

Adaptive Geo-Replication Strategy

Amadeo Ascó
Trifork Leeds, Leeds, UK

ABSTRACT

The amount of data being processed in Data Centres (DCs) keeps growing at enormous rate so that full replication may start being impractical. The use of replication between DCs is used to increase data availability in the presence of site failures and to reduce access costs by accessing the data locally, if possible. This means that replicating the data only in some of the DCs is becoming more critical to reduce the costs of keeping the data consistent or eventually consistent and still maintaining a high availability (scalability) and low access costs. For changing patterns of read and write data requests the data to be replicated and its location has to be obtained dynamically. Given that the problem of finding an optimal replication schema in a general network has been shown to be NP-complete for the static case, it is unlikely to be able to generate a general algorithm to find efficient solutions to the dynamic problem.

So a new adaptive bio-inspired replication strategy is presented here, which is completely decentralised, adaptive, inspired on the Ant Colony algorithm, and event-driven. Also the replication protocol is independent of the strategy implemented but it is guided by the strategy.

Keywords

Adaptive Replication, Geo-replication, Multiple DCs, Large-Scale Database Replication, Partial Replication, Scalability

1. INTRODUCTION

The amount of data being processed in DCs keeps growing at enormous rate [21, 6, 5]. Some of the areas where the amount of stored data already reach terabytes (TBs) and even petabytes (PBs) are data mining, particle physics, climate modelling, high energy physics and astrophysics, to site few, data which needs to be shared and analysed [13, 16, 17]. DCs are able to ensure that stored data is highly accessible and scalable. But the location of a DC in respect of the

client accessing the data has an impact on availability, access times (latency - accessibility) and costs derived from providing the data. Replicating some of the data at multiple sites is a possible solution to reduce some of these undesirable effects, [4, 1, 22]. An increase in the number of replications may result in a large bandwidth savings and lead to a reduction in user response time on reads or writes depending on the replication type, i.e. replication on read or write. But keeping too many replicas of the data incurs in extra costs, such as extra replication traffic to keep all versions of the data coherent, extra required storage and extra computational power [20, 9].

This means that finding an optimal replication distribution that minimises the amount of network traffic given certain read and write frequencies for various objects should alleviate these extra costs when replicating, [25, 4]. Given the volume of operations considered, which are predominately higher in reads than writes, and the speed of the access expected then any algorithm suitable to be applied to this constrain optimisation problem must have a very fast execution time or it will have a detrimental effect in the latency.

The replications may be grouped, based on its existence, into two types; **static replication** where a replica persists until it is deleted by a user or its duration expires, and **dynamic replication** where the creation and deletion of a replica are managed automatically and normally directed by the access pattern of the data used by the users [7]. In static replication the major drawback is their inability to adapt to changes in the access pattern of the data used by the users.

Also there are two types of replication based on their effect on the data; **partial replication** is concerned in the number of parts the full data is composed, all of which may be located in different parts of the overall system, i.e. DCs, within a DC in different nodes or at the client-side [19, 4, 20], whereas **adaptive geo-replication** is concerned in what data and where the data or part of the data is located within the overall system of DCs and how many replicas exist simultaneously, [11, 13, 24, 1, 2, 15].

The problem of finding an optimal geo-replication schema in a general network has been shown to be NP-complete for the static case [3, 27, 26] so there is no known efficient geo-replication algorithm to locate a convergent optimal solution dynamically. Given this it is proposed an algorithm based on Wolfson's algorithm, [25], which proposes an adap-

tive algorithm for replicated data between processors which takes into account the changes in the read-write pattern of the processors in the network, and it is also based on the principles of the Ant Colony Optimisation algorithms, which are inspired in the behaviour of ant colonies when deciding which path to follow when foraging, [8].

It should be noted that the main propose of the replication is not to recover from disasters as this is the responsibility of the recovery centres which are data centres sufficiently close to the operational DC, they are associated to, so copies can be processed quickly enough but at the same time sufficiently far apart as to avoid any potential geographical issues that may happen to the DCs, i.e. earthquakes, fire that destroys the DC, the shutting down of main power plant which provides energy to the DC. Neither it is the DC responsibility to provide analytical services as these are provided by the data warehouse(s). The main propose of the considered DCs is to provide operational access to clients (operational DCs), which corresponds to intensive read/write operations to the client's most recent data. The following references to DCs in this document refer to operational data centres.

2. ALGORITHM

The general idea of this algorithm is to decide without the need of human intervention where and when to replicate with the main objective of reducing the latency and network traffic (reduce usage of bandwidth).

In general terms any read operation in a DC reinforces the need for a replica of the data in such DC, similarly but perhaps with a different degree it happens with the write operations, but write operations decrease the need for a replica of the data in the other DCs with replicas of the data, so eventually these DCs will not have any replica of the data. Given that we do not want to keep replicas if it is not necessary then the need for such replica will decay as time pass, but always making sure that the data is present (replicated) at least in one of the DCs. This algorithm is further explained below together with a mathematical representation. The variables and constants used in the algorithm are summarised in Table 1.

Equation 1 represents the existence of a replica of data k in DC d , $X_{dk} = 1$, or its absence $X_{dk} = 0$. R_k is the number of replicas of data k , as expressed in Equation 2.

$$X_{kd} = \begin{cases} 1 & \text{if } D_k \in DC_d \ (\exists D_{kd}) \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

$$R_k = \sum_{d=1}^{|DC|} X_{dk} \quad (2)$$

F_{kd} represents the full strength of the replication of the data k in DC d , which it is expressed in Equation 3.

$$F_{kd} = \max(0, \min(L_k, r_{kd} * \Delta r_k + w_{kd} * \Delta w_k - \sum_{i=1, i \neq d}^{|DC|} w_{ki} * \Delta w_{kdi} - X_{kd} * \Gamma)) \quad (3)$$

Equation 3 shows that the replication strength for the data k in DC d is strengthened by the reads and writes requested through DC d , with intensities Δr_k and Δw_k respectively, and weakened by the writes requested through other DCs than DC d , with intensity Δw_{kdi} , and it is furthermore weakened by a temporal decay on the replication strength (last term in the equation). Where Equation 3 states that only DCs with a replica of the data, D_k , will be penalised by writes and time, which in practice it is achieved by not sending the write to those DCs so it does not incur in any extra cost when implementing it. Similarly Γ will not be applied to DCs without replicas of the data with the exceptions of those candidate DCs where the data may later be replicated. Reads may increase the number of replicas where writes may strengthen the replication in a DC but may also potentially remove a replica from one of the other DCs, effect that it is strengthen by the temporal decay of the replication strength.

T_k^+ is the threshold of the replication strength of data k which determines when to keep an existing replica in a DC, as expressed in Equation 4.

$$F_{kd} \leq T_k^+, \nexists D_{kd}, t < t_0 \wedge F_{kd} > T_k^+, t_0 \implies \exists D_{kd}, t_0 \ (D_{kd} \in DC_d) \quad (4)$$

T_k^- is the threshold of the replication strength of data k which determines when to remove an existing replica in a DC when the other constraints are still kept. i.e. minimum number of replicas, as expressed in Equation 5. Also a simple view of the thresholds is shown in Figure 1.

$$F_{kd} > T_k^- \wedge \exists D_{kd}, t < t_0 \wedge F_{kd} \leq T_k^-, t_0 \implies \nexists D_{kd}, t_0 \ (D_{kd} \in DC_d) \quad (5)$$

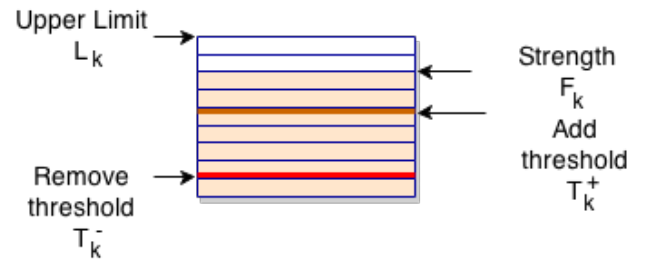


Figure 1: Constraints.

The Equations 6 and 7 show the cases where a data k is not replicated in a DC d . A replica in a DC, with a replication strength not bigger than the reduction replication strength T_k^- , will only be destroyed if the number of replicas is bigger than a minimum, N_k , as represented in Equation 6. Also for a DC which does not have a replica of the data the replication strength must be higher than a pre-set threshold T_k^+ to create a new replica, as represented in Equation 7.

$$\exists D_{kd}, t < t_0, R_k > N_k, F_{kd} > T_k^- \implies \nexists D_{kd}, t_0, F_{kd} \leq T_k^- \quad (6)$$

Variable	Description	Type
DC	It is the set of all DCs. d identifies one of the DCs, $d \in \{1, \dots, DC \}$. DC_d represents the DC d which holds some data, D_d .	$d \in \mathbb{N}^+$
D_{kd}	It represents the replica of data k in DC d , $k \in \{1, \dots, D \}$.	$k \in \mathbb{N}^+$
r_{kd}	It is the number of reads for data k requested on DC d .	$r_{kd} \in \mathbb{N}_0$
Δr_k	It represents the strengthening of the replication in the DC used to execute one read.	$\Delta r_k \in \mathbb{R}^+$
w_{kd}	It is the number of writes for data k directly requested on DC d .	$w_{kd} \in \mathbb{N}_0$
Δw_k	It represents the strengthening of the replication in the DC directly used to execute a write.	$\Delta w_k \in \mathbb{R}^+$
w_{kdi}	It represents the decay of the replication in the DC i consequence of the write request in DC d . This value will depend of both DCs d and i and may also depend on the time of the day or other useful information available at the time it is used.	$w_{kdi} \in \mathbb{R}^+$
Γ	It is the decay of the replication strength with time. A simple example corresponds to a constant decay (τ) of the replication strength from the time the replica was created, $\Gamma = \Delta t * \tau$.	$\Gamma \in \mathbb{R}^+$
N_k	Minimum number of replicas of data k with $1 \leq N_k \leq DC $, default $N_k = 1$. The number of DCs that can be safely removed from the system without information loss is $N_k - 1$.	$N_k \in \mathbb{N}^+$
T_k^+	It is the replication strength required to start the replication of data k in a DC which does not currently contain a replica of the data.	$T_k^- \in \mathbb{R}^+$
T_k^-	It is the replication strength from where the replication of the data is removed k in a DC which contains currently a replica of the data, default $T_k^- = 0$.	$T_k^- \in \mathbb{R}^+$
L_k	The maximum replication strength for data k .	$L_k \in \mathbb{R}^+$
F_k	The replication strength for data k , $F_k \leq L_k$.	$F_k \in \mathbb{R}^0$
R_k	It is the number of replicas for data k , D_k , with $1 \leq N_k \leq R_k \leq DC $.	$R_k \in \mathbb{N}^+$
X_{kd}	It refers to the existence of a replica of data k in DC d , with value 1 if the replicas exists or 0 otherwise.	$X_{kd} \in (0, 1)$

Table 1: List of variables and constants.

$$\nexists D_{kd}, t < t_0, F_{kd} \leq T_k^+ \implies \nexists D_{kd}, t_0, F_{kd} \leq T_k^+ \quad (7)$$

Each DC with a replica of the data must know about the other DCs that have replicas of the same data in order to manage the replication of such data.

If data k is only replicated in one DC d and a write is requested using a different DC j than the one it is correctly replicated in, so that its replication strength in DC d is reduced to zero or under ($F_{kd} < 0$), then the data will continue to be replicated in that DC until a replica is place in another DC or its replication is increased over zero.

For the time being, it is considered the case where w_{kdi} increases with the distance (network distance) between both DCs d and i . Also it is assumed that the value is symmetrical, such that $w_{kdi} = w_{kid}$, so there is the same cost of transferring the data from DC d to DC i than from DC i to DC d . If to transfer the data between both DCs d and i requires the use of an intermediate DC j then $w_{kdi} \geq w_{kdj} + w_{kji}$, similarly if many intermediate nodes are used the decay will be at least the sum of the intermediate decays.

The effect of Γ is required to ensure that in absence of writes some replications will still vanish and if the reads are concentrated in a few DCs then the replicas in other DCs will be eventually removed. There is an extra requirement in the case that the data only exists in a DC, in which case, the temporal effect should be ignored, so the data exists at least in one DC.

A read request to a DC, which does not have a replica of the

data, will be forwarded to the closest DC with a replica. The DC with a replica will not gain strength from this read operation as the read was not initiated directly from itself. This DC will then have a knowledge of the data but not a replica, so this knowledge will be used in subsequent reads/writes to the DC which eventually may keep a replica. Once the replication strength is higher than the threshold (T_k^+), this DC will notify all the DCs with a replica of the existence of the new DC with a replica.

It may also be required to use another temporal effect, the Time To Live (TTL), to make sure that eventually the data will be fully removed. This value may not be applied to data stored in recovery centres and data warehouses which may have their own TTL. Also it would be desirable for data that expires to be copied into those data centres before it is removed from all the DCs.

Overall it is assumed that reads and writes are not treated differently (not using Command Query Responsibility Segregation (CQRS)), so they are not directed to different DCs, otherwise it may be needed some adjustment to the approach presented here or may even invalidate it.

The algorithm is optimal in the sense that when the replication scheme stabilises, the total number of replicas required for the reads and writes is minimal.

The read sequence diagrams for this algorithm are shown in Table 2 with its flowchart presented in Figure 2, and the write ones are shown in Table 3 with its flowchart presented in Figure 3. The second figure in Table 2 may also have a notification for all DCs with a replica of the data to notify of the new replica which will be sent by the direct DC for

when the threshold has been passed and a replica is created in the direct DC.

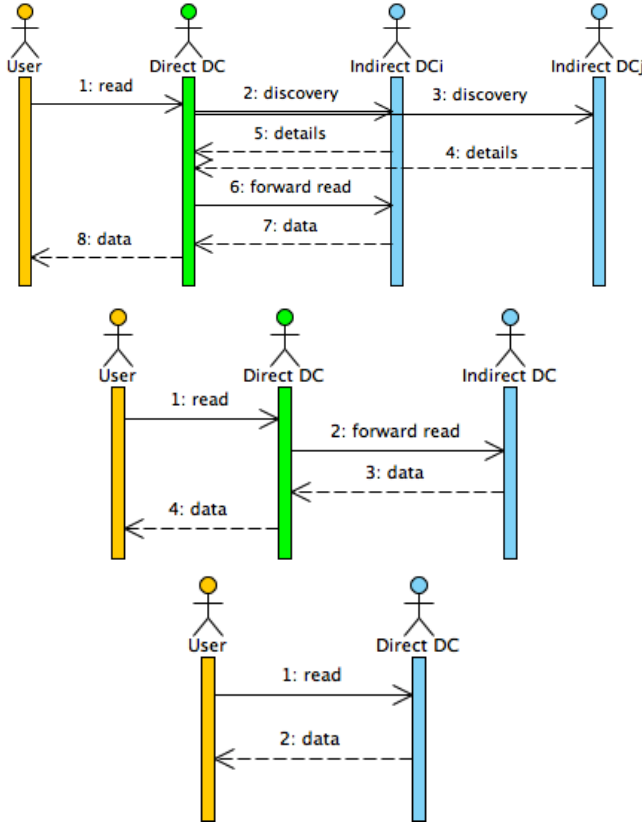


Table 2: Read Sequence Diagrams on write replication.

The creation of data will generate a replica in the directly accessed DC, as the number of replicas must be at least one ($M_k \geq 1$). If the minimum number of replicas is bigger than 1 then the algorithm should generate and distribute the other replicas between the vicinity DCs. It is proposed that the approach to distribute created replicas to be implemented as a pluggable distribution approach, i.e. distribute replicas within the vicinity. The creation of the data will require to check if the data exists already, i.e. in a key value tuple that the key does not already exist in any of the other DCs. Then only in the case where the data does not already exist in any of the DCs the data will be replicated in the direct accessed DC and distributed replicas within other of the DCs based on the current creation approach and the minimum number of replicas.

Some notes on parameters:

Given that the number of read is expected to be higher than the number of writes, and normally a read will be executed before a write, it may make sense for Δw_k to be smaller than Δr_k ($\Delta w_k < \Delta r_k$) so that more writes would be required to maintain or create a replica in a DC than when using reads.

Some of the parameters may be generalised even further by

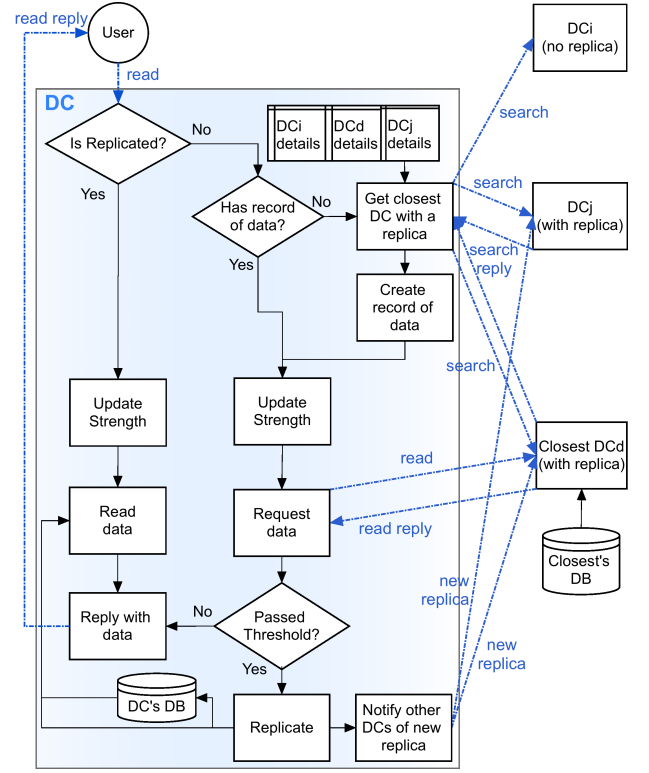


Figure 2: Flowchart for read requests.

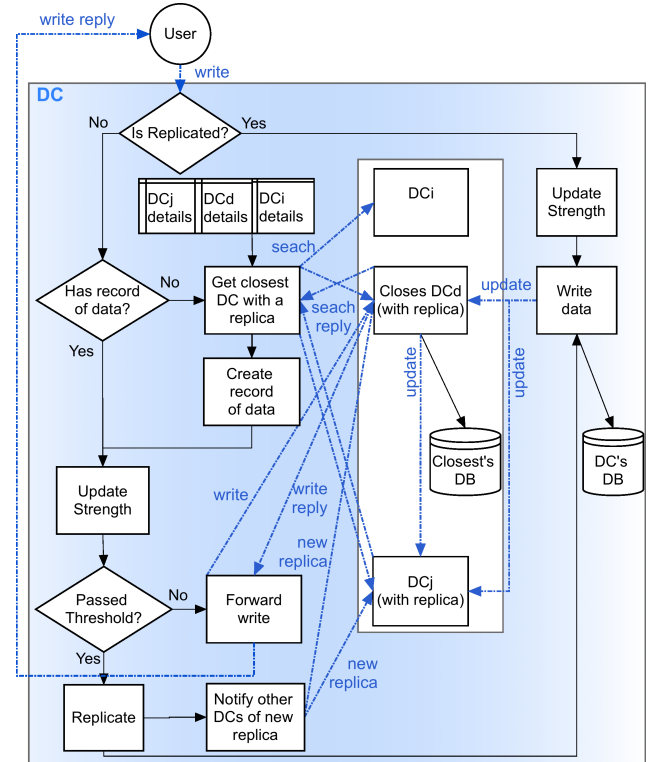


Figure 3: Flowchart for write requests.

allowing them to depend on the DC the calculation executed

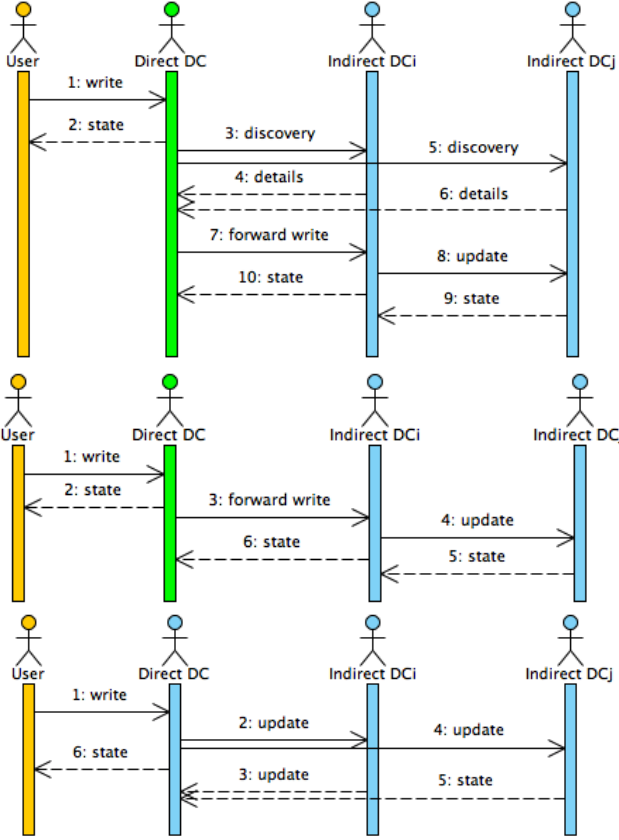


Table 3: Write Sequence Diagrams on write replication.

in, e.g. Δr_{kd} instead of Δr_k and Δw_{kd} instead of Δw_k . Even the terms $r_k * \Delta r_k$ and $w_k * \Delta w_k$ could be generalised to consider other factors like bandwidth used and current storage capacity available in the DC, or new terms could be added to take account of those new factors.

Paiva et al. [18] use a very simple representation of the available storage capacity, which is taking into account in the replication of the data. We have assumed that the capacity in a DC is sufficient and when more storage capacity is required then extra storage is added to the DC. If this is not the case and S^d is the total storage capacity in DC d , and S_k is the size of the data to store then Equation 8 provides the free storage available in DC d , s^d .

$$s^d = S^d - \sum_{k=1}^{|D|} S_k * X_{kd} \quad (8)$$

Accessibility in general is increased for Adaptive replication when compared to no replication. But also accessibility in general decreases for Adaptive replication when compared to full replication, given that when less replicas exist there is more chances that, on partition, a user is in the side of the partition where there is not any replica. Of course, this can be reduced by increasing the minimum number of replicas that must exist at each time and even adding extra con-

straints in relation to where those replicas must be in relation to each other.

2.1 Mathematical Cost Representation

It is considered that there is a direct cost associated to the reads (c_{kd}^r) and writes (c_{kd}^w) executed in the system for the data k in the DC d . Also there is a cost associate to the reads (c_{kdj}^r , $c_{kdd}^r = 0$) and writes (c_{kdj}^w , $c_{kdd}^w = 0$) between DCs d and j for data k , shown in Table 4.

The cost between time t and $t + dt$ is expressed by Equation 9 and the minimum number of replicas is represented by Inequality 10 for time t .

$$\begin{aligned} Cost_k(t, t + dt) = & \sum_{d=1}^{|DC|} \left(\underbrace{r_{kd}(t, t + dt) * c_{kd}^r(t) + w_{kd}(t, t + dt) * c_{kd}^w(t)}_{\text{direct cost}} \right. \\ & + \underbrace{\sum_{j=1}^{|DC|} w_{kd}(t, t + dt) * c_{kdj}^w(t) * X_{kj}(t)}_{\text{update propagation cost}} \\ & + \underbrace{(1 - X_{kd}(t)) * r_{kd}(t, t + dt) * \min_{j \in R_k} (c_{kdj}^r(t))}_{\text{redirect reads cost}} \\ & + \underbrace{(1 - X_{kd}(t)) * w_{kd}(t, t + dt) * \min_{j \in R_k} (c_{kdj}^w(t))}_{\text{redirect writes cost}} \\ & \left. + \underbrace{\sum_{j=1, j \neq d}^{|DC|} (1 - X_{kd}(t)) * X_{kj}(t) * \rho_{kd}(t, t + dt) * c_{*dj}^r(t)}_{\text{search cost}} \right) \quad (9) \end{aligned}$$

$$\sum_{d=1}^{|DC|} X_{kd}(t) \geq N_k \quad (10)$$

The cost expresses in Equation 9 corresponds to the cost of the operations for the reads and writes of data k in DC d , ‘direct cost’, plus the cost of propagating the updates, ‘update propagation cost’, plus the cost of redirecting the reads to the less costly DC with a replica, ‘redirect reads cost’, and finally plus the cost of searching for a DC with a replica of the data k to forward the request if the directly access DC does not have a replica of such data, ‘search cost’, from time t to time $t + dt$. The Equation 9 has a quadratic term that corresponds to the ‘search cost’. So the total time up to a time T is expressed in Equation 11.

$$Cost_k = \frac{\int_{t=0}^T Cost_k(t, t + dt) * dt}{T} \quad (11)$$

The algorithm proposed does not incur the full ‘search cost’ as only the first reply is processed and DCs with no replica of the data do not reply to the search request. Also all the reads received could be combined into one to the closes DC with a replica.

There is a cost incurred when a replica is placed in a new DC and when the data is removed from a DC which are

Cost	Description	Type
$X_{kd}(t)$	It refers to the existence of a replica of data k in DC d for time t , with value 1 if the replicas exists or 0 otherwise.	$\{0, 1\}$
$r_{kd}(t, t + dt)$	It is the number of reads for data k requested on DC d from time t to $t + dt$. It could be also represented as $r_{kd}(t, t + dt) = r_{kd}(t + dt) - r_{kd}(t)$, where $r_{kd}(t)$ is the number of reads from the beginning up to time t .	\mathbb{N}_0
$w_{kd}(t, t + dt)$	It is the number of writes for data k directly requested on DC d from time t to $t + dt$. It could be also represented as $w_{kd}(t, t + dt) = w_{kd}(t + dt) - w_{kd}(t)$, where $w_{kd}(t)$ is the number of writes from the beginning up to time t .	\mathbb{N}_0
$\rho_{kd}(t, t + dt)$	The existence of an operation from time t to $t + dt$, witch is 1 if $(r(t, t + dt)_{kd} + w(t, t + dt)_{kd}) > 0$ or zero otherwise.	$\{0, 1\}$
$c_{kd}^r(t)$	The cost of a read request of the data k to DC d for time t .	\mathbb{R}^+
$c_{kd}^w(t)$	The cost of a write request of the data k to DC d for time t .	\mathbb{R}^+
$c_{kdj}^r(t)$	The cost of a read request of data k from DC d to DC j , and $c_{kdd}^r = 0$ for time t .	\mathbb{R}^+
$c_{*dj}^r(t)$	The cost of a search request from DC d to DC j , and $c_{*dd}^r = 0$ for time t .	\mathbb{R}^+
$c_{kdj}^w(t)$	The cost of a read request of data k from DC d to DC j , and $c_{kdd}^w = 0$ for time t .	\mathbb{R}^+

Table 4: Costs and variables time dependent.

not present in Equation 9. If X'_{kd} is the new replication state after all the current operations have been executed with a cost per notification of c_{*kdj}^r then the cost of a new replica can be represented as in Equation 12, and the cost of removing a replica from a DC, with cost per notification of c_{*kdj}^r , can be represented as in Equation 13.

$$\sum_{j=1, X_{kj}(t)=1}^{|DC|} X'_{kd}(t) * (1 - X_{kd}(t)) * c_{*kdj}^r(t) \quad (12)$$

$$\sum_{j=1, j \neq d, X_{kj}(t)=1}^{|DC|} (1 - X'_{kd}(t)) * X_{kd}(t) * c_{*kdj}^r(t) \quad (13)$$

The extended cost is shown in Equation 14 where the time dependency has been omitted for simplicity and clarity.

$$\begin{aligned}
Cost_k(t, t + dt) = & \sum_{d=1}^{|DC|} \left(\underbrace{r_{kd} * c_{kd}^r + w_{kd} * c_{kd}^w}_{\text{direct cost}} + \underbrace{\sum_{j=1}^{|DC|} w_{kd} * c_{kdj}^w * X_{kj}}_{\text{update propagation cost}} \right. \\
& + \underbrace{(1 - X_{kd}) * r_{kd} * \min_{j \in R_k} (c_{kdj}^r)}_{\text{redirect reads cost}} \\
& + \underbrace{\sum_{j=1, j \neq d}^{|DC|} (1 - X_{kd}) * X_{kj} * c_{*dj}^r}_{\text{search cost}} \\
& + \underbrace{\sum_{j=1, X_{kj}=1}^{|DC|} X'_{kd} * (1 - X_{kd}) * c_{*kdj}^r}_{\text{new replica cost}} \\
& \left. + \underbrace{\sum_{j=1, j \neq d, X_{kj}=1}^{|DC|} (1 - X'_{kd}) * X_{kd} * c_{*kdj}^r}_{\text{remove replica cost}} \right) \quad (14)
\end{aligned}$$

Equation 14 could be used by another adaptive algorithm to set the parameters of this algorithm at running time in a dynamic way, instead of using fixed parameter values.

3. SOME COMPARISONS (DRAFT)

Here the following implementations are compared:

- Only one fixed replica, so no replication as such.
- Proposed algorithm, adaptive location of replicas.
- Full replication, so a replica in all DCs.

In the proposed algorithm, second implementation, by changing the minimum number of replicas we moved from the first case to the last one.

Table 5 shows the operations for each sequence diagram for the implementations considered. Looking at the storage requirements having replicates uses more storage than when the data is located only in one DC. Regarding the operations executed can be seen that for reads the existence of the data in many DCs reduces the total number of operation to execute as it is more likely that the accessed DC to get the data already have a replica of the data compared to when the data only exists in one DC. Whereas for writes having the data replicated in multiple DCs increases the total number of operations proportionally to the number of replicas. One technic to reduce the number of write executed from a DC with a replica to the other DCs with a replica is to combine multiple write in one before forwarding it to other DCs.

The algorithm proposed may behave as one fixed replica or as full replication by changing the value of some of its parameters, see below:

- **One fixed replica (A):**

$$N_k = 1, T_k^+ = \infty, T_k^- = -1, \Delta r_{kd} = 0, \Delta w_k = 0, \forall d \in |DC|.$$

Sequence	One fixed replica (A)	Algorithm (B)	Full replica (C)
Storage	1* data	$R_k * (\text{data} + \text{info}) + \Delta \text{info}$	$ DC * \text{data}$
First read without replica	$2 * \text{reads} + (DC - 1) * \text{discoveries}$	$2 * \text{reads} + (DC - 1) * \text{discoveries}$	
Read without replica	$2 * \text{reads} + (DC - 1) * \text{discoveries}$	$2 * \text{reads}$	$1 * \text{reads}$
Read with replica	$1 * \text{reads}$	$1 * \text{reads}$	
First write without replica	$2 * \text{writes} + (N_k - 1) * \text{discoveries}$	$(R_k + 1) * \text{writes} + (N_k - 1) * \text{discoveries}$	
Write without replica	$2 * \text{writes} + (N_k - 1) * \text{discoveries}$	$(R_k + 1) * \text{writes}$	$ DC * \text{writes}$
Write with replica	$1 * \text{write}$	$R_k * \text{writes}$	
Create data	$1 * \text{write}$	$N_k * \text{writes}$	$ DC * \text{writes}$

Table 5: Required operations by the different implementations for the specified sequences, $1 \leq N_k \leq R_k \leq |DC|$, on write replication.

• **Full replication (C):**

$$N_k = |DC|, T_k^+ = 0, T_k^- = -1, \Delta r_{kd} = 0, \Delta w_k = 0, \forall d \in |DC|.$$

This means that systems where the number of read is significantly higher than the number of writes will benefit from replicating the data in multiple DCs, whereas systems with similar or higher number of write will benefit from low or no replication (only one replica). This will be true for write replication, but if the replication (update) is conducted on the reads then increase in the number of replicas will be advantageous for writes and detrimental to reads. Future implementations of the proposed algorithm may take this into account so the replication type (on read or write) is also part of the adaptive mechanism.

There are also many more requirements that have an important influence on the selection of the system configuration, i.e. number of DCs and replication. Some of the most common requirements are Scalability, Accessibility, Latency and Security.

- **Scalability:** (A) does not provide scalability, where replication provides scalability. [12] provides an analytical study which shows the scalability limits of full replication as updates have to be sent and executed at all replicated sites (symmetric processing). To reduce this it can be used asymmetric processing where transactions are processed first at the originating site then collected and eventually propagated and applied to the other sites, which improves scalability. From the point of view of scalability the processing power in a DC, when a write is received, is invested in processing the write and the updates. As the number of replicas increases, there is a point at which the increase on the number of DCs, so replicas, does not increase any more the system capacity. The main reason is that most of the system processing power is used in processing the updates.
- **Availability/Accessibility:** (B) and (C) improve accessibility, where (A) provides a limited accessibility, [14].
- **Latency:** Given that (B) and (C) provide replicas in multiple DCs then the data can be accessed from any of those DCs and choosing the one with less latency will improve the latency seen by the customer, improve responsiveness.
- **Security/Fault-tolerance:** Security normally is achieved by using recovery centres but it is also improved/achievable by the provision of replicas (redundancy) as if a DC suffer a catastrophic event where all data is lost if the data was replicated then the data can be retrieved from the other DC(s) where the data is also replicated, [10, 23].

4. SUMMARY

This algorithm implies:

1. For reads on a DC, the DC does not require to communicate any information to any of the other DCs. A

read only needs to use resources in the DC, where the read is initially requested from, so no extra network traffic is imposed on the systems.

In the particular case where a read is executed on a DC, without a replica of the data, storage and processing power in the DC will be required and the request will be forwarded to its closest DC with a replica. This operation incurs in extra network traffic but if it is repeated too often the extra network traffic will be removed by replicating the data in the requesting DC where the data was initially requested.

2. On a write the DC receiving the original request will (eventually) transmit it to the other DCs, which have a replica of the data, so no need to add extra network traffic as this is the normal approach. But extra data will be sent to the other DCs to notify them of the number of writes the changes refer to, which will depend of the type of Conflict-free Replicated Data Type (CRDT) approach used, i.e. op-based or state-based. In some cases not every write is transmitted to the other DCs, such is the case of the state-base approach, so it would be required to keep some track of the number of merged writes.

Also the merging of the data should only be executed after the replication strength has been calculated and when it is still higher than zero, which will reduce unnecessary operations.

3. On data without any reads and writes on any DC, the number of replicas will be reduced as time passes by the temporal effect (Γ), until the data is only replicated in the specified minimum number of DCs.

This approach requires simple and fast operations to adapt to the changing access pattern of the data. Furthermore, some or all of the parameters could be determined at run time by some other adaptive approach which has into account the cost function provided in 2.1.

The replication type (on read or on write) has an important influence in the costs of replicating data and should be incorporated into the strategy described above to even improve more its performance.

5. ACKNOWLEDGEMENTS

This work has been supported by the project [SyncFree](#) (co-financed by the European Commission through the grant agreement number 609551).

References

- [1] Cristina L. Abad, Yi Lu, and Roy H. Campbell. Dare: Adaptive data replication for efficient cluster scheduling. In *Proceedings of the 2011 IEEE International Conference on Cluster Computing*, CLUSTER '11, pages 159–168, Washington, DC, USA, 2011. IEEE Computer Society. ISBN 978-0-7695-4516-5. doi: 10.1109/CLUSTER.2011.26. URL <https://wiki.engr.illinois.edu/download/attachments/194990492/cluster11.pdf>.
- [2] S. Abdul-Wahid, R. Andonie, J. Lemley, J. Schwing, and J. Widger. Adaptive distributed database replication through colonies of pogo ants. In *Parallel and Distributed Processing Symposium, 2007. IPDPS 2007. IEEE International*, pages 1–8, March 2007. doi: 10.1109/IPDPS.2007.370575. URL <http://www.cwu.edu/~andonie/MyPapers/IPDPS%20Long%20Beach%202007%20final.pdf>.
- [3] Peter M. G. Apers. Data allocation in distributed database systems. *ACM Transactions on Database Systems*, 13:263–304, 1988.
- [4] Iwan Briquemont. Optimising client-side geo-replication with partially replicated data structures. Master's thesis, Louvain-la-Neuve, September 2014. URL <http://www.info.ucl.ac.be/~pvr/MemoireIwanBriquemont.pdf>.
- [5] Aimee Chanthadavong. Internet of things to drive explosion of useful data: Emc. Technical report, ZDNet, April 2014. URL <http://www.zdnet.com/internet-of-things-to-drive-explosion-of-useful-data-en>.
- [6] Cisco. The zettabyte era-trends and analysis. Technical report, Cisco, June 2014. URL http://www.cisco.com/c/en/us/solutions/collateral/service-provider/visual-networking-index-vni/VNI_Hyperconnectivity_WP.pdf.
- [7] Xiaohua Dong, Ji Li, Zhongfu Wu, Dacheng Zhang, and Jie Xu. On dynamic replication strategies in data service grids. In *Object Oriented Real-Time Distributed Computing (ISORC), 2008 11th IEEE International Symposium on*, pages 155–161, May 2008. doi: 10.1109/ISORC.2008.66.
- [8] M. Dorigo. *Optimization, Learning and Natural Algorithms*. PhD thesis, Politecnico di Milano, Italy, 1992.
- [9] Sushant Goel and Rajkumar Buyya. Data replication strategies in wide area distributed systems. In Robin G. Qiu, editor, *Enterprise Service Computing: From Concept to Deployment*, pages 211–241. Idea Group Inc, 2006. URL <http://www.cloudbus.org/papers/DataReplicationInDSChapter2006.pdf>.
- [10] Rachid Guerraoui and André Schiper. Fault-tolerance by replication in distributed systems. In Alfred Strohmeier, editor, *Reliable Software Technologies – Ada Europe 96*, volume 1088 of *Lecture Notes in Computer Science*, pages 38–57. Springer Berlin Heidelberg, 1996. ISBN 978-3-540-61317-6. doi: 10.1007/BFb0013477. URL <http://dx.doi.org/10.1007/BFb0013477>.
- [11] Junsang Kim; Won Joo Lee; Changho Jeon. A priority based adaptive data replication strategy for hierarchical cluster grids. *International Journal of Multimedia & Ubiquitous Engineering*, 9(6):127–140, 2014. URL http://www.sersc.org/journals/IJMUE/vol9_no6_2014/13.pdf.
- [12] Ricardo Jiménez-Peris, M. Patiño Martínez, Gustavo Alonso, and Bettina Kemme. Are quorums an alternative for data replication? *ACM Trans. Database Syst.*, 28(3):257–294, September 2003. ISSN 0362-5915. doi: 10.1145/937598.937601. URL <http://lsd.ls.fi.upm.es/lsd/papers/2003/tods03.pdf>.

- [13] R. Kingsy Grace and R. Manimegalai. Dynamic replica placement and selection strategies in data grids- a comprehensive survey. *J. Parallel Distrib. Comput.*, 74(2):2099–2108, February 2013. ISSN 0743-7315. doi: 10.1016/j.jpdc.2013.10.009. URL <http://dx.doi.org/10.1016/j.jpdc.2013.10.009>.
- [14] Rivka Ladin, Barbara Liskov, Liuba Shrira, and Sanjay Ghemawat. Providing high availability using lazy replication. *ACM Trans. Comput. Syst.*, 10(4):360–391, November 1992. ISSN 0734-2071. doi: 10.1145/138873.138877. URL <http://doi.acm.org/10.1145/138873.138877>.
- [15] Thanasis Loukopoulos and Ishfaq Ahmad. Static and adaptive distributed data replication using genetic algorithms. *J. Parallel Distrib. Comput.*, 64(11):1270–1285, November 2004. ISSN 0743-7315. doi: 10.1016/j.jpdc.2004.04.005. URL http://pdf.aminer.org/000/297/337/static_and_adaptive_data_replication_algorithms_for_fast_information_access.pdf.
- [16] Noriyani Mohd. Zin, A. Noraziah, AinulAzila Che Fauzi, and Tutut Herawan. Replication techniques in data grid environments. In Jeng-Shyang Pan, Shyi-Ming Chen, and NgocThanh Nguyen, editors, *Intelligent Information and Database Systems*, volume 7197 of *Lecture Notes in Computer Science*, pages 549–559. Springer Berlin Heidelberg, 2012. ISBN 978-3-642-28489-2. doi: 10.1007/978-3-642-28490-8_57. URL http://dx.doi.org/10.1007/978-3-642-28490-8_57.
- [17] Shaik Naseera and K.V. Madhu Murthy. Agent based replica placement in a data grid environment. *Computational Intelligence, Communication Systems and Networks, International Conference on*, 0:426–430, 2009. doi: <http://doi.ieeecomputersociety.org/10.1109/CICSYN.2009.77>.
- [18] João Paiva, Pedro Ruivo, Paolo Romano, and Luís Rodrigues. Autoplacer: scalable self-tuning data placement in distributed key-value stores. In *Proceedings of the 10th International Conference on Autonomic Computing, ICAC’13*, San Jose, CA, USA, June 2013. USENIX. URL <https://www.usenix.org/conference/icac13/technical-sessions/presentation/paiva>.
- [19] *Conflict-free Partially Replicated Data Types*, 2015. PPOPP.
- [20] D. Serrano, M. Patino-Martinez, R. Jimenez-Peris, and B. Kemme. Boosting database replication scalability through partial replication and 1-copy-snapshot-isolation. In *Dependable Computing, 2007. PRDC 2007. 13th Pacific Rim International Symposium on*, pages 290–297, Dec 2007. doi: 10.1109/PRDC.2007.39. URL http://www.researchgate.net/publication/200023090_Boosting_Database_Replication_Scalability_through_Partial_Replication_and_1-Copy-Snapshot-Isolation/links/0deec520a3cdf6504e000000.
- [21] K.M. Tolle, D. Tansley, and A.J.G. Hey. The fourth paradigm: Data-intensive scientific discovery [point of view]. *Proceedings of the IEEE*, 99(8):1334–1337, Aug 2011. ISSN 0018-9219. doi: 10.1109/JPROC.2011.2155130. URL <http://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=5958175>.
- [22] Srikumar Venugopal, Rajkumar Buyya, and Kotagiri Ramamohanarao. A taxonomy of data grids for distributed data sharing, management, and processing. *ACM Comput. Surv.*, 38(1), June 2006. ISSN 0360-0300. doi: 10.1145/1132952.1132955. URL <http://www.cloudbus.org/reports/DataGridTaxonomy.pdf>.
- [23] J. von Neumann. *Automata Studies*, chapter Probabilistic Logics and Synthesis of Reliable Organisms from Unreliable Components, pages 43–98. Princeton University Press, 1956. URL https://ece.uwaterloo.ca/~ssundara/courses/prob_logics.pdf.
- [24] Zhe Wang, Tao Li, Naixue Xiong, and Yi Pan. A novel dynamic network data replication scheme based on historical access record and proactive deletion. *J. Supercomput.*, 62(1):227–250, October 2012. ISSN 0920-8542. doi: 10.1007/s11227-011-0708-z. URL <http://dx.doi.org/10.1007/s11227-011-0708-z>.
- [25] Ouri Wolfson. A distributed algorithm for adaptive replication of data. Technical report, Department of Computer Science, Columbia University, 1990. URL <http://hdl.handle.net/10022/AC:P:21285>.
- [26] Ouri Wolfson and Amir Milo. The multicast policy and its relationship to replicated data placement. *ACM Trans. Database Syst.*, 16(1):181–205, March 1991. ISSN 0362-5915. doi: 10.1145/103140.103146. URL <http://academiccommons.columbia.edu/catalog/ac/3A142996>.
- [27] Ouri Wolfson, Sushil Jajodia, and Yixiu Huang. An adaptive data replication algorithm. *ACM Trans. Database Syst.*, 22(2):255–314, June 1997. ISSN 0362-5915. doi: 10.1145/249978.249982. URL http://www.cs.uic.edu/~wolfson/mobile_ps/tods-adaptive-replication.pdf.

Glossary

CQRS Command Query Responsibility Segregation

CRDT Conflict-free Replicated Data Type

DC Data Centre

PB petabytes

TB terabyte

TTL Time To Live