliterator.from_yaml(yaml_)
213 gt.transliterate("ab")

```

```

WARNING:graphtransliterator:No matching transliteration rule at token pos 2_
↳of [' ', 'a', 'b', ' ']

```

```
NoMatchingTransliterationRuleException
```

```

214 gt.ignore_errors = True
215 gt.transliterate("ab")

```

```

WARNING:graphtransliterator:No matching transliteration rule at token pos 2_
↳of [' ', 'a', 'b', ' ']

```

```
'A'
```

Additional Methods

Graph Transliterator also offers a few additional methods that may be useful for particular tasks.

Serialization and Deserialization

The settings of a Graph Transliterator can be serialized using `GraphTransliterator.dump()`, which returns a dictionary of native Python data types. A JSON string of the same can be accessed using `GraphTransliterator.dumps()`. Validation is not performed during a dump.

By default, `GraphTransliterator.dumps()` will use compression level 2, which removes the internal graph and indexes tokens and graph node labels. Compression level 1 also indexes tokens and graph node labels and contains the graph. Compression level 0 is human readable and includes the graph. No information is lost during compression. Level 2, the default, loads the fastest and also has the smallest file size.

A GraphTransliterator can be loaded from serialized settings, e.g. in an API context, using `GraphTransliterator.load()` and from JSON data as `GraphTransliterator.loads()`. Because they are intended to be quick, neither method performs ambiguity checks or strict validation checking by default.

Serialization can be useful if providing an API or making the configured Graph Transliterator available in other programming languages, e.g. Javascript.

Matching at an Index

The method `match_at()` is also public. It matches the best transliteration rule at a particular index, which is the rule that contains the largest number of required tokens. The method also has the option `match_all` which, if set, returns all possible transliteration matches at a particular location:

```

216 gt = GraphTransliterator.from_yaml(''
217     tokens:
218     a: []
219     a a: []
220     ' ': [wb]

```

(continues on next page)

(continued from previous page)

```

221     rules:
222         a: <A>
223         a a: <AA>
224     whitespace:
225         default: ' '
226         consolidate: True
227         token_class: wb
228 '''
229 tokens = gt.tokenize("aa")
230 tokens # whitespace added to ends

```

```
[' ', 'a', 'a', ' ']
```

```
231 gt.match_at(1, tokens) # returns index to rule
```

```
0
```

```
232 gt.rules[gt.match_at(1, tokens)] # actual rule
```

```
TransliterationRule(production='<AA>', prev_classes=None, prev_tokens=None, tokens=['a
→', 'a'], next_tokens=None, next_classes=None, cost=0.41503749927884376)
```

```
233 gt.match_at(1, tokens, match_all=True) # index to rules, with match_all
```

```
[0, 1]
```

```
234 [gt.rules[_] for _ in gt.match_at(1, tokens, match_all=True)] # actual rules, with_
→match_all
```

```
[TransliterationRule(production='<AA>', prev_classes=None, prev_tokens=None, tokens=[
→'a', 'a'], next_tokens=None, next_classes=None, cost=0.41503749927884376),
TransliterationRule(production='<A>', prev_classes=None, prev_tokens=None, tokens=['a
→'], next_tokens=None, next_classes=None, cost=0.5849625007211562)]
```

Details of Matches

Each Graph Transliterator has a property `last_matched_rules` which returns a list of `TransliterationRule` of the previously matched transliteration rules:

```
235 gt.transliterate("aaa")
```

```
'<AA><A>'
```

```
236 gt.last_matched_rules
```

```
[TransliterationRule(production='<AA>', prev_classes=None, prev_tokens=None, tokens=[
→'a', 'a'], next_tokens=None, next_classes=None, cost=0.41503749927884376),
TransliterationRule(production='<A>', prev_classes=None, prev_tokens=None, tokens=['a
→'], next_tokens=None, next_classes=None, cost=0.5849625007211562)]
```

The particular tokens matched by those rules can be accessed using `last_matched_rule_tokens`:

```
237 gt.last_matched_rule_tokens
```

```
[[ 'a', 'a'], [ 'a']]
```

Pruning of Rules

In particular cases, it may be useful to remove certain transliteration rules from a more robustly defined Graph Transliterator based on the string output produced by the rules. That can be done using `pruned_of()`:

```
238 gt.rules
```

```
[TransliterationRule(production='<AA>', prev_classes=None, prev_tokens=None, tokens=[
→ 'a', 'a'], next_tokens=None, next_classes=None, cost=0.41503749927884376),
  TransliterationRule(production='<A>', prev_classes=None, prev_tokens=None, tokens=['a
→'], next_tokens=None, next_classes=None, cost=0.5849625007211562)]
```

```
239 gt.pruned_of('<AA>').rules
```

```
[TransliterationRule(production='<A>', prev_classes=None, prev_tokens=None, tokens=['a
→'], next_tokens=None, next_classes=None, cost=0.5849625007211562)]
```

```
240 gt.pruned_of(['<A>', '<AA>']).rules
```

```
[]
```

Internal Graph

Graph Transliterator creates a directed tree during its initialization. During calls to `transliterate()`, it searches that graph to find the best transliteration match at a particular index in the tokens of the input string.

DirectedGraph

The tree is an instance of `DirectedGraph` that can be accessed using `GraphTransliterator.graph`. It contains: a list of nodes, each consisting of a dictionary of attributes; a dictionary of edges keyed between the head and tail of an edge that contains a dictionary of edge attributes; and finally an edge list.

```
241 gt = GraphTransliterator.from_yaml(
242     """
243     tokens:
244       a: []
245       ' ': [wb]
246     rules:
247       a: b
248       <wb> a: B
249       ' ': ' '
250     whitespace:
251       token_class: wb
252       default: ' '
253       consolidate: false
```

(continues on next page)

(continued from previous page)

```
254     """)
255 gt.graph
```

```
<graphtransliterator.graphs.DirectedGraph at 0x7fd57c312b00>
```

Nodes

The tree has nodes of three types: *Start*, *token*, and *rule*. A single *Start* node, the root, is connected to all other nodes. A *token* node corresponds to a token having been matched. Finally, *rule* nodes are leaf nodes (with no outgoing edges) that correspond to matched transliteration rules:

```
256 gt.graph.node
```

```
[{'type': 'Start', 'ordered_children': {'a': [1], ' ': [4]}},
 {'type': 'token', 'token': 'a', 'ordered_children': {'__rules__': [2, 3]}},
 {'type': 'rule', 'rule_key': 0, 'accepting': True, 'ordered_children': {}},
 {'type': 'rule', 'rule_key': 1, 'accepting': True, 'ordered_children': {}},
 {'type': 'token', 'token': ' ', 'ordered_children': {'__rules__': [5]}},
 {'type': 'rule', 'rule_key': 2, 'accepting': True, 'ordered_children': {}}]
```

Edges

Edges between these nodes may have different constraints in their attributes:

```
257 gt.graph.edge
```

```
{0: {1: {'token': 'a', 'cost': 0.41503749927884376},
      4: {'token': ' ', 'cost': 0.5849625007211562}},
 1: {2: {'cost': 0.41503749927884376, 'constraints': {'prev_classes': ['wb']}},
      3: {'cost': 0.5849625007211562}},
 4: {5: {'cost': 0.5849625007211562}}}
```

Before the *token* nodes, there is a *token* constraint on the edge that must be matched before the transliterator can visit the token node:

```
258 gt.graph.edge[0][1]
```

```
{'token': 'a', 'cost': 0.41503749927884376}
```

On the edges before rules there may be other *constraints*, such as certain tokens preceding or following tokens of the corresponding transliteration rule:

```
259 gt.graph.edge[1][2]
```

```
{'cost': 0.41503749927884376, 'constraints': {'prev_classes': ['wb']}}
```

An edge list is also maintained that consists of a tuple of (head, tail):

```
260 gt.graph.edge_list
```

```
[(0, 1), (1, 2), (1, 3), (0, 4), (4, 5)]
```

Search and Preprocessing

Graph Transliterator uses a best-first search, implemented using a stack, that finds the transliteration with the the lowest cost. The cost function is:

$$\text{cost}(\text{rule}) = \log_2 \left(1 + \frac{1}{1 + \text{count_of_tokens_in}(\text{rule})} \right)$$

It results in a number between 1 and 0 that lessens as more tokens must be matched. Each edge on the graph has a cost attribute that is set to the lowest cost transliteration rule following it. When transliterating, Graph Transliterator will try lower cost edges first and will backtrack if the constraint conditions are not met.

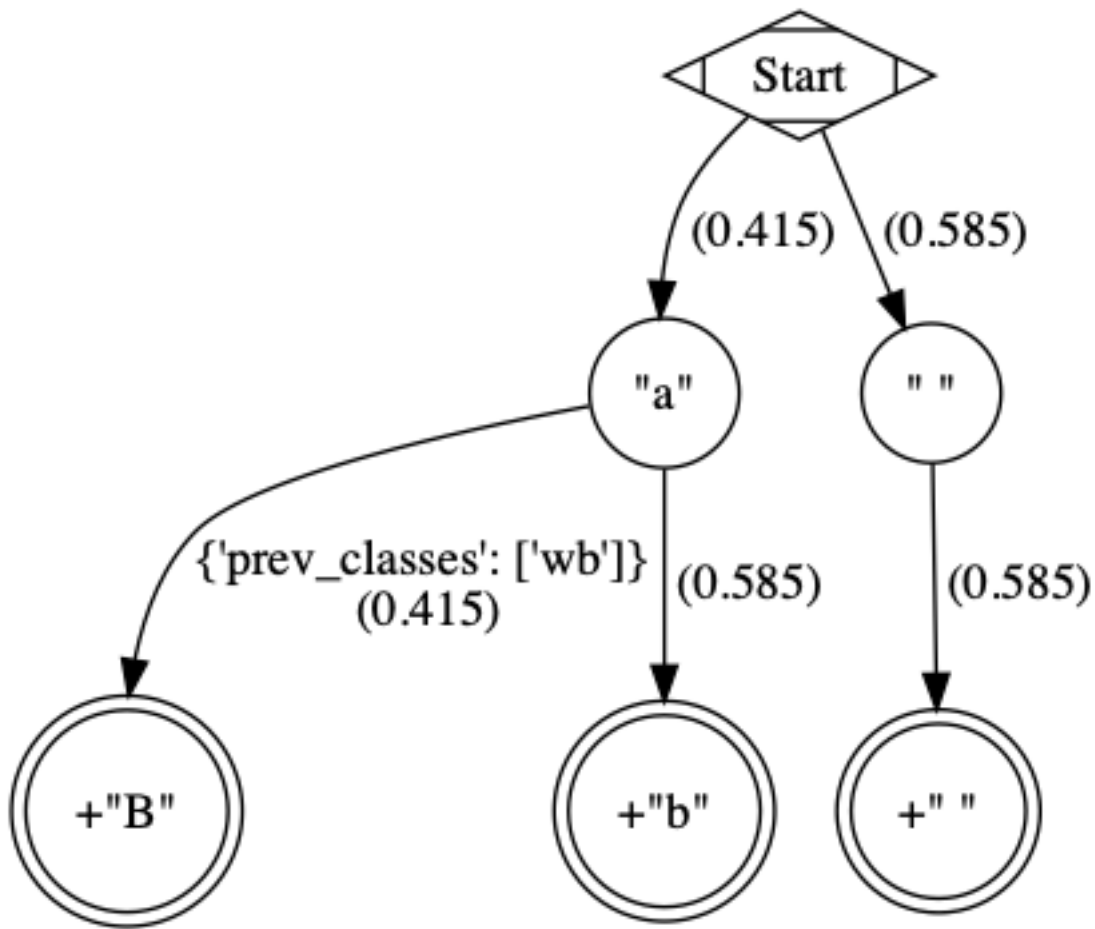


Fig. 2: An example graph created for the simple case of a Graph Transliterator that takes as input two token types, `a` and `" "` (space), and renders `" "` as `" "`, and `a` as `b` unless it follows a token of class `wb` (for wordbreak), in which case it renders `a` as `B`. The *rule* nodes are in double circles, and *token* nodes are single circles. The numbers are the cost of the particular edge, and less costly edges are searched first. Previous token class (`prev_classes`) constraints are found on the edge before the leftmost leaf rule node.

To optimize the search, during initialization an `ordered_children` dictionary is added to each non-leaf node. Its values are a list of node indexes sorted by cost and keyed by the following *token*:

```
261 gt.graph.node[0]
```

```
{'type': 'Start', 'ordered_children': {'a': [1], ' ': [4]}}
```

Any *rule* connected to a node is added to each *ordered_children*. Any rule nodes immediately following the current node are keyed to `__rules__`:

```
262 gt.graph.node[1]
```

```
{'type': 'token', 'token': 'a', 'ordered_children': {'__rules__': [2, 3]}}
```

Because of this preprocessing, Graph Transliterator does not need to iterate through all of the outgoing edges of a node to find the next node to search.

3.2.3 Bundled Transliterators

Note: Python code on this page: `bundled.py` Jupyter Notebook: `bundled.ipynb`

Graph Transliterator includes bundled transliterators in a `Bundled` subclass of `GraphTransliterator` that can be used as follows:

```
1 import graphtransliterator.transliterators as transliterators
2 example_transliterator = transliterators.Example()
3 example_transliterator.transliterate('a')
```

```
'A'
```

To access transliterator classes, use the iterator `transliterators.iter_transliterators()`:

```
4 bundled_iterator = transliterators.iter_transliterators()
5 next(bundled_iterator)
```

```
<example.Example at 0x7f34c0678210>
```

To access the names of transliterator classes, use the iterator `transliterators.iter_names()`:

```
6 bundled_names_iterator = transliterators.iter_names()
7 next(bundled_names_iterator)
```

```
'Example'
```

The actual bundled transliterators are submodules of `graphtransliterator.transliterators`, but they are loaded into the namespace of `transliterators`:

```
8 from graphtransliterator.transliterators import Example
```

Each instance of `Bundled` contains a `directory` attribute:

```
9 transliterator = Example()
10 transliterator.directory
```

```
'/home/docs/checkouts/readthedocs.org/user_builds/graphtransliterators/checkouts/
↳latest/graphtransliterators/transliterators/example'
```

Each will contain an easy-reading YAML file that you can view:

```
tokens:
  a: [vowel]
  ' ': [whitespace]
  b: [consonant]
rules:
  a: A
  b: B
  ' ': ' '
  (<consonant> a) b (a <consonant>): "!B!"
onmatch_rules:
  - <vowel> + <vowel>: ", "
whitespace:
  consolidate: False
  default: " "
  token_class: whitespace
metadata:
  name: example
  version: 1.0.0
  description: "An Example Bundled Transliterators"
  url: https://github.com/seanpue/graphtransliterators/tree/master/transliterators/
↳sample
  author: Author McAuthorson
  author_email: author_mcauthorson@msu.edu
  license: MIT License
  keywords:
    - example
  project_urls:
    Documentation: https://github.com/seanpue/graphtransliterators/tree/master/
↳graphtransliterators/transliterators/example
    Source: https://github.com/seanpue/graphtransliterators/tree/graphtransliterators/
↳transliterators/example
    Tracker: https://github.com/seanpue/graphtransliterators/issues
```

There is also a JSON dump of the transliterators for quick loading:

```
{ "graphtransliterators_version": "1.2.0", "compressed_settings": [ [ "consonant", "vowel",
↳ "whitespace" ], [ " ", "a", "b" ], [ [ 2 ], [ 1 ], [ 0 ] ], [ [ "!B!", [ 0 ], [ 1 ], [ 2 ], [ 1 ], [ 0 ], -5 ], [ "A", 0, 0,
↳ [ 1 ], 0, 0, -1 ], [ "B", 0, 0, [ 2 ], 0, 0, -1 ], [ " ", 0, 0, [ 0 ], 0, 0, -1 ], [ " ", "whitespace", 0 ], [ [ 1 ],
↳ [ 1 ], " ", " ] ], { "name": "example", "version": "1.0.0", "description": "An Example Bundled_
↳ Transliterators", "url": "https://github.com/seanpue/graphtransliterators/tree/master/
↳ transliterators/sample", "author": "Author McAuthorson", "author_email": "author_
↳ mcauthorson@msu.edu", "license": "MIT License", "keywords": [ "example" ], "project_urls": {
↳ "Documentation": "https://github.com/seanpue/graphtransliterators/tree/master/
↳ graphtransliterators/transliterators/example", "Source": "https://github.com/seanpue/
↳ graphtransliterators/tree/graphtransliterators/transliterators/example", "Tracker":
↳ "https://github.com/seanpue/graphtransliterators/issues" } }, null ] }
```

Test Coverage of Bundled Transliterators

Each bundled transliterator requires rigorous testing: every node and edge, as well as any onmatch rules, if applicable, must be visited. A separate subclass `CoverageTransliterator` of `GraphTransliterator` is used during testing.

It logs visits to nodes, edges, and onmatch rules. The tests are found in a subdirectory of the transliterator named “tests”. They are in a YAML file consisting of a dictionary keyed from transliteration input to correct output, e.g.:

```
# YAML declaration of tests for bundled Graph Transliterator
# These are in the form of a dictionary.
# The key is the source text, and the value is the correct transliteration.
' ': ' '
a: A
aa: A,A
babab: BA!B!AB
b: B
```

Once the tests are completed, Graph Transliterator checks that all components of the graph and all of the onmatch rules have been visited.

Class Structure and Naming Conventions

Each transliterator must include a class definition in a submodule of `transliterators`.

The class name of each transliterator must be unique and follow camel-case conventions, e.g. *SourceToTarget*. File and directory names should, if applicable, be lowercased as *source_to_target*.

The bundled files should follow this directory structure, where `{{source_to_target}}` is the name of the transliterator:

```
transliterators
├── {{source_to_target}}
│   ├── __init__.py
│   ├── {{source_to_target}}.json
│   └── {{source_to_target}}.yaml
└── tests
    ├── test_{{source_to_target}}.py
    └── {{source_to_target}}_tests.yaml
```

The bundled transliterator will:

- include both an easy-reading YAML file `{{source_to_target}}.yaml` and a JSON file `{{source_to_target}}.json`.
- have tests in a YAML format consisting of a dictionary keyed from transliteration to correct output in `{{source_to_target}}_tests.yaml`. It must include complete test coverage of its graph. Every node and edge of the graph must be visited during the course of the tests, as well as every on-match rule. Each on-match rule must be utilized during the course of the tests.
- include metadata about the transliterator in its easy-reading YAML file.
- have an optional custom test file `test_{{source_to_target}}.py`. This is useful during development.

Metadata Requirements

Each `Bundled` transliterator can include the following metadata fields. These fields are a subset of the metadata of `setuptools`.

name (*str*)

Name of the transliterator, e.g. “source_to_target”.

version (*str*, optional)

Version of the transliterator. Semantic versioning (<https://semver.org>) is recommended.

url (*str*, optional)

URL for the transliterator, e.g. github repository.

author (*str*, optional)

Author of the transliterator

author_email (*str*, optional)

E-mail address of the author.

maintainer (*str*, optional)

Name of the maintainer.

maintainer_email (*str*, optional)

E-mail address of the maintainer.

license (*str*, optional)

License of the transliterator. An open-source license is required for inclusion in this project.

keywords (*list of str*, optional)

List of keywords.

project_urls (*dict of {str: str}*, optional)

Dictionary of project URLs, e.g. *Documentation*, *Source*, etc.

Metadata is validated using a `BundledMetadataSchema` found in `transliterators.schemas`.

To browse metadata, you can use `iter_transliterators()`:

```
11 import pprint
12 transliterator = next(transliterators.iter_transliterators())
13 pprint.pprint(transliterator.metadata)
```

```
{'author': 'Author McAuthorson',
 'author_email': 'author_mcauthorson@msu.edu',
 'description': 'An Example Bundled Transliterator',
 'keywords': ['example'],
 'license': 'MIT License',
 'name': 'example',
 'project_urls': {'Documentation': 'https://github.com/seanpue/graphtransliterator/
↳tree/master/graphtransliterator/transliterators/example',
                  'Source': 'https://github.com/seanpue/graphtransliterator/tree/
↳graphtransliterator/transliterators/example',
                  'Tracker': 'https://github.com/seanpue/graphtransliterator/issues'},
 'url': 'https://github.com/seanpue/graphtransliterator/tree/master/transliterator/
↳sample',
 'version': '1.0.0'}
```


3.2.4 Command Line Interface

Graph Transliterator has a simple command line interface with six commands: `dump`, `dump-tests`, `generate-tests`, `list-bundled`, `make-json`, `test`, and `transliterate`.

```
$ graphtransliterator --help
```

```
Usage: main [OPTIONS] COMMAND [ARGS]...
```

Options:

```
--version  Show the version and exit.
--help     Show this message and exit.
```

Commands:

```
dump                Dump transliterator as JSON.
dump-tests          Dump BUNDLED tests.
generate-tests      Generate tests as YAML.
list-bundled        List BUNDLED transliterators.
make-json           Make JSON rules of BUNDLED transliterator(s).
test               Test BUNDLED transliterator.
transliterate       Transliterate INPUT.
```

Dump

The `dump` command will output the specified transliterator as JSON:

```
$ graphtransliterator dump --help
```

```
Usage: dump [OPTIONS]
```

```
Dump transliterator as JSON.
```

Options:

```
-f, --from <CHOICE TEXT>...  Format (bundled/yaml_file) and source (name or
                               filename) of transliterator [required]
-ca, --check-ambiguity / -nca, --no-check-ambiguity
                               Check for ambiguity. [default: no-check-
                               ambiguity]
-cl, --compression-level INTEGER
                               Compression level, from 0 to 2 [default: 2]
--help                        Show this message and exit.
```

It require a `--from` or `-f` option with two arguments. The first argument specifies the format of the transliterator (*bundled* or *yaml_file*) and the second a parameter for that format (the name of the bundled transliterator or the name of a YAML file).

To load a bundled transliterator, used *bundled* as the first parameter and give its (class) name, which will be in CamelCase, as the second:

```
$ graphtransliterator dump --from bundled Example
```

```
{"graphtransliterator_version":"1.2.4","compressed_settings":[[{"consonant","vowel",
↪ "whitespace"],[" ","a","b"],[[2],[1],[0]],[["!B!",[0],[1],[2],[1],[0],-5],["A",0,0,
```

(continues on next page)

(continued from previous page)

```

↪ [1],0,0,-1],[ "B",0,0,[2],0,0,-1],[ " ",0,0,[0],0,0,-1],[ " ", "whitespace",0],[[1],
↪ [1], " "],{"name":"example","version":"1.0.0","description":"An Example Bundled
↪ Transliterator","url":"https://github.com/seanpue/graphtransliterator/tree/master/
↪ transliterator/sample","author":"Author McAuthorson","author_email":"author_
↪ mcauthorson@msu.edu","license":"MIT License","keywords":["example"],"project_urls":{"
↪ Documentation":"https://github.com/seanpue/graphtransliterator/tree/master/
↪ graphtransliterator/transl iterators/example","Source":"https://github.com/seanpue/
↪ graphtransliterator/tree/graphtransliterator/transl iterators/example","Tracker":
↪ "https://github.com/seanpue/graphtransliterator/issues"}},null]}

```

To load from a YAML file, give `yaml_file` as the first and the the name of the file as the second parameter:

```

$ graphtransliterator dump --from yaml_file ../graphtransliterator/transl iterators/
↪ example/example.yaml

```

```

{"graphtransliterator_version":"1.2.4","compressed_settings":[{"consonant","vowel",
↪ "whitespace"],[ " ","a","b"],[[2],[1],[0]],[["!B!",[0],[1],[2],[1],[0],-5],[ "A",0,0,
↪ [1],0,0,-1],[ "B",0,0,[2],0,0,-1],[ " ",0,0,[0],0,0,-1],[ " ", "whitespace",0],[[1],
↪ [1], " "],{"name":"example","version":"1.0.0","description":"An Example Bundled
↪ Transliterator","url":"https://github.com/seanpue/graphtransliterator/tree/master/
↪ transliterator/sample","author":"Author McAuthorson","author_email":"author_
↪ mcauthorson@msu.edu","license":"MIT License","keywords":["example"],"project_urls":{"
↪ Documentation":"https://github.com/seanpue/graphtransliterator/tree/master/
↪ graphtransliterator/transl iterators/example","Source":"https://github.com/seanpue/
↪ graphtransliterator/tree/graphtransliterator/transl iterators/example","Tracker":
↪ "https://github.com/seanpue/graphtransliterator/issues"}},null]}

```

If you want to check for ambiguity in the transliterator before the dump, use the `--check-ambiguity` or `-ca` option:

```

$ graphtransliterator dump --from bundled Example --check-ambiguity # human readable

```

```

{"graphtransliterator_version":"1.2.4","compressed_settings":[{"consonant","vowel",
↪ "whitespace"],[ " ","a","b"],[[2],[1],[0]],[["!B!",[0],[1],[2],[1],[0],-5],[ "A",0,0,
↪ [1],0,0,-1],[ "B",0,0,[2],0,0,-1],[ " ",0,0,[0],0,0,-1],[ " ", "whitespace",0],[[1],
↪ [1], " "],{"name":"example","version":"1.0.0","description":"An Example Bundled
↪ Transliterator","url":"https://github.com/seanpue/graphtransliterator/tree/master/
↪ transliterator/sample","author":"Author McAuthorson","author_email":"author_
↪ mcauthorson@msu.edu","license":"MIT License","keywords":["example"],"project_urls":{"
↪ Documentation":"https://github.com/seanpue/graphtransliterator/tree/master/
↪ graphtransliterator/transl iterators/example","Source":"https://github.com/seanpue/
↪ graphtransliterator/tree/graphtransliterator/transl iterators/example","Tracker":
↪ "https://github.com/seanpue/graphtransliterator/issues"}},null]}

```

The compression level can of the JSON be specified using the `--compression-level` or `-cl` command. Compression level 0 is human readable; compression level 1 is not human readable and includes the generated graph; compression level 2 is not human readable and does not include the graph. Compression level 2, which is the fastest, is the default. There is no information lost during these compressions:

```

$ graphtransliterator dump --from bundled Example --compression-level 0 # human
↪ readable, with graph

```

```
{
  "tokens": {
    "a": ["vowel"],
    " ": ["whitespace"],
    "b": ["consonant"]
  },
  "rules": [
    {
      "production": "!B!",
      "prev_classes": ["consonant"],
      "prev_tokens": ["a"],
      "tokens": ["b"],
      "next_classes": ["consonant"],
      "next_tokens": ["a"],
      "cost": 0.22239242133644802,
      {"production": "A", "tokens": ["a"], "cost": 0.5849625007211562},
      {"production": "B", "tokens": ["b"], "cost": 0.5849625007211562},
      {"production": " ", "tokens": [" "], "cost": 0.5849625007211562}],
    "whitespace": {
      "default": " ",
      "token_class": "whitespace",
      "consolidate": false
    },
    "onmatch_rules": [
      {
        "prev_classes": ["vowel"],
        "next_classes": ["vowel"],
        "production": ",",
        "metadata": {
          "name": "example",
          "version": "1.0.0",
          "description": "An Example Bundled Transliterator",
          "url": "https://github.com/seanpue/graphtransliterator/tree/master/transliterator/sample",
          "author": "Author McAuthorson",
          "author_email": "author_mcauthorson@msu.edu",
          "license": "MIT License",
          "keywords": ["example"],
          "project_urls": {
            "Documentation": "https://github.com/seanpue/graphtransliterator/tree/master/graphtransliterator/transliterator/example",
            "Source": "https://github.com/seanpue/graphtransliterator/tree/graphtransliterator/transliterator/example",
            "Tracker": "https://github.com/seanpue/graphtransliterator/issues"
          },
          "ignore_errors": false,
          "onmatch_rules_lookup": {
            "a": {
              "a": [0]
            },
            "tokens_by_class": {
              "vowel": ["a"],
              "whitespace": [" "],
              "consonant": ["b"]
            },
            "graph": {
              "node": [
                {
                  "type": "Start",
                  "ordered_children": {
                    "b": [1],
                    "a": [3],
                    " ": [6]
                  },
                  {"token": "b", "type": "token", "ordered_children": {
                    "__rules__": [2, 5]
                  }, {"type": "rule", "accepting": true, "rule_key": 0}, {"token": "a", "type": "token", "ordered_children": {
                    "__rules__": [4]
                  }, {"type": "rule", "accepting": true, "rule_key": 1}, {"type": "rule", "accepting": true, "rule_key": 2}, {"token": " ", "type": "token", "ordered_children": {
                    "__rules__": [7]
                  }, {"type": "rule", "accepting": true, "rule_key": 3}],
              "edge": {
                "0": {
                  "1": {
                    "token": "b",
                    "cost": 0.22239242133644802
                  },
                  "3": {
                    "token": "a",
                    "cost": 0.5849625007211562
                  },
                  "6": {
                    "token": " ",
                    "cost": 0.5849625007211562
                  },
                  "1": {
                    "2": {
                      "cost": 0.22239242133644802,
                      "constraints": {
                        "prev_classes": ["consonant"],
                        "prev_tokens": ["a"],
                        "next_tokens": ["a"],
                        "next_classes": ["consonant"]
                      },
                      "5": {
                        "cost": 0.5849625007211562
                      },
                      "3": {
                        "4": {
                          "cost": 0.5849625007211562
                        },
                        "6": {
                          "7": {
                            "cost": 0.5849625007211562
                          }
                        }
                      },
                      "edge_list": [[0, 1], [0, 3], [0, 6], [1, 2], [1, 5], [3, 4], [6, 7]]
                    },
                    "tokenizer_pattern": "(b|a|\\s)",
                    "graphtransliterator_version": "1.2.4"
                  }
                }
              }
            }
          }
        }
      }
    ]
  }
}
```

```
$ graphtransliterator dump --from bundled Example --compression-level 1 # not human-readable, with graph
```

```
{
  "graphtransliterator_version": "1.2.4",
  "compressed_settings": [
    ["consonant", "vowel", "whitespace"],
    [" ", "a", "b"],
    [[2], [1], [0]],
    [{"!B!", 0, [1], [2], [1], [0], -5}, {"A", 0, 0, [1], 0, 0, -1}, {"B", 0, 0, [2], 0, 0, -1}, {" ", 0, 0, [0], 0, 0, -1}],
    [" ", "whitespace", 0],
    [[1], [1], ",", ""],
    {"name": "example", "version": "1.0.0", "description": "An Example Bundled Transliterator", "url": "https://github.com/seanpue/graphtransliterator/tree/master/transliterator/sample", "author": "Author McAuthorson", "author_email": "author_mcauthorson@msu.edu", "license": "MIT License", "keywords": ["example"], "project_urls": {
      "Documentation": "https://github.com/seanpue/graphtransliterator/tree/master/graphtransliterator/transliterator/example",
      "Source": "https://github.com/seanpue/graphtransliterator/tree/graphtransliterator/transliterator/example",
      "Tracker": "https://github.com/seanpue/graphtransliterator/issues"
    }},
    [{"Start", "rule", "token"},
    [[0, 0, {"2": [1], "1": [3], "0": [6]}], [2, 0, 2, {"-1": [2, 5]}], [1, 1, 0], [2, 0, 1, {"-1": [4]}], [1, 1, 1], [1, 1, 2], [2, 0, 0, {"-1": [7]}], [1, 1, 3], {"0": {"1": [0, -5, 2], "3": [0, -1, 1], "6": [0, -1, 0]}, "1": {"2": [[0], [1], [1], [0]], -5, -1}, "5": [0, -1, -1], "3": {"4": [0, -1, -1], "6": {"7": [0, -1, -1]}}]]
  ]
}
```

```
$ graphtransliterator dump --from bundled Example --compression-level 2 # default; not human readable, no graph
```

```
{ "graphtransliterators_version": "1.2.4", "compressed_settings": [ [ "consonant", "vowel",  
→ "whitespace" ], [ " ", "a", "b" ], [ [ 2 ], [ 1 ], [ 0 ] ], [ [ "!B!", [ 0 ], [ 1 ], [ 2 ], [ 1 ], [ 0 ], -5 ], [ "A", 0, 0,  
→ [ 1 ], 0, 0, -1 ], [ "B", 0, 0, [ 2 ], 0, 0, -1 ], [ " ", 0, 0, [ 0 ], 0, 0, -1 ] ], [ " ", "whitespace", 0 ], [ [ 1 ],  
→ [ 1 ], " " ] ], { "name": "example", "version": "1.0.0", "description": "An Example Bundled  
→ Transliterators", "url": "https://github.com/seanpue/graphtransliterators/tree/master/  
→ transliterators/sample", "author": "Author McAuthorson", "author_email": "author_  
→ mcauthorson@msu.edu", "license": "MIT License", "keywords": [ "example" ], "project_urls": {  
→ "Documentation": "https://github.com/seanpue/graphtransliterators/tree/master/  
→ graphtransliterators/transliterators/example", "Source": "https://github.com/seanpue/  
→ graphtransliterators/tree/graphtransliterators/transliterators/example", "Tracker":  
→ "https://github.com/seanpue/graphtransliterators/issues" } }, null ] }
```

Dump Tests

The `dump-tests` command dumps the tests of a bundled transliterators:

```
$ graphtransliterators dump-tests --help
```

```
Usage: dump-tests [OPTIONS] BUNDLED
```

Dump BUNDLED tests.

Options:

```
-t, --to [json|yaml]  Format (json/yaml) in which to dump  [default: yaml]  
--help                Show this message and exit.
```

By default, it outputs the original YAML tests file, preserving any comments:

```
$ graphtransliterators dump-tests Example
```

```
# YAML declaration of tests for bundled Graph Transliterators  
# These are in the form of a dictionary.  
# The key is the source text, and the value is the correct transliteration.  
' ': ' '  
a: A  
aa: A,A  
babab: BA!B!AB  
b: B
```

To output as JSON, use the `--to` or `-t` flag:

```
$ graphtransliterators dump-tests --to json Example
```

```
{ " ": " ", "a": "A", "aa": "A,A", "babab": "BA!B!AB", "b": "B" }
```

Generate Tests

The `generate-tests` command generates YAML tests keyed from input to desired output covering the entire internal graph. This command can be used to view the output of the transliterator in Unicode. It can also be used to generate starter tests for bundled transliterators:

```
$ graphtransliterator generate-tests --help
```

```
Usage: generate-tests [OPTIONS]
```

```
Generate tests as YAML.
```

```
Options:
```

```
-f, --from <CHOICE TEXT>...    Format (bundled/json/json_file/yaml_file) and
                                source (name, JSON, or filename) of
                                transliterator [required]
-ca, --check-ambiguity / -nca, --no-check-ambiguity
                                Check for ambiguity. [default: no-check-
                                ambiguity]
--help                          Show this message and exit.
```

It also require a `--from` or `-f` option with two arguments. The first argument specifies the format of the transliterator (*bundled*, *json*, *json_file*, *yaml_file*), and the second a parameter for that format (the name of the bundled transliterator, the actual JSON, or the name of a YAML file). Ambiguity checking can be turned on using `--check_ambiguity` or `-ca`:

```
$ graphtransliterator generate-tests --from bundled Example
```

```
' ': ' '
a: A
aa: A,A
b: B
babab: BA!B!AB
```

List Bundled Transliterators

The `list-bundled` command provides a list of bundled transliterators:

```
$ graphtransliterator test --help
```

Make JSON of Bundled Transliterator(s)

The `make-json` command makes new JSON files of bundled transliterators:

```
$ graphtransliterator make-json --help
```

It also allows regular-expression matching using the `--reg-ex` or `-re` flag. Matching starts at the start of the string. This command is for people creating new bundled transliterators.

Test

The `test` command tests a bundled transliterator:

```
$ graphtransliterator test --help
```

```
Usage: test [OPTIONS] BUNDLED

Test BUNDLED transliterator.

Options:
  -ca, --check-ambiguity / -nca, --no-check-ambiguity
                                Check for ambiguity. [default: no-check-
                                ambiguity]
  --help                        Show this message and exit.
```

It can only be used with bundled transliterators, so it only needs the name of the transliterator as its argument. This feature is useful when developing a transliterator. You can write the tests first and then begin developing the transliterator:

```
$ graphtransliterator test Example
```

```
True
```

Transliterate

The `transliterate` command will transliterate any following arguments:

```
$ graphtransliterator transliterate --help
```

```
Usage: transliterate [OPTIONS] [INPUT]...

Transliterate INPUT.

Options:
  -f, --from <CHOICE TEXT>...  Format (bundled/json/json_file/yaml_file) and
                                source (name, JSON, or filename) of
                                transliterator [required]
  -t, --to [json|python]        Format in which to output [default: python]
  -ca, --check-ambiguity / -nca, --no-check-ambiguity
                                Check for ambiguity. [default: no-check-
                                ambiguity]
  -ie, -nie, --ignore-errors / --no-ignore-errors
                                Ignore errors. [default: no-ignore-errors]
  --help                        Show this message and exit.
```

It also requires a `--from` or `-f` option with two arguments. The first argument specifies the format of the transliterator (*bundled*, *json*, *json_file*, *yaml_file*), and the second a parameter for that format (the name of the bundled transliterator, the actual JSON, or the name of a YAML file).

The *transliterate* command will transliterate every argument that follows. If there is only one input string, it will return a string:

```
$ graphtransliterator transliterate --from bundled Example a
```

```
A
```

```
$ graphtransliterator transliterate -f json_file ../graphtransliterator/  
↳transliterators/example/example.json a
```

```
A
```

```
$ graphtransliterator transliterate -f yaml_file ../graphtransliterator/  
↳transliterators/example/example.yaml a
```

```
A
```

Otherwise, it will return a list:

```
$ graphtransliterator transliterate -f bundled Example a a
```

```
['A', 'A']
```

The *transliterate* command also an optional `--to` or `-t` command that specifies the output format, a ``python` string (default) or a `json` string:

```
$ graphtransliterator transliterate --from bundled Example a
```

```
A
```

```
$ graphtransliterator transliterate --from bundled Example --to json a
```

```
"A"
```

```
$ graphtransliterator transliterate --from bundled Example --to python a a
```

```
['A', 'A']
```

```
$ graphtransliterator transliterate --from bundled Example --to json a a
```

```
["A", "A"]
```

3.2.5 Tutorial: Using GraphTransliterator

Note: Python code on this page: `tutorial.py` Jupyter Notebook: `tutorial.ipynb`

Graph Transliterator is designed to allow you to quickly develop rules for transliterating between languages and scripts. In this tutorial you will use a portion of Graph Transliterator's features, including its token matching, class-based matching, and on match rules, using the `GraphTransliterator` class.

Tutorial Overview

The task for this tutorial will be to design a transliterator between the [ITRANS \(Indian languages TRANSliteration\)](#) encoding for [Devanagari](#) (Hindi) and standard [Unicode](#). ITRANS developed as a means to transliterate Indic-language using the latin alphabet and punctuation marks before there were Unicode fonts.

The Devanagari alphabet is an abugida (alphasyllabary), where each “syllable” is a separate symbol. Vowels, except for the default अ (“a”) have a unique symbol that connects to a consonant. At the start of the words, they have a unique shape. Consonants in sequence, without intermediary vowels, change their shape and are joined together. In Unicode, that is accomplished by using the [Virama](#) character.

Graph Transliterator works by first converting the input text into a series of tokens. In this tutorial you will define the tokens of ITRANS and necessary token classes that will allow us to generate rules for conversion.

Graph Transliterator allows rule matching by preceding tokens, tokens, and following tokens. It allows token classes to precede or follow any specific tokens. For this task, you will use a preceding token class to identify when to write vowel signs as opposed to full vowel characters.

Graph Transliterator also allows the insertion of strings between matches involving particular token classes. This transliterator will need to insert the virama character between transliteration rules ending with consonants in order to create consonant clusters.

Configuring

Here you will parameterize the Graph Transliterator using its “easy reading” format, which uses [YAML](#). It maps to a dictionary containing up to five keys: `tokens`, `rules`, `onmatch_rules` (optional), `whitespace`, and `metadata` (optional).

Token Definitions

Graph Transliterator tokenizes its input before transliterating. The `tokens` section will map the input tokens to their token classes. The main class you will need is one for consonants, so you can use `consonant` as the class. Graph Transliterator also requires a dedicated whitespace class, so you can use `whitespace`.

Graph Transliterator allows the use of Unicode character names in files using `\N{UNICODE CHARACTER NAME HERE}}` notation. You can enter the Unicode characters using that notation or directly. YAML will also unescape `\u####`, where `####` is the hexadecimal notation for a character.

Here is a subsection of that definition:

```
tokens:
  k: [consonant]
  kh: [consonant]
  "\N{LATIN SMALL LETTER N WITH DOT ABOVE}": [consonant]
  a: [vowel]
```

(continues on next page)

(continued from previous page)

```

aa: [vowel]
A: [vowel]
' ': [wb, whitespace]
"\t": [wb, whitespace]
.N: [vowel_sign]

```

Transliteration Rule Definitions

The rule definitions in Graph Transliterator in “easy reading” format are also a dictionary where the rules are the key and the production—what should be outputted by the rule—is the value. For this task, you just need to match individual tokens and also any preceding token classes:

```

rules:
  b: \N{DEVANAGARI LETTER B}
  <consonant> A: \N{DEVANAGARI LETTER AA}
  A: \N{DEVANAGARI LETTER AA}

```

These rules will replace “b” with the devanagari equivalent (ब), and “A” with with a full letter आ if it is at a start of a word (following a token of class “wb”, for wordbreak) or otherwise with a vowel sign ा if it is not, presumably following a consonant. Graph Transliterator automatically sorts rules by how many tokens are required for them to be matched, and it picks the one with that requires the most tokens. So the “A” following a consonant would be matched before an “A” after any other character. Graph Transliterator will also check for ambiguity in these rules, unless `check_ambiguity` is set to `False`.

While not necessary for this tutorial, Graph Transliterator can also require matching of specific previous or following tokens and also classes preceding and following those tokens, e.g.

```

k a r (U M g A <wb>): k,a,r_followed_by_U,M,g,A_and_a_wordbreak
s o (n a): s,o_followed_by_n,a
(<wb> p y) aa r: aa,r_preceded_by_a_wordbreak,p,and_y

```

Here is a subsection of the rules:

```

rules:
  "\t": "\t"
  ' ': ' '
  ',': ','
  .D: "\N{DEVANAGARI LETTER DDDHA}"
  <consonant> A: "\N{DEVANAGARI VOWEL SIGN AA}"
  "\N{LATIN SMALL LETTER N WITH DOT ABOVE}": "\N{DEVANAGARI LETTER NGA}"

```

On Match Rule Definitions

You will want to insert the Virama character between consonants so that they will join together in Unicode output. To do so, add an “onmatch_rules” section:

```

onmatch_rules:
  - <consonant> + <consonant>: "\N{DEVANAGARI SIGN VIRAMA}"

```

Unlike the tokens and rules, the *onmatch rules are ordered*. The first rule matched is applied. In YAML, they consist of a list of dictionaries each with a single key and value. The value is the production string to be inserted between matches. The ``+`` represents that space. So in the input string `kyA`, which would tokenize as `[' ', 'k', 'y', 'A', ' ']`, a virama

character would be inserted when `y` is matched, as it is of class “consonant” and the previously matched transliteration rule for “`k`” ends with a “consonant”.

Whitespace Definitions

The final required setup parameter is for whitespace. These include the `default` whitespace token, which is temporarily added before and after the input tokens; the `consolidate` option to replace sequential whitespace characters with a single default whitespace character; and the `token_class` of whitespace tokens:

```
whitespace:
  consolidate: false
  default: ' '
  token_class: whitespace
```

Metadata Definitions

Graph Transliterator also allows metadata to be added to its settings:

```
metadata:
  title: "ITRANS Devanagari to Unicode"
  version: "0.1.0"
```

Creating a Transliterator

Now that the settings are ready, you can create a Graph Transliterator. Since you have been using the “easy reading” format, you can use `GraphTransliterator.from_yaml_file()` to read from a specific file or the `GraphTransliterator.from_yaml()` to read from a YAML string. You read from the loaded contents of an “easy reading” YAML file using `GraphTransliterator.from_dict()`. Graph Transliterator will convert those settings into basic Python types and then return a `GraphTransliterator`:

```
1 from graphtransliterator import GraphTransliterator
2 easyreading_yaml = """
3 tokens:
4   k: [consonant]
5   kh: [consonant]
6   g: [consonant]
7   gh: [consonant]
8   ~N: [consonant]
9   "\N{LATIN SMALL LETTER N WITH DOT ABOVE}": [consonant]
10  ch: [consonant]
11  chh: [consonant]
12  Ch: [consonant]
13  j: [consonant]
14  jh: [consonant]
15  ~n: [consonant]
16  T: [consonant]
17  Th: [consonant]
18  D: [consonant]
19  Dh: [consonant]
20  N: [consonant]
21  t: [consonant]
22  th: [consonant]
```

(continues on next page)

(continued from previous page)

```

23 d: [consonant]
24 dh: [consonant]
25 n: [consonant]
26 ^n: [consonant]
27 p: [consonant]
28 ph: [consonant]
29 b: [consonant]
30 bh: [consonant]
31 m: [consonant]
32 y: [consonant]
33 r: [consonant]
34 R: [consonant]
35 l: [consonant]
36 ld: [consonant]
37 L: [consonant]
38 zh: [consonant]
39 v: [consonant]
40 sh: [consonant]
41 Sh: [consonant]
42 s: [consonant]
43 h: [consonant]
44 x: [consonant]
45 kSh: [consonant]
46 GY: [consonant]
47 j~n: [consonant]
48 dny: [consonant]
49 q: [consonant]
50 K: [consonant]
51 G: [consonant]
52 J: [consonant]
53 z: [consonant]
54 .D: [consonant]
55 .Dh: [consonant]
56 f: [consonant]
57 Y: [consonant]
58 a: [vowel]
59 aa: [vowel]
60 A: [vowel]
61 i: [vowel]
62 ii: [vowel]
63 I: [vowel]
64 ee: [vowel]
65 u: [vowel]
66 uu: [vowel]
67 U: [vowel]
68 RRi: [vowel]
69 R^i: [vowel]
70 LLi: [vowel]
71 L^i: [vowel]
72 RRI: [vowel]
73 LLI: [vowel]
74 a.c: [vowel]
75 ^e: [vowel]
76 e: [vowel]
77 ai: [vowel]
78 A.c: [vowel]
79 ^o: [vowel]

```

(continues on next page)

(continued from previous page)

```

80  o: [vowel]
81  au: [vowel]
82  ' ': [wb, whitespace]
83  "\t": [wb, whitespace]
84  ',': [wb]
85  .h: [wb]
86  H: [wb]
87  OM: [wb]
88  AUM: [wb]
89  '|': [wb]
90  '||': [wb]
91  '0': [wb]
92  '1': [wb]
93  '2': [wb]
94  '3': [wb]
95  '4': [wb]
96  '5': [wb]
97  '6': [wb]
98  '7': [wb]
99  '8': [wb]
100 '9': [wb]
101 Rs.: [wb]
102 ~Rs.: [wb]
103 .a: [wb]
104 a.e: [vowel_sign]
105 .N: [vowel_sign]
106 .n: [vowel_sign]
107 M: [vowel_sign]
108 .m: [vowel_sign]
109 rules:
110   "\t": "\t"
111   ' ': ' '
112   ',': ','
113   .D: "\N{DEVANAGARI LETTER DDDHA}"
114   .Dh: "\N{DEVANAGARI LETTER RHA}"
115   .N: "\N{DEVANAGARI SIGN CANDRABINDU}"
116   .a: "\N{DEVANAGARI SIGN AVAGRAHA}"
117   .h: "\N{DEVANAGARI SIGN VIRAMA}\N{ZERO WIDTH NON-JOINER}"
118   .m: "\N{DEVANAGARI SIGN ANUSVARA}"
119   .n: "\N{DEVANAGARI SIGN ANUSVARA}"
120   '0': "\N{DEVANAGARI DIGIT ZERO}"
121   '1': "\N{DEVANAGARI DIGIT ONE}"
122   '2': "\N{DEVANAGARI DIGIT TWO}"
123   '3': "\N{DEVANAGARI DIGIT THREE}"
124   '4': "\N{DEVANAGARI DIGIT FOUR}"
125   '5': "\N{DEVANAGARI DIGIT FIVE}"
126   '6': "\N{DEVANAGARI DIGIT SIX}"
127   '7': "\N{DEVANAGARI DIGIT SEVEN}"
128   '8': "\N{DEVANAGARI DIGIT EIGHT}"
129   '9': "\N{DEVANAGARI DIGIT NINE}"
130   <consonant> A: "\N{DEVANAGARI VOWEL SIGN AA}"
131   <consonant> A.c: "\N{DEVANAGARI VOWEL SIGN CANDRA O}"
132   <consonant> I: "\N{DEVANAGARI VOWEL SIGN II}"
133   <consonant> LLi: "\N{DEVANAGARI VOWEL SIGN VOCALIC LL}"
134   <consonant> LLi: "\N{DEVANAGARI VOWEL SIGN VOCALIC L}"
135   <consonant> L^i: "\N{DEVANAGARI VOWEL SIGN VOCALIC L}"
136   <consonant> RRI: "\N{DEVANAGARI VOWEL SIGN VOCALIC RR}"

```

(continues on next page)

(continued from previous page)

```

137 <consonant> RRi: "\N{DEVANAGARI VOWEL SIGN VOCALIC R}"
138 <consonant> R^i: "\N{DEVANAGARI VOWEL SIGN VOCALIC R}"
139 <consonant> U: "\N{DEVANAGARI VOWEL SIGN UU}"
140 <consonant> ^e: "\N{DEVANAGARI VOWEL SIGN SHORT E}"
141 <consonant> ^o: "\N{DEVANAGARI VOWEL SIGN SHORT O}"
142 <consonant> a: ' '
143 <consonant> a.c: "\N{DEVANAGARI VOWEL SIGN CANDRA E}"
144 <consonant> aa: "\N{DEVANAGARI VOWEL SIGN AA}"
145 <consonant> ai: "\N{DEVANAGARI VOWEL SIGN AI}"
146 <consonant> au: "\N{DEVANAGARI VOWEL SIGN AU}"
147 <consonant> e: "\N{DEVANAGARI VOWEL SIGN E}"
148 <consonant> ee: "\N{DEVANAGARI VOWEL SIGN II}"
149 <consonant> i: "\N{DEVANAGARI VOWEL SIGN I}"
150 <consonant> ii: "\N{DEVANAGARI VOWEL SIGN II}"
151 <consonant> o: "\N{DEVANAGARI VOWEL SIGN O}"
152 <consonant> u: "\N{DEVANAGARI VOWEL SIGN U}"
153 <consonant> uu: "\N{DEVANAGARI VOWEL SIGN UU}"
154 A: "\N{DEVANAGARI LETTER AA}"
155 A.c: "\N{DEVANAGARI LETTER CANDRA O}"
156 AUM: "\N{DEVANAGARI OM}"
157 Ch: "\N{DEVANAGARI LETTER CHA}"
158 D: "\N{DEVANAGARI LETTER DDA}"
159 Dh: "\N{DEVANAGARI LETTER DDHA}"
160 G: "\N{DEVANAGARI LETTER GHHA}"
161 GY: "\N{DEVANAGARI LETTER JA}\N{DEVANAGARI SIGN VIRAMA}\N{DEVANAGARI LETTER NYA}"
162 H: "\N{DEVANAGARI SIGN VISARGA}"
163 I: "\N{DEVANAGARI LETTER II}"
164 J: "\N{DEVANAGARI LETTER ZA}"
165 K: "\N{DEVANAGARI LETTER KHHA}"
166 L: "\N{DEVANAGARI LETTER LLA}"
167 LLI: "\N{DEVANAGARI LETTER VOCALIC LL}"
168 LLi: "\N{DEVANAGARI LETTER VOCALIC L}"
169 L^i: "\N{DEVANAGARI LETTER VOCALIC L}"
170 M: "\N{DEVANAGARI SIGN ANUSVARA}"
171 N: "\N{DEVANAGARI LETTER NNA}"
172 OM: "\N{DEVANAGARI OM}"
173 R: "\N{DEVANAGARI LETTER RRA}"
174 RRI: "\N{DEVANAGARI LETTER VOCALIC RR}"
175 RRi: "\N{DEVANAGARI LETTER VOCALIC R}"
176 R^i: "\N{DEVANAGARI LETTER VOCALIC R}"
177 Rs.: "\N{INDIAN RUPEE SIGN}"
178 Sh: "\N{DEVANAGARI LETTER SSA}"
179 T: "\N{DEVANAGARI LETTER TTA}"
180 Th: "\N{DEVANAGARI LETTER TTHA}"
181 U: "\N{DEVANAGARI LETTER UU}"
182 Y: "\N{DEVANAGARI LETTER YYA}"
183 ^e: "\N{DEVANAGARI LETTER SHORT E}"
184 ^n: "\N{DEVANAGARI LETTER NNNA}"
185 ^o: "\N{DEVANAGARI LETTER SHORT O}"
186 a: "\N{DEVANAGARI LETTER A}"
187 a.c: "\N{DEVANAGARI LETTER CANDRA E}"
188 a.e: "\N{DEVANAGARI LETTER CANDRA A}"
189 aa: "\N{DEVANAGARI LETTER AA}"
190 ai: "\N{DEVANAGARI LETTER AI}"
191 au: "\N{DEVANAGARI LETTER AU}"
192 b: "\N{DEVANAGARI LETTER BA}"
193 bh: "\N{DEVANAGARI LETTER BHA}"

```

(continues on next page)

(continued from previous page)

```

194 ch: "\N{DEVANAGARI LETTER CA}"
195 chh: "\N{DEVANAGARI LETTER CHA}"
196 d: "\N{DEVANAGARI LETTER DA}"
197 dh: "\N{DEVANAGARI LETTER DHA}"
198 dny: "\N{DEVANAGARI LETTER JA}\N{DEVANAGARI SIGN VIRAMA}\N{DEVANAGARI LETTER NYA}"
199 e: "\N{DEVANAGARI LETTER E}"
200 ee: "\N{DEVANAGARI LETTER II}"
201 f: "\N{DEVANAGARI LETTER FA}"
202 g: "\N{DEVANAGARI LETTER GA}"
203 gh: "\N{DEVANAGARI LETTER GHA}"
204 h: "\N{DEVANAGARI LETTER HA}"
205 i: "\N{DEVANAGARI LETTER I}"
206 ii: "\N{DEVANAGARI LETTER II}"
207 j: "\N{DEVANAGARI LETTER JA}"
208 jh: "\N{DEVANAGARI LETTER JHA}"
209 j~n: "\N{DEVANAGARI LETTER JA}\N{DEVANAGARI SIGN VIRAMA}\N{DEVANAGARI LETTER NYA}"
210 k: "\N{DEVANAGARI LETTER KA}"
211 kSh: "\N{DEVANAGARI LETTER KA}\N{DEVANAGARI SIGN VIRAMA}\N{DEVANAGARI LETTER SSA}"
212 kh: "\N{DEVANAGARI LETTER KHA}"
213 l: "\N{DEVANAGARI LETTER LA}"
214 ld: "\N{DEVANAGARI LETTER LLA}"
215 m: "\N{DEVANAGARI LETTER MA}"
216 n: "\N{DEVANAGARI LETTER NA}"
217 o: "\N{DEVANAGARI LETTER O}"
218 p: "\N{DEVANAGARI LETTER PA}"
219 ph: "\N{DEVANAGARI LETTER PHA}"
220 q: "\N{DEVANAGARI LETTER QA}"
221 r: "\N{DEVANAGARI LETTER RA}"
222 s: "\N{DEVANAGARI LETTER SA}"
223 sh: "\N{DEVANAGARI LETTER SHA}"
224 t: "\N{DEVANAGARI LETTER TA}"
225 th: "\N{DEVANAGARI LETTER THA}"
226 u: "\N{DEVANAGARI LETTER U}"
227 uu: "\N{DEVANAGARI LETTER UU}"
228 v: "\N{DEVANAGARI LETTER VA}"
229 x: "\N{DEVANAGARI LETTER KA}\N{DEVANAGARI SIGN VIRAMA}\N{DEVANAGARI LETTER SSA}"
230 y: "\N{DEVANAGARI LETTER YA}"
231 z: "\N{DEVANAGARI LETTER ZA}"
232 zh: "\N{DEVANAGARI LETTER LLLA}"
233 '|': "\N{DEVANAGARI DANDA}"
234 '||': "\N{DEVANAGARI DOUBLE DANDA}"
235 ~N: "\N{DEVANAGARI LETTER NGA}"
236 ~Rs.: "\N{INDIAN RUPEE SIGN}"
237 ~n: "\N{DEVANAGARI LETTER NYA}"
238 "\N{LATIN SMALL LETTER N WITH DOT ABOVE}": "\N{DEVANAGARI LETTER NGA}"
239 onmatch_rules:
240 - <consonant> + <consonant>: "\N{DEVANAGARI SIGN VIRAMA}"
241 whitespace:
242   consolidate: false
243   default: ' '
244   token_class: whitespace
245 metadata:
246   title: ITRANS to Unicode
247   version: 0.1.0
248   """
249 gt = GraphTransliterator.from_yaml(easyreading_yaml)

```

Transliterating

With the transliterator created, you can now transliterate using `GraphTransliterator.transliterate()`:

```
250 gt.transliterate("aaj mausam ba.Daa beiimaan hai, aaj mausam")
```

```
'ॐ ॐॐॐ ॐॐ ॐॐॐॐॐ ॐ, ॐ ॐॐॐ'
```

Other Information

Graph Transliterator has a few other tools built in that are for more specialized applications.

If you want to receive the details of the most recent transliteration, access `GraphTransliterator.last_matched_rules` to get this list of rules matched:

```
251 gt.last_matched_rules
```

```
[TransliterationRule(production=' ', prev_classes=None, prev_tokens=None, tokens=['aa
→'], next_tokens=None, next_classes=None, cost=0.5849625007211562),
 TransliterationRule(production='j', prev_classes=None, prev_tokens=None, tokens=['j
→'], next_tokens=None, next_classes=None, cost=0.5849625007211562),
 TransliterationRule(production=' ', prev_classes=None, prev_tokens=None, tokens=['
→'], next_tokens=None, next_classes=None, cost=0.5849625007211562),
 TransliterationRule(production='m', prev_classes=None, prev_tokens=None, tokens=['m
→'], next_tokens=None, next_classes=None, cost=0.5849625007211562),
 TransliterationRule(production='au', prev_classes=['consonant'], prev_tokens=None,
→tokens=['au'], next_tokens=None, next_classes=None, cost=0.41503749927884376),
 TransliterationRule(production='s', prev_classes=None, prev_tokens=None, tokens=['s
→'], next_tokens=None, next_classes=None, cost=0.5849625007211562),
 TransliterationRule(production='', prev_classes=['consonant'], prev_tokens=None,
→tokens=['a'], next_tokens=None, next_classes=None, cost=0.41503749927884376),
 TransliterationRule(production='m', prev_classes=None, prev_tokens=None, tokens=['m
→'], next_tokens=None, next_classes=None, cost=0.5849625007211562),
 TransliterationRule(production=' ', prev_classes=None, prev_tokens=None, tokens=['
→'], next_tokens=None, next_classes=None, cost=0.5849625007211562),
 TransliterationRule(production='b', prev_classes=None, prev_tokens=None, tokens=['b
→'], next_tokens=None, next_classes=None, cost=0.5849625007211562),
 TransliterationRule(production='', prev_classes=['consonant'], prev_tokens=None,
→tokens=['a'], next_tokens=None, next_classes=None, cost=0.41503749927884376),
 TransliterationRule(production='D', prev_classes=None, prev_tokens=None, tokens=['.D
→'], next_tokens=None, next_classes=None, cost=0.5849625007211562),
 TransliterationRule(production='aa', prev_classes=['consonant'], prev_tokens=None,
→tokens=['aa'], next_tokens=None, next_classes=None, cost=0.41503749927884376),
 TransliterationRule(production=' ', prev_classes=None, prev_tokens=None, tokens=['
→'], next_tokens=None, next_classes=None, cost=0.5849625007211562),
 TransliterationRule(production='b', prev_classes=None, prev_tokens=None, tokens=['b
→'], next_tokens=None, next_classes=None, cost=0.5849625007211562),
 TransliterationRule(production='e', prev_classes=['consonant'], prev_tokens=None,
→tokens=['e'], next_tokens=None, next_classes=None, cost=0.41503749927884376),
 TransliterationRule(production='ii', prev_classes=None, prev_tokens=None, tokens=['ii
→'], next_tokens=None, next_classes=None, cost=0.5849625007211562),
 TransliterationRule(production='m', prev_classes=None, prev_tokens=None, tokens=['m
→'], next_tokens=None, next_classes=None, cost=0.5849625007211562),
 TransliterationRule(production='aa', prev_classes=['consonant'], prev_tokens=None,
→tokens=['aa'], next_tokens=None, next_classes=None, cost=0.41503749927884376),
 TransliterationRule(production='n', prev_classes=None, prev_tokens=None, tokens=['n
```

(continues on next page)

(continued from previous page)

```

→'], next_tokens=None, next_classes=None, cost=0.5849625007211562),
TransliterationRule(production=' ', prev_classes=None, prev_tokens=None, tokens=['
→'], next_tokens=None, next_classes=None, cost=0.5849625007211562),
TransliterationRule(production='h', prev_classes=None, prev_tokens=None, tokens=['h
→'], next_tokens=None, next_classes=None, cost=0.5849625007211562),
TransliterationRule(production='h', prev_classes=['consonant'], prev_tokens=None,
→tokens=['ai'], next_tokens=None, next_classes=None, cost=0.41503749927884376),
TransliterationRule(production=' ', prev_classes=None, prev_tokens=None, tokens=['
→'], next_tokens=None, next_classes=None, cost=0.5849625007211562),
TransliterationRule(production=' ', prev_classes=None, prev_tokens=None, tokens=['
→'], next_tokens=None, next_classes=None, cost=0.5849625007211562),
TransliterationRule(production='aa', prev_classes=None, prev_tokens=None, tokens=['aa
→'], next_tokens=None, next_classes=None, cost=0.5849625007211562),
TransliterationRule(production='j', prev_classes=None, prev_tokens=None, tokens=['j
→'], next_tokens=None, next_classes=None, cost=0.5849625007211562),
TransliterationRule(production=' ', prev_classes=None, prev_tokens=None, tokens=['
→'], next_tokens=None, next_classes=None, cost=0.5849625007211562),
TransliterationRule(production='m', prev_classes=None, prev_tokens=None, tokens=['m
→'], next_tokens=None, next_classes=None, cost=0.5849625007211562),
TransliterationRule(production='m', prev_classes=['consonant'], prev_tokens=None,
→tokens=['au'], next_tokens=None, next_classes=None, cost=0.41503749927884376),
TransliterationRule(production='s', prev_classes=None, prev_tokens=None, tokens=['s
→'], next_tokens=None, next_classes=None, cost=0.5849625007211562),
TransliterationRule(production=' ', prev_classes=['consonant'], prev_tokens=None,
→tokens=['a'], next_tokens=None, next_classes=None, cost=0.41503749927884376),
TransliterationRule(production='m', prev_classes=None, prev_tokens=None, tokens=['m
→'], next_tokens=None, next_classes=None, cost=0.5849625007211562)]

```

Or if you just want to know the tokens matched by each rule, check `GraphTransliterator.last_matched_rule_tokens`:

```
252 gt.last_matched_rule_tokens
```

```

[['aa'],
 ['j'],
 [' '],
 ['m'],
 ['au'],
 ['s'],
 ['a'],
 ['m'],
 [' '],
 ['b'],
 ['a'],
 ['.D'],
 ['aa'],
 [' '],
 ['b'],
 ['e'],
 ['ii'],
 ['m'],
 ['aa'],
 ['n'],
 [' '],
 ['h'],
 ['ai'],

```

(continues on next page)

(continued from previous page)

```
[','],
[' '],
['aa'],
['j'],
[' '],
['m'],
['au'],
['s'],
['a'],
['m']]
```

You can access the directed tree used by GraphTransliterator using `GraphTransliterator.graph`:

253 `gt.graph`

```
<graphtransliterator.graphs.DirectedGraph at 0x7f8aac66be40>
```

3.2.6 Advanced Tutorial: Bundling a Transliterator

This advanced tutorial builds upon the original tutorial to show you how to bundle a transliterator for inclusion in Graph Transliterator.

Contributions to Graph Transliterator are strongly encouraged!

You will make a very simple transliterator while going through the steps of bundling it into Graph Transliterator.

Git Basics: Fork, Branch, Sync, Commit

Fork

The first thing to do, if you have not already, is to create a fork of Graph Transliterator. See <https://help.github.com/en/articles/fork-a-repo>

(From here on out, we will be using the command line.)

After creating a fork, clone your forked repo:

```
git clone https://github.com/YOUR-USERNAME/graphtransliterator
```

Branch

Once you have done that, go into that directory and create a new branch:

```
cd graphtransliterator
git checkout -b [name_of_your_transliterator_branch]
```

For this example, you can use the branch `a_to_b`:

```
cd graphtransliterator
git checkout -b a_to_b
```

Then, push that branch to the origin (your personal github fork):

```
git push origin [name_of_your_transliterator_branch]
```

Here that would be:

```
.. code-block:: bash
```

```
git push origin a_to_b
```

Next, add a remote upstream for Graph Transliterator (the official Graph Transliterator repo):

```
git remote add upstream https://github.com/seanpue/graphtransliterator.git
```

Sync

To update your local copy of the the remote (official Graph Transliterator repo), run:

```
git fetch upstream
```

To sync your personal fork with the remote, run:

```
git merge upstream/master
```

See <https://help.github.com/en/articles/syncing-a-fork> for more info. You can run the previous two commands at any time.

Commit

You can commit your changes by running:

```
git commit -m 'comment here about the commit'
```

Adding A Transliterator

To add a transliterator, the next step is to create a subdirectory in `transliterators`. For this tutorial, you can make a branch named `a_to_b`.

Note that this will be under `graphtransliterator/transliterators`, so from the root directory enter:

```
cd graphtransliterator/transliterators
mkdir [name_of_your_transliterator]
cd [name_of_your_transliterator]
```

For this example, you would enter:

```
cd graphtransliterator/transliterators
mkdir a_to_b
cd a_to_b
```

In the `graphtransliterator/transliterators/[name_of_your_transliterator]` directory, you will add:

- an `__init__.py`
- a YAML file in the “easy reading format”

- a JSON file that is a serialization of the transliterator (optional)
- a `tests` directory including a file named `[name_of_your_transliterator]_tests.yaml`
- a Python test named `test_[name_of_your_transliterator].py` (optional)

Here is a tree showing the file organization:

```
transliterators
├── {{source_to_target}}
│   ├── __init__.py
│   ├── {{source_to_target}}.json
│   └── {{source_to_target}}.yaml
└── tests
    ├── test_{{source_to_target}}.py
    └── {{source_to_target}}_tests.yaml
```

YAML File

The YAML file should contain the “easy reading” version of your transliterator. For this example, create a file called `a_to_b.yaml`. Add a `metadata` field to the YAML file, as well, following the guidelines.

```
tokens:
  a: [a_class]
  ' ': [whitespace]
rules:
  a: A
onmatch_rules:
  - <a_class> + <a_class>: ",", "
whitespace:
  default: ' '
  token_class: whitespace
  consolidate: false
metadata:
  name: A to B
  version: 0.0.1
  url: http://website_of_project.com
  author: Your Name is Optional
  author_email: your_email@is_option.al
  maintainer: Maintainer's Name is Optional
  maintainer_email: maintainers_email@is_option.al
  license: MIT or Other Open Source License
  keywords: [add, keywords, here, as, a, list]
  project_urls:
    Documentation: https://link_to_documentation.html
    Source: https://link_to_sourcecode.html
    Tracker: https://link_to_issue_tracker.html
```

For most use cases, the `project_urls` can link to the Graph Transliteritor Github page.

JSON File

To create a JSON file, you can use the command line interface:

```
$ graphtransliterator dump --from yamL_file a_to_b.yamL > a_to_b.json
```

Alternatively, you can use the `make-json` command:

```
$ graphtransliterator make-json AToB
```

The JSON file loads more quickly than the YAML one, but it is not necessary during development.

`__init__.py`

The `__init__.py` will create the bundled transliterator, which is a subclass of *GraphTransliterator* named *Bundled*.

Following convention, you need to name your transliterator's class is CamelCase. For this example, it would be `AToB`:

```
from graphtransliterator.transl iterators import Bundled

class AToB(Bundled):
    """
    A to B Bundled Graph Transliterator
    """

    def __init__(self, **kwargs):
        """Initialize transliterator from YAML."""
        self.from_YAML(
            **kwargs
        ) # defaults to check_ambiguity=True, check_coverage=True
        # When ready, remove the previous lines and initialize more quickly from JSON:
        # self.init_from_JSON(**kwargs) # check_ambiguity=False, check_coverage=False
```

When you load the bundled transliterator from YAML using `from_YAML` it will check for ambiguity as well as check the coverage of the tests. You can turn those features off temporarily here.

When a transliterator is added into Graph Transliterator, it will likely be set to load from JSON by default. Tests will check for ambiguity and coverage.

Tests

Graph Transliterator requires that all bundled transliterators have tests that visit every edge and node of the internal graph and that use all on-match rules. The test file should be a YAML file defining a dictionary keyed from input to correct output.

You can test the transliterator as you are developing it by adding YAML tests and running the command:

```
graphtransliterator test [name_of_your_transliterator]
```

Tests can be generated using the command line interface:

```
mkdir tests
graphtransliterator generate-tests --from bundled [name_of_your_transliterator] >_
↪ tests/[name_of_your_transliterator]
```

Testing the Transliterator

You should test the transliterator to make sure everything is correct, including its metadata. To do that, navigate back to the root directory of *graphtransliterator* and execute the command:

```
py.test tests/test_transliterators.py
```

You can also run the complete suite of tests by running:

```
tox
```

Pushing Your Transliterator

When you are finished with a version of your transliterator, you should once again commit it to your github branch after syncing your branch with the remote. Then you can make a pull request to include the transliterator in Graph Transliterator. You can do that from the Graph Transliterator Github page. See <https://help.github.com/en/articles/creating-a-pull-request-from-a-fork>.

3.2.7 Contributing

Contributions are welcome, and they are greatly appreciated! Every little bit helps, and credit will always be given.

Contributor Code of Conduct

Please note that this project is released with a Contributor Code of Conduct. By participating in this project you agree to abide by its terms.

Types of Contributions

You can contribute in many ways:

Report Bugs

Report bugs at <https://github.com/seanpue/graphtransliterator/issues>.

If you are reporting a bug, please include:

- Your operating system name and version.
- Any details about your local setup that might be helpful in troubleshooting.
- Detailed steps to reproduce the bug.

Fix Bugs

Look through the GitHub issues for bugs. Anything tagged with “bug” and “help wanted” is open to whoever wants to implement it.

Implement Features

Look through the GitHub issues for features. Anything tagged with “enhancement” and “help wanted” is open to whoever wants to implement it.

Write Documentation

Graph-based Transliterator could always use more documentation, whether as part of the official Graph-based Transliterator docs, in docstrings, or even on the web in blog posts, articles, and such.

Submit Feedback

The best way to send feedback is to file an issue at <https://github.com/seanpue/graphtransliterator/issues>.

If you are proposing a feature:

- Explain in detail how it would work.
- Keep the scope as narrow as possible, to make it easier to implement.
- Remember that this is a volunteer-driven project, and that contributions are welcome :)

Add Transliterators

We welcome new transliterators to be added to the bundled transliterators!

See the documentation about Bundled Transliterators and look at Example as a model.

Raise an issue on Github, <https://github.com/seanpue/graphtransliterator/issues>

Then create a new branch with the new transliterator. Make sure the transliterator passes all of these requirements:

- is a submodule of `graphtransliterator.transliterators`
- has a unique name, preferably in format `source_to_target`
- has the following files: - `__init__.py` - `{{source_to_target}}.yaml` - `{{source_to_target}}.json` - `tests/{{source_to_target}}_tests.yaml` - `tests/test_{{source_to_target}}.py` (optional)
- has a classname in camel case, e.g. `SourceToTarget`
- has complete test coverage of all nodes and edges of generated graph and all onmatch rules, if present
- has required metadata in the YAML file.

When all the requirements are fulfilled, submit a pull request, and it will be reviewed for inclusion in a near-future release.

Get Started!

Ready to contribute? Here's how to set up *graphtransliterator* for local development.

1. Fork the *graphtransliterator* repo on GitHub.
2. Clone your fork locally:

```
$ git clone git@github.com:your_name_here/graphtransliterator.git
```

3. Install your local copy into a virtualenv. Assuming you have virtualenvwrapper installed, this is how you set up your fork for local development:

```
$ mkvirtualenv graphtransliterator
$ cd graphtransliterator/
$ python setup.py develop
```

4. Create a branch for local development:

```
$ git checkout -b name-of-your-bugfix-or-feature
```

Now you can make your changes locally.

5. When you're done making changes, format your code using the Black code formatter. (You can do that in your editor, as well). Then check that your changes pass flake8 and the tests, including testing other Python versions with tox:

```
$ black graphtransliterator
$ flake8 graphtransliterator tests
$ python setup.py test or py.test
$ tox
```

To get black, flake8, and tox, just pip install them into your virtualenv.

You should also test your coverage using make:

```
$ make coverage
```

6. Commit your changes and push your branch to GitHub:

```
$ git add .
$ git commit -m "Your detailed description of your changes."
$ git push origin name-of-your-bugfix-or-feature
```

7. Submit a pull request through the GitHub website.

Pull Request Guidelines

Before you submit a pull request, check that it meets these guidelines:

1. The pull request should include tests.
2. If the pull request adds functionality, the docs should be updated. Put your new functionality into a function with a docstring, and add the feature to the list in README.rst.
3. The pull request should work for Python 3.7 and 3.8 for PyPy. Check https://travis-ci.org/seanpue/graphtransliterator/pull_requests and make sure that the tests pass for all supported Python versions.

Tips

To run a subset of tests:

```
$ py.test tests.test_graphtransliterator
```

Deploying

A reminder for the maintainers on how to deploy. Make sure all your changes are committed (including an entry in HISTORY.rst). Then run:

```
$ bumpversion patch # possible: major / minor / patch
$ git push
$ git push --tags
```

The module uses Github Actions to deploy to TestPyPI and to PyPI.

3.2.8 API Reference

A list of the full API reference of all public classes and functions is below.

Public members can (and should) be imported from *graphtransliterator*:

```
from graphtransliterator import GraphTransliterator
```

Bundled transliterators require that *graphtransliterator.translitterators*: be imported:

```
import graphtransliterator.translitterators
translitterators.iter_names()
```

Core Classes

```
class graphtransliterator.GraphTransliterator (tokens, rules, whitespace, onmatch_rules=None,
                                                metadata=None, ignore_errors=False,
                                                check_ambiguity=True,
                                                onmatch_rules_lookup=None,
                                                tokens_by_class=None, graph=None,
                                                tokenizer_pattern=None,
                                                graphtransliterator_version=None, **kwargs)
```

A graph-based transliteration tool that lets you convert the symbols of one language or script to those of another using rules that you define.

Transliteration of tokens of an input string to an output string is configured by: a set of input token types with classes, pattern-matching rules involving sequences of tokens as well as preceding or following tokens and token classes, insertion rules between matches, and optional consolidation of whitespace. Rules are ordered by specificity.

Note: This constructor does not validate settings and should typically not be called directly. Use *from_dict()* instead. For “easy reading” support, use *from_easyreading_dict()*, *from_yaml()*, or *from_yaml_file()*. Keyword parameters used here (*ignore_errors*, *check_ambiguity*) can be passed from those other constructors.

Parameters

- **tokens** (*dict* of {*str*: *set* of *str*}) – Mapping of input token types to token classes
- **rules** (*list* of *TransliterationRule*) – *list* of transliteration rules ordered by cost
- **onmatch_rules** (*list* of *OnMatchRule*, or *None*) – Rules for output to be inserted between tokens of certain classes when a transliteration rule has been matched but before its production string has been added to the output
- **whitespace** (*WhitespaceRules*) – Rules for handling whitespace
- **metadata** (*dict* or *None*) – Metadata settings
- **ignore_errors** (*bool*, optional) – If true, transliteration errors are ignored and do not raise an exception. The default is false.
- **check_ambiguity** (*bool*, optional) – If true (default), transliteration rules are checked for ambiguity. *load()* and *loads()* do not check ambiguity by default.
- **onmatch_rules_lookup** (*dict* of {*str*: *dict* of {*str*: *list* of *int*}}, optional) – On-MatchRules lookup, used internally, will be generated if not present.
- **tokens_by_class** (*dict* of {*str*: *set* of *str*}, optional) – Tokens by class, used internally, will be generated if not present.
- **graph** (*DirectedGraph*, optional) – Directed graph used by Graph Transliterator, will be generated if not present.
- **tokenizer_pattern** (*str*, optional) – Regular expression pattern for input string tokenization, will be generated if not present.
- **graphtransliterator_version** (*str*, optional) – Version of graphtransliterator, added by *dump()* and *dumps()*.

Example

```

1 from graphtransliterator import GraphTransliterator, OnMatchRule, \
  ↪ TransliterationRule, WhitespaceRules
2 settings = {'tokens': {'a': {'vowel'}, ' ': {'wb'}}, 'onmatch_rules': \
  ↪ [OnMatchRule(prev_classes=['vowel'], next_classes=['vowel'], production=',')],
  ↪ 'rules': [TransliterationRule(production='A', prev_classes=None, prev_
  ↪ tokens=None, tokens=['a'], next_tokens=None, next_classes=None, cost=0.
  ↪ 5849625007211562), TransliterationRule(production=' ', prev_classes=None, prev_
  ↪ tokens=None, tokens=[' '], next_tokens=None, next_classes=None, cost=0.
  ↪ 5849625007211562)], 'metadata': {'author': 'Author McAuthorson'}, 'whitespace': \
  ↪ WhitespaceRules(default=' ', token_class='wb', consolidate=False)}
3 gt = GraphTransliterator(**settings)
4 gt.transliterate('a')
```

```
'A'
```

See also:

from_dict

Constructor from dictionary of settings

from_easyreading_dict

Constructor from dictionary in “easy reading” format

from_yaml

Constructor from YAML string in “easy reading” format

from_yaml_file

Constructor from YAML file in “easy reading” format

dump (*compression_level=0*)

Dump configuration of Graph Transliterator to Python data types.

Compression is turned off by default.

Parameters

compression_level (*int*) – A value in 0 (default, no compression), 1 (compression including graph), and 2 (compression without graph)

Returns

GraphTransliterator configuration as a dictionary with keys:

"tokens"

Mappings of tokens to their classes (*OrderedDict* of {*str*: list of *str*})

"rules"

Transliteration rules in direct format (*list* of *dict* of {*str*: *str*})

"whitespace"

Whitespace settings (*dict* of {*str*: *str*})

"onmatch_rules"

On match rules (*list* of *OrderedDict*)

"metadata"

Dictionary of metadata (*dict*)

"ignore_errors"

Ignore errors in transliteration (*bool*)

"onmatch_rules_lookup"

Dictionary keyed by current token to previous token containing a list of indexes of applicable *OnmatchRule* to try (*dict* of {*str*: *dict* of {*str*: list of *int*}})

"tokens_by_class"

Tokens keyed by token class, used internally (*dict* of {*str*: list of *str*})

"graph"

Serialization of *DirectedGraph* (*dict*)

"tokenizer_pattern"

Regular expression for tokenizing (*str*)

"graphtransliterator_version"

Module version of *graphtransliterator* (*str*)

Return type

OrderedDict

Example

```

5 yam1_ = '''
6 tokens:
7   a: [vowel]
8   ' ': [wb]
9 rules:
10  a: A
11  ' ': ' '
12 whitespace:
13   default: " "
14   consolidate: false
15   token_class: wb
16 onmatch_rules:
17   - <vowel> + <vowel>: ', ' # add a comma between vowels
18 metadata:
19   author: "Author McAuthorson"
20 '''
21 gt = GraphTransliterator.from_yaml(yam1_)
22 gt.dump()

```

```

OrderedDict([('tokens', {'a': ['vowel'], ' ': ['wb']}),
             ('rules',
              [OrderedDict([('production', 'A'),
                           ('tokens', ['a']),
                           ('cost', 0.5849625007211562)]),
               OrderedDict([('production', ' '),
                           ('tokens', [' ']),
                           ('cost', 0.5849625007211562)]))],
             ('whitespace',
              {'default': ' ', 'token_class': 'wb', 'consolidate': False}),
             ('onmatch_rules',
              [OrderedDict([('prev_classes', ['vowel']),
                           ('next_classes', ['vowel']),
                           ('production', ',')])]),
             ('metadata', {'author': 'Author McAuthorson'}),
             ('ignore_errors', False),
             ('onmatch_rules_lookup', {'a': {'a': [0]}},
             ('tokens_by_class', {'vowel': ['a'], 'wb': [' ']}),
             ('graph',
              {'node': [{'type': 'Start',
                          'ordered_children': {'a': [1], ' ': [3]}},
                       {'token': 'a',
                        'type': 'token',
                        'ordered_children': {'__rules__': [2]}},
                       {'type': 'rule', 'accepting': True, 'rule_key': 0},
                       {'token': ' ',
                        'type': 'token',
                        'ordered_children': {'__rules__': [4]}},
                       {'type': 'rule', 'accepting': True, 'rule_key': 1}],
              'edge': {0: {1: {'token': 'a', 'cost': 0.5849625007211562},
                          3: {'token': ' ', 'cost': 0.5849625007211562}},
                       1: {2: {'cost': 0.5849625007211562}},
                       3: {4: {'cost': 0.5849625007211562}}},
              'edge_list': [(0, 1), (0, 3), (1, 2), (3, 4)]},
             ('tokenizer_pattern', '(a|\\ )'),
             ('graphtransliterator_version', '1.2.4')])

```

See also:

dumps

Dump Graph Transliterator configuration to JSON string

load

Load Graph Transliteration from configuration in Python data types

loads

Load Graph Transliteration from configuration as a JSON string

dumps (*compression_level=2*)

Parameters

- **compression_level** (*int*) – A value in 0 (no compression), 1 (compression including graph), and 2 (default, compression without graph)
- **separators** (*tuple of str*) – Separators used by `json.dumps()`, default is compact
- **(JSON)** . (*Dump settings of Graph Transliterator to Javascript Object Notation*) –
- **default** . (*Compression is turned on by*) –

Returns

JSON string

Return type

str

Examples

```
23 yaml_ = '''
24     tokens:
25         a: [vowel]
26         ' ': [wb]
27     rules:
28         a: A
29         ' ': ' '
30     whitespace:
31         default: " "
32         consolidate: false
33         token_class: wb
34     onmatch_rules:
35         - <vowel> + <vowel>: ',', # add a comma between vowels
36     metadata:
37         author: "Author McAuthorson"
38     '''
39 gt = GraphTransliterator.from_yaml(yaml_)
40 gt.dumps()
```

```
'{"graphtransliterator_version":"1.2.4","compressed_settings":[[{"vowel","wb"},
↪ [{" ","a"],[[1],[0]],[["A",0,0,[1],0,0,-1],[[" ",0,0,[0],0,0,-1],[[" ","wb",
↪ 0],[[["0],[0],[""],{"author":"Author McAuthorson"},null]]]'
```

See also:

dump

Dump Graph Transliterator configuration to Python data types

load

Load Graph Transliteration from configuration in Python data types

loads

Load Graph Transliteration from configuration as a JSON string

static from_dict (*dict_settings*, ***kwargs*)

Generate GraphTransliterator from *dict* settings.

Parameters

dict_settings (*dict*) – Dictionary of settings

Returns

Graph transliterator

Return type

GraphTransliterator

static from_easyreading_dict (*easyreading_settings*, ***kwargs*)

Constructs *GraphTransliterator* from a dictionary of settings in “easy reading” format, i.e. the loaded contents of a YAML string.

Parameters

easyreading_settings (*dict*) – Settings dictionary in easy reading format with keys:

"tokens"

Mappings of tokens to their classes (*dict* of {*str*: *list* of *str*})

"rules"

Transliteration rules in “easy reading” format (*list* of *dict* of {*str*: *str*})

"onmatch_rules"

On match rules in “easy reading” format (*dict* of {*str*: *str*}, optional)

"whitespace"

Whitespace definitions, including default whitespace token, class of whitespace tokens, and whether or not to consolidate (*dict* of {'default': *str*, 'token_class': *str*, consolidate: *bool*}, optional)

"metadata"

Dictionary of metadata (*dict*, optional)

Returns

Graph Transliterator

Return type

GraphTransliterator

Note: Called by *from_yaml()*.

Example

```
41 tokens = {
42     'ab': ['class_ab'],
43     ' ': ['wb']
44 }
45 whitespace = {
46     'default': ' ',
47     'token_class': 'wb',
48     'consolidate': True
49 }
50 onmatch_rules = [
51     {'<class_ab> + <class_ab>': ', '}
52 ]
53 rules = {'ab': 'AB',
54         ' ': '_'}
55 settings = {'tokens': tokens,
56            'rules': rules,
57            'whitespace': whitespace,
58            'onmatch_rules': onmatch_rules}
59 gt = GraphTransliterator.from_easyreading_dict(settings)
60 gt.transliterate("ab abab")
```

```
'AB_AB,AB'
```

See also:

from_yaml

Constructor from YAML string in “easy reading” format

from_yaml_file

Constructor from YAML file in “easy reading” format

static from_yaml (*yaml_str*, *charnames_escaped*=True, ***kwargs*)

Construct GraphTransliterator from a YAML str.

Parameters

- **yaml_str** (*str*) – YAML mappings of tokens, rules, and (optionally) onmatch_rules
- **charnames_escaped** (*boolean*) – Unescape Unicode during YAML read (default True)

Note: Called by *from_yaml_file()* and calls *from_easyreading_dict()*.

Example

```

61 yaml_ = '''
62 tokens:
63     a: [class1]
64     ' ': [wb]
65 rules:
66     a: A
67     ' ': ' '
68 whitespace:
69     default: ' '
70     consolidate: True
71     token_class: wb
72 onmatch_rules:
73     - <class1> + <class1>: "+"
74 '''
75 gt = GraphTransliteritor.from_yaml(yaml_)
76 gt.transliterate("a aa")

```

```
'A A+A'
```

See also:

from_easyreading_dict

Constructor from dictionary in “easy reading” format

from_yaml

Constructor from YAML string in “easy reading” format

from_yaml_file

Constructor from YAML file in “easy reading” format

static *from_yaml_file* (yaml_filename, **kwargs)

Construct GraphTransliteritor from YAML file.

Parameters

yaml_filename (*str*) – Name of YAML file, containing tokens, rules, and (optionally) onmatch_rules

Note: Calls *from_yaml* ().

See also:

from_yaml

Constructor from YAML string in “easy reading” format

from_easyreading_dict

Constructor from dictionary in “easy reading” format

property graph

Graph used in transliteration.

Type

DirectedGraph

property graphtransliterator_version

Graph Transliterator version.

Type

str

property ignore_errors

Ignore transliteration errors setting.

Type

bool

property last_input_tokens

Last tokenization of the input string, with whitespace at start and end.

Type

list of str

property last_matched_rule_tokens

Last matched tokens for each rule.

Type

list of list of str

property last_matched_rules

Last transliteration rules matched.

Type

list of TransliterationRule

static load (*settings*, ***kwargs*)

Create GraphTransliterator from settings as Python data types.

Parameters

settings – GraphTransliterator configuration as a dictionary with keys:

"tokens"

Mappings of tokens to their classes (*dict of {str: list of str}*)

"rules"

Transliteration rules in direct format (*list of OrderedDict of {str: str}*)

"whitespace"

Whitespace settings (*dict of {str: str}*)

"onmatch_rules"

On match rules (*list of OrderedDict, optional*)

"metadata"

Dictionary of metadata (*dict, optional*)

"ignore_errors"

Ignore errors. (*bool, optional*)

"onmatch_rules_lookup"

Dictionary keyed by current token to previous token containing a list of indexes of applicable OnmatchRule to try (*dict of {str: dict of {str: list of int}}*, optional)

tokens_by_class

Tokens keyed by token class, used internally (*dict of {str: list of str}*, optional)

graph

Serialization of *DirectedGraph* (*dict, optional*)

"tokenizer_pattern"

Regular expression for tokenizing (*str*, optional)

"graphtransliterator_version"

Module version of *graphtransliterator* (*str*, optional)

Returns

Graph Transliterator

Return type

GraphTransliterator

Example

```

77 from collections import OrderedDict
78 settings = {'tokens': {'a': ['vowel'], ' ': ['wb']}},
79 'rules': [OrderedDict([('production', 'A'),
80                        # Can be compacted, removing None values
81                        # ('prev_tokens', None),
82                        ('tokens', ['a']),
83                        ('next_classes', None),
84                        ('next_tokens', None),
85                        ('cost', 0.5849625007211562)])],
86 OrderedDict([('production', ' '),
87               ('prev_classes', None),
88               ('prev_tokens', None),
89               ('tokens', [' ']),
90               ('next_classes', None),
91               ('next_tokens', None),
92               ('cost', 0.5849625007211562)])],
93 'whitespace': {'default': ' ', 'token_class': 'wb', 'consolidate': False},
94 'onmatch_rules': [OrderedDict([('prev_classes', ['vowel']),
95                                ('next_classes', ['vowel']),
96                                ('production', ' ')])],
97 'metadata': {'author': 'Author McAuthorson'},
98 'onmatch_rules_lookup': {'a': {'a': [0]}},
99 'tokens_by_class': {'vowel': ['a'], 'wb': [' ']},
100 'graph': {'edge': {0: {1: {'token': 'a', 'cost': 0.5849625007211562},
101                          3: {'token': ' ', 'cost': 0.5849625007211562}},
102            1: {2: {'cost': 0.5849625007211562}},
103            3: {4: {'cost': 0.5849625007211562}}},
104 'node': [{ 'type': 'Start', 'ordered_children': {'a': [1], ' ': [3]}},
105           { 'type': 'token', 'token': 'a', 'ordered_children': {'__rules__': [2]}},
106           { 'type': 'rule',
107             'rule_key': 0,
108             'accepting': True,
109             'ordered_children': {}},
110           { 'type': 'token', 'token': ' ', 'ordered_children': {'__rules__': [4]}},
111           { 'type': 'rule',
112             'rule_key': 1,
113             'accepting': True,
114             'ordered_children': {}},
115           'edge_list': [(0, 1), (1, 2), (0, 3), (3, 4)],
116           'tokenizer_pattern': '(a| )',
117           'graphtransliterator_version': '0.3.3'}
118 gt = GraphTransliterator.load(settings)
119 gt.transliterate('aa')

```

```
'A, A'
```

```
120 # can be compacted
121 settings.pop('onmatch_rules_lookup')
122 GraphTransliterator.load(settings).transliterate('aa')
```

```
'A, A'
```

See also:

dump

Dump Graph Transliterators configuration to Python data types

dumps

Dump Graph Transliterators configuration to JSON string

loads

Load Graph Transliteration from configuration as a JSON string

static loads (*settings*, ***kwargs*)

Create GraphTransliterators from JavaScript Object Notation (JSON) string.

Parameters

settings – JSON settings for GraphTransliterators

Returns

Graph Transliterators

Return type

GraphTransliterators

Example

```
123 JSON_settings = '''{"tokens": {"a": ["vowel"], " ": ["wb"]}, "rules": [{
↪ "production": "A", "prev_classes": null, "prev_tokens": null, "tokens": ["a
↪ "], "next_classes": null, "next_tokens": null, "cost": 0.5849625007211562},
↪ {"production": " ", "prev_classes": null, "prev_tokens": null, "tokens": ["
↪ "], "next_classes": null, "next_tokens": null, "cost": 0.5849625007211562}},
↪ "whitespace": {"default": " ", "token_class": "wb", "consolidate": false},
↪ "onmatch_rules": [{"prev_classes": ["vowel"], "next_classes": ["vowel"],
↪ "production": ",,"}], "metadata": {"author": "Author McAuthorson"}, "ignore_
↪ errors": false, "onmatch_rules_lookup": {"a": {"a": [0]}}, "tokens_by_class
↪ ": {"vowel": ["a"], "wb": [" "]}, "graph": {"node": [{"type": "Start",
↪ "ordered_children": {"a": [1], " ": [3]}}, {"type": "token", "token": "a",
↪ "ordered_children": {"__rules__": [2]}}, {"type": "rule", "rule_key": 0,
↪ "accepting": true, "ordered_children": {}}, {"type": "token", "token": " ",
↪ "ordered_children": {"__rules__": [4]}}, {"type": "rule", "rule_key": 1,
↪ "accepting": true, "ordered_children": {}}], "edge": {"0": {"1": {"token":
↪ "a", "cost": 0.5849625007211562}, "3": {"token": " ", "cost": 0.
↪ 5849625007211562}}, "1": {"2": {"cost": 0.5849625007211562}, "3": {"4": {
↪ "cost": 0.5849625007211562}}, "edge_list": [[0, 1], [1, 2], [0, 3], [3,
↪ 4]]}, "tokenizer_pattern": "(a| )", "graphtransliterators_version": "1.2.2"}'
↪ '''
124
125 gt = GraphTransliterators.loads(JSON_settings)
126 gt.transliterate('a')
```

```
'A'
```

See also:*dump*

Dump Graph Transliterator configuration to Python data types

dumps

Dump Graph Transliterator configuration to JSON string

load

Load Graph Transliteration from configuration in Python data types

match_at (*token_i*, *tokens*, *match_all=False*)

Match best (least costly) transliteration rule at a given index in the input tokens and return the index to that rule. Optionally, return all rules that match.

Parameters

- **token_i** (*int*) – Location in *tokens* at which to begin
- **tokens** (*list of str*) – List of tokens
- **match_all** (*bool*, optional) – If true, return the index of all rules matching at the given index. The default is false.

Returns

Index of matching transliteration rule in *GraphTransliterator.rules* or None. Returns a *list of int* or an empty *list* if *match_all* is true.

Return type

int, *None*, or *list of int*

Note: Expects whitespaces token at beginning and end of *tokens*.

Examples

```

127 gt = GraphTransliterator.from_yaml('''
128     tokens:
129         a: []
130         a a: []
131         ' ': [wb]
132     rules:
133         a: <A>
134         a a: <AA>
135     whitespace:
136         default: ' '
137         consolidate: True
138         token_class: wb
139 ''')
140 tokens = gt.tokenize("aa")
141 tokens # whitespace added to ends

```

```
[' ', 'a', 'a', ' ']
```

```
142 gt.match_at(1, tokens) # returns index to rule
0

143 gt.rules[gt.match_at(1, tokens)] # actual rule
TransliterationRule(production='<AA>', prev_classes=None, prev_tokens=None, ↵
↵tokens=['a', 'a'], next_tokens=None, next_classes=None, cost=0.
↵41503749927884376)

144 gt.match_at(1, tokens, match_all=True) # index to rules, with match_all
[0, 1]

145 [gt.rules[_] for _ in gt.match_at(1, tokens, match_all=True)]
[TransliterationRule(production='<AA>', prev_classes=None, prev_tokens=None, ↵
↵tokens=['a', 'a'], next_tokens=None, next_classes=None, cost=0.
↵41503749927884376),
TransliterationRule(production='<A>', prev_classes=None, prev_tokens=None, ↵
↵tokens=['a'], next_tokens=None, next_classes=None, cost=0.5849625007211562)]
```

property metadata

Metadata of transliterator

Type

dict

property onmatch_rules

Rules for productions between matches.

Type

list of OnMatchRules

property onmatch_rules_lookup

On Match Rules lookup

Type

dict

property productions

List of productions of each transliteration rule.

Type

list of str

pruned_of (*productions*)

Remove transliteration rules with specific output productions.

Parameters

productions (*str*, or *list of str*) – list of productions to remove

Returns

Graph transliterator pruned of certain productions.

Return type

graphtransliterator.GraphTransliterator

Note: Uses original initialization parameters to construct a new *GraphTransliterators*.

Examples

```

146 gt = GraphTransliterators.from_yaml(''
147     tokens:
148         a: []
149         a a: []
150         ' ': [wb]
151     rules:
152         a: <A>
153         a a: <AA>
154     whitespace:
155         default: ' '
156         consolidate: True
157         token_class: wb
158 '')
159 gt.rules

```

```

[TransliterationRule(production='<AA>', prev_classes=None, prev_tokens=None,
↳tokens=['a', 'a'], next_tokens=None, next_classes=None, cost=0.
↳41503749927884376),
 TransliterationRule(production='<A>', prev_classes=None, prev_tokens=None,
↳tokens=['a'], next_tokens=None, next_classes=None, cost=0.5849625007211562)]

```

```

160 gt.pruned_of('<AA>').rules

```

```

[TransliterationRule(production='<A>', prev_classes=None, prev_tokens=None,
↳tokens=['a'], next_tokens=None, next_classes=None, cost=0.5849625007211562)]

```

```

161 gt.pruned_of(['<A>', '<AA>']).rules

```

```

[]

```

property rules

Transliteration rules sorted by cost.

Type

list of TransliterationRule

tokenize (input)

Tokenizes an input string.

Adds initial and trailing whitespace, which can be consolidated.

Parameters

input (*str*) – String to tokenize

Returns

List of tokens, with default whitespace token at beginning and end.

Return type

list of str

Raises

ValueError – Unrecognizable input, such as a character that is not in a token

Examples

```
162 tokens = {'ab': ['class_ab'], ' ': ['wb']}
163 whitespace = {'default': ' ', 'token_class': 'wb', 'consolidate': True}
164 rules = {'ab': 'AB', ' ': '_'}
165 settings = {'tokens': tokens, 'rules': rules, 'whitespace': whitespace}
166 gt = GraphTransliterator.from_easyreading_dict(settings)
167 gt.tokenize('ab ')
```

```
[' ', 'ab', ' ']
```

property tokenizer_pattern

Tokenizer pattern from transliterator

Type

str

property tokens

Mappings of tokens to their classes.

Type

dict of {*str*

Type

set of *str*}

property tokens_by_class

Tokenizer pattern from transliterator

Type

dict of {*str*

Type

list of *str*}

transliterate (*input*)

Transliterate an input string into an output string.

Parameters

input (*str*) – Input string to transliterate

Returns

Transliteration output string

Return type

str

Raises

ValueError – Cannot parse input

Note: Whitespace will be temporarily appended to start and end of input string.

Example

```

168 GraphTransliterator.from_yaml (
169     '''
170     tokens:
171         a: []
172         ' ': [wb]
173     rules:
174         a: A
175         ' ': '_'
176     whitespace:
177         default: ' '
178         consolidate: True
179         token_class: wb
180     ''').transliterate("a a")

```

```
'A_A'
```

property whitespace

Whitespace rules.

Type

WhiteSpaceRules

class graphtransliterator.CoverageTransliterator (*args, **kwargs)

Subclass of GraphTransliterator that logs visits to graph and on_match rules.

Used to confirm that tests cover the entire graph and onmatch_rules.

check_coverage (raise_exception=True)

Check coverage of graph and onmatch rules.

First checks graph coverage, then checks onmatch rules.

check_onmatchrules_coverage (raise_exception=True)

Check coverage of onmatch rules.

clear_visited ()

Clear visited flags from graph and onmatch_rules.

Bundled Transliterators

graphtransliterator.transliterators

Bundled transliterators are loaded by explicitly importing `graphtransliterator.transliterators`. Each is an instance of `graphtransliterator.bundled.Bundled`.

class graphtransliterator.transliterators.Bundled (*args, **kwargs)

Subclass of GraphTransliterator used for bundled Graph Transliterator.

property directory

Directory of bundled transliterator, used to load settings.

from_JSON (check_ambiguity=False, coverage=False, **kwargs)

Initialize from bundled JSON file (best for speed).

Parameters

- **check_ambiguity** (*bool*,) – Should ambiguity be checked. Default is *False*.
- **coverage** (*bool*) – Should test coverage be checked. Default is *False*.

from_YAML (*check_ambiguity=True, coverage=True, **kwargs*)

Initialize from bundled YAML file (best for development).

Parameters

- **check_ambiguity** (*bool*,) – Should ambiguity be checked. Default is *True*.
- **coverage** (*bool*) – Should test coverage be checked. Default is *True*.

generate_yaml_tests (*file=None*)

Generates YAML tests with complete coverage.

Uses the first token in a class as a sample. Assumes for onmatch rules that the first sample token in a class has a unique production, which may not be the case. These should be checked and edited.

load_yaml_tests ()

Iterator for YAML tests.

Assumes tests are found in subdirectory *tests* of module with name *NAME_tests.yaml*, e.g. *'source_to_target/tests/source_to_target_tests.yaml'*.

property name

Name of bundled transliterator, e.g. 'Example'

classmethod new (*method='json', **kwargs*)

Return a new class instance from method (json/yaml).

Parameters

method (*str* (*json* or *yaml*)) – How to load bundled transliterator, JSON or YAML.

run_tests (*transliteration_tests*)

Run transliteration tests.

Parameters

transliteration_tests (*dict* of {*str:str*}) – Dictionary of test from source -> correct target.

run_yaml_tests ()

Run YAML tests in *MODULE/tests/MODULE_tests.yaml*

property yaml_tests_file

Metadata of transliterator

Type

dict

class graphtransliterator.transl iterators.**Example** (***kwargs*)

Example Bundled Graph Transliterator.

class graphtransliterator.transl iterators.**ITRANSDevanagariToUnicode** (***kwargs*)

ITRANS Devanagari to Unicode Transliterator.


```
class graphtransliterator.transl iterators.MetadataSchema (*, only: Sequence[str] |
    AbstractSet[str] | None = None,
    exclude: Sequence[str] |
    AbstractSet[str] = (), many: bool
    = False, context: dict | None =
    None, load_only: Sequence[str] |
    AbstractSet[str] = (),
    dump_only: Sequence[str] |
    AbstractSet[str] = (), partial:
    bool | Sequence[str] |
    AbstractSet[str] | None = None,
    unknown: str | None = None)
```

Schema for Bundled metadata.

```
graphtransliterator.transl iterators.iter_names()
```

Iterate through bundled transliterator names.

```
graphtransliterator.transl iterators.iter_transl iterators (**kws)
```

Iterate through instances of bundled transl iterators.

Graph Classes

```
class graphtransliterator.DirectedGraph (node=None, edge=None, edge_list=None)
```

A very basic dictionary- and list-based directed graph. Nodes are a list of dictionaries of node data. Edges are nested dictionaries keyed from the head -> tail -> edge properties. An edge list is maintained. Can be exported as a dictionary.

node

List of node data

Type

list of dict

edge

Mapping from head to tail of edge, holding edge data

Type

dict of {int: dict of {int: dict}}

edge_list

List of head and tail of each edge

Type

list of tuple of (int, int)

Examples

```
181 from graphtransliterator import DirectedGraph
182 DirectedGraph()
```

```
<graphtransliterator.graphs.DirectedGraph at 0x7f250c24c500>
```

add_edge (*head*, *tail*, *edge_data*=None)

Add an edge to a graph and return its attributes as dict.

Parameters

- **head** (*int*) – Index of head of edge
- **tail** (*int*) – Index of tail of edge
- **edge_data** (*dict*, default {}) – Edge data

Returns

Data of created edge

Return type

`dict`

Raises

ValueError – Invalid head or tail, or edge_data is not a dict.

Examples

```
183 g = DirectedGraph()
184 g.add_node()
```

```
(0, {})
```

```
185 g.add_node()
```

```
(1, {})
```

```
186 g.add_edge(0,1, {'data_key_1': 'some edge data here'})
```

```
{'data_key_1': 'some edge data here'}
```

```
187 g.edge
```

```
{0: {1: {'data_key_1': 'some edge data here'}}}
```

add_node (*node_data*=None)

Create node and return (*int*, *dict*) of node key and object.

Parameters

node_data (*dict*, default {}) – Data to be stored in created node

Returns

Index of created node and its data

Return type

tuple of (*int*, *dict*)

Raises

ValueError – node_data is not a dict

Examples

```
188 g = DirectedGraph()
189 g.add_node()
```

```
(0, {})
```

```
190 g.add_node({'datakey1': 'data value'})
```

```
(1, {'datakey1': 'data value'})
```

```
191 g.node
```

```
[{}, {'datakey1': 'data value'}]
```

class graphtransliterator.**VisitLoggingDirectedGraph**(*graph*)

A DirectedGraph that logs visits to all nodes and edges.

Used to measure the coverage of tests for bundled transliterators.

check_coverage (*raise_exception=True*)

Checks that all nodes and edges are visited.

Parameters

raise_exception (*bool*, default) – Raise IncompleteGraphCoverageException (default, *True*)

Raises

IncompleteGraphCoverageException – Not all nodes/edges of a graph have been visited.

clear_visited()

Clear all visited attributes on nodes and edges.

Rule Classes

class graphtransliterator.**TransliterationRule**(*production, prev_classes, prev_tokens, tokens, next_tokens, next_classes, cost*)

A transliteration rule containing the specific match conditions and string output to be produced, as well as the rule's cost.

production

Output produced on match of rule

Type

str

prev_classes

List of previous token classes to be matched before tokens or, if they exist, *prev_tokens*

Type

list of str, or None

prev_tokens

List of tokens to be matched before *tokens*

Type

list of str, or None

tokens

List of tokens to match

Type

list of str

next_tokens

List of tokens to match after *tokens*

Type

list of str, or None

next_classes

List of tokens to match after *tokens* or, if they exist, *next_tokens*

Type

list of str, or None

cost

Cost of the rule, where less specific rules are more costly

Type

float

class graphtransliterator.**OnMatchRule** (*prev_classes, next_classes, production*)

Rules about adding text between certain combinations of matched rules.

When a translation rule has been found and before its production is added to the output, the productions string of an OnMatch rule is added if previously matched tokens and current tokens are of the specified classes.

prev_classes

List of previously matched token classes required

Type

list of str

next_classes

List of current and following token classes required

Type

list of str

production

String to added before current rule

Type

str

class graphtransliterator.**WhitespaceRules** (*default, token_class, consolidate*)

Whitespace rules of GraphTransliterator.

default

Default whitespace token

Type

str

token_class

Whitespace token class

Type

str

consolidate

Consolidate consecutive whitespace tokens and render as a single instance of the specified default whitespace token.

Type

bool

Exceptions

exception graphtransliterator.**GraphTransliteratorException**

Base exception class. All Graph Transliterator-specific exceptions should subclass this class.

exception graphtransliterator.**AmbiguousTransliterationRulesException**

Raised when multiple transliteration rules can match the same pattern. Details of ambiguities are given in a `logging.warning()`.

exception graphtransliterator.**NoMatchingTransliterationRuleException**

Raised when no transliteration rule can be matched at a particular location in the input string's tokens. Details of the location are given in a `logging.warning()`.

exception graphtransliterator.**UnrecognizableInputTokenException**

Raised when a character in the input string does not correspond to any tokens in the GraphTransliterator's token settings. Details of the location are given in a `logging.warning()`.

Schemas

class graphtransliterator.**DirectedGraphSchema** (*, only: *Sequence[str] | AbstractSet[str] | None = None*, exclude: *Sequence[str] | AbstractSet[str] = ()*, many: *bool = False*, context: *dict | None = None*, load_only: *Sequence[str] | AbstractSet[str] = ()*, dump_only: *Sequence[str] | AbstractSet[str] = ()*, partial: *bool | Sequence[str] | AbstractSet[str] | None = None*, unknown: *str | None = None*)

Schema for *DirectedGraph*.

Validates graph somewhat rigorously.

class graphtransliterator.**EasyReadingSettingsSchema** (*, only: *Sequence[str] | AbstractSet[str] | None = None*, exclude: *Sequence[str] | AbstractSet[str] = ()*, many: *bool = False*, context: *dict | None = None*, load_only: *Sequence[str] | AbstractSet[str] = ()*, dump_only: *Sequence[str] | AbstractSet[str] = ()*, partial: *bool | Sequence[str] | AbstractSet[str] | None = None*, unknown: *str | None = None*)

Schema for easy reading settings.

Provides initial validation based on easy reading format.

```
class graphtransliterator.GraphTransliteratorSchema (*, only: Sequence[str] | AbstractSet[str] |  
None = None, exclude: Sequence[str] |  
AbstractSet[str] = (), many: bool =  
False, context: dict | None = None,  
load_only: Sequence[str] |  
AbstractSet[str] = (), dump_only:  
Sequence[str] | AbstractSet[str] = (),  
partial: bool | Sequence[str] |  
AbstractSet[str] | None = None,  
unknown: str | None = None)
```

Schema for Graph Transliterator.

```
class graphtransliterator.OnMatchRuleSchema (*, only: Sequence[str] | AbstractSet[str] | None =  
None, exclude: Sequence[str] | AbstractSet[str] = (),  
many: bool = False, context: dict | None = None,  
load_only: Sequence[str] | AbstractSet[str] = (),  
dump_only: Sequence[str] | AbstractSet[str] = (),  
partial: bool | Sequence[str] | AbstractSet[str] | None  
= None, unknown: str | None = None)
```

Schema for *OnMatchRule*.

```
class graphtransliterator.SettingsSchema (*, only: Sequence[str] | AbstractSet[str] | None = None,  
exclude: Sequence[str] | AbstractSet[str] = (), many: bool  
= False, context: dict | None = None, load_only:  
Sequence[str] | AbstractSet[str] = (), dump_only:  
Sequence[str] | AbstractSet[str] = (), partial: bool |  
Sequence[str] | AbstractSet[str] | None = None, unknown:  
str | None = None)
```

Schema for settings in dictionary format.

Performs validation.

```
class graphtransliterator.TransliterationRuleSchema (*, only: Sequence[str] | AbstractSet[str] |  
None = None, exclude: Sequence[str] |  
AbstractSet[str] = (), many: bool =  
False, context: dict | None = None,  
load_only: Sequence[str] |  
AbstractSet[str] = (), dump_only:  
Sequence[str] | AbstractSet[str] = (),  
partial: bool | Sequence[str] |  
AbstractSet[str] | None = None,  
unknown: str | None = None)
```

Schema for *TransliterationRule*.

```
class graphtransliterator.WhitespaceDictSettingsSchema (*, only: Sequence[str] |
    AbstractSet[str] | None = None,
    exclude: Sequence[str] |
    AbstractSet[str] = (), many: bool =
    False, context: dict | None = None,
    load_only: Sequence[str] |
    AbstractSet[str] = (), dump_only:
    Sequence[str] | AbstractSet[str] =
    (), partial: bool | Sequence[str] |
    AbstractSet[str] | None = None,
    unknown: str | None = None)
```

Schema for Whitespace definition as a *dict*.

```
class graphtransliterator.WhitespaceSettingsSchema (*, only: Sequence[str] | AbstractSet[str] |
    None = None, exclude: Sequence[str] |
    AbstractSet[str] = (), many: bool = False,
    context: dict | None = None, load_only:
    Sequence[str] | AbstractSet[str] = (),
    dump_only: Sequence[str] |
    AbstractSet[str] = (), partial: bool |
    Sequence[str] | AbstractSet[str] | None =
    None, unknown: str | None = None)
```

Schema for Whitespace definition that loads as `WhitespaceRules`.

3.2.9 Credits

Development Lead

- A. Sean Pue @seanpue <pue@msu.edu>

Contributors

- Valentino Constantinou @vc1492a
- Rebecca Sutton Koeser @rlskoeser

3.2.10 Acknowledgements

Software development was supported by an Andrew W. Mellon Foundation New Directions Fellowship (Grant Number 11600613) and by matching funds provided by the College of Arts and Letters, Michigan State University.

3.2.11 Kudos

Graph Transliterator's developers acknowledge the following open-access projects, which have been particularly helpful in Graph Transliterator's development. These include: [astropy](#) (guide for documentation style expanding on [numpy](#)), [click](#) (command line interface), [contributor_covenant](#) (basis for the code of conduct), [cookiecutter-pypackage](#) (initial Python module template), [jupyter-sphinx](#) (renderer of live code results in the docs), and [marshmallow](#) (object serializer/deserializer).

Those from which code/text has been adopted are mentioned in [NOTICE](#).

[~ Dependencies scanned by PyUp.io ~]

3.2.12 History

[Unreleased - Maybe]

- save match location in tokenize using token_details
- allow insertion of transliteration error messages into output
- fix Devanagari output in Sphinx-generated Latex PDF
- add translated messages
- add static typing with mypy
- adjust IncorrectVersionException to only consider major, minor versioning not patch
- Adjust CSS for CLI output in docs
- add doc making test to commit

[To do]

- Add on/off switch characters
- Update module publication

1.2.4 (2023-10-15)

- switched to poetry for module publishing, based on cookiecutter-poetry
- fixed click version command
- added code coverage

1.2.3 (2023-10-09)

- added python 3.10, 3.11, removed <=3.8
- updated dependencies (used pur)
- updated jupyter-download syntax
- reformatted with black
- adjusted flake8 line length
- removed collect_ignore for pytest
- updated Github actions

1.2.2 (2021-08-11)

- updated CONTRIBUTING.rst for new Python versions
- added github actions to publish to pypi and testpypi
- shifted to github CI
- updated dependencies
- fixed tox.ini
- updated schema.py error message
- updated docs/conf.py for jupyter_sphinx

1.2.1 (2020-10-29)

- updated docs/conf.py for jupyter_sphinx

1.2.0 (2020-05-13)

- changes to bundled.py and cli.py with dump-tests command
- updated cli.rst

1.1.2 (2020-04-29)

- updated LICENSE, minor code updates, security updates

1.1.1 (2020-04-21)

- Added test to check compressed dump is uniform
- Fixed sorting of class id in compressed dump to make JSON output uniform
- Added Python 3.8 support

1.1.0 (2020-01-10)

- Added pre-commit hook to rebuild bundled transliterators with bump2version
- remove to_dict from DirectedGraph, since it is handled through Marshmallow schemas.
- Adjust documentation to mention compression.
- added list-bundled CLI command
- added --regex/-re flag to graphtransliterator make-json CLI command to allow regular expressions
- removed coverage keyword from GraphTransliterator
- reorganized core.py
- converted from_dict, from_easyreading_dict, from_yaml, and from_yaml_file to static methods from class methods
- moved ambiguity-checking functions to ambiguity.py and tests to test_ambiguity.py

- set three levels of compression: 0 (Human-readable), 1 (no data loss, includes graph), 2 (no data loss, and no graph); 2 is fastest and set to default.
- set check_ambiguity to read keyword during JSON load
- allowed empty string productions during JSON compression
- added compression.py with decompress_config() and compress_config() to compress JSON
- added tests/test_compression.py to test compression.py
- added sorting of edge_list to DirectedGraph to allow dumped JSON comparison in tests
- adjusted _tokenizer_string_from() to sort by length then string for JSON comparison

1.0.7 (2019-12-22)

- added IncorrectVersionException, if serialized version being loaded is from a later version than the current graph-transliterator version
- added automatic edge_list creation if edge parameter in DirectedGraph
- added fields to and started using NodeDataSchema
- added pre_dump to GraphTransliteratorSchema, NodeDataSchema to remove empty values to compress Serialization
- removed rule from graph leaves and updated docs accordingly

1.0.6 (2019-12-15)

- fixed serialization of graph node indexes as integer rather than strings

1.0.5 (2019-12-14)

- added JOSS citation to README
- added --version to cli
- removed some asserts
- removed rule dictionaries from graph leaves to compress and simplify serialization

1.0.4 (2019-11-30)

- updates to docs

1.0.3 (2019-11-30)

- update to paper

1.0.2 (2019-11-30)

- updates for Zenodo

1.0.1 (2019-11-29)

- updated requirements_dev.txt

1.0.0 (2019-11-26)

- removed extraneous files
- updated development status in setup.py
- set to current jupyter-sphinx

0.4.10 (2019-11-04)

- fixed typo in requirements_dev.txt

0.4.9 (2019-11-04)

- quick fix to requirements_dev.txt due to readthedocs problem with not reading changes

0.4.8 (2019-11-04)

- twine update to 2.0

0.4.7 (2019-11-04)

- temp switch back to dev version of jupyter-sphinx for overflow error
- Dropped Python 3.5 support for twine 2.0 update

0.4.6 (2019-11-04)

- switched to latest jupyter-sphinx
- travis adjustments

0.4.5 (2019-10-31)

- Adjusted make-json CLI test to restore original example.json

0.4.4 (2019-10-24)

- moved README.rst to include in index.rst
- fixed error in advanced_tutorial.rst

0.4.3 (2019-10-24)

- fixed requirements_dev.txt

0.4.2 (2019-10-24)

- fixed README.rst for PyPI

0.4.1 (2019-10-24)

- fixed links to code in docs
- fixed link to NOTICE
- added acknowledgements

0.4.0 (2019-10-24)

- added bundled transliterators to api.rst
- adjustments to usage.rst
- adjustments to tutorial.rst
- fixes to docs (linking module)
- adjustments to advanced_tutorial.rst
- adjustments to README.rst
- fixes to AUTHORS.rst
- added kudos.rst to docs to acknowledge inspirational projects
- added advanced tutorial on bundling a transliterator.
- added cli.rst to docs
- fixed regex in get_unicode_char to allow hyphen
- added cli.py and adjusted setup.py
- updated tutorial
- added statement of need to README. Thanks [@rlskoesser](#).
- Removed continue-linenos jupyter-sphinx directive in favor of configuration settings
- added preface to documentation source files with links to production version, etc. Thanks [@rlskoesser](#).

- added custom css for jupyter-sphinx cells
- added jupyter-sphinx documentation with line numbering
- removed pkg_resources as source for version due to problem with loading from pythonpath for jupyter-sphinx in readthedocs, instead used __version__
- adjust path in docs/conf.py to fix docs error
- added bundled/schemas.py with MetadataSchema for bundled transliterator metadata
- added coverage to from_dict()
- added allow_none in onmatch_rules in GraphTransliteratorSchema
- adjusted core.py so that all edges are visited during search, even if no constraints
- removed _count_of_tokens() in favor of cost
- added IncompleteGraphCoverageException to exceptions.py
- added VisitLoggingDirectedGraph to graphs.py
- added tests/test_transliterator.py
- partially updated transliterators/README.rst
- removed transliterators/sample/*
- added yaml and json to package_data in setup.py
- Added to core.py class CoverageTransliterator, which tracks visits to edges, nodes, and onmatch rules, and allows clearing of visits and checking of coverage, used to make sure tests are comprehensive
- created test/test_coverage.py to test CoverageTransliterator
- created transliterators/bundled.py with class Bundled for bundled transliterators
- added load_from_YAML() and load_from_JSON() initializers to Bundled to load from bundled YAML (for development) and JSON (for speed)
- added load_yaml_tests(), run_yaml_tests(), and run_tests() to Bundled
- created transliterators/__init__.py that finds bundled transliterators in subdirectory and adds them to graphtransliterator.transliterator namespace
- added iter_names() and iter_transl iterators() to transliterators/__init__.py
- created test/test_transliterator.py to check bundled transliterator loading and functions
- created in transliterators/example/ __init__.py, example.json, example.yaml
- created in transliterators/example/tests test_example.py and example_tests.yaml

0.3.8 (2019-09-18)

- fixed load() docstring example
- updated check_ambiguity() to use cost

0.3.7 (2019-09-17)

- Adjusted docs to show readme as first page
- Added sample graph and code to README.rst
- moved images in docs to _static

0.3.6 (2019-09-17)

- adjusted installation.rst renaming libraries to modules
- updated paper and bibliography.

0.3.5 (2019-09-15)

- flake8 fix for core.py
- fixed bug in schemas.py whereby, during load(), DirectedGraphSchema() was modifying input settings
- added tests for modifications to settings by load()
- adjusted DirectedGraphSchema to allow for compacted transliteration rule settings
- adjusted GraphTransliteratorSchema to allow for compacted settings
- added tests to confirm all optional fields passed to load() are really optional
- added ValidationError if onmatch_rules_lookup present without onmatch_rules
- adjusted DirectedGraphSchema edge definition to remove str if loading from JSON
- added more rigorous schema definitions for edge_list and node in DirectedGraphSchema
- fixed flake8 warning in graphs.py
- adjusted docstrings in core.py for dump(), dumps(), load(), and loads()

0.3.4 (2019-09-15)

- added sphinx-issues and settings to requirements_dev.txt, docs/conf.py
- added .readthedocs.yml configuration file to accommodate sphinx-issues
- removed history from setup.py due to sphinx-issues
- fixed GraphTransliteratorException import in __init__.py
- added docs/_static directory
- fixed emphasis error and duplicate object description in docs/usages.rst
- fixed docstring in core.py
- added python versions badge to README.rst ([openjournals/joss-reviews#1717](https://openjournals.org/joss-reviews#1717)). Thanks @vc1492a.
- added NOTICE listing licenses of open-source text and code
- added Dependencies information to docs/install.rst ([openjournals/joss-reviews#1717](https://openjournals.org/joss-reviews#1717)). Thanks @vc1492a.
- updated AUTHORS.rst
- minor updates to README.rst

0.3.3 (2019-09-14)

- fixed missing marshmallow dependency (#47). Thanks @vc1492a.
- removed unused code from test (#47). Thanks @vc1492a.
- removed cerberus dependency

0.3.2 (2019-08-30)

- fixed error in README.rst

0.3.1 (2019-08-29)

- adjustments to README.rst
- cleanup in initialize.py and core.py
- fix to docs/api.rst
- adjusted setup.cfg for bumpversion of core.py
- adjusted requirements.txt
- removed note about namedtuple in dump docs
- adjusted docs (api.rst, etc.)

0.3.0 (2019-08-23)

- Removed `_tokens_of()` from `init`
- Removed `serialize()`
- Added `load()` to `GraphTransliterator`, without ambiguity checking
- Added `dump()` and `dumps()` to `GraphTransliterator` to export configuration
- renamed `_tokenizer_from()` to `_tokenizer_pattern_from()`, and so that `regex` is compiled on load and passed as pattern string (`tokenizer_pattern`)
- added settings parameters to `DirectedGraph`
- added `OnMatchRule` as `namedtuple` for consistency
- added new `GraphTransliterator.from_dict()`, which validates `from_yaml()`
- renamed `GraphTransliterator.from_dict()` to `GraphTransliterator.from_easyreading_dict()`
- added `schemas.py`
- removed `validate.py`
- removed `cerberus` and added `marshmallow` to `validate.py`
- adjusted tests
- Removed `check_settings` parameter

0.2.14 (2019-08-15)

- minor code cleanup
- removed yaml from validate.py

0.2.13 (2019-08-03)

- changed setup.cfg for double quotes in bumpversion due to Black formatting of setup.py
- added version test

0.2.12 (2019-08-03)

- fixed version error in setup.py

0.2.11 (2019-08-03)

- travis issue

0.2.10 (2019-08-03)

- fixed test for version not working on travis

0.2.9 (2019-08-03)

- Used Black code formatter
- Adjusted tox.ini, contributing.rst
- Set development status to Beta in setup.py
- Added black badge to README.rst
- Fixed comments and minor changes in initialize.py

0.2.8 (2019-07-30)

- Fixed ambiguity check if no rules present
- Updates to README.rst

0.2.7 (2019-07-28)

- Modified docs/conf.py
- Modified equation in docs/usage.rst and paper/paper.md to fix doc build

0.2.6 (2019-07-28)

- Fixes to README.rst, usage.rst, paper.md, and tutorial.rst
- Modifications to core.py documentation

0.2.5 (2019-07-24)

- Fixes to HISTORY.rst and README.rst
- 100% test coverage.
- Added draft of paper.
- Added graphtransliterator_version to serialize().

0.2.4 (2019-07-23)

- minor changes to readme

0.2.3 (2019-07-23)

- added xenial to travis.yml

0.2.2 (2019-07-23)

- added CI

0.2.1 (2019-07-23)

- fixed HISTORY.rst for PyPI

0.2.0 (2019-07-23)

- Fixed module naming in docs using `__module__`.
- Converted DirectedGraph nodes to a list.
- Added Code of Conduct.
- Added GraphTransliterator class.
- Updated module dependencies.
- Added requirements.txt
- Added `check_settings` parameter to skip validating settings.
- Added tests for ambiguity and `check_ambiguity` parameter.
- Changed name to Graph Transliterator in docs.
- Created core.py, validate.py, process.py, rules.py, initialize.py, exceptions.py, graphs.py
- Added `ignore_errors` property and setter for transliteration exceptions (UnrecognizableInputToken, NoMatchingTransliterationRule)

- Added logging to graphtransliterator
- Added positive cost function based on number of matched tokens in rule
- added metadata field
- added documentation

0.1.1 (2019-05-30)

- Adjusted copyright in docs.
- Removed Python 2 support.

0.1.0 (2019-05-30)

- First release on PyPI.

3.3 Indices and tables

- [genindex](#)
- [modindex](#)
- [search](#)

PYTHON MODULE INDEX

g

`graphtransliterator`, [52](#)

`graphtransliterator.transliterators`, [67](#)

A

`add_edge()` (*graphtransliterator.DirectedGraph* method), 69
`add_node()` (*graphtransliterator.DirectedGraph* method), 70
AmbiguousTransliterationRulesException, 73

B

Bundled (class in *graphtransliterator.transliterators*), 67

C

`check_coverage()` (*graphtransliterator.CoverageTransliterator* method), 67
`check_coverage()` (*graphtransliterator.VisitLoggingDirectedGraph* method), 71
`check_onmatchrules_coverage()` (*graphtransliterator.CoverageTransliterator* method), 67
`clear_visited()` (*graphtransliterator.CoverageTransliterator* method), 67
`clear_visited()` (*graphtransliterator.VisitLoggingDirectedGraph* method), 71
`consolidate` (*graphtransliterator.WhitespaceRules* attribute), 73
`cost` (*graphtransliterator.TransliterationRule* attribute), 72
CoverageTransliterator (class in *graphtransliterator*), 67

D

`default` (*graphtransliterator.WhitespaceRules* attribute), 72
DirectedGraph (class in *graphtransliterator*), 69
DirectedGraphSchema (class in *graphtransliterator*), 73
`directory` (*graphtransliterator.transliterators.Bundled* property), 67
`dump()` (*graphtransliterator.GraphTransliterator* method), 54
`dumps()` (*graphtransliterator.GraphTransliterator* method), 56

E

EasyReadingSettingsSchema (class in *graphtransliterator*), 73
`edge` (*graphtransliterator.DirectedGraph* attribute), 69
`edge_list` (*graphtransliterator.DirectedGraph* attribute), 69
Example (class in *graphtransliterator.transliterators*), 68

F

`from_dict()` (*graphtransliterator.GraphTransliterator* static method), 57
`from_easyreading_dict()` (*graphtransliterator.GraphTransliterator* static method), 57
`from_JSON()` (*graphtransliterator.transliterators.Bundled* method), 67
`from_yaml()` (*graphtransliterator.GraphTransliterator* static method), 58
`from_YAML()` (*graphtransliterator.transliterators.Bundled* method), 68
`from_yaml_file()` (*graphtransliterator.GraphTransliterator* static method), 59

G

`generate_yaml_tests()` (*graphtransliterator.transliterators.Bundled* method), 68
`graph` (*graphtransliterator.GraphTransliterator* property), 59
graphtransliterator module, 52
GraphTransliterator (class in *graphtransliterator*), 52
graphtransliterator.transliterators module, 67
`graphtransliterator_version` (*graphtransliterator.GraphTransliterator* property), 59
GraphTransliteratorException, 73
GraphTransliteratorSchema (class in *graphtransliterator*), 74

I

`ignore_errors` (*graphtransliterator.GraphTransliterator* property), 60

`iter_names()` (in module `graphtranslitera-
tor.transliterators`), 69
`iter_transliterators()` (in module `graphtranslitera-
tor.transliterators`), 69
`ITRANSDevanagariToUnicode` (class in `graph-
transliterator.transliterators`), 68

L

`last_input_tokens` (`graphtranslitera-
tor.GraphTransliterator` property), 60
`last_matched_rule_tokens` (`graphtranslitera-
tor.GraphTransliterator` property), 60
`last_matched_rules` (`graphtranslitera-
tor.GraphTransliterator` property), 60
`load()` (`graphtransliterator.GraphTransliterator` static
method), 60
`load_yaml_tests()` (`graphtranslitera-
tor.transliterators.Bundled` method), 68
`loads()` (`graphtransliterator.GraphTransliterator` static
method), 62

M

`match_at()` (`graphtransliterator.GraphTransliterator`
method), 63
`metadata` (`graphtransliterator.GraphTransliterator` prop-
erty), 64
`MetadataSchema` (class in `graphtranslitera-
tor.transliterators`), 68
`module`
 `graphtransliterator`, 52
 `graphtranslitera-
tor.transliterators`, 67

N

`name` (`graphtransliterator.transliterators.Bundled` prop-
erty), 68
`new()` (`graphtransliterator.transliterators.Bundled` class
method), 68
`next_classes` (`graphtransliterator.OnMatchRule` at-
tribute), 72
`next_classes` (`graphtransliterator.TransliterationRule`
attribute), 72
`next_tokens` (`graphtransliterator.TransliterationRule`
attribute), 72
`node` (`graphtransliterator.DirectedGraph` attribute), 69
`NoMatchingTransliterationRuleException`,
73

O

`onmatch_rules` (`graphtranslitera-
tor.GraphTransliterator` property), 64
`onmatch_rules_lookup` (`graphtranslitera-
tor.GraphTransliterator` property), 64
`OnMatchRule` (class in `graphtransliterator`), 72

`OnMatchRuleSchema` (class in `graphtransliterator`), 74

P

`prev_classes` (`graphtransliterator.OnMatchRule` at-
tribute), 72
`prev_classes` (`graphtransliterator.TransliterationRule`
attribute), 71
`prev_tokens` (`graphtransliterator.TransliterationRule`
attribute), 71
`production` (`graphtransliterator.OnMatchRule` at-
tribute), 72
`production` (`graphtransliterator.TransliterationRule` at-
tribute), 71
`productions` (`graphtransliterator.GraphTransliterator`
property), 64
`pruned_of()` (`graphtransliterator.GraphTransliterator`
method), 64

R

`rules` (`graphtransliterator.GraphTransliterator` property),
65
`run_tests()` (`graphtranslitera-
tor.transliterators.Bundled` method), 68
`run_yaml_tests()` (`graphtranslitera-
tor.transliterators.Bundled` method), 68

S

`SettingsSchema` (class in `graphtransliterator`), 74

T

`token_class` (`graphtransliterator.WhitespaceRules` at-
tribute), 73
`tokenize()` (`graphtransliterator.GraphTransliterator`
method), 65
`tokenizer_pattern` (`graphtranslitera-
tor.GraphTransliterator` property), 66
`tokens` (`graphtransliterator.GraphTransliterator` prop-
erty), 66
`tokens` (`graphtransliterator.TransliterationRule` attribute),
72
`tokens_by_class` (`graphtranslitera-
tor.GraphTransliterator` property), 66
`transliterate()` (`graphtranslitera-
tor.GraphTransliterator` method), 66
`TransliterationRule` (class in `graphtransliterator`),
71
`TransliterationRuleSchema` (class in `graph-
transliterator`), 74

U

`UnrecognizableInputTokenException`, 73

V

VisitLoggingDirectedGraph (*class in graph-transliterator*), [71](#)

W

whitespace (*graphtransliterator.GraphTransliterator property*), [67](#)

WhitespaceDictSettingsSchema (*class in graph-transliterator*), [74](#)

WhitespaceRules (*class in graphtransliterator*), [72](#)

WhitespaceSettingsSchema (*class in graph-transliterator*), [75](#)

Y

yaml_tests_filen (*graphtranslitera-tor.transliterators.Bundled property*), [68](#)