

This is an electronic reprint of the original article. This reprint may differ from the original in pagination and typographic detail.

Cost-Efficient, Utility-Based Caching of Expensive Computations in the Cloud

Byholm, Benjamin; Jokhio, Fareed; Ashraf, Adnan; Lafond, Sebastien; Lilius, Johan; Porres Paltor, Ivan

Published in:

23rd Euromicro International Conference on Parallel, Distributed, and Network-Based Processing

DOI:

[10.1109/PDP.2015.49](https://doi.org/10.1109/PDP.2015.49)

Publicerad: 01/01/2015

[Link to publication](#)

Please cite the original version:

Byholm, B., Jokhio, F., Ashraf, A., Lafond, S., Lilius, J., & Porres Paltor, I. (2015). Cost-Efficient, Utility-Based Caching of Expensive Computations in the Cloud. In M. Daneshtalab, M. Aldinucci, V. Leppänen, J. Lilius, & M. Brorsson (Eds.), *23rd Euromicro International Conference on Parallel, Distributed, and Network-Based Processing* (pp. 505–513). IEEE Computer Society Conference Publishing Services (CPS).
<https://doi.org/10.1109/PDP.2015.49>

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

Copyright Notice

The document is provided by the contributing author(s) as a means to ensure timely dissemination of scholarly and technical work on a non-commercial basis. This is the author's version of the work. The final version can be found on the publisher's webpage.

This document is made available only for personal use and must abide to copyrights of the publisher. Permission to make digital or hard copies of part or all of these works for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage. This works may not be reposted without the explicit permission of the copyright holder.

Permission to reprint/republish this material for advertising or promotional purposes or for creating new collective works for resale or redistribution to servers or lists, or to reuse any copyrighted component of this work in other works must be obtained from the corresponding copyright holders. It is understood that all persons copying this information will adhere to the terms and constraints invoked by each copyright holder.

IEEE papers: © IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works. The final publication is available at <http://ieeexplore.ieee.org>

ACM papers: © ACM. This is the author's version of the work. It is posted here by permission of ACM for your personal use. Not for redistribution. The final publication is available at <http://dl.acm.org/>

Springer papers: © Springer. Pre-prints are provided only for personal use. The final publication is available at <link.springer.com>

Cost-Efficient, Utility-Based Caching of Expensive Computations in the Cloud

Benjamin Byholm^{*†}, Fareed Jokhio^{*§}, Adnan Ashraf^{*†‡}, Sébastien Lafond^{*†}, Johan Lilius^{*†} and Ivan Porres^{*†}
Email: bbyholm@abo.fi, fjokhio@abo.fi, aashraf@abo.fi, slafond@abo.fi, jolilius@abo.fi, iporres@abo.fi

^{*}*Åbo Akademi University, Department of Information Technologies, Turku, Finland*

[†]*Turku Centre for Computer Science (TUCS), Turku, Finland*

[‡]*Department of Software Engineering, International Islamic University, Islamabad, Pakistan*

[§]*Quaid-e-Awam University of Engineering, Science & Technology, Nawabshah, Pakistan*

Abstract—We present a model and system for deciding on computing versus storage trade-offs in the Cloud using von Neumann-Morgenstern lotteries. We use the decision model in a video-on-demand system providing cost-efficient transcoding and storage of videos. Video transcoding is an expensive computational process that converts a video from one format to another. Video data are large enough to cause concern over rising storage costs. In the general case, our work is of interest when dealing with expensive computations that generate large results that can be cached for future use. Solving the decision problem entails solving two sub-problems: how long to store cached objects and how many requests we can expect for a particular object in that duration. We compare the proposed approach to always storing and to our previous approach over one year using discrete-event simulations. We observe a 72 % cost reduction compared to always storing and a 13 % reduction compared to our previous approach. This reduction in cost stems from the proposed approach storing fewer unpopular objects when it does not regard it as cost-efficient to do so.

Keywords—Cache storage; Decision theory; Markov processes; Simulation; Transcoding; Utility theory; Web services

I. INTRODUCTION

In this paper we present a decision model for caching expensive computations that produce large amounts of data in a cloud. This model is applicable to cloud services that produce vast quantities of data that can be reused in subsequent requests. As a concrete example, we shall study its application to a cloud-based video transcoding service.

A video transcoding service converts a digital video from one format to another. The need for such a service arises due to the existence of a large number of video compression techniques as well as packaging formats. On the other hand, client devices, and especially mobile devices, can decode and play only a limited number of video formats. If a video is not supported at the client-side device, it needs to be converted into one of the supported formats before playing it. This process is known as video transcoding [1] and it is a CPU-intensive operation. Transcoded videos need to be stored in the server while they are being streamed to the client but they can be deleted from the storage once the streaming operation is completed. Still, there may be new requests for a previously transcoded video in the near future. By storing a transcoded video for an additional amount of time, we can avoid repeating CPU-intensive transcoding operations,

thereby saving relatively large amounts of money. After this additional time, we may reevaluate the circumstances and make a new decision on whether or not to store the video.

In the context of a pay-per-use cloud computing infrastructure, each transcoding operation has a monetary cost due to use of CPU resources, while video storage has a cost based on the amount of data and time to be stored. In order to reduce the operating costs of the service, we need to decide when and for how long each transcoded video should be cached in the storage. A service that stores data that will not be requested in the future will incur unnecessary storage costs. On the other hand, a service that discards data too eagerly is susceptible to incur unnecessary computing costs.

In this paper we study this problem and propose a decision model for cloud-based caches with the objective to reduce operating costs. In previous work on video transcoding we developed a transcoding–storage cost trade-off strategy called cost and popularity score (CPS) [2], which resulted in significantly lower operating costs compared to always storing. This paper improves the decision process by applying utility theory through von Neumann-Morgenstern lotteries, which we have previously made use of for cost-efficient, reliable, utility-based session management in the Cloud [3].

Our utility model for decision making requires three unknown parameters: the storage duration t , the mean number of arrivals $m(t)$ over the storage duration and the popularity distribution p_i of cached objects o_i in the system. We present a natural way of obtaining a good value for the storage duration t , having nice properties that help evaluate the performance of the decision algorithm. We obtain the number of arrivals $m(t)$ over the storage duration t by solving a subproblem consisting of predicting future arrival counts through singular value decomposition. Finally, we employ the Simple Good-Turing frequency estimator to estimate the relative popularity p_i of each cacheable object in the system.

We evaluate the decision making approaches through discrete-event simulations and find that the proposed approach offers 72 % lower cost compared to always storing all requested objects. Compared to our previous approach [2], we see 13 % less cost. This cost reduction stems from the proposed approach storing fewer unpopular objects when it determines that doing so would lead to unnecessary costs.

A. Related Work

There are only a few works in the area of computation and storage trade-off analysis for cost-efficient usage of cloud resources. One of the earlier attempts is Adams et al. [4], which addressed the problem of maximizing efficiency by trading storage for computation. It highlighted some of the important issues involved in constructing a cost-benefit model, which can be used to analyze the trade-offs between computation and storage. However, it did not propose a strategy to find the proper balance between the two resources.

Deelman et al. [5] studied cost and performance trade-offs for an astronomy application using Amazon Elastic Compute Cloud (EC2)¹ and Amazon Simple Storage Service (S3)² cost models. It also examined the trade-offs between three different data management models for cloud storage, namely remote I/O, regular, and dynamic cleanup. The paper concluded that, based on the likelihood of reuse, storing popular datasets in the Cloud can be cost-efficient. However, it did not provide a concrete strategy for cost-efficient computation and storage of scientific datasets in an actual, cloud-based environment.

The Nectar system [6] is designed to automate the management of data and computation in a data center. It initially stores all the derived datasets when they are generated. However, when the available disk space falls below a threshold, all obsolete or least valued datasets are garbage collected to improve resource utilization. Nectar makes use of the usage history of datasets to perform cost-benefit analysis, which determines the usefulness of each dataset. The cost-benefit analysis considers the size of the dataset, the elapsed time since it was last used, the number of times it has been used, and its cumulative computation time. The datasets with the largest cost-to-benefit ratios are deleted. Although Nectar provides a computation and storage trade-off strategy, it is not designed to reduce the total cost of computation and storage in a cloud-based service which makes use of infrastructure as a service (IaaS) resources.

Yuan et al. [7] proposed two strategies for cost-efficient storage of scientific datasets in the Cloud, which compare the computation cost and the storage cost of the datasets, and a cost transitive tournament shortest path (CTT-SP) algorithm to find the best trade-off between the computation and the storage resources. The strategies are called cost rate based storage strategy [8], [9] and local-optimization based storage strategy [10]. The cost rate based storage strategy compares computation cost rate and storage cost rate to decide storage status of a dataset. Whereas, the local-optimization based storage strategy partitions a data dependency graph (DDG) of datasets into linear segments and applies the CTT-SP algorithm to achieve a localized optimization. The local-optimization based storage strategy tends to be more cost-efficient than the cost rate based storage strategy. However,

¹<http://aws.amazon.com/ec2/>

²<http://aws.amazon.com/s3/>

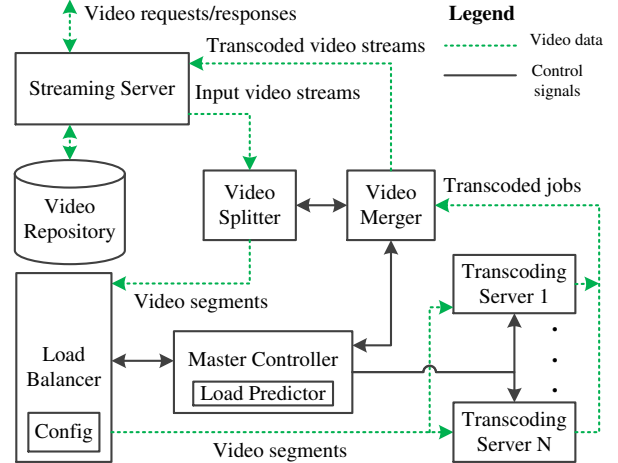


Figure 1. Architecture of the video transcoding service.

due to the overhead introduced by the CTT-SP algorithm, it is less efficient and less scalable. The DDG-based local-optimization based storage strategy of Yuan et al. [10], which provides cost-efficient results for scientific datasets, has little use in video transcoding, which has few data dependencies.

Jokhio et al. [2], [11] presented a computation and storage trade-off strategy called CPS. It estimates an equilibrium point on the time axis where the computation cost and the storage cost of a transcoded video become equal. It also estimates the popularity of the individual transcoded videos to differentiate among videos based on their individual popularity levels.

Kathpal et al. [12] analyzed compute versus storage trade-off for transcoded videos. It proposed an elimination metric to decide which transcoded videos can be removed from the video repository. However, in contrast to the cost and popularity score based strategy of Jokhio et al. [2], it did not account for the video popularity score. Moreover, although the results are also based on Amazon EC2 and Amazon S3, it used rather short videos, with durations no more than 60s.

II. A VIDEO ON-DEMAND SYSTEM

The architecture of the cloud-based, on-demand video transcoding service is shown in Figure 1. It consists of a *streaming server*, a *video splitter*, a *video merger*, a *video repository*, a dynamically scalable cluster of *transcoding servers*, a *load balancer*, a *master controller*, and a *load predictor*. Video requests and responses flow through the *streaming server*. Since our focus in this paper is on computation and storage trade-off for video transcoding, we assume that the *streaming server* will not be a bottleneck.

The video streams are stored in the *video repository* in various compressed formats. The *streaming server* accepts video requests from users and checks if the required video is available in the *video repository*. If it finds the video in the desired format and resolution, it starts streaming the video. However, if it finds that the requested video is stored

only in another format or resolution than the one desired by the user, it sends the video for segmentation and subsequent transcoding. As soon as it receives the transcoded video from the *video merger*, the *streaming server* begins streaming it.

After each transcoding operation, the computation and storage trade-off strategy determines if the transcoded video should be stored in the *video repository* or not. Moreover, if a transcoded video is stored, then the trade-off strategy also determines the duration for which the video should be stored. Therefore, it allows us to trade computation for storage or vice versa in order to reduce the total operating costs and to improve the performance of the video transcoding service.

The *video splitter* splits the video streams into smaller segments called jobs, which are placed into the job queue. Further discussion on video segmentation at the group of pictures (GOP) level is provided in Jokhio et al. [13], [14].

The *load balancer* distributes load on the transcoding servers. In other words, it routes and load balances transcoding jobs on the *transcoding servers*. It maintains a configuration list of active *transcoding servers*. This list is updated often as a result of dynamic virtual machine (VM) allocation and deallocation. The *load balancer* serves the jobs in FIFO (First In, First Out) order. It implements one or more job scheduling policies, such as, the *shortest queue length* policy, which selects a *transcoding server* with the shortest queue length and the *shortest queue waiting time* policy, which selects the *transcoding server* that currently has the shortest available queue waiting time of all the servers.

The actual transcoding is performed by the transcoding servers. They get compressed video segments, perform the required transcoding operations, and return the transcoded video segments for merging. A *transcoding server* runs on a dynamically provisioned VM. Each *transcoding server* processes one or more simultaneous jobs. When a transcoding job arrives at a *transcoding server*, it is placed in the server's queue, from where it subsequently will be processed further.

The *master controller* acts as the main controller and resource allocator. It implements prediction-based dynamic resource allocation and deallocation algorithms [15] and one or more computation and storage trade-off strategies. The resource allocation and deallocation is mainly based on the target play rate of the video streams and the predicted transcoding rate of the transcoding servers. The *master controller* uses the *load predictor* for load prediction [16]. The *video merger* merges the transcoded jobs into video streams, which form video responses. Our resource allocation and load prediction algorithms are described in detail in Jokhio et al. [15] and Ashraf et al. [16]. In this paper, we focus on a cost-efficient computation and storage trade-off.

III. UTILITY MODEL FOR COMPUTING VERSUS STORING

In this section we present the utility model that will be used to govern the caching strategy of the service. We proceed by introducing the basics of von Neumann-Morgenstern lotteries.

A. Von Neumann-Morgenstern Lotteries

A von Neumann-Morgenstern lottery [17] consists of mutually exclusive outcomes that may occur with a given probability. The sum of probabilities in a lottery should be equal to one. For example, the simple lottery L described by

$$L = 0.20A + 0.80B, \quad (1)$$

denotes a scenario where the probability of event A , $P(A) = 0.20$, the probability of event B , $P(B) = 0.80$, and exactly one of the possible outcomes will occur. In general, a lottery L with n outcomes A_i and probabilities p_i is expressed as:

$$L = \sum_{i=1}^n p_i A_i \quad (2)$$

$$\text{subject to } \sum_{i=1}^n p_i = 1.$$

According to the von Neumann-Morgenstern utility theorem [17], an agent faced with the problem of choosing between a set of lotteries has a utility function, provided that the four axioms of the theorem are satisfied. The four axioms of the utility theorem on lotteries L , M and N are:

- *completeness* (L or M is preferred, or they are equal)
 $L \leq M \vee M \leq L$
- *transitivity* (consistent preference across 3 operations)
 $(L \leq M \wedge M \leq N) \rightarrow L \leq N$
- *continuity* (transitive preference is continuous)
 $(L \leq M \wedge M \leq N) \rightarrow \exists p \in [0, 1] pL + (1-p)N = M$
- *independence* (independence of irrelevant alternatives)
 $L < M \rightarrow \forall N \forall p \in (0, 1] pL + (1-p)N < pM + (1-p)N.$

If an agent satisfies these axioms, it has a utility function u , assigning a real value $u(A)$ to every possible outcome A , so that for any two lotteries L and M , $Eu(L)$ is the expected value of u in L , $Eu(M)$ is the expected value of u in M and

$$L < M \leftrightarrow Eu(L) < Eu(M). \quad (3)$$

By using the utility function we can determine which lotteries to play. In this paper we will model the choice between computing and storing as von Neumann-Morgenstern lotteries. Choosing among lotteries, we can make the best decisions.

There are, however, some limitations to von Neumann-Morgenstern utility. Von Neumann and Morgenstern [17] acknowledged that nested gambling is ignored. An example of nested gambling with lotteries L and M would be $pL + (1-p)M$, which gets treated as a lottery itself. Another limitation is that utilities cannot be compared between agents X and Y with different utility functions u_X and u_Y . Expressions like $u_X(L) + u_Y(L)$ are undefined. As we use neither nested gambling nor multiple agents, these limitations do not affect us. We may design a utility function that incorporates risk aversion or diminishing returns, which could be beneficial in an environment with a relatively high degree of uncertainty.

We assume that we can model requests arriving to the system as an inhomogeneous Poisson process $N(t)$ with rate $\lambda(t)$, where $N(t)$ is the number of arrivals by time t . The mean number of arrivals $m(t)$ by time t is obtained through

$$m(t) = \int_0^t \lambda(u) du. \quad (4)$$

$N(t)$ has a Poisson distribution with mean parameter $m(t)$:

$$P(N(t) = k) = \frac{m(t)^k}{k!} e^{-m(t)}. \quad (5)$$

Using Poisson splitting, we can model the requests arriving for each video as independent Poisson processes $N_i(t)$ with mean number of arrivals $p_i m(t)$, where p_i is the relative frequency of requests for video i . Each process $N_i(t)$ then has a Poisson distribution with parameter $p_i m(t)$, defined as

$$P(N_i(t) = k) = \frac{(p_i m(t))^k}{k!} e^{-p_i m(t)}. \quad (6)$$

Requests may arrive to the system at any time. Whenever a video is requested, we check if a cached copy is available. If a cached copy is available, we serve the request from the cache at no extra cost. If there is no cached copy available, we need to transcode the corresponding source video into the correct format, this will incur a fixed cost c_t of transcoding. After the video has been transcoded, we have the option of storing the result in the cache for duration t at cost $c_s t$, after which we can continue storing it for a new duration t at cost $c_s t$. Alternatively, we may delete it at no further cost.

If we decide to cache a transcoded video for duration t it will always cost $c_s t$, regardless of whether any further requests arrive for the video or not. If we decide not to cache the transcoded video, one of two possible outcomes will occur: either further requests for that video arrive, say $n_i \geq 1$ requests, at which point we will have to transcode the source video again at cost $n_i c_t$, or no further requests arrive, costing us nothing. Thus, we may consider the following outcomes:

- A: Delete, no requests arrive
- B: Delete, requests arrive
- C: Store, no requests arrive
- D: Store, requests arrive.

Using (6), we can compute the probability of 0 arrivals, $P(N_i(t) = 0) = e^{-p_i m(t)}$. Conversely, the probability of more than 0 arrivals is $1 - P(N_i(t) = 0) = 1 - e^{-p_i m(t)}$. Assuming that n_i is approximately equal to the mean of the corresponding zero-truncated Poisson process $n_i \approx \frac{p_i m(t)}{1 - e^{-p_i m(t)}}$, we obtain the following utilities for each possible outcome:

$$\begin{aligned} u(A) &= 0 \\ u(B) &= -C_{ti} \frac{p_i m(t)}{1 - e^{-p_i m(t)}} \\ u(C) &= -C_{si} t \\ u(D) &= -C_{si} t. \end{aligned}$$

We can then formulate the alternatives as von Neumann-Morgenstern lotteries, L_d for deleting and L_s for storing:

$$L_d = P(N_i(t) = 0)A + (1 - P(N_i(t) = 0))B \quad (7)$$

$$L_s = P(N_i(t) = 0)C + (1 - P(N_i(t) = 0))D. \quad (8)$$

The two lotteries then have the following expected utilities:

$$Eu(L_d) = -p_i m(t) C_{ti} \quad (9)$$

$$Eu(L_s) = -C_{si} t. \quad (10)$$

Always choosing the lottery with the highest expected utility should give us the best result in the long run. However, if we are to compute the expected utility of each lottery, we need actual values for $m(t)$ as well as p_i , none of which are directly observable. In Sections IV and V we present possible methods for estimating these parameters. It would also be possible to use other estimators for specific problems.

B. Determining the Duration for Caching

To decide whether to cache the results of a computation, we first need to determine when to make this decision. If we decide not to store a transcoded video, we naturally cannot make further decisions for that video until it has been transcoded again, as the data are discarded. If we decide to store a transcoded video, we also need to determine how long to store it. Storing indefinitely is not viable. Jokhio et al. [2], [11] determined that there is an equilibrium point τ , where the cost of storing over time C_s becomes equal to the cost of transcoding C_t . Figure 2 illustrates how storage cost and transcoding cost are related. For example, if the cost of transcoding a given video v_i is $C_{ti} = \$1$ and the cost of storing that video is $C_{si} = \$0.5 \text{ d}^{-1}$, storing for $\frac{\$1}{\$0.5 \text{ d}^{-1}} = 2 \text{ d}$ results in the same cost as transcoding once: $C_{ti} = \tau C_{si} = C_i$. This means that if we store a video v_i for duration τ and get at least one request for it, say n_i requests, we will break even or even save money compared to transcoding it n_i times, as $C_i \leq n_i C_t, n_i \geq 1$. Setting $t = \tau$ in (9) results in a simplified expression, in which the costs have now been made equal:

$$Eu(L_d) = -p_i m(\tau) C_i \quad (11)$$

$$Eu(L_s) = -C_i. \quad (12)$$

Thus, we can conclude that we should decide to store when

$$p_i m(\tau) \geq 1. \quad (13)$$

IV. ARRIVAL RATE PREDICTION

Obtaining $m(\tau)$ in (13) requires knowledge of the future arrival rate up to time τ . We have decided to use a time-series prediction approach based on truncated singular value decomposition, as outlined by Shen and Huang [18], which presented methods for predicting arrivals to a call center.

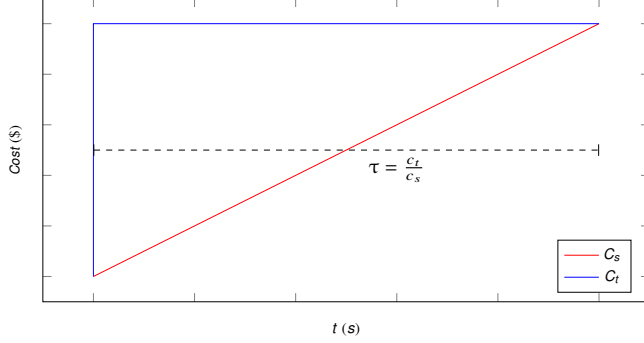


Figure 2. Storage cost C_s approaches transcoding cost C_t .

A. Truncated Singular Value Decomposition

Let X be an $n \times m$ matrix recording the number of requests for n days, each day having m time periods. The rows of this matrix constitute a vector-valued time series in \mathbb{R}^m . We wish to build a time-series model and forecast future values. The dimensionality of this time series is large. We can significantly reduce the dimensionality with a decomposition

$$\mathbf{x}_i = \beta_{i1}\mathbf{f}_1 + \dots + \beta_{iK}\mathbf{f}_K + \boldsymbol{\epsilon}_i, \quad i = 1, \dots, n, \quad (14)$$

where $\mathbf{f}_1, \dots, \mathbf{f}_K \in \mathbb{R}^m$ are basis vectors, $\boldsymbol{\epsilon}_1, \dots, \boldsymbol{\epsilon}_n \in \mathbb{R}^m$ are the corresponding error terms, and $\beta_{ik} \in \mathbb{R}$ are scalars.

The singular value decomposition (SVD) of X is given by

$$X = USV^T \quad (15)$$

and gives the solution for minimizing the error terms in (14):

$$\mathbf{x}_i \approx s_1 u_{i1} \mathbf{v}_1 + \dots + s_K u_{iK} \mathbf{v}_K, \quad (16)$$

where s_1, \dots, s_K are the K largest singular values from the diagonal of S , u_{i1}, \dots, u_{iK} , $i = 1, \dots, n$ are the entries of column i in U and $\mathbf{v}_1, \dots, \mathbf{v}_K$ are the corresponding columns of V . We can now forecast each series $\{\beta_{ik}\}$ separately, which we do by extrapolating from linear regression. We can then predict the number of arrivals in future days accordingly:

$$\hat{\mathbf{x}}_{n+h} = \hat{\beta}_{n+h,1}\mathbf{f}_1 + \dots + \hat{\beta}_{n+h,K}\mathbf{f}_K. \quad (17)$$

However, to keep predictions non-negative, we employ the root-unroot method of Brown et al. [19]: we transform X to $\sqrt{X + \frac{1}{4}}$ before the SVD and adjust the prediction by $\hat{\mathbf{x}}^2 - \frac{1}{4}$.

SVD generally requires $O(m^2n + n^3)$ operations on an $m \times n$ real matrix. However, as we only want the K most important singular values, we can use efficient algorithms for truncated singular value decomposition. We therefore use the implicitly restarted Lanczos bidiagonalization method of Baglama and Reichel [20], implemented in the `irlbpy` library [21]. In practice, this method requires as little as $O(mnK)$ operations [21], scaling linearly with size of data.

V. FREQUENCY ESTIMATION

Suppose there are five kinds of videos: $\{\mathcal{A}, \mathcal{B}, \mathcal{C}, \mathcal{D}, \mathcal{E}\}$. You observe the number of times each video is requested and find 2 requests for \mathcal{A} , 2 requests for \mathcal{B} , 3 requests for \mathcal{C} and 1 request for \mathcal{D} . A naïve frequency estimator, like the maximum likelihood frequency estimator, will assign a probability p_i of .25 to \mathcal{A} , .25 to \mathcal{B} , .375 to \mathcal{C} and .125 to \mathcal{D} . Notice that no requests for \mathcal{E} were observed, resulting in the maximum likelihood frequency estimator assigning a probability of 0 for requests to the video. However, this must be false, because we know that it is possible to request the video. Video \mathcal{E} should therefore have a probability greater than zero. The maximum likelihood estimator does not account for any missing samples. This is the reason why it tends to underestimate rare videos and why the maximum likelihood frequency estimator may not be the best option.

A. The Good-Turing Frequency Estimator

Since we are most concerned about low expected numbers of requests, we are dealing with rare events of which we might not have made any observations. Similarly, when new videos are made available, we will again be dealing with a lack of information. This is why we propose to use the Simple Good-Turing frequency estimator [22], which accounts for unobserved events. The estimator tends to underestimate frequent items, but this can be mitigated by using the empirical estimate for them [23]. As we are really only interested in infrequent items where $p_i m(t) \approx 1$, it is not necessary for our particular case. For the five videos, the Simple Good-Turing frequency estimator assigns a probability of .22 to \mathcal{A} , .22 to \mathcal{B} , .285 to \mathcal{C} , .15 to \mathcal{D} , and .125 to \mathcal{E} .

The Good-Turing frequency estimators assume that the observed items follow a binomial distribution [22], but we assume that they follow a Poisson distribution. However, this is not a problem in our case, as the binomial distribution converges to the Poisson distribution as $np = \lambda, n \rightarrow \infty, p \rightarrow 0$. With a large number of observable items, the probabilities will be small for infrequent items. The Good-Turing frequency estimators also assume that the underlying frequency distribution is static. We postulate that we can relax this assumption by using a sliding window. We only need to ensure that the window is large enough, compared to the mean number of arrivals $m(t)$, to give a suitable number of observations n within the duration of the time window.

VI. EVALUATION USING DISCRETE-EVENT SIMULATIONS

We used SimPy³ to develop a discrete-event simulation of a video transcoding service. We chose to compare the proposed, utility-based approach with the previously developed CPS policy and a reference policy based on always opting to store.

³<http://simpy.readthedocs.org/>

A. Setting up the Experiment

We simulated a video transcoding service with 10000 videos over one year. Video sizes were randomly assigned according to a double Pareto-lognormal distribution [24] with parameters $\alpha = 2, \beta = 4, \mu = 0, \sigma = 1$ scaled by a factor of 256 MiB. The probability of a request to the system belonging to a particular video is given by a truncated Pareto distribution with parameters $x_m = 1, \alpha = 2$. Transcoding cost $\$1.7 \times 10^{-5} \text{ s}^{-1}$, the cost of a medium instance in Amazon EC2. Transcoding rate was 2.4 MiB s^{-1} . Storage cost $\$3.6 \times 10^{-11} \text{ MiB}^{-1} \text{ s}^{-1}$, as in Amazon S3. Thus, the time to store each video was a constant, calculated through

$$\tau = \frac{\$1.7 \times 10^{-5} \text{ s}^{-1}}{2.4 \text{ MiB s}^{-1} \times \$3.6 \times 10^{-11} \text{ MiB}^{-1} \text{ s}^{-1}} = 55 \text{ h.} \quad (18)$$

The arrival process was a randomly generated inhomogeneous Poisson process. We constructed the arrival process by thinning a homogeneous Poisson process with rate $\lambda = 20 \text{ s}^{-1}$ using a Bernoulli trial with a time-dependent probability $p(t)$. We generated the probability vector by dividing the simulation duration into a random number X of parts according to a Poisson distribution with mean $\mu = 26$. We then generated an exponentially distributed duration with mean $\omega = \frac{t_{\text{total}}}{X}$ for each interval. With this information we finally generated a linear spline with random coefficients, assuming values between 0 and 1. Figure 3 shows the resulting mean rate.

We sought to compare the proposed, utility-based approach with the established CPS approach. We also included the always store policy and a version of the utility-based approach with perfect knowledge of the mean number of arrivals $m(t)$ and popularity p_i in the benchmarks, acting as references. To enable direct comparison of the approaches by classifying each decision as good, bad or neutral, we used a minimum storage duration of $SD_{\tau_i} = 55 \text{ h}$ for the CPS approach. This resulted in marginally higher cost for this approach, compared to when using a minimum duration of 24 h, but this increase was not significant enough to alter the outcome of the experiments. The reference policy always chose storing.

For the utility-based approach, the rate predictor used a sliding past window of 7 d. The popularity estimator also used a sliding window of 7 d. Because the system cannot make good predictions before historical data have been collected, all requested videos were always stored during the first 7 d. Samples of request counts were collected every hour: $t_s = 1 \text{ h}$.

B. Results

Figure 3 also shows the predicted arrival rate obtained through truncated singular value decomposition. The prediction accuracy appears good considering the simplicity of the used approach. The total operating cost over time is shown in Figure 4. Always storing cost \$605, CPS cost \$196 and the utility-based approach cost \$171 (\$164 with perfect information). The utility-based approach clearly operates at a lower cost than the other approaches. Compared to always

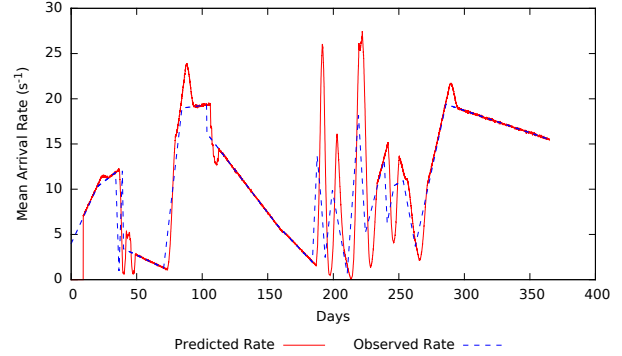


Figure 3. Observed and predicted arrival rate for the system.

storing, the utility based approach operated at $1 - \frac{\$171}{\$605} = 72 \%$ less cost. Compared to the CPS approach, the utility based approach operated at $1 - \frac{\$171}{\$196} = 13 \%$ less cost. With perfect information, the utility-based approach operated at $1 - \frac{\$164}{\$171} = 4 \%$ less cost than the actual implementation. This is a small difference, indicating that the proposed implementation performed close to its theoretical optimum. The reason why the proposed approach fares better than the previously developed approach can be seen by classifying each decision as good, bad or neutral, as shown in Figure 5.

C. Analysis

The utility-based approach operates at less cost than the CPS approach because it did fewer bad decisions, as seen in Figures 5a and 5b. A bad decision is either a bad decision to store or a bad decision to delete. A bad decision to store means that no requests arrived for the video in question during duration τ , meaning that it cost us C_i when it could have cost us \$0 if we had chosen the other alternative. Conversely, a bad decision to delete means that more than one request arrived for the video in duration τ , which cost us more than storing at cost C_i would have done. When only one request arrives in duration τ , it does not matter whether we store or delete, as storing for duration τ costs the same as transcoding

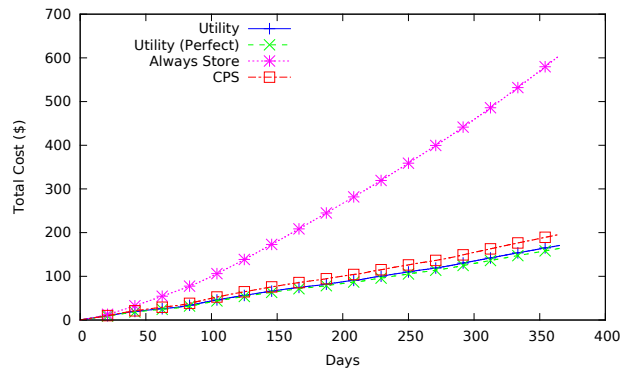
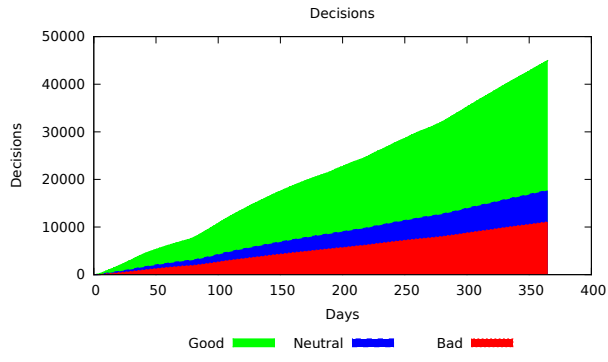
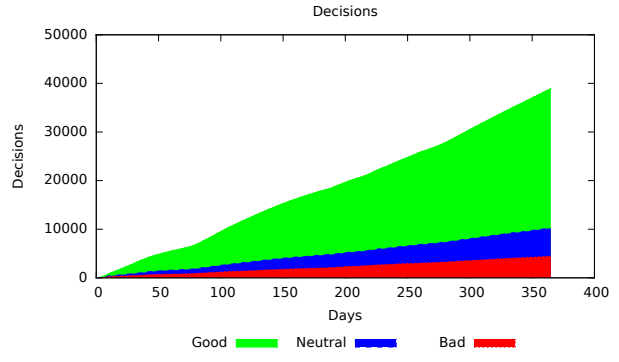


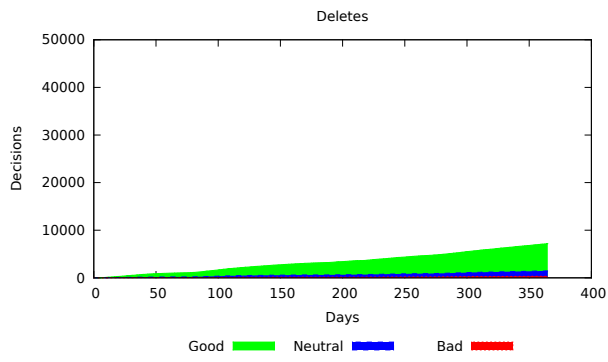
Figure 4. Cost of operation for the approaches over one year.



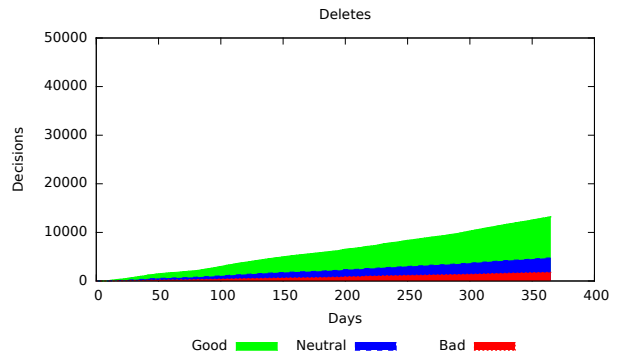
(a) Decisions by the CPS approach.



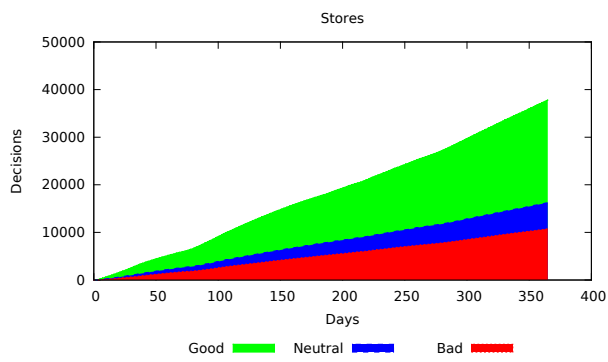
(b) Decisions by the utility-based approach.



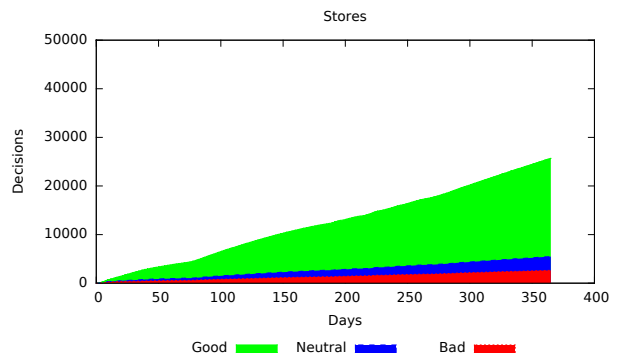
(c) Deletes by the CPS approach.



(d) Deletes by the utility-based approach.



(e) Stores by the CPS approach.



(f) Stores by the utility-based approach.

Figure 5. Classification of decisions made by the two actual approaches as good, bad or neutral.

Table I
BAD, NEUTRAL AND GOOD OPERATIONS PER APPROACH

Operation	Approach	Bad	Neutral	Good
Delete	CPS	4 %	16 %	80 %
	Utility	14 %	22 %	64 %
Store	CPS	28 %	15 %	57 %
	Utility	10 %	11 %	79 %
All	CPS	25 %	15 %	60 %
	Utility	11 %	15 %	74 %

once: C_i . These decisions are thus neutral. Table I shows the proportion of bad, neutral, and good delete and store operations for each approach. CPS made 44 973 decisions, comprising 7172 deletes and 37 801 stores. The utility-based approach made 38 894 decisions, comprising 13 214 deletes and 25 680 stores. While the utility-based approach did slightly more erroneous deletes than the CPS approach, it still did far fewer erroneous stores. CPS deletes too rarely.

VII. CONCLUSIONS

We presented a utility-based decision strategy for caching expensive computations in a cloud-computing setting. We exemplified our approach with a video transcoding service, but we believe that a similar approach can be used in other services that transcode, uncompress or index large amounts of data based on user requests. Our model requires three unknown parameters: the storage duration t , the mean number of arrivals $m(t)$ over the storage duration and the popularity distribution p_i of objects to store in the cache o_i . We presented a formally justified way of obtaining a good value for the storage duration t , having good properties that helped evaluate the performance of the decision algorithm. We obtained the mean number of arrivals $m(t)$ over the storage duration t by solving a subproblem consisting of predicting future arrival counts through singular value decomposition. Finally, we used the Simple Good-Turing frequency estimator to estimate the relative popularity p_i of each available video.

We evaluated our approach using discrete-event simulations. The utility-based approach incurred 72 % less cost than always storing and 13 % less cost than the CPS approach [2] over one simulated year. Prediction accuracy was high, as evidenced by only obtaining 4 % less cost when using perfect information. This made the utility-based approach better. We only evaluated the system with a static popularity distribution, but have accounted for non-static distributions by using a sliding window for popularity estimation. Determining the sizes of the sliding windows used by the arrival rate predictor and popularity estimator is a domain-specific problem which we have only touched upon lightly. Given a good estimator for non-static popularity distributions, the utility-based approach should work well. We have assumed independent arrivals, which is not always true of a real-world service. However, our approach can also be evaluated against actual load traces.

ACKNOWLEDGMENTS

This work was supported by the N4S research project. Adnan Ashraf, Benjamin Byholm and Fareed Jokhio were supported by the Foundation of Nokia Corporation. Adnan Ashraf and Fareed Jokhio were also supported by scholarships from the Higher Education Commission (HEC) of Pakistan.

REFERENCES

- [1] A. Vetro, C. Christopoulos, and H. Sun, "Video transcoding architectures and techniques: An overview," *Signal Processing Magazine, IEEE*, vol. 20, no. 2, pp. 18–29, mar 2003.
- [2] F. Jokhio, A. Ashraf, S. Lafond, and J. Lilius, "A computation and storage trade-off strategy for cost-efficient video transcoding in the cloud," in *39th EUROMICRO Conference on Software Engineering and Advanced Applications*. IEEE Computer Society, 2013, pp. 365–372.
- [3] B. Byholm and I. Porres, "Cost-efficient, reliable, utility-based session management in the cloud," in *14th IEEE/ACM International Symposium on Cluster, Cloud, and Grid Computing*, ser. IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing. IEEE Computer Society, 2014, pp. 102–111.
- [4] I. F. Adams, D. D. E. Long, E. L. Miller, S. Pasupathy, and M. W. Storer, "Maximizing efficiency by trading storage for computation," in *Proceedings of the 2009 conference on Hot topics in cloud computing*, ser. HotCloud'09. USENIX Association, 2009.
- [5] E. Deelman, G. Singh, M. Livny, B. Berriman, and J. Good, "The cost of doing science on the cloud: The montage example," in *Proceedings of the 2008 ACM/IEEE conference on Supercomputing*, ser. SC '08. IEEE Press, 2008, pp. 50:1–50:12.
- [6] P. K. Gunda, L. Ravindranath, C. A. Thekkath, Y. Yu, and L. Zhuang, "Nectar: automatic management of data and computation in datacenters," in *Proceedings of the 9th USENIX conference on Operating systems design and implementation*, ser. OSDI'10. USENIX Association, 2010, pp. 1–8.
- [7] D. Yuan, Y. Yang, X. Liu, and J. Chen, "Computation and storage trade-off for cost-effectively storing scientific datasets in the cloud," in *Handbook of Data Intensive Computing*. Springer New York, 2011, pp. 129–153.
- [8] —, "A cost-effective strategy for intermediate data storage in scientific cloud workflow systems," in *Parallel Distributed Processing (IPDPS), 2010 IEEE International Symposium on*, 2010, pp. 1–12.
- [9] D. Yuan, Y. Yang, X. Liu, G. Zhang, and J. Chen, "A data dependency based strategy for intermediate data storage in scientific cloud workflow systems," *Concurrency and Computation: Practice and Experience*, vol. 24, no. 9, pp. 956–976, 2012.
- [10] D. Yuan, Y. Yang, X. Liu, and J. Chen, "A local-optimisation based strategy for cost-effective datasets storage of scientific applications in the cloud," in *Cloud Computing (CLOUD), 2011 IEEE International Conference on*, 2011, pp. 179–186.

- [11] F. A. Jokhio, A. Ashraf, S. Lafond, I. Porres, and J. Lilius, "Cost-efficient dynamically scalable video transcoding in cloud computing," Turku Centre for Computer Science (TUCS), Tech. Rep. 1098, 2013.
- [12] A. Kathpal, M. Kulkarni, and A. Bakre, "Analyzing compute vs. storage tradeoff for video-aware storage efficiency," in *Proceedings of the 4th USENIX Conference on Hot Topics in Storage and File Systems*, ser. HotStorage'12. USENIX Association, 2012, pp. 13–13.
- [13] F. Jokhio, T. Deneke, S. Lafond, and J. Lilius, "Bit rate reduction video transcoding with distributed computing," in *Parallel, Distributed and Network-Based Processing (PDP), 2012 20th Euromicro International Conference on*, 2012, pp. 206–212.
- [14] F. A. Jokhio, T. Deneke, S. Lafond, and J. Lilius, "Analysis of video segmentation for spatial resolution reduction video transcoding," in *Intelligent Signal Processing and Communications Systems (ISPACS), 2011 International Symposium*. IEEE, 2011.
- [15] F. Jokhio, A. Ashraf, S. Lafond, I. Porres, and J. Lilius, "Prediction-based dynamic resource allocation for video transcoding in cloud computing," in *Parallel, Distributed and Network-Based Processing (PDP), 21st Euromicro International Conference on*. IEEE Computer Society, 2013, pp. 254–261.
- [16] A. Ashraf, B. Byholm, and I. Porres, "A session-based adaptive admission control approach for virtualized application servers," in *Utility and Cloud Computing (UCC), 5th IEEE/ACM International Conference on*. IEEE Computer Society, 2012, pp. 65–72.
- [17] J. von Neumann and O. Morgenstern, *Theory of Games and Economic Behavior*, 3rd ed. Princeton, NJ: Princeton University Press, 1953.
- [18] H. Shen and J. Z. Huang, "Interday forecasting and intraday updating of call center arrivals," *Manufacturing and Service Operations Management*, vol. 10, no. 3, pp. 391–410, 2008.
- [19] L. Brown, N. Gans, A. Mandelbaum, A. Sakov, H. Shen, S. Zeltyn, and L. Zhao, "Statistical analysis of a telephone call center: A queueing-science perspective," *Journal of the American Statistical Association*, vol. 100, no. 469, pp. 36–50, 2005.
- [20] J. Baglama and L. Reichel, "Augmented implicitly restarted Lanczos bidiagonalization methods," *SIAM Journal on Scientific Computing*, vol. 27, no. 1, pp. 19–42, 2005.
- [21] J. Baglama, M. Kane, B. Lewis, and L. Reichel, *irlbpy: Truncated SVD by Implicitly Restarted Lanczos Bidiagonalization for NumPy*, 2013.
- [22] W. A. Gale and G. Sampson, "Good-Turing frequency estimation without tears," *Journal of Quantitative Linguistics*, vol. 2, no. 3, pp. 217–237, 1995.
- [23] A. Orlitsky, N. P. Santhanam, and J. Zhang, "Always Good Turing: Asymptotically optimal probability estimation," *Science*, vol. 302, no. 5644, pp. 427–431, 2003.
- [24] W. J. Reed and M. Jorgensen, "The double Pareto-lognormal distribution—A new parametric model for size distributions," *Communications in Statistics-Theory and Methods*, vol. 33, no. 8, pp. 1733–1753, 2004.