# An NDN Testbed for Large-scale Scientific Data

**Huhnkuk Lim**

**Korea Institute of Science & Technology Information (KISTI)**

**NDNComm 2015**
**Sep. 28, 2015**

# Motivations on NDN for Large-scale Scientific Application

- **As the data volumes and complexity increase, data-intensive science cannot rely on extension in the storage infrastructure.**

- **It needs to investigate new methods of intelligent processing and data distribution over networks.**

- **Use of caching technique changes traffic pattern in the network and improves corrupted data rate.**

- **NDN based large-scale scientific application**
  - **Climate modeling application as an initial focus**
  - **Extension of NDN architecture to various data-intensive science application such as HEP and astronomy with hierarchical naming strategies**

- **Innovative data management lead to traffic pattern change**

# Backgrounds on NDN for Climate Modeling Application

## Why climate data transfer using NDN Architecture

- Current CMIP5 data transfer using ESGF, long time latency and corrupted data occur
- To provide innovative transfer, management, and security function for scientific big data using the NDN architecture
- Movement of traffic pattern in data-intensive science and reduction of data explosion on it

## R&D on NDN based data-intensive science application

- NDN testbed for climate modeling application (CSU univ.)

- NDN architecture design, development, and deployment for LHC big data transfer (Fermi Lab)

- ESnet for research networks in US



data discovery and fast retrieval

UCSD (planned)

**Climate modeling NDN testbed in US**

## Data-intensive science applications

1. Climate Modeling

2. HEP (LHC, CMS)

3. Astronomy

3

# NDN Testbed for Climate Modeling Application

**Graphic User Interface (Web browser)**



| **Functions of front-end system in consumer** | | **Functions of back-end system in producer** |
|---|---|---|
| ▪ To provide GUI for climate modeling application based on NDN architecture<br>▪CMIP5 data search using controlled vocabulary<br>▪NDN name based CMIP5 data downloading | **Kisti-ndn-atmos package** | ▪ To translate .nc file names to NDN names<br>▪ NDN based repository establishment for CMIP5 data management<br>▪ NDN name database establishment, in order to search a CMIP5 data of interest in producer |

# Key Components in the NDN Testbed

**Category based search**

**Keyword based search**

**GUI to support NDN based climate modeling application**

**NDN Name Translator for climate modeling application**



Query result

Results

◈ Name lists sorting
◈ To show meta data corresponding to each searched CMIP5 data
◈ Search results is changed to CMIP5 file name following DRS syntax

**Works to support NDN based Climate Modeling Application**

◈ To translate CMIP5 data files stored in NDN repository to NDN names and to store them in DB
◈ NDN name translation following DRS structure

**NDN network for climate modeling in Korea**

◈ Forwarding and caching of interest/data packets
◈ **Synchrinized FIB table management in the NDN testbed**
◈ NDN platform (ver O.3.4)
  – NDN-cxx, NFD
  – NDN-js (one of NDN-ccl)
  – NDNfs-port

5

# Features of GUI (1)

- **Reflection of the ESGF system workflow**
- **CMIP5 climate data searching following climate DRS structure**
  - To show original CMIP5 nc file names changed from NDN names, together with meta data sets corresponding to .nc file names
  - Key word based CMIP5 data search and user-friendly sorting for search results



<MetaData for the above nc file>

- **CMIP5 data downloading in metadata window**
  - **Download button have the address corresponding to an NDN name of interest in producer side**
    - Address: NDN name based URI
    - "ndn:/catalog/myUniqueName/*<CMOR fiflename.nc>*"
      - ex) ndn:/catalog/myUniqueName/*psl_amip_MIROC5_historical_r1i1p1_1950010100-xx.nc*



**\<Downloading of CMIP5 climate data\>**

# Features of Name Translator

- **To translate all nc file names stored in repository to NDN names**
  - **Parsing of each name component**
  - **To check time variable in an nc file has the same value in metadata**
    - Sometimes, time in metadata is slightly different from one in real data.
    - For allowable error range, name translation for an nc file name.
    - If they are outside from it, no translation for that one.

**6 nc files in NDN file system (repository)**

```
uns@ubuntu:/tmp/ndnfs$ ls
pr_day_GFDL-ESM2M_historical_r1i1p1_20010101-20051231.nc
psl_6hrPlev_MIROC5_historical_r1i1p1_1950010100-1950123118.nc
snoToIce_OImon_bcc-csm1-1-m_rcp85_r1i1p1_210001-210012.nc
snw_LImon_GFDL-ESM2M_1pctCO2_r1i1p1_019601-020012.nc
tasmin_day_GFDL-CM3_historical_r1i1p1_20050101-20051231.nc
va_6hrPlev_GFDL-CM3_historical_r1i1p1_2005010100-2005123123.nc
uns@ubuntu:/tmp/ndnfs$
```

*name translation*

**6 CMIP5 NDN names translated in Mysql DB repository**

```
name                                                                    | activity | product | organization | model       | experiment | freq

/CMIP5/output/BCC/bcc-csm1-1-m/rcp85/mon/seaIce/snoToIce/r1i1p1/210001-210012/     | CMIP5 | output  | BCC          | bcc-csm1-1-m | rcp85      | mon

/CMIP5/output1/NOAAGFDL/GFDL-CM3/historical/6hr/atmos/va/r1i1p1/2005010100-2005123123/  | CMIP5 | output1 | NOAAGFDL     | GFDL-CM3     | historical | 6hr

/CMIP5/output1/NOAAGFDL/GFDL-ESM2M/historical/day/atmos/pr/r1i1p1/20010101-20051231/    | CMIP5 | output1 | NOAAGFDL     | GFDL-ESM2M   | historical | day

/CMIP5/output1/NOAAGFDL/GFDL-ESM2M/rcp45/day/atmos/ta/r1i1p1/20960101-21001231/         | CMIP5 | output1 | NOAAGFDL     | GFDL-ESM2M   | rcp45      | day

                                                                          GFDL-CM3     | historical | day

                                                                          MIROC5       | historical | 6hr
```

| name | sha256 | activity | product | organization | model | experiment | frequency | modeling_realm | variable_name | ensemble | time |
|------|--------|----------|---------|--------------|-------|------------|-----------|----------------|---------------|----------|------|
| Full name | Hash value | CMIP5 | output | MIROC | MIROC5 | historical | 6hr | atmos | psl | r1i1p1 | 1968 ..... |

Database schema => http://redmine.named-data.net/projects/ndn-atmos/wiki/Schema

8

# Summary of kisti-ndn-atmos SW Package

## Summary of kisti-ndn-atmos SW package

| Key function | | kisti-ndn-atmos |
|---|---|---|
| User Interface | Data search | To show .nc file name lists following DRS structure |
| | Metadata | Supported |
| | File downloading | Supported |
| | User-friendly functions | Sorting and key word based searching |
| Name translator | | NDN name translation for valid climate data |
| Repository for NDN | | To provide a repository using ndnfs-port |

**There have been significant code sharing between KISTI and CSU project, in order to develop each ndn-atmos SW package for climate application**

# Climate Data Transfer by Federated NDN Testbed in Korea and US

- **Transfer by the Earth System Grid Federation (ESGF) infrastructure**
  - ESGF: Distributed CMIP5 data management protocol in current IP based networks
  - Data explosion for duplicate big data requests results in BW waste

- **Transfer by federated NDN testbeds**
  - Smart transfer for duplicate big data requests
  - Change of traffic pattern results in traffic reduction in networks
  - Prevention of data explosion in networks



ESGF architecture based CMIP5 delivery



NDN based CMIP5 delivery

- **Current works on federated NDN Testbed in Korea and US**
  - Interoperability for front and back-end systems in each doman
  - To create synchronized FIB tables to search for all CMIP5 data sets at each producer using NLSR
  - Caching scheme for large scale scientific data

# Summary and Future Works

- **Current climate data transfer by ESGF results in long time latency and high corrupted data rate.**

- **To provide large-scale scientific data with innovative transfer and management.**

- **To change traffic pattern in data-intensive science and to prevent data explosion in networks.**

- **NDN testbed with kisti-ndn-atmos package for climate application**
  - Front-end system in consumer and back-end system in producer
  - To show original climate .nc file names following DRS and corresponding metadata sets
  - Key word based climate data search and downloading
  - To translate all .nc file names stored in the NDN repository to NDN names
  - Forwarding and caching of interest/data packets on climate modeling application

- **Future works**
  - Federated NDN testbed in Korea and US for climate modeling application
  - Performance analysis for ESGF and NDN based transfer
  - Caching and mobility to consider characteristics of large-scale scientific data