

# 하이브리드 클라우드 환경에서의 SLA기반 오토-스케일링 기법

강혜정 고정인<sup>○</sup> 김윤희<sup>☆</sup> 조혜영<sup>†</sup>

숙명여자대학교 컴퓨터과학부, 한국과학기술정보연구원<sup>†</sup>

{hjkang, jungin, yulan}@sookmyung.ac.kr, chohy@kisti.re.kr<sup>†</sup>

## SLA-Aware Auto-Scaling in Hybrid Cloud Environment

Hyejeong Kang Jungin Koh<sup>○</sup> Yoonhee Kim<sup>☆</sup> Hyeyoung Cho<sup>†</sup>

Dept. of Computer Science, Sookmyung Women's University

Korea Institute of Science and Technology Information<sup>†</sup>

### 요 약

최근 계산 과학 분야에서 대용량 컴퓨팅 자원의 요구를 능동적으로 만족시키기 위해, 자원을 필요한 만큼 빌려 쓸 수 있는 클라우드 기술을 적용하여 과학 클라우드(Science Cloud)를 구축하는 연구가 세계 여러 곳에서 진행되고 있다. 클라우드 컴퓨팅 환경을 이용하는 데 있어 가장 중요한 이슈 중 하나는 자원을 사용함에 있어 실제 필요한 양만큼 할당하여 사용하는 것이다. 현대 응용의 가변적인 작업 부하에 맞춰 필요한 만큼만 자원을 제공하기 위해 기존 상용 클라우드들에서는 규칙 기반의 메커니즘을 이용하여 자원 할당의 자동화를 시도하고 있다. 그러나 대부분의 오토 스케일링 기법들은 단순히 자원의 성능 지표만을 지원할 뿐 응용의 데드라인이나 비용과 같은 Service Level Agreement(SLA)를 고려하지 않는다. 본 논문에서는 사설/상용 클라우드로 이루어진 하이브리드 클라우드 환경에서 가변적인 자원 요구에 따라 자원을 할당하고 SLA 위반을 최소화하는 오토-스케일링 기법을 제안한다.

## 1. 서 론

최근 계산 과학 분야에서 대용량 컴퓨팅 자원의 요구를 능동적으로 만족시키기 위해, 자원을 필요한 만큼 빌려 쓸 수 있는 클라우드 기술을 적용하여 과학 클라우드(Science Cloud)를 구축하는 연구가 세계 여러 곳에서 진행되고 있다. 빅데이터 처리 응용이나 계산 과학 응용(예, 과학 워크플로우)은 고성능 컴퓨팅(HPC: High Performance Computing) 또는 다처리 컴퓨팅(HTC: High Throughput Computing)을 요구하는 대규모 사이즈 응용(MTC: Many Task Computing)으로 정의할 수 있다. 이런 응용은 고성능의 자원을 장시간 활용해야 하며 안정적이고 지속적인 자원의 공급이 응용 실행의 필수이다. 이는 엑사스케일(Exascale) 컴퓨팅, 페타스케일(Petascale) 컴퓨팅의 연구로 이어지고 있으며 이런 환경에서 성공적인 작업 수행을 위한 실행관리, 자원 관리 등을 제공하는 계산 문제 풀이 환경(Computational Problem Solving Environment)에 대한 연구가 더욱 절실하다.

이런 문제 풀이 환경에 필요한 클라우드 자원을 효율적으로 활용하기 위한 방법으로 온-디맨드 자원 가상화 특성을 활용한 오토-스케일링 기법을 고려할 수 있다. 오토-스케일링 기법은 응용의 수행 중에 작업 수행을 위한 가상 자원의 수를 줄이거나 늘리는 방법으로 AWS

(Amazon Web Service)[1]의 “Auto-scaling” 서비스가 대표적인 예이다. AWS나 Scalr[2]와 같은 오토-스케일링의 서비스들은 사용자가 정의한 규칙(예, CPU 이용량이 60% 이상일 때, 가상머신 2개 추가 할당)을 기반으로 자원의 규모를 결정한다. 단순한 방법으로 오토-스케일링을 가능케 하는 장점이 있지만 현대 응용의 가변적인 워크로드로 인한 동적인 자원 요구 패턴을 만족하지 못하는 문제점이 있다. 특히, 응용의 자원 요구가 충족되지 않을 경우 응용의 데드라인 위반, 수행 실패의 증가와 같은 SLA(Service Level Agreement) 위반으로 이어질 수 있는 심각한 문제점이 존재한다.

본 논문에서는 하이브리드 클라우드 컴퓨팅 환경에서 자원을 효율적으로 활용하기 위한 오토-스케일링 기법을 제안한다. 제안하는 오토-스케일링 기법은 데드라인, 비용/성능 정책과 같은 SLA를 충족한다. 또한 계산 집약, I/O 집약, 데이터 집약 등 응용의 특성과 함께 상용 클라우드의 인스턴스 타입을 고려함으로써 비용 효율적인 오토-스케일링이 가능하다.

본 논문의 구성은 다음과 같다. 1장의 서론에 이어 2장에서는 관련 연구들을 살펴보고, 3장에서는 본 논문에서 제안하는 SLA기반 오토-스케일링 기법을 적용한 프레임워크 구조에 대하여 설명한다. 4장에서는 SLA 기반 오토-스케일링 기법에 대해 설명하고 마지막으로 5장에서 결론을 맺는다.

## 2. 관련 연구

클라우드 컴퓨팅은 가상화 기술을 통해 자원의 무한

<sup>☆</sup> 교신저자

“본 연구는 한국과학기술정보연구원의 국가슈퍼컴퓨팅 서비스 기반 강화의 연구결과로 수행되었음” (K-13-L01-C01)

한 제공이 가능하고, 자원 규모의 확장과 축소가 용이하다. 이로 인해 효율적인 자원 관리 방법으로 오토-스케일링 기법에 대한 연구가 활발히 이루어지고 있다.

대표적인 오토-스케일링 기술로는 AWS[1]의 “Auto-scaling” 서비스가 존재함을 앞서 언급하였다. 이 서비스는 사용자가 정의한 규칙을 통해 자원의 규모를 확장하고 축소하는 규칙 기반의 자원 할당 자동화 기법으로 CPU 사용량, 디스크 사용량과 같은 하드웨어의 성능과 관련된 수치를 기준으로 한다. 기준이 되는 성능치에 상한 값과 하한 값을 정하고 그 값을 기준으로 미리 정해진 수만큼 가상머신을 추가하거나 제거하여 자원의 규모를 확장하거나 축소시키는 방법이다. 비슷한 서비스로 Windows Azure[3]에서 사용되는 Paraleap[4]과 Scalr[2], RightScale[5]이 존재한다. 규칙 기반의 스케일링 기법은 비교적 단순한 방법으로 동적인 자원 할당이 가능하지만 복잡한 수행 패턴을 갖는 작업에 대해서 실제로 필요한 만큼 가상머신이 추가/제거되지 않을 수 있다. 이러한 문제는 작업 수행 데드라인이나 자원 사용 비용과 같은 SLA를 갖는 작업들에서는 SLA 위반이라는 큰 문제를 야기할 수 있다.

작업 수행 데드라인이나 자원 사용 비용 등을 고려한 오토-스케일링 기법 연구로는 [6], [7]이 존재한다. [6]에서는 사용자가 요청한 작업의 데드라인 내에서 작업 수행이 가능하도록 오토-스케일링을 수행한다. 이 연구에서는 작업의 데드라인뿐만 아니라 자원 할당에 따른 비용까지 고려하였다. 그러나 상용 클라우드를 단독으로 사용하였을 경우의 자원 할당만을 고려하였다는 제한사항이 존재한다. [7]은 자원 사용 비용을 최소화하기 위해 물리머신(Physical Machine)의 사용량과 자원 사용료 간에 트레이드-오프가 최적인 물리머신 사용량을 찾아 해당하는 만큼 가상머신을 스케일링하는 기법을 제안한 연구이다. 가상머신의 개수를 추가/제거하여 자원 규모를 확장하고 축소하는 수평적(Horizontal) 스케일링과 가상머신의 하드웨어 사양을 조절하는 수직적(Vertical) 스케일링 기법을 혼합 활용한다. 그러나 이 연구에서도 수평적 스케일링 시에는 고정된 개수의 가상머신만을 추가/제거하여 복잡한 수행 패턴의 작업에 대하여 동적인 자원 할당이 가능하다고 하기엔 부족하다.

본 논문에서는 하이브리드 클라우드 환경에서 데드라인과 자원 사용 비용과 같은 SLA를 고려하여 자원 할당이 가능한 오토-스케일링 기법을 제안한다. 지속적인 모니터링을 통하여 자원 할당이 이루어지며 자원 할당 시에 필요한 가상머신의 개수를 계산하므로 그때그때 적절한 양의 가상머신을 추가/제거할 수 있다.

### 3. 오토-스케일링 프레임워크 구조

본 논문에서 제안한 하이브리드 클라우드 환경에서 오토-스케일링 기법을 지원하는 프레임워크의 구조는 그림 1과 같다. ‘Scaling Decision Service’는 핵심 서비

스 모듈로서 ‘Job Execution Monitoring Service’를 통해 작업 실행 정보를 얻어 오토-스케일링에 이용한다. 최종 자원 추가/제거 여부를 ‘Dynamic Resource Mgmt. Service’에 전달하여 실제 가상머신을 추가/제거하도록 한다. ‘Job Mgmt. Service’는 ‘Queue’에 대기중인 작업들의 각 요구 수행시간 및 요구 자원 성능 정보를 관리하고 ‘Scaling Decision Service’에서 결정된 자원 할당 스케줄에 따라 작업을 자원에 제출하는 서비스를 제공한다. ‘Queue’에는 사용자가 제출한 작업이 자원 할당이 완료될 때까지 머물게 된다. ‘VM Catalog’는 사설 클라우드에 존재하는 가상머신의 성능과 상용 클라우드에서 제공하는 인스턴스 타입 및 가격에 대한 정보를 담고 있다.

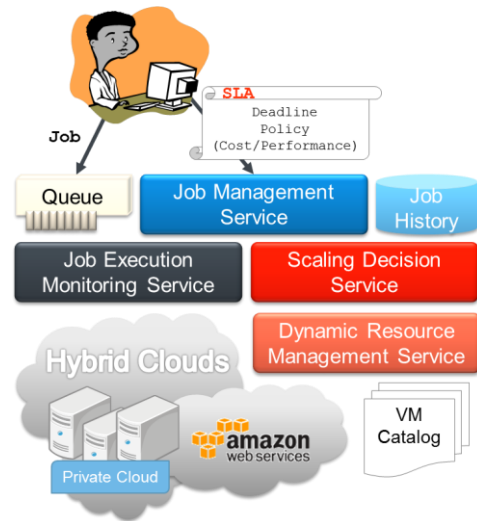


그림 1 오토-스케일링 프레임워크 구조

### 4. SLA 기반 오토-스케일링 기법

본 논문에서는 하이브리드 클라우드 환경에서의 오토-스케일링 기법을 제안한다. 작업 부하가 동적으로 변하는 현대 응용의 특성상 작업 수행에 대한 예측 및 자원 사용량에 관한 예측이 힘들기 때문에 주기적인 모니터링을 통해 시시각각 변하는 자원 요구를 충족시켜줄 필요가 있다. 자원 요구를 만족시키는 것과 동시에 데드라인, 자원 사용 비용과 같은 SLA로 정의되는 사용자 요구사항을 충족시키는 것 또한 중요하다. 따라서 제안하는 기법에서는 비용/성능 두 가지 정책과 함께 데드라인을 SLA로 정의하고 데드라인을 기준으로 사용자가 선택한 정책기반의 자원 할당과 오토-스케일링을 수행한다. 비용 정책은 사용자가 상용 클라우드 자원을 사용함으로써 발생하는 자원 이용료를 고려하여 합리적인 비용으로 데드라인 내에서 작업을 수행하도록 자원을 할당하는 정책이고, 성능 정책은 최대한 빨리 작업을 완료하도록 자원을 할당하는 정책이다. 성능 정책을 선택할 경우 데드라인은 최상 경로(Critical Path) 이상의 값을 가져야 하며, 최소 성능 요구도 SLA로 함께 포함된다. 즉 자원 할당 시에 사용자가 요구한 최소 성능 이상의 사양을 갖는 자원을 할당하는 것이다. 자원 요구

에 따라 자원 규모를 확장하는 스케일 아웃의 경우 정책에 따른 자원 할당 시에 이루어지며 자원 규모의 축소인 스케일 인의 경우 자원 할당이 끝난 후 일괄적으로 이루어지게 된다.

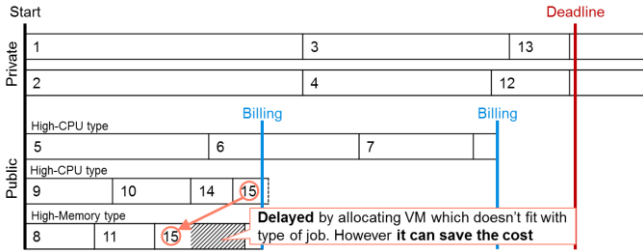


그림 2 최초 자원 할당 결과

그림 2는 제안하는 오토-스케일링 기법의 최초 자원 할당 결과를 보여준다. 작업들은 수행시간을 기준으로 내림차순으로 정렬한 후, 자원 이용료가 들지 않는 사실 클라우드 자원에 최대한 많은 작업을 할당한다. 이후 데드라인을 위반할 것이라 예상되면 상용 클라우드에서 자원을 구매하여 작업을 수행한다. 상용 클라우드의 자원에 작업을 할당할 경우, 이미 Running 중인 자원을 우선 고려하는데 이때 새로운 가상머신을 샀을 때의 가격, 인스턴스 타입을 고려하여 새로운 가상머신의 구입 여부를 결정한다. 인스턴스 타입을 고려함으로써 성능대비 비용 효율적인 자원을 사용하도록 하였다. 여기서 성능대비 비용 효율적이라 함은 계산 집약적인 작업을 수행한다고 했을 때, Amazon EC2의 인스턴스 타입을 기준으로 같은 사양의 CPU 성능을 갖는 일반 Standard 인스턴스와 High-CPU 인스턴스의 가격 비교 시, 후자의 가격이 싸기 때문이다. 그림에서 작업 15번의 경우 작업 타입에 맞는 가상머신에 작업을 그대로 할당하면 가격 책정(Billing) 시간을 초과하여 추가 금액이 발생하나, 다른 타입의 자원을 할당하여 비록 수행시간이 길어지더라도 가격 책정 시간 내에 수행되어 추가적인 비용 발생을 방지한 것을 볼 수 있다.

작업 수행시간으로 데드라인 위반(그림 2의 작업 12)이 일어날 것이라 판단되면 아직 대기 상태에 있는 작업들(작업 4, 7, 12, 13)에 대하여 자원 할당을 새롭게 수행한다. 작업 길이에 따라 4, 7, 12, 13의 순으로 자원 할당이 이루어진다. 사실 클라우드에 가용 자원이 존재하므로 우선 고려 되고 7, 12의 경우 사실 클라우드에서 수행할 경우 데드라인을 위반하므로 상용 클라우드 자원에 할당된다. 여기서 작업 7이 그대로 원래 할당되었던 자원에 할당될 경우 가격 책정 시간 초과로 추가 비용( $\alpha$ )이 발생한다. 이 때, 작업 12를 7 뒤에 할당하면 데드라인 위반이므로 새로운 가상머신을 구매( $\beta$ )하여 할당하여야 하고, 작업 12의 경우 작업 길이가 짧아 단위 가격 책정 시간 동안 유휴 상태가 길어진다. 그러므로 작업 7을 다시 새로운 가상머신에 할당하여 결과적으로 새로운 가상머신 구매 금액인  $\beta$ 만큼만 자원 사용 비용에 추가된다. 마지막으로 작업 13이 사실 클라우드에 할당되고 자원 할당이 끝난 후 유휴 자원이 존재할 경우 해당 자원을 Release하여 자원 규모를 축소한다.

## 5. 결론 및 향후 연구

본 논문에서는 하이브리드 클라우드 컴퓨팅 환경에서 효율적으로 자원을 관리하기 위한 오토-스케일링 기법을 제안하였다. 제안하는 오토-스케일링 기법을 통하여 데드라인 내에서 작업 수행에 필요한 자원을 필요할 때에 필요한 만큼만 구매하여 사용할 수 있으며, 이를 통해 사용자의 자원 사용 비용을 최소화할 수 있다. 또한 제공되는 비용과 성능 정책 중, 성능 정책 선택 시 비용에 대한 고려를 최소화하고 작업의 수행시간을 단축하여 시간적인 이득을 볼 수도 있다.

향후에는 데이터센터의 전력 소모량을 고려하여 동적인 자원 스케일링이 가능한 에코-스케일링 기법을 추가 연구할 예정이다.

## 6. 참고 문헌

- [1] Amazon Web Service, <http://aws.amazon.com/>
- [2] Scalr, <http://scalr.com/>
- [3] Windows Azure, <http://www.windowsazure.com/>
- [4] Paraleap, <https://www.paraleap.com/>
- [5] RightScale, <http://www.rightscale.com/>
- [6] Ming Mao and Marty Humphrey, "Auto-Scaling to Minimize Cost and Meet Application Deadlines in Cloud Workflows," Proceedings of 2011 International Conference for High Performance Computing, Networking, Storage and Analysis, Nov. 12-18, 2011.
- [7] Dutta, S., Gera, S., Verma, A., and Viswanathan, B. "SmartScale: Automatic Application Scaling in Enterprise Clouds," Proceedings of 2012 IEEE 5th International Conference on Cloud Computing (CLOUD), pp. 221-228, June. 2012.

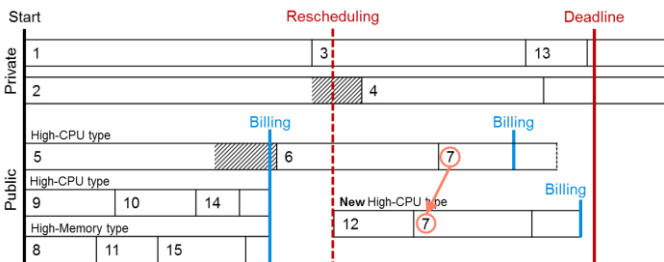


그림 3 모니터링 및 오토-스케일링 결과

그림 3은 지속적인 모니터링을 통해 응용 수행의 데드라인 위반 여부를 확인하고 자원 요구량 변화에 따른 오토-스케일링을 수행한 결과를 보여준다. 현재 자원에서 수행되는 작업들에 대해서 예상 수행 시작 시간과 실제 수행 시작 시간을 비교하여, 이전 자원 할당 시에 예상한 것보다 작업의 시작이 지체되었을 경우(작업 4, 6) 앞서 수행된 작업의 지연(작업 2, 5)으로 판단하고 지연 정도를 계산한다. 만약 이 지연으로 인해 늘어난