

# 클라우드 컴퓨팅에서의 대용량 데이터 처리 모델에 관한 조사

진아연\*, 박영호\*

\*숙명여자대학교 멀티미디어학과

e-mail: einbilden@naver.com

yhpark@sookmyung.ac.kr

## A Survey on Massive Data Processing Model in Cloud Computing

Ah-Yeon Jin\*, Young-Ho Park\*

\*Dept of Multimedia Science, Sookmyung Women's University

### 요 약

클라우드 컴퓨팅은 세계적인 시장조사기관인 가트너사의 10대전략기술에서 2년 연속 1위를 할 정도로 많은 각광을 받고 있다. 클라우드 컴퓨팅이란 인터넷 기술을 활용하여 가상화된 컴퓨팅 자원을 서비스로 제공하는 것으로, 사용자는 IT자원을 필요한 만큼 빌려서 사용하고 사용한 만큼 비용을 지불하는 컴퓨팅을 지칭한다. 이러한 클라우드 컴퓨팅 상에서 폭발적으로 증가하고 있는 데이터를 효율적으로 병렬 처리할 수 있는 방법에 대하여 많은 연구가 활발히 이루어지고 있다. 이러한 대용량 데이터 처리를 위한 대표적인 모델에는 MapReduce와 Dryad가 있으며, 서로간에 많은 공통점이 있지만 MapReduce는 범용 프로그래밍 언어를 기반으로 쉬운 병렬 프로그래밍을 가능하게 했다는 점에서 많이 사용되고 있으며 Dryad는 재사용이 쉽고 데이터 처리 흐름을 유연하게 작성할 수 있다는 점에서 장점을 가지고 있다.

### 1. 서론

클라우드 컴퓨팅은 세계적인 시장조사기관인 가트너사가 선정한 10대 전략기술에서 2010년에 이어 2011년에도 2년 연속 1위로 선정되었고[1], 진 홀 가트너 최고경영자는 2014년엔 클라우드 기반 사업자 중 적어도 한 곳이 글로벌 톱10 IT서비스 사업자에 이름을 올리게 될 것[2]이라고 확언을 할 정도로 클라우드 컴퓨팅의 중요성은 더욱 높아지고 있다.

이러한 관심이 지속되는 가운데, 폭발적으로 증가하고 있는 데이터를 클라우드 컴퓨팅 상에서 효율적으로 처리할 수 있는 방법에 대하여 많은 연구가 활발히 이루어지고 있다.

본 논문에서는 클라우드 컴퓨팅을 정의하고 클라우드 컴퓨팅에서의 MPP 방식(massively parallel processing workload)의 대용량 데이터 처리를 위한 대표적인 모델인 MapReduce[3], Dryad[4]에 대해 조사한다.

### 2. 클라우드 컴퓨팅

클라우드 컴퓨팅은 정의 자체가 모호하고 범위가 매우 넓기 때문에 다양한 정의들이 존재한다. 그 가운데 아래의 두 정의를 선별하여 본 논문의 클라우드 컴퓨팅의 정의로 차용한다.

- IBM : 웹 기반 애플리케이션을 활용하여 대용량 데이터베이스를 인터넷 가상공간에서 분산 처리하고 이 데이

터를 데스크톱PC, 휴대 전화, 노트북PC, PDA 등 다양한 단말기에서 불러오거나 가공할 수 있게 하는 환경

- UC Berkeley: “데이터 센터가 인터넷을 통해 서비스 형태로 제공하는 응용과 이를 위한 하드웨어와 시스템 소프트웨어”

위의 정의들을 요약하자면 클라우드 컴퓨팅이란 인터넷 기술을 활용하여 가상화된 컴퓨팅 자원을 서비스로 제공하는 것으로, 사용자는 IT자원을 필요한 만큼 빌려서 사용하고 사용한 만큼 비용을 지불하는 컴퓨팅을 지칭한다.

클라우드 컴퓨팅은 대규모 데이터, 처리시간이 극단적으로 높은 작업 등 여러 요청들을 효율적으로 처리하기 위해 비공유 구조로 연결된 수많은 노드들을 이용한 병렬 처리가 가장 큰 특징이다.

다음 장에서는 이러한 대용량 데이터 처리를 위한 대표적인 모델인 MapReduce와 Dryad에 대해 설명한다.

### 3. 대용량 데이터 처리를 위한 모델

#### 3.1 MapReduce

MapReduce는 대용량 데이터 셋을 처리하기 위한 프로그래밍 모델로, 추상화에 초점이 맞추어져 있다. MapReduce는 시스템의 분산 구조를 감추면서 범용 프로그래밍 언어를 기반으로 쉬운 병렬 프로그래밍을 가능하게 한다. MapReduce의 프로그래밍 모델은 Map과 Reduce 함수의 단순한 구조로 되어 있으며, 이 함수의 정의는 사

용자가 지정하게 한다.

MapReduce 프레임워크는 Map과 Reduce의 2 단계로 질의를 처리한다. 사용자가 작성한 map과 reduce 함수는 모든 워커들에게 자동으로 배포되어 각 단계에 각 함수들이 수행된다. 마스터에 의해 스케줄링된 작업은 다수의 태스크로 분할되고 분할된 태스크는 각 워커들에게 할당된다.

MapReduce는 분산 파일시스템인 Google File System(GFS)를 하부 저장 구조로 활용하며 하나의 마스터와 다수의 워커(worker)로 구성[6]된다. 마스터는 작업 스케줄링과 태스크(task) 할당 및 진행상황 모니터링을 전담한다. 사용자가 어떠한 작업을 제출하면 마스터는 해당 작업을 M 개의 Map 태스크로 분할하고, 이 태스크들을 각 워커들에게 할당한다. 각 워커들은 설정된 최대 태스크량만큼의 태스크를 동시에 수행할 수가 있는데, 이러한 워커들의 상황을 마스터가 파악하여 균등하게 태스크를 할당시킨다. Map 태스크를 수행 한 결과는 Reduce 단계에서의 입력 데이터로 이용된다. 이때 Map 워커들은 중간결과를 로컬 디스크에 저장하고 마스터에게 저장된 중간결과와 정보를 전달하고, 이 정보를 Reduce 태스크에 담아 Reduce 태스크 워커들에게 전달하여 작업을 최종 수행한다. 이러한 과정은 MapReduce가 각 워커들의 태스크 완료 시간의 균형을 맞추는 것을 목적으로 하는 런타임 스케줄링 전략을 가지기 때문에, 클러스터를 구성하는 노드들의 성능 차이가 크더라도 평균적으로 좋은 성능을 보여준다.

또한 클라우드 기반에서 처리하는 대부분의 데이터들은 장시간의 분석적 처리를 필요로 하는데 MapReduce는 단순하면서 강력한 내고장성을 가진다. 마스터가 주기적으로 보내는 핑에 대한 워커들의 응답으로 고장 발생 여부를 판단하게 되며, 문제가 있는 태스크를 다시 스케줄링 하고 재시작하는 단순한 방법으로 내고장성을 보장한다.

### 3.2 Dryad

Dryad는 MapReduce와 마찬가지로 사용자의 프로그램을 자동으로 분배 및 스케줄링하여 병렬 실행한다. 그렇기 때문에 개발자는 동시성 제어, 통신 등 분산 환경에 대한 고려 없이도 프로그램을 개발할 수 있다. Dryad는 기존에 작성된 프로그램들의 재사용과 유연한 데이터처리 흐름 지원을 목적으로 하기 때문에 데이터 흐름이 고정된 MapReduce와 달리 데이터 처리 흐름을 정의하기 위한 다양한 요소를 프로그래밍 모델에 포함시키고 있다. 첫 번째는 데이터 처리 연산자를 작성하는 것으로 정점(vertex) 프로그래밍 모델이다. 정점 프로그램은 빈번하게 일어나는 통신 유형들을 따르는 클래스들을 미리 만들어놓은 C++ 클래스 파일을 상속받아 작성하게 된다. 두 번째는 데이터 흐름을 작성하는 그래프 프로그래밍 모델이다. 그래프 프로그램은 데이터 흐름을 기술하는 역할을 하며 질의 처리 계획을 사용자가 직접 작성한다. 이렇게 작성된 프로그램

은 비순환 방향 그래프(directed acyclic graph)로 나타나는데, 여기에서 정점(vertex)은 연산자를 의미하며 간선(edge)는 연산자 간의 데이터 전달을 의미한다. 이는 자유로운 처리 흐름을 만들 수 있어 다양한 문제를 다룰 수 있지만 저수준의 언어였기 때문에 알고리즘을 기술하는데 있어 매우 복잡하였다. 따라서 이를 개선하기 위해 .NET 기반의 언어인 LINQ(.NET Language INtegrated Query)를 기반으로 개발된 고수준 프로그래밍 언어를 이용한 DryadLINQ를 제안하였다.

## 4. 결론

지금까지 살펴본 대용량 처리 모델은 프로그래밍을 통해 작업을 결정하므로 비정형적인 데이터의 처리에 특화되어있기 때문에 클라우드 환경에 적합한 특징을 가진다. 조사한 두 모델은 프로그래밍 언어를 이용하여 어떻게 데이터를 처리 할 것인지를 기술한 것으로, 서로 간에 많은 공통점이 있다. 하지만 MapReduce는 범용 프로그래밍 언어를 기반으로 쉬운 병렬 프로그래밍을 가능하게 하여 많은 곳에서 사용이 되고 있으며, Dryad는 기존에 작성된 프로그램들의 재사용이 쉽고 데이터 처리 흐름을 유연하게 작성할 수 있다는 점이 장점이다. 이러한 각 모델들의 특징을 알고 이에 맞춰 클라우드 컴퓨팅 시스템을 구성한다면 더 효율성이 높아질 것이라 판단된다.

이 논문은 2011년도 정부(교육과학기술부)의 재원으로 한국연구재단의 지원을 받아 수행된 기초연구사업임(2011-0002707)

## 참고문헌

- [1] Gartner, "Gartner Identifies the Top 10 Strategic Technologies for 2011", <http://www.gartner.com/it/page.jsp?id=1454221>
- [2] 안호천, "[인터뷰]진 홀 가트너 CEO", <http://www.ciobiz.co.kr/news/articleView.html?idxno=6306>
- [3] Jeffrey Dean and Sanjay Ghemawat. MapReduce: Simplified data processing on large clusters. In OSDI '04: Proceedings of USENIX Symposium on Operating System Design and Implementation, pp.137-150, 2004.
- [4] M. Isard and et al. Dryad: Distributed data- parallel programs from sequential building blocks. In Proceedings of the 2nd ACM SIGOPS/EuroSys European Conference on Computer Systems 2007, p.72. ACM, 2007.
- [5] 이경하, 최현식, 정연돈, "클라우드 컴퓨팅에서의 대용량 데이터 처리와 관리 기법에 관한 조사", 데이터베이스 제 38권 제 2호, pp.104-125, 2011
- [6] 민영수, 김홍연, 김영균, "클라우드 컴퓨팅을 위한 분산 파일 시스템 기술", 정보과학회지 제 27권 제 5호, pp.86-94, 2009