# GATECloud.net:
# Cloud Infrastructure for
# Large-Scale, Open-Source Text Processing

Valentin Tablan    Ian Roberts    Hamish Cunningham
Kalina Bontcheva

University of Sheffield

28 September 2011

# GATE is. . .

- Infrastructure for language engineering
- Mature (1995 -)
- Open source (LGPL, AGPL)
- Developed (mainly) at Sheffield
- GATE Family includes desktop and server software
- ≈20..30,000 downloads / year
- http://gate.ac.uk

# Once Upon a Time. . .

### The year was 2007. . .

- **MATRIXWARE**™
  INFORMATION SERVICES
- Semantic Annotation of patents
- "We can do that!"
- "We've done large collections before (like 60,000 documents)!"

# Once Upon a Time. . .

### The year was 2007. . .

- **MATRIXWARE** INFORMATION SERVICES
- Semantic Annotation of patents
- "We can do that!"
- "We've done large collections before (like 60,000 documents)!"

### Let's start small. . .

- ". . . with about 20 million documents"

### Did you say "20 **Millions**"?!

- $\approx$ 5 sec / document
- 20 Million docs $\approx$ 3 years processing time

# Fast Forward to 2009

After. . .

- two years
- a few '00,000s €

We now had. . .

- some new software:
  the *GATE parallelizer*
- one IBM compute cluster (10 nodes)

. . . and were able to process the 20 million collection. . .

- in under 1 week
- using:
  - a *SUN Grid Engine* install
  - several ad-hoc shell scripts
  - and a lot of operator time

# Fast Forward to 2009

After. . .

- two years
- a few '00,000s €

We now had. . .

- some new software: the *GATE parallelizer*
- one IBM compute cluster (10 nodes)

. . . and were able to process the 20 million collection. . .

- in under 1 week
- using:
  - a *SUN Grid Engine* install
  - several ad-hoc shell scripts
  - and a lot of operator time

## Big data is hard!

# Processing Large Collections is Useful for Scientists Too!

E.g.:

- Social Media: millions of new tweets every day
- Life Sciences: millions of Medline abstracts
- Web Research: focused web crawling can still return millions of pages
- . . .

# Processing Large Collections is Useful for Scientists Too!

E.g.:

- Social Media: millions of new tweets every day
- Life Sciences: millions of Medline abstracts
- Web Research: focused web crawling can still return millions of pages
- . . .

Know any scientists that have a few '00,000s € to spare?

# 2011: Introducing GATECloud.net

## Design Ideas / Requirements

- Use our experience to automate large processing jobs
- Maximum automation, minimal operator intervention
- Scalability by distributing effort on a IaaS platform
- Should work with unmodified GATE pipelines

## Offering

- Large processing jobs
- On-demand dedicated servers with specialised server software for text processing

Motivation and Background
GATECloud.net
Evaluation

**Design**
On-demand Servers
Annotation Jobs

# GATECloud.net

### Decisions & Compromises

- No AMIs (use parameterised cloud-init scripts instead)
    - (+) simpler upgrade path
    - (+) no storage costs for root partitions
    - (-) slower start-up ($\approx$5 minutes)
- No 'elastic IPs' (not enough addresses available); use a proxy server instead.
    - (+) unlimited names (SSL capable)
    - (-) slower throughput (proxy in the cloud could mitigate)
    - (-) single point of failure
- Sensible defaults and centralised administration instead of user configuration
    - (+) much more user-friendly
    - (-) potentially sub-optimal (but user **mis-**configuration is also a risk!)

Motivation and Background
GATECloud.net
Evaluation

**Design**
On-demand Servers
Annotation Jobs

# GATECloud.net: Architectural Elements

## Web App

- On-line store (item descriptions, pricing rules)
- User support (user accounts, security)
- Admin (defining new shop items, managing server swarms)
- Workflow management (orchestrating the platform)

## "DNS" Proxy Server

- Secure traffic routing

## The Cloud

- EC2: compute instances on demand
- SimpleDB: management of workflow tasks
- S3: software distros, job data, job results

Motivation and Background
GATECloud.net
Evaluation

Design
On-demand Servers
Annotation Jobs

# On-demand Dedicated Servers

## What

- GATE Teamware
- GATE Mímir

## Why

- Complex build process (multiple web apps, programming languages, setup parameters)
- Complex dependencies (runtime, database setup, application server, etc.)

Motivation and Background
GATECloud.net
Evaluation

Design
On-demand Servers
Annotation Jobs

# On-demand Servers: Making a Reservation

Motivation and Background
GATECloud.net
Evaluation

Design
On-demand Servers
Annotation Jobs

# On-demand Servers: Starting a Server

Motivation and Background    Design
GATECloud.net    On-demand Servers
Evaluation    Annotation Jobs

# Annotation Jobs

- Process (large) document collections
- User friendly process:
    - upload documents
    - upload processing pipeline (or use pre-configured one)
    - press 'Go'
    - download results
- Maximum throughput
    - Multi-threading on each node
    - Scale up by spinning up more nodes
- Automated workflow management
- Web-based front end

Motivation and Background
GATECloud.net
Evaluation

Design
On-demand Servers
Annotation Jobs

# Annotation Jobs: Workflow

## Evaluation

| | | Time (hh:mm:ss) | | | Speed | Cost (GBP) | |
|---|---|---|---|---|---|---|---|
| | | CPU Time | Computer Time | Clock Time | (Kb/s) | £/GB | Total |
| **Experiment 1: Patents** | Desktop | N/A | N/A | N/A | | | |
| 100,000 patent documents | Server | 91:42:00 | 18:39:54 | 18:39:54 | 85.33 | | |
| | Cloud | 162:38:00 | 16:56:37 | 02:03:32 | 773.6 | £3.08 | £16.83 |
| **Experiment 2: News** | Desktop | 05:20:19 | 05:20:19 | 05:20:19 | 71.52 | | |
| 20,000 documents | Server | 04:43:00 | 03:08:00 | 03:08:00 | 121.86 | | |
| | Cloud | 07:47:00 | 01:21:20 | 00:35:31 | 645.04 | £1.51 | £1.98 |
| **Experiment 3: Tweets** | Desktop | 32:28:46 | 32:28:46 | 32:28:46 | 52.80 | | |
| 50,000,000 tweets | Server | 22:16:15 | 03:19:12 | 03:19:12 | 516.53 | | |
| | Cloud | 40:08:00 | 07:00:14 | 01:25:46 | 1199.69 | £1.35 | £7.92 |

### Did it work?

20,000,000 patents:

| | **2007** | **2009** | **2011** |
|---|---|---|---|
| **Time** | 3 years | ≈7 days | ≈16 days (could be made arbitrarily shorter with a larger swarm) |
| **Cost** | N/A | 00,000s € investment | ≈3,500 € |

# Thank you!

### Questions?

### More Information

- http://gatecloud.net
- http://gate.ac.uk

### Future Work

- Investigate the use of FTL neutrinos for getting results <u>before</u> the job is submitted?