

GSDC 구조생물학분야 TEM 분석 팜 소개

(‘20.2. 한국과학기술정보연구원 국가슈퍼컴퓨팅본부 대용량데이터허브센터)

□ 글로벌 대용량데이터허브센터(GSDC) 소개



- 데이터집약형 기초과학 실험지원을 위한 한국형 통합 데이터센터 역할
- 약 10,000+ 코어, 10+ PB 스토리지 및 첨단연구망 연계 기초과학 데이터 통합 인프라 구축·운영
- 세계 주요 가속기 및 대형 검출기 등에서 생성된 대용량의 실험데이터에 대한 저장·공유·분석 환경 구축 및 연구자 맞춤형 서비스를 위한 연구개발 수행

국내 7개 연구 커뮤니티 공동활용 데이터집약형 다양한 연구도메인으로 확대

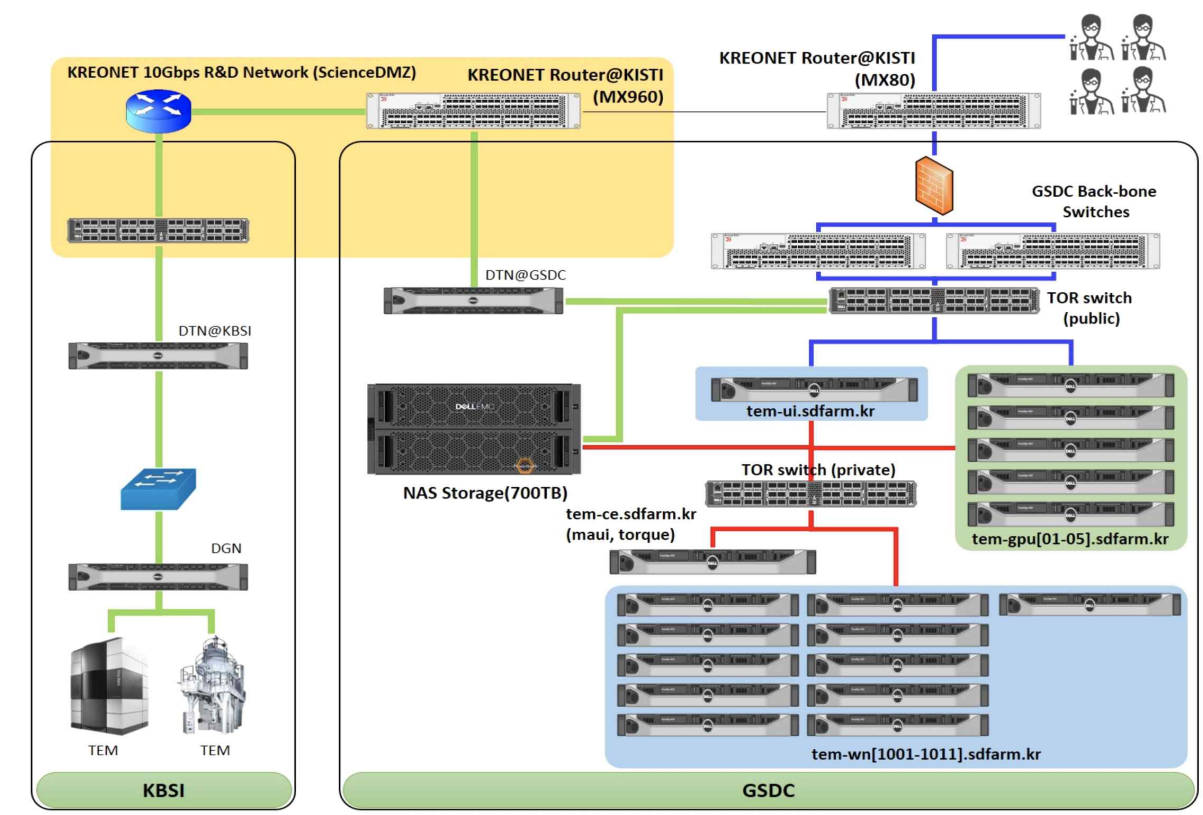


□ 국내 대형연구장비 연계 데이터 집약형 연구지원 확대

- ALICE, CMS, BELLE, LIGO, 유전체, RENO 실험 이외에 국내 대형연구장비 (전자현미경, 포항방사광가속기 등)와 연계 데이터 집약형 연구지원 확대
- 오창 KBSI에 있는 초극저온전자현미경과 연계한 TEM 데이터 분석 팜을 구축·운영 시범 서비스 (2017년~) 및 공식 서비스 (2019년 09월 ~)
- 포항가속기연구소 4세대 방사광가속기와 연계한 데이터 분석 환경을 구축·운영 시범 서비스 (2018년~)

□ GSDC 구조생물학분야 TEM 데이터 분석 팜 현황

- KBSI ↔ GSDC 간 전용의 10Gbps 고속데이터 전송망 제공 (KREONET, KISTI 첨단연구망 ScienceDMZ 활용)
- 408 CPU 코어, 10 GPGPU, 800 TB 스토리지 제공
- Relion, cisTEM, cryoSPARC 등 다양한 데이터 분석 도구 활용 환경 제공



□ GSDC TEM 분석 팜 컴퓨팅/스토리지 자원 (2020년 2월 기준)

- TEM 분석 팜 구조 및 시스템 명세

Category	Name	Specification	Resources size
Login	tem-ui.sdfarm.kr	<ul style="list-style-type: none"> CPU : Intel(R) Xeon(R) CPU E5-2697v3 @ 2.60GHz 14Core * 2 CPUs RAM : DDR4 8GB * 24 (192GB) HDD : 12G SAS HDD 1.2TB * 2EA (RAID-1) 	28 cores
Computing (master)	tem-ce.sdfarm.kr	<ul style="list-style-type: none"> CPU : Intel(R) Xeon(R) CPU E5-2697v3 @ 2.60GHz 14Core * 2 CPUs RAM : DDR4 8GB * 24 (192GB) HDD : 12G SAS HDD 1.2TB * 2EA (RAID-1) 	28 cores
Computing (workers)	tem-wn[1001-1011].sdfarm.kr	<ul style="list-style-type: none"> CPU : Intel(R) Xeon(R) CPU E5-2697v3 @ 2.60GHz 14Core * 2 CPUs RAM : DDR4 8GB * 24 (192GB) HDD : 12G SAS HDD 1.2TB * 2EA (RAID-1) 	308 cores
	tem-gpu[01-05].sdfarm.kr	<ul style="list-style-type: none"> CPU : Intel® Xeon® CPU E5-2690v4 @ 2.60GHz 14Core * 2 CPUs RAM : DDR4 16GB * 24 (384GB) SSD : 6G SATA SSD 800GB * 2EA (RAID-1) GPU : NVIDIA P100 * 2ea (tem-gpu[01-03]) GPU : NVIDIA P40 * 2ea (tem-gpu[04-05]) 	140 cores
Storage	Dell EMC Isilon NAS	Network attached storage 700 TB	
Total		504 CPU cores, 10 GPGPUs, 700TB Storage	

※ 2020년 100 CPU 코어, 100TB 스토리지 증설 예정 (총 408 CPU 계산 코어, 10 GPGPU, 800TB 스토리지)

- TEM 분석 팜 관리 소프트웨어 (배치 시스템, OpenMPI, Nvidia CUDA 등)

Category	Name	Description	Version (module path)
OS	Scientific Linux	Operating system	6.x
System M/W	Environment module	<ul style="list-style-type: none"> Module environment https://modules.readthedocs.io/en/latest 	v3.2.10
	OpenPBS(torque)	<ul style="list-style-type: none"> Cluster resources management http://www.adaptivecomputing.com/products/torque 	v6.1.2
	OpenMPI	<ul style="list-style-type: none"> Messaging Pass Interface(MPI) Reference implementation for MPI standard https://www.open-mpi.org 	v1.8.8 (mpi/gcc/openmpi/1.8.8)
	cuda	<ul style="list-style-type: none"> Compute Unified Device Architecture(CUDA) NVIDIA CUDA Runtime & Toolkit https://developer.nvidia.com/cuda-toolkit 	9.1 (cuda/9.1)
	python	<ul style="list-style-type: none"> Python runtime 	v2.6.6

- 데이터셋 분석 도구

Category	Name	Description	Version (module path)
Tools	Relion	A stand-alone computer program that employs an empirical Bayesian approach to refinement of (multiple) 3D reconstructions or 2D class averages in electron cryo-microscopy (cryo-EM). • https://www3.mrc-lmb.cam.ac.uk/relion/index.php	v3.0.7 (apps/gcc/4.4.7/relion/cpu/3.0.7) (apps/gcc/4.4.7/relion/gpu/3.0.7)
	cisTEM	User-friendly software to process cryo-EM images of macromolecular complexes and obtain high-resolution 3D reconstructions. • https://cistem.org	v1.0.0 (apps/gcc/4.4.7/cistem/1.0.0)
	CryoSPARC	CryoSPARC is the state-of-the-art platform used globally for obtaining 3D structural information from single particle cryo-EM data. • https://cryosparc.com	Not deployed yet (TBD)

※ 2020년 웹 기반 CryoSPARC 시범 서비스 예정. Relion, cisTEM, cryoSPARC 등 지속적인 업데이트

□ GSDC TEM 분석 팜 사용자 계정 (account) 별 신청 절차

1. 첨부 TEM 계정신청서 작성 후, 이메일로 계정신청서 제출 (연구자 → GSDC TEM 서비스 담당자)
2. 계정 및 접속 정보 (GSDC TEM 서비스 담당자 → 연구자)

□ GSDC TEM 분석 팜 사용자 계정 (account) 별 지원 현황

- tem-ui.sdfarm.kr(로그인 서버)을 통한 SSH 터미널 접속 지원
- 데이터 분석 작업 실행을 위한 CPU 계산자원 공용 큐, GPU 계산자원 공용 큐
- 원본 데이터셋 저장 및 분석용 스토리지 20TB
- 데이터셋 분석 도구 실행을 위한 사용자 기술 지원

□ GSDC 구조생물학분야 TEM 서비스 사용자 매뉴얼

```
GSDC TEM Relion

* Official GSDC TEM users guide : https://tem-docs.readthedocs.io
=====
* Hostname.....: tem-ui.sdfarm.kr
* OS Release.....: Scientific Linux release 6.10 (Carbon)
* System uptime...: 217 days 6 hours 5 minutes 40 seconds
* Users.....: Currently 8 user(s) logged on
* Processes.....: 1263 running
* CPU usage.....: 0.87, 0.48, 0.36 (1, 5, 15 min)
* Memory (used/total)...: 3293 MB / 193583 MB
* Swap in use.....: 181 MB
-----
* TEM disk (used/total)..: 150 TB / 800 TB (19%)
* Current user.....: tem
* Home directory.....: /tem/home/tem
* Disk Quota limit.....: 20480 GB
* Disk usage.....: 877 GB (4.2838 %)
* # of Files.....: 1571741
=====
[tem@tem-ui ~]$
```

<https://tem-docs.readthedocs.io>

- GSDC TEM 분석 팜 개요 및 소개
- TEM 분석 팜 접속 방법 (리눅스/맥, 윈도우 환경)
- 데이터 분석 도구 활용을 위한 모듈 환경 (module environment)
- 배치시스템 작업 스크립트 작성 방법 및 예시
- 배치시스템 작업 공용 큐 (Queue) 소개 (CPU Queue, GPU Queue)
- Relion : CPU/GPU 계산자원을 활용한 데이터 분석 방법
- cisTEM : CPU 계산자원을 활용한 데이터 분석 방법 등

□ 연락처 및 문의

- 유정록 (junglok.yu@kisti.re.kr, 042-869-0622 KISTI GSDC)
- 여일연 (ilyeon9@kisti.re.kr, 042-869-0658, KISTI GSDC)

□ 추가 안내사항

- KBSI가 아닌 연구그룹 연구실에서의 원활한 실험데이터셋 전송 (GSDC ↔연구실)를 위해서는 첨단연구망 사용 신청 필요
- 매년 초 신청, 평가 진행됨

