
Emit Documentation

Release 0.2.0

Brian Hicks

January 17, 2013

CONTENTS

Emit is a library that hooks into Celery (or just local memory) to provide subscriptions and notifications for your functions. It makes stream processing a whole lot easier.

Supported Pythons:

- CPython 2.6
- CPython 2.7
- CPython 3.2
- PyPy 1.9

Contents:

GETTING STARTED

After reading this document you'll have a good idea of what you need to do to get your first emit graph running, both in local memory and Celery.

1.1 Installing

Eventually, you'll be able to install via *pip*. However, while the library is under initial development you'll need to install via git:

```
pip install emit
```

1.2 Quickstart

For a sampler, we're going to make a simple command-line application that will take and count all the words in a document, giving you the top 5.

Put the following into *graph.py*

```
1 from __future__ import print_function
2 from collections import Counter
3 from emit import Router
4 import sys
5
6 router = Router()
7
8
9 def prefix(name):
10     return '%s.%s' % (__name__, name)
11
12
13 @router.node(('word',), entry_point=True)
14 def words(msg):
15     print('got document')
16     for word in msg.document.strip().split(' '):
17         yield word
18
19
20 WORDS = Counter()
21
22
```

```
23 @router.node(('word', 'count'), prefix('words'))
24 def count_word(msg):
25     print('got word (%s)' % msg.word)
26
27     global WORDS
28     WORDS.update([msg.word])
29
30     return msg.word, WORDS[msg.word]
31
32 if __name__ == '__main__':
33     router(document=sys.stdin.read())
34
35     print()
36     print('Top 5 words:')
37     for word, count in WORDS.most_common(5):
38         print('      %s: %s' % (word, count))
```

(incidentally, this file is available in the project directory as `examples/simple/graph.py`.)

Now on the command line: `echo "the rain in spain falls mainly on the plain" | python graph.py`. You should get some output that looks similar to the following

```
got document
got word (the)
got word (rain)
got word (in)
got word (spain)
got word (falls)
got word (mainly)
got word (on)
got word (the)
got word (plain)
```

```
Top 5 words:
the: 2
on: 1
plain: 1
mainly: 1
rain: 1
```

1.2.1 Breaking it Down

First, we need to construct a router:

```
router = Router()
```

Since we're keeping everything in-memory, we don't need to specify anything to get this to work properly. It should "Just Work(TM)".

Next, we define a function to split apart a document on spaces to get words:

```
@router.node(('word',), entry_point=True)
def words(msg):
    print('got document')
    for word in msg.document.strip().split(' '):
        yield word
```

`Router` provides a decorator. The first argument is the fields that the decorated function returns. These are wrapped in a message and passed around between functions.

We don't specify any subscriptions on this function, since it really doesn't need any. In fact, it's an entry point, so we specify that instead. This specifically means that if you call the router directly it will delegate to this function. There can be multiple functions with `entry_point` set to true on a given Router.

If the decorated function is a generator, each yielded value is treated as a separate input into the next nodes in the graph.

Splitting the document into parts is only as useful as what we can do with the words, so let's count them now:

```
WORDS = Counter()
@router.node('word', 'count', prefix('words'))
def count_word(msg):
    print('got word (%s)' % msg.word)

    global WORDS
    WORDS.update([msg.word])

    return msg.word, WORDS[msg.word]
```

There's a little less going on in this function. We just update a `Counter` builtin, and then return the word and the count to be passed down the graph. In real life, you'd probably persist this value in a database to allow multiple workers to process different parts of the stream.

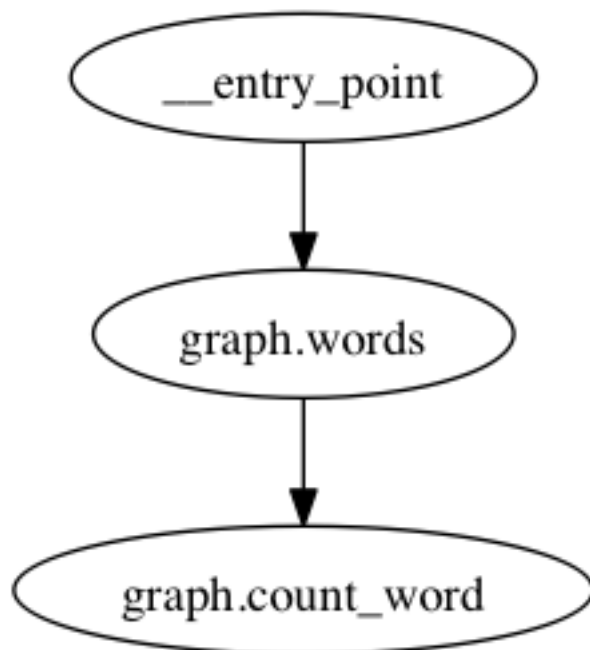
In non-entry nodes, the second argument of `router.node` is a string or list of functions to subscribe to. These need to be fully qualified when you're using Celery, but for now they're fine.

Now that we've defined both functions, it's time to send some data into our graph:

```
router(document=sys.stdin.read())
```

Calling this graph is easy, since we defined a function as an entry point. You can call any of the functions (or the router itself) by using keyword arguments or passing a dictionary.

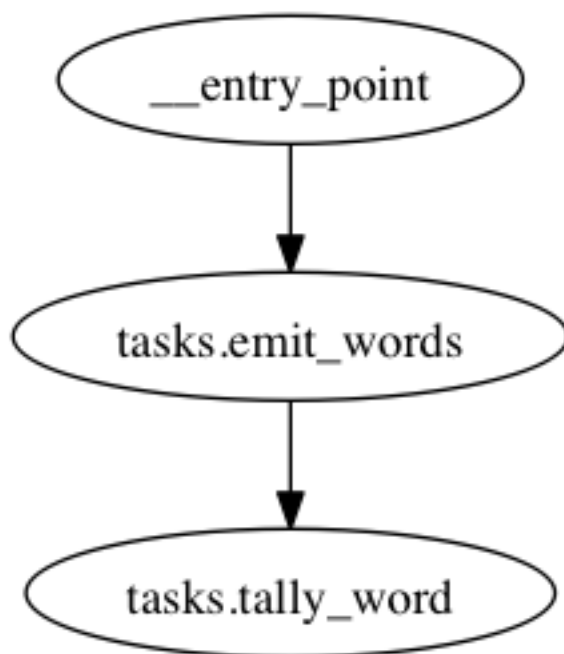
In the end, data flows through the graph like this:



USING CELERY TO DISTRIBUTE PROCESSING

Emit makes it simple to use celery to distribute realtime processing across many worker nodes. To demonstrate this, we'll be scaling our quickstart example

We'll be making, in essence, this graph:



2.1 Setting up Celery

Create an `app.py` file for your celery initialization code to live in:

```
1 'simple celery app'
2 from celery import Celery
3 from emit import CeleryRouter
4
5 import logging
6
```

```
7 app = Celery(
8     'celery_emit_example',
9     broker='redis://'
10 )
11 app.conf.update(
12     CELERY_IMPORTS=('tasks',)
13 )
14
15 router = CeleryRouter(celery_task=app.task, node_modules=['tasks'])
16
17 logging.basicConfig(format='%(levelname)s: %(message)s', level=logging.DEBUG)
```

Take note that Router is initialized using the default celery task in this case. This is probably the best way to do it, since per-task settings should belong in the task (possible in Emit's decorator), and app-level configuration should be on the app object (as on line 10).

Next we'll define (in `tasks.py`) a function to take a document and emit each word:

```
@router.node(('word',), entry_point=True)
def emit_words(msg):
    for word in msg.document.strip().split(' '):
        yield word
```

We don't have to give any special syntax to get these tasks to work with celery: since we specified it in the router, they just do.

However, if you want to give special celery attributes to a particular function, you can do that too:

```
@router.node(('word', 'count'), subscribe_to='tasks.emit_words', celery_task=app.task(rate_limit='5/s'))
def tally_word(msg):
    redis = Redis()
    return msg.word, redis.zincrby('celery_emit_example', msg.word, 1)
```

Obviously rate limiting to 5 per second in this case is a bit contrived, but you get the general idea: it's easy to configure tasks within the decorator by passing in the celery decorator.

2.2 Running the Graph

We'll need to boot up the celery daemon:

```
celery worker -A app.app -l INFO -E
```

And enter the following on the command line to start something fun processing (if you'd like, the relevant code is in `examples/celery/kickoff.py` in the project directory, start it and get a prompt with `ipython -i kickoff.py`):

```
from app import router
import random
words = 'the rain in spain falls mainly on the plain'.split(' ')
router(document=' '.join(random.choice(words) for i in range(50)))
```

You should get something like the following:

```
{'word': 'the'},
{'word': 'spain'},
{'word': 'in'},
# ...
{'word': 'falls'},
```

```
{'word': 'falls'},  
{'word': 'mainly'})
```

And you should see the celery window quickly scrolling by with updated totals. Run the command a couple more times, if you like, and you'll see the totals keep going up.

USING EMIT IN OTHER LANGUAGES

You can use Emit in other languages through the multilang API.

3.1 Defining Tasks

We're going to define a node that takes a number and emits each integer in that range. Let's do it with Ruby! (why not?)

```
require "json"
message = JSON.parse(STDIN.read)

message["count"].times do |i|
  puts i.to_json
end
```

(the equivalent in Python is in `examples/multilang/test.py`)

The messages passed in and out are expected to be in JSON format. Output from the functions should be json strings separated by newlines.

3.2 Creating a Node

We'll be subclassing `emit.multilang.ShellNode` to tell emit how to execute our task:

```
@router.node(('n',))
```

After that, you can call your node and subscribe as normal.

REGEX ROUTING

You can subscribe to arbitrary node’s output by providing a regular expression. In this example, we’ll use Redis’ pubsub capabilities to notify an external receiver of all tasks passing through the graph.

The product of this example is in `examples/regex/graph.py`.

First, we’ll create a function that yields multiple return values. In this case, we’re going to naively parse a HTTP querystring.

```
@router.node(('key', 'value'), entry_point=True)
def parse_querystring(msg):
    'parse a querystring into keys and values'
    for part in msg.querystring.strip().lstrip('?').split('&'):
        key, value = part.split('=')
        yield key, value
```

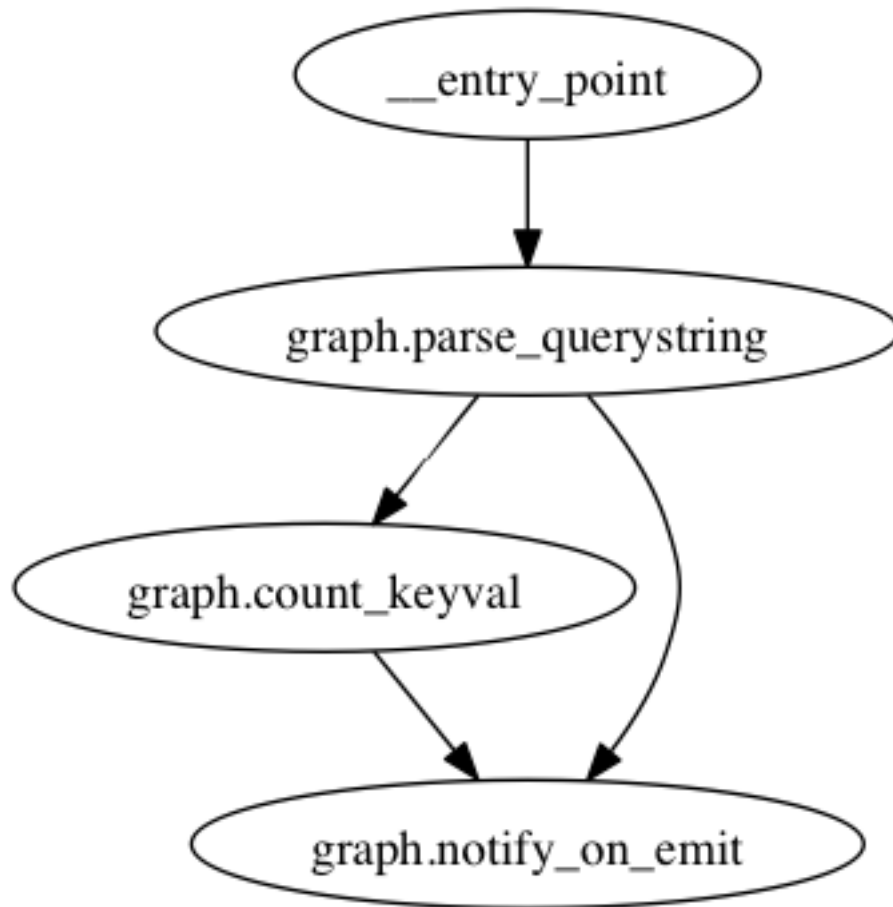
Now we’re going to count keys and values:

```
@router.node(('key', 'value', 'count'), prefix('parse_querystring'))
def count_keyval(msg):
    count = redis.zincrby('querystring_count.%s' % msg.key, msg.value, 1)
    return msg.key, msg.value, count
```

Next, we’ll make a function that publishes to Redis on every message:

```
@router.node(tuple(), '.*')
def notify_on_emit(msg):
    redis.publish('notify_on_emit.%s' % msg._origin, msg.as_json())
    return NoResult
```

Now, when you call `router(querystring='?a=1&b=2&c=3')`, `notify_on_emit` will publish seven messages: three with origin “`graph.parse_querystring`”, three with origin “`graph.count_keyval`”, and one with origin “`__entry_point`”. The graph ends up looking like this:

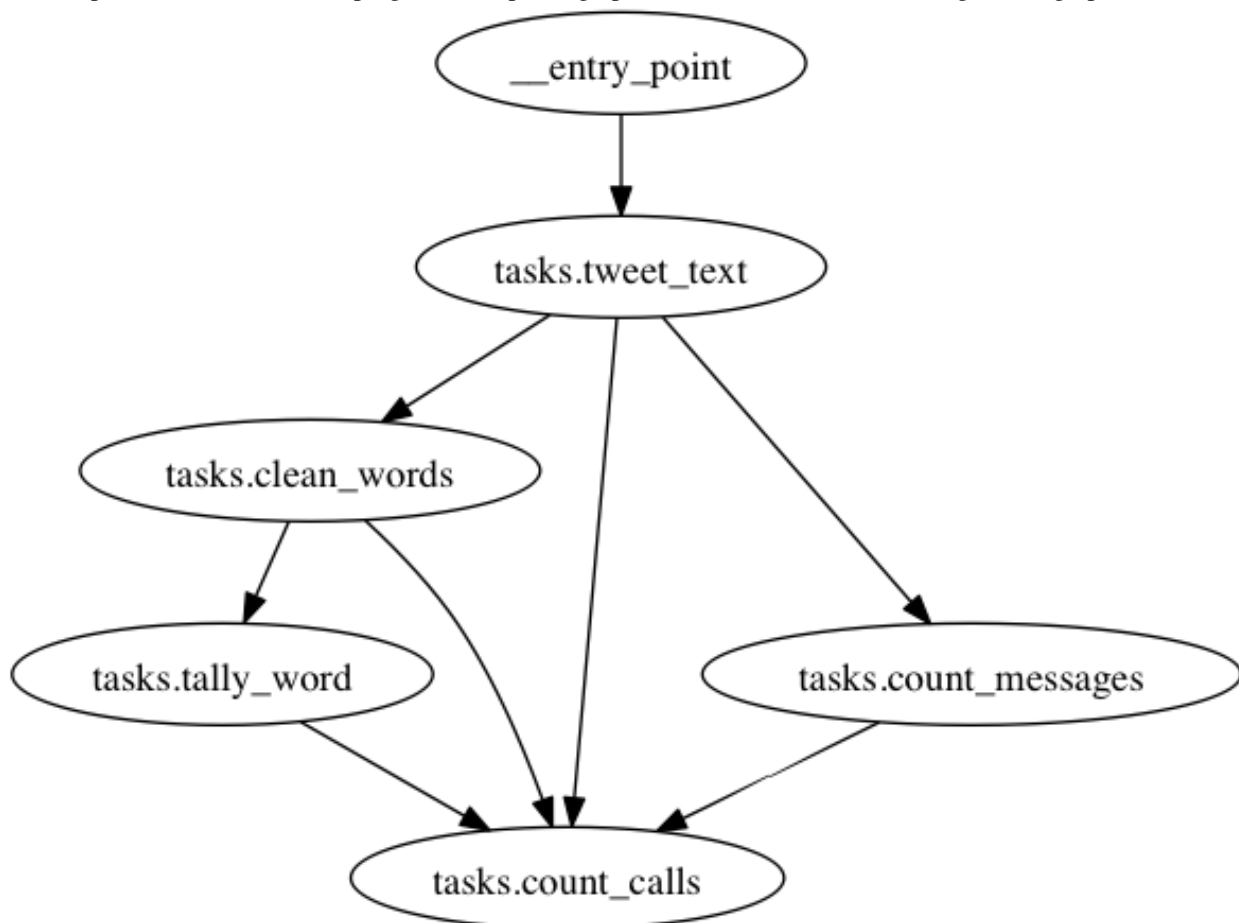


You can also specify `ignores` in `Router.node`, which can cut a little fat out of an otherwise greedy regex.

COMMAND LINE UTILITIES

5.1 emit_digraph - Generate Graph Images

Emit ships with a command-line program to inspect a graph: `emit_digraph`. Use it to generate graphs like this:



`emit_digraph` will output the code `graphviz` needs to properly generate the graph. (You'll need `graphviz` installed on your machine for this to render properly.) To use it, pass it the path of your router. (for example, `emit_digraph app.router` in the Celery example.) The output should look something like this:

```
digraph router {
"tasks.clean_words" -> "tasks.tally_words";
```

```
"tasks.clean_text" -> "tasks.clean_words";  
"tasks.tweet_text" -> "tasks.count_messages";  
"__entry_point" -> "tasks.tweet_text";  
}
```

to make graphviz generate a PNG of the graph, pipe it into the following command:

```
emit_digraph app.router | dot -T png -o graph.png
```

LOGGING

Emit is set up to handle logging using Python's standard logger. It currently uses the following levels:

- `DEBUG`: task registration and calls - very verbose
- `INFO`: route registration, receipts

So far there's been no need for anything above `INFO`, but that may change in the future.

6.1 Setting Up Logging

In some file (I recommend the file where the router is initialized, but your project may vary) insert the following lines:

```
import logging
logging.basicConfig(format='%(levelname)s: %(message)s', level=logging.DEBUG) # or INFO etc.
```

6.2 Setting Up Logging in Django

Set up a logger for `emit`. Examples to come.

API DOCUMENTATION

7.1 Router

class `emit.router.Router` (*message_class=None, node_modules=None, node_package=None*)

A router object. Holds routes and references to functions for dispatch

__init__ (*message_class=None, node_modules=None, node_package=None*)

Create a new router object. All parameters are optional.

Parameters

- **message_class** (`emit.message.Message` or subclass) – wrapper class for messages passed to nodes
- **node_modules** (a list of `str`, or `None`.) – a list of modules that contain nodes
- **node_package** (`str`, or `None`.) – if any `node_modules` are relative, the path to base off of.

Exceptions `None`

Returns `None`

__call__ (***kwargs*)

Route a message to all nodes marked as entry points.

Note: This function does not optionally accept a single argument (dictionary) as other points in this API do - it must be expanded to keyword arguments in this case.

node (*fields, subscribe_to=None, entry_point=False, ignore=None, **wrapper_options*)

Decorate a function to make it a node.

Parameters

- **fields** (ordered iterable of `str`) – fields that this function returns
- **subscribe_to** (`str` or iterable of `str`) – functions in the graph to subscribe to. These indicators can be regular expressions.
- **ignore** (`str` or iterable of `str`) – functions in the graph to ignore (also uses regular expressions.) Useful for ignoring specific functions in a broad regex.
- **entry_point** (`bool`) – Set to `True` to mark this as an entry point - that is, this function will be called when the router is called directly.

In addition to all of the above, you can define a `wrap_node` function on a subclass of Router, which will need to receive node and an options dictionary. Any extra options passed to node will be passed down to the options dictionary. See `emit.router.CeleryRouter.wrap_node` as an example.

Returns decorated and wrapped function, or decorator if called directly

Examples

Multiple fields:

```
@router.node(['quotient', 'remainder'])
def division_with_remainder(msg):
    return msg.numer / msg.denom, msg.numer % msg.denom
```

This function would end up returning a dictionary that looked something like:

```
{'quotient': 2, 'remainder': 1}
```

The next node in the graph would receive a `emit.message.Message` with “quotient” and “remainder” fields.

Emitting multiple values:

```
@router.node(['word'])
def parse_document(msg):
    for word in msg.document.clean().split(' '):
        yield word
```

If the function returns a generator, Emit will gather the values together and make sure the generator exits cleanly before returning (but this may change in the future via a flag.) Therefore, the return value will look like this:

```
({'word': "I've"},
 {'word': 'got'},
 {'word': 'a'},
 {'word': 'lovely'},
 {'word': 'bunch'},
 {'word': 'of'},
 {'word': 'coconuts'})
```

Each message in the tuple will be passed on individually in the graph.

dispatch (*origin, destination, message*)
dispatch a message to a named function

Parameters

- **destination** (`str`) – destination to dispatch to
- **message** (`emit.message.Message` or subclass) – message to dispatch

get_message_from_call (**args, **kwargs*)
Get message object from a call.

Raises `TypeError` (if the format is not what we expect)

This is where arguments to nodes are turned into Messages. Arguments are parsed in the following order:

- A single positional argument (a dict)
- No positional arguments and a number of keyword arguments

get_name (*func*)
Get the name to reference a function by

Parameters **func** (*callable*) – function to get the name of

regenerate_routes ()

regenerate the routes after a new route is added

register (*name, func, fields, subscribe_to, entry_point, ignore*)

Register a named function in the graph

Parameters

- **name** (*str*) – name to register
- **func** (*callable*) – function to remember and call

fields, subscribe_to and *entry_point* are the same as in `Router.node()`.

register_ignore (*origins, destination*)

Add routes to the ignore dictionary

Parameters

- **origins** (*str* or iterable of *str*) – a number of origins to register
- **destination** (*str*) – where the origins should point to

Ignore dictionary takes the following form:

```
{ 'node_a': set(['node_b', 'node_c']),
  'node_b': set(['node_d']) }
```

register_route (*origins, destination*)

Add routes to the routing dictionary

Parameters

- **origins** (*str* or iterable of *str* or *None*) – a number of origins to register
- **destination** (*str*) – where the origins should point to

Routing dictionary takes the following form:

```
{ 'node_a': set(['node_b', 'node_c']),
  'node_b': set(['node_d']) }
```

route (*origin, message*)

Using the routing dictionary, dispatch a message to all subscribers

Parameters

- **origin** (*str*) – name of the origin node
- **message** (`emit.message.Message` or subclass) – message to dispatch

wrap_result (*name, result*)

Wrap a result from a function with it's stated fields

Parameters

- **name** (*str*) – fields to look up
- **result** (*anything*) – return value from function. Will be converted to tuple.

Returns *dict*

class `emit.router.CeleryRouter` (*celery_task, *args, **kwargs*)

Router specifically for Celery routing

__init__ (*celery_task*, *args, **kwargs)
Specifically route when celery is needed

Parameters **celery_task** (*A celery task decorator, in any form*) – celery task to apply to all nodes
(can be overridden in `Router.node()`.)

dispatch (*origin, destination, message*)
enqueue a message with Celery

Parameters

- **destination** (`str`) – destination to dispatch to
- **message** (`emit.message.Message` or subclass) – message to dispatch

wrap_node (*node, options*)
celery registers tasks by decorating them, and so do we, so the user can pass a celery task and we'll wrap our code with theirs in a nice package celery can execute.

7.2 Message

class `emit.message.Message` (*args, **kwargs)
Convenient wrapper around a dictionary to provide attribute access

as_dict ()
representation of this message as a dictionary

Returns dict

as_json ()
representation of this message as a json object

Returns str

7.3 Multilang

class `emit.multilang.ShellNode`
callable object to wrap communication to a node in another language

to use this, subclass `ShellNode`, providing “command”. Decorate it however you feel like.

Messages will be passed in on lines in msgpack format. This class expects similar output: msgpack messages separated by a newline.

__call__ (*msg*)
call the command specified, processing output

deserialize (*msg*)
deserialize output to a Python object

get_command ()
get the command as a list

get_cwd ()
get directory to change to before running the command

GLOSSARY

Some of the terms in this project may be a little rough. I'm still considering what things should really be called here, but here's a primer:

- *Router*: an object (implemented in `emit.router.Router`) that keeps references to functions and their names and handles dispatch. It knows where everything is.
- *Node*: a function or callable class that receives messages, processes them in it's own way, and passes them on down the graph.
- *Subscription*: an edge in the graph - only flows one way.
- *Graph*: a directed graph, like graph theory. A collection of nodes connected by subscriptions.

CHANGELOG

9.1 0.2.0

- New argument for `node: ignores`. Pass it some regex to ignore items in otherwise broad subscriptions.
- Add support for Python 2.6

9.2 0.1.0

- Initial Release to PyPI

INDICES AND TABLES

- *genindex*
- *modindex*
- *search*

LICENSE

Copyright (c) 2012–2013 Brian Hicks

Permission is hereby granted, free of charge, to any person obtaining a copy of this software and associated documentation files (the "Software"), to deal in the Software without restriction, including without limitation the rights to use, copy, modify, merge, publish, distribute, sublicense, and/or sell copies of the Software, and to permit persons to whom the Software is furnished to do so, subject to the following conditions:

The above copyright notice and this permission notice shall be included in all copies or substantial portions of the Software.

THE SOFTWARE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS IN THE SOFTWARE.

PYTHON MODULE INDEX

e

`emit.message, ??`
`emit.multilang, ??`
`emit.router, ??`