

クラスタ分割されたネットワークのノード故障対策について

5X-4

石田 賢治 宮尾 淳一 菊野 亨 吉田 典可

広島大学

1. まえがき

本稿では、クラスタ分割されたコンピュータネットワーク[1]-[2]上のコントロールノードにおける故障対策の中で、特に障害回復に関する基礎的考察を行う。

2. クラスタ分割されたネットワーク

コンピュータネットワークの大規模化に伴い、ネットワークの効率良い管理は益々困難になりつつある。1つの解決策に、ネットワークをクラスタに分割し、クラスタごとに管理する方法がある(図1参照)。

各クラスタ i にはコントロールノード C_i と幾つかのボーダノードが存在する。クラスタ i 内のノードは C_i によって集中的に管理される。但し、各クラスタ間は分散的に管理さる。

これらのネットワーク管理に必要な情報は(1)広域的なネットワーク情報を含む表 TG と、(2)クラスタ i に固有なネットワーク情報を含む表 TL_i として保持される。表の保持方法に関して次の仮定をおく。

H1 クラスタ i のコントロールノード C_i は表 $TG \cup TL_i$ を保持している。

H2 クラスタ i の各ボーダノードは表 TL_i を保持している。

3. 故障回復

最近、ソフトウェアの誤りやオペレータの誤動作が問題となりつつある[3]。ここでは、コントロールノード上の表 $T = TG \cup TL_i$ が消失した場合を想定し、表の回復方法について議論する。

今の場合、他のコントロールノードから TG を転送し、クラスタ i 内のボーダノード

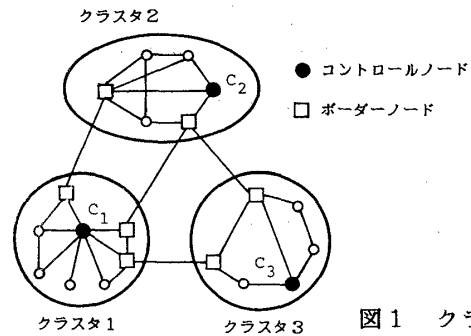


図1 クラスタ分割

から TL_i を転送すれば、 c_i では表 $TG \cup TL_i$ を再構成できる。ここでは、 TG (あるいは TL_i)それぞれ自体も複数個の表から構成されていると仮定する。

4. 問題の定式化

以降、ネットワーク上のノード v_0 の表 T が消失し、ノード v_1, v_2, \dots, v_n 上に(コピーとして)重複して置かれている T を利用して回復するものとする。

ノード v_0 と v_i の間のリンクを l_i で表し、 l_i に関し次のパラメータを導入する。

$r(l_i) \dots l_i$ を設定するのに要するコスト(時間)

$d(l_i) \dots$ 単位パケットを l_i 上で転送するのに要するコスト

更に、表 T を $T = \{t_1, t_2, \dots, t_m\}$ と仮定する。ここで、各 t_j はいずれも表である。各 t_j に対し、次のパラメータを導入する。

$s(t_j) \dots$ 表 t_j の転送に必要なパケットの総数

ここでは、各ノード v_i への表(の集合)の割当て $A(v_i) (i \in 2^T)$ を決定する問題として定式化する。すなわち、各 v_i から v_0

に表 $A(v_i)$ をリンク l_i 上を転送することで、 v_0 上に表 T を再構成する。ここでは、 A を評価するため、次の評価関数 $D(A)$ を導入する。

$$D(A) = \max_{1 \leq i \leq n} \{ d(l_i) \cdot \sum_{t_j \in A(v_i)} s(t_j) + r(l_i) \}$$

次に、回復問題 RT の定式化を与える。

[回復問題 RT] 次の①～④が与えられたとき、評価関数 $D(A)$ の値が最小となる割当て A を求めよ。

- ① $L = \{ l_i \mid 1 \leq i \leq n \}$
- ② $r(l_i), d(l_i) \in Z^+$
- ③ $T = \{ t_j \mid 1 \leq j \leq m \}$
- ④ $s(t_j) \in Z^+$

[例 1] 回復問題 RT の例として次の①～④について考えてみる。

- ① $L = \{ l_1, l_2, l_3, l_4 \}$
- ② $r(l_1) = 5, r(l_2) = 6,$
 $r(l_3) = 1, r(l_4) = 3,$
 $d(l_1) = 3, d(l_2) = 4,$
 $d(l_3) = 1, d(l_4) = 2.$
- ③ $T = \{ t_1, t_2, t_3, t_4, t_5 \}$
- ④ $s(t_1) = 9, s(t_2) = 10,$
 $s(t_3) = 8, s(t_4) = 12,$
 $s(t_5) = 20.$

これに対し最適解として $A = (\{ t_2 \}, \phi, \{ t_4, t_5 \}, \{ t_1, t_3 \})$ が求まる。このとき $D(A) = 37$ となる。

5. ヒューリスティック解法

回復問題 RT の時間計算複雑さに関し、次の命題が導ける。

[命題 1] 回復問題 RT は NP -困難である。

命題 1 より、回復問題 RT に対するヒューリスティック解法の開発が重要になる。以降では、簡単化のため、 $d(l_1) < d(l_j)$, $j \neq 1$ と仮定する。更に、 $A(v_i)$ を系列として考え、例えば $A(v_i) = (t_s, t_{s+1}, \dots, t_p)$ と表す。この系列に対し、先頭の要素を $TOP(A(v_i))$ 、要素 t_s を削除する操作を $A(v_i) - t_s$ 、要素 t_s を系列の最後に追加する操作を $A(v_i) + t_s$ で表す。

[アルゴリズム RT]

ステップ 1 (初期解) ... 集合 T 上の系列 $\bar{T} = (t_{i1}, t_{i2}, \dots, t_{im})$ を構成する。但し、 $s(t_{ij}) \leq s(t_{ik})$, $j < k$ とする。次に、初期解として $A = (\bar{T}, \lambda, \dots, \lambda)$ とおく。

ステップ 2 (逐次改良) ... この時点での解を $A' = (a'_1, a'_2, \dots, a'_n)$ とする。

2.1 (候補の決定) 関数 $H(j) = r(l_j) + d(l_j) * (\sum_{t \in A'(v_j)} s(t) + s(TOP(a'_1)))$, $2 \leq j \leq n$ を計算し、最小の値を与える添字を 1 つ選ぶ(今、それを j^* とする)。

2.2 (解の更新) 新しく $A = (a'_1 - TOP(a'_1), a'_2, \dots, a_{j^*} + TOP(a'_1), \dots, a'_n)$ を解とする。

2.3 (判定) $D(A) < D(A')$ ならステップ 2 を繰り返す、それ以外なら A' を最終解として停止する。

[例 2] 例 1 と同じ問題に対し、上述のアルゴリズム RT を適用すると、解として $A = (\{ t_1 \}, \phi, \{ t_4, t_5 \}, \{ t_2, t_3 \})$ が求まる。このとき $D(A) = 39$ (最適解は 37) となる。

6. むすび

アルゴリズム RT を日本・データゼネラル社の ECLIPSE MV/4000 上で C 言語を用いて実現した。40 個のデータについてシミュレーション実験を行った結果、最適解(分岐限定法による)との相対誤差の平均値は約 19% であった。

現在、耐故障性を向上させるための表の分散配置方法、故障時のネットワーク再構成の方法について検討中である。

文献

- [1] Kikuno, T., et al.: "A cluster-based approach to fault-tolerant computer networks", Proc. ISCAS-86, pp.630-633 (1986).
- [2] Perlman, R.: "Hierarchical networks and subnetwork partition problem", Computer Networks and ISDN Systems, 9, 4, pp.297-303 (1985).
- [3] 当麻喜弘: "フォールトトレラント技術の新フロンティア", 信学技報, FTS85-1, pp.1-8 (1985).