

# Boosting command line experience

Python meets AWK

Kirill Pavlov

Technical Recruiter, Terminal 1

October 22, 2017

***CODECONF 2017***  
***22 Oct @ ICC, Kowloon***

- ▶ A lot of Python/AWK examples here.
- ▶ Source code and slides [available online](#).
- ▶ At the end: build a stock trading system and check NYSE:CS.

# Table of content

1. Problem Background
2. AWK Bootcamp in 5 min
3. Tabtools architecture and features
4. Stock example with NYSE:CS

# Table of content

1. Problem Background
2. AWK Bootcamp in 5 min
3. Tabtools architecture and features
4. Stock example with NYSE:CS

- ▶ [Yandex](#), year 2010. Hadoop was not widely adopted.
- ▶ 10Gb of archived ads data daily: *time*, *ad\_id*, *site\_id*, *clicks*.
- ▶ Task: daily data aggregation (simple functions: group by, sum, join) and feature generation for further machine learning classification.
- ▶ Solution: released a set of command line scripts.

# Example

This presentation uses UCI machine learning [Higgs boson](#) data: 11M objects, 28 attributes, 7.5Gb unarchived.

## Questions:

1. What is the maximum value of *lepton\_eta*?
2. What is the average value of *lepton\_phi* by class 0 and 1?
3. Filter objects with  $m_{jj} > 0.75$  (8.9M objects) and sort them by *m\_wbb*.

## Solutions:

1. In-memory Python with Pandas.
2. Database SQL queries (PostgreSQL and Docker).
3. Command line with AWK.

# Demo Time

# Reality Check

It's not as agile as it seems. You work inside the company network.

1. You **don't have sudo** rights and your admin does not want to install anything for you. Like no database or user privileges, etc.
2. The **server does not have GitHub/Internet access** and the only deployment possible is Java JARs or C/C++/etc. So, no NodeJS/Python packages. And of course no R/Matlab/Excel.
3. Get better at command line tools ;)



# Table of content

1. Problem Background
2. AWK Bootcamp in 5 min
3. Tabtools architecture and features
4. Stock example with NYSE:CS

# Basic concepts

1. **AWK**<sup>1</sup>— language for streaming columnar data processing. Standard in unix-like OS.
2. Actual AWK is outdated, use mawk (fast) or gawk (flexible).
3. Limited data structures: strings, **associative arrays (hash maps)** and regexps.
4. Built-in variables:
  - ▶ \$1, \$2, ... (\$0 is entire record)
  - ▶ NR - number of processed lines (records)
  - ▶ NF - number of columns (fields)
5. Use vars without declaration. Default values are 0s. One liners. **Hipster friendly**.

---

<sup>1</sup>Tutorial by Bruce Barnett. Careful, he [writes his blog in txt](#)

# AWK Examples 1 & 2

1. Count number of words and lines at [codeconf.hk](https://codeconf.hk):

```
cat codeconf.md | awk '{w += NF}END{print NR, w}'
```

```
370 1445
```

2. Most popular words on [codeconf.hk](https://codeconf.hk) website:

```
cat codeconf.md \  
| awk '{for(i=1; i<=NF; i++) words[tolower(\$i)]++}  
      END{for(w in words) print w, words[w]}' \  
| sort -k2 -nr
```

Most popular non stop-words: "Serverles" and "Android".

SEO winners: Davide Benvegnù and Richard Cohen.

10:15 - 10:45

Davide Benvegnù



## Go **Serverless** - Design Patterns and Best Practices

**Serverless** compute makes it sooooo easy to create an http endpoint or just run arbitrary code in the cloud.

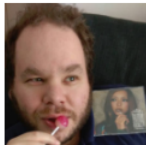
But with great power comes great responsibility and often users make fundamental design mistakes that end up effecting their **Serverless** performance.

After a brief Azure **Serverless** services introduction, will dive into **Serverless** design principles and architecture considerations effecting performance and overall functionality.

17:00 - 17:30

Richard Cohen

Kotlin for Android



The biggest cheer at I/O 2017 was for the announcement that Google would be supporting Kotlin as a first class language for Android app dev.

We're very close to that being delivered in Android Studio 3.0, so - what is Kotlin?

Where does it fit in the Android world?

Did they do it just so they could finally hire Jake Wharton?

Why would you want to use it?

How does it compare either Java?

Why **wouldn't** you want to use it?

What does the future hold?

In summary:

A language overview, comparison with Java, and how to do Android dev with it

## 3. Find the longest line in the text (if-then-else example):

```
cat codeconf.md | awk '{
    l = length > length(1) ? \ $0 : 1
}END{
    print length(1), 1
}'
```

```
146 * We believe that the Hong Kong developer community is skilled and
diverse, but that often these skills end up hidden away in big organisations.
```

## Demo Time

# Table of content

1. Problem Background
2. AWK Bootcamp in 5 min
3. Tabtools architecture and features
4. Stock example with NYSE:CS

# Basic concepts

1. Special files format: tsv + header (meta information). Easy to convert and autogenerate headers.

#	<i>Date</i>	<i>Open</i>	<i>High</i>	<i>Low</i>	<i>Close</i>	<i>Volume</i>
	2014-02-21	84.35	84.45	83.9	83.45	17275.0

2. Python script manages file descriptors headers, convert column names to column numbers and executes command line command, e.g. cat/tail/sort.
3. Heavy lifting goes to awk: tawk (map) and tgrp (map-reduce).
4. Based on command line expressions, it generates awk command and executes it with incoming stream.
5. Visual sugar: tpretty and tplot.



# Features

1. Streaming expressions: parametrized running/total sum/average/maximum<sup>2</sup>.
2. Aggregators: first, last, min, max, count.
3. Modules: `deque`.
4. Build to self-contained 2k LOC portable python (2.7, 3.3+) scripts.
5. All together: zero-configuration extensible sql in command line. It is readable and faster than a generic python/cython code (even after shedskin) and perl.

---

<sup>2</sup>moving maximum in `linear time` with `deque` implemented on top of awk associative arrays.

# Solutions comparison

Dell xps 15, 16Gb RAM, 8 CPUs:

	Python	PostgreSQL	gawk	mawk	Tabtools
Read time	104.4	180.3	0	0	0
Q1: "max" time	0	15.2	22.8	12.2	12.8
Q2: "group + avg" time	0	5.8	30.5	12.6	26.6 <sup>3</sup>
Q3: "filter + sort" time	21.3	33.6	174.2	36.3	33.5
<b>Total, sec.</b>	125.7	243.9	227.5	<b>61.1</b>	<b>72.9</b>

---

<sup>3</sup>Uses  $\Omega(n \log(n))$  complexity instead of  $\Omega(n)$ . Could be improved.

# Table of content

1. Problem Background
2. AWK Bootcamp in 5 min
3. Tabtools architecture and features
4. Stock example with NYSE:CS

# Data description

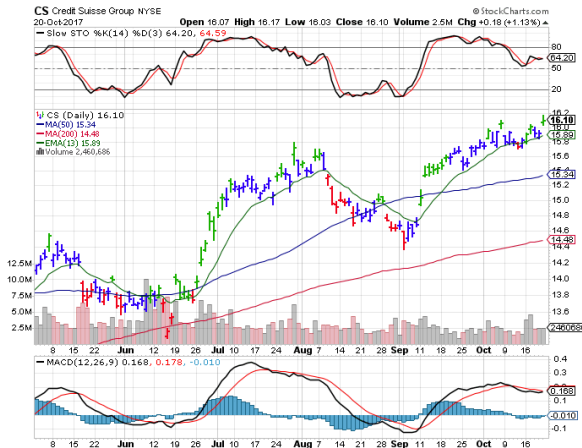
Credit Suisse (NYSE:CS) daily stock data from [Yahoo Finance](#): 'CS.csv' + 'cs.tsv'.

```
cat cs.tsv | tgrp \  
  -k "Week=strftime(\"%U\", DateEpoch(Date))" \  
  -g "Date=FIRST(Date)" \  
  -g "Open=FIRST(Open)" \  
  -g "High=MAX(High)" \  
  -g "Low=MIN(Low)" \  
  -g "Close=LAST(Close)" \  
  -g "Volume=SUM(Volume)" \  
| ttail \  
| tsrt -k Date:desc \  
| tpretty
```

# Demo Time

1. Moving Average for window size 200 and 50.
2. Exponential moving average for window size 26 and 13.
3. MACD(26, 12, 9) histogram.
4. Moving maximum and minimum for window size 14.
5. Fast and Slow Stochastics.

# Demo: plot (expected and actual)



# Thank you!

Kirill Pavlov <k@p99.io>, Recruiter, [Terminal 1](#).

GitHub: [@pavlov99](#) | Presentation: [2017-10-22-codeconf](#) | [tabtools](#)