

Performance Improvement of Distributed Systems by Autotuning of the Configuration Parameters^{*}

ZHANG Fan (张帆)¹, CAO Junwei (曹军威)^{2,3,**},
LIU Lianchen (刘连臣)^{1,3}, WU Cheng (吴澄)^{1,3}

1. National CIMS Engineering and Research Center, Tsinghua University, Beijing 100084, China;
2. Research Institute of Information Technology, Tsinghua University, Beijing 100084, China;
3. Tsinghua National Laboratory for Information Science and Technology, Beijing 100084, China

Abstract: The performance of distributed computing systems is partially dependent on configuration parameters recorded in configuration files. Evolutionary strategies, with their ability to have a global view of the structural information, have been shown to effectively improve performance. However, most of these methods consume too much measurement time. This paper introduces an ordinal optimization based strategy combined with a back propagation neural network for autotuning of the configuration parameters. The strategy was first proposed in the automation community for complex manufacturing system optimization and is customized here for improving distributed system performance. The method is compared with the covariance matrix algorithm. Tests using a real distributed system with three-tier servers show that the strategy reduces the testing time by 40% on average at a reasonable performance cost.

Key words: distributed systems; performance evaluation; autotune configuration parameters; ordinal optimization; covariance matrix algorithm

Introduction

Performance improvements, such as high throughput and reducing end-to-end response times, are critical in distributed computing systems^[1] such as cluster computing^[2], grid computing^[3], cyber-infrastructure^[4], and cloud computing^[5]. These mainstream distributed computing technologies are implemented in various computationally intensive areas to provide online services, such as real time on-line e-commerce transactions, collaborative games, and large scale scientific

workflow analyses. One major difficulty in improving the performance of distributed computing systems is that the services are deployed on different servers, which requires system performance improvement at a global level.

Distributed systems usually have some common configuration files and many configurable parameters, such as the session time, maximum number of connected clients, and cache pool size. Their values are recorded in a plain text file. The session time specifies how long the conversational state with one client is maintained across multiple requests. Maximum number of clients specifies how many concurrent clients can be connected to the server at one time. The cache pool is used in database systems to store temporary variables, tables, and results.

The performance of distributed systems is partially dependent on these system configuration parameters, so there is an escalating interest in how to make best

Received: 2011-01-26; revised: 2011-06-22

* Supported by the National Natural Science Foundation of China (No. 60803017) and the National Key Basic Research and Development (973) Program of China (Nos. 2011CB302505 and 2011CB302805)

** To whom correspondence should be addressed.

E-mail: jcao@tsinghua.edu.cn; Tel: 86-10-62772260

use of these configuration parameters for system performance improvements for the following three reasons:

(1) Parameter tuning is easy since their values are recorded as plain text or as XML text, which is commonly used in mainstream servers such as those by Apache, JBoss, and MySQL.

(2) Parameter tuning requires little extra cost compared with hardware investments.

(3) Parameter values significantly impact system performance. Tests show that optimal parameters can increase system throughput by 24.5% and reduce average response time by 28.1% compared with the default settings.

Many companies tune these parameters based on their experience. However, this study provides a scientific approach applicable to different application scenarios.

Although configuration parameter autotuning can significantly improve system performance, there are many difficulties.

For example, if the session time is set too short, even though this reduces the potential for malicious attacks, clients have to more frequently connect to the server, which leads to connection overhead. Increasing the number of maximum clients and cache pool size increases system throughput, but cannot guarantee satisfactory response time since the server has to deal with more client data, conversational states, and temporary information. Thus the optimal configuration parameters to run applications on web servers and/or server farms to improve system performance are difficult to identify.

Dynamic and random events in these complex systems, whose performance is affected by a combination of many factors, are difficult to predict using precise mathematical models, so the effect of traditional heuristic optimization methods is limited. Many researchers have proposed use of the black box model^[6,7] to optimize performance, but the proper combination of configuration parameters to maximize system throughput as well as minimize average response time is still a challenge.

In addition, determination of the proper configuration parameters is quite time consuming in the black box search model. Searches generally take 25 minutes or more to find an optimal value^[6,7], which is not

acceptable in many real-time scenarios. Thus quickly finding a configuration value is a great challenge.

This work focuses on fast autotuning of configuration parameters in dynamic distributed computing systems to address these challenges.

1 Preliminaries

Several concepts used in this work are defined here.

Definition 1 Policy

A policy, denoted by θ , is one vector of configuration parameters. For example, $\theta = [200, 300, 50 \text{ s}, 200 \text{ MB}, 20 \text{ MB}]$ is one policy for the configuration parameters: [MaxClients, MaxConnections, SessionTime, KeyBufferSize, MaxPoolSize].

Definition 2 Performance

Performance, defined as a cost performance, is $\Omega = \text{throughput}/\text{response time}$. The throughput is the average number of requests processed per second. The response time is the average end-to-end time from initiating a request to receiving the result.

Definition 3 Ideal performance

The ideal performance, $J(\theta)$, is used but there are many uncertainties affecting performance autotuning. Thus, tests should be replicated N (generally a large number) times for each policy with the arithmetic mean used to get Ω for each policy. The exact relationship J is then $\Omega = J(\theta)$.

Definition 4 Measured performance

The measured performance, $\hat{J}(\theta)$, is a relatively simple but computationally fast model. $\hat{J}(\theta)$ is used to approximate the ideal performance. ω type and noise level analytical methods defined in Definitions 7 and 8 are used to bridge the gap between the ideal performance and the measured performance as much as possible. The simple model used here is a neural network.

This model significantly reduces the computational cost to get a good-enough policy, but not necessarily the best one. The measured performance can be seen as the ideal performance plus some noise:

$$\hat{J}(\theta) = J(\theta) + \text{noise} \quad (1)$$

Definition 5 Good-enough set, G

The best g policies of the ideal performance are defined as the good-enough set G . The value of g is preset by the users.

Definition 6 Selected set, S

The selection set is the best s policies of the

measured performance.

Instead of finding the optimal policy θ^* in the whole search space, the ordinal optimization (OO) method searches for a small set S , which contains k good-enough policies in G . The success probability of such a search was set as α (e.g. 98%). The good-enough policies are the top g ($g \geq k$) in the search space Θ . The numbers k and g are preset by the users. They follow the condition in Eq. (2) and illustrated in Fig. 1.

$$P[|G \cap S| \geq k] \geq \alpha \quad (2)$$

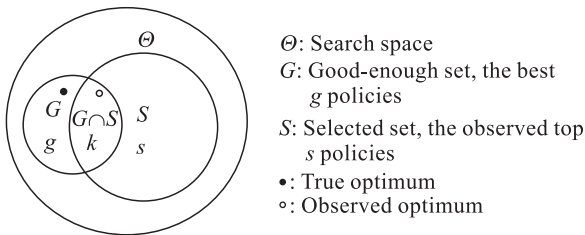


Fig. 1 Graphical illustration of the notation describing the key concepts used in ordinal optimization. G and S represent the good-enough and selected set with sizes g and s . S should have at least k overlapped policies in G . This narrows the search space from a very large $|\Theta|$ to a much smaller s .

The simple illustration in Fig. 2 shows how the OO method works. Suppose the problem is

$$\min_{\theta \in \Theta} J(\theta) = \theta, \quad \Theta = \{1, \dots, 9\} \quad (3)$$

The measured performance can be seen as Eq. (1) for the ideal performance plus noise. The noise is assumed to be a random variable uniformly distributed in $[0, 4]$.

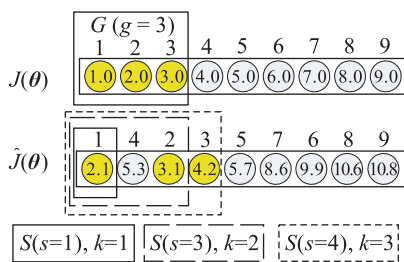


Fig. 2 Example to illustrate how ordinal optimization works. The search space consists of 9 policies in ascending order. The good-enough set G is shown by the left (best) 3 ($g=3$) shaded policies. The first, second, and third policies in the measured performance set ($s=1, 3$, or 4) are selected to get at least 1, 2, or 3 policies ($k=1, 2$, or 3) in G . Different sizes of s will give different k . The search is narrowed from 9 to 1/3/4 depending on the number of good-enough policies.

After generating selection set S , the ideal performance model is used to simulate the s policies to get the

best policy. In this way, the simulation set is reduced from the original $|\Theta|$ to s , which greatly reduces the computational time.

The size of s is determined by the size of the good-enough set, g , the overlap number k , and the success probability α . Two other concepts are also important.

Definition 7 ω type

The ω type is used to define the distribution of all the policies. There are generally five ω types as shown in Fig. 3. Suppose $\{\theta_1, \theta_2, \dots, \theta_n\}$ is the policy set Θ with their measured performance being $\{\hat{J}(\theta_1), \hat{J}(\theta_2), \dots, \hat{J}(\theta_n)\}$. Their performances are then listed in ascending order $\{\hat{J}(\theta_{[1]}), \hat{J}(\theta_{[2]}), \dots, \hat{J}(\theta_{[n]})\}$, with x as $[0, n-1]$, and $F(x)$ as

$$[\hat{J}(\theta_{[i]}) - \hat{J}(\theta_{[1]})] / [\hat{J}(\theta_{[n]}) - \hat{J}(\theta_{[1]})] \quad (4)$$

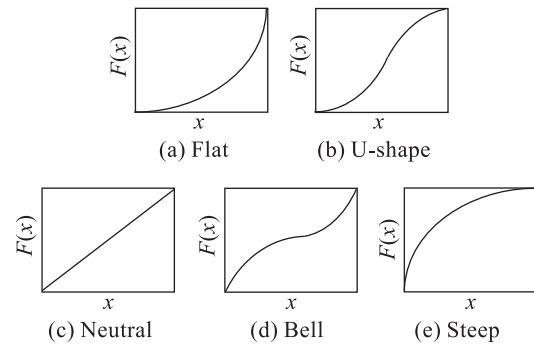


Fig. 3 Five ω types based on the measured performance of each policy. x represents the x -th smallest value of the measured performance while $F(x)$ is the normalized value from Eq. (4).

For a minimization problem, ω is like Fig. 3a, then many policies are good since they are very small. The selection set S could then be relatively small and still guarantee many good enough results. Similar rules can be applied to the other figures to get relationships between the problem type and the selection set S .

Definition 8 Noise level (NL)

The noise level is another factor that affects the size of selection set S . NL describes how “rough” the rough model $\hat{J}(\theta)$ is compared with the true performance model $J(\theta)$.

Again using Fig. 2 as an example, the standard deviation of the noise is 3. Thus, $NL = 3/(9-0) = 0.33$, which is a small noise level. The noise level can be calculated by normalizing the measured performance $\{\hat{J}(\theta_1), \hat{J}(\theta_2), \dots, \hat{J}(\theta_n)\}$ into $[-1, 1]$, and finding the maximum standard deviation, σ . These can be categorized into small noise levels ($0 < \sigma \leq 0.5$),

medium noise levels ($0.5 < \sigma \leq 1$), large noise levels ($1 < \sigma \leq 2.5$), and very large noise levels ($\sigma > 2.5$).

A small noise example is given in Fig. 2, three noise

levels are illustrated in Fig. 4 with the noise on a continuous interval $[-3, 3]$.

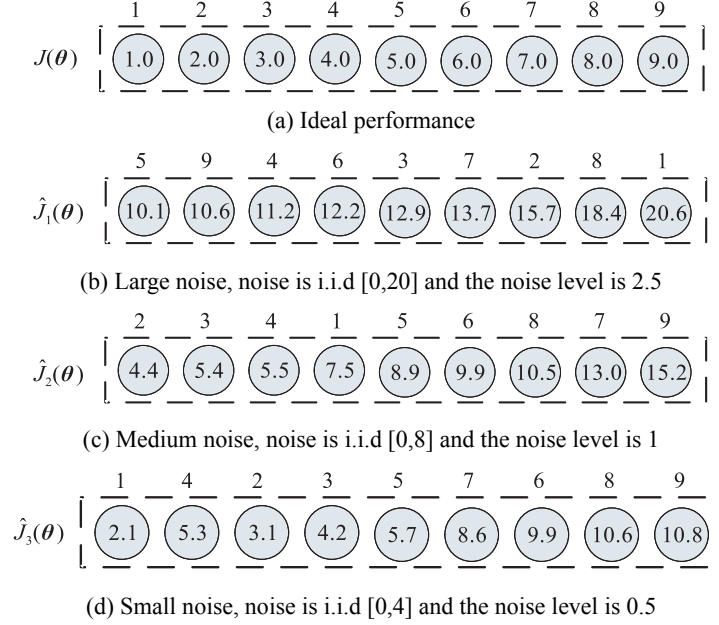


Fig. 4 Noise level can affect the measured performance. (a) ideal performance (values in the circles) and their order (figures above the dashed rectangular box), (b) measured performance with large noise with uniformly distributed random noise between $[0,20]$ and a noise level of $(20-0)/(9.0-1.0)=2.5$. The numbers above the values show their original order. In this example, the best policy (1.0 in (a)) becomes the worst (20.6 in (b)) after the large noise evaluation. Similar rules are applied in (b) and (c). The disorder in (d) is intuitively less than that in (b) and (c).

The characteristics of the autotuning configuration problem perfectly match the OO characteristics. Each configuration parameter θ would require a time consuming simulation to get the performance, which is not practical with a large search space. Finding the best θ^* in the good-enough set G , perhaps the top 2%, would be much easier. The OO method guarantees finding at least k policies in the good-enough set with a probability larger than α (e.g., 98%).

2 OO-Based Configuration Parameter Autotuning

2.1 When and why OO works

The configuration parameter must be specified for the distributed system, for example for Apache+MySQL servers or Tomcat+JBoss+Oracle servers. The parameters of the mainstream JSP/Servlet system are [MaxClients, MaxConnections, SessionTime, KeyBufferSize, MaxPoolSize], which typically have values between $[10, 1000]$, $[100, 500]$, $[30 \text{ s}, 200 \text{ s}]$, $[3 \text{ MB},$

$500 \text{ MB}]$, and $[10, 150]$. Different scenarios may have different intervals and configuration parameters. The objective is to determine which configuration parameters are used and their value space for a specific scenario. The OO method can then be used to identify the proper configuration parameters to make the application run smoothly and efficiently.

The method can effectively find good-enough configuration parameters because it digs out structural information to find the ranges for the good enough results, based on the ω type, noise level, and neural network models. This rough model fast records the structural information to reduce the number of options.

The OO method greatly reduces computational cost because the neural network model is only rough, with precision loss to reduce the overhead, to further reduce the number of searches.

Rather than taking much time to search for the optimal policy like some heuristic method does, the OO method searches for a proper balance between performance and simulation time.

2.2 Number of training samples

The OO method reduces the very large search space, $|\Theta|$ to a much smaller set s . Since some configuration parameters are continuous, such as SessionTime in [30 s, 200 s], a uniform and random sample in Θ is sufficient for the optimization

The problem is to set the training sample size, T , to enable the neural network model to represent the performance of the true model. Thus, the system needs to balance the neural network model effectiveness and the test time.

One extreme would do 1000 tests in a brute-force search of the entire space, Θ_N to get the best one, but this is quite time consuming. The other extreme would do only 1 real test, but the resulting rough model is too rough to be of use, even though the time cost is small. Traditional evolutionary strategies, such as SHC, CMA, are close to the brute-force end of the spectrum, while OO will be towards to the fast end. The tests will illustrate how to decide how many tests are needed.

The number of good-enough policies with the OO method is very difficult to determine. This depends on how much test time is available. Neural network model rules, such as avoiding underfitting and overfitting, can be used to determine the number of training samples. In our experiment, the number was initially 40 with better values found in tradeoffs. The tests show what happens to the system performance if this number is changed.

2.3 Configuration autotuning procedure

Given n parameters $\{p_1, p_2, \dots, p_n\}$ with value parameters p_i range of $[\underline{p}_i, \bar{p}_i]$, the OO method is then:

(1) For $i=1, \dots, n$, linearly quantize $[\underline{p}_i, \bar{p}_i]$ into $[0, 144]$;

$$\text{Map } p_i \in [\underline{p}_i, \bar{p}_i] \text{ to } \bar{p}_i + \frac{144}{\bar{p}_i - \underline{p}_i} p_i.$$

(2) Randomly choose 1000 groups uniform n -dimensional Gaussian distribution of configuration parameters, $\{\theta_1, \theta_2, \dots, \theta_{1000}\}$, in $[0, 144]^n$ and quantize p_i back to its original space $[\underline{p}_i, \bar{p}_i]$. Here θ_j represents the n -dimensional vector $\{v_{1j}, v_{2j}, \dots, v_{ij}, \dots, v_{nj}\}$, $j \in [1, 1000]$.

$$\text{Map } v_{ij} \in [0, 144] \text{ to } \frac{\bar{p}_i - \underline{p}_i}{144} v_{ij}.$$

(3) Use $\{\theta_1, \theta_2, \dots, \theta_{40}\}$ as training samples. Test the real performance for each $\theta_i (i \in [1, 40])$ and record the cost performance as $\{c(\theta_1), c(\theta_2), \dots, c(\theta_{40})\}$.

(4) Use the three-tier BP neural network model^[8] shown in Fig. 5 to evaluate the 1000 inputs, θ_i , with their related outputs $c(\theta_i)$. This rough model is computationally fast.

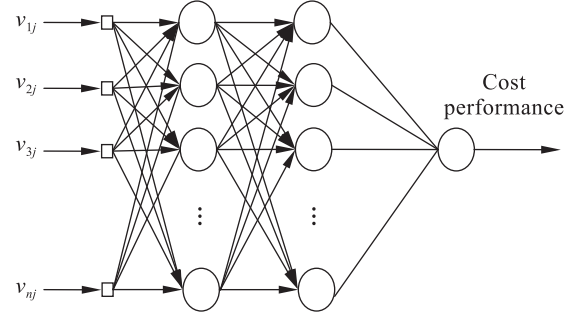


Fig. 5 Three-tier neural network (rough model) for autotuning configuration parameters with θ_j for the n configuration parameters and the output as the cost performance. 40 training samples from real experiment results are used to train the neural network.

The BP neural network roughly learns the relationship between each configuration parameter, θ_i , and its cost performance $c(\theta_i)$ by defining weight on each edge. The cost performances of the remaining samples (θ_{41} to θ_{1000}) are derived from the neural network output, so this is called rough model. These are denoted as $\{\hat{c}(\theta_i), i = 41, 42, \dots, 1000\}$.

(5) Use the model to roughly evaluate the ω type (Flat, U-shape, Neutral, Bell or Steep). Simulations suggest that the majority of applications have a Bell shape, which implies that the good-enough data distribution to search is neutral (neither too good nor too bad). The effect of the Bell shape on the final selection set S is described by Ho et al.^[9]

(6) Specify the noise level. Normalize the observed performance in $\{\hat{c}(\theta_i), i = 41, 42, \dots, 1000\}$ into $[-1, 1]$, and find the maximum standard deviation σ . The noise level is then found based on σ .

(7) Specify the size, g , of the good-enough set G , the alignment level size, k , and the alignment probability α , as shown in Fig. 1.

(8) Use a look-up table^[9] to calculate the size, s , of the selection set S . OO theory guarantees that S contains at least k good-enough configuration parameter vectors with a probability no less than α .

(9) Find the best s θ_i based on the rough model with the least cost performance based on the neural network model output.

(10) Use tests to evaluate the s θ_i to get the best vector for configuration parameters.

These steps greatly reduce the computational cost from \mathcal{O} to \mathcal{O}_N and then to s .

3 Performance Evaluations of Actual Systems

3.1 System testbed

Tests of an actual running system used a typical mainstream three-tier JSP/Servlet system using the MVC framework. This system is used to organize different workflow tasks in a company to cooperate with each

other in task related events. The first tier is a web (HTTP/HTML) server to handle requests from remote clients, using Apache servers. The application server, Tomcat, has a controller component to identify the request type and target the person to handle the request. Then the request is transferred to the related servlet component to invoke its model and find its execution source. During the invoking process, a MySQL database system is also used to manage the relational data tables, the contents file, and documents. The system structure and workflow are illustrated in Fig. 6.

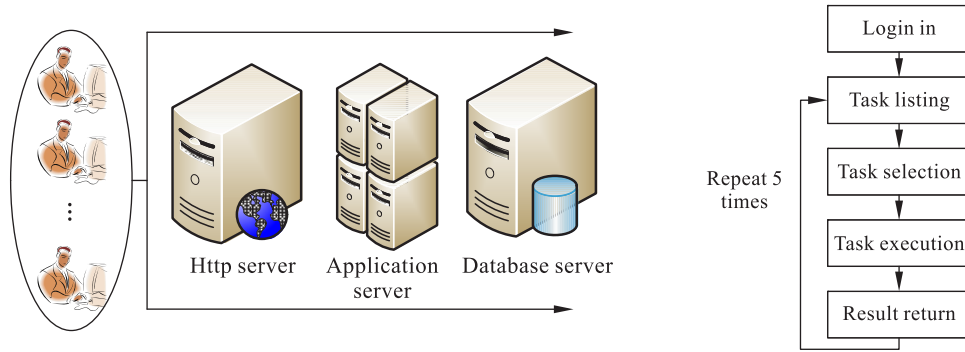


Fig. 6 Test system structure and workflow

The system has 5 action categories. Each user has to login only once, with all of the user's tasks then listed. After that, the tasks are selected and executed and the system returns the result. The user then returns to the task list and continues the procedure. The workflow is simplified by having each user repeat the process five times before logging out.

Seven parameters are used to represent the system performance, [MaxKeepAliveRequests, KeepAliveTimeout, ThreadCacheSize, MaxInactiveInterval, MaxConnections, KeyBufferSize, SortBufferSize]. All these parameters impact system performance when the system is saturated. The default values are $X^0 = [100, 5\text{ s}, 8, 2\text{ s}, 400, 20\text{ MB}, 256\text{ KB}]$. Parameter spaces are $[10, 200] \times [10\text{ s}, 200\text{ s}] \times [5, 100] \times [5\text{ s}, 50\text{ s}] \times [100, 500] \times [8\text{ MB}, 256\text{ MB}] \times [128\text{ KB}, 1024\text{ KB}]$. The parameter spaces were chosen for this specific application with some default settings not necessarily in these intervals in the software. Some of the settings depend on the application type, such as ThreadCacheSize and MaxInactiveInterval. Some of the settings also depend on the intrinsic characteristics of the server, such as the SortBufferSize which is dependent on MySQL. All these configuration parameters are inter-related so the optimal set is difficult to find.

3.2 Test results and analysis

A neural network was used with 40 groups of input/response pairs to generate the rough model. The method described in Section 2 was used to estimate the noise level, which is quite high (large noise level) due to many unpredictable factors in the dynamic system behavior. The Bell curve ω type was used to get a good-enough set g with 100 vectors (actually, with many more because of the large search space), with the required alignment level k equal to 5 and the alignment probability α as 98%. Those parameters were set before the tests to give a selected set S size of 20. The probability is $P[|G_\theta \cap S| \geq 5] \geq 98\%$.

The tests show the impact on three performance metrics, throughput (average number of requests processed per second), response time (average time to finish a servlet execution), and cost performance. 410 virtual users were used to saturate the three-tier system. To make the test more applicable to a real scenario, the workload was not added simultaneously. Instead, the workload was gradually increased and lasted for a specific period of time (60 s) after each user was initiated. The performance results are illustrated in Fig. 7 with the default parameters: default throughput $t(X^0)$

(requests/s): 149.21; default response time $r(X^0)$ (ms): 501.8; default cost performance $c(X^0)$: 0.2695.

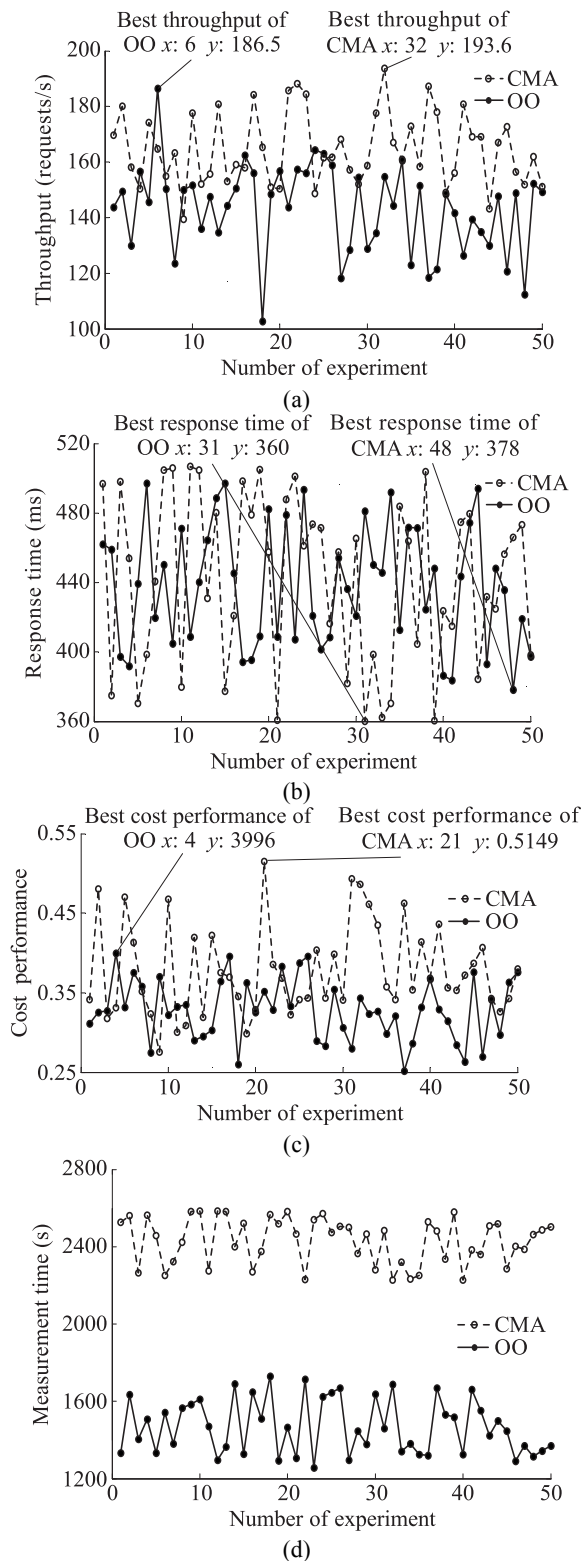


Fig. 7 Performance with the OO and CMA systems in (a) throughput, (b) response time, (c) cost performance, and (d) measurement time

The CMA method had the best cost performance in the 21st test. $X^{\text{CMA}} = [126, 28.19, 16, 29, 89, 136, 323.51]^T$, which resulted in a throughput, $t(X^{\text{CMA}}) = 185.73$ requests/s and cost performance $c(X^{\text{CMA}}) = 0.5149$, 24.5% and 91.1% better than for the default setting. The response time, $r(X^{\text{CMA}}) = 361$ ms was 28.1% less than for the default settings. The OO optimal result was in test 4 for $X^{\text{OO}} = [194, 46.28, 46, 35, 101, 188, 445.95]^T$. The throughput was $t(X^{\text{OO}}) = 156.6$ (requests/s) and the cost performance was $c(X^{\text{OO}}) = 0.3996$, 5.0% and 48.3% increases compared with the default settings. The response time was reduced by 21.9% to $r(X^{\text{OO}}) = 392$ ms.

A comparison of the CMA and OO methods shows the CMA throughput is 15.7% better and the response time is 8.0% better. The OO method tends to explore more of the estimated best value space and with the BP neural network model used here, the performance is strongly determined by the accuracy of the rough neural network model. On the contrary, the CMA method tries both exploration and exploitation, so it more easily picks out better configuration settings in the search space. The tradeoff then is that CMA gives a slightly more optimal result at a cost of a longer simulation time, while the OO method gets a reasonably optimal results in a dramatically shorter time. The test times are illustrated in Fig. 7d.

The measurement time comparison used the (4,9)-CMA algorithm, which on average required five generations to arrive at the optimal result, which is also reasonable for evaluating the simulation time. The time consumed of the OO method is much less than for CMA. The OO testing time includes establishing the neural network rough model, defining the ω type and noise level, and looking up the table to get the selection set S . Actually these parts can be neglected since they are done only once before the tests. Thus, the time depends only on the selected set, S .

The average measurement time with CMA was $t(X^{\text{CMA}}) = 2432.25$ s compared to only $t(X^{\text{OO}}) = 1470.43$ s with OO. The average OO time was 60% less than CMA. Thus, the OO method gives a significantly reduced measurement time for autotuning the configuration parameters.

3.3 Test results with different number of training samples

The number of training samples for the neural network

rough model is quite important. If this number is too large, the system will find more information about good policies at the cost of more simulation time.

The BP neural network in the previous tests used 40 training samples. This number will affect the OO performance. Thus, different numbers of training samples were used to show their impacts on the throughput, response time, cost performance and, most importantly, the test time. The results are also compared with the CMA method.

110 tests were made using the OO method as in Section 2. Additional tests did not lead to noticeable performance improvements. More tests would require more measurement time, which reduces the advantage of the OO method. All the tests used the same

configuration parameter set and value set $\{\theta_1, \theta_2, \dots, \theta_{1000}\}$. The test time was reduced by uniformly sampling a number of training samples (10, 20, etc).

The results of Fig. 8 show the use of more samples to train the BP neural network finds a better configuration parameter set because the method goes deeper into the structural information, which compensates for the drawback of the OO method. However, the measurement time increase is almost linear with the number of samples, so the number of training samples depends on how much measurement time is available. For the same number of training samples, 110 in these tests, the throughput, response time, and cost performance did not improve.

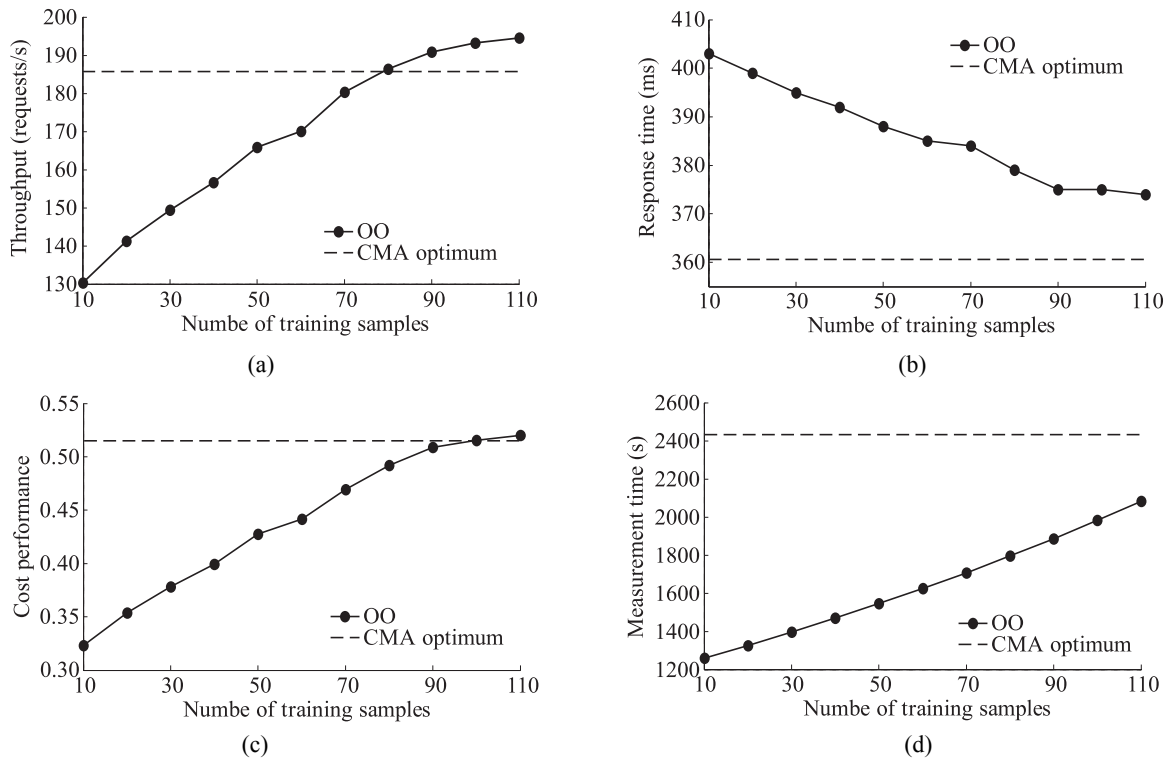


Fig. 8 Impact of neural network training samples on (a) throughput, (b) response time, (c) cost performance, and (d) measurement time. With increase numbers of training samples, the OO performance is better at the cost of more measurement time. The CMA optimum is the best result.

4 Conclusions and Future Work

This paper describes an ordinal optimization based strategy to improve the performance with less measurement time for autotuning configuration parameters. This method significantly reduces the measurements in real distributed systems with a slight optimization performance decrease compared with the traditional

evolutionary optimization strategy. Future work can be categorized into four directions.

The “no free lunch theorem” suggests that the algorithm needs to carefully search the intrinsic characteristics of the autotuning configuration problem to dig out more information to support the decision-making. The iterative rough model described is an important extension of the OO algorithm.

Web servers have too many parameters to be

thoroughly searched. Some of these are application dependent; thus, limiting traditional optimization methods. The configuration parameters can be selected according to the application type for specific scenarios.

There have been many extensions of the OO method introduced in related works. Some, such as Optimal computing budget allocation (OCBA)^[10-13] and breadth & depth (B&D)^[14], can be used together to further reduce the test time. The vector ordinal optimization (VOO)^[15] is also effective in dealing with multiple objective problems and other performance metrics.

The number of training samples used here was specified based on experience and some rules of thumb for neural networks, such as to avoid over fitting or under fitting. These could be further improved for the characteristics of distributed computing systems to further reduce the measurement time.

Acknowledgement

Junwei Cao thanks Professor Erik Katsavounidis of the MIT LIGO Laboratory for his collaboration on the LIGO gravitational-wave research. Fan Zhang was supported by 2010-2011 and 2011-2012 IBM Ph.D. Fellowships.

References

- [1] Hwang K, Xu Z. Scalable Parallel Computing. USA: McGraw-Hill, 1998.
- [2] Jia Q, Zhao Q. A SVM-based method for engine maintenance strategy optimization. In: Proceedings of the 2006 IEEE International Conference on Robotics and Automation. Orlando, FL, USA, 2006: 1066-1071.
- [3] Foster I, Kesselman C. The Grid: Blueprint for a New Computing Infrastructure. USA: Morgan-Kaufmann, 1998.
- [4] Atkins D, Droegemeier K, Feldman S, et al. Revolutionizing science and engineering through cyberinfrastructure. Technical report. National Science Foundation Blue - Ribbon Advisory Panel on Cyberinfrastructure, 2003.
- [5] Boss G, Malladi P, Quan D. IBM high performance on demand solutions. Technical report. IBM. Oct. 2007.
- [6] Xi B, Liu Z, Raghavachari M. A smart hill-climbing algorithm for application server configuration. In: Proceedings of the 3rd International Conference on World Wide Web. New York, NY, USA, 2004: 287-296.
- [7] Saboori A, Jiang G, Chen H. Autotuning configurations in distributed systems for performance improvements using evolutionary strategies. In: Proceedings of the 28th International Conference on Distributed Computing Systems. Beijing, China. 2008: 765-772.
- [8] Haykin S. Neural Networks: A Comprehensive Foundation (2nd Edition). USA: Prentice Hall, 1998.
- [9] Ho Y C, Zhao Q C, Jia Q S. Ordinal Optimization, Soft Optimization for Hard Problems. Germany: Springer, 2007.
- [10] Teng S Y, Lee L, Chew P. Integration of indifference-zone with multi-objective computing budget allocation. *European Journal of Operational Research*, 2010, **203**(2): 419-429.
- [11] He D H, Lee L, Chen C H, et al. Simulation optimization using the cross-entropy method with optimal computing budget allocation. *ACM Transactions on Modeling and Computer Simulation*, 2010, **20**(1): 1-22.
- [12] Hsieh B, Chen C H, Chang S H. Efficient simulation-based composition of dispatching policies by integrating ordinal optimization with design of experiment. *IEEE Transactions on Automation Science and Engineering*, 2007, **4**(4): 553-568.
- [13] Chen J, Lee L. A multi-objective selection procedure of determining a Pareto set. *Computers and Operations Research*, 2009, **36**(6): 1872-1879.
- [14] Chen C H, Lin Ji W, Yücesan E, et al. Simulation budget allocation for further enhancing the efficiency of ordinal optimization. *Journal of Discrete Event Dynamic Systems: Theory and Applications*, 2000, **10**: 251-270.
- [15] Li D, Lee H, Ho Y C. Vector ordinal optimization – A new heuristic approach and its application to computer networks routing design problems. *International Journal of Operations and Quantitative Management*, 1999, **5**: 211-230.