



## Load-balancing on swapped or OTIS networks

Chenggui Zhao<sup>a,b</sup>, Wenjun Xiao<sup>b</sup>, Behrooz Parhami<sup>c,\*</sup>

<sup>a</sup> Department of Computer Science, Yunnan University of Finance and Economics, Kunming, 650221, China

<sup>b</sup> School of Computer Science and Engineering, South China University of Technology, Guangzhou, 510640, China

<sup>c</sup> Department of Electrical and Computer Engineering, University of California, Santa Barbara, CA 931106-9560, USA

### ARTICLE INFO

#### Article history:

Received 12 January 2008

Received in revised form

27 June 2008

Accepted 12 January 2009

Available online 30 January 2009

#### Keywords:

Diffusion

Dimension exchange

Hierarchical network

Load balancing

Optical transpose interconnect system

(OTIS) network

Swapped network

### ABSTRACT

Existing local iterative algorithms for load-balancing are ill-suited to many large-scale interconnection networks. The main reasons are complicated Laplace spectrum computations and flow scheduling strategies. Many large-scale networks are modular and/or hierarchically structured, a prime example being the class of swapped or OTIS networks that have received much attention in recent years. We propose a new local scheme, called DED-X, for load-balancing on homogeneous and heterogeneous swapped/OTIS networks. Our scheme needs spectral information only for the much smaller basis or factor graph, which is of size  $O(n)$  rather than  $O(n^2)$ , and it schedules load flow on intragroup and intergroup links separately. We justify the improvements offered by DED-X schemes over traditional X schemes analytically and verify the advantages of our approach, in terms of efficiency and stability, via simulation.

© 2009 Elsevier Inc. All rights reserved.

### 1. Introduction

Load-balancing is the process of redistributing workloads among computational nodes in a parallel or distributed computing environment, when static or a priori distribution fails to achieve near-perfect balance, thus leading to suboptimal efficiency and speedup. In some cases, dynamic, poorly predictable task characteristics bring about a highly uneven load distribution, causing severe performance degradation. Examples abound in cluster and grid computing, which are characterized by a broad collection of dynamically generated loads on processing nodes. Workload equalization is achieved by moving tasks and/or finer chunks of work between nodes and their neighbors via communication links that connect them.

Load-balancing algorithms typically assume that the node workload consists of equally sized items and that the workload is infinitely divisible. The goal of load-balancing is to design scheduling algorithms to migrate load across network links, with each node ideally ending up with a load that matches its capabilities. A node communicates with one neighbor at a time in *dimension exchange* algorithms and with all neighbors simultaneously in *diffusion* algorithms. In parallel processing terminology, these are

known as single-port and all-port communication models, respectively. Commonly, load-balancing entails two distinct phases of balance calculation and choice of items to transfer. Balance calculation yields the amount of load that should be migrated between a node and its neighbors to achieve a balanced status. Selection of load items for actual transfer entails a number of criteria relating to workload characteristics. A good load balancing algorithm has a numerically stable iterative process, with low computational complexity and a small flow on communication links for achieving a balanced state.

A number of significant algorithms has been developed for load-balancing on general networks. For homogeneous networks, Cybenko [4] presented a local diffusion load-balancing scheme. Muthukrishnan et al. [13] refer to Cybenko's method as first-order scheme (FOS) and endeavor to speed up its iteration process by using an overrelaxation iterative method in their second-order scheme (SOS). Diekmann et al. [6] developed an iterative algorithm, dubbed the optimal polynomial scheme (OPS), to balance loads among nodes within a finite number of iterations. They also provided a theoretical analysis for their algorithm. Elsässer et al. [7] presented an optimal diffusion scheme, OPT, which balances the node loads in a finite number of iterative steps, once the graph's Laplace spectrum becomes known. They subsequently extended several polynomial diffusion schemes for load-balancing from homogeneous to heterogeneous networks [8], with the resulting "balanced" node loads being proportional to their given weights. The schemes cited above produce  $l_2$ -minimal

\* Corresponding author.

E-mail addresses: [zhaochenggui@126.com](mailto:zhaochenggui@126.com) (C. Zhao), [wjxiao@scut.edu.cn](mailto:wjxiao@scut.edu.cn) (W. Xiao), [parhami@ece.ucsb.edu](mailto:parhami@ece.ucsb.edu) (B. Parhami).

balancing flows. In fact, it is the case that all local iterative algorithms lead to a minimal flow, independent of the algorithms and parameters used.

Motivated by load-balancing research in other contexts, Qin [16] designed a data-aware load-balancing strategy to achieve high performance for data-intensive jobs in data grid environments. This was accomplished via a model for estimating the job response time to calculate slowdowns imposed on jobs to balance the load of a data grid in such way that computation and storage resources in each site are simultaneously utilized. Harchol-Balter and Downey [11] proposed a functional form to fit the distribution of lifetimes for Unix processes and derived a preemptive migration load-balancing strategy. They also showed that their policy reduces the mean delay by 35%–50%, compared with other preemptive migration policies.

General diffusion algorithms are ill-suited to load-balancing on large-scale networks, owing to their complicated Laplace spectrum computation. Swapped or OTIS networks constitute a case in point. Optical transpose interconnect system (OTIS) networks, named *swapped* networks by Parhami [15] (who provides a historical review and cites many references to original papers on these networks), are built of  $n$  copies of an  $n$ -node factor or basis network, and thus have a total of  $n^2$  nodes. Nodes are linked according to the connectivity of the basis network in clusters and via intercluster links to other clusters. The original proposal aimed to realize these intercluster links as optical channels, hence the name “OTIS”. Computing the Laplace spectrum of a swapped/OTIS network must take about  $O(n^2)$  steps. Interestingly, many algorithms for swapped networks can be based on the respective algorithms/properties of the much smaller basis network. It is thus natural to ask whether spectrum computation for load-balancing can likewise be limited to the basis network in order to make the algorithm much more efficient. Before showing that this simplification can indeed be achieved, we review several studies that deal with load-balancing in large-scale networks.

Several diffusion schemes, described as *alternating direction iteration* (ADI), have been proposed by Elsässer et al. [7] to deal with load equilibrium problems for large and scalable networks. The same research group also introduced *mixed direction iteration* (MDI) in [9] to obtain a smaller flow than ADI, with the same number of iterations. They present ADI-FOS and ADI-OPT, that is, ADI versions of the general diffusion schemes FOS and OPT. When applied to product graphs, the MDI method converges to balanced status faster than the corresponding general diffusion schemes, and the number of iterations is always smaller than the latter. However, these schemes are applicable only to networks modeled as a Cartesian product of two graphs, and thus cannot be used for OTIS architectures.

The diffusion algorithms discussed thus far assume the all-port communication model. For the single-port model, Arndt [1] has constructed the dimension-exchange (DE) algorithms DE-OPT and DE-OPS (which use the same iteration as in OPT and the same iterative polynomial as in OPS diffusion algorithms), but divided each iterative step in both diffusive algorithms into substeps corresponding to edge colors on the underlying graph. DE-OPS and DE-OPT use the most recent information, meaning that a node exchanges its load information with one of its neighbors in each substep. The diffusion matrices of OPT and OPS are replaced with their DE counterparts, and the eigenvalues of the original Laplacian matrix are likewise modified. The convergence speed of DE still depends on the number of distinct eigenvalues of the Laplacian. For product graphs, Arndt [2] developed a new diffusion algorithm ADI-OPS, together with two DE algorithms DE-ADI-OPT and DE-ADI-OPS.

Our aim here is to construct new diffusion algorithms based on general diffusion schemes, such as FOS, SOS, and OPT, so that they

**Table 1**

Terminology and abbreviations, listed for ready reference.

Term	Meaning or interpretation
ADI	Alternating direction iteration
Adjacency	Matrix representation of a graph
Basis	Component graph from which a swapped/OTIS network is built
Convergence	Speed with which a balanced load status is achieved
CS	Weight assignment with all nodes of weight 1, except one with weight $n + 1$
DE	Dimension exchange; communication with one neighbor at a time
DED-X	Three-phase diffusion-exchange-diffusion scheme based on X
Error	Difference between an achieved load distribution and the ideal balanced load
Factor	Same as “basis”
Flow	Extent of workload transfers in a network
FOS	First-order scheme
HOMO	Homogeneous weight assignment: all nodes are given weight 1
Intragroup	Links that connect nodes located in the same basis network
Intergroup	Links that connect nodes located in different basis networks
Laplacian	A particular matrix representation of a graph
Load	Units of work assigned to a particular entity (node, subnetwork, or network)
MDI	Mixed direction iteration
Migration	Transfer of workload among the network nodes
Network	Set of nodes interconnected by links; used interchangeably with “graph”
Norm	Parameter characterizing a diffusion scheme
OPT	Optimal diffusion scheme
OTIS	Optical transpose interconnect system; used interchangeably with “swapped”
OTIS-G	OTIS network with graph $G$ as basis network
OTIS- $H_d$	OTIS network with a $dD$ hypercube as basis network; generically, OTIS-Cube
OTIS- $M_{k \times m}$	OTIS network with a $k \times m$ mesh as basis network; generically, OTIS-Mesh
PEAK	A highly skewed load distribution where a single node holds the entire load
Quality	Inversely related to extent of flow: the smaller the flow, the higher the quality
RANL	Random load distribution
RANW	Weight assignment, with all node weights being random integers in [1, 64]
SEMI	Node assignment in which nodes are given weight 1 or 2, in equal numbers
SOS	Second-order scheme
Stability	Property of an algorithm whose error decreases smoothly and monotonously
Swapped	Used interchangeably with “OTIS”
Weight	Value assigned to a node, reflecting its computational capacity

can perform load -balancing on OTIS networks with the same level of computational overhead as would be needed for load-balancing on their much smaller basis networks. Accordingly, we propose several hybrid load-balancing schemes and show them to possess a simple iteration process, as well as high efficiency, when applied to a wide array of OTIS networks whose basis networks have regular topologies. Table 1 contains a list of key terms and abbreviations used in concert with other standard graph-theory and parallel-computing terms [10,14].

The rest of this paper is organized as follows. After presenting basic definitions pertaining to load-balancing, diffusion algorithms, and OTIS networks in Section 2, we review the application of several existing general diffusion schemes to homogenous OTIS networks in Section 3. We present our local iterative algorithms for load-balancing on homogeneous swapped/OTIS networks in Section 4, extending the proposed schemes to heterogeneous OTIS networks in Section 5. In Section 6, we analyze the performance of these schemes and present simulation results to show the viability of our approach.

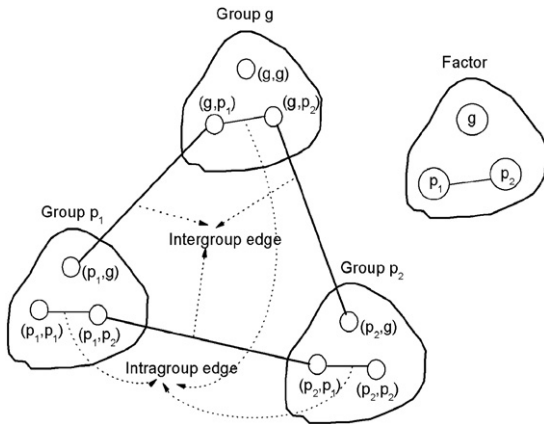


Fig. 1. The structure of swapped/OTIS network.

2. Definitions and background

The swapped/OTIS architecture (see, e.g., [5,15]) derived from a general graph  $G$  is denoted as OTIS- $G$  or  $S(G)$ . A formal definition is given below. Throughout this paper, we use swapped and OTIS networks interchangeably.

**Definition 1 (Swapped/OTIS Graph).** Let  $G = (V_G, E_G)$  be an undirected graph. The swapped or OTIS graph associated with  $G$ , OTIS- $G = S(G) = (V, E)$ , is an undirected graph with the vertex set  $V = \{(g, p) \mid g, p \in V_G\}$  and the edge set  $E = E_b \cup E_s$ , where  $E_b = \{((g, p_1), (g, p_2)) \mid g \in V_G, (p_1, p_2) \in E_G\}$  and  $E_s = \{((g, p), (p, g)) \mid g, p \in V_G, g \neq p\}$ . □

Informally, a swapped/OTIS network, derived from an  $n$ -node network  $G$ , is composed of  $n$  clusters, each of which is internally connected as  $G$ . Additionally, node  $i$  of cluster  $j$  (where  $i \neq j$ ) is connected externally to node  $j$  of cluster  $i$ .

The graph  $G$  is the *factor* or *basis* network of OTIS- $G$ . If  $G$  has  $n$  nodes, then OTIS- $G$  is composed of  $n$  node-disjoint subnetworks  $G_i$ ,  $i = 0, 1, \dots, n - 1$ , which constitute the groups or clusters. Each of these groups is isomorphic to  $G$ . Denote the vertex set of  $G_i$  as  $V_i = \{v_{ij} \mid 0 \leq j \leq n - 1\}$  and its edge set as  $E_i = \{(v_{ik}, v_{il}) \mid (v_k, v_l) \in E_G\}$ . The vertex set  $V$  of OTIS- $G$  is  $V = \cup_{0 \leq i \leq n-1} V_i$ . The edge set  $E$  of OTIS- $G$  can be partitioned into two subsets: The intragroup or basis edge set  $E_b$ , and the intergroup or swap edge set  $E_s$ . Clearly,  $E_b = \cup_{0 \leq i \leq n-1} E_i$  and  $E_s = \{(v_{ij}, v_{ji}) \mid i < j\}$ .

Fig. 1 depicts the structure of a swapped/OTIS network, along with the terminology used to refer to its various parts. Fig. 2(a) contains an example OTIS network formed with the 4-node cycle  $C_4$  as its basis or factor network. The example OTIS network shown in Fig. 2(b) is based on the 6-node complete graph  $K_6$ .

Let  $w_i = (w_{i1}, w_{i2}, \dots, w_{in})^T$  and  $c_i = (c_{i1}, c_{i2}, \dots, c_{in})^T$  represent the load and weight on the  $i$ th factor  $G_i$  of OTIS- $G$ . Similarly, let  $w = (w_0, w_1, \dots, w_{n-1})^T$  and  $c = (c_0, c_1, \dots, c_{n-1})^T$  denote the load and weight on OTIS- $G$ , where the  $i$ th components are the load  $w_i$  and the weight  $c_i$  of  $G_i$ . The  $j$ th node  $v_{ij}$  of the  $i$ th basis network has initial load  $w_{ij}^0 \geq 0$  and weight  $c_{ij} > 0$ . Notationally,  $C$  and  $C_i$  are taken to be diagonal matrices with elements of the vectors  $c$  and  $c_i$  as their diagonal entries, respectively. That is:

$$C = \text{diag}(c_{01}, \dots, c_{0n}, c_{11}, \dots, c_{1n}, \dots, c_{(n-1)1}, \dots, c_{(n-1)n})$$

$$C_i = \text{diag}(c_{i1}, c_{i2}, \dots, c_{in}).$$

Let  $B_b$  and  $B$  be the node-edge incidence matrices of the basis graph  $G$  and OTIS- $G$  respectively; take  $B_s$  to be the matrix specifying the incidence of the intergroup edges in  $E_s$  to nodes of OTIS- $G$ . Matrices  $B_b$ ,  $B_s$  and  $B$  all have in each column exactly two nonzero entries 1 and  $-1$ , which represent the nodes incident to the corresponding edge. The signs of these nonzero entries imply directions of the flows produced in the process of load-balancing on the corresponding edges. The Laplacian  $L$  of a graph is  $L = BB^T$ . Let  $L$  and  $L_b$  be the Laplace matrices of OTIS- $G$  and its basis network  $G$ , respectively. Let  $A_{ij}$  denote the  $n \times n$  matrix with only the  $ij$ th entry being 1 and other entries being 0. Let  $A_s$  be a matrix with the  $ij$ th entry being  $A_s(i, j) = A_{ij}$ . Then,  $L = I_n \otimes (L_b + I_n) - A_s$ , where  $\otimes$  represents the Kronecker product.

We denote the distinct eigenvalues of  $L$  with  $\lambda_i$  ( $0 \leq i \leq m$ ) and those of  $L_b$  with  $\lambda_i^b$  ( $0 \leq i \leq m_b$ ), arranged in increasing order. Let  $\alpha \in (0, 1)$  be a constant edge weight for OTIS- $G$  and  $\alpha_b$  for  $G$ . Take  $M = I_{n^2} - \alpha L$  and  $M_b = I_n - \alpha_b L_b$  to be the corresponding diffusion matrices of polynomial-based diffusion schemes. Then,  $M$  and  $M_b$  have the eigenvalues  $\mu_i = 1 - \alpha \lambda_i$  and  $\mu_i^b = 1 - \alpha_b \lambda_i^b$ . Denote the second largest eigenvalue of  $M$  and  $M_b$  according to their absolute values with  $\gamma = \max(|\mu_2|, |\mu_m|)$ ,  $\gamma^b = \max(|\mu_2^b|, |\mu_{m_b}^b|)$ . The workload  $w^k$  in step  $k$  for polynomial based diffusion schemes can be commonly expressed as a general iteration form  $w^k = p_k(M)w^0$ . The convergence of this iteration depends on whether the error term  $e^k = w^k - w^l$  tends to zero when  $k$  increases, where  $w^l$  is the node load vector when the network achieves a balanced status.

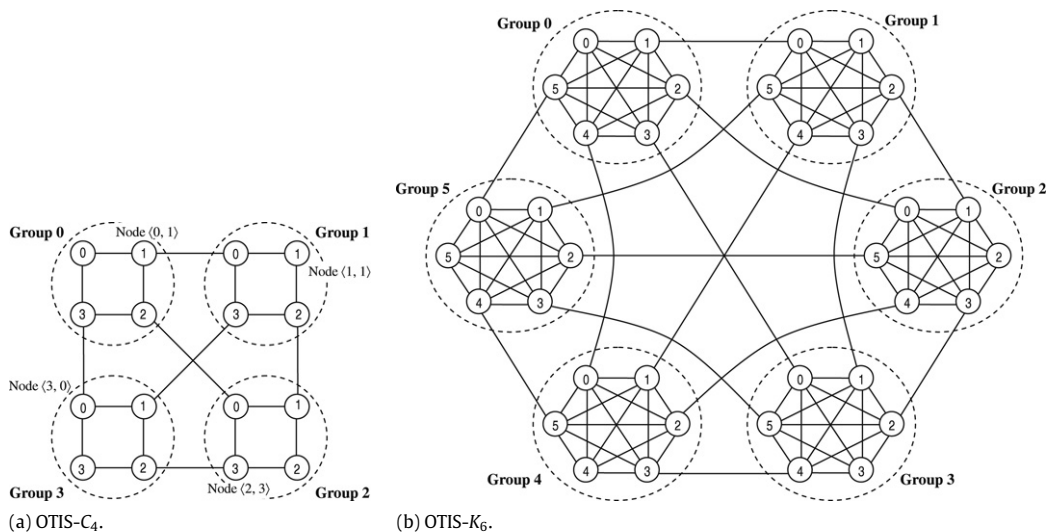


Fig. 2. Two example OTIS networks, one built of the 4-node cycle  $C_4$  as the basis network and the other based on the 6-node complete graph  $K_6$ .

Based on the discussion above, the error term  $e^k$  is an indicator of the quality of load-balancing after  $k$  steps, and it can be used for comparing distinct algorithms with comparable computational complexities. The iteration error term  $e^k$  satisfies (see [6]):

$$\|e^k\|_2 \leq \max\{|p_k(\mu_i)| \cdot \|e^0\|_2, 2 \leq i \leq m\}. \quad (1)$$

In particular, the first-order scheme (FOS) satisfies  $w^k = Mw^{k-1}$  and yields error results in  $\|e^k\|_2 \leq \gamma^k \cdot \|e^0\|_2$ . Since  $\gamma(M) = \max(|1 - \alpha\lambda_2(L)|, |1 - \alpha\lambda_m(L)|)$ , the minimum of  $\gamma$  is achieved for  $1 - \alpha\lambda_2(L) = -1 + \alpha\lambda_m(L)$ . Thus, the optimal value of  $\alpha$  is  $\alpha = 2/(\lambda_2 + \lambda_m)$ . Consequently, we have  $\gamma = (1 - \rho)/(1 + \rho)$ , with  $\rho = \lambda_2(L)/\lambda_m$  constituting the condition number of  $L$ .

For OTIS-G, let  $y^k$  and  $x^k$  be two flow vectors whose entries corresponding to edge  $e$  represent the amount of load migrated along  $e$  in step  $k$  and the total amount of load until step  $k$ , respectively. For  $G$ , let  $y_b^k$  and  $x_b^k$  represent the corresponding parameters of  $G$ . The directions of the flows are determined by the directions of the edges in the incidence matrix. The flow  $x$  is called a balancing flow if and only if  $Bx = w^0 - w^l$ .

We now proceed to describe several known algorithms, including FOS [4], OPS [6], and OPT [7]. These all constitute diffusion schemes for load-balancing on networks with general topologies. The FOS scheme can be expressed as a local iterative scheme. It changes the workload vector  $w^k$  of nodes and schedules the flow vector  $x$  of edges according to  $w^k = Mw^{k-1}$ ,  $x^k = x^{k-1} - \alpha B^T w^{k-1}$ ,  $k \geq 1$ . The error term  $e^k$  of FOS [6] satisfies  $\|e^k\|_2 \leq \gamma^k \cdot \|e^0\|_2$ . To improve the relatively slow convergence of FOS, another polynomial-based iterative method, called a second-order scheme (SOS), was devised [12]. The latter is based on the polynomials:

$$p_0(t) = 1, \quad p_1(t) = t,$$

$$p_k(t) = \beta t p_{k-1}(t) + (1 - \beta) p_{k-2}(t) \quad \text{for } k \geq 2.$$

The iterative process of SOS migrates workloads according to  $w^1 = Mw^0$ ,  $w^k = \beta Mw^{k-1} + (1 - \beta)w^{k-1}$ ,  $k \geq 2$ . It is known that  $w^k$  converges to  $w^l$  whenever  $\beta \in (0, 2)$ , with the fastest convergence occurring for  $\beta = 2/(1 + \sqrt{1 - \gamma^2})$ . Following [6], we denote this optimal value of  $\beta$  as  $\beta_0$ . Then, the error term  $e^k$  in the  $k$ th iteration satisfies:

$$\|e^k\|_2 \leq (\beta_0 - 1)^{k/2} (1 + k\sqrt{1 - \gamma^2}) \|e^0\|_2. \quad (2)$$

After these parameters have been computed, the SOS algorithm can be expressed as a general framework, as suggested in [6].

The optimal scheme OPT [7] has the following simpler iteration process:  $2 \leq k \leq r - 1$ ,  $y^{k-1} = (1/\lambda_{k+1})B^T w^{k-1}$ ,  $x^k = x^{k-1} + y^{k-1}$ ,  $w^k = [I - (1/\lambda_{k+1})L]w^{k-1}$ .

### 3. Hybrid diffusion schemes for homogeneous OTIS-networks

For an OTIS-network with  $n$  vertices on its basis network, all eigenvalues of an  $n^2 \times n^2$  matrix have to be computed before the load-balancing process starts. This is sometimes impractical, motivating us to pursue a hybrid scheme called DED-X, that combines diffusion and dimension exchange, for OTIS networks. Its basic idea is to divide the load-balancing process into three stages of diffusion, exchange, and diffusion. To describe our DED-X approach, let the symbol  $X$  denote any known general load-balancing scheme. In the first stage, DED-X iteratively diffuses node loads until the initial load  $w_i^0$  of the  $i$ th basis network achieves a balanced status  $w_i^l$  locally within the basis network. In this stage, the workload  $w^k$  of step  $k$  can be expressed as  $w^k = (I_n \otimes p_k(M_b))w^0$ .

In the second stage, a dimension exchange strategy is applied over all intergroup links. In this stage, basis networks interchange their balanced node load by a way of swapping the load of node  $(u, v)$  with that of node  $(v, u)$ . At the end of this stage, the total

load on each basis network is the same. Given the status at the end of the first stage, the load of the entire network after this second stage is given by  $w^{l+1} = A_s w^l$ .

In the third stage, we proceed with diffusion using the same iterative polynomial-based scheme as in the first stage. However, given that all basis networks have the same initial load vector, we only compute the load migration on one of the basis networks, using the resulting common flow on all other basis networks.

Fig. 3 illustrates the preceding three-stage load balancing process. Note that the key property responsible for DED-X's extreme efficiency is the ability of an OTIS network to disperse the balanced basis network loads uniformly over all clusters, making the load vector in each cluster identical after a single load exchange step via the intercluster (swap) links.

We now proceed to prove that any polynomial-based scheme used in the DED-X framework must force  $w^k$  to the average of node loads in the entire network after the completion of DED-X's diffusion-exchange-diffusion process.

**Theorem 1.** For any polynomial-based scheme  $X$  that takes at most  $l$  steps to iteratively balance the load within a basis network, from the initial loads  $w_i^0$  for the  $i$ th basis network to the common load  $\frac{1}{n}Jw_i^0$ , DED-X scheme balances the load  $w^0$  to the common load  $\frac{1}{n^2}(J \otimes J)w^0$  in at most  $2l + 1$  steps.

**Proof.** By the condition of this proposition,  $p_l(M_b) = (1/n)J$ . Applying the DED-X scheme, we have:

$$\begin{aligned} w^{2l+1} &= [I_n \otimes p_l(M_b)]A_s[I_n \otimes p_l(M_b)]w^0 \\ &= \frac{1}{n^2}(I_n \otimes J)A_s(I_n \otimes J)w^0 = \frac{1}{n^2}(J \otimes J)w^0. \end{aligned}$$

The equation above shows that  $w^0$  will tend to the final balanced load of  $\frac{1}{n^2}(J \otimes J)w^0$  within  $2l + 1$  steps.  $\square$

The DED-X scheme, as applied here to homogeneous systems, will be expressed in the form of a local iterative algorithm in Section 4 (see Fig. 4). Such a DED-X algorithm for heterogeneous networks readily yields DED-X for a homogeneous network as a special case. The performance of these algorithms will be discussed in detail in Sections 5 and 6.

Note that one can select any previously known (e.g., FOS, SOS, OPT) or newly-proposed load-balancing scheme to replace  $X$  in DED-X. Thus, DED-X leads to a variety of practical algorithms with diverse attributes.

The flow calculated by DED-X is not minimal in  $l_2$ -norm. In the following discussion,  $\lambda_i^b$ ,  $0 \leq i \leq m_b$ , denotes the Laplacian eigenvalues of the basis graph  $G$  and  $\lambda_i$ ,  $0 \leq i \leq m$ , denotes those of OTIS-G. Let  $z_i^b$  be the orthogonal eigenvectors corresponding to  $\lambda_i$  satisfying  $\sum_{1 \leq i \leq n} z_i^b = w_k^0$ , where  $w_k^0$  is the part of initial load on the  $k$ th basis network. Let  $z_i$  be the orthogonal eigenvectors corresponding to  $\lambda_i$ . Take  $x^g$  and  $x^{\text{DED}}$  to represent the flows on OTIS-G links resulting from general diffusion schemes and DED-X schemes, respectively. By Theorem 7 of reference [1], the flow  $x^g$  satisfies:

$$x^g = B^T \sum_{i=2}^m \frac{1}{\lambda_i} z_i = \begin{bmatrix} (I \otimes B_b^T) \sum_{i=2}^m \frac{1}{\lambda_i} z_i \\ B_s^T \sum_{i=2}^m \frac{1}{\lambda_i} z_i \end{bmatrix}. \quad (3)$$

If  $x^{D1}$ ,  $x^E$ , and  $x^{D2}$  represent the flows produced at the first diffusion, exchange, and the second diffusion stages of DED-X schemes, then:

$$x^{\text{DED}} = \begin{bmatrix} x^{D1} + x^{D2} \\ x^E \end{bmatrix}. \quad (4)$$

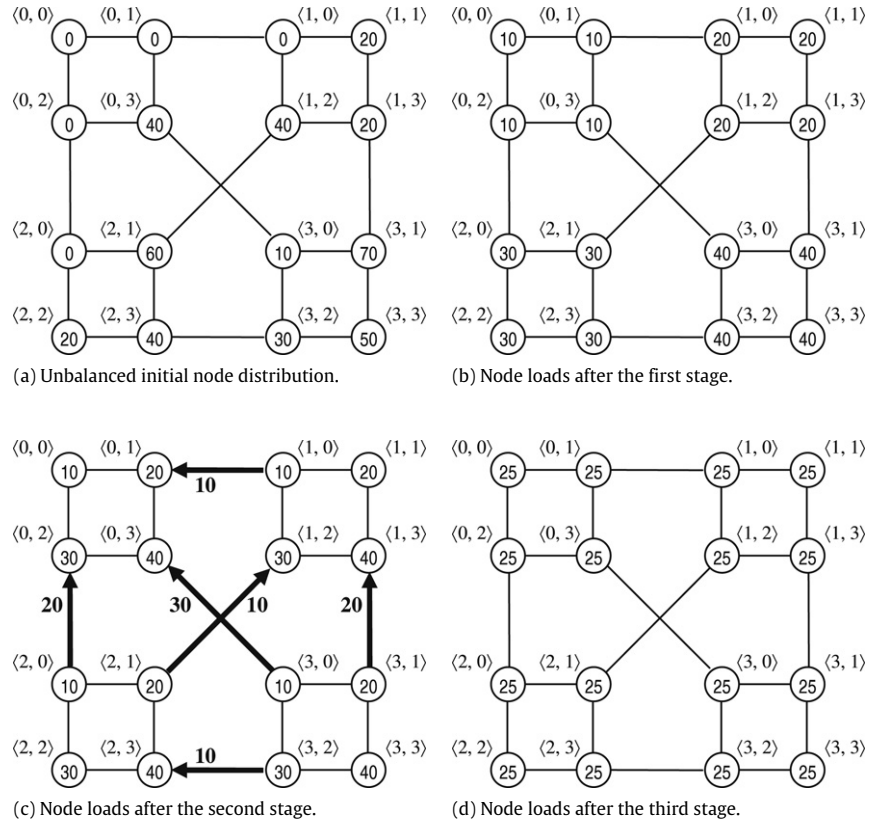


Fig. 3. A simple example for the hybrid DED-X schemes, with loads shown at the beginning and after each of the three stages.

```

Algorithm DED-X

for all groups  $G_i$  of OTIS-G do
    run the procedure X on  $G_i$ ,  $0 \leq i \leq n - 1$ ,
        and input the balanced load vector as  $w_i^1$ ;
end for

for all intergroup edges  $e=((i,j),(j,i))$ ,  $i \neq j$  do
     $y_e^E = \frac{(\sum_{k=1}^{n-1} c_{ik} \sum_{k=1}^{n-1} c_{jk})}{\sum_{i=1}^{n-1} \sum_{j=1}^{n-1} c_{ij}^2} \left( \frac{w_{ij}^{k1}}{c_{ij}} - \frac{w_{ji}^{k1}}{c_{ji}} \right)$ ;
     $w_{ij}^E = w_{ij}^1 - y_e^E$ ;
end for

for all groups  $G_i$  of OTIS-G do
    if load error of  $G_i >$  threshold then
        run algorithm X on  $i$ th group with initial load vector  $w_i^E$ ;
        output the balanced load vector as  $w_i^1$ ;
    end if
end for
    
```

Fig. 4. The structure of the DED-X algorithm.

Note that we have the following two equalities:

$$x^{D_1} = (I \otimes B_b^T) \begin{bmatrix} \sum_{i=1}^{m_b} \frac{1}{\lambda_i^b} z_{1,i}^b \\ \vdots \\ \sum_{i=1}^{m_b} \frac{1}{\lambda_i^b} z_{n,i}^b \end{bmatrix} = (I \otimes B_b^T) \sum_{i=1}^n \sum_{j=1}^{m_b} \frac{1}{\lambda_j^b} (e_i \otimes z_{*,j}^b) \quad (5)$$

$$x^{D_2} = (I \otimes B_b^T) \begin{bmatrix} \sum_{i=1}^{m_b} \frac{1}{\lambda_i^b} z_{*,i}^b \\ \vdots \\ \sum_{i=1}^{m_b} \frac{1}{\lambda_i^b} z_{*,i}^b \end{bmatrix} = (I \otimes B_b^T) \sum_{i=1}^n \sum_{j=1}^{m_b} \frac{1}{\lambda_j^b} (e_i \otimes z_{*,j}^b). \quad (6)$$

Substituting  $x^E = B_s^T \sum_{1 \leq i \leq n} e_i \otimes z_{i,1}^b$  as well as the expressions for  $x^{D_1}$  and  $x^{D_2}$  from Eqs. (5) and (6) into (4), we obtain:

$$x^{DED} = \begin{bmatrix} (I \otimes B_b^T) \sum_{i=1}^n \sum_{j=1}^{m_b} \frac{1}{\lambda_j^b} (e_i \otimes (z_{i,j}^b + z_{*,j}^b)) \\ B_s^T \sum_{i=1}^n e_i \otimes z_{i,1}^b \end{bmatrix}. \quad (7)$$

Because  $z_{i,j}^b + z_{*,j}^b$  is also an eigenvector of  $\lambda_j^b$  letting  $\bar{z}_{i,j}^b = z_{i,j}^b + z_{*,j}^b$ , we get:

$$x^{DED} = \begin{bmatrix} (I \otimes B_b^T) \sum_{i=1}^n \sum_{j=1}^{m_b} \frac{1}{\lambda_j^b} (e_i \otimes \bar{z}_{i,j}^b) \\ B_s^T \sum_{i=1}^n e_i \otimes z_{i,1}^b \end{bmatrix}. \quad (8)$$

Comparing Eqs. (3) and (8), we can conclude that the flow by DED-X scheme will approach the optimal flow when OTIS-G has relatively complicated eigenvalues with respect to those of G, a condition that is true for most regular basis networks.

#### 4. DED-X schemes for heterogeneous OTIS-networks

The DED-X scheme described thus far cannot achieve the balanced state by only local balancing on basis networks, along with load transposition in the second stage, when the OTIS network is not homogeneous. To generalize the DED-X scheme for use with heterogeneous OTIS networks, it is necessary to revise the load exchange strategy and flow schedule components. In this section, we accomplish this goal by assigning weights to the intergroup links.

With a heterogeneous OTIS network, the X scheme used in the first diffusion phase must be tailored for individual basis networks. Next, we must revise the exchange strategy over intergroup links in the second stage. Using  $c_{ij}$  to denote the weight of node  $(i, j)$ , the weight  $l_{ij}$  for the edge  $((i, j), (j, i))$ ,  $i \neq j$ , linking the basis networks  $i$  and  $j$ , is assigned as follows:

$$l_{ij} = \frac{\left(\sum_{k=1}^{n-1} c_{ik}\right) \left(\sum_{k=1}^{n-1} c_{jk}\right)}{\sum_{i=1}^{n-1} \sum_{j=1}^{n-1} c_{ij}}. \quad (9)$$

In the third stage, we proceed with diffusion on each of the basis networks by means of the same iterative polynomial as in the first stage, with flow-scheduling on intragroup edges of each basis network.

We can prove that the DED-X scheme must converge from the initial load  $w^k$  to the balanced status in the case of heterogeneous networks.

**Theorem 2.** For any polynomial based scheme X, if X takes at most  $k_1$  and  $k_2$  steps to redistribute the initial loads  $w_i^0$  and  $w_i^l$  of each of basis networks to achieve balanced status, respectively, then the DED-X scheme can lead from the initial load  $w^0$  to the balanced status in at most  $k_1 + k_2 + 1$  steps.

**Proof.** For any polynomial-based scheme X, and for any  $0 \leq i \leq n - 1$ , the scheme X leads the basis network  $i$  to local balanced status after the first stage. It follows that:

$$w_{ij}^{k_1} = \frac{c_{ij}}{\sum_{k=1}^{n-1} c_{ik}} \sum_{k=1}^{n-1} w_{ik}^0. \quad (10)$$

In the second stage, based on Eqs. (9) and (10), exchanging loads on intergroup links results in the new loads:

$$w_{ij}^{k_1+1} = w_{ij}^{k_1} - \frac{\left(\sum_{k=1}^{n-1} c_{ik}\right) \left(\sum_{k=1}^{n-1} c_{jk}\right)}{\sum_{i=1}^{n-1} \sum_{j=1}^{n-1} c_{ij}} \left(\frac{w_{ij}^{k_1}}{c_{ij}} - \frac{w_{ji}^{k_1}}{c_{ji}}\right) \\ = \left(\frac{c_{ij}}{\sum_{k=1}^{n-1} c_{ik}} - \frac{\sum_{k=1}^{n-1} c_{jk}}{\sum_{i=1}^{n-1} \sum_{j=1}^{n-1} c_{ij}}\right) \sum_{k=1}^{n-1} w_{ik}^0 - \frac{\sum_{k=1}^{n-1} c_{ik}}{\sum_{i=1}^{n-1} \sum_{j=1}^{n-1} c_{ij}} \sum_{k=1}^{n-1} w_{jk}^0. \quad (11)$$

In the third stage, the same iterative polynomial is used, but with different initial loads on the nodes. We thus get, after  $k_2$  additional steps:

$$w_{ij}^{k_1+k_2+1} = \frac{c_{ij}}{\sum_{k=1}^{n-1} c_{ik}} \sum_{j=1}^{n-1} w_{ij}^{k_1+1} = \frac{c_{ij}}{\sum_{k=1}^{n-1} c_{ik}} \sum_{j=1}^{n-1} \left[\frac{\sum_{k=1}^{n-1} c_{ik}}{\sum_{i=1}^{n-1} \sum_{j=1}^{n-1} c_{ij}} \sum_{k=1}^{n-1} w_{jk}^0\right] \\ = \frac{c_{ij}}{\sum_{i=1}^{n-1} \sum_{j=1}^{n-1} c_{ij}} \sum_{j=1}^{n-1} \sum_{k=1}^{n-1} w_{jk}^0 = w_{ij}^l. \quad (12)$$

Eq. (12) establishes that  $w^0$  tends to  $w^l$  within  $k_1 + k_2 + 1$  steps.  $\square$

The structure of the DED-X algorithm is outlined in Fig. 4 as a local iterative process. As in the case of the homogeneous version of the algorithm, one can select any load-balancing scheme, such as FOS, SOS, or OPT, to replace X. Fig. 5 depicts a simple example of the application of the DED-X algorithm for heterogeneous load balancing. The performance of DED-X schemes will be discussed in detail in Sections 5 and 6.

#### 5. Algorithm analysis

A difference between the DED-X and the X algorithm is that the X scheme can run on any general network, whereas DED-X is specific to OTIS networks. However, when the X algorithm is applied to OTIS-G, Laplacian eigenvalues of the entire OTIS-G graph must be known and iterations are executed on all nodes by flowing loads over all edges synchronously. But with DED-X, only the Laplacian eigenvalues of the basis graph G are necessary and iterations proceed only within groups, separated

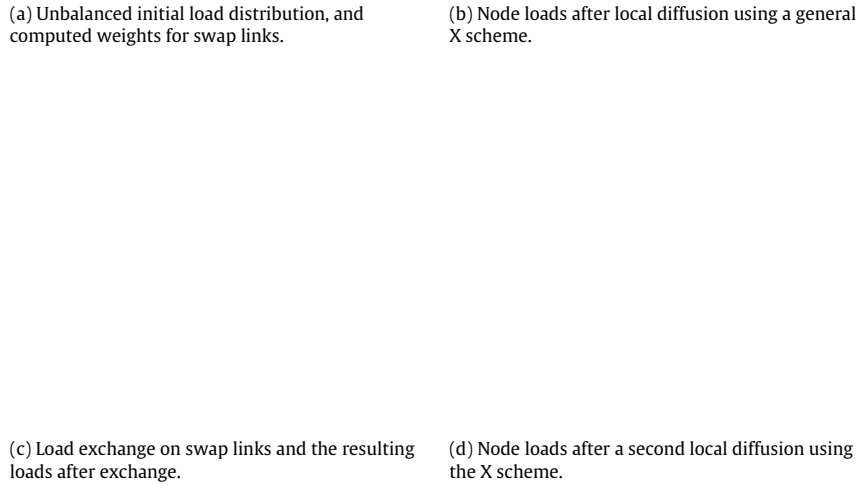


Fig. 5. An example of the DED-X scheme applied to a heterogeneous OTIS network. Node weights in {0.25, 0.5, 1.0} are given immediately below node indices.

by a dimension exchange process over intergroup edges. The most important parameters characterizing the performance of the proposed algorithm include load-balancing accuracy, number of iterative steps, and the amount of flow on communication links. The load-balancing accuracy of DED-X follows from Theorems 1 and 2. Quickness of convergence follows directly from the structure of DED-X, based on the convergence of X. As for the amount of flow, we have analytic results to characterize the extent, but given that the presentation of these results requires the introduction of new notations and extensive derivations, we will report them separately in our future work.

Let  $\mathbf{R}^V$  be the set of functions from the vertex set  $V$  of graph  $G$  to the set  $\mathbf{R}$  of real numbers, that is,  $\mathbf{R}^V = \{f : V \rightarrow \mathbf{R}\}$ . We proceed to prove that OTIS-G has a smaller diffusion norm than  $G$ , when using a polynomial based diffusion scheme X.

**Theorem 3.** Let  $\lambda_1^b$  and  $\lambda_1$  represent the second smallest Laplace eigenvalue of the basis network  $G$  and OTIS-G, respectively. Take  $\lambda_{\max}^b$  and  $\lambda_{\max}$  to denote their largest eigenvalue, respectively. Then it is true that

$$\lambda_1 \leq \lambda_1^b \tag{13}$$

$$\lambda_{\max}^b + 1 \leq \lambda_{\max}. \tag{14}$$

**Proof.** Let  $f \in \mathbf{R}^V$  be a function defined by the eigenvector corresponding to  $\lambda_1^b$ . Then,  $f \perp e$ , where  $e = (1, 1, \dots, 1)^T$ . For any  $g \in \mathbf{R}^V$ , we have  $(g \otimes f)^T(e \otimes e) = 0$ , so we have  $(g \otimes f) \perp e$ . Now considering the Rayleigh quotient expression of matrix eigenvalue, and noting that  $L_b f = \lambda_1^b f$ , we have:

$$\lambda_1 \leq \frac{\langle L_s(g \otimes f), g \otimes f \rangle}{\langle g \otimes f, g \otimes f \rangle} = \frac{(g \otimes f)^T (I_n \otimes L_b + I_{n^2} - A_s)(g \otimes f)}{(g \otimes f)^T (g \otimes f)}$$

$$= \lambda_1^b + 1 - \frac{(g^T \otimes f^T) A_s (g \otimes f)}{\langle g, g \rangle \langle f, f \rangle}. \tag{15}$$

On the other hand, the following two equations hold:

$$\frac{(g^T \otimes f^T) A_s (g \otimes f)}{\langle g, g \rangle \langle f, f \rangle} = \frac{1}{\langle g, g \rangle \langle f, f \rangle} \sum_{i=1}^n g_i f^T \left( \sum_{k=1}^n A_{ki} g_k f \right) \tag{16}$$

$$\sum_{i=1}^n g_i f^T \left( \sum_{k=1}^n A_{ki} g_k f \right) = \sum_{i=1}^n g_i f_i^T g = \langle g, f \rangle^2. \tag{17}$$

Combining Eqs. (15)–(17) we conclude that:

$$\lambda_1 \leq \lambda_1^b + 1 - \frac{\langle g, f \rangle^2}{\langle g, g \rangle \langle f, f \rangle}. \tag{18}$$

By choosing  $g$  to equal  $f$ , we have  $\frac{\langle g, f \rangle^2}{\langle g, g \rangle \langle f, f \rangle} = 1$  and  $\lambda_1 \leq \lambda_1^b$ . Eq. (13) is thus satisfied. Similarly, for any  $g \in \mathbf{R}^V$ , we choose  $f$  as an eigenvector of the maximal eigenvalue  $\lambda_{\max}^b$  of  $L_b$ . Because  $g \otimes f \in \mathbf{R}^V$ , for  $h \in \mathbf{R}^V$  the largest eigenvalue  $\lambda_{\max}$  satisfies:

$$\lambda_{\max} = \max_h \frac{\langle L_s h, h \rangle}{\langle h, h \rangle}. \tag{19}$$

Replacing  $h$  with  $g \otimes f$ , we get:

$$\lambda_{\max} \geq \max_g \frac{\langle L_s(g \otimes f), g \otimes f \rangle}{\langle g \otimes f, g \otimes f \rangle}. \tag{20}$$

In a manner similar to the proof of Eq. (13), we can establish that:

$$\frac{\langle L_s(g \otimes f), g \otimes f \rangle}{\langle g \otimes f, g \otimes f \rangle} = \lambda_{\max}^b + 1 - \frac{\langle g, f \rangle^2}{\langle g, g \rangle \langle f, f \rangle}. \tag{21}$$









