†                                    †    Corinne Touati†

Pareto

Nash
Pareto

# Fair and Pareto optimal load balancing in distributed computer systems

Atsushi Inoie,† Hisao Kameda† and Corinne Touati†

In this paper, we examine a new criterion of fairness for sharing resources among users. The Nash Equilibria can be perceived as fair as they correspond to an equilibrium obtained by fair competition between users. However, they may not be Pareto optimal. We therefore consider, as a fair point the Pareto optimal point where the utilities of users are proportional to those of a Nash Equilibrium point. We study the properties of the new fair point in a simple load balancing system consisting of two servers and two set of users and compare it with previously proposed fair criteria.

## 1. Introduction

Load balancing among computers in a distributed system is a means of efficiently sharing resources among users. We can consider many optimizations in load balancing. An important objective of a load balancing system is Pareto optimality. Indeed, there exists no other state where all users have better benefits simultaneously than in a Pareto optimum situation. There exist innumerably many Pareto-optimal situations, and we may not have any absolute preference among Pareto optima. The weighted sum optimization (which minimizes the weighted sum of the users' costs with different weights) is one of the means to obtain a Pareto optimum.

The choice of one to achieve can be controversial among users. One selection criterion is fairness among users. Various fairness concepts that achieve Pareto optima have been already proposed[3),5)~8)].

There also exist many systems where multiple independent users, or players, strive to unilaterally optimize their own ulitily or cost. They are regarded as noncooperative games, and the equilibrium is called a Nash equilibrium[9)].

† Graduate School of Systems and Information Engineering, University of Tsukuba

Nash equilibria may be Pareto inefficient. However, each Nash equilibrium is fair on all users in the sense that it is achieved by the fair competition (with no coalition) among users. Then, among the Pareto-optima, only those that are strongly Pareto-superior to the Nash equilibrium could satisfy all users. In particular, as the situations that would make all users to feel fairness similar to that of the Nash equilibrium, we consider a group of situations where each user's utility is proportionately larger than that of the Nash equilibrium. We say that such situations are *Nash proportionately fair* to the Nash equilibrium. If we identify a Nash-proportionately-fair Pareto optimum, the resulting situation will satisfy all users since it reflects the competitive fairness given by the Nash equilibrium and is Pareto optimal, at the same time.

By the Pareto set of a system, we mean the set of all Pareto optima of the system. We are quite interested in the positions that the already proposed and Nash proportionate fairness objectives occupy in the Pareto set. It may seem difficult to study this problem in a general framework from this beginning stage. Therefore, in this paper, we deal with simple static load balancing model with two identical servers (computers) each of which has an identical arrival and its own queue, and numerically obtain the cost (the mean response time) of each user at

the points in the Pareto set, at the solutions that achieve various fairness objectives, and at the Nash equilibrium (which is unique in this case[1]).

The rest of this paper is organized as follows. Section 2 describes our model and formulates as various types of fair and optimal load balancing problems. Section 3 shows some numerical results. Section 4 concludes this article.

## 2. Model and Assumptions

We consider a distributed computer system (shown in Fig. 1) consisting of two servers (computers), 1 and 2, with two flows of demands $\phi_1$ and $\phi_2$ arriving from users 1 and 2 at servers 1 and 2, respectively. Service times of servers 1 and 2 are according to exponential distribution at mean $1/\mu_1$ and $1/\mu_2$, respectively. Let a fraction $x_i$ $(0 \leq x_i \leq \phi_i)$ of a flow of jobs be forwarded from server $i$ to the other server $j$ $(\neq i)$. Denote by $x$ the vector $(x_1, x_2)$, and by $\beta_1$ and $\beta_2$, respectively, the resulting loads on nodes 1 and 2. Then, we have

$$\beta_i = \phi_i - x_i + x_j, \quad i, j = 1, 2 \ (i \neq j).$$

We assume that the processing time at server $i$ for the load of rate $\beta_i$ is given by $(\mu_i - \beta_i)^{-1}$. For simplicity, we assume that forwarding a job requires a fixed delay $t$. Therefore, the cost of user $i$, that is, the delay of each flow arriving from the user $i$, can be written:

$$T_i(\boldsymbol{x}) = \frac{1}{\phi_i}\Big[ \frac{\phi_i - x_i}{\mu_i - \phi_i + x_i - x_j} + x_i(t + \frac{1}{\mu_j - \phi_j + x_j - x_i}) \Big],$$
$$(1)$$

for $i, j = 1, 2 (j \neq i)$.

Denote by $C$ the feasible region of $x$, that is, $C = \{x \mid 0 \leq x_i \leq \phi_i, \ i = 1, 2 \text{ and } \mu_i - \phi_i + x_i - x_j > 0, \ i, j = 1, 2 \ (i \neq j)\}$. Clearly, $C$ is a convex set. However, note that $T_i(\boldsymbol{x})$, $i = 1, 2$, $\boldsymbol{x} \in C$, is not convex in $\boldsymbol{x}$ while $T_i(\boldsymbol{x})$, $i = 1, 2$, is convex in each of $x_1$ and $x_2$ for $(x_1, x_2) \in C$.

For example, we may consider that the utilities of the users are inversely proportional to their costs, that is, $U_i(T_i) = 1/T_i, i = 1, 2$.

### 2.1 Pareto set

Denote by $\Pi$ the Pareto set defined as follows:

$\Pi = \{(T_1(\boldsymbol{x}), T_2(\boldsymbol{x})) \mid \boldsymbol{x} \in C,$ and for any $\boldsymbol{x}'$ $(\boldsymbol{x}' \in C)$
if $T_i(\boldsymbol{x}') < T_i(\boldsymbol{x})(i = 1, 2)$ then $T_j(\boldsymbol{x}') > T_j(\boldsymbol{x})(j \neq i)\}$

Minimization of a weighted sum of costs of users is one of the objectives to obtain the Pareto set. The objective is represented as follows:

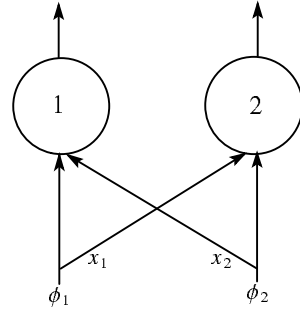$$\Omega(\bar{\boldsymbol{x}}) = \min_{x \in C} \Omega(\boldsymbol{x}) \qquad (2)$$

**Fig.1**  Load balancing in a distributed system consisting of two servers

where $\Omega(\boldsymbol{x}) = \sum_p \xi_p T_p(\boldsymbol{x}), \xi_i \geq 0, i = 1, 2$, and $\sum_p \xi_p > 0$. Then, the cost of user $i$ is given by $T_i(\bar{\boldsymbol{x}})$, $i = 1, 2$. Clearly, $\bar{\boldsymbol{x}}$ gives a Pareto optimum.

This objective has the following property[2]. If the strategy set $C$ is convex, and if the cost functions $T_i$, $i = 1, 2$, are convex, any Pareto optimum is given by a solution $\bar{\boldsymbol{x}}$ of the minimization of weighted sums of costs (2). In fact, however, $T_i$, $i = 1, 2$, are not convex in $\boldsymbol{x}$ as noted above. Therefore, there may exist Pareto optimal points that are not a solution of minimizing a weighted sum of costs, which is actually shown in the numerical examples given in the later section.

### 2.2 Nash proportionate fairness

In this system, a Nash equilibrium $\tilde{\boldsymbol{x}}$ is given as follows:

$$T_i(\tilde{\boldsymbol{x}}) = \min_{x_i} T_i(x_i, \tilde{x}_j), \quad \text{s.t.} \ (x_i, \tilde{x}_j) \in C, \ i, j = 1, 2 \ (i \neq j).$$
$$(3)$$

For the model in question, there exists a unique Nash equilibrium[1].

The Nash equilibrium, $T_i(\tilde{\boldsymbol{x}})$, $i = 1, 2$, may be Pareto inefficient[4]. Consider $P_i = \eta T_i(\tilde{\boldsymbol{x}})$, $i = 1, 2$, By decreasing $\eta$, if $(P_1, P_2)$ hits the curve $\Pi$, and reaches a Pareto optimal point, $(\bar{P}_1, \bar{P}_2)$, it is the Nash-proportionate-fair Pareto optimum.

### 2.3 Already proposed fairness

Already proposed lines of general parameterized fairness objectives are expressed, for example, in the following form[8].

$$F(\hat{\boldsymbol{x}}) = \min_{x \in C} F(\boldsymbol{x}), \qquad (4)$$

where $F(\boldsymbol{x}) = (1 - \alpha)^{-1} \sum_p \{T_p(\boldsymbol{x})\}^{1-\alpha}$.

Similarly as (4), we consider maximization of the utilities of the users, and formulate the following fairness objectives:

$$\check{F}(\hat{\boldsymbol{x}}) = \max_{x \in C} \check{F}(\boldsymbol{x}), \qquad (5)$$

where $\check{F}(\boldsymbol{x}) = (1 - \alpha)^{-1} \sum_p \{U_p(\boldsymbol{x})\}^{1-\alpha}$ and $U_p(\boldsymbol{x}) =$
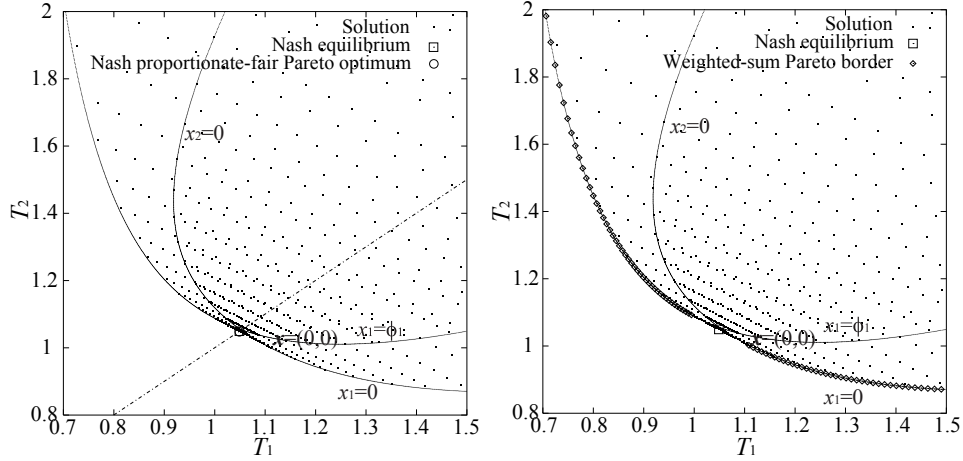
**Fig. 2** Combinations of response times, respectively, $T_1$ and $T_2$, of users 1 and 2 in the case where the values of system parameters are $\phi_1 = 2.1$, $\mu_1 = 3$, $\phi_2 = 2.7$, $\mu_2 = 3.7$, and $t = 0.001$. Note that the left and right graphs of this figure and the left and right graphs of the following Fig. 3 show the same case of the system. Only differences lie in that the points shown may be what achieve objectives different among the four graphs. This applies also to the set of 4 graphs in Figs. 4 and 5 and to the set of 4 graphs in Figs. 6 and 7. We show the points that achieve **[Left]** Nash proportionate fairness and **[Right]** minimization of weighted sums of costs.
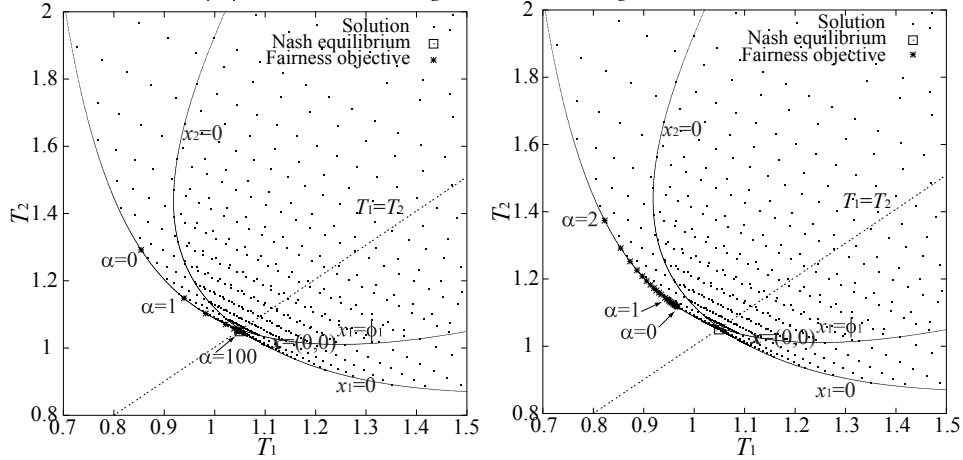


**Fig. 3** Combinations of response times, respectively, $T_1$ and $T_2$, of users 1 and 2 in the case where $\phi_1 = 2.1$, $\mu_1 = 3$, $\phi_2 = 2.7$, $\mu_2 = 3.7$, and $t = 0.001$. We show the points that achieve the already proposed fairness objectives (5) **[Left]** and (4) **[Right]**.

$1/T_p(\boldsymbol{x})$.

The case of $\alpha = 0$ shows the simple sum of the costs of the users. The case of $\alpha \to 1$ presents a Nash bargaining solution and, in particular, proportional fairness in this system. The case of $\alpha \to \infty$ corresponds to the Max-Min fairness[3),5)~8)]. Obviously, (5) also follows the lines of already proposed general parameterized fairness objectives.

## 3. Numerical Results

We characterize fair and optimal load balancing prob-

lem through some numerical results. For convenience, we add the constraint $\xi_1 + \xi_2 = 1$ in minimization of weighted-sums without losing generality.

### 3.1 A case where weighted-sum objective does not cover all the Pareto optima

Figs. 2 and 3 show the Nash equilibrium, the part of Pareto set obtained by the weighted-sum optimization and the fairness solutions that achieve (4) and (5). The values of the system parameters are $\phi_1 = 2.1$, $\mu_1 = 3$, $\phi_2 = 2.7$, $\mu_2 = 3.7$, and $t = 0.001$.

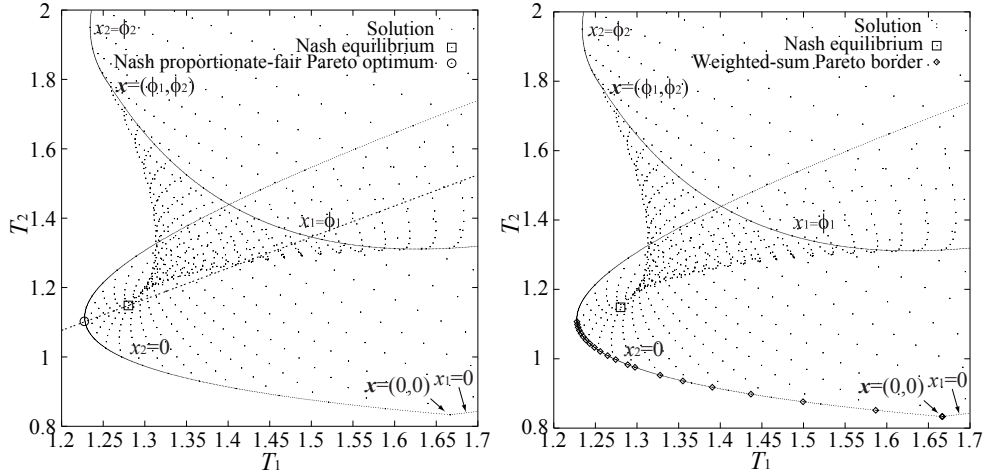In the right graph of Fig. 2, we observe that the parts

**Fig. 4**  Combinations of response times, respectively, $T_1$ and $T_2$, of users 1 and 2 in the case where $\phi_1 = 0.9$, $\mu_1 = 1.5$, $\phi_2 = 0.8$, $\mu_2 = 2$, and $t = 0.35$. We show **[Left]** the point that achieves the Nash proportionate fairness, and **[Right]** the points that achieve the minimization of weighted sums of costs.
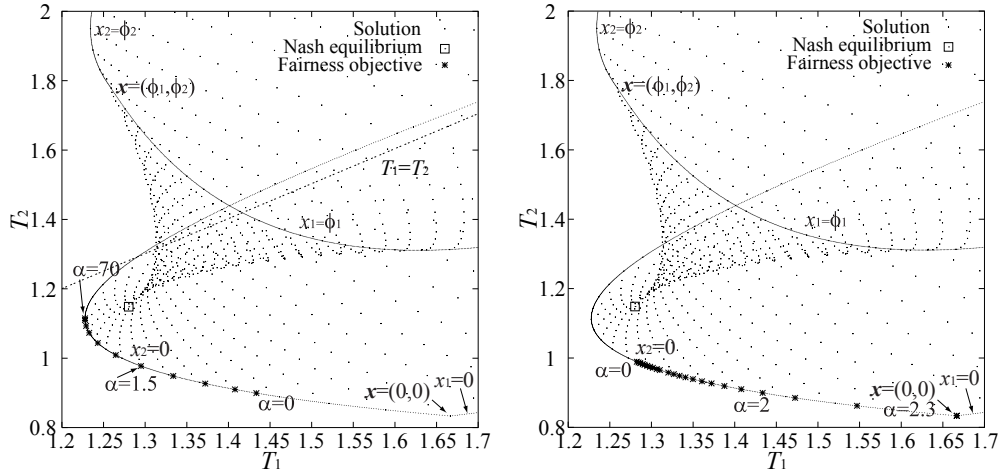


**Fig. 5**  Combinations of response times, respectively, $T_1$ and $T_2$, of users 1 and 2 in the case where $\phi_1 = 0.9$, $\mu_1 = 1.5$, $\phi_2 = 0.8$, $\mu_2 = 2$, and $t = 0.35$. We show the points that achieve the already proposed fairness objectives (5) **[Left]** and (4) **[Right]**.

of the Pareto set obtained by the weighted-sum objectives are divided into two parts and do not cover all the Pareto set. We note that $T_1$ and $T_2$ are nonconvex in $x$, and, therefore, the optimal solutions to the weighted-sum objectives may not cover the Pareto set since the conditions of the Aubin's theorem are not satisfied. Thus, the above result presents a counter-example that shows that if a condition of the Aubin's theorem is not satisfied, the theorem does not hold.

In Fig. 3, we observe that, as to the solutions obtained by the already proposed fairness objectives (5), the optimal points converges to the point of the Pareto set satisfy-

ing $T_1 = T_2$ as the value of $\alpha$ increases. (Note, however, that, in this case, we were unable to obtain numerically the optimal values for very large values of $\alpha$ ($>$ 100). This is perhaps because of accumulation of round-off errors.) In particular, some such optimal points are in the part of the Pareto set which the weighted-sum objective cannot cover. On the other hand, the points optimal for the objectives (4) diverge from the Max-Min fair point as the value of $\alpha$ increases. Superficially thinking, both the objectives (5) and (4) would be anticipated to show similar behaves, but, in fact, the objective (4) is not good as a general fairness objective.
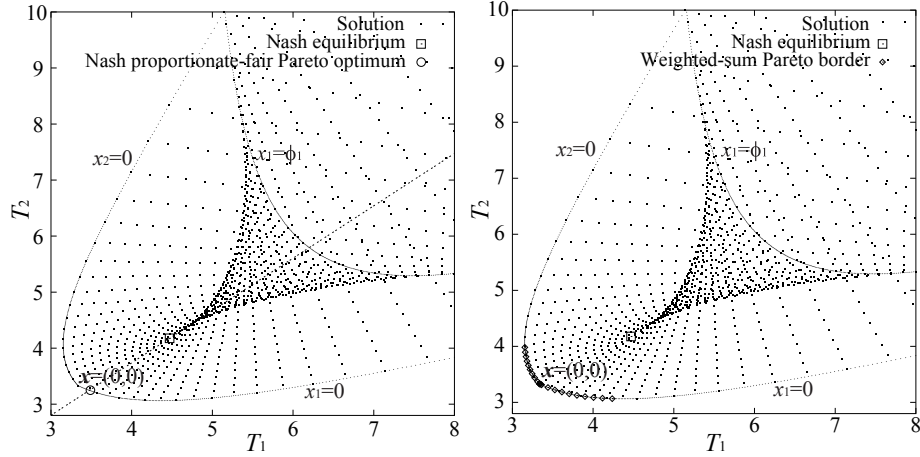
18

**Fig. 6** Combinations of response times, respectively, $T_1$ and $T_2$, of users 1 and 2 in the case where $\phi_1 = 0.7$, $\mu_1 = 1.0$, $\phi_2 = 0.9$, $\mu_2 = 1.2$, and $t = 3$. We show the points that achieve **[Left]** Nash proportionate fairness and **[Right]** minimization of weighted sums of costs.
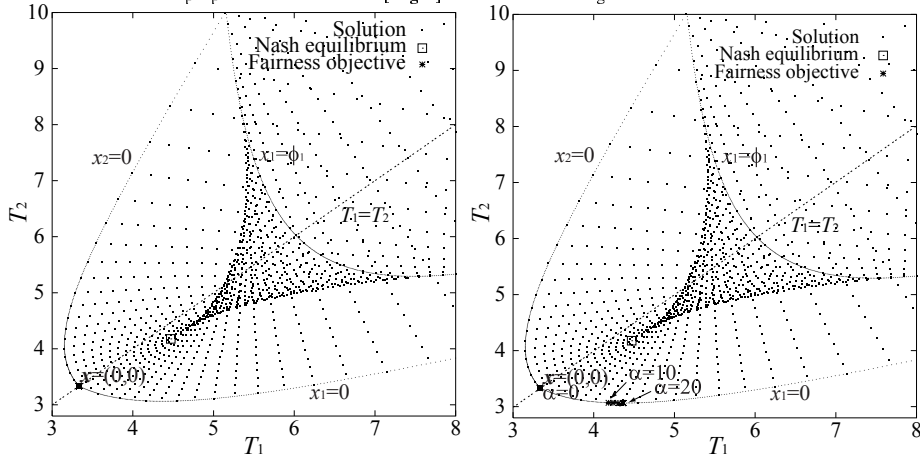


**Fig.7** Combinations of response times, respectively, $T_1$ and $T_2$, of users 1 and 2 in the case where $\phi_1 = 0.7$, $\mu_1 = 1.0$, $\phi_2 = 0.9$, $\mu_2 = 1.2$, and $t = 3$. The points achieve the fairness objectives (5) **[Left]** and (4)**[Right]**.

It is seen in the above figures that the Nash equilibrium is almost Pareto optimal, and almost identical to the Nash proportionate-fair Pareto optimum. In this case, however, the Nash proportionate-fair Pareto-optimal point is not in the part of the Pareto set obtained by weighted-sum objectives and any fairness objectives.

### 3.2 A case where the Nash equilibrium is not Pareto optimal

Figs. 4 and 5 show a case where the values of system parameters are $\phi_1 = 0.9$, $\mu_1 = 1.5$, $\phi_2 = 0.8$, $\mu_2 = 2$, and $t = 0.35$.

In Fig. 4, we observe that the Nash equilibrium is not on the Pareto set. The Pareto set and the straight line passing through the origin (0,0) and the Nash equilibrium

intersect at a point, which is the Nash proportionate-fair Pareto-optimal point.

In this case, Pareto optimal points that achieve the weighted-sum optimization cover all the Pareto set. The Pareto optimum corresponding to the Nash proportionate fairness is given by $\xi_1 \simeq 0.934$ and $\xi_2 \simeq 0.066$. In this case, the Nash-proportionate-fair optimal point happens to be the point that achieves the fairness objective (5) for $\alpha \simeq 26.4$. On the other hand, in this case, no points that achieve the fairness objectives (4) with any values of $\alpha$ can be identical with it.
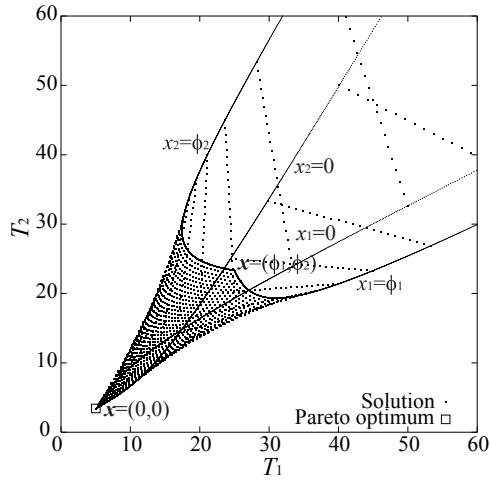
19

**Fig. 8** Combinations of response times, respectively, $T_1$ and $T_2$, of users 1 and 2 in the case where $\phi_1 = 0.5$, $\mu_1 = 0.7$, $\phi_2 = 0.4$, $\mu_2 = 0.7$, and $t = 20$.

### 3.3 A case where only one Pareto optimum point achieves the fairness objectives (5) with various values of $\alpha$

Figs. 6 and 7 show a case where only one Pareto optimum point achieves the fairness objective (5) at $T_1 \simeq 3.333$ and $T_2 \simeq 3.333$, that is, the case of no load balancing ($x_1 = 0$ and $x_2 = 0$). The values of the system parameters are $\phi_1 = 0.7$, $\mu_1 = 1.0$, $\phi_2 = 0.9$, $\mu_2 = 1.2$, and $t = 3$. Note that in this case, the following relation is satisfied: $\phi_1 - \mu_1 = \phi_2 - \mu_2$. Note that, in this case also, the Nash proportionate-fair Pareto-optimal point is different from the Pareto optimum that achieves the fairness objective (5).

### 3.4 A case where only one Pareto optimum point exists

Fig. 8 shows a case where only one Pareto optimum point exists. The value of the system parameters are $\phi_1 = 0.5$, $\mu_1 = 0.7$, $\phi_2 = 0.4$, $\mu_2 = 0.7$, and $t = 20$. We note that load balancing must be ineffective when job forwarding time $t$ has a large value. In Fig. 8, all optimal points that achieve the weighted-sum optimization for any combinations of the values of $\xi_1$ and $\xi_2$, the points that achieve both fairness objectives (4) and (5) with any values of $\alpha$, and the Nash equilibrium point happens to be the Pareto optimal point.

### 4. Concluding Remarks

We have numerically examined the generally parameterized fairness objectives and the Nash-proportionate

fairness recently introduced. The platform of this research has been simple static load balancing model with two identical servers (computers) each of which has an identical arrival and its own queue.

The points that achieve the general parameterized fairness objectives generally cover a part of the Pareto set, and at times, do not covered the Nash-proportionate-fair Pareto optimal point. Since each Pareto optimum may have its own significance, we may wish to have a more generally parameterized fairness objective that all the Pareto optimal points may be achieved with a certain choice of the values of the parameters.

We have observed that careful consideration is needed in establishing the concrete form of the fairness objective along the lines of the generally parameterized fairness objectives. Otherwise, we may have an inappropriate objective that would not give us truly fair assignment of resources to users.

### References

1) Altman, E., Kameda, H. and Hosokawa, Y.: Nash equilibria in load balancing in distributed computer systems, *International Game Theory Review*, Vol. 4, No. 2, pp. 91–100 (2002).

2) Aubin, J.-P.: *Optima and Equilibria: An Introduction to Nonlinear Analysis, 2nd Ed*, Springer-Verlag, Berlin (1998).

3) Bertsekas, D. and Gallager, R.: *Data Networks, 2nd Ed*, Prentice-Hall, Englewood Cliffs (1992).

4) Kameda, H., Altman, E., Kozawa, T. and Hosokawa, Y.: Braess-like paradoxes in distributed computer systems, *IEEE Trans. Automatic Contr.*, Vol. 45, No. 9, pp. 1687–1691 (2000).

5) Kelly, F. P.: Charging and rate control for elastic traffic, *European Transactions on Telecommunications*, Vol. 8, pp. 33–37 (1997).

6) Maulloo, A., Kelly, F. P. and Tan, D.: Rate control in communication networks: Shadow prices, proportional fairness and stability, *Journal of the Operational Research Society*, Vol.49, pp.237–252 (1998).

7) Mazumdar, R., Mason, L. G. and Doulgligeris, C.: Fairness in network optimal flow control: Optimality of product forms, *IEEE Trans. Communications*, Vol. 39, No. 5, pp. 237–252 (1991).

8) Mo, J. and Walrand, J.: Fair end-to-end window-based congestion control, *IEEE/ACM Trans. Networking*, Vol. 8, No. 5, pp. 556–567 (2000).

9) Nash, J. F.: Non-cooperative games, *Ann. Math.*, Vol. 54, pp. 286–295 (1951).