

Designing K-Means for the cloud

Matthieu Durut, Fabrice Rossi

March 6, 2012

Few reminders about distributed K-Means

$$\begin{aligned} \text{KMeans Sequential Time} &= (3Knd + Kn + Kd + nd) * I * T^{flop} \\ &\simeq 3Knd * I * T^{flop} \end{aligned}$$

$$\begin{aligned} \text{KMeans SMP Distributed Cost} &= T_P^{\text{comp}} \\ &= \frac{(3Knd + Kn + Kd + nd) * I * T^{\text{flop}}}{P} \\ &\approx \frac{3Knd * I * T^{\text{flop}}}{P} \end{aligned}$$

KMeans DMM Cost

$$\begin{aligned} &= T_P^{comp} + T_P^{comm} \\ &= \frac{(3Knd + Kn + Kd + nd) * I * T^{flop}}{P} + T_P^{comm} \\ &\approx \frac{3Knd * I * T^{flop}}{P} + O(\log(P)) \end{aligned}$$

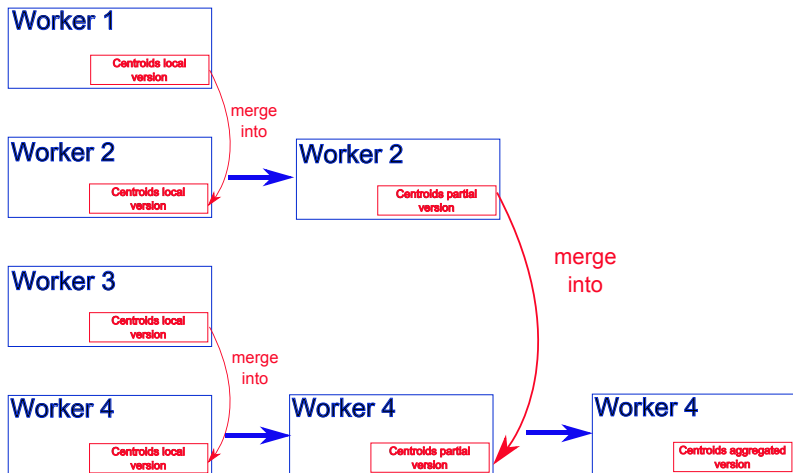


Figure: Merging centroids logic

The size of the biggest pipe (on top of the tree) and the aggregated bandwidth.

modélisation of T_P^{comm} in DMM

$$T_P^{comm} = \sum_{p=1}^{\lceil \log_2(P) \rceil} \left(2 \frac{IKd * SoD}{B_{\frac{P}{2^p}}^{DMM, MPI}} + 5IKd * T^{flop} \right)$$

modélisation de T_P^{comm} in DMM (2)

$$T_P^{comm} = \left(2 \frac{IKd * SoD}{B_{P,avg}^{DMM,MPI}} + 5IKd * T^{flop} \right) * \lceil \log_2(P) \rceil$$

$$\begin{aligned}
 \text{Speedup} &= \frac{\text{KMeans Sequential Cost}}{\text{KMeansDMMDistributedCost}} \\
 &= \frac{\text{KMeansSequentialCost}}{T_P^{\text{comp}} + T_P^{\text{comm}}} \\
 &= \frac{3nKdIT^{\text{flop}}}{\frac{3nKdIT^{\text{flop}}}{P} + \left(2\frac{IKd*SoD}{B_P^{\text{DMM,MPI}}} + 5IKdT^{\text{flop}}\right) \lceil \log_2(P) \rceil} \\
 &= \frac{3nT^{\text{flop}}}{\frac{3nT^{\text{flop}}}{P} + \left(2\frac{SoD}{B_P^{\text{DMM,MPI}}} + 5T^{\text{flop}}\right) \lceil \log_2(P) \rceil}
 \end{aligned}$$

- ▶ RAM is not modelize \Rightarrow speed up might be greater
- ▶ speedup independant of K, l and d
- ▶ winner paper model !

$$\lim_{n \rightarrow +\infty} SpeedUp = P \quad (1)$$

(2)

First case study : EDF

- ▶ $n = 20\,000\,000$ time series
- ▶ $d = 24 \text{ hours} * 365 \text{ days} * 10 \text{ years} = 87600$ values per series
- ▶ $K = \sqrt{n} = 4472$ clusters
- ▶ Note : volume of timeseries is 12.3 TeraBytes.

1 monocore

37 years

Yahoo terasort winner

- ▶ 910 nodes
- ▶ 4 dual core Xeons @ 2.0ghz per a node
- ▶ 1 gigabit ethernet on each node
- ▶ 40 nodes per a rack
- ▶ 8 gigabit ethernet uplinks from each rack to the core

Yahoo terasort winner(2)

$$\begin{aligned}T_P^{comp} &= \frac{(3nKd + nK + dK + nd) * I * T^{flop}}{P} \\ &= 322893 \text{seconds}\end{aligned}$$

$$\begin{aligned}T_P^{comm} &= (2 * \text{Broadcast time between 2 processors} + \text{Merging time}) * [\log_2(P)] \\ &= \left(2 \frac{IKd * SoD}{B_{P,avg}^{DMM,MPI}} + 5IKd * T^{flop}\right) * [\log_2(P)] \\ &= 225866 \text{seconds}\end{aligned}$$

Yahoo terasort winner(2) bounded by aggregated bandwidth

$$B_{P,avg}^{DMM,MPI} = 109 \text{ Mbits}$$

$$\begin{aligned} T_P^{comm} &= (2 * \text{Broadcast time between 2 processors} + \text{Merging time}) * [\log_2(P)] \\ &= \left(2 \frac{1Kd * SoD}{B_{P,avg}^{DMM,MPI}} + 51Kd * T^{flop}\right) * [\log_2(P)] \\ &= 434155 \text{ seconds} \end{aligned}$$

- ▶ Storage price : 1845 dollars/Month (0,15 dollars /GB/Month)
- ▶ No charge for internal transfers ?
- ▶ update : 10 GBytes/Month : 3 dollars (0.3dollars /GB)
- ▶ cost of running the KMeans : (very rough) 0.12dollars/hour *
324120 = 40 000dollars
- ▶ storage requests : ? from 0.01 dollars to 11 millions dollars

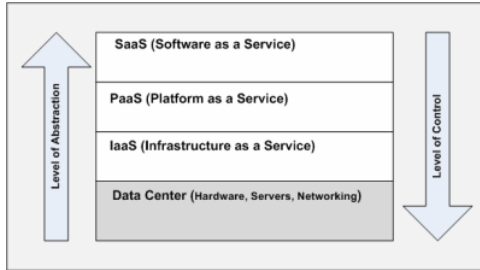


Figure: Paas, Iaas, SaaS

Paas and Iaas : some tradeoffs

- ▶ abstraction/control tradeoff
- ▶ scalability/cost to develop tradeoff

Introducing Lokad.Cloud

- ▶ higher level abstractions

- ▶ Azure Queues
- ▶ Azure BlobStorage
- ▶ Azure TableStorage
- ▶ Azure SQLStorage

- ▶ loosely coupled asynchronous communication mechanism
- ▶ not scalable yet (in size and number of messages)
- ▶ FIFO when not sought
- ▶ messages can be processed several times if worker fails (see idempotency)

- ▶ objects up to 50 GB
- ▶ consistency ?
- ▶ Key/value system

- ▶ Tables/Entities/Properties
- ▶ key/value pair also

- ▶ SOA
- ▶ Queue consumers

ACID/BASE constraints on Cloud

- ▶ ACID
- ▶ BASE
- ▶ CAP Theorem
- ▶ CLAP system

ACID/BASE constraints on Azure

- ▶ Availability : No
- ▶ Internet and Electricity neither
- ▶ Consistency : ?

Consistency on Azure : replicas and consistency

- ▶ Block commits
- ▶ Etag
- ▶ very few overwrites

- ▶ The SQL trick

- ▶ read/write tradeoff
- ▶ Responsivity/efficiency
- ▶ Synchronous replication/Asynchronous replication ?

BlobStorage or TableStorage ?

- ▶ Partition Key and locality
- ▶ size of objects
- ▶ efficiency of Blobstorage for objects larger than 4Kbytes
- ▶ Price

- ▶ BlobStorage Read : 13MBytes/sec/worker. Aggregated bounded to 393 MBytes/sec
- ▶ BlobStorage Write : 5MBytes/sec/worker. Aggregated bounded to 124 MBytes/sec
- ▶ Direct worker communication : 55MBytes/sec on very large messages
- ▶ VM : from 1 to 8 cores
- ▶ Queues throughput : 500 messages/sec

Redesign of the algorithm

- ▶ SOA
- ▶ No Master
- ▶ All workers can be switched
- ▶ No worker/data per affinity but ...
- ▶ Scaling up as an initiative in PaaS
- ▶ No exact workers count
- ▶ Pinging queues as a tradeoff between cost and simplicity
- ▶ Queue Messages count
- ▶ idempotent jobs

Redesign of the algorithm (2)

- ▶ Idempotent Counters
- ▶ Workers communication ?

Proposed Algorithm

- ▶ Scheme

Difficulty to anticipate speedUp

- ▶ Serialization/Deserialization.DataContract ? BinaryFormat ? GoogleProtoBuf ?
- ▶ variability of the bandwidth
- ▶ latency due to pinging queues
- ▶ a worker can fail, a VM can be killed
- ▶ Storage unavailability
- ▶ Stragglers issues ?

Is time to load data important ?

$$T_P^{Load} = \frac{nd * SoD}{P * \max(Bandwidth^{read}, \frac{MaxReadThroughput}{P})}$$

$$T_P^{Load} \ll T_P^{Comp}$$

Implicit supposition : loading is run once.

$$\frac{SoD}{3IK * T^{flop} * \max(Bandwidth^{read}, \frac{MaxReadThroughput}{P})} \ll 1$$

$$T_P^{comm,periteration} = \begin{aligned} & \text{Time to read 1 blob "simultaneously" by } P \text{ workers} \\ & + \text{Time to write } P \text{ blobs "simultaneously"} \\ & + \text{Time to read } P \text{ blobs by 1 worker} \\ & + \text{Time to merge the different centroids version} \\ & + \text{Time to write 1 blob per 1 worker} \end{aligned}$$

Communication modelisation(2)

$$\begin{aligned} T_P^{comm,periteration} = & \frac{Kd * SoD}{\text{Min}(\text{Bandwidth}^{read}, \frac{\text{BlobReadThroughput}}{P})} \\ & + \frac{Kd * SoD}{\text{Min}(\text{Bandwidth}^{write}, \frac{\text{MaxWriteThroughput}}{P})} \\ & + \frac{PKd * SoD}{\text{Bandwidth}^{read}} \\ & + (P - 1)Kd * 5\text{floatOperations} * T^{flop} \\ & + \frac{Kd * SoD}{\text{Bandwidth}^{write}} \end{aligned}$$

Communication modelisation(3)

$$T_P^{comm} \simeq \frac{IPKd * SoD}{Bandwidth^{read}}$$

$$\begin{aligned} \text{SpeedUp} &= \frac{T_1^{\text{comp}}}{T_P^{\text{Comp}} + T_{P\text{comm}}} \\ &\approx \frac{3IKdnT^{\text{flop}}}{\frac{3IKdnT^{\text{flop}}}{P} + \frac{IPKd*SoD}{\text{Bandwidth}^{\text{read}}}} \\ &\approx \frac{3nT^{\text{flop}}}{\frac{3nT^{\text{flop}}}{P} + \frac{P*SoD}{\text{Bandwidth}^{\text{read}}}} \end{aligned}$$

Optimal Number of workers

$$P^* = \sqrt{\frac{3n T^{flop} \text{Bandwidth}^{read}}{SoD}}$$

Numerical Example

- ▶ $n = 432\ 000$
- ▶ $K = 657$
- ▶ $d = 168$
- ▶ $l = 100$
- ▶ $P_{Max} = 64$
- ▶ $Bandwidth^{read} = 20\text{ Mbytes/sec (resp. }5\text{ Mbytes/sec)}$

Sequential KMeans time = 14300seconds

Optimal number of workers = 57(resp28)

SpeedUp = 28,5(resp14)

Improvements : more reducers ?

scheme

Direct worker communication ?