

A Learning Architecture for Scheduling Workflow Applications in the Cloud

Enda Barrett

School of Engineering & Informatics
National University of Ireland, Galway
Galway, Ireland

Email: Enda.Barrett@nuigalway.ie

Enda Howley

School of Engineering & Informatics
National University of Ireland, Galway
Galway, Ireland

Email: Enda.Howley@nuigalway.ie

Jim Duggan

School of Engineering & Informatics
National University of Ireland, Galway
Galway, Ireland

Email: James.Duggan@nuigalway.ie

Abstract—The scheduling of workflow applications involves the mapping of individual workflow tasks to computational resources, based on a range of functional and non-functional quality of service requirements. Workflow applications such as scientific workflows often require extensive computational processing and generate significant amounts of experimental data. The emergence of cloud computing has introduced a utility-type market model, where computational resources of varying capacities can be procured on demand, in a pay-per-use fashion. In workflow based applications dependencies exist amongst tasks which requires the generation of schedules in accordance with defined precedence constraints. These constraints pose a difficult planning problem, where tasks must be scheduled for execution only once all their parent tasks have completed. In general the two most important objectives of workflow schedulers are the minimisation of both cost and makespan. The cost of workflow execution consists of both computational costs incurred from processing individual tasks, and data transmission costs. With scientific workflows potentially large amounts of data must be transferred between compute and storage sites. This paper proposes a novel cloud workflow scheduling approach which employs a Markov Decision Process to optimally guide the workflow execution process depending on environmental state. In addition the system employs a genetic algorithm to evolve workflow schedules. The overall architecture is presented, and initial results indicate the potential of this approach for developing viable workflow schedules on the Cloud.

I. INTRODUCTION

The recent advancement of Cloud computing and its pay per use model enable the procurement of large amounts of computational resources on demand. Cloud providers leveraging on large economies of scale are capable of delivering increasing amounts of computational resources at lower costs and with greater reliability. Cloud computing delivers computational resources by means of virtualisation technologies. Through the instantiation of virtual machines users can deploy their applications on resources with varying performance and cost levels. Most recently numerous scientific experiments in fields such as astronomy [1], epigenomics [2] and neuroscience involve complex computational analyses on large data sets. These communities rely greatly on insights gained through computational analysis in order to advance the state of the art in their respective communities. Many of these computational experiments can be characterised neatly as workflows, with large segments capable of being executed in parallel.

Initially the processing of many large scientific models and application workflows involved a significant investment in high performance computing infrastructures. The initial monetary outlay, coupled with setup and on-going maintenance costs, rendered it prohibitive to most. Apart from a small number of well funded research projects most didn't have the resources. Advancements in grid computing technologies led to the development of large computational grids, where through virtual organisations the scientific community could share resources among each other in a large distributed grid network. In response to the sharing of their computational resources, members of the virtual organisation were allocated pre-reserved slots for their own experiments. However these platforms possessed a number of limitations both technical and bureaucratic [3]. Often minor research groups could not gain access to the global grids for political reasons. Furthermore, many of these grids use pre-defined platforms with specific operating systems, application suites and API's. Platform limitations present a clear barrier to enabling workflows to execute without a significant technical overhead. Recently there has been significant high profile attention regarding the possibilities of executing scientific workflows on the cloud [4]. A recent study highlighted the cost effectiveness of executing scientific workflows on Amazon EC2¹ [5]. However with the market oriented model of clouds, existing workflow management tools need to be adapted to better utilise cloud platforms. Previous workflow scheduling approaches for grids focussed on resources which were reserved in advance. However through virtual machines greater amounts of resources can be procured in real time and on demand. To address this, scheduling approaches for clouds should be dynamic (online) and adaptive to system behaviours. They should also be capable of adjusting to fluctuating costs for communication and computation. To address these issues we propose a novel cloud workflow scheduling architecture capable of scheduling application workflows in a cloud computing environment. We schedule workflow applications based on user specified QoS constraints, namely cost and makespan. We adopt a multifaceted approach to scheduling where a genetic algorithm is used to generate optimal schedules and a Markov Decision

¹<http://aws.amazon.com/ec2/>

Process (MDP) chooses from among them based on the observed state of the environment. The virtualised nature of cloud computing (Disk and Cpu), with multiple virtual machines deployed on the same physical hosts, poses problems predicting task execution and data transfer times. Temporal variability such as peak usage times with large network traffic should be taken into consideration for any planning approach scheduling workflows on the cloud. The main contribution of this work is the formalism of the scheduling problem as a MDP, creating a mapping from states to evolved schedules. This is achieved through the evolution feasible schedules to the planning problem, and then selecting the appropriate solution based on environmental state. In conjunction with this we adopt a Bayesian Model learning algorithm enabling the learning agent to reason about its choices, based on actual environmental feedback.

The rest of this paper is structured as follows: Background Research provides an overview of relevant and related work in this field. A number of aspects of workflow scheduling and cloud computing are discussed. We also detail the algorithms we use to solve the scheduling problem. Workflow architecture outlines the specific components of our cloud workflow scheduling approach. Cloud workflow scheduling details the problem of scheduling workflows onto cloud resources. Initial Results details our preliminary findings, leading to Conclusions & Future Work.

II. BACKGROUND RESEARCH

Extensive research has taken place on the scheduling of workflow applications onto distributed resources in grid computing. Yu et al [6] proposed a workflow scheduling approach based on using MDPs to schedule tasks on utility grids. Their approach partitions the workflow into branches, where they allocate sub-deadlines to each branch based on the overall deadline. A branch containing sequential tasks is then mapped to services using a MDP. The main objective is to meet the overall deadline, concurrently minimising cost. In this context the MDP binds resources to tasks based on the policy obtained through the solution of the MDP via value iteration. A scheduling approach based on genetic algorithms was also developed by Yu et al. [7]. Depending on user preference the approach was constrained to optimising either execution cost or the overall makespan. This work was extended further [8] with the inclusion of multiobjective evolutionary algorithms. Our evolutionary scheduling approach is similar, where a genetic algorithm is used to evolve solutions to the workflow scheduling problem. However our approach differs to these in two regards, firstly we schedule on clouds where virtual machines of varying capacity can be created on demand and without advanced reservation. Secondly we employ a MDP to optimally choose from amongst the evolved schedules and not scheduling individual tasks to resources. Our MDP approach then chooses from amongst these evolved schedules, choosing the approach which maximises the expected reward in the long run.

A number of grid and cluster workflow management tools have been proposed to address the scheduling problem. Many of these have been devised to facilitate the execution of scientific workflows. Tools such as Pegasus [9], Kepler [10] and Taverna workbench [11] have all successfully executed scientific workflows on computational grids and clusters. These generally schedule tasks based on earliest finish time, earliest starting time or high processing capabilities. Collectively these can be considered as "best resource selection" (BRS) [12] approach, where a resource is selected solely based on its performance.

Pandey et al. recently developed a cloud scheduling approach based on Particle Swarm Optimisation (PSO) [12]. Experiments have shown that PSOs can provide a good distribution of tasks across the available services, reducing resource overloading and performing better than BRS based approaches. The improvement in performance is also due to the fact that the PSO scheduler factors in communication costs between all tasks in the workflow when developing a feasible schedule. Wu et al. [13] recently proposed a market oriented hierarchical scheduling approach to cloud workflow scheduling. They adopt a two stage scheduling approach, aiming to optimise scheduling at the overall service level and virtual machine level. They employ a number of meta-heuristic algorithms to evolve different schedules allowing the user to choose a schedule from amongst those returned.

Whilst the above work shows promise in improving the performance of workflow scheduling on clouds, their focus on the development of heuristic and meta-heuristic algorithms to improve performance, is limited by the inherent variability in the environment. There reliance on advanced modelling and task execution prediction, reduces their effectiveness in truly dynamic environments like clouds. We contend that a workflow management system should possess a certain amount of environmental knowledge which can facilitate the scheduling process.

To reason about environmental knowledge we adopt a Bayesian model learning algorithm, a subset of machine learning closely allied with reinforcement learning. Learning techniques such as this allow the system to "keep up" with the constantly changing underlying environment. A number of studies have looked at applying reinforcement learning to resource allocation problems [14]. This work presented a framework using reinforcement learning, capable of dynamically allocating resources in a distributed system. Reinforcement learning methods are capable of closely approximating true action value functions in scenarios where a complete model of the environment is not available. They allow the learning agent to learn through observed rewards as they explore the state space. While resource allocation is more concerned with low level scheduling of tasks at the virtual machine level, the parallels between them still merit their inclusion. Tesauro investigated the use of a hybrid reinforcement learning technique for autonomic resource allocation [15]. He applied this research to optimizing server allocation in data centers. Germain-Renaud et al. [16] looked at similar resource allocation issues. Here a

workload demand prediction technique was used to predict the resource allocation required each time. Reinforcement learning has also been successfully applied to grid computing as a job scheduler. Here the scheduler can seamlessly adapt its decisions to changes in the distributions of inter-arrival time, QoS requirements, and resource availability [17]. The objective of applying reinforcement learning to resource allocation in these works is to schedule tasks independent of constraints across a number of resources. Our work differs from these in that we schedule tasks with precedence constraints and our objective is to optimise both costs and makespan.

A. Markov Decision Processes

Markov Decision Processes are a particular mathematical framework suited to modelling decision making under uncertainty. They can be represented as a four tuple consisting of states, actions, transition probabilities and rewards (S, A, P, R). In general the learning agent interacts with its environment through a sequence of discretized time steps. At the end of each time period t the agent occupies state $s_t \in S$, where S represents the set of all possible states. The agent then chooses an action $a_t \in A(s_t)$, where $A(s_t)$ is the set of all possible actions within state s_t . The transition probability P defines the probabilities of moving from one state to the next.

$$P_{s,s'}^a = Pr\{s_{t+1} = s' | s_t = s, a_t = a\} \quad (1)$$

R is the expected value of the next reward [18]. Figure 1 depicts a simple two state MDP with two possible actions for each state. In Figure 1, choosing action A_1 in *State 1* will lead you to *State 2* with a probability of 0.7 and back to *State 1* with a probability of 0.3. Choosing A_2 will lead you to *State 2* with a probability of 0.4 and back to *State 1* with a probability of 0.6. Clearly an agent wishing to transition to *State 2* is far more likely should they choose action A_1 . The key principle of MDPs is that state transitions are nondeterministic, meaning that choosing an action might not transition you to your intended next state. The overall objective of the learning agent is to optimize its value function. The agent makes decisions based on its value estimates of states and actions. $V^\pi(s)$ is called the state-value function. It is the expected value of being in a particular state under policy π .

$$V^\pi(s, a) = E_\pi\left\{\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} | s_t = s, a_t = a\right\} \quad (2)$$

The goal of the learning agent is to maximize its returns in the long run often forgoing short term gains in place of long term benefits. By introducing a discount factor γ , ($0 < \gamma < 1$), an agents degree of myopia can be controlled. A value close to 1 for γ assigns a greater weight to future rewards, while a value close to 0 considers only the most recent rewards.

B. Genetic Algorithms

Genetic algorithms are stochastic search and optimization techniques based on evolution. In their simplest form, a set of possible solutions to a particular problem are evaluated

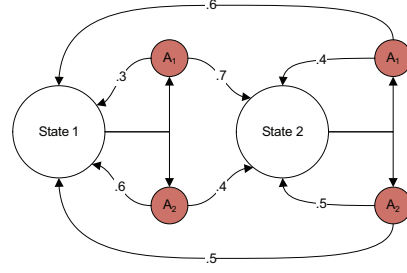


Fig. 1. Markov Decision Process with two states

in an iterative manner. From the fittest of these solutions, the next generation is created and the evaluation process begins once more. A solutions suitability to its environment is determined using a fitness function. By iterating through successive generations good approximate solutions can be found for the given environment.

In creating feasible workflow schedules, a task to resource mapping $\{T_i, C_j\}$ is represented as a single gene in the chromosome. A valid chromosome contains a sequence of genes, mapping every task in the workflow to a corresponding resource. The order of the genes represents the schedule execution order on the chosen resources. A feasible solution to the scheduling problem must maintain the precedence constraints between the tasks specified in the workflow.

The genetic algorithm is outlined by Algorithm 1.

Algorithm 1 Genetic Algorithm

Initialise population of feasible schedules

repeat

Evaluate chromosome fitness

Rank chromosomes according to overall population fitness

Select parents using roulette wheel selection

if (random > crossoverRate) **then**

 Apply single point crossover

end if

if (random > mutationRate) **then**

 Apply mutation

end if

Create next generation

until end

III. WORKFLOW ARCHITECTURE

Figure 2 details the architecture of the cloud workflow management system. Workflow execution is a multi stage process consisting of a number of steps leading up to and after workflow execution. Firstly users submit their application workflows through the user interface, along with their non-functional QoS constraints. The next stage of the process involves performance estimation on the submitted workflow. Performance estimation calculates the average execution time

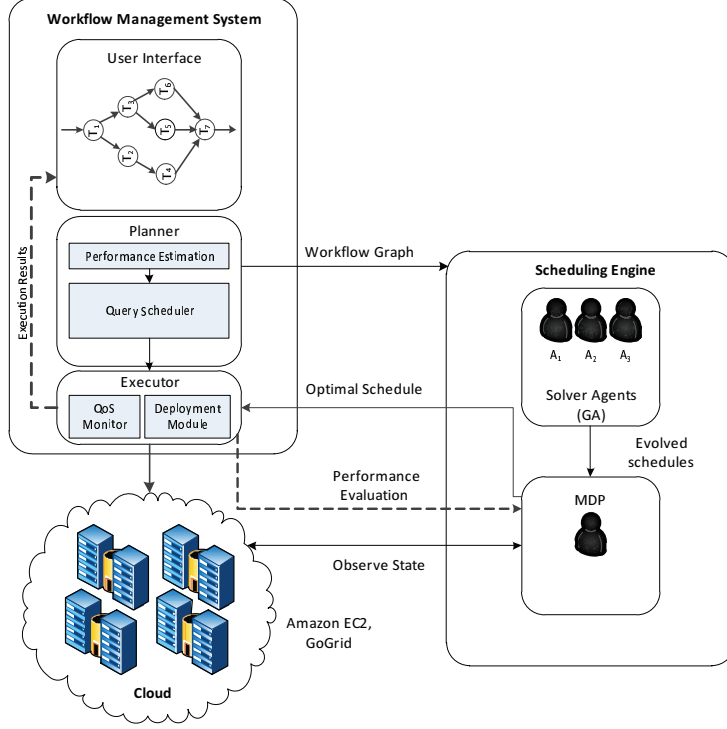


Fig. 2. Workflow Management System Architecture

that a given task will take to run on a specified resource. In cloud environments resources can be virtualised based on CPU, memory and storage. As the cost of cloud services is inversely proportional to their processing capabilities, once one knows the average execution time for a specific resource configuration, one can infer across all resources. Techniques such as analytical modelling [19] or historical data [20], can be used to determine average execution times. One of the principle advantages of scheduling with our learning approach, is that we are less reliant than other approaches on the performance of these techniques. For instance in the case of historical data, workflow applications such as scientific workflows are unlikely to be repeating sufficiently to allow for accurate estimates. Coupled with this the dynamic nature of the virtualised resources will also add to the uncertainty surrounding the estimates of task processing times. With our approach depending on environmental state, the learning agent will bias its decisions towards solutions less reliant on these estimates as they will have proven more successful in the past. These values are used by the solvers to estimate the makespan of the workflow. A number of solvers with ranging configurations are instantiated to produce schedules of varying cost and makespan. From these schedules an agent utilising a MDP computes the optimal schedule based on the current state of the cloud environment. The scheduling plan is then executed on the cloud via the Executor module. This module queues the tasks from the schedule with the associated resources. It also

monitors the tasks and returns the results to the user interface. Once a processing schedule has completed on the cloud, the QoS monitor returns the actual cost and makespan incurred by the schedule. This is then returned to the MDP agent who updates the transition probabilities accordingly.

IV. CLOUD WORKFLOW SCHEDULING

This section describes our algorithmic approaches to scheduling workflow applications in the cloud. It also outlines Directed Acyclic Graphs which can be neatly used to represent the execution order and data dependencies between tasks in a workflow.

A. Directed Acyclic Graph

Workflow applications can generally be modelled as a Directed Acyclic Graph (DAG) $G = \{V, E\}$. The set of vertices $V = \{T_1, \dots, T_n\}$ denotes each individual task in the workflow. E the set of directed edges represents precedence constraints between task nodes. A directed edge $E_{i,j}$ states that T_i is the parent task of T_j . Child tasks can only be executed once all the parent tasks have completed. $F_{i,j}$ denotes a data dependency between tasks T_i and T_j . We have a finite set of compute services $C = \{C_1, \dots, C_n\}$ which are capable of executing T_i . The execution of task T_i on a compute service C_j incurs a cost. This cost is inversely proportional to time taken to process it, where the greater the expense, the faster the resource. We also have a set of storage sites $S = \{S_1, \dots, S_n\}$. The cost of data transfer per unit of data between C_i and S_j is

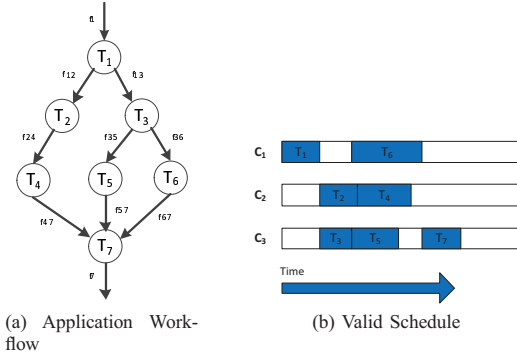


Fig. 3. Sample workflow application with feasible schedule

fixed and known in advance. The total data costs incurred by tasks executing on a given compute service C_i is $D_{total}(C)_i$. This includes all data costs between tasks executing on service C_i and those that are not. The overall makespan M_{total} of the workflow application is defined as the latest finished time across all virtual machines. Let M_{C_i} be equal to the total makespan of compute service C_i . Then the total makespan for the entire workflow is

$$M_{total} = \max(C_i) \forall i \in C \quad (3)$$

Data transfer costs between compute services can be calculated by the file size of the output of the parent task. Generally for two tasks executing on the same resource there is no data transfer cost. The total processing costs P for a given compute service is P_{total} . Total task execution costs are

$$Ex_{total} = D_{total}(C)_i + P_{total}(C)_i \forall i \in C \quad (4)$$

Figure 3a depicts a sample workflow DAG containing 7 task nodes. The edges between tasks indicate the input files and output files between task nodes. Figure 3b depicts a valid schedule for the adjacent workflow. Each individual task is assigned to a specific resource and can only execute once its parent tasks have completed. The objective of the cloud workflow scheduler is to find a mapping from tasks to compute services that minimises the cost of execution and the overall makespan of the workflow.

$$\text{Min}(Ex_{total} + M_{total}) \quad (5)$$

In general, workflow scheduling onto distributed resources is an NP-Complete problem [21]. To generate near optimal solutions we apply a metaheuristic algorithm to solve the workflow scheduling problem.

B. Scheduling Using a Genetic Algorithm

In our scheduling architecture each solver agent employs a Genetic Algorithm in order to create valid schedules for the workflow scheduling task. In our approach evolution occurs over the entire population, with offspring from the fittest solutions replacing only the weakest schedules in the population. Firstly an initial population of feasible schedules

is created. Then each solution is assigned a fitness value based on Equation 5. A solutions total cost and execution time infers its fitness in the environment. For selection purposes we normalise fitness according to the following equation.

$$F(i) = \frac{f_{max} - f_i}{f_{max} - f_{min}} + \mu \quad (6)$$

The fitness $F(i)$ is the normalised fitness of individual i . f_{max} and f_{min} are the maximum and minimum fitness values in the population, with $\mu = 0.3$. Normalising the fitness in this manner ensures that the highest probability is apportioned to the fittest solutions during the selection process. Individuals are selected for reproduction using roulette wheel selection, based on their fitness. Roulette wheel selection involves ranking chromosomes in terms of their fitness and probabilistically selecting them. The selection process is weighted in favour of chromosomes possessing a higher fitness. To ensure that agents, already optimal for their environment are not lost in the evolutionary process elitism is applied. Elitism involves the selection of a certain percentage of the fittest chromosomes and moving them straight into the next generation, avoiding the normal selection process. In creating offspring for the next generation, the selection of two parents is required. Each pairing results in the reproduction of two offspring. Figure 4 shows the crossover of two parents and subsequent production of two offspring. Crossover involves taking certain aspects/traits of both parents' chromosomes and creating a new chromosome. There are a number of ways to achieve this including, single point crossover, two point crossover and uniform crossover. Our crossover function employs single point crossover, where a point in the bit string is randomly selected, at which crossover is applied. Crossover generally has a high probability of occurrence. Figure 4 shows two valid schedules combining to create two valid offspring. However crossover can also result in the creation of schedules which are invalid. To convert an invalid schedule into a feasible solution we apply an adjustment after crossover has completed [22]. Mutation involves randomly altering the bit string altering aspects of a chromosome. Mutation occurs on the assigned service in a given task-service mapping as seen in Figure 5.

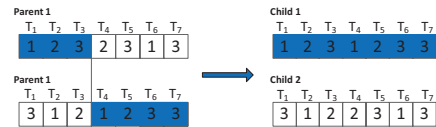


Fig. 4. Single point crossover

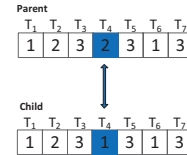


Fig. 5. Mutation

Once the required number of offspring have been created they form the next generation and the process begins once more. The algorithm terminates when the stopping condition is met.

C. Learning Agent - Choosing Schedules Using a MDP

The key novelty proposed by our architecture is the modelling of the scheduling problem as a finite MDP. This allows a learning agent to solve this MDP and reason about the possible schedules generated by the solver agents which was outlined in the previous sections. In solving a MDP two algorithms, value iteration or policy iteration algorithms from dynamic programming can be used. In this work we choose the value iteration algorithm to calculate the state value function. Approximations of $V^\pi(s)$ which are indicative as to the benefit of being state s under policy π are calculated after each time interval. Actions are chosen based on π the policy being followed. The policy denotes the optimal mapping from states to actions.

Algorithm 2 outlines the value iteration algorithm used to compute the value function estimates. Firstly all states are initialised arbitrarily. Our state signal is comprised of three variables, which are **Time**, **Resource Load** and **Execution Result**. The first variable *Time* is necessary in order to analyse the effects of temporal fluctuations on the executing workflows. At peak times (e.g between 7-9pm) in any given region there could be significant effect on the workflow execution performance. By including time into the state space the learning agent can reason about the possible effects of this. The second variable is the overall *Resource Load* of the virtual machines which the workflow system has instantiated. It coarsely evaluates its own resources by summing up all workflow tasks currently executing and divides this by the number of virtual machines. Resource Load is split into three distinctions, *light*, *moderate* and *heavy*. This coarse interpretation is surprisingly effective at enabling a good estimate of the system load. The final variable *Execution Result* defines whether or not a given selection was successful in executing the workflow. This information is returned via the QoS monitor. There are many other possible metrics that could, if added to the state signal, give even greater insights (CPU utilisation, disk reads, bandwidth). However for the sake of simplicity and to enable a manageable state set we focus on these three. If we assume time to transition on an hourly basis (0-23), combined with the three distinctions for resource load and the two terminal states (Success and Failure), then the total state space consists of 144 states. The set of all possible actions in the MDP are defined by the number of solver agents present in the system. The MDP may choose a schedule from amongst this set of solvers. The rewards are defined as the cost incurred as a result of the execution of a given workflow. If a schedule resulted in the violation of the specified deadline, then additional penalties are added to the reward. The goal of the MDP is to ensure the successful completion of the workflow execution within the budget and deadline constraints.

Algorithm 2 Value Iteration

Initialize V arbitrarily

repeat

$\Delta \leftarrow 0$

For each $s \in S$

$v \leftarrow V(s)$

$V(s) \leftarrow \min_a \sum_{s'} P_{s,s'}^a [R_{s,s'}^a + \gamma V(s')]$

$\Delta \leftarrow \max(\Delta, |v - V(s)|)$

until $\Delta < \Theta$

Output policy π

For all $s \in S$

$\pi(s) = \operatorname{argmin}_a \sum_{s'} P_{s,s'}^a [R_{s,s'}^a + \gamma V(s')]$

1) *Bayesian Model Learning*: One important consideration is that approximating the true values of $V(s)$ is largely dependent on the number of times state s is visited. In fact the true values will only converge on the limit. Value iteration generally only guarantees asymptotic convergence. We adopt a Bayesian model learning approach similar to what was proposed by Doshi et al. [23]. The MDP updates its estimates of the transition probabilities based on experience. When submitting an evolved schedule for execution, the MDP agent needs to have an estimate of the certainty with which the executed schedule will result in the desired outcome. These estimates are a measure of the uncertainty of the underlying processes and are manifested in the transition function P . The optimal policy π^* is naturally dependent on these estimates. In order to achieve better policies one needs to update approximate transition probabilities based on experience. In a MDP, choosing an action a in state s results in a state transition from s to s' . The experience we have of this transition is a measure of how certain our actions will result in the predicted outcome. To update these probabilities we initialise an experience counter *Expc* to 1 and each time we observe a particular state transition the experience associated with it is incremented by one. The updated transition probability is given by Equation

$$P'(s = s'|a, s = s) = \frac{P(s = s'|a, s = s) \times \text{Expc} + 1}{\text{Expc}'} \quad (7)$$

where Expc' is the incremented counter. In order to make the probabilities sum up to 1 we update the remainder of the probabilities using the following equation

$$P'(s = s'|a, s = s) = \frac{P(s = s'|a, s = s) \times \text{Expc}}{\text{Expc}'} \quad (8)$$

Doshi's work proposes an Execute&Learn [23] algorithm to learning the true transition probabilities, where model learning and workflow execution are interleaved. Actual rewards observed by the agent are used to update the probabilities. This allows for the continuous updating the MDP's estimates enabling closer approximation of the true transition probabilities.

TABLE I
COMMUNICATION COSTS (\$CENTS) BETWEEN REGIONS

	C_1	C_2	C_3	C_4
C_1	0	0.15	0.19	0.20
C_2	0.15	0	0.19	0.20
C_3	0.15	0.15	0	0.20
C_4	0.15	0.15	0.19	0

V. INITIAL RESULTS

This section presents our initial findings with our cloud workflow system. We carry out an empirical analysis of our learning architecture highlighting the potential of the proposed system.

A. Experimental Setup

To evaluate our MDP approach we used Cloudsim, a simulator developed by the CLOUDS group at the University of Melbourne [24]. Cloudsim is cloud simulation tool which allows you to model typical cloud scenarios on a single processor. It facilitates the creation of data centres and the deployment of multiple hosts and virtual machines of varying capacities. You can also specify detailed pricing policies with regard to virtual machines and data transmission costs. In our experiments we simulate four separate data centers in four geographic locations. We place one in the US, one in Europe, one in Hong Kong/Singapore and one in Japan. Data transfers between compute sites are serviced by a content delivery network (CDN) such as Amazon CloudFront². Each data centre supports a single host capable of instantiating multiple virtual machines. Table I depicts Amazons current costs per unit data transfer between the four geographic regions at the time of writing.

For each experiment we use three separate solvers to produce different schedules for the workflow scheduling problem. Each solver is configured with the following values for population size, elitism, cross-over and mutation 50, 5%, 85% and 3%. Each solver is given its own unique configuration, to ensure a diversity of schedule solutions. Importantly each solver applies a weighting to the average values associated with task processing times. This value is in the range of $[0, 1]$. A value close 1 associates high degree of confidence in the predicted task execution times, whereas a value close to 0 indicates low confidence. These parameters ensure a diverse range of solutions with respect to cost and makespan are returned. Therefore the goal of the learning agent is to identify the solver which is most likely to achieve the best performance given the current environmental state.

B. Increasing Data Size

Our first experiment analyses the cost savings of the MDP approach as the amount of data increases across the system. Figure 6 shows the growth in costs as the data sizes increase in the range of [64mb - 2048mb]. The BRS approach performs

badly in this regard as it does not take the costs of data transfer between resources into account. It myopically assigns a task to a resource based on the most expensive resource. As the total amount of data across the system increases, the BRS approach incurs additional costs as tasks are executed on different resources in different regions, at greater expense. The costs for executing the MDP schedule over BRS increases at a much more slower rate. The learning agent chooses from amongst the various solvers enabling a better performance. In conjunction with this the schedules evolved with the Genetic Algorithm perform much better individually than the BRS, as the schedule it evolves makes considerations over the entire workflow and does not simply choose the best resource.

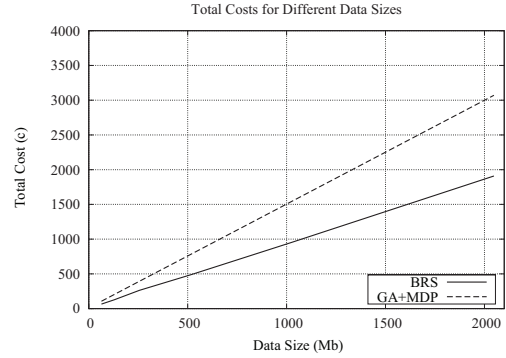


Fig. 6. Overall cost as data increases

C. Load Variance

To simulate a variable work load on the system we fix the arrival rate of workflows to be processed. We vary the number of tasks in each workflow application from 20 to 50. The number of compute services is constrained from 5 – 15. By controlling the arrival rate of workflow applications and the number of tasks at any time we can vary execution load across the virtual machines over time. Figure 7 shows the performance of the MDP as the load varies. We compare the performance of the MDP approach against the BRS approach and a GA. The genetic algorithm parameter configuration was determined experimentally through identifying the best settings to encourage the best solutions across the widest range of system loads. Initially as the MDP assigns equal probabilities to all actions and doesn't perform as well as the GA. It spends time exploring suboptimal actions and get penalised as a result. As the MDP explores the state space, it outperforms both the GA and BRS. Since the data size used for the experiment was fixed, BRS does not incur any significant loss of costs due to load variances.

VI. CONCLUSIONS & FUTURE WORK

We presented an initial architecture for a cloud workflow scheduler and some preliminary results. We have presented a continuous state action space formalism to schedule workflows on a cloud computing environment. Our results show that our

²<http://aws.amazon.com/cloudfront/>

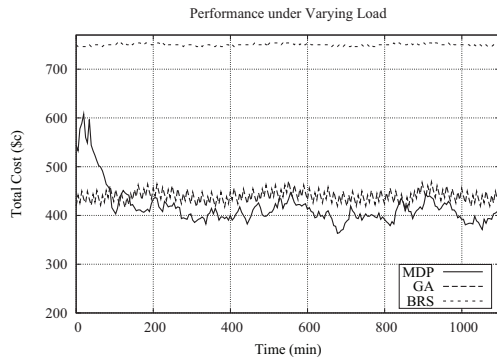


Fig. 7. Performance of MDP as system load varies

MDP agent can optimally choose from a set of evolved workflow schedules. This is achieved in spite of an environment that has varying computational loads and data sizes. We presented a Bayesian model learning approach capable of quickly learning approximate transition probabilities with no prior knowledge. It is more challenging to learn optimal value functions for states that are less frequently visited. These states are most likely to occur when system load approaches its maximum or minimum values. Therefore, MDP performance is heavily dependant on it gaining a greater experience across the widest range of possible states. This means the performance of the MDP will suffer unless it can gain greater experience. In future work we hope to investigate a solution based on developing policies offline using sample data from previous executions. This will facilitate the development of better policies online and produce more optimal solutions. Furthermore, we hope to introduce a range of optimisation algorithms that can provide a broader range of schedules and thereby improve overall system performance. Finally we hope to perform a qualitative evaluation of each workflow application and incorporate this into the state space.

ACKNOWLEDGMENT

The authors would like to gratefully acknowledge the continued support of Science Foundation Ireland.

REFERENCES

- [1] Montage: An astronomical image engine, <http://montage.ipac.caltech.edu>.
- [2] Illumina: , <http://www.illumina.com>.
- [3] C. Vecchiola, S. P. and R. Buyya, "High-performance cloud computing: A view of scientific applications."
- [4] G. Juve and E. Deelman, "Scientific workflows and clouds," *ACM Crossroads*, vol. 16, no. 3, pp. 14–18, 2010.
- [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. Piscataway, NJ, USA: IEEE Press, 2008, pp. 50:1–50:12. [Online]. Available: <http://portal.acm.org/citation.cfm?id=1413370.1413421>
- [6] J. Yu, R. Buyya, and C. K. Tham, "Cost-based scheduling of scientific workflow application on utility grids," in *Proceedings of the First International Conference on e-Science and Grid Computing*. Washington, DC, USA: IEEE Computer Society, 2005, pp. 140–147. [Online]. Available: <http://portal.acm.org/citation.cfm?id=1107836.1107867>
- [7] J. Yu and R. Buyya, "Scheduling scientific workflow applications with deadline and budget constraints using genetic algorithms," in *Scientific Programming*. IOS Press, 2006, pp. 217–230.
- [8] J. Yu, M. Kirley, and R. Buyya, "Multi-objective planning for workflow execution on grids," in *Proceedings of the 8th IEEE/ACM International Conference on Grid Computing*, ser. GRID '07. Washington, DC, USA: IEEE Computer Society, 2007, pp. 10–17. [Online]. Available: <http://dx.doi.org/10.1109/GRID.2007.4354110>
- [9] E. Deelman, G. Singh, M. hui Su, J. Blythe, A. Gil, C. Kesselman, G. Mehta, K. Vahi, G. B. Berriman, J. Good, A. Laity, J. C. Jacob, and D. S. Katz, "Pegasus: a framework for mapping complex scientific workflows onto distributed systems," *Scientific Programming Journal*, vol. 13, pp. 219–237, 2005.
- [10] B. Ludscher, I. Altintas, C. Berkley, D. Higgins, E. Jaeger, M. Jones, E. A. Lee, J. Tao, and Y. Zhao, "Scientific workflow management and the kepler system," in *Concurr. Comput. : Pract. Exper.*, 2005, p. 2006.
- [11] T. Oinn, M. Addis, J. Ferris, D. Marvin, T. Carver, M. R. Pocock, and A. Wipat, "Taverna: A tool for the composition and enactment of bioinformatics workflows," *Bioinformatics*, vol. 20, p. 2004, 2004.
- [12] S. P. L. Wu, S. M. Guru, and R. Buyya, "A particle swarm optimization-based heuristic for scheduling workflow applications in cloud computing environments."
- [13] Z. Wu, X. Liu, Z. Ni, D. Yuan, and Y. Yang, "A market-oriented hierarchical scheduling strategy in cloud workflow systems," *The Journal of Supercomputing*, pp. 1–38, 2011, 10.1007/s11227-011-0578-4. [Online]. Available: <http://dx.doi.org/10.1007/s11227-011-0578-4>
- [14] D. Vengerov, "A reinforcement learning approach to dynamic resource allocation," *Eng. Appl. Artif. Intell.*, vol. 20, no. 3, pp. 383–390, 2007.
- [15] G. Tesauro, N. K. Jong, R. Das, and M. N. Bannani, "On the use of hybrid reinforcement learning for autonomic resource allocation," *Cluster Computing*, vol. 10, no. 3, pp. 287–299, 2007.
- [16] D. Gmach, J. Rolia, L. Cherkasova, and A. Kemper, "Workload analysis and demand prediction of enterprise data center applications," in *Workload Characterization, 2007. IISWC 2007. IEEE 10th International Symposium on*, Sept. 2007, pp. 171–180.
- [17] J. Perez, C. Germain-Renaud, B. Kégl, and C. Loomis, "Grid differentiated services: A reinforcement learning approach," in *CCGRID '08: Proceedings of the 2008 Eighth IEEE International Symposium on Cluster Computing and the Grid*. Washington, DC, USA: IEEE Computer Society, 2008, pp. 287–294.
- [18] R. S. Sutton and A. G. Barto, "Reinforcement learning :an introduction," 1998.
- [19] G. R. Nudd, D. J. Kerbyson, E. Papaefstathiou, S. C. Perry, J. S. Harper, and D. V. Wilcox, "Pace - a toolset for the performance prediction of parallel and distributed systems," 1999.
- [20] W. Smith, I. Foster, and V. Taylor, "Predicting application run times using historical information," 1997.
- [21] J. D. Ullman, "Np-complete scheduling problems," *J. Comput. Syst. Sci.*, vol. 10, pp. 384–393, June 1975. [Online]. Available: [http://dx.doi.org/10.1016/S0022-0000\(75\)80008-0](http://dx.doi.org/10.1016/S0022-0000(75)80008-0)
- [22] J. Oh and C. Wu, "Genetic-algorithm-based real-time task scheduling with multiple goals," *J. Syst. Softw.*, vol. 71, pp. 245–258, May 2004. [Online]. Available: <http://portal.acm.org/citation.cfm?id=993017.993022>
- [23] P. Doshi, "Dynamic workflow composition using markov decision processes," *International Journal of Web Services Research*, vol. 2, pp. 1–17, 2005.
- [24] R. N. Calheiros, R. Ranjan, C. A. F. D. Rose, and R. Buyya, "Cloudsim: A novel framework for modeling and simulation of cloud computing infrastructures and services," *CoRR*, vol. abs/0903.2525, 2009.