

分散ファイル・システムにおけるオンデマンド・レプリケーションによる  
ファイルアクセス効率の改善について

芦原 評、 清水 謙多郎、 前川 守、 浜野 純  
東京大学理学部情報科学科

ファイル配置問題は多様に研究されてきたが、過去の多くの解はシステムの完全な知識と集中制御に基づくもので分散システム上での実現には不適であり、それゆえ極めて特殊な用途か、理論的研究の範囲を出なかった。本論文では、我々がダイナミック・レプリカと名付けた動的・部分的なレプリカファイルとその生成を決定する分散的アルゴリズムを示した。決定は各ノードにより動的、自律的に行われ、局所的判断のみによるにもかかわらず、例によっては過去のアルゴリズムによる最適解にまさる場合さえあることが示された。

ダイナミック・レプリカの概念の定義の後、ダイナミック・レプリカを生成、維持するための体系を検討し、性能改善のための条件を求める。その結果多くの現実的場面においてダイナミックレプリカがシステムの性能を改善することが示された。

"File Access Improvements in Distributed Systems  
By On-Demand Replication of Files" (in Japanese)  
by Hyo ASHIHARA, Kentaro SHIMIZU, Mamoru MAEKAWA  
and Jun HAMANO (Department of Information Science,  
Faculty of Science, University of Tokyo, 7-3-1  
Bunkyo-ku, Tokyo, 153, Japan)

Technical Report 87-26, Department of Information  
Science, University of Tokyo (in English)

## 1. はじめに

分散ファイル・システムにおいて、ファイル配置の決定は理論・実際の両面で極めて重要な問題である。

与えられた条件のもとで、最適配置、すなわち、平均応答時間あるいは通信コストを最小にするような、分散システム内の各ノードへのファイルの割り当て（もしくはその近似）を求める必要がある。その配置は、システムの設計時、あるいは個々のファイルの生成時などに外部から指定することもできるが、ユーザーがシステムの物理的構造にとらわれず、すべてのファイルがローカルに存在するかのように扱える（透過的）システムを、効率を犠牲にすることなく構成するためには、システム自身が最適配置を決定できることが望ましい。

このとき、同一ファイルのレプリカ（コピー）を複数のノードに置くことがありうる。レプリカはシステムの信頼性を確保すると同時に、性能の向上のためにも重要である。しかしながら、コピーファイルを許した場合のファイル配置問題はNP困難であることが知られている。

レプリケーションの特殊な形と考えられるものとして、キャッシング、すなわちプロセスによる遠隔ファイルの一時的取込みがある。本論文では、リモートファイルを必要に応じてローカルにコピーする機構（ダイナミック・レプリカと呼ぶ）を提案し、その効果を解析する。ダイナミック・レプリカは、最初ページテーブルのみが存在し、実体のレプリケーションはページごとにオンデマンドで行われる。

### 1.1 分散システム

分散ファイル・システム、又は分散データベース・システムについての多くの研究の中で、たとえばSDD-1 [1]やLOCUS [2]は信頼性、性能の双方を指向したレプリカを扱うものであるが、信頼性・可用性についての面をより強調している。性能向上を主目的としてもいくつかの方式が提供されており、ITC [13]ではリモート・アクセスに際して常にファイル全体をキャッシュするという方式をとっている。また読み出し専用のファイルのレプリケーションも行う。CFS [14]も、読み出し専用ファイルへのリモート・アクセスに限ってキャッシングを行う。

このほか分散ファイル・システムに関する最近の試みについてはSvobodova [19]やTanenbaum [20]に概観されているが、ファイルのレプリケーションをサポートしたものは比較的少なく、特にレプリカの生成、存在とその位置に関する透過性を実現したものは筆者の知る限りない。これは、レプリケーションに際して何等かの人間の判断が要求されることを意味する。

### 1.2 レプリケーションと性能

レプリケーションによる性能向上は、原理的には、読み出しの時に最も近いレプリカにアクセスすればよいかわり、書き込みの際にはすべてのレプリカを更新しなければならぬとすれば、そのファイルに対する読み出しが更新に比して十分多いとき期待できる。ファイルが読み出し専用の場合は更新に関連した複雑な制御機構が不要なため特に有利である。

その定量的な評価としてはCasey [2]やMuro [10]のものがある。これらは一般のファイル配置問題の一要素としてレプリカの存在を考慮したものである。

Caseyは読み出しが更新より多い場合コスト上レプリケーションが有利となりうるが、完全結合ネットワークの各ノードにおいて読み出しに対する更新の割合が一定ならば、その逆数がオリジナルを除くレプリカの数の限界であることを示す一方、ARPAネットワークのデータを使用した例によりこの比が0.4の場合のレプリケーションを含む最適配置を示している。Fisherも同じくARPAの例で、更新比率0.4でレプリケーションが有利となる結果を得ている [6]。さらにFisherの結果の一つは、レプリカが少なすぎるよりは多すぎる方がよいことを

も示唆している。

一方Muroは乱数を用いて生成したネットワークモデル上の実験により、細かい条件をレプリケーションに有利に単純化した場合でも、最適配置にレプリカが含まれるのは更新比率0.2程度以下の場合に限り、ほとんどの場合多数のレプリカの配置は現実的でないとして述べている。しかしながらMuroのモデルは平面上のノードの分布を一律、各ノードからのファイル参照の頻度をすべて等しいと仮定したもので、これが両者の結論の差の主たる原因と思われる。

### 1.3 ファイル配置

ファイル配置問題についても、Chu [3]以来多くの研究がなされている [5, 21]が、考慮すべき要素があまりにも多く、一般性が現実性的一方を犠牲にすることになりがちである。複雑な条件をすべて組み入れた場合、ケース・スタディの形をとらざるを得ない。しかも理論的にはきわめて単純なモデルでさえNP困難であるため、現実のシステムにおけるファイル配置決定のため、さまざまな「実用的」アルゴリズムが考案されている。この際、レプリケーションは特例的であって現実の有効性は低く、問題を徒らに複雑化すると、考慮から外されることも多い。

たとえば山崎はレプリケーションを認めなくともファイル相互の関係によるコストの変化を考えた場合同じくNP困難であるとしている [25]。

これら過去に行われた研究はネットワークをグローバルに見ての議論が多い。ネットワーク全域にわたるデータの交換と集中的な配置決定は困難であり、実用に供するファイル・アクセス決定のアルゴリズムは効率的かつ分散向きでなければならぬ。これらはネットワーク・トポロジや各ノードのファイルへのアクセス頻度を既知の固定したものとしているが、実際には、ハードウェアによる制限など、さまざまな「現実的」条件を考慮しつつ、ダイナミックな変動に追随しなければならぬ。そして、ファイル配置をそのように変更する場合、再配置自体に要するコストも考慮しなくてはならない。単一のサービスのためでなく、公共、グループ、個人を含むさまざまなユーザーのためのさまざまな目的を持つ大システムの場合これらは顕著であって、Muroの仮定とは逆に、ノードごとのアクセス頻度に大きな偏りがありうる場合が重要である。レプリケーションは、アクセス頻度による制限など、容易にわかれる。特にほとんどのファイルが（結果的に）ローカルなアクセスしかされない（しかし無視できない例外は常に存在する）ようなシステムにおいても、アクセスの効率と配置決定の効率を両立させねばならない。

他方一般に資源管理・負荷分散と呼ばれる問題も概念上同じ目的を持つが、逆に往々にして議論がプロセス単位に限定され、また与えられた資源と要求をマッチさせるという形を取るためレプリケーションは無視されがちである。Kuroseはレプリケーションを認めないか、ただかレプリカの数があらかじめ与えられた場合を考察している [8]。

また、1つのファイルの部分によっても条件の偏りがあり、性能上レプリケーションを考えるとき常にファイル全体をレプリケートすることは無駄、時には有害になる。

### 1.4 ダイナミック・レプリケーション

本論文においては、リモートファイルに対するアクセスが必然的にページ転送を伴うタイプのシステムを対象に、ダイナミックにページ単位で生成、消滅するレプリカ（ダイナミック・レプリカと呼ぶ）を考えることにより、この問題の一つの実用的な解決を計る。

すなわちダイナミック・レプリカ（以下DR）は生成時には構造のみが存在し、実体はページ単位でアクセス要求があるご

とに転送されてレプリカにとどまり、一定時間アクセスがなければ消滅する。これはキャッシュの拡張と見ることもできるが、キャッシュと異なりプロセス間で共有することができる。(図1)

各ノードの立場に立った場合、ファイルへのリモートアクセスはレプリカ生成の要求であると考え、その点にレプリカを作ることが作れない場合より有利かどうかを判断するという形になる。

この方式は現在開発中の分散オペレーティング・システム G A L A X Y [9]において採用したものである。

その有効性の議論が本稿の目的である。任意のファイル、ノードにおいて、DRを作るか否かの0-1問題として、単純化されたアクセスモデルに基づき通信コストの評価を行い、いくつかの場合についてその可否を論ずる。

## 2. モデル

ネットワーク上のノードに、PページのオリジナルファイルFが存在するとする。オリジナルの位置を移動させることはないとし、以後、このノードとFを同一視する。このFを、距離L(=ページあたりの伝送コスト、一定とする)の位置D付近にあるノードまたはノード群D'のプロセスからアクセスする。このとき、Dにレプリカを作ることによりコストが改善されるかどうかを検討する。

DとD'の距離はいずれもL'であり、原則としてL' < Lである。そのためD'とFの距離はLとする。(図1)

各プロセスはキャッシュあるいはバッファを持ち得るが、それらはプロセスの一部とみなす。外部から観測されるのは、プロセスからのページ要求、すなわちプロセスが必要とする(読み出し、書き込みを問わず)ファイルのページがキャッシュになかった場合の要求であり、これが単位時間にfの割合で起こる。またファイルへのページの書き戻しがwの割合で起こる。

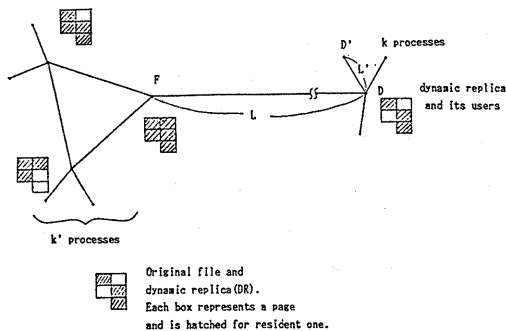


図1 ダイナミック・レプリカ概念図

箱はファイルと各ページを、斜線部はそのうちDRに実在するものを示す

### 2.1 変数

以下の変数を定義する。

- L : FとD、およびFとD'の距離
- L' : DとD'の距離 (L' < L)
- f : 単位時間中のページ要求
- w : 単位時間中の書き込み要求
- t : 時間
- k : D'においてファイルにアクセスするプロセスの数
- k' : 同じファイルにアクセスする上記以外のプロセスの数
- P : ファイルのページ数
- v(p) : ページ要求において、第pページが要求される確率
- v\_w(p) : 書き込みにおいて、第pページに書き込まれる確率
- a(p, t) : ページpが時刻tまでにアクセスされる確率
- S(t) : 時刻tにおいてDRに実在するページ数

### 2.2 レプリカのない場合

k個のプロセスがそれぞれ時間あたりfの比率でFからページを読み出し、wの比率で書き戻す。t時間の間の総通信コストは

$$C_1(t) = k(f+w)tL$$

となる。

### 2.3 レプリカのある場合

次にDにDRを置く。(DとDRを同一視する)これを使用するプロセス(群)D'からの距離はすべてL'とする。D=D'ならL'=0としてよい。Dに実在しないページが要求されるとFから読み込みを行い、そのままDにとどまる。(図1) ページpが確率v(p)でランダムにアクセスされると仮定する。

以後はpを連続変数とみなすと

$$\int_0^P v(p) dp = 1 \quad (0 \leq p \leq P)$$

一度アクセスされたページはDにとどまるので、ファイルの第pページが時刻tまでにアクセスされている確率

$$a(p, t) = 1 - \text{EXP}(-v(p)kft)$$

より、

時刻tにおいてDに存在するページの総計は

$$S(t) = \int_0^P a(p, t) dp$$

$$S(t) \leq kft \quad (t \geq 0)$$

は直観的にも明らか。

[例]

1)

v(p) = 1/P (すなわち一様)と仮定すれば

$$S(t) = P(1 - \text{EXP}(-fkt/P))$$

この場合を以後一様アクセスと呼ぶ。固定したtに対し一様アクセスのときS(t)は最大となる。

$$v(p) = \begin{cases} 2/P & 0 \leq p \leq P/2 \\ 0 & P/2 < p \leq P \end{cases}$$

(全ページの半分だけが実際にアクセスされる)

$$S(t) = (P/2)(1 - \text{EXP}(-2fk t/P))$$

$$3) \quad v(p) = (2/P^2)p \quad (\text{ページ数に比例})$$

$$S(t) = P - (P^2/2fk t)(1 - \text{EXP}(-2fk t/P))$$

(図2)

通信コストは

$$C_2(t) = S(t)(L+L') + kwt(L+L') + (fkt - S(t))L'$$

$L' = 0$  (ローカル)の場合を考えると

$$C_3(t) = S(t)L + kwtL \leq C_1$$

ファイルのアクセス特性に関する研究[11,12,15,18,23,24]は、ここでのいくつかの仮定はやや単純化がすぎること示唆しているが、我々の目的には十分な近似を与える。以下の議論において、一般性を保持するため、 $t \rightarrow \infty$ での長期的条件と、一様アクセスなどの代表的ケースを論ずるにとどめ、関数 $S(t)$ の正確な形に依存することは極力避けたが、個々のファイルの特性が予め知られていればもの小さな状況ではより現実的な判定が可能になる。SheltzerはOSのネーム・サービスに限ったアクセス特性とキャッシングの効率を論じている[16]。

一般に $v(p)$ は未知であり、 $f, w$ はファイル全体にわたる平均であって、ページごとの相当する値の変動を含む。実際には、ファイルがクラスタに分割され、同一ファイル内で性質の異なる部分に分かれる場合、クラスタごとにレプリケーションを行う。(すなわち、クラスタ単位で構造が、ページ単位で実体が生成される)クラスタのサイズについては、実現上のオーバーヘッドとの兼ね合いになる。部分的に読み出し専用が指定されている( $\rho = \infty$ )場合は特に重要である。逆に、特に書き込みの多いことがわかっているページにレプリケーションを禁止するフラグを立てるなども考えられる。

$v(p)$ に偏りがあると、部分的にレプリケーションが不利なページが現われるが、全体としては一様のときよりよくなることに注意する。

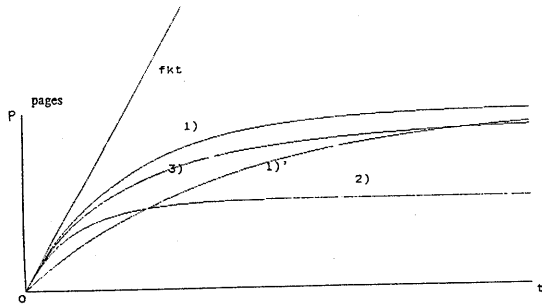


図2 【例】の3種類のアクセス特性による $S(t)$ の変化

1) ' は  $k f$  が半分の場合の一様アクセスを示す

## 2.4 更新

次にFが更新される場合を考える。F側にも同様のプロセスが $k'$ あるとする。(他のレプリカの利用者も含む) それらのプロセスの更新アクセスが第pページに落ちる確率を $v_w(p)$ とすれば、D上に当該ページが存在するときそれを書き換える通信コストは

$$\begin{aligned} & (\text{転送されるページ数}) \\ & = (\text{全更新の頻度}) \times \\ & \int \int (\text{時刻 } t \text{ にページ } p \text{ が更新される確率}) \\ & \quad \times (\text{ } t \text{ に } p \text{ が } D \text{ に存在する確率}) dt dp \end{aligned}$$

$$C_4(t) = k' w \int_0^t \int_0^P v_w(p) a(p, t) dt dp L$$

$$= k' wt - \frac{k' w}{kf} \int_0^t \frac{v_w(p)}{v(p)} (1 - e^{-v(p)kf}) dp L$$

$$= (k' w / kf) (fkt - S^*(t)) L \quad \text{とおく。}$$

$$S^*(t) = \int (v_w(p) / v(p)) a(p, t) dp$$

$$v_w(p) = v(p) \text{ と仮定すれば } S^*(t) = S(t)$$

$$\begin{aligned} & \text{一様アクセスの場合} \\ & v_w \text{ によらず、同じく} \\ & S^*(t) = S(t) \end{aligned}$$

$$(\therefore \int v_w dp = 1)$$

これは一般には保証されないが、最悪の場合でも $S^*(t) > 0$ であるから、この性質が以下の議論において致命的になることはない。

## 2.5 利得

ここで、DRの有効性について考察する。コストの差は、

$$\begin{aligned} \Delta_1 &= C_1 - (C_2 + C_4) \\ &= \left( \frac{k'w}{kf} S^*(t) - S(t) \right) L - \left\{ wL' + fL' + \left( \frac{k'w}{kf} \right) fL - fL \right\} kt \\ & \quad k' / k = \kappa, f / w = \rho, L' / L = \lambda \text{ とおくと、第二項は} \\ & \quad -(\lambda + \rho \lambda + \kappa - \rho) kw t L \end{aligned}$$

ここから、

$$\begin{aligned} (1 + \kappa) &< (1 + \rho)(1 - \lambda) \quad \dots\dots\dots(1) \\ \text{または} \\ (w / (w + f)) &< (k / (k + k')) ((L - L')) \end{aligned}$$

の時、 $t \rightarrow \infty$ でレプリカが有利になる。(長期条件) 特に、 $\alpha = \kappa / \rho$  とおけば、 $\lambda = 0$  ( $L' = 0$ ) では

$$\alpha < 1$$

これらのパラメータ中、直感的に最も理解しにくいのは  $\rho$  であるが、これは更新比率の逆数で、ファイルの性質によって決まる定数と見てよい。ただし、機構上書き込みの際にもキャッシングが行われることから、全アクセスに対する更新比率に近い意味を持ち、 $\rho > 1$  が期待されるが、保証はされない。

さらに  $\alpha = \kappa / \rho$  は D の読み出しに対する F の書き換えの比率と見ることができるから、上の結果は極めて当然のことである。

## 2.6 遷移

$t$  が小さく  $S(t)$  が卓越する場面について、いくつかの特別な場合を論ずる。これらは現実にもしばしば重要である。

$\lambda = 0$ 、 $S^*(t) \geq S(t)$  では、 $\kappa / \rho < 1$  のとき  $t$  によらず  $\Delta_1 > 0$ 。

(1) が満たされるとして、 $t = t_0$  において  $\Delta_1 = 0$  とおくと、

$$S^*(t) = S(t) \text{ の場合} \\ T_0 = (\rho + 1)\lambda / (\rho - \kappa)$$

$$(T(t) = 1 - S(t)) / \kappa f t, T_0 = T(t_0) \\ 0 < T < 1, t \rightarrow \infty \text{ で } T \rightarrow 1$$

また時刻  $t$  において DR の存在が有利である条件は  $(1 + \kappa) < (1 + \rho)(1 - \lambda / T)$  ……(1t)

$$S^*(t) = 0 \text{ の場合} \\ (1 + \kappa / T) < (1 + \rho)(1 - \lambda / T) \text{ ……}(1t')$$

逆に  $S^* > S$  であれば (1t) より良くなる。

$\rho = \infty (w = 0)$  であれば、(1) は常に成立し、 $S^*$  によらず

$$T_0 = \lambda$$

[例]  
一様アクセスでは  
 $\lambda = 1 - P(1 - \text{EXP}(-f \kappa t_0 / P)) / f \kappa t_0$

$T(P / \kappa f) \approx 0$ 、3.7 であるから、 $\lambda < 0$ 、3.7 であれば、P ページが伝送されるまでには DR が有利になっている。他の  $v(p)$  においてはこれより良くなる。

$\kappa = 0 (\kappa' = 0)$  では、同様に

$$T_0 = (1 + 1 / \rho)\lambda$$

F 側にプロセスがない ( $\kappa = 0$ ) の場合、安全上の配慮を別とすれば、DR の更新をただちに F に書き戻す必要はない。このケースでは、最終的に書き戻すべきページ数

$$S_w(t) = \int (1 - \text{EXP}(-v_w \kappa w t)) dp \text{ とおくと} \\ S_w(t) \leq \kappa w t$$

$$\Delta_1 + \kappa w t L - S_w(t)$$

$$= -S(t)L - S_w(t)L - (wL' + fL' - fL - wL)\kappa t \\ = -S(t)L - S_w(t)L + (1 + \rho)(1 - \lambda)\kappa w t L$$

長期条件は常に成立し、 $t$  との関係も  $\rho = \infty$  の場合とほぼ同じになる。

## 2.7 拡張

ここまで、D 側に  $k$  個、F 側に  $k'$  個のプロセスが独立・並列に走っていると見てきた。あるノードで全く同じプロセスが同時にあるいは繰り返し使われる場合、アクセスされるページはほとんど同じであり、理想的には更新分を除いて新たな転送は不要になると考えられる。

別な見方として、すべてのプロセスについては  $v(p) = 0 \rightarrow v_w(p) = 0$  であるが、D 側と F 側のプロセスが違えばこの制限は必要でない。

これらを念頭に、  
 $Q = \lim_{t \rightarrow \infty} \int v_w(p) a(p, t) dp$   
すなわち  $v(p) > 0$  となる  $p$  についての  $v_w(p)$  の積分 ( $0 \leq Q \leq 1$ ) とおく。

$$\text{DR の飽和以後、} \\ \Delta_2 = \kappa f t L - Q \kappa' w t L - \kappa w t L' - f \kappa t L' \\ = -(\lambda + \rho \lambda + Q \kappa - \rho) \kappa w t L$$

$\kappa$  を  $Q \kappa$  におきかえたものとなる。

$$\text{条件は、} \\ (1 + Q \kappa) < (1 + \rho)(1 - \lambda) \text{ ……}(2)$$

であって、当然 (1) より有利である。

次に、ハードウェアの制限その他により、レプリカの容量を制限した場合を考える。

現実にはリプレースメント・アルゴリズム等考える必要が生ずる。

1)  $S^*(t) = S(t)$  の場合

$t = t_1$  以降 D に実在するページ数が増えないとすると、以後のアクセスに関して

$$\Delta_3 = (d \Delta_1 / dt) |_{t=t_1} \cdot t \\ = \left[ \left( \frac{\kappa' w}{\kappa f} - 1 \right) S'(t_1) L - \left\{ wL' + fL' + \left( \frac{\kappa' w}{\kappa f} - 1 \right) fL \right\} \kappa \right] \cdot t \\ = R f \kappa t L - w \kappa t L' - f \kappa t L' - R w \kappa' t L \\ = (R \rho - R \kappa - \lambda - \rho \lambda) \kappa w t L$$

$$(R \equiv 1 - S'(t_1) / \kappa f = \int v(p) a(p, t_1) dp, \\ 0 < R \leq 1, \text{ 単調増加})$$

$$\text{(一様アクセスでは } R = S(t_1) / P \text{)}$$

$\Delta_3 > 0$  の条件は

$$(1 + \kappa) < (1 + \rho)(1 - \lambda / R) \text{ ……}(3)$$

すなわち (1) に対し  $\lambda$  を  $\lambda / R$  で置き換えた結果となる。 $\lambda < R$  である必要があるのをはじめ一般には生成の条件・生成後の効率共に不利になるが、 $\lambda = 0$  の場合は条件には影響を受けない。

2) 一般の場合

$S(t) \neq S^*(t)$  でありうるから、

$$\Delta_3 = (R \rho - R^* \kappa - \lambda - \rho \lambda) \kappa w t L$$

$$\text{条件は} \\ (1 + R^* \kappa / R) < (1 + \rho)(1 - \lambda / R) \text{ ……}(3^*)$$

$$(R^* \equiv Q - S^*(t_1) / \kappa f)$$

となって、複雑である。

$v$ 、 $v_w$ が未知であるから、あらかじめ対応することは困難であるが、前出のクラスタリングによって問題をかなり避けることができる。

【例】

$$v(p) = \begin{cases} \frac{2}{3P} \\ \frac{4}{3P} \end{cases} \quad v_w(p) = \begin{cases} \frac{2}{P} & (0 \leq p < \frac{P}{2}) \\ 0 & (\frac{P}{2} \leq p \leq P), \end{cases}$$

なる場合

$$S(t) = P \left( 1 - \frac{e^{-\frac{2kt}{3P}}}{2} - \frac{e^{-\frac{4kt}{3P}}}{2} \right)$$

$$S^*(t) = \frac{3}{2} P \left( 1 - e^{-\frac{2kt}{3P}} \right)$$

( $T \equiv f k t$ )

$$1 - R = \text{EXP}(-2T/3P)/3 + 2 \text{EXP}(-4T/3P)/3$$

$$1 - R^* = \text{EXP}(-2T/3P) \equiv X$$

$$\therefore R = 1 - X/3 - 2X^2/3$$

$$R^* = 1 - X$$

$\lambda = 0$ であれば

$$(3^*): \quad \alpha = \frac{\kappa}{\rho} < \frac{R}{R^*} = \frac{2X}{3} + 1$$

すなわち、 $1 < \alpha < 5/3$ のとき、全体ではレプリケーション不可だが、容量を

$$X = 3(\alpha - 1)/2$$

$$S(t)/P = 1 - X/2 - X^2/2$$

$$= -(9\alpha^2 - 12\alpha - 5)/8$$

までに制限した場合は有利になる。最善となるのは

$$X = (3\alpha - 1)/4$$

$$S/P = -(9\alpha^2 + 6\alpha - 35)/32$$

のとき。

$\alpha = 4/3$ ならば容量約  $1/3$  の時最善、 $5/8$  以上では不利になる。

定性的には、 $v$ の比較的小さい  $p$  において  $v_w/v$  が大きいと、このような逆転現象が生ずる。逆に、 $v$ の大きい時  $v_w/v$  が大きいと、容量の制限は(3)よりさらに不利となる。

## 2.8 ページの無効化

以上の方式では、インプリメントを考えた場合、オリジナル側が、DRの位置と、存在するページに関する情報を把握しておかねばならないという欠点がある。

信頼性を目的とする「通常の」レプリカの場合、原則として全レプリカを同時に書き換えて整合性を保つが、DRはキャッシュに近い「軽い」存在であるから必ずしも同じ方法を執る必要はない。しかしこの場合でも更新同士の衝突を避けるほかプロセス間通信などによる干渉に対し矛盾が生じないようにしなければならない。

Fで書き換えられたページが無効化される場合を一種の理想的状态として考えると、同一ページに続けて書き込みが行われた時に無駄な転送をしなくてよい。またオリジナル側がDRを無視できるという特長があり、ページごとのオンデマンド転送という概念との整合性もよい。

キャッシュに対しても同種の機構は必要であるが、性能上ここでは問題としない。(なおこれはいわゆるコンカレンシ・コントロールとは別の問題である)[7]

各ページが時刻  $t$  に  $D$  に存在する確率  $a(t)$  について

$$\frac{da}{dt} = (1-a)kfv - ak'vw_w$$

( $v_w = v$  と仮定)  $= kfv - (kf + k'w)va$   
これを解いて

$$a(t) = \frac{kf}{kf + k'w} (1 - e^{-(kf + k'w)vt})$$

ページがアクセスされた回数の期待値を  $b$  とおくと

$$b = \int_0^t (1-a)vkf dt$$

$$= \frac{kk'fwv}{kf + k'w} t + \left( \frac{kf}{kf + k'w} \right)^2 (1 - e^{-(kf + k'w)vt})$$

時刻  $t$  までにアクセスされるページの総和

$$S_2(t) = \int_0^t b d p$$

$$= \frac{kk'fw}{kf + k'w} t + \left( \frac{kf}{kf + k'w} \right)^2 S \left( \frac{kf + k'w}{kf} t \right)$$

【例】

一様アクセスの時

$$S_2(t) = \frac{kk'fw}{kf + k'w} t + \left( \frac{kf}{kf + k'w} \right)^2 P \left( 1 - e^{-\frac{kf + k'w}{P} t} \right)$$

$$C_6(t) = S_2(t)L + kfL' + kwL(L+L')$$

$$= \frac{kk'fw}{kf + k'w} tL + \left( \frac{kf}{kf + k'w} \right)^2 S \left( \frac{kf + k'w}{kf} t \right) L + kfL' + kwL(L+L')$$

$$\Delta_4 = C_1 - C_6$$

$$= kft(L-L') - S_2(t)L - kwL'$$

$$= - \left( \frac{\rho}{\rho + \kappa} \right)^2 S \left( \frac{\rho + \kappa}{\rho} t \right) L + \left( \rho - \frac{\kappa\rho}{\rho + \kappa} - \rho\lambda - \lambda \right) kwL$$

より、 $t \rightarrow \infty$  で  $\Delta_4 > 0$  となるのは

$$(1 + \rho)(\kappa + \rho)\lambda < \rho^2 \quad \dots\dots(4)$$

特に  $\lambda = 0$  であれば  $\kappa$ 、 $\rho$ 、 $t$  によらず常にレプリカが存在が有利になる。

D側の書き込みの際もFのみ更新してDを無効化する(ライトスルー)方式も考えられる。定常状態で

$$\Delta_4(\kappa = 1) + kwL' - \frac{kk'fw}{kf + k'w} tL > 0$$

$$(1 + \rho)(\kappa + \rho)\lambda < \rho^2 - \kappa \quad \dots\dots(5)$$

(4)より悪いが、 $\kappa > 1$  では(1)より良くなる。

以上の二方式を同時に使用すると、 $\Delta I + k w L' > 0$ において  $\kappa \rightarrow \kappa + 1$  と置き換えれば  $(1 + \rho + \kappa)\lambda < \rho$  あるいは  $(1 + \kappa) < (1 + \rho + \kappa)(1 - \lambda)$  ……(6) であって、(1)、(4)より有利になる。

(4)、(5)、(6)の条件によって  $v = v_0$  が最悪の仮定であることは容易にわかる。ライトスルー方式では、書き込みの多いページが無効化されることから、一般にはさらに有効である。

### 3. 応用

機構としてのレプリケーションを論じた以上の結果は、そのまま実際のシステムにおいて個々のレプリカの生成を決定するアルゴリズムにつながる。 $\kappa$ 、 $\rho$ などのパラメータについては、システムがある程度自動的にデータを収集することができる。過去のデータを近い将来のアクセス頻度の予測として用いることにより、レプリカ生成の可否を決定することができる。また逆に使用されていないレプリカは削除する。さらに長期にわたる使用の「実績」によってDRと通常のレプリカの間の昇・降格を行うことも考えられる。これによって、「見えざる手」に要求の変動に自動的に対応することができる。

この議論の結果形成されるファイル配置がグローバルな意味で最適に近いとは限らないが、二点間のみの問題として各ノードが自律的に考えることができるため、分散システム上での意志決定は単純で効率的である。他のレプリカやネットワークのグローバルな構造についてノードが知る必要はない。(オリジナル点に関する情報に含まれる)。逆に言えばノードごとのローカルな条件を直ちに反映できる。ローカルに生成され、ローカルに使用されるようなファイルについてははじめから問題としないですむ。ただしシステムを理解しているユーザーが明示的に指定することを妨げるものではない。

プロセス終了後のレプリカの寿命の設定は将来の研究の対象である。これは非分散システムでのページング、キャッシング[4,17]、またファイルマイグレーションとも類似の問題を含む。もちろん、より深い分析と、きめ細かいデータによれば、それだけ効果的な判定ができるはずだが、それは個々のシステムの特性にも大きく依存したものととなる。

実現時にまず問題となるのは、実際には仮定の場合ほどはつきりリモートグループが分離できないことである。そのため、まず  $\lambda = 0$  の場合に限って具体的な方法の例を示す。ノードRにおいて、ある程度以上速方のファイルFを呼び出す必要が初めて生じたとき、Rに小さなファイルdを作り、以後のFに対するアクセスの総数( $fkt$ 相当)を記録する。一方Fには、すべてのノードからの読み出し、書き込みそれぞれの要求の総和が記録される。その、他のノードからの全要求に対する自己の要求が2、において求めた基準を満たせばその次のオープン要求の際Rにdにかわりダイナミック・レプリカDを生成する。

ネットワーク内の近い距離にあるいくつかのノードがグループ(物理的には、たとえばLAN)をなしている場合、グループ全体として上記の手続きを行うこともできる。その際、こうして一旦生成されたDRをグループ外から使用することを許すか、またそれをレプリカ維持のためのデータにどう算入するか、などの問題が生ずる。慎重に行わないと、状況が変化した後もDRが惰性的に存続して不利益を生ずる可能性がある。より大きな入に対しては適用可能な実現法の開発が求められる。

### 4. 例

すでに述べたように、本方式では、“それ以前より改善される”ことだけが保証され、最適配置にどの程度近づくかはわからない。この点を検証するため、Caseyの用いたネットワークのデータ(詳細は[2]を参照)により試算を行う。(1)式の判定で、query/update比をそのまま  $\rho$  に使用。ページごとの効果は含まれていない。

#### 4.1 5ノード

(図3)は、5つのノードからなるネットワークの各ノード間の距離及び各ノードからのファイル参照頻度を示している。ここで、ファイルを単一のノードに配置するならばノード4または5が最善であり、その時のコストは915である。

レプリカを認めた場合の最適配置は(1, 4, 5)で、コストは705。

上のアルゴリズムを各ノードについて適用すると、すべてに条件が成立する。従って、ノード5にオリジナル、1、2、3、4にDRを置くと、コストは759。

同じ全ノードにレプリカを置く場合でも、Caseyの示す完全結合のコスト(771)と異なる(オリジナルが適した位置に置かれていれば、通常改善される)ことに注意。これは本方式の持つもう一つの利点であって、次の例ではさらに劇的に示される。

#### 4.2 ARPA

更新比率0.1のとき、オリジナルの配置は10が最適(コスト  $2.4 \times 10^5$ )。Caseyによる最適配置は(2, 10, 14)で、コストは  $1.2 \times 10^5$ 。

2-5、13-15、17-19をそれぞれグループとしてまとめて計算すると、2-5グループ、12、13-15グループが条件を満たす(更新比率が0.25以上であれば、この方式ではDRは作られない)。2、12、14にそれぞれDRを生成し、DRはグループ外から使わないとして、コスト  $1.6 \times 10^5$ 。各ノードが常に最も近いレプリカにアクセスするならば、コストはCaseyの最適解より良い  $1.1 \times 10^5$  となる。このことから、複数ノードへのメッセージの最適伝達経路と、それをも考慮した最適配置という新しい問題が導かれる。これについては、他の論文でさらに検討する予定である。

INPUT PARAMETERS		COST PER MEGABYTE SHIPPED														
QUERY UPDATE FILE		TRAFFIC														
U	R	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	1	24	2	4	8	8	12	9	6	12	9	6	6	12	9	6
0	2	24	4	12	6	8	12	6	12	6	12	6	12	6	12	6
0	3	24	4	8	8	8	8	8	8	8	8	8	8	8	8	8
0	4	24	4	8	8	8	8	8	8	8	8	8	8	8	8	8
0	5	24	4	8	8	8	8	8	8	8	8	8	8	8	8	8

(a) The five-node example

QUERY FILE		COST PER MEGABYTE SHIPPED														
U	R	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	1	0	75	75	75	75	75	75	75	75	75	75	75	75	75	75
0	2	0	75	75	75	75	75	75	75	75	75	75	75	75	75	75
0	3	0	75	75	75	75	75	75	75	75	75	75	75	75	75	75
0	4	0	75	75	75	75	75	75	75	75	75	75	75	75	75	75
0	5	0	75	75	75	75	75	75	75	75	75	75	75	75	75	75

(b) The ARPA example

図3 ネットワークの例 (Casey[2]による)

ノード間の通信コストと通信量を示す

## 5. おわりに

以上で、オリジナル側の更新に際して余分な通信を必要とする1)方式であっても、 $\kappa$ に比して $\rho$ が大きい(特に、 $\kappa \ll 1$  (maybe = 0)、あるいは $w \ll 1$ )場合、ダイナミックレプリケーションが効率を高めることが示された。またレプリカの整合性に関する制約を緩めればさらに効率を上げることができる。

効率を論じるに際し問題となる要素には、通信コストの他、ストレージコスト、計算コスト、反応時間などがある。実際の運用にあたってはこれらを考慮に入れねばならないが、個々のノードにおけるストレージコストやストレージ制限を含めた計算は容易である。

通信時間についても、距離の定義を変更すれば上記の結果はほぼ妥当する。ただし、通信コストの場合、リモート/ローカル、LAN内/遠距離のアクセスについては十分小さいと見てよいが、時間においてはレプリカを経由することによる遅れは無視できないことが多い。またコンカレンシコントロールに要する時間も大きく、実現時のメカニズム、アルゴリズムに強く依存する。

ファイル配置問題は多様に研究されてきたが、過去の多くの解はシステムの完全な知識と集中制御に基づくもので分散システム上での実現には不適であり、それゆえ極めて特殊な用途か、理論的研究の範囲を出なかった。本論文では、我々がダイナミック・レプリカと名付けた動的・部分的なレプリカファイルとその生成を決定する分散的アルゴリズムを示した。決定は各ノードにより動的、自律的に行われ、局所的判断のみであるにもかかわらず、例によっては過去のアルゴリズムによる最適解にまさる場合さえあることが示された。

## 謝辞

本論文をまとめるにあたり、熱心な討論と多くの貴重な助言をいただいたGALAXY計画のメンバーに深く感謝する。

## 参考文献

- [1] Bernstein, P. A., et al. : "The concurrency control mechanism of SDD-1: a system for distributed databases (the fully redundant case)", IEEE Trans. Softw. Eng., Vol. SE-4, No. 3, (May 1978).
- [2] Casey, R. G. : "Allocation of copies of a file in an information network", in Proc. AFIPS Spring Joint Comput. Conf., (1972).
- [3] Chu, W. W. : "Optimal file allocation in a multiple computer system" IEEE Trans. on Comput. Vol. C-18, NO. 10 (Oct. 1969).
- [4] Denning, P. J. : "Virtual memory" Comput. Surv., Vol. 2 No. 3 (Sep. 1970).
- [5] Dowdy, L. W., D. V. Foster : "Comparative models of the file assignment problem" Comput. Surv. Vol. 14, No. 2 (June 1982).
- [6] Fisher, M. L., et al. : "Database location in computer networks" JACH. Vol. 27, No. 4, (Oct. 1980).
- [7] Hamano, J. : Master's Thesis, To appear.
- [8] Kurose, J. F. : "A microeconomic approach to optimal file allocation" in Proc. of the 6th Intern. Conf. on Distrib. Comput. Sys., IEEE, (May 1986).
- [9] Maekawa, M., et al. : "Object Management and Address Space of GALAXY Holonic Processing System", Tech. Rep. 86 15, Dept. of Information Science, Univ. of Tokyo (Oct. 1986).
- [10] Muro, S., et al. : "Evaluation of the file redundancy in distributed database systems", IEEE Trans. Softw. Eng., Vol. SE-11 No. 2, (Feb. 1985).
- [11] Ousterhout, J. K., et al. : "A Trace-Driven Analysis of the UNIX 4.2 BSD File System", in Proc. of the 10th Symp. on Oper. Sys. Principles. (1985).
- [12] Rodriguez-Rosell, J. : "Empirical Data Reference Behavior in Data Base Systems" IEEE Computer (Nov. 1987)
- [13] Satyanarayanan, M., et al. : "The ITC Distributed File System: Principles and Design" in Proc. of the 10th Symp. on Oper. Sys. Principles. (1985).
- [14] Schroeder, M. D., et al. : "A Caching File System For a Programmer's Workstation" in Proc. of the 10th Symp. on Oper. Sys. Principles. (1985).
- [15] Shedler, G. S., C. Tung : "Locality in Page Reference Strings" SIAM J. Comput., Vol. 1, No. 3 (Sept. 1972).
- [16] Sheltzer, A. B., et al. : "Name Service Locality and Cache Design in a Distributed Operating System" in Proc. of the 6th Intern. Conf. on Distrib. Comp. Sys., IEEE, (May 1986).
- [17] Smith, A. J. : "Analysis of Long Term File Reference Patterns for Application to File Migration Algorithms", IEEE Trans. on Softw. Eng., Vol. SE-7, No. 4, (July 1981).
- [18] Smith, A. J. : "Cache Memories", Comput. Surv., Vol. 14, No. 3, (Sep. 1982).
- [19] Svobodova, L. : "File Servers for Network-based Distributed Systems", Comput. Surv., Vol. 16, No. 4, (Dec. 1984).
- [20] Tanenbaum, A. S., and R. van Renesse : "Distributed Operating Systems", Comput. Surv., Vol. 17, No. 4, (Dec. 1985).
- [21] Wah, B. W. : "File Placement on Distributed Computer Systems", IEEE Computer, (Jan. 1984).
- [22] Walker, B., et al. : "The LOCUS Distributed Operation System", ACM SIGOPS Oper. Sys. rev. 17, 5, (Oct. 1983).
- [23] 畑下 豊仁 他 : UNIXワークステーションにおけるディスク・アクセス特性とディスク・キャッシュの考察 : 情報処理学会論文誌 Vol. 28, No. 6, (June 1987).
- [24] 菅 隆志 他 : ディスク・キャッシュ装置のシミュレーションによる効果測定 : 情報処理学会論文誌 Vol. 22, No. 1, (Jan. 1981).
- [25] 山崎 浩 他 : 分散型データベースシステムにおけるファイル配置問題の計算の複雑さについて 電子情報通信学会論文誌 D Vol. J70-D No. 4, (Apr. 1987).