# A high throughput cloud computation architecture for 'deep' parsing

Alexandre Rademaker and Henrique Muniz

IBM Research, Brazil

July 15, 2019

# What is it?

It is our first steps on the use of emerging technologies for distributed cloud computing for building scalable and high-performance architecture for 'deep' parsing with DELPHI-IN tools.
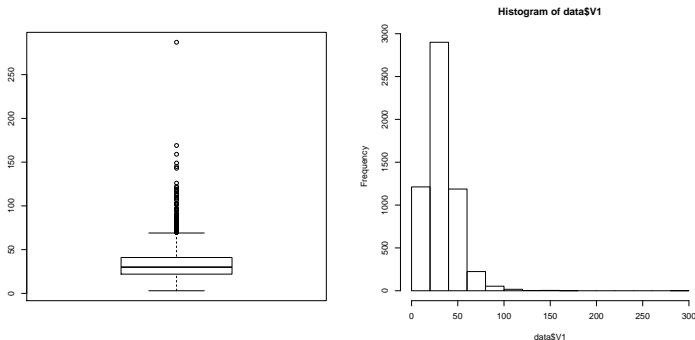
## Our goals

- a high throughput architecture
- as simple as possible
- scalable, pay what you need
- flexibility

# The available options

- LOGON's batch parsing script (pvm library last release 2005)
- Heart of Gold middleware
- ACE and Arbiter

# Data for experiments



A corpus with 5602 sentences obtained from 155 text passages relevant to petroleum systems extracted from documents randomly selected from a corpus of 1298 publicly available **English** language geological reports, published by the United States Geological Survey (USGS), Geological Survey of Canada (GSC), and British Geological Survey (BGS).

# Cluster 1

One node with 40 cores and 126 GB RAM.

## LOGON parsing script

- `-count 4` takes 8 hours 30 min
- `-count 8` takes 6 hours 30 min
- `-count 10` takes 4 hours 30 min
- `-count 20` takes 4 hours 40 minutes

# Cloud Architecture



Libraries and tools: ACE, PyDephin, Kubernets and Docker, Python RQ (Redis Queues).
... but we are lisp programmers! ;-)
https://github.com/own-pt/k8s-delphin-parsing

# IBM Cloud Kubernets Cluster Service

IBM Cloud Kubernets Cluster (RIS):
15 workers, 56 cores, 242 GB RAM.

2273/5602 sentences (many results
PyDelphin could not read?!) in 3
hours using 8 workers.

# IBM Cloud Kubernets Cluster Service

IBM Cloud Kubernets Cluster (RIS): 15 workers, 56 cores, 242 GB RAM.

2273/5602 sentences (many results PyDelphin could not read?!) in 3 hours using 8 workers.

# IBM Cloud Kubernets Cluster Service

IBM Cloud Kubernets Cluster (RIS): 15 workers, 56 cores, 242 GB RAM.

2273/5602 sentences (many results PyDelphin could not read?!) in 3 hours using 8 workers.

# Future work

- It should run in any cloud environment that supports Kubernets! But we need to try more cloud environments.
- Standard protocols from DELPH-IN (i.e. ErgApi?)
- Part of an internal text processing pipeline for IE of scientific articles from the O&G domain. More experiments.