Contents lists available at SciVerse ScienceDirect



Sustainable Computing: Informatics and Systems

journal homepage: www.elsevier.com/locate/suscom



Faster exploration of data centre cooling using thermal influence indices

Harshad Bhagwat^a, Amarendra Singh^{a,*}, Arunchandar Vasan^b, Anand Sivasubramaniam^{c,1}

^a TCS Innovation Labs-TRDDC, Tata Consultancy Services Limited, Pune, India

^b TCS Innovation Labs-Chennai, Tata Consultancy Services Limited, Chennai, India

^c Department of Computer Science and Engineering, The Pennsylvania State University, PA, USA

ARTICLE INFO

ABSTRACT

Article history: Received 11 September 2012 Accepted 29 January 2013

Keywords: Data centre Efficient cooling Computational fluid dynamics (CFD) Cooling is an important issue in data centre design and operation. Accurate evaluation of a design or operational parameter choice for cooling is difficult as it requires several runs of computationally intensive computational fluid dynamics (CFD) based models. Therefore there is need for an exploration method that does not incur enormous computation. In addition, the exploration should also provide insights that enable informed decision making. Given these twin goals of reduced computation and improved insights, we present a novel approach to data centre cooling exploration. The key idea is to do a local search around the current design or operation of a data centre to obtain better design or operation parameters subject to the desired constraints. To do this, all the microscopic information about airflow and temperature in data centre available from a single run of CFD computation is converted into macroscopic metrics called influence indices. The influence indices, which characterize the causal relationship between heat sources and sinks, are used to refine the design or operation of the data centre either manually or programmatically. New designs are evaluated with further CFD runs to compute new influence indices and the process is repeated to yield improved design or operation as per the computation budget available. We have carried out guided explorations of design and then operation of a realistic data centre using this methodology. Specifically, we considered maximization of the heat load (design parameter) and supply temperatures of CRAC (operating parameter) in the data centre subject to the constraints that: (1) servers are kept at appropriate temperatures and (2) overloading of CRACs is avoided. Our evaluation for a 1500sq.ft. data centre shows that the use of influence indices cuts down the exploration time for design by 80% as compared to unguided explorations. It is demonstrated quantitatively that this solution is close to optimal. A guideline was evolved to reduce the effect of initial configuration on the final solution. Proposed methodology provides an additional 10% reduction in operational cost over existing methods of unguided explorations.

© 2013 Elsevier Inc. All rights reserved.

1. Introduction

With the rapid growth of computing requirements, data centres have become ubiquitous. Data centres are estimated to consume up to 2% of the entire energy consumption of the world and this number is only poised to increase [2]. This has led to research in the Greening of IT. Several works have focused on computing issues like managing sever workload, consolidation, virtualization, energy efficiency in hardware design, etc. Because IT equipment need to operate at an appropriate temperature to ensure that reliability requirements are met, cooling a data centre adequately is essential

* Corresponding author.

E-mail addresses: harshad.bhagwat@tcs.com (H. Bhagwat),

amarendra.singh@tcs.com (A. Singh), arun.vasan@tcs.com (A. Vasan),

anand.sivasubramaniam@tcs.com (A. Sivasubramaniam).

[3]. Thus, data centre power consumption has several non-IT components in addition to the IT power and cooling the IT equipment takes significant power. For instance, as a rule of thumb, for every watt consumed by IT components, one may need an additional watt of non-IT power, of which cooling constitutes a significant portion. Indeed, cooling accounts for nearly 80% of non-IT power in a typical data centre [4]. Therefore, efficient cooling of a data centre is an important issue.

There are several challenges in efficiently cooling a data centre both at the design and operation stages. At the design stage, data centre architect or design engineer needs to explore various parameters such as placements of racks, perforated tiles, computer room air-conditioner (CRAC) units, room dimensions, and plenum height. This leads to a combinatorial explosion of the search space. At the operations stage, data centre manager has to dynamically control the set-points of CRAC to minimize the power consumption. Also, managers deploy software solutions such as virtualization for consolidation and dynamically move around workloads without

¹ The work was done during the author's sabbatical at Tata Consultancy Services.

^{2210-5379/\$ -} see front matter © 2013 Elsevier Inc. All rights reserved. http://dx.doi.org/10.1016/j.suscom.2013.01.004

understanding the cooling implications of such transitions. This could potentially create hotspots in a data centre if cooling is insufficient, or lose in (excessive) cooling whatever energy is saved in computing.

The key issue in these cases is that a reasonably accurate evaluation of a design or operational parameter choice for cooling is computationally hard. This is because one needs to typically run a CFD model, which provides a three-dimensional temperature and flow profile inside a data centre, to quantify the impact of a choice before making the choice. CFD models are based on the numerical solutions of non-linear, coupled partial differential equations and are computationally intensive. For instance, one instance of CFD model of a typical 1500 sq.ft. data centre takes 21 h time on a machine with an Intel Xeon 3 GHz processor and 4 GB RAM. Though the time of computation can be brought down by use of parallel computing, this calls for expensive resources. Therefore there is a need to explore options in cooling both at design and operations stages of a data centre in a manner that does not incur enormous computational time or resources. Another reason for avoiding repeated CFD computations is limited financial budgets available for such evaluations. In addition to CFD modelling and evaluation of a data center's cooling, the exploration should also provide valuable insights to the designers and managers into the current operation of the data centre, so as to enable him or her to take decisions which are better than the norm

Given the twin goals of reduced computation and improved insights, we present a novel approach to data centre cooling exploration. Intuitively, the idea is to do a local search around the current design or operation of a data centre to obtain better design or operating parameters subject to the desired constraints. To do this, we distill all the microscopic information about airflow and temperature in data centre available from a single run of CFD computation into macroscopic metrics called influence indices. The influence indices characterize the causal relationship between heat sources and sinks (e.g., Rack₈ receives 50% of its cold air from CRAC₃) and provide deep insight into the operation of the data centre (e.g., setting the supply temperature of this CRAC₄ will affect Rack₃ and Rack₄ the most).

Because the indices make use of microscopic information, they can be composed into causal relationships at various granularities (e.g., servers, rack, cluster of racks, all racks in a row, etc.) thereby providing the data centre designer and manager various levers with a quantification of their impacts. The influence indices can be used to refine the design or operation of the data centre either manually or programmatically. The new design or operating point can be evaluated with a further CFD run to compute new influence indices and the process can be repeated to yield improved design or operation as per the computation budget available. Note that the desired levels of influence indices themselves can be used as criteria for stopping the iterations (explained in Section 3 in detail). We review related work in Section 2. To the best of our knowledge, this is the first work that examines how to explore design or operation choices with limited computation in the context of data centres.

Specific contributions of our work include:

- We propose simple metrics that quantify causal heat relationships between sources and sinks, thereby providing deep insights into the practical operation of a data centre which an IT manager can easily understand and make use of in real time.
- We propose a systematic approach to explore the design space and operation space for cooling a data centre given its cooling budget (i.e., in terms of CRAC capacities and locations) through a process of iterative refinement, which can be either guided or automated.

- We prove the effectiveness of the methodology by quantifying how close the results are to the optimal. We also evaluate the sensitivity of the procedure to the initial starting configuration and present effective guidelines to reduce this sensitivity.
- We reduce significantly the number of CFD computation runs required to explore design or operation choices, thereby enabling quick and reasonably accurate thermal decision making.
- We show significant reduction in operating costs after exploration of operation space carried out using influence indices compared to existing methods.

We introduced the guided exploration of data centre cooling using influence indices in Ref. [1]. We evaluate our approach with exploration of design space and operation space of a realistic data centre. In exploration of design space, we choose to explore the effect of server placements in racks, which is a design parameter and dictates heat load distribution in a data centre. We wish to maximize total heat load in the data centre that can be cooled, subject to following constraints: (1) servers are kept below temperatures to meet the reliability requirements of servers and (2) overloading of CRACs is avoided. Our evaluation shows that the use of influence indices cuts down the exploration time by 80%. We show that the methodology is successful in locating one of the best designs in spite of reduced number of CFD runs and local search employed. As the final solution is dependent on initial state, we analyze the sensitivity of the solution to the initial starting point, and identify guidelines to ensure that the search ends at one of the best designs.

We then explore operation space for reducing operational costs. We aim at maximizing supply temperatures of CRAC subject to the same constraints stated above. We found that use of influence indices to guide the exploration results in more efficient operation: reducing operating costs further by 10% on top of the efficient operation resulted from the existing approach.

The rest of the paper is organized as follows. Section 2 surveys related work. Section 3 explains the preliminaries regarding influence indices and CFD model of data centre. Section 3 also presents the problem statement and explains the solution methodology employed to solve the above problem with results presented in Section 4. Finally, Section 5 presents our conclusions.

2. Related work

Related work can be broadly categorized into four categories: (1) optimization of cooling (both design and operation) in data centre; (2) reduction of computation for evaluation of cooling; (3) performance metrics and (4) importance of heat load distribution.

• Optimization of cooling: there have been numerous efforts to optimize the cooling in a data centre. The scale considered for the optimization is ranging from inside computing equipment [5,6] to rack housing these equipment [7,8] to data centre as whole [9]. Efforts for optimization considering data centre as whole can be broadly categorized further into design stage approaches and operational approaches. Design stage approaches include tile placement, CRACs placement, server placement, rack placement, room dimensions, and plenum heights [10-12]; while operational stage approaches include as dynamic controlling of CRACs, temperature aware workload scheduling [13–17]. For example, Bhopte et al. [18] have used full scale CFD model to explore design changes which invoke large number of instances of CFD computations for couple of design parameters. Full design space exploration through CFD alone would be computationally prohibitive. Our work is complementary to any such optimization strategy in that it can help guide any optimization strategy using Download English Version:

https://daneshyari.com/en/article/493828

Download Persian Version:

https://daneshyari.com/article/493828

Daneshyari.com