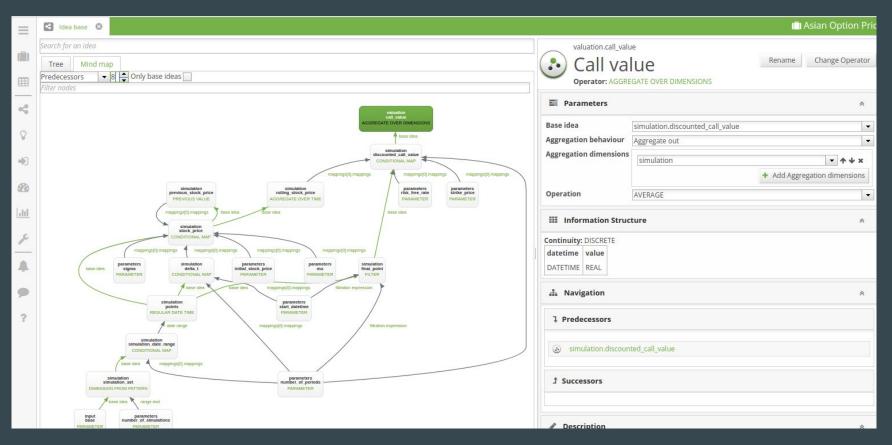# Managing infrastructure on EC2 Spot
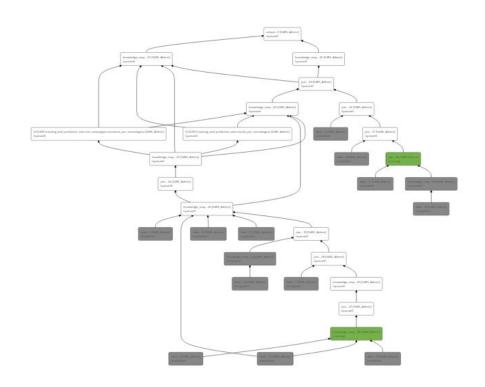
Peter Juriz

# What we do

- Modelling Platform for various forms of data:
- Financial / Insurance / Logistical / Medical/ Marketing
- Compute, Simulation, Machine Learning and Data Visualization
- Modelling has a visual component

Graph based compute engine (C++)

# Running

Compile these graphs into dependent Hadoop Map-Reduce jobs

Calculations can get pretty big

**hadoop02**
```
1 [|                1.3%]   5 [                0.0%]   9 [                0.0%]  13 [                0.0%]
2 [                0.0%]   6 [                0.0%]  10 [                0.0%]  14 [                0.0%]
3 [                0.0%]   7 [                0.0%]  11 [                0.0%]  15 [                0.0%]
4 [                0.0%]   8 [                0.0%]  12 [                0.0%]  16 [                0.0%]
Mem[|||            981/122952MB]   Tasks: 31, 7 thr; 1 running
Swp[               0/0MB]          Load average: 0.00 0.01 0.05
                                  Uptime: 16:27:10
```

**hadoop03**
```
1 [||||||100.0%]   5 [||||||100.0%]   9 [||||||100.0%]  13 [||||||100.0%]
2 [||||||100.0%]   6 [||||||100.0%]  10 [||||||99.4%]   14 [||||||99.4%]
3 [||||||100.0%]   7 [||||||99.4%]   11 [||||||100.0%]  15 [||||||100.0%]
4 [||||||100.0%]   8 [||||||100.0%]  12 [||||||100.0%]  16 [||||||99.4%]
Mem[||||||  11275/122952MB]   Tasks: 71, 1672 thr; 26 running
Swp[               0/307199MB]  Load average: 24.70 15.43 11.82
                               Uptime: 16:09:35
```

**hadoop04**
```
1 [||||||61.5%]    5 [||||||55.7%]   9 [||||||58.9%]   13 [||||||58.2%]
2 [||||||100.0%]   6 [||||||93.0%]  10 [||   14.6%]    14 [||||  36.1%]
3 [||||||55.5%]    7 [||||||88.6%]  11 [||||||55.4%]   15 [||||||54.3%]
4 [||||||5.1%]     8 [||||||100.0%] 12 [||||||67.3%]   16 [||||||41.4%]
Mem[||||||  4774/122952MB]    Tasks: 50, 957 thr; 12 running
Swp[               278/307199MB]  Load average: 9.99 9.89 9.79
                               Uptime: 1 day, 20:23:53
```

**hadoop05**
```
1 [||||||100.0%]   5 [||||||100.0%]   9 [||||||100.0%]  13 [||||||100.0%]
2 [||||||100.0%]   6 [||||||100.0%]  10 [||||||100.0%]  14 [||||||100.0%]
3 [||||||100.0%]   7 [||||||100.0%]  11 [||||||100.0%]  15 [||||||100.0%]
4 [||||||100.0%]   8 [||||||99.4%]   12 [||||||100.0%]  16 [||||||99.4%]
Mem[|||||||||  7366/122952MB]  Tasks: 59, 1454 thr; 41 running
Swp[               0/307199MB]  Load average: 25.42 14.76 12.31
                               Uptime: 16:09:34
  PID USER      PRI  NI  VIRT   RES   SHR S CPU% MEM%   TIME+  Command
74952 mapred     20   0 1856M  135M 16192 S 115.  0.1  0:03.87 /usr/lib/jvm/ja
F1Help F2Setup F3Search F4Filter F5Tree F6SortBy F7Nice - F8Nice + F9Kill F10Quit
```

**hadoop06**
```
               99.4%]   5 [||||||99.4%]   9 [||||||98.2%]  13 [||||||98.1%]
               98.8%]   6 [||||||99.4%]  10 [||||||99.4%]  14 [||||||99.4%]
               97.5%]   7 [||||||99.4%]  11 [||||||100.0%]  15 [||||||96.9%]
               99.4%]   8 [||||||100.0%] 12 [||||||100.0%]  16 [||||||99.4%]
               9730/122952MB]   Tasks: 64, 1649 thr; 23 running
               0/307199MB]      Load average: 25.09 15.27 11.87
                                Uptime: 16:01:31
```

**hadoop07**
```
1 [||||||99.4%]    5 [||||||99.4%]    9 [||||||99.4%]   13 [||||||100.0%]
2 [||||||99.4%]    6 [||||||99.4%]   10 [||||||99.4%]   14 [||||||95.7%]
3 [||||||99.4%]    7 [||||||98.8%]   11 [||||||99.4%]   15 [||||||99.4%]
4 [||||||98.2%]    8 [||||||99.4%]   12 [||||||99.4%]   16 [||||||99.4%]
Mem[|||||||   11592/122952MB]   Tasks: 71, 1747 thr; 24 running
Swp[               0/307199MB]   Load average: 27.89 14.38 10.11
                                Uptime: 16:01:33
  PID USER      PRI  NI  VIRT   RES   SHR S CPU% MEM%   TIME+  Command
124094 mapred    20   0 1853M  140M 16252 S 136.  0.1  0:04.03 /usr/lib/jvm/ja
```

**hadoop08 (left)**
```
1 [||||||99.4%]    5 [||||||100.0%]   9 [||||||99.4%]   13 [||||||98.8%]
2 [||||||100.0%]   6 [||||||99.4%]   10 [||||||100.0%]  14 [||||||100.0%]
3 [||||||100.0%]   7 [||||||100.0%]  11 [||||||100.0%]  15 [||||||98.8%]
4 [||||||98.2%]    8 [||||||99.4%]   12 [||||||99.4%]   16 [||||||100.0%]
Mem[|||||||||  9439/122952MB]   Tasks: 62, 1497 thr; 26 running
Swp[               0/307199MB]   Load average: 24.69 13.75 10.99
                                Uptime: 16:01:33
  PID USER      PRI  NI  VIRT   RES   SHR S CPU% MEM%   TIME+  Command
111755 mapred    20   0 1839M  129M 16184 S 109.  0.1  0:03.33 /usr/lib/jvm/ja
F1Help F2Setup F3Search F4Filter F5Tree F6SortBy F7Nice - F8Nice + F9Kill
```

**hadoop09**
```
               100.0%]   5 [||||||100.0%]   9 [||||||100.0%]  13 [||||||100.0%]
               100.0%]   6 [||||||100.0%]  10 [||||||100.0%]  14 [||||||100.0%]
               100.0%]   7 [||||||100.0%]  11 [||||||100.0%]  15 [||||||100.0%]
               100.0%]   8 [||||||100.0%]  12 [||||||100.0%]  16 [||||||100.0%]
               11399/122952MB]   Tasks: 71, 1833 thr; 39 running
               0/307199MB]       Load average: 35.66 19.63 12.69
                                 Uptime: 16:01:34
               1857M  473M 16844 S 100.  0.4  1:54.51 /usr/lib/jvm/ja
```

**hadoop10**
```
1 [||||||100.0%]   5 [||||||99.4%]    9 [||||||98.2%]   13 [||||||100.0%]
2 [||||||100.0%]   6 [||||||99.4%]   10 [||||||99.4%]   14 [||||||100.0%]
3 [||||||99.4%]    7 [||||||100.0%]  11 [||||||100.0%]  15 [||||||100.0%]
4 [||||||99.4%]    8 [||||||98.2%]   12 [||||||100.0%]  16 [||||||100.0%]
Mem[|||||||||  18067/122952MB]   Tasks: 76, 1930 thr; 29 running
Swp[               0/307199MB]    Load average: 32.65 16.08 11.61
                                 Uptime: 16:01:32
  PID USER      PRI  NI  VIRT   RES   SHR S CPU% MEM%   TIME+  Command
124494 mapred    20   0 1839M  133M 16204 S 119.  0.1  0:03.71 /usr/lib/jvm/ja
```

**hadoop11 (left)**
```
1 [||||||97.5%]    5 [||||||98.1%]    9 [||||||98.8%]   13 [||||
2 [||||||94.5%]    6 [||||||93.8%]   10 [||||||95.1%]   14 [||||
3 [||||||100.0%]   7 [||||||96.9%]   11 [||||||98.1%]   15 [
4 [||||||97.5%]    8 [||||||97.5%]   12 [||||||98.2%]   16 [
Mem[||||||   17477/122952MB]   Tasks: 79, 1933 thr; 21 runn
Swp[               0/307199MB]   Load average: 31.17 15.25 10
                                Uptime: 16:01:33
  PID USER      PRI  NI  VIRT   RES   SHR S CPU% MEM%   TIME+  Command
122604 mapred    20   0 1919M  227M 16028 S 128.  0.2  0:05.00 /usr/
```

**hadoop12**
```
1 [||||||100.0%]   5 [||||||99.4%]    9 [||||||100.0%]  13 [||||||98.8%]
2 [||||||100.0%]   6 [||||||98.8%]   10 [||||||99.4%]   14 [||||||100.0%]
3 [||||||97.5%]    7 [||||||99.4%]   11 [||||||99.4%]   15 [||||||99.4%]
4 [||||||100.0%]   8 [||||||100.0%]  12 [||||||100.0%]  16 [||||||99.4%]
Mem[||||||   17521/122952MB]   Tasks: 77, 1930 thr; 27 running
Swp[               0/307199MB]   Load average: 33.11 16.68 11.05
                                Uptime: 16:01:30
  PID USER      PRI  NI  VIRT   RES   SHR S CPU% MEM%   TIME+  Command
115665 mapred    20   0 1853M  146M 16260 S 124.  0.1  0:04.12 /usr/lib/jvm/ja
F1Help F2Setup F3Search F4Filter F5Tree F6SortBy F7Nice - F8Nice + F9Kill F10Quit
```

**hadoop13**
```
1 [||||||99.4%]    5 [||||||100.0%]   9 [||||||100.0%]  13 [||||||100.0%]
2 [||||||100.0%]   6 [||||||100.0%]  10 [||||||100.0%]  14 [||||||98.2%]
3 [||||||99.4%]    7 [||||||100.0%]  11 [||||||99.4%]   15 [||||||99.4%]
4 [||||||99.4%]    8 [||||||100.0%]  12 [||||||100.0%]  16 [||||||100.0%]
Mem[||||||   8426/122952MB]    Tasks: 61, 1569 thr; 32 running
Swp[               0/307199MB]   Load average: 24.68 13.79 10.98
                                Uptime: 16:01:32
  PID USER      PRI  NI  VIRT   RES   SHR S CPU% MEM%   TIME+  Command
89978 mapred     20   0 1843M  152M 16292 S 122.  0.1  0:04.62 /usr/lib/jvm/ja
F1Help F2Setup F3Search F4Filter F5Tree F6SortBy F7Nice - F8Nice + F9Kill F10Quit
```

**bottom-left node**
```
1 [||||||99.4%]    5 [||||||94.5%]    9 [||||||97.5%]   13 [||||||98.8%]
2 [||||||97.0%]    6 [||||||98.8%]   10 [||||||98.2%]   14 [||||||95.6%]
3 [||||||95.1%]    7 [||||||96.9%]   11 [||||||98.8%]   15 [||||||98.8%]
4 [||||||100.0%]   8 [||||||98.2%]   12 [||||||96.3%]   16 [||||||95.7%]
Mem[||||||   17021/122952MB]   Tasks: 81, 1910 thr; 22 running
Swp[               0/307199MB]   Load average: 31.77 16.30 11.10
                                Uptime: 16:01:28
  PID USER      PRI  NI  VIRT   RES   SHR S CPU% MEM%   TIME+  Command
124649 mapred    20   0 1847M  103M 16176 S 132.  0.1  0:02.35 /usr/lib/jvm/ja
```

**bottom-center node**
```
               98.8%]   5 [||||||100.0%]   9 [||||||100.0%]  13 [||||||95.7%]
               98.8%]   6 [||||||98.8%]   10 [||||||95.0%]   14 [||||||95.7%]
               99.4%]   7 [||||||94.5%]   11 [||||||99.4%]   15 [||||||100.0%]
               98.8%]   8 [||||||99.4%]   12 [||||||99.4%]   16 [||||||98.1%]
               8400/122952MB]   Tasks: 66, 1556 thr; 26 running
               0/307199MB]      Load average: 21.36 11.74 9.17
                                Uptime: 16:01:33
               1891M  124M 15772 S 124.  0.1  0:03.14 /usr/lib/jvm/ja
```

**bottom-right node**
```
1 [||||||100.0%]   5 [||||||100.0%]   9 [||||||100.0%]  13 [||||||99.4%]
2 [||||||100.0%]   6 [||||||100.0%]  10 [||||||98.8%]   14 [||||||100.0%]
3 [||||||100.0%]   7 [||||||100.0%]  11 [||||||100.0%]  15 [||||||100.0%]
4 [||||||100.0%]   8 [||||||99.4%]   12 [||||||100.0%]  16 [||||||100.0%]
Mem[||||||   16675/122952MB]   Tasks: 75, 1860 thr; 38 running
Swp[               0/307199MB]   Load average: 37.64 19.01 13.65
                                Uptime: 16:00:13
  PID USER      PRI  NI  VIRT   RES   SHR S CPU% MEM%   TIME+  Command
99693 mapred     20   0 1839M  141M 16248 S 102.  0.1  0:03.97 /usr/lib/jvm/ja
F1Help F2Setup F3Search F4Filter F5Tree F6SortBy F7Nice - F8Nice + F9Kill F10Quit
```

# Motivations for using EC2 Spot

- Big calculations require many instances
- Much cheaper ~ $0.24 spot vs $1.06 on demand for r4.4xlarge per hour (4.4 times cheaper) - when not surging
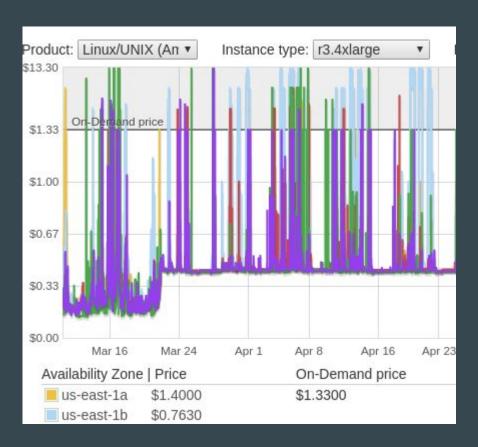- Can handle losing some machines

# How Spot Pricing Works

- Quasi market
- Set a bid, pay the market price
- If market price goes above your bid => your instance gets terminated
- Spot price can go 10X above the on demand price
- Different prices for different instance types

# Different volatility for different instance types



Product: Linux/UNIX (An ▾)  Instance type: r4.4xlarge ▾  Date range: 3 hours ▾

| Availability Zone | Price | On-Demand price | Date |
|---|---|---|---|
| us-east-1a | $0.2434 | $1.0640 | 6/9/2017, 12:50:31 PM UTC+0200 |
| us-east-1b | $0.2529 | | |
| us-east-1c | $0.2435 | | |
| us-east-1d | $0.2362 | | |
| us-east-1e | $0.2445 | | |

Product: Linux/UNIX (An ▾)  Instance type: r3.4xlarge ▾  Date range: 3 hours ▾

| Availability Zone | Price | On-Demand price | Date |
|---|---|---|---|
| us-east-1a | $0.4268 | $1.3300 | 6/9/2017, 12:50:50 PM UTC+0200 |
| us-east-1b | $0.4365 | | |
| us-east-1c | $0.4358 | | |
| us-east-1d | $0.4250 | | |
| us-east-1e | $0.4998 | | |

# Jumps in market price

# Problems we've encountered

- What if price stays up?
- Multi-AZ - incurs interzone traffic costs
- Losing local disk can be a problem
- Some services don't handle node churn

# Having insights is important

# Build internal monitoring tools

Gryphon

Elephant Cluster

Monitored Tasks

Query Monitor

Condor

Cron

**Cluster Size:** Minimum `1` Desired `24` Maximum `24` `Set size`

| 10.10.10.13 | RUNNING |
|---|---|
| r4.4xlarge | |
| uptime: 0:33:45.351762 | |

cores (16/16)

memory (65536mb/119952mb)

---

| 10.10.10.136 | STARTUP |
|---|---|
| r4.4xlarge | |
| uptime: 0:03:16.351762 | |

---

| 10.10.10.147 | RUNNING |
|---|---|
| r4.4xlarge | |
| uptime: 0:42:33.351762 | |

cores (16/16)

memory (65536mb/119952mb)

---

| 10.10.10.186 | RUNNING |
|---|---|
| r4.4xlarge | |
| uptime: 0:42:33.351762 | |

cores (16/16)

memory (63488mb/119952mb)

---

| 10.10.10.208 | RUNNING |
|---|---|
| r4.4xlarge | |
| uptime: 2:29:43.351762 | |

cores (16/16)

memory (65536mb/119952mb)

---

| 10.10.10.210 | STARTUP |
|---|---|
| r4.4xlarge | |
| uptime: 0:03:16.351762 | |

---

| 10.10.10.221 | STARTUP |
|---|---|
| r4.4xlarge | |
| uptime: 0:03:16.351762 | |

---

| 10.10.10.230 | RUNNING |
|---|---|
| r4.4xlarge | |
| uptime: 0:33:45.351762 | |

cores (16/16)

memory (65536mb/119952mb)

---

| 10.10.10.47 | RUNNING |
|---|---|
| r4.4xlarge | |
| uptime: 0:33:45.351762 | |

cores (16/16)

memory (65536mb/119952mb)

---

| 10.10.10.53 | RUNNING |
|---|---|
| r4.4xlarge | |
| uptime: 0:42:33.351762 | |

cores (16/16)

memory (65536mb/119952mb)

---

| 10.10.10.7 | RUNNING |
|---|---|
| r4.4xlarge | |
| uptime: 0:33:45.351762 | |

cores (16/16)

memory (65536mb/119952mb)

---

| 10.10.10.86 | STARTUP |
|---|---|
| r4.4xlarge | |
| uptime: 0:03:16.351762 | |

---

| 10.10.20.135 | RUNNING |
|---|---|
| r4.4xlarge | |
| uptime: 0:33:45.351762 | |

cores (16/16)

memory (65536mb/119952mb)

---

| 10.10.20.142 | RUNNING |
|---|---|
| r4.4xlarge | |
| uptime: 0:42:35.351762 | |

cores (16/16)

memory (65536mb/119952mb)

---

| 10.10.20.144 | STARTUP |
|---|---|
| r4.4xlarge | |
| uptime: 0:03:19.351762 | |

---

| 10.10.20.176 | RUNNING |
|---|---|
| r4.4xlarge | |
| uptime: 0:42:35.351762 | |

cores (16/16)

memory (65536mb/119952mb)

---

| 10.10.20.198 | STARTUP |
|---|---|
| r4.4xlarge | |
| uptime: 0:03:19.351762 | |

---

| 10.10.20.207 | STARTUP |
|---|---|
| r4.4xlarge | |
| uptime: 0:03:19.351762 | |

---

| 10.10.20.22 | RUNNING |
|---|---|
| r4.4xlarge | |
| uptime: 0:33:45.351762 | |

cores (16/16)

memory (65536mb/119952mb)

---

| 10.10.20.41 | RUNNING |
|---|---|
| r4.4xlarge | |
| uptime: 0:33:45.351762 | |

cores (16/16)

memory (65536mb/119952mb)

---

| 10.10.20.53 | RUNNING |
|---|---|
| r4.4xlarge | |
| uptime: 0:42:35.351762 | |

cores (16/16)

memory (65536mb/119952mb)

---

| 10.10.20.56 | RUNNING |
|---|---|
| r4.4xlarge | |
| uptime: 0:33:45.351762 | |

cores (16/16)

memory (65536mb/119952mb)

---

| 10.10.20.68 | RUNNING |
|---|---|
| r4.4xlarge | |
| uptime: 0:42:35.351762 | |

cores (16/16)

memory (65536mb/119952mb)

---

| 10.10.20.74 | STARTUP |
|---|---|
| r4.4xlarge | |
| uptime: 0:03:19.351762 | |

Cluster Size: Minimum `1` Desired `24` Maximum `24` `Set size`

| 10.10.10.13 RUNNING |
| r4.4xlarge |
| uptime: 0:52:10.051000 |
| cores (16/16) |
| memory (65536mb/119952mb) |

| 10.10.10.136 RUNNING |
| r4.4xlarge |
| uptime: 0:21:41.051000 |
| cores (16/16) |
| memory (65536mb/119952mb) |

| 10.10.10.147 RUNNING |
| r4.4xlarge |
| uptime: 1:00:58.051000 |
| cores (16/16) |
| memory (65536mb/119952mb) |

| 10.10.10.186 RUNNING |
| r4.4xlarge |
| uptime: 1:00:58.051000 |
| cores (16/16) |
| memory (65536mb/119952mb) |

| 10.10.10.208 RUNNING |
| r4.4xlarge |
| uptime: 2:48:08.051000 |
| cores (16/16) |
| memory (65536mb/119952mb) |

| 10.10.10.210 RUNNING |
| r4.4xlarge |
| uptime: 0:21:41.051000 |
| cores (16/16) |
| memory (43008mb/119952mb) |

| 10.10.10.221 ERROR |
| r4.4xlarge |
| uptime: 0:21:41.051000 |

| 10.10.10.230 RUNNING |
| r4.4xlarge |
| uptime: 0:52:10.051000 |
| cores (16/16) |
| memory (65536mb/119952mb) |

| 10.10.10.47 RUNNING |
| r4.4xlarge |
| uptime: 0:52:10.051000 |
| cores (16/16) |
| memory (65536mb/119952mb) |

| 10.10.10.53 RUNNING |
| r4.4xlarge |
| uptime: 1:00:58.051000 |
| cores (16/16) |
| memory (65536mb/119952mb) |

| 10.10.10.7 RUNNING |
| r4.4xlarge |
| uptime: 0:52:10.051000 |
| cores (16/16) |
| memory (65536mb/119952mb) |

| 10.10.10.86 RUNNING |
| r4.4xlarge |
| uptime: 0:21:41.051000 |
| cores (16/16) |
| memory (65536mb/119952mb) |

| 10.10.20.135 RUNNING |
| r4.4xlarge |
| uptime: 0:52:10.051000 |
| cores (16/16) |
| memory (65536mb/119952mb) |

| 10.10.20.142 RUNNING |
| r4.4xlarge |
| uptime: 1:01:00.051000 |
| cores (16/16) |
| memory (65536mb/119952mb) |

| 10.10.20.144 RUNNING |
| r4.4xlarge |
| uptime: 0:21:44.051000 |
| cores (16/16) |
| memory (65536mb/119952mb) |

| 10.10.20.176 RUNNING |
| r4.4xlarge |
| uptime: 1:01:00.051000 |
| cores (16/16) |
| memory (63488mb/119952mb) |

| 10.10.20.198 RUNNING |
| r4.4xlarge |
| uptime: 0:21:44.051000 |
| cores (16/16) |
| memory (65536mb/119952mb) |

| 10.10.20.207 RUNNING |
| r4.4xlarge |
| uptime: 0:21:44.051000 |
| cores (16/16) |
| memory (65536mb/119952mb) |

| 10.10.20.22 RUNNING |
| r4.4xlarge |
| uptime: 0:52:10.051000 |
| cores (16/16) |
| memory (65536mb/119952mb) |

| 10.10.20.41 RUNNING |
| r4.4xlarge |
| uptime: 0:52:10.051000 |
| cores (16/16) |
| memory (65536mb/119952mb) |

| 10.10.20.53 RUNNING |
| r4.4xlarge |
| uptime: 1:01:00.051000 |
| cores (16/16) |
| memory (65536mb/119952mb) |

| 10.10.20.56 RUNNING |
| r4.4xlarge |
| uptime: 0:52:10.051000 |
| cores (16/16) |
| memory (65536mb/119952mb) |

| 10.10.20.68 RUNNING |
| r4.4xlarge |
| uptime: 1:01:00.051000 |
| cores (16/16) |
| memory (65536mb/119952mb) |

| 10.10.20.74 RUNNING |
| r4.4xlarge |
| uptime: 0:21:44.051000 |
| cores (16/16) |
| memory (65536mb/119952mb) |

Cluster Size: Minimum 1    Desired 24    Maximum 24    Set size

| 10.10.10.13 | RUNNING |
| --- | --- |
| r4.4xlarge | |
| uptime: 1:01:26.495599 | |

cores (0/16)

memory (0mb/119952mb)

| 10.10.10.136 | RUNNING |
| --- | --- |
| r4.4xlarge | |
| uptime: 0:30:57.495599 | |

cores (0/16)

memory (0mb/119952mb)

| 10.10.10.147 | RUNNING |
| --- | --- |
| r4.4xlarge | |
| uptime: 1:10:14.495599 | |

cores (1/16)

memory (2048mb/119952mb)

| 10.10.10.186 | RUNNING |
| --- | --- |
| r4.4xlarge | |
| uptime: 1:10:14.495599 | |

cores (0/16)

memory (0mb/119952mb)

| 10.10.10.208 | RUNNING |
| --- | --- |
| r4.4xlarge | |
| uptime: 2:57:24.495599 | |

cores (0/16)

memory (0mb/119952mb)

| 10.10.10.210 | RUNNING |
| --- | --- |
| r4.4xlarge | |
| uptime: 0:30:57.495599 | |

cores (0/16)

memory (0mb/119952mb)

| 10.10.10.221 | RUNNING |
| --- | --- |
| r4.4xlarge | |
| uptime: 0:30:57.495599 | |

cores (0/16)

memory (0mb/119952mb)

| 10.10.10.230 | RUNNING |
| --- | --- |
| r4.4xlarge | |
| uptime: 1:01:26.495599 | |

cores (0/16)

memory (0mb/119952mb)

| 10.10.10.47 | RUNNING |
| --- | --- |
| r4.4xlarge | |
| uptime: 1:01:26.495599 | |

cores (0/16)

memory (0mb/119952mb)

| 10.10.10.53 | RUNNING |
| --- | --- |
| r4.4xlarge | |
| uptime: 1:10:14.495599 | |

cores (0/16)

memory (0mb/119952mb)

| 10.10.10.7 | RUNNING |
| --- | --- |
| r4.4xlarge | |
| uptime: 1:01:26.495599 | |

cores (0/16)

memory (0mb/119952mb)

| 10.10.10.86 | RUNNING |
| --- | --- |
| r4.4xlarge | |
| uptime: 0:30:57.495599 | |

cores (0/16)

memory (0mb/119952mb)

| 10.10.20.135 | RUNNING |
| --- | --- |
| r4.4xlarge | |
| uptime: 1:01:26.495599 | |

cores (1/16)

memory (4096mb/119952mb)

| 10.10.20.142 | RUNNING |
| --- | --- |
| r4.4xlarge | |
| uptime: 1:10:16.495599 | |

cores (0/16)

memory (0mb/119952mb)

| 10.10.20.144 | RUNNING |
| --- | --- |
| r4.4xlarge | |
| uptime: 0:31:00.495599 | |

cores (0/16)

memory (0mb/119952mb)

| 10.10.20.176 | RUNNING |
| --- | --- |
| r4.4xlarge | |
| uptime: 1:10:16.495599 | |

cores (0/16)

memory (0mb/119952mb)

| 10.10.20.198 | RUNNING |
| --- | --- |
| r4.4xlarge | |
| uptime: 0:31:00.495599 | |

cores (0/16)

memory (0mb/119952mb)

| 10.10.20.207 | RUNNING |
| --- | --- |
| r4.4xlarge | |
| uptime: 0:31:00.495599 | |

cores (0/16)

memory (0mb/119952mb)

| 10.10.20.22 | RUNNING |
| --- | --- |
| r4.4xlarge | |
| uptime: 1:01:26.495599 | |

cores (0/16)

memory (0mb/119952mb)

| 10.10.20.41 | RUNNING |
| --- | --- |
| r4.4xlarge | |
| uptime: 1:01:26.495599 | |

cores (0/16)

memory (0mb/119952mb)

| 10.10.20.53 | RUNNING |
| --- | --- |
| r4.4xlarge | |
| uptime: 1:10:16.495599 | |

cores (0/16)

memory (0mb/119952mb)

| 10.10.20.56 | RUNNING |
| --- | --- |
| r4.4xlarge | |
| uptime: 1:01:26.495599 | |

cores (0/16)

memory (0mb/119952mb)

| 10.10.20.68 | RUNNING |
| --- | --- |
| r4.4xlarge | |
| uptime: 1:10:16.495599 | |

cores (0/16)

memory (0mb/119952mb)

| 10.10.20.74 | RUNNING |
| --- | --- |
| r4.4xlarge | |
| uptime: 0:31:00.495599 | |

cores (0/16)

memory (0mb/119952mb)

Cluster Size: Minimum 1   Desired 24   Maximum 24   Set size

| | | | | |
|---|---|---|---|---|
| **10.10.10.13** RUNNING<br>r4.4xlarge<br>uptime: 1:12:53.577803<br>cores (16/16)<br>memory (65536mb/119952mb) | **10.10.10.136** RUNNING<br>r4.4xlarge<br>uptime: 0:42:24.577803<br>cores (16/16)<br>memory (65536mb/119952mb) | **10.10.10.147** RUNNING<br>r4.4xlarge<br>uptime: 1:21:41.577803<br>cores (16/16)<br>memory (65536mb/119952mb) | **10.10.10.186** RUNNING<br>r4.4xlarge<br>uptime: 1:21:41.577803<br>cores (16/16)<br>memory (65536mb/119952mb) | **10.10.10.208** RUNNING<br>r4.4xlarge<br>uptime: 3:08:51.577803<br>cores (16/16)<br>memory (63488mb/119952mb) |
| **10.10.10.210** RUNNING<br>r4.4xlarge<br>uptime: 0:42:24.577803<br>cores (16/16)<br>memory (65536mb/119952mb) | **10.10.10.221** RUNNING<br>r4.4xlarge<br>uptime: 0:42:24.577803<br>cores (16/16)<br>memory (65536mb/119952mb) | **10.10.10.230** RUNNING<br>r4.4xlarge<br>uptime: 1:12:53.577803<br>cores (16/16)<br>memory (65536mb/119952mb) | **10.10.10.47** RUNNING<br>r4.4xlarge<br>uptime: 1:12:53.577803<br>cores (16/16)<br>memory (65536mb/119952mb) | **10.10.10.53** RUNNING<br>r4.4xlarge<br>uptime: 1:21:41.577803<br>cores (16/16)<br>memory (65536mb/119952mb) |
| **10.10.10.7** RUNNING<br>r4.4xlarge<br>uptime: 1:12:53.577803<br>cores (16/16)<br>memory (65536mb/119952mb) | **10.10.10.86** RUNNING<br>r4.4xlarge<br>uptime: 0:42:24.577803<br>cores (16/16)<br>memory (65536mb/119952mb) | **10.10.20.135** RUNNING<br>r4.4xlarge<br>uptime: 1:12:53.577803<br>cores (16/16)<br>memory (65536mb/119952mb) | **10.10.20.142** RUNNING<br>r4.4xlarge<br>uptime: 1:21:43.577803<br>cores (16/16)<br>memory (65536mb/119952mb) | **10.10.20.144** RUNNING<br>r4.4xlarge<br>uptime: 0:42:27.577803<br>cores (16/16)<br>memory (65536mb/119952mb) |
| **10.10.20.176** RUNNING<br>r4.4xlarge<br>uptime: 1:21:43.577803<br>cores (16/16)<br>memory (65536mb/119952mb) | **10.10.20.198** RUNNING<br>r4.4xlarge<br>uptime: 0:42:27.577803<br>cores (16/16)<br>memory (65536mb/119952mb) | **10.10.20.207** RUNNING<br>r4.4xlarge<br>uptime: 0:42:27.577803<br>cores (16/16)<br>memory (65536mb/119952mb) | **10.10.20.22** RUNNING<br>r4.4xlarge<br>uptime: 1:12:53.577803<br>cores (16/16)<br>memory (65536mb/119952mb) | **10.10.20.41** RUNNING<br>r4.4xlarge<br>uptime: 1:12:53.577803<br>cores (16/16)<br>memory (65536mb/119952mb) |
| **10.10.20.53** RUNNING<br>r4.4xlarge<br>uptime: 1:21:43.577803<br>cores (16/16)<br>memory (65536mb/119952mb) | **10.10.20.56** RUNNING<br>r4.4xlarge<br>uptime: 1:12:53.577803<br>cores (16/16)<br>memory (65536mb/119952mb) | **10.10.20.68** RUNNING<br>r4.4xlarge<br>uptime: 1:21:43.577803<br>cores (16/16)<br>memory (65536mb/119952mb) | **10.10.20.74** RUNNING<br>r4.4xlarge<br>uptime: 0:42:27.577803<br>cores (16/16)<br>memory (65536mb/119952mb) | |

# Profile the code



**Query Performance Profile**
Total wall time: 16,754,000ms

Idea Profiling | Operator Profiling | Activity Profiling | Overhead Profiling | Knowledge Map Profiling

| Operator | Time Spent (ms) - Total: 228,821,645ms |
|---|---|
| dimensions_from_expressions | 81,840,081 (35.77%) |
| filter | 38,992,312 (17.04%) |
| aggregate_over_dimension_ranges | 35,581,190 (15.55%) |
| conditional_map | 27,207,787 (11.89%) |
| input | 12,892,003 (5.63%) |
| lookup | 10,775,539 (4.71%) |
| ideas_from_dimensions | 6,613,394 (2.89%) |
| aggregate_over_dimensions | 4,170,770 (1.82%) |
| predict | 3,005,895 (1.31%) |
| zip | 2,767,238 (1.21%) |
| unzip | 2,230,198 (0.97%) |
| rank | 1,942,622 (0.85%) |
| join | 658,194 (0.29%) |
| bucketing | 119,121 (0.05%) |
| rename_dimensions | 25,301 (0.01%) |

# Insights

- Monitoring can makes your life easier
- Compute coming up faster saves $$$ (docker vs puppet): image download vs package install
- Match cluster size to resource requirements
- Having the right autoscaling conf matters

# Autoscaling issues

- Don't flip on and off
  - Two Causes for us:
    - new job => scale up => complete => scale down => new job (within an instance hour)
    - Scale up too high too quickly => low resource utilization => scale down => (repeat)
- Scale up on load metrics (high cpu), gate scaledown on # messages in alive queue

```python
while True:
    try:
        conn = boto.sqs.connect_to_region('us-east-1')
        q = conn.get_queue(heartbeat_queue_name)
        m = Message()
        m.set_body("I am alive")
        q.write(m)
        logger.info("Sending heartbeat message")
        time.sleep(60)
    except Exception:  # pylint: disable=broad-except
        logger.exception("Unknown error occurred in heartbeat thread")
        time.sleep(5)
```

# Some recent spot autoscaling usage

# Questions