

解説



定性推論のパラメータチューニングへの応用†*

伊藤 潔†† 本位田 真一†††

1. はじめに

動的なシステムの性能を改善するパラメータチューニングを厳密に行おうとすると、複雑な数式処理や数値計算が必要で計算量も多くなる。また、一般的な解法が存在しないこともあるため、さまざまな近似モデルや近似解法を用いることもある。この近似モデルや近似解法が現実の対象システムから遊離してしまうこともある。このため、専門家は、問題の性質に応じて種々の解法を巧みに組み合わせたり直観的な経験則や知識を用いていることが多い。

調整可能なパラメータが n 個ある場合、チューニング可能な増減パターンが 2^n 通りある。しかし、改善するために必要なパラメータの組み合わせは、 n 個より少ないパラメータの組み合わせである場合が多い。また、必ずしも、1 通りではなく複数の組み合わせが存在し得る。一つの組み合わせをチューニングプラン (tuning plan) と呼ぶ。チューニングプランを前もって求めずにむやみにチューニングしようとする、 2^n 通り試みる必要がある。

このようなチューニングプランを求めるために、専門家の知識を整理する。この知識はかなり定性的であり定性的な表現の枠組みでうまく定式化できる。パラメータチューニングの第1段階として、定性的な推論で、チューニングに有効であるパラメータの増減を指示した複数のチューニングプランを求める。

この定性推論によるチューニングプランはパラメータの増減の指定のみであるから、具体的な増減量が不明なため、このままでは実用上不十分である。パラメータチューニングの第2段階として、前段の定性的なチューニングプランごとに、定量的な推論によってパラメータの定量的なチューニングプランを求める。

定性推論 [BOB85], [APT86], [DEK85], [MIZ89], [NIS88, 89] は、経済システムや電気回路などの定性的な挙動を表すモデルとして使われている。本稿で解説する待ち行列ネットワークへの適用は、定性推論の応用としては新規なものである。パラメータチューニングへの定性推論の主な適用例として、Rajagopalan, R. のターボジェットエンジンの対気速度と絞り弁の設定の増減関係を表す定性モデル [RAJ84] がある。本解説の方法では、動的なシステムのパラメータチューニングに、定性推論ばかりでなく相補的に定量推論を用いる。

一般に、定性推論が対象とするシステムは、

- ① その挙動の詳細が不明であり、かつ挙動を表す微分方程式などの式が不明なシステム、あるいは、
- ② 挙動がある程度分かっているが、その挙動を表す式が複雑かつ膨大、その解法が明確ではない、あるいは解くのに時間がかかるシステム、である。

定性推論では、このようなシステムの挙動について、数少ないパラメータ間の因果関係でモデル化し、数少ない定式で説明や推論を行う。

3.2 で述べるとおり、パラメータチューニングを行うために待ち行列ネットワークの挙動をそのまま表すと非線形の連立等式・不等式となる。この連立式は膨大であり解くのに時間がかかる。すなわち、待ち行列ネットワークのパラメータチューニングは②の範疇に属し、定性推論の適用が有効であると考えられる。

† Application of Qualitative Reasoning to Parameter Tuning Process by Kiyoshi ITOH (Laboratory of Information Science, Natural Sciences Center, Faculty of Science and Technology, Sophia Univ.) and Shinichi HONIDEN (Systems and Software Engineering Laboratory, Toshiba Corporation).

†† 上智大学理工学部一般科学研究室情報科学部門

††† (株)東芝システム・ソフトウェア技術研究所

* 本解説は、待ち行列ネットワークのパラメータチューニングに関する著者の論文 [ITO*ab, **ab], [SAW**ab] の内容を、待ち行列ネットワークの基礎事項も含めて定性推論の応用としての観点からできるだけ分かりやすく解説したものである。

2. 対象の待ち行列ネットワーク

2.1 単一のサーバ

単一のサーバ (server) の5つのパラメータを定義する。

λ : 単位時間当りにサーバに到着するエンティティの平均個数 (到着率: arrival rate)

μ : 単位時間当りにサーバで処理できるエンティティの平均個数 (サービス率: servicing rate)

t : 単位時間当りにサーバから出てくるエンティティの平均個数 (スループット: throughput)

ρ : サーバが実際にエンティティを処理する時間の割合 (稼働率: utilization rate)

q : サーバの前で処理を待つエンティティの平均個数 (待ち行列長: queue length)

図-1と図-2に示すとおり、 λ と μ が与えられるとその大小関係によりそのサーバの性能を表す ρ , t が決まる。

図-1は、到着率がサービス率より小さい場合、すなわち、平均到着時間間隔(到着率の逆数)が平均サービス時間(サービス率の逆数)より大きい場合、サーバはよどみなくエンティティを処理できる。サーバから出るエンティティの量(スループット)は、サーバへの到着率に等しくなる。稼働率は、

$$\rho = (\text{平均サービス時間}) / (\text{平均到着時間間隔})$$

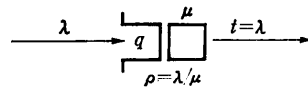


図-1 $\lambda < \mu$ の場合

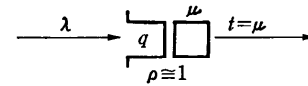


図-2 $\lambda \geq \mu$ の場合

$$= (1/\mu) / (1/\lambda) = \lambda/\mu \text{ となる。}$$

図-2は、到着率がサービス率より大きい場合、すなわち、平均到着時間間隔が平均サービス時間より小さい場合、サーバはよどみなくエンティティを処理できない。すなわち、あるエンティティに対するサービスを行っている間に次のエンティティがそのサーバに到着する。このサーバからのスループットは、サーバのサービス率に等しくなる。稼働率は、サーバが常に働いていることになるので、 $\rho \equiv 1$ となる。

2.2 待ち行列ネットワーク

待ち行列ネットワーク (queueing network; 以下QNと略記)は、複数のサーバが結合してネットワーク構造になったもの(たとえば図-3の“QN4”)である。QN内で二つのサーバが直列に結合する場合には、前段のサーバからのスループット t が、後段のサーバへの到着率 λ となる。二つのサーバからのスループット t_1 と t_2 が合流し

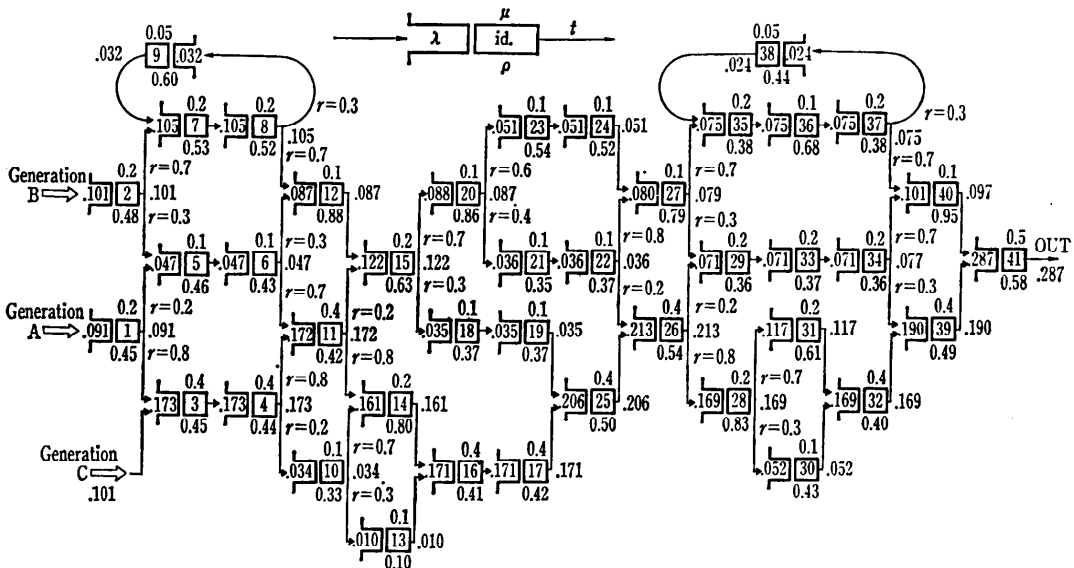


図-3 例題待ち行列ネットワーク“QN4”

て後段のサーバに入る場合には、 t_1+t_2 が後段のサーバへの到着率となる。あるサーバ s_1 のスループット t_1 が二つのサーバ s_2, s_3 に分岐する場合には、その分岐確率 (branching probability) を r , $1-r$ とすると、 s_2 と s_3 への到着率は rt_1 , $(1-r)t_1$ となる。QN 4 の各サーバの入出力関係を表す定式を図-4 に示す。

対象 QN は、エンティティが全て外部から到着し、処理済みになると外部へ出るオープン型の QN である。エンティティは外部から複数のサーバに到着できるが、内部ではそれらを区別しない。すなわち、内部ではエンティティを1種類と考えた単一フロー型の QN である。

2.3 稼働率とスループットの順次的計算

QN の構造 (サーバ同士の結合の仕方) が固定であり、外部からのエンティティの入力量、サーバからの分岐確率、及び個々のサーバのサービス率が既知であれば、外部からエンティティが入るサーバを起点として、順次、下流に向かって、個々のサーバの λ, ρ, t を、図-1 と図-2 に従って決定できる。すなわち、 $\lambda < \mu$ の場合は図-1 よりその到着率 λ をスループット t と決め、 $\lambda \geq \mu$ の場合をサーバのサービス率 μ をスループットと決める。その際、稼働率 ρ も決まる。このスループットは直下流のサーバへの到着率の一部一直列の場合はそのままの量 (図-4 (b)), 合流の場合は他との和 (図-4 (d)), 分岐の場合は分岐確率

を乗じた量 (図-4 (c)(e)) 一となり、順次、下流に向かって λ, ρ, t が決定できる。

実際のシステムを QN とみなしてモニタリング (monitoring) を行い、パラメータを実測する場合、ばらつきの存在や立ち上がり時の測定データの棄却方法次第により、正確には一致しないがほぼ上の議論に従う。

3. QN のパラメータチューニング

3.1 ボトルネック

QN 内で稼働率が1にきわめて近いサーバはボトルネック (bottleneck) である。このボトルネックサーバの待ち行列 (queue) は無限に成長する恐れがある。ボトルネックサーバが存在する場合、QN は非定常 (過負荷) 状態 (unstable or overloaded state) である。

稼働率が1に近くないが過大なサーバはボトルネックとなる可能性がある。サーバへの到着時間間隔やサーバのサービス時間が大きなばらつきをもつ場合、稼働率が1に近くなくても待ち行列長が時には急激に増加する危険性がある。フェールセーフ (failure to safety) の考え方を導入して、稼働率が0.7以上のサーバにボトルネックの可能性があると経験的に診断していることが多い。

また、稼働率が過大ではないが待ち行列長が過大であるサーバには一測定がそれほど長時間行われていないかもしれない、測定終了以降にエン

(a) 外部入力	$\lambda_{1,0} = t_{1,10} + t_{1,11}$	$t_{1,0,10} = t_{1,0} * r_{1,0,10}$
$\lambda_1 = ga, \lambda_2 = gb$	$\lambda_{1,0} = t_{1,10} + t_{1,11}$	$t_{1,0,11} = t_{1,0} * r_{1,0,11}$
(b) 直列	$\lambda_{2,0} = t_{1,7} + t_{1,8}$	$t_{1,1,10} = t_{1,1} * r_{1,1,10}$
$\lambda_4 = t_2, \lambda_6 = t_3, \lambda_8 = t_7$	$\lambda_{2,0} = t_{2,20} + t_{2,21}$	$t_{1,1,11} = t_{1,1} * r_{1,1,11}$
$\lambda_{17} = t_{1,6}, \lambda_{19} = t_{1,8}$	$\lambda_{27} = t_{2,27} + t_{2,28}$	$t_{1,1,10} = t_{1,1} * r_{1,1,10}$
$\lambda_{22} = t_{2,1}, \lambda_{24} = t_{2,2}$	$\lambda_{28} = t_{2,20} + t_{2,21}$	$t_{1,1,11} = t_{1,1} * r_{1,1,11}$
$\lambda_{22} = t_{2,9}, \lambda_{24} = t_{2,2}$	$\lambda_{27} = t_{2,27} + t_{2,28}$	$t_{1,1,10} = t_{1,1} * r_{1,1,10}$
$\lambda_{26} = t_{2,6}, \lambda_{27} = t_{2,8}$	$\lambda_{28} = t_{2,20} + t_{2,21}$	$t_{1,1,11} = t_{1,1} * r_{1,1,11}$
(c) 分岐	$\lambda_{29} = t_{2,29} + t_{2,30}$	$t_{1,0,20} = t_{1,0} * r_{1,0,20}$
$\lambda_6 = t_{2,9}, \lambda_{10} = t_{2,10}$	$\lambda_{30} = t_{2,27} + t_{2,28}$	$t_{1,0,21} = t_{1,0} * r_{1,0,21}$
$\lambda_{12} = t_{1,0,12}, \lambda_{13} = t_{1,1,13}$	$\lambda_{39} = t_{2,39} + t_{2,40}$	$t_{1,0,20} = t_{1,0} * r_{1,0,20}$
$\lambda_{20} = t_{1,9,20}, \lambda_{21} = t_{1,0,21}$	$\lambda_{40} = t_{2,40} + t_{2,37,40}$	$t_{1,0,21} = t_{1,0} * r_{1,0,21}$
$\lambda_{22} = t_{2,0,22}, \lambda_{23} = t_{2,0,23}$	$\lambda_{41} = t_{2,0} + t_{2,0}$	$t_{1,0,20} = t_{1,0} * r_{1,0,20}$
$\lambda_{20} = t_{2,0,20}, \lambda_{21} = t_{2,0,21}$	(e) スループットと分岐確率	$t_{1,0,21} = t_{1,0} * r_{1,0,21}$
$\lambda_{26} = t_{2,7,26}$	$t_{1,0} = t_{1,0} * r_{1,0}$	$t_{1,0,20} = t_{1,0} * r_{1,0,20}$
(d) 合流	$t_{1,0} = t_{1,0} * r_{1,0}$	$t_{1,0,21} = t_{1,0} * r_{1,0,21}$
$\lambda_2 = gc + t_{1,3}$	$t_{2,7} = t_{2,7} * r_{2,7}$	$t_{1,0,20} = t_{1,0} * r_{1,0,20}$
$\lambda_4 = t_{1,4} + t_{2,5}$	$t_{4,10} = t_{4,10} * r_{4,10}$	$t_{1,0,21} = t_{1,0} * r_{1,0,21}$
$\lambda_7 = t_{2,7} + t_8$	$t_{4,11} = t_{4,11} * r_{4,11}$	$t_{1,0,20} = t_{1,0} * r_{1,0,20}$
$\lambda_{11} = t_{4,11} + t_{6,11}$	$t_{4,11} = t_{4,11} * r_{4,11}$	$t_{1,0,21} = t_{1,0} * r_{1,0,21}$
$\lambda_{12} = t_{6,12} + t_{8,12}$	$t_{6,10} = t_{6,10} * r_{6,10}$	$t_{1,0,20} = t_{1,0} * r_{1,0,20}$
$\lambda_{14} = t_{1,0,14} + t_{11,14}$	$t_{6,11} = t_{6,11} * r_{6,11}$	$t_{1,0,21} = t_{1,0} * r_{1,0,21}$
	$t_{6,12} = t_{6,12} * r_{6,12}$	(f) 分岐確率
		$r_i, j + r_i, k = 1 (j \neq k)$

図-4 各サーバの入出力の定式

ティティが過大に到着して稼働率が大きくなる危険性がある。この過大な待ち行列長は後に徐々に解消されるかもしれないが、ボトルネックの発生の事前防止というフェールセーフの考え方で、測定時に平均待ち行列長が1以上のサーバにボトルネックの可能性があると診断する。

3.2 パラメータチューニングの難しさ

ボトルネックを改善するチューニングプラン(通常一つ以上存在)は、QNの構造のみを固定としほかの全てのパラメータを可変とした場合に、稼働率 ρ が0.7を下回るためにどのパラメータを変更すればよいか、という問題の、択一的に選択可能な解である。

図-3のQN4で、その構造は固定であり、外部からのエンティティの入力量、サーバからの分岐確率、及び各サーバのサービス率が既知であるとする。実測、シミュレーション、上流からの順次的な計算(2.3)などの方法があるが、この例ではシミュレーションにより各サーバの稼働率、スループット、到着率が求められた。この後、ボトルネックを改善するために、構造のみを固定としてパラメータチューニングする。

QN4ではボトルネックの可能性のあるサーバとして、 $\rho \geq 0.7$ のs12, s14, s20, s27, s28, s40がある。改善の例としてs20をあげる。このs20の稼働率 ρ_{20} が0.7を下回るようにするためには、s20にのみ着目すると、

- (1) サービス率 μ_{20} を増加する、あるいは
- (2) 到着率 λ_{20} を減少する という方法をとる、

(新しい ρ_{20}) = (新しい λ_{20}) / (新しい μ_{20}) < 0.7とする。

単一サーバではなく複数のサーバが結合しているQNでは、この(1), (2)の適用のみでは不十分である。

(1)の方法を実施すると、これまで稼働率が高いためs20の待ち行列に滞留する傾向のあったエンティティが、s20のサービスを以前より多く受けることが可能となり、s20のスループットが増加する。s20の直下流のs21, s23への到着率を増加させ、さらにその下流のサーバへの到着率も増加させる。s20の下流にありすでにボトルネックの可能性のあったサーバ(s27, s28, s40)の稼働率をさらに増加させるとともに、ボトルネックでは

なかったサーバもその可能性をもつものにしてしまう危険性がある。

(2)の方法を実施するためには、s20から上流に遡行し外部からの入口までの経路上の

(2-1) 分岐確率を減少する、あるいは、

(2-2) 外部からの入力量を減少する、方法をとる。

(2-1)の方法を実施するためには、上流への経路上に複数の分岐点が存在する場合、どの分岐点の分岐確率を減少するかを決定しなければならない。s20から上流への経路には、3個の分岐点(s15, s11, s4)が存在する。また、その分岐点の他方は分岐確率が増加するので下流の到着率が増加することになり、(1)と同様に下流にあるボトルネックに留意しなければならない。

チューニングプランを得るために2.3の順次的計算式を使うと、上記の(1), (2), (2-1), (2-2)に現れたパラメータを変数とした式が立式できる。この式には、(2-1)の複数の分岐確率を変数とした積項が含まれる。また、改善後の稼働率が0.7を下回るから全てのサーバについて $\rho = \lambda/\mu < 0.7$ という式を与える必要がある。立式は非線形の連立等式・不等式となる。

この連立式は膨大であり解くのに時間がかかる。すなわち、待ち行列ネットワークのパラメータチューニングは、1.の②の範疇に属し、定性推論の適用が有効であると考えられる問題である。

待ち行列理論(queueing theory)(たとえば文献[KLE 75], [GEL 80])では、平均値だけではなく、到着時間間隔やサービス時間にさまざまな分布を仮定し、稼働率、待ち時間などの性能パラメータの厳密解や近似解を求めるものが多い。性能を改善するためにパラメータ群をチューニングする問題を対象とすることはあまり多くない。

非線形の連立式を解かないで、択一的に選択可能なチューニングプランを可能なかぎり効率的に得たい。このために、QNの構造と性能パラメータに関する経験則を用いて定性推論の枠組みでチューニングプランを列挙する方法を4.で述べる。

3.3 チューニングプランの意味

一つのチューニングプランは、①外部からの入力量を減少、②分岐確率を減少、③サーバのサービス率を増加、などのパラメータ変更の組み合わせ

となる。これらのパラメータを全て可変としており、そのため、どのチューニングプランも定性的には選択可能である。

①はシステムに対する負荷の減少を意味する。
 ②は、エンティティがこれまで流れていたパスの部分的な変更を意味する。分岐後が並列なパスであれば、分岐確率の変更は各パスへの負荷を変更することになり、②は実現可能である。並列なパスではない場合、分岐確率を変更してよいかどうかはシステムによる。

③の説明として、図-5 (a) のボトルネックサーバ α を改善するためにサーバの稼働率を 0.7 に減少させる例を述べる。図-5 (b) では、サービス率が 2 倍 (サービス時間が 1/2) のサーバに置き換えて同じ仕事をさせる。このようなサーバの置き換えが現実的ではない場合、図-5 (c) のサーバのパイプライン化や図-5 (d) のサーバの並列化が、実際的な方法である。図-5 (c) は、14h かかる仕事を 7h と 7h に直列に分割できる場合、二つのサーバをパイプライン化して並べる方法である。図-5 (d) は、14h で行うサーバを

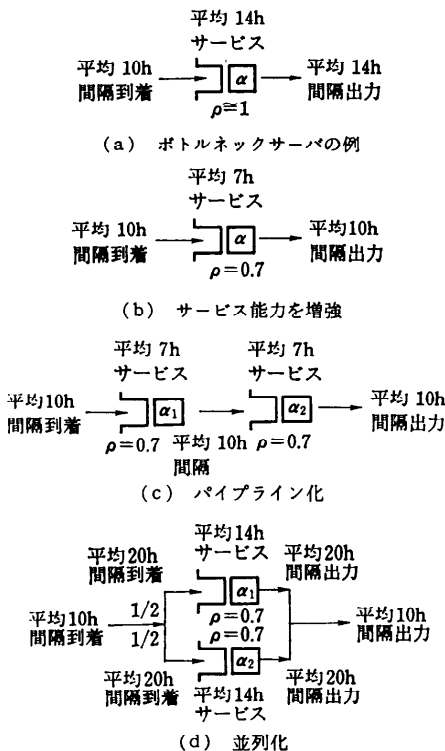


図-5 ボトルネックサーバのサービス能力を増強する方法
 到着時間間隔、出力時間間隔、サービス時間の逆数は、到着率、スループット、サービス率

もう一つ用意しておのおのへ同じ分岐確率で分岐させる方法である。パイプライン化と並列化についての一般的な例は本解説では省略する。

4. 定性的なパラメータチューニング

4.1 ボトルネックランドマーク

専門家の経験則として、稼働率 0.7 をボトルネックランドマーク (bottleneck landmark: BL) とし、また、行列長 1 を待ち行列長に関するボトルネックランドマーク (QBL) とし、 $\rho \geq BL$, あるいは、 $q \geq QBL$ ならばボトルネックの可能性があると看做する。

定性推論記法で表現する。[ρ] と [q] は、おのおの、BL と QBL を原点として評価した定性値 (+, 0, -) である。

[ρ] = + : ボトルネックの可能性ある。この場合フェールセーフの考え方で $d\rho = +$ と考える。

[q] = + : ボトルネックの可能性ある。この場合 $dq = +$ と考える。

4.2 単一サーバの定性式

サーバの挙動を定性式で表す。ボトルネックの可能性の有無で用いる式の種類が異なることに特長がある。以下の式をサーバの定性挙動式 (qualitative behavior expression: QLBE) と呼ぶ。ボトルネックの可能性ある場合 (下記の b)) を特に定性改善式 (qualitative bottleneck improvement expression: QL-BIE) と呼ぶ。

a) $\rho < BL$ の場合 ($\lambda < \mu$ の場合: 図-1 参照) すなわち

[ρ] = - (ρ が BL より小さい) の場合

a-1) λ について

$$d\rho = \pm \leftarrow d\lambda = \pm (\lambda \text{ と } \rho \text{ は同じ方向に増減})$$

$$dt = \pm \leftarrow d\lambda = \pm (\lambda \text{ と } t \text{ は同じ方向に増減})$$

a-2) μ について

$$d\rho = \mp \leftarrow d\mu = \pm (\mu \text{ と } \rho \text{ は反対方向に増減})$$

$$dt = 0 \leftarrow d\mu = \pm (\mu \text{ の増減に } t \text{ は追従せず})$$

b) $\rho \geq BL$ の場合 ($\lambda \geq \mu$ の場合: 図-2 参照)

すなわち

[ρ] = + or 0 (ρ が BL 以上、すなわちボトルネックの可能性ある)

または、 $q \geq QBL$, すなわち [q] = + or 0 の場合

b-1) λ について

$$d\rho = - \leftarrow d\lambda = - (\lambda \text{ の減少に } \rho \text{ は追従})$$

$$dt = 0 \leftarrow d\lambda = - (\lambda \text{ の減少に } t \text{ は追従せず})$$

b-2) μ について

$$d\rho = - \leftarrow d\mu = + (\mu \text{ が増加すると } \rho \text{ は減少})$$

$$dt = + \leftarrow d\mu = + (\mu \text{ が増加すると } t \text{ は増加})$$

b-3) q について

$$dq = - \leftarrow d\rho = - (\rho \text{ の減少に } q \text{ は追従})$$

$$d\rho = - \leftarrow dq = - (q \text{ の減少に } \rho \text{ は追従})$$

4.3 経験則を用いない定性挙動推論

3.2 で述べたとおり、サーバの λ の変化のためには、上流のサーバからの t の変化が必要であ

る。あるサーバの t の増加は下流のサーバへの λ の増加につながるので、下流のボトルネックへの留意が必要である。チューニングの一つの方法は、全てのサーバに対して定性式で立式して、それらを組み合わせて解く方法である。

図-3 のQN4で、ボトルネックの可能性のあるサーバの定性改善式(QL-BIE)及びそのほかのサーバの定性挙動式(QLBE)を図-6に示す。QN4に対する式の個数は、図-6に示すとおり

for $i = \{s12, s14, s20, s27, s28, s40\}$
(which may be bottleneck)

• qualitative state: $[p_i] = +$
 $d\rho_i = +$

• QL-BIE;

$$d\rho_i = - \leftarrow d\lambda_i = -$$
$$dt_i = 0 \leftarrow d\lambda_i = -$$
$$d\rho_i = - \leftarrow d\mu_i = +$$
$$dt_i = + \leftarrow d\mu_i = +$$

subtotal: 24 expressions
= 4 expressions \times 6 (servers)
total: 304 expressions

for $j = \{s1, s2, s3, s4, s5, s6, s7, s8, s9, s10, s11, s13,$
 $s15, s16, s17, s18, s19, s21, s22, s23, s24, s25,$
 $s26, s29, s30, s31, s32, s33, s34, s35, s36, s37,$
 $s38, s39, s41\}$

• qualitative state: $[p_j] = -$

• QLBE: $d\rho_j = \pm \leftarrow d\lambda_j = \pm$

$$dt_j = \pm \leftarrow d\lambda_j = \pm$$
$$d\rho_j = \mp \leftarrow d\mu_j = \pm$$
$$dt_j = 0 \leftarrow d\mu_j = \pm$$

subtotal: 280 expressions
= 8 expressions \times 35 (servers)

図-6 QN4の各サーバの定性挙動式と定性改善式

(a) 外部入力

$$d\lambda_i = \pm \leftarrow dga = \pm$$

$$d\lambda_j = \pm \leftarrow dgb = \pm$$

subtotal: 4 expressions
= 2 expressions \times 2 (servers)

(b) 直列

$$d\lambda_s = \pm \leftarrow dt_s = \pm$$

$$d\lambda_t = \pm \leftarrow dt_t = \pm$$

$$d\lambda_r = \pm \leftarrow dt_r = \pm$$

$$d\lambda_{17} = \pm \leftarrow dt_{17} = \pm$$

$$d\lambda_{19} = \pm \leftarrow dt_{19} = \pm$$

$$d\lambda_{23} = \pm \leftarrow dt_{23} = \pm$$

$$d\lambda_{25} = \pm \leftarrow dt_{25} = \pm$$

$$d\lambda_{27} = \pm \leftarrow dt_{27} = \pm$$

$$d\lambda_{31} = \pm \leftarrow dt_{31} = \pm$$

$$d\lambda_{33} = \pm \leftarrow dt_{33} = \pm$$

$$d\lambda_{35} = \pm \leftarrow dt_{35} = \pm$$

$$d\lambda_{37} = \pm \leftarrow dt_{37} = \pm$$

subtotal: 22 expressions
= 2 expressions \times 11 (servers)

(c) 分岐

$$d\lambda_s = \pm \leftarrow dt_{1,s} = \pm$$

$$d\lambda_{10} = \pm \leftarrow dt_{10,10} = \pm$$

$$d\lambda_{13} = \pm \leftarrow dt_{13,13} = \pm$$

$$d\lambda_{18} = \pm \leftarrow dt_{18,18} = \pm$$

$$d\lambda_{20} = \pm \leftarrow dt_{20,20} = \pm$$

$$d\lambda_{21} = \pm \leftarrow dt_{21,21} = \pm$$

$$d\lambda_{22} = \pm \leftarrow dt_{22,22} = \pm$$

$$d\lambda_{28} = \pm \leftarrow dt_{28,28} = \pm$$

$$d\lambda_{30} = \pm \leftarrow dt_{30,30} = \pm$$

$$d\lambda_{31} = \pm \leftarrow dt_{31,31} = \pm$$

$$d\lambda_{37} = \pm \leftarrow dt_{37,37} = \pm$$

total: 266 expressions

subtotal: 22 expressions
= 2 expressions \times 11 (servers)

(d) 合流

$$d\lambda_s = \pm \leftarrow dgc = \pm; dt_{1,s} = \pm$$

$$d\lambda_t = \pm \leftarrow dt_{t,s} = \pm; dt_{t,s} = \pm$$

$$d\lambda_r = \pm \leftarrow dt_{r,s} = \pm; dt_{r,s} = \pm$$

$$d\lambda_{11} = \pm \leftarrow dt_{1,11} = \pm; dt_{1,11} = \pm$$

$$d\lambda_{12} = \pm \leftarrow dt_{1,12} = \pm; dt_{1,12} = \pm$$

$$d\lambda_{14} = \pm \leftarrow dt_{1,14} = \pm; dt_{1,14} = \pm$$

$$d\lambda_{15} = \pm \leftarrow dt_{1,15} = \pm; dt_{1,15} = \pm$$

$$d\lambda_{16} = \pm \leftarrow dt_{1,16} = \pm; dt_{1,16} = \pm$$

$$d\lambda_{23} = \pm \leftarrow dt_{1,23} = \pm; dt_{1,23} = \pm$$

$$d\lambda_{26} = \pm \leftarrow dt_{1,26} = \pm; dt_{1,26} = \pm$$

$$d\lambda_{27} = \pm \leftarrow dt_{1,27} = \pm; dt_{1,27} = \pm$$

$$d\lambda_{29} = \pm \leftarrow dt_{1,29} = \pm; dt_{1,29} = \pm$$

$$d\lambda_{30} = \pm \leftarrow dt_{1,30} = \pm; dt_{1,30} = \pm$$

$$d\lambda_{32} = \pm \leftarrow dt_{1,32} = \pm; dt_{1,32} = \pm$$

$$d\lambda_{37} = \pm \leftarrow dt_{1,37} = \pm; dt_{1,37} = \pm$$

$$d\lambda_{41} = \pm \leftarrow dt_{1,41} = \pm; dt_{1,41} = \pm$$

subtotal: 68 expressions
= 4 expressions \times 17 (servers)

(e) スループットと分岐確率

$$dt_{1,s} = \pm \leftarrow dt_s = \pm; dr_{1,s} = \pm$$

$$dt_{1,t} = \pm \leftarrow dt_t = \pm; dr_{1,t} = \pm$$

$$dt_{1,r} = \pm \leftarrow dt_r = \pm; dr_{1,r} = \pm$$

$$dt_{1,10} = \pm \leftarrow dt_{10} = \pm; dr_{1,10} = \pm$$

$$dt_{1,11} = \pm \leftarrow dt_{11} = \pm; dr_{1,11} = \pm$$

$$dt_{1,11} = \pm \leftarrow dt_{11} = \pm; dr_{1,11} = \pm$$

$$dt_{1,11} = \pm \leftarrow dt_{11} = \pm; dr_{1,11} = \pm$$

$$dt_{1,11} = \pm \leftarrow dt_{11} = \pm; dr_{1,11} = \pm$$

total: 266 expressions

$$dt_{1,s} = \pm \leftarrow dt_s = \pm; dr_{1,s} = \pm$$

$$dt_{1,12} = \pm \leftarrow dt_{12} = \pm; dr_{1,12} = \pm$$

$$dt_{1,13} = \pm \leftarrow dt_{13} = \pm; dr_{1,13} = \pm$$

$$dt_{1,14} = \pm \leftarrow dt_{14} = \pm; dr_{1,14} = \pm$$

$$dt_{1,15} = \pm \leftarrow dt_{15} = \pm; dr_{1,15} = \pm$$

$$dt_{1,16} = \pm \leftarrow dt_{16} = \pm; dr_{1,16} = \pm$$

$$dt_{1,18} = \pm \leftarrow dt_{18} = \pm; dr_{1,18} = \pm$$

$$dt_{1,20} = \pm \leftarrow dt_{20} = \pm; dr_{1,20} = \pm$$

$$dt_{1,21} = \pm \leftarrow dt_{21} = \pm; dr_{1,21} = \pm$$

$$dt_{1,22} = \pm \leftarrow dt_{22} = \pm; dr_{1,22} = \pm$$

$$dt_{1,23} = \pm \leftarrow dt_{23} = \pm; dr_{1,23} = \pm$$

$$dt_{1,25} = \pm \leftarrow dt_{25} = \pm; dr_{1,25} = \pm$$

$$dt_{1,27} = \pm \leftarrow dt_{27} = \pm; dr_{1,27} = \pm$$

$$dt_{1,28} = \pm \leftarrow dt_{28} = \pm; dr_{1,28} = \pm$$

$$dt_{1,30} = \pm \leftarrow dt_{30} = \pm; dr_{1,30} = \pm$$

$$dt_{1,31} = \pm \leftarrow dt_{31} = \pm; dr_{1,31} = \pm$$

$$dt_{1,33} = \pm \leftarrow dt_{33} = \pm; dr_{1,33} = \pm$$

$$dt_{1,35} = \pm \leftarrow dt_{35} = \pm; dr_{1,35} = \pm$$

$$dt_{1,37} = \pm \leftarrow dt_{37} = \pm; dr_{1,37} = \pm$$

$$dt_{1,41} = \pm \leftarrow dt_{41} = \pm; dr_{1,41} = \pm$$

$$dt_{1,41} = \pm \leftarrow dt_{41} = \pm; dr_{1,41} = \pm$$

$$dt_{1,41} = \pm \leftarrow dt_{41} = \pm; dr_{1,41} = \pm$$

$$dt_{1,41} = \pm \leftarrow dt_{41} = \pm; dr_{1,41} = \pm$$

$$dt_{1,41} = \pm \leftarrow dt_{41} = \pm; dr_{1,41} = \pm$$

$$dt_{1,41} = \pm \leftarrow dt_{41} = \pm; dr_{1,41} = \pm$$

$$dt_{1,41} = \pm \leftarrow dt_{41} = \pm; dr_{1,41} = \pm$$

$$dt_{1,41} = \pm \leftarrow dt_{41} = \pm; dr_{1,41} = \pm$$

$$dt_{1,41} = \pm \leftarrow dt_{41} = \pm; dr_{1,41} = \pm$$

$$dt_{1,41} = \pm \leftarrow dt_{41} = \pm; dr_{1,41} = \pm$$

$$dt_{1,41} = \pm \leftarrow dt_{41} = \pm; dr_{1,41} = \pm$$

$$dt_{1,41} = \pm \leftarrow dt_{41} = \pm; dr_{1,41} = \pm$$

$$dt_{1,41} = \pm \leftarrow dt_{41} = \pm; dr_{1,41} = \pm$$

$$dt_{1,41} = \pm \leftarrow dt_{41} = \pm; dr_{1,41} = \pm$$

$$dt_{1,41} = \pm \leftarrow dt_{41} = \pm; dr_{1,41} = \pm$$

$$dt_{1,41} = \pm \leftarrow dt_{41} = \pm; dr_{1,41} = \pm$$

$$dt_{1,41} = \pm \leftarrow dt_{41} = \pm; dr_{1,41} = \pm$$

$$dt_{1,41} = \pm \leftarrow dt_{41} = \pm; dr_{1,41} = \pm$$

$$dt_{1,41} = \pm \leftarrow dt_{41} = \pm; dr_{1,41} = \pm$$

$$dt_{1,41} = \pm \leftarrow dt_{41} = \pm; dr_{1,41} = \pm$$

$$dt_{1,41} = \pm \leftarrow dt_{41} = \pm; dr_{1,41} = \pm$$

$$dt_{1,41} = \pm \leftarrow dt_{41} = \pm; dr_{1,41} = \pm$$

$$dt_{1,41} = \pm \leftarrow dt_{41} = \pm; dr_{1,41} = \pm$$

図-7 各サーバの入出力の定性式

図4 (a)~(f)の定量式に対応, (d)の+, (e)の*は; (or)になる

304本である。ボトルネックの改善を行う場合には、これらの式に加えてサーバ間の入出力情報(図-7の入出力の定性式): 266本を与えなければならない。しかし、これらを全て組み合わせて定性挙動推論により解くと、状態数の爆発が起きて効率的ではない。

4.4 経験則に基づく定性挙動推論

全ての立式(計570本)を組み合わせて解くのではなく、専門家の経験則(浅い知識)を導入する。それは、サーバの結合による部分形状ごとの性能パラメータ間の増減関係を表した定性改善式である。

著者らは、この定性改善式を扱い定性推論によりチューニングプランを列挙するエキスパートシステムBDES(Bottleneck Diagnosis Expert System)を開発した。

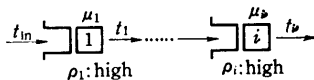
BDESには、ボトルネック改善のための9つの知識(下流への影響, 直列, 合流, 分岐, ループなど)がある。この知識により、使用する性能パラメータの個数が、パラメータ全てに着目する場合よりかなり減少する。

図-8は、ボトルネックサーバの改善のためその μ を増加するとスループットが増加するため、下流の過大な ρ を減少させる知識である。

図-9の直列型(tandem)のサーバs2の ρ_2 を改善するために、 t_1 と μ_2 のバランスを変える必要がある。そのために接続する他のサーバのパラメータの大きさを変える。図-9は、一つ上流のサーバs1がボトルネックでない場合、 t_{in} を減少して ρ_2 を改善する知識である。

図-10は、合流の場合、合流入力の大きさにより改善すべきパラメータを決める知識である。

図-11の分岐の知識は、分岐確率を修正する場

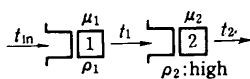


ボトルネックサーバのサービス率を増加する場合、下流のサーバの稼働率を減少する必要あり

$$d\rho_1 = - \leftarrow d\mu_1 = +$$

$$d\rho_2 = - \leftarrow d\rho_1 = -$$

図-8 下流への影響



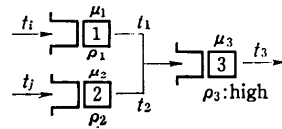
$$d\rho_2 = - \leftarrow dt_{in} = -$$

図-9 直列型の知識

合と、直上流のスループットを修正する場合に分かれる。

図-12の知識はまず、ループの機能を変えないために、ループからの分岐確率の修正を禁止する。また、ループからのスループットを受け取るボトルネックサーバs3を改善するためには、ループ内部のパラメータではなくループへの到着量のみを減少する。すなわち、 $d\rho_3 = - \leftarrow dt_{in} = -$ という定性改善式のみを使用する知識である。

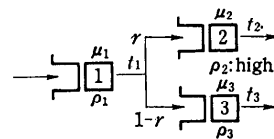
図-13は、ループ内のボトルネックを改善するために、ループへの到着率を減少する知識である。ボトルネックサーバのサービス率を増加するとスループットが増加し、これがループして自身



$t_i > t_j$ の場合
 $d\rho_3 = - \leftarrow dt_i = -$

$t_i = t_j$ の場合
 $d\rho_3 = - \leftarrow dt_i = -$
 $d\rho_3 = - \leftarrow dt_j = -$

図-10 合流型の知識

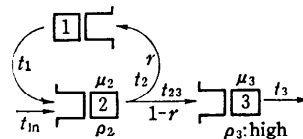


分岐確率が修正可能な場合
 $d\rho_3 = - \leftarrow dr = -$

s_3 とその下流のサーバの稼働率にも留意

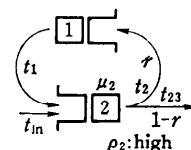
修正不可能な場合
 $d\rho_3 = - \leftarrow dt_i = -$

図-11 分岐型の知識



ループからの分岐確率は修正禁止
 $d\rho_3 = - \leftarrow dt_{in} = -$

図-12 ループ直後のボトルネック



サービス率の増加は禁止
 $d\rho_3 = - \leftarrow dt_{in} = -$

図-13 ループ内のボトルネック

への入力の一部となり、ボトルネックの改善にはつながらない。

以上の知識と定性改善式は、専門家のもついわゆる浅い知識である。4.2で述べたサーバごとの定性挙動式、定性改善式及び入出力の定性式はいわゆる深い知識である。一般に浅い知識の妥当性は簡単には証明ができないといわれているが、文献 [ITO 90 a] に示すとおり、本節の浅い知識は深い知識により証明可能である。

図-3のQN4のボトルネックサーバs20に対する、BDESの定性挙動推論プロセスを図-14に示す。図-6と図-7に比べ式の本数はかなり減少し、状態数の爆発は抑制される。

図-14の最上部のブロック1は、ボトルネックs20を改善したいというチューニングのゴール

($d\rho_{20} = -$)である。ブロック1からブロック2・3への分岐はs20に対する定性改善式(4.2)の適用である。ブロック2は、図-8の下流への影響に関する知識である。ブロック3は、図-11の分岐型の知識によりブロック4と5に分岐する。以下同様に知識が適用される。

以上の推論の結果、ブロック2, 4, 9, 15, 20, 22, 24の下に、改善のゴール($d\rho_{20} = -$)を実現するための択一的に選択可能な定性的なチューニングプラン(plan 1~7)が示されている。たとえば、プラン1は、s20の改善のために、s20自体という局所的なサービス率の増加のほかに、広域的にs27, s28, s40の稼働率の減少を指示している。

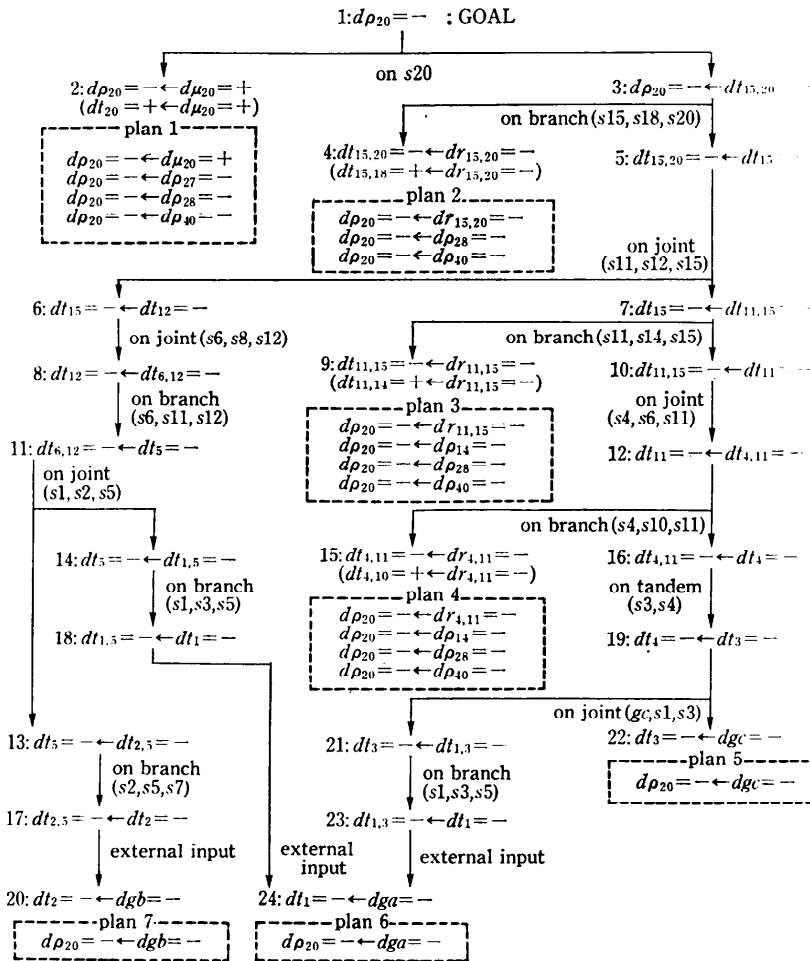


図-14 QN4の改善のための定性推論

5. 定量的なパラメータチューニング

著者らは、ボトルネック改善のための性能パラメータの定量推論を行うエキスパートシステム BIES (Bottleneck Improvement Expert System) を開発した。図-14 に示したとおり、一つのボトルネックに対して定性的なチューニングプランが複数得られたが、これらは択一的に選択可能である。ユーザが選んだ一つの定性チューニングプランに対して、BIES は改善後のパラメータの定量値を推論する。すなわち低いサービス率を定量的に増加させたり、高い到着率や高い分岐確率を定量的に減少させる。さらに、このような一つのサーバのパラメータチューニングの後、ほかのサーバの新しい到着率や新しいスループットを自動的に算出する。

図-15 は、図-14 の BDES の 7 つの定性的な

チューニングプラン全てに対して定量的なチューニングを行ったプロセスである。定性チューニングプラン 1 に対して、BIES がきめ細かな定量推論を行うと定量チューニングプランが 8 個提示される。また、定性チューニングプラン 5 は外部入力 gc を減少することを示しているが、BIES の定量チューニングプランでは、きめ細かな定量推論を行うと、この gc の改善のみでは $s20$ のボトルネックを解消できないことが示されている。定性推論を行わないで直接、複数の定量チューニングプランを列挙することは必ずしも容易ではない。このように定性推論と定量推論を相補的に使用することが有効である。

図-14 の全ての定性チューニングプランや、このおののに対する図-15 の全ての定量チューニングプランを、BDES と BIES は順次全て出力する。たとえばあるパラメータが変更禁止の固定パ

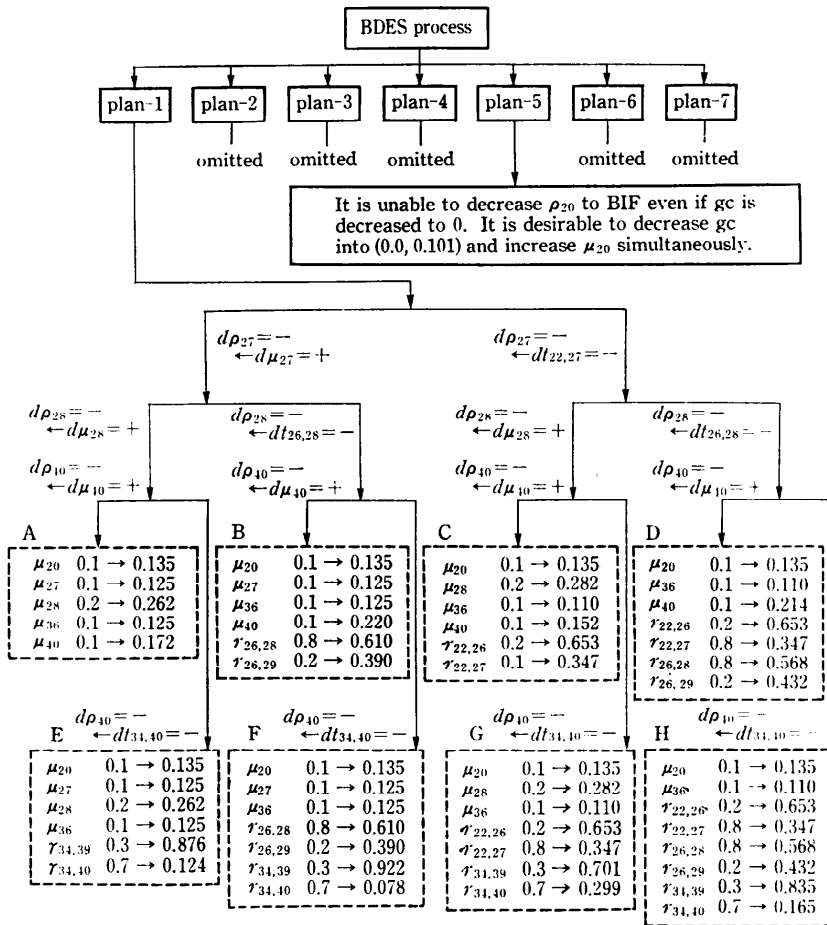


図-15 QN 4 の改善のための定量推論

ラメータであれば、そのような改善プランをユーザは選択せず、次のプランの提示を促す。また選択したプランが不適切であるとユーザが判断したら、BDES と BIES はパラメータを元の値に自動的に復帰させるバックトラックを行う。

6. おわりに

本解説では、定性推論と定量推論を組み合わせることで QN を例にあげてそのボトルネックを改善するためのパラメータチューニングの方法を解説した。この方法は、複数種類のエンティティが流れる多重フロー型 QN、エンティティ同士が同期を取る同期型 QN、時間/確率ペトリネットなどのボトルネック改善にも適用した^{[MAT90],[SHI90]}。定性推論に基づいて並列オブジェクトプログラムの性能デバガにも適用した^[HON90ab]。

待ち行列ネットワーク理論では、外部からの到着率、サーバのサービス率、分岐確率を与えて、各サーバの稼働率や待ち行列長を求めるという順問題 (direct problem: 初期条件を与えても対象を支配する方程式を解く問題) を解くことが多い。パラメータチューニングとは、本解説のようにたとえば稼働率を 0.7 より小さくするためには、到着率、サービス率、分岐確率などをどのように設定すればよいかという逆問題 (inverse problem)^{[SUU86],[SIM90]} を解くことである。待ち行列ネットワークの逆問題を解析的な最適化手法で解く方法も数少ない(たとえば [WAD90])。本解説の手法は、この逆問題を解析的ではない定性推論で解くヒューリスティックな解法である。パラメータチューニングを逆問題としてとらえ定性推論により解く方法は、ほかの分野でもあまり報告されていない。しかし、システムの構造がネットワークやグラフで表現されその上を何かが伝播するシステム (力学/電気回路, 経済システムなど数多い) のパラメータチューニングにきわめて有効であると考えている。

謝辞 BDES と BIES の共同研究者である上智大学院生沢村淳氏(現, ソニー(株)), 志田圭介氏に感謝します。

参考文献

- [APT 86] Apte, C. et al.: Using Qualitative Reasoning to Understand Financial Arithmetic Proc. AAAI (1986).
- [BOB 85] Bobrow, D.G. et al. ed.: Qualitative Reasoning about Physical Systems, MIT Press (1985).
- [DEK 85] De Kleer, J.: How Circuits Work, in [BOB 85].
- [GEL 80] Gelenbe, E. et al.: Analysis and Synthesis of Computer Systems, Academic Press (1980).
- [HON 90 a] 本位田, 内平, 伊藤: 定性推論と定量推論による待ち行列ネットワークの性能改善エキスパートシステム BDES & BIES—その 5 並列プログラムへの適用一, 4th NC JSAI (July 1990).
- [HON 90 b] Honiden, S. and Itoh, K. et al.: Rapid Prototyping Method for Performance Design in Realtime Systems, InfoJapan '90, Vol. 1, pp. 103-110 (Oct. 1990).
- [ITO 89 a] Itoh, K. et al.: Knowledge-Based Parameter Tuning for Queueing Network Type System—A New Application of Qualitative Reasoning, Proc. CAPE '89, pp. 209-216 (Oct. 1989).
- [ITO 89 b] Itoh, K. et al.: Knowledge-Based Diagnosis on Simulation Result for Queueing Network, Proc. Beijing I.C. on System Simulation and Scientific Computing (Oct. 1989).
- [ITO 90 a] 伊藤, 本位田, 沢村, 志田: 定性推論と定量推論を導入した待ち行列ネットワークのボトルネック診断と改善法, 人工知能学会誌, Vol. 5, No. 1 (1990).
- [ITO 90 b] Itoh, K., Honiden, S. et al.: Role of Qualitative and Quantitative Reasoning in Diagnosis and Improvement for Queueing Network Bottleneck, InfoJapan '90, Vol. 2, pp. 171-178 (Oct. 1990).
- [KLE 75] Kleinrock, L.: Queueing Systems, John Wiley & Sons, Inc. (1975).
- [MAT 90] 松永, 志田, 本位田, 伊藤: 定性推論と定量推論による待ち行列ネットワークの性能改善エキスパートシステム BDES & BIES—その 4 多重フロー型待ち行列ネットワークのエキスパートシステム—4th NC JSAI (July 1990).
- [MIZ 89] 溝口, 古川, 安西編: 定性推論, 共立 (1989).
- [NIS 88] 西田: 定性推論に関する最近の研究動向, 情報処理, Vol. 29, No. 9 (Sep. 1988).
- [NIS 89] 西田: 定性推論の基礎, 人工知能学会誌, Vol. 4, No. 5, pp. 522-527 (1988).
- [RAJ 84] Rajagopalan, R.: Qualitative Modeling in the Turbojet Engine Domain, Proc. AAAI (1984).
- [SAW 89 a] 沢村, 本位田, 伊藤: 定性推論を導入した待ち行列ネットワークのボトルネック診断, 情報処理学会知識工学と人工知能研究会 (Jan. 1989).
- [SAW 89 b] 沢村, 本位田, 志田, 伊藤: 知識工学的手法を用いた待ち行列ネットワークのボトルネック診断, 情報処理学会論文誌, Vol. 30, No. 8, pp. 990-1002 (Aug. 1989).
- [SHI 90] 志田, 伊藤, 本位田, 早瀬: 定性/定量推論による同期型待ち行列ネットワークのボトルネック診断と改善, 情報処理学会知識工学と人工知能研究会 (Jan. 1990).
- [SIM 90] コンピュータ利用による逆問題解析, シミュレーション, 9.1 (Mar. 1990).
- [SUU 86] 逆問題, 数理科学, 274 (Apr. 1986).
- [WAD 90] 和田, 米田: 待ち行列網の逆問題とその設備計画への応用, SSOR 90 (1990).

(平成 2 年 9 月 10 日受付)



伊藤 謙 (正会員)

1951年生。1974年京都大学工学部情報工学科卒業。1979年同大学院情報工学専攻博士課程修了。京都大学工学博士。1979年より上智大学に勤務。1985年より助教授。現在、同理工学部一般科学研究室情報科学部門所属。ソフトウェア工学、シミュレーション手法、定性理論の待ち行列ネットワークへの応用、三面図からのソリッドモデルの自動合成法の研究に従事。情報処理学会学会誌編集委員。「仕様記述の効率的適用と評価」研究グループ幹事。ISO/TC 184/SC 2 WG 3 (産業用ロボットの安全性)議長団。著書「ソフトウェア開発のためのプロトタイプングツール」(共著)、「システムプログラム」、訳書「並行処理と Unix」(共訳)。人工知能学会、IEEE、ACM など各会員。



本位田真一 (正会員)

1953年生。1976年早稲田大学理工学部電気工学科卒業。1978年同大学院理工学研究科電気工学専攻修士課程修了。工学博士。同年(株)東芝入社。現在、同社システム・ソフトウェア技術研究所研究主務。早稲田大学非常勤講師。主として、ソフトウェア工学、人工知能の研究に従事。ソフトウェアの基礎理論に興味をもつ。1986年情報処理学会論文賞受賞。著書「ソフトウェア開発のためのプロトタイプング・ツール」(共著)、「KE 養成講座② エキスパートシステム基礎技術」(共著)、「オブジェクト指向システム分析」(共訳)、人工知能学会、日本ソフトウェア科学会、電気学会、AAAI 各会員。

