

ScienceDMZ 기반 빅데이터 고속도로 체계 및 AI 컴퓨팅 체계 연계 구축 방안 연구

문정훈, 김기현, 석우진

한국과학기술정보연구원

jhmoon@kisti.re.kr, kkh1258@kisti.re.kr, wjseok@gmail.com

A Study on the big data super highway and AI computing environment integrate based on ScienceDMZ

Moon Jeong Hoon, Kim Ki Hyeon, Seok Woo Jin

Korea Institute of Science and Technology Information

요 약

본 논문은 데이터 집약형 과학의 발달로 대용량 데이터의 생산이 급증하고 있으며, 또한 IT 기술의 발달로 인하여 실험, 관측, 계산 장비들의 고도화에 따른 연구결과와 대용량과 복잡화 되어 가는 추세에 있다. 이러한 동향은 4차 산업 시대에 전세계적인 현상이며, 선진국을 중심으로 여러 나라들이 컨소시엄을 구성한 국제 프로젝트들의 등장으로 엑사바이트 규모 이상의 데이터가 생산되고 있다. 이러한 거대 규모의 데이터는 한 두 곳의 데이터 센터에서 처리하기 어려운 수준이며, 이것은 전송, 저장, 공유의 문제를 의미하기도 한다. 따라서 이러한 과학 빅데이터들은 전세계 주요 데이터 센터들에 분산되어 저장되며, 실제 연구자들은 연구에 투자하는 시간보다 이렇게 분산되어 있는 데이터들을 모으고, 분류하는데 더 많은 시간을 사용하고 있다. 따라서 본 논문에서는 데이터 집약형 과학분야의 빅데이터를 위한 빅데이터 고속도로 체계인 ScienceDMZ의 구축과 이러한 빅데이터를 기반으로 AI 연구를 위한 AI 컴퓨팅 체계의 연계 구축 방안을 제안한다.

1. 서 론

본 논문에서는 데이터 집약형 과학분야의 비약적인 발전과 IT 기술의 발전으로 인하여 과거에 비해 수백배 이상의 과학 데이터들이 생성되고 있다. 이런 대용량 데이터의 생성은 실험, 관측 장비로부터 저장을 위한 데이터 센터로의 전송과 데이터 센터에서 연구자에게 전송을 효과적으로 해야 하지만 기존 네트워크의 전송 성능의 한계로 인하여 연구자에게 적절한 전송 환경이 제공되지 못하고 있다. 특히 미국 과학재단에서는 향후 과학분야에서 중점적으로 지원해야 하는 데이터 집약형 과학 분야를 7개로 분류하고 있는데 입자물리, 천문 관측 및 천체물리, 바이오분야의 게놈 데이터, 의료데이터, 지진, 기상 분야를 포함하는 지구과학 데이터, 인공지능분야, 가상현실 분야이다. 이들 분야 중 입자물리는 스위스 CERN에서 발생하는 데이터를 전세계 데이터 센터로의 전송과 분석을 통해 새로운 입자를 발견하고 있으며, 천문분야의 LSST 프로젝트 [1]는 남미 칠레의 관측 망원경으로부터 관측된 자료를 전세계 연구자들과 공유하는 프로젝트로서 하루 15TB~30TB가 발생하고 있다. 이러한 거대 데이터는 한 두 곳의 데이터 센터에서 처리할 수 있는 수준이 아니기 때문에 전세계의 빅데이터 센터로 분산 저장되고 있다. 또한 분산 저장된 데이터들을 연구자의 입장에서 하나의 데이터 셋으로 구성하여 원하는 연구를 할려면 고속의 전송과 분석 작업이 효과적으로 연계되어 진행되어야 한다. 이러한 과학분야의 빅데이터를 효과적으로 처리하기 위해서는 전송 성능의 향상과 이와 연계된 AI 컴퓨팅 작업이 함께 이루어지는 것이 바람직한데 기존 TCP/IP 프로토콜에서 TCP는 아직까지 데이터 전송에서 가장 적합한 프로토콜이다[2]. 그러나 망에서 여러 가지 요인들로 인하여 전송 성능의 한계가 있다. 그중에서 특히 불특정한 이유로 발생하는 패킷로스는 그림 1에서와 같이 고대역망에 대해서 최저 성능까지 떨어졌다가 다시 회복되는 과정을 거쳐 기존 성능으로 복구가 되는데

이러한 현상은 망 자체가 리부팅 되는 수준으로 성능저하가 일어난다.[3]

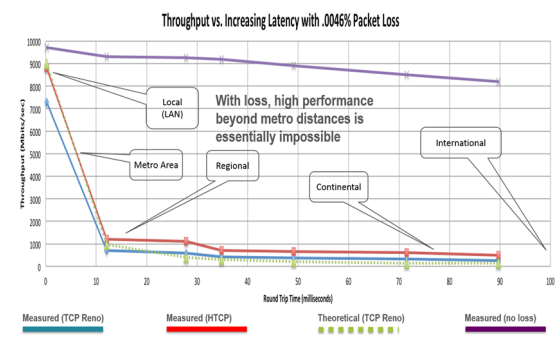


그림 1. 패킷로스 발생에 따른 네트워크 전송 성능 [3]

그러나 대부분의 과학분야 응용프로그램들은 TCP 기반으로 발전해 왔기 때문에 새로운 기술 체계로의 전환은 어렵다. 따라서 이러한 문제를 개선하기 위해 기존 네트워크 환경의 변화 없이 빅데이터 고속도로 체계를 구축하는 기술인 ScienceDMZ 기술의 적용을 통해 이러한 전송 문제를 해결한다. 또한 인공지능의 경우 GPU 기반의 컴퓨팅 환경이 필수적이지만 한 연구자 또는 대학에서 보유하고 있는 GPU 자원의 한계로 인하여 고성능의 계산은 어려운 실정이다. 이러한 연구 환경에 대해 ScienceDMZ기반의 빅데이터 고속도로를 통한 고성능 전송 환경을 분산 환경에 적용하여 분산되어 있는 전산 자원들을 오케스트레이션 기법을 통하여 연동함으로써 분산되어 있는 GPU자원들을 단일 자원화하여 사용하게 된다.

ACKNOWLEDGMENT

국가과학기술연구회 “출연연 중심 AI 융합을 위한 빅데이터 Super Highway 융합클러스터” 사업

참 고 문 헌

- [1] 지구 남반구 관측 천문 국제 프로젝트, URL, <https://www.lsst.org/>
- [2] J. Postel. “Transmission Control Protocol. Request for Comments(Standard) 793”, Internet Engineering Task Force, September, 1981
- [3] E. Dart, L. Rotman, B. Tierney, “The ScienceDMZ: A Network Design Pattern for Data-Intensive Science”, Scientific Programming Volume 22, Issue 2, Pages 173-185
- [4] 미국국가연구망 ScienceDMZ, URL, ScienceDMZ<https://fasterdata.es.net/science-dmz/DTN/>