

Don't step on the Distribution's Tail!

Investigating the impact of random fluctuations on efficient resource utilization

Fabrice Saffre
BT Research and Innovation
Ipswich, United Kingdom
fabrice.saffre@bt.com

Hanno Hildmann
NEC Laboratories Europe
Heidelberg, Germany
hanno.hildmann@neclab.eu

Abstract—It is sometimes assumed that the total amount of a resource being consumed is the key consideration when attempting to devise the most efficient management strategy. We explore a case in which the rate and timing of resource utilization are also susceptible to impact on performance and illustrate how bringing random fluctuations under control can help maximize efficiency in a simple client-server scenario. The role of self-organization and distributed control methods in achieving this goal is briefly discussed.

Index Terms—random fluctuations; resource management; performance optimization; client-server architecture

I. INTRODUCTION

MUCH has been said about demand side management [1], the so called smart grid [2] and about intelligent infrastructure [3]. Intuitive results are often generalised and the important little exceptions are sometimes overlooked. This short paper aims to invite the reader to follow the authors on a interesting, empirical result based investigation.

We simulate a server farm under a capping constraint, as it is commonly discussed in the literature when e.g. considering the use of renewable energies to power some infrastructure. By comparing two scenarios and assuming a cost for the deviations from the norm for both, we can compare the impact either strategy has on the cost efficiency of the operation. This led us to a result which we consider worth sharing.

II. BACKGROUND

Efficient resource management is fast becoming one of the most critical features of almost any human activity. Whether this is the case seems clear to us, why this is the case we would rather leave to a philosophical debate. On a very abstract level, this is arguably because we have been so successful as a species that we have reached a limit above which the planet is no longer capable to sustain our wasteful habits. There was a time when we could afford to be oblivious of how much water, food, energy etc. we were using, simply because our impact was negligible, the environment was able to replenish the stocks of whatever resource we needed as quickly as it was consumed. But is no longer the case, and has not been the case for some time.

A. The problem

For better or worse, we have entered an age in which we will have to act more responsibly, to use only what is needed when it is needed, under penalty of seeing our global society

descend into conflicts and service disruption as we are forced to compete ever more aggressively for dwindling supplies. The *we* here can be understood as the societal *all of us*, but also as the managerial *our company* or even the personal and cost aware *I* in every day live.

The most visible aspect of this requirement is that we must seek to limit the total amount of a resource that is consumed in the process of achieving a certain goal. There are countless examples of this today, from designing more fuel-efficient vehicles to engineering crops that require less water or nutrients. There is however another, perhaps less obvious angle to this quest for efficiency: notwithstanding how much of a resource is consumed in total, when and at which rate it is being used can also have a critical impact on its availability (and cost). While this seems an obvious fact, it is often overlooked, or lost in abstraction.

B. Example

A good example of such a resource is renewable power. Imagine that a solar plant producing 1 MW for 10 hours a day is used to power a community that consumes a total of 10 MWh over 24 hours. If some power is used at night or if the load ever exceeds the available output from the solar plant during the day, even if the aggregate supply as well as consumption total for the day is still exactly identical (i.e. 10 MWh), some provision for this “mismatch” between supply and demand will need to be made, in the form of an additional power source or storage facility. On the contrary, if the demand could be exactly mapped to the output of the solar plant over time, then it could be met locally and efficiently (with fewer losses). This is of course widely discussed in the recent literature, and techniques to shape the load so as to better approximate the available supply are collectively known as “Demand-Side Management” or DSM [4].

There is a slightly different but related and very common special case of this problem: what if the average demand for a resource is constant and known but instantaneous consumption fluctuates randomly and unpredictably? If a facility (or supply) is dimensioned so as to accommodate the average requirement, then it will necessarily be sometimes under-, sometimes over-used. Moreover, depending on the characteristics (amplitude, symmetry, . . .) of the fluctuations, the time spent in either state (over-utilization and under-utilization), as well as the deviation from the average, may vary greatly.

C. The issue with statistics

In practice, it is often assumed that, because fluctuations statistically cancel each other out over time, this is no cause for concern and so not a relevant field of study. But this view is contradicted by the realization that the rate and timing of resource utilization can impact on efficiency (as illustrated by the above example). Similarly, the commonly held view that in a large enough population of consumers, deviations from the global average will be negligible is also a simplistic one. Indeed, there is no guarantee that the costs incurred by such deviations grow linearly. If they don't, then even small deviations can have a significant impact and large ones, although extremely rare, could have a disastrous effect.

D. Aim of this paper

In this paper, we experiment with a set of conceptual tools designed to quantitatively measure the influence of random fluctuations on performance in a simple client-server scenario. Specifically, we compare the case in which no upper bound is imposed on the number of active servers to that in which there is such a constraint. In effect, we propose a cost/benefit analysis of two strategies: cutting the upper-end tail of the distribution (which creates execution delays) versus "stepping on it" (which may incur extra costs).

III. MODEL AND SIMULATIONS

A. The model

We used Monte Carlo simulations to approximate the dynamics of a group of servers. On every time-step, all identical servers have a fixed probability P to receive a new job, the duration of which is comprised between 1 and $(2 \times \text{avg} - 1)$ time-steps (where avg is the average duration of a job). P is chosen in such a way that, statistically, slightly less than 1 out of 8 servers is expected to be busy at any time.

On every time-step, the number of servers in the "busy" state is recorded. In the default scenario (no constraints), this value is allowed to exceed the 1:8 ratio. In the other scenario, no more than 1 out of 8 servers are allowed to be active simultaneously. If random fluctuations in job arrivals / duration would cause the system to exceed this limit, the execution of a corresponding number of jobs (randomly selected) is temporarily suspended (i.e. the servers processing them are put on stand-by).

In the default scenario, performance degradation is measured by the number of server-time-steps falling above the 1:8 target and incurring extra costs. In the constrained scenario, it is the cumulative delay (total number of time-steps) suffered as a result of imposing a cap on the number of active servers (QoS (quality of service) penalty).

B. Simulations

Data-centers were modelled to operate between 1000 and 4000 servers (by increments of 500) and were simulated for 512 time-steps. The average job duration was 8 time-steps (i.e. flat distribution between 1 and 15). There were 1000

realizations for every size and for each scenario, totaling 14000 independent simulation runs.

A control simulation (simulating batches up to 1024 time steps) was run to confirm the results presented in the next section. The results from this investigation confirmed that the system was at, or very close to, its steady state at $t=512$.

IV. RESULTS

A. Frequency distribution

Fig. 1 shows the frequency distribution of simulation outcome for 4 different server population sizes. As expected, for the "constrained" scenario, there is a sharp peak corresponding to the 1:8 ratio (125, 250, 375 and 500 servers respectively) since the whole tail of the distribution has "collapsed" onto this single value. Note that the remainder of the distribution closely follows the Gaussian profile found for the default scenario, apart from in the case of smaller populations, where capping also seems to negatively affect the height of the normal peak.

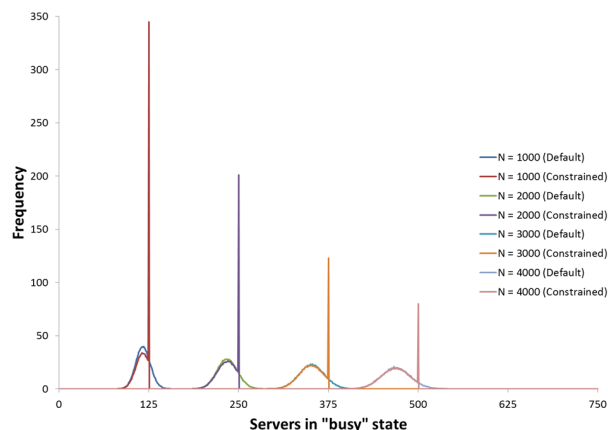


Fig. 1. Frequency distribution of system state as a function of the number of active servers, for 4 different population sizes and for the two possible scenarios ("Default" vs. "Constrained"). 1000 runs per set of parameter values.

The values reported in Fig. 1 consider only the values for the last 128 (out of 512) time-steps, when the system is at or very close to steady state. The motivation is that the initial values would be too much affected by the build up, but after $3/4^{th}$ of the steps we can assume that the overspill from previous steps has normalised (especially since the individual duration can not exceed 15 and we are not including the first 384 steps).

B. Aggregated QoS violation

Fig. 2 shows the evolution of the average delay (QoS penalty) incurred over all processed jobs, for increasing system size and for both scenarios. As expected, this value is unaffected by system size in the default case, where a delay only occurs when a job is submitted to a server before it has completed the execution of its predecessor. By contrast, in the constrained scenario, a QoS penalty is also incurred when the overall workload exceeds the processing capacity of $1/8^{th}$ of the population.

The main finding is that the average delay converges quickly for both cases as system size increases. This is due to the amount of time spent above the cut-off limit in the default scenario being inversely proportional to the number of servers involved; thus capping becomes proportionally less frequent in the constrained scenario (illustrated by the decreasing height of the sharp peak corresponding to 1:8 limit, see Fig. 1).

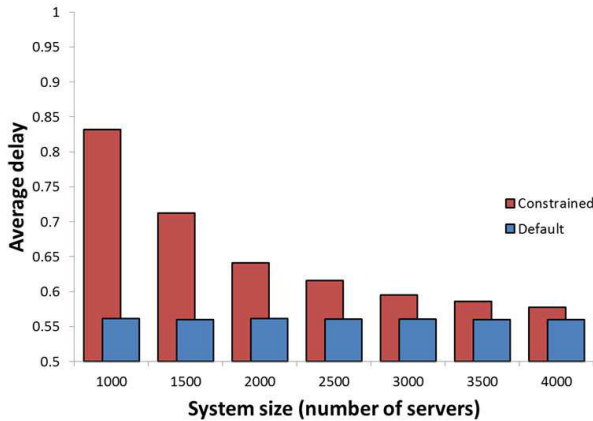


Fig. 2. Evolution of the average delay per job as a function of population size (workload proportional to the number of servers). The larger the system, the lower the QoS penalty resulting from enforcing the 1:8 cap.

C. Comparative QoS violations

Although to some extent these results are merely intuitive and can be anticipated from known statistical properties [5], they can also be interpreted in a different and more surprising way. For instance, the fact that, in the default scenario, the absolute number of server-time-steps falling inside the high-end tail of the frequency distribution shown on Fig. 1 (i.e. above the 1:8 limit) exhibits a maximum (see Fig. 3) clearly suggests that the capping strategy would yield most benefits within a finite range of system sizes.

D. Potential benefit analysis

Table I illustrates how the potential benefits of using the capping strategy would vary as a function of system size, under the arbitrary assumption that every server-time-step over the 1:8 target incurs an extra £0.01 cost (e.g. because of higher server rental cost) and every time-step delay incurs a £0.01 penalty fee (e.g. as compensation for breaching the service-level agreement).

TABLE I
POTENTIAL SAVINGS FROM APPLYING A CAPPING STRATEGY (ARBITRARY PENALTY COSTS). NOTE THE PRESENCE OF A MAXIMUM FOR N = 3000.

N	Cost Breakdown		
	Cost Differential (default-constraint)	QoS Penalty (constraint-default)	Net Savings
1000	£33,792	£21,251	£12,541
2000	£41,238	£12,687	£28,552
3000	£43,052	£8,160	£34,892
4000	£39,187	£5,562	£33,625

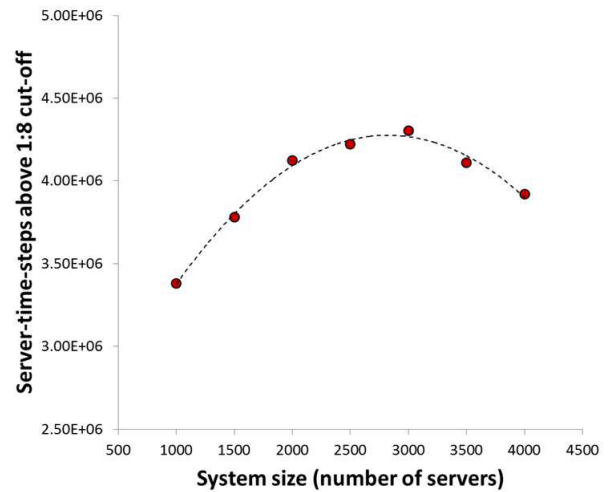


Fig. 3. Total (absolute) number of server-time-steps falling above the 1:8 cut-off limit in the default scenario, as a function of population size.

V. CONCLUSIONS

These preliminary findings confirm that there are realistic applications for which limiting the total amount of a resource being consumed over a period of time is not enough. Preventing the instantaneous load from exceeding a predetermined value can be beneficial and contribute to improving efficiency.

Moreover, relying on statistical effects, either over time or over a large population, to ensure a smooth demand profile may be a high-risk strategy as it is possible that even small or rare deviations from the average incur a severe penalty.

In this context, it seems important to investigate control methods that would permit to regulate the aggregated demand from a population of resource consumers so as to prevent them from exceeding a certain target. Because a central mechanism for achieving such a goal may not be feasible (due to scalability problems, limits on information availability or even ownership boundaries), we argue that a form of collective intelligence capable of supporting self-management is required. The investigation of a candidate technology for implementing such a distributed resource controller will be the subject of future work.

REFERENCES

- [1] A. A. Garcia, "Demand side management integration issues a case history," *Power Systems, IEEE Transactions on*, vol. 2, no. 3, pp. 772–778, Aug. 1987.
- [2] T. T. Kim and H. V. Poor, "Scheduling power consumption with price uncertainty," *Smart Grid, IEEE Transactions on*, vol. PP, no. 99, p. 1, 2011.
- [3] F. Saffre, H. Hildmann, and S. Nicolas., "The adaptive grid." The Institute of Telecommunications Professionals (ITP) - ITP Journal, vol. 6, no. 3, pp. 31–38, 2012.
- [4] G. Strbac, "Demand-side management: benefits and challenges," *Energy Policy*, vol. 36, pp. 4419–4426, 2008.
- [5] D. Bini, G. Latouche, and B. Meini, *Numerical Methods for Structured Markov Chains*, ser. Numerical Mathematics and Scientific Computation, 2005.