

Stream Processing with Python in Wallaroo

DataEngConf NYC, Oct. 2017

Andrew Turley, Lead Engineer andy@WallarooLabs.com

Wallaroo Labs

Wallroo Labs is a small startup that making it easy to develop distributed data processing applications

- **★**NYC-based
- **★**VC funded
- ★Working on Wallaroo for 18 months
- ★Extensive experience working on distributed data processing applications

The Challenge (in Tweets)

Finding the Trending Topics in 1M tweets is easy to prototype on a laptop (in between coffees)

...but months to productionize at large-scale



...and then the volume spikes

...and then you add sentiment analysis ...and then (aargh!)

Focus on Algorithms, not Infrastructure

Wallaroo is an open-source framework for building and operating fast data, big data, and machine learning applications

- ★ Deploy rapidly, scale automatically, and operate at very low-cost
- ★ Makes the infrastructure virtually disappear
- ★ Native for Python and GoLang, no JVM needed
- ★ Ultrafast and accurate processing



Do What You're Good At

Wallaroo lets you do what you're good at, and takes care of the rest

- **★**Wallaroo
 - **★**Scaling
 - **★**Message Delivery
 - **★**Error Recovery
- **★**Application Developer
 - **★**Business Logic

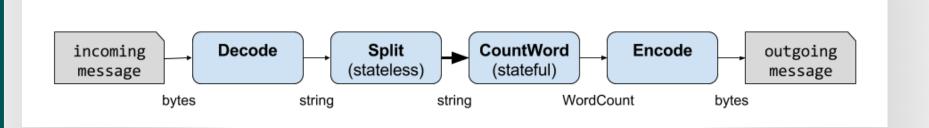
Another Challenge (in Words)

The canonical WORD COUNT example

- ★Receive messages that consist of one or more words
- ★Split the words into individual words
- ★Keep a running total of how many times each word occurs
- ★Output a message for each received word with the total occurrences of the word

Let's Count Some Words!

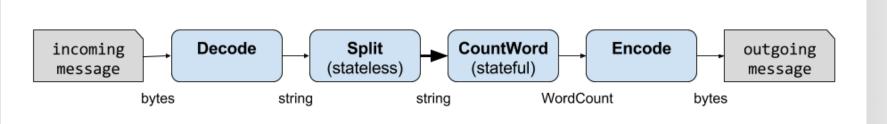
The canonical WORD COUNT example



- ★Get a message containing one or more words
- ★Keep a tally of how many times each word has been seen
- ★For each incoming word, output the current tally for that word

Coffee

Let's get some coffee





Word Count

Single machine prototype

```
class WordTotals:
    def init (self):
        self.counts = {}
                                                                  import socket
    def update(self, word):
                                                                  def main():
        self.counts[word] = self.counts.get(word, 0) + 1
                                                                       word_totals = WordTotals()
    def get_count(self, wo
                                                                                                     INET, socket.SOCK_STREAM)
        return WordCount(word
class WordCount(object):
                                                                      listener = socket.socket(socket.AF_INET, socket.SOCK_STREAM)
    def __init__(self, word, count):
                                                                      <u>listene</u>r.bind(('', 5000))
        self.word = word
                                   listener.listen(1)
(source, addr) = listener.accept()
while True:
        self.count = count
def decode(bs):
                                                                          raw words = decode(source.recv(512))
    return bs.encode("utf-8")
                                                                          words = split(raw words)
                                                                          for w in words:
def split(words):
                                                                              word count = count_words(w, word_totals)
    return words.split()
                                                                              sink.send(encode(word_count))
def count_words(word, word_totals):
                                                                  if __name__ == "__main ":
    word_totals.update(word)
                                                                      main()
    return word totals.get count(word)
def encode(word count):
    return "{} => {}\n".format(word_count.word, word_count.count)
```

Wallaroo Handles Scaling

Run at scale and autoscale

```
class Split(object):
                                                              class WordCount(object):
   def name(self):
                                                                   def init (self, word, count):
       return "split into words"
                                                                       self.word = word
                                                                       self.count = count
   def compute_multi(self, data):
       punctuation = " !\"#$%&'()*+,-./:;<=>?@[\]^_`{|}~"
                                                              class WordTotalsBuilder(object):
       words = []
                                                                   def build(self):
                                                                        return WordTotals()
       for line in data.spli
           clean line = line.lower().strip(punctuation)
           for word in clean line.split(
               clean word = word.strip(pun
               words.append(clean word)
                                                                       if data[0] >= 'a' or data[0] <= 'z':</pre>
                                                                          return data[0]
       return words
                                                                       else:
class CountWord(object)
   def name(self):
       return "Count
   def compute(self, word, word totals):
       word_totals.update(word)
       return (word totals.get count(word), True)
                                                                   def payload length(self, bs):
                                                                       return struct.unpack(">I", bs)[0]
class WordTotals(object):
   def __init__(self):
       self.word totals = {}
                                                                   def decode(self, bs):
                                                                       return bs.decode("utf-8")
   def update(self, word):
       if self.word_totals.has_key(word):
                                                              class Encoder(object):
           self.word_totals[word] = self.word_totals[word] + 1
                                                                   def encode(self, data):
                                                                       output = data.word + " => " + str(data.count) + "\n"
           self.word totals[word] = 1
                                                                       print output
   def get count(self, word):
                                                                       return output
       return WordCount(word, self.word_totals[word])
```

Plain old Python

Wallaroo lets you describe your business logic in plain Python

```
class Split(object):
    def name(self):
        return "split into words"
    def compute_multi(self, data):
        punctuation = " !\"#$%&'()*+,-./:;<=>?@[\]^ `{|}~"
       words = []
        for line in data.split("\n"):
            clean_line = line.lower().strip(punctuation)
            for word in clean_line.split(' '):
                clean_word = word.strip(punctuation)
                words.append(clean word)
        return words
```

Plain old Python

Wallaroo lets you operate on Python objects

```
class CountWord(object):
   def name(self):
        return "Count Word"
   def compute(self, word, word_totals):
       word_totals.update(word)
       return (word totals.get count(word), True)
class WordTotals(object):
   def __init__(self):
       self.word_totals = {}
   def update(self, word):
       if self.word totals.has key(word):
            self.word totals[word] = self.word totals[word] + 1
       else:
            self.word_totals[word] = 1
   def get count(self, word):
       return WordCount(word, self.word_totals[word])
```

Wallaroo Handles Scaling

Wallaroo handles scaling for you automatically

Python!

Leverage your existing knowledge, use your favorite libraries





NumPy









NLTK

Wallaroo

Wallaroo is a framework for building and operating fast data, big data, and machine learning applications that deploy rapidly, scale automatically and operate at very low-cost



wallaroolabs.com

github.com/wallaroolabs/wallaroo