TSINGHUA UNIVERSITY
CENTER FOR EARTH SYSTEM SCIENCE
清华大学地球系统科学研究中心

# CAFE: a Collaborative Analysis Framework for distributed Environmental data

Hao Xu
On behalf of the CAFE Team

Dec.7th 2016

Outline

# Introduction-about CAFE

CAFE is a dedicated software package for **collaborative analysis** of large volumes of **distributed environmental data.**

**Key features:**

1) Computing near the data;

2) data is logically grouped, while physically distributed;

3) Analytic tasks are divided as subtasks, and then fulfilled on corresponding nodes;

4) Easy way to enrich the built-in analytic functions;

5) Open source projects on github

# Introduction-why CAFE?

## Typical workflow

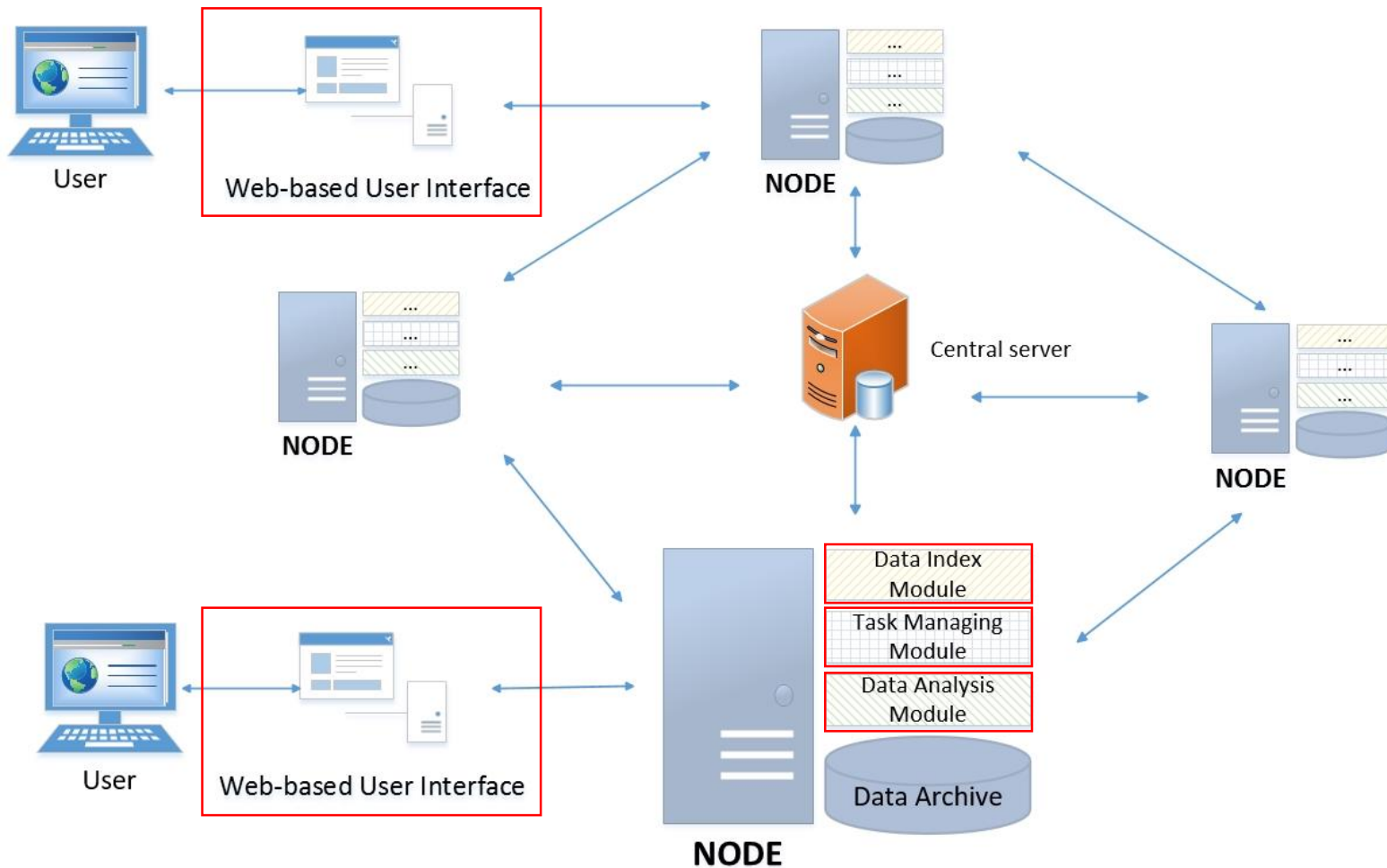- ❖ Download data from multiple nodes
- ❖ Data subsetting
- ❖ Data regridding
- ❖ Data averaging
- ❖ Writing analytic scripts
- ❖ Executing codes locally
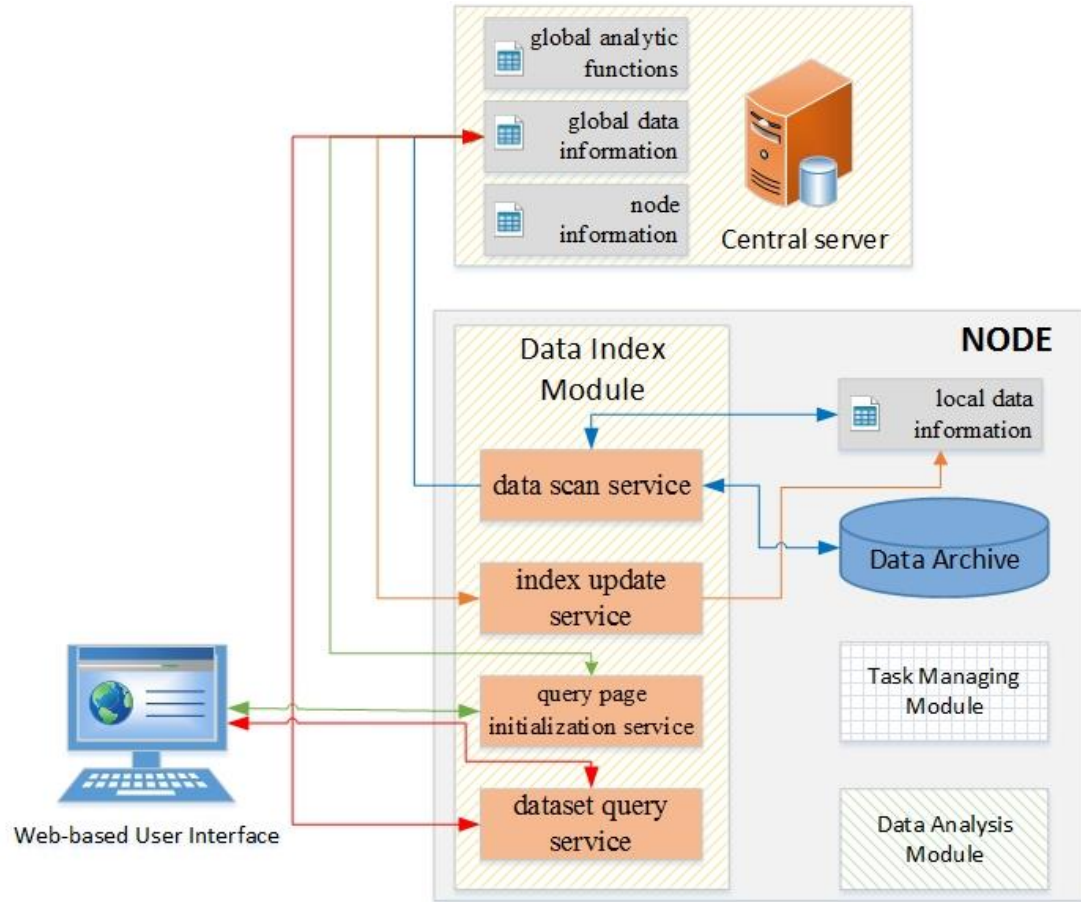- ❖ Result visulization
- ❖ Result analysis
- ❖ *......*

## CAFE Features

- ❖ Web-based UI provided
- ❖ REST APIs availbale
- ❖ One-stop service for data discovery, visualization, and analysis
- ❖ User-transparent
- ❖ Task management
- ❖ Multi-node collaboration
- ❖ Support for data intercomparison
- ❖ Multiple built-in analytic functions
- ❖ Easy way for extensions
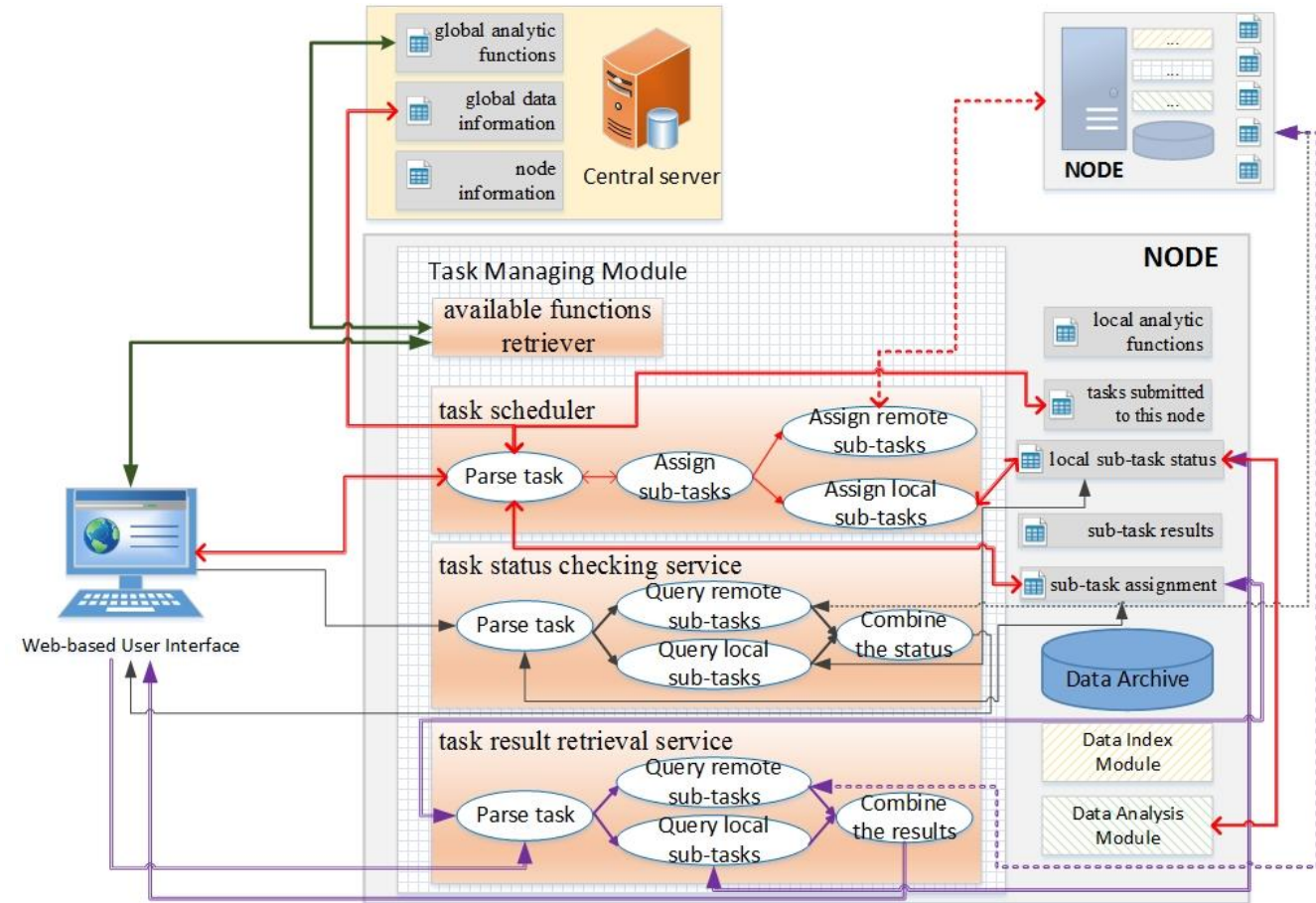- ❖ Batch analysis

    *......*

A p2p architecture

# System design-modularity

**Data Index Module**
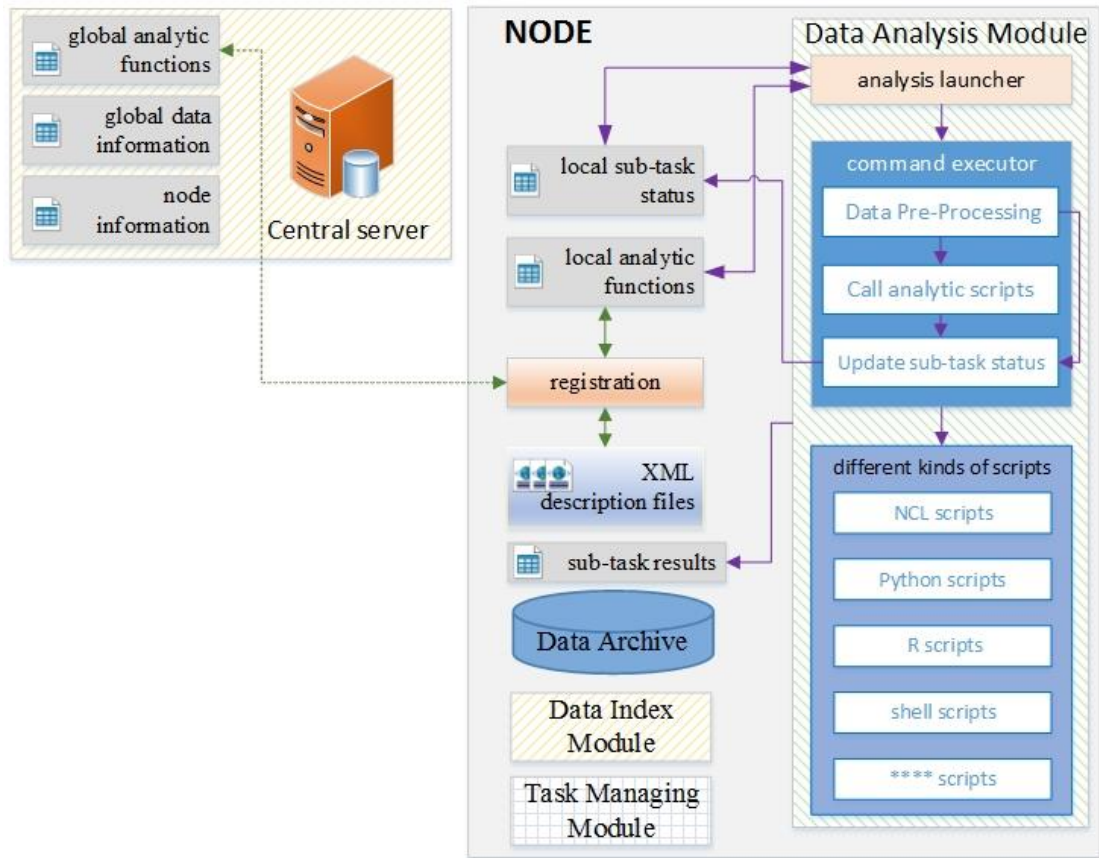


**Task Managing Module**

Data Analysis Module

Web-based User Interface

TSINGHUA UNIVERSITY
CENTER FOR EARTH SYSTEM SCIENCE
清华大学地球系统科学研究中心

```xml
<?xml version="1.0"?>
<function lang="NCL">  <!--defining the language of the analytic function-->
    <name>PolarNHEOF</name>  <!--the name of the analytic function-->
    <script>PolarNHEOF.ncl</script> <!--the file name the script-->
    <description>xxx</description> <!--description of the function-->
    <isGlobalFunction>true</isGlobalFunction>
    <!--defining if the function can be distributed to all the nodes-->
    <MultiInputFiles>false</MultiInputFiles>
    <!--defining if the function can be distributed to all the nodes-->
    <InputFileFormat>netCDF</InputFileFormat> <!--setting format of the input file-->
    <properties>
        <!--defining which datasets can use this function, type can be "include" or "exclude"-->
        <!--the contributor can define model,frequency,modelingRealm and variableName values for filtering-->
        <Model type="include"><!--refer to PCMDI documents to get acceptable values-->
            <value>xxx</value>
            ...<!--multiple values-->
        </Model>
        ... <!--frequency,modelingRealm,variableName-->
    </properties>
    <Controls>  <!--setting the parameters that input from the webpage-->
        <parameter>
            <!--defining name,description,tag,type and value range of the parameter-->
        </parameter>
        ... <!--multiple parameters-->
    </Controls>
    <InputFileParameters>  <!--setting the parameters about the input file(s)-->
        <parameter>
            <!--defining name,description,tag,type and value rangeof the parameter-->
        </parameter>
        ...
    </InputFileParameters>
    <OutputFileParameters>  <!--setting the parameters about the output file(s)-->
        <result>
            <filetype>xxx</filetype> <!--png/nc/txt...-->
            <filecount>1</filecount>
            <parameter>
                <!--defining name,description,tag,type and value range of the parameter-->
            </parameter>
            ...
        </result>
        ... <!--multiple kinds of results-->
    </OutputFileParameters>
    <pre-processing> <!--defining pre-processing type before invoking the script-->
        <pre-processing-type>Latlon</pre-processing-type> <!--Origin/Latlon/YearAvg/SeasAvg/LTM/Subset-->
        <AddToResults>false</AddToResults>
        <!--defining if the pre-processing result need to be added to the result files-->
    </pre-processing>
</function>
```

**properties**

**controls**

**Input/Output related**

**Pre-processing**

Easy way for extensions. Only the analytic script and its XML description are needed.
The script should have an I/O interface for command line and can be invoked by Java.

```
#Example: for NCL
if(.not.isvar("date")) then
    date=19000101
end if
if(.not.isvar("filename")) then
    filename="test.nc"
end if    print(date)
print(filename)
#command invoking: ncl test.ncl
date=19000120'filename="test1.nc"'
```

8

```xml
<?xml version="1.0"?>
<function lang="NCL">  <!--defining the language of the analytic function-->
    <name>PolarNHEOF</name>  <!--the name of the analytic function-->
    <script>PolarNHEOF.ncl</script> <!--the file name the script-->
    <description>xxx</description> <!--description of the function-->
    <isGlobalFunction>true</isGlobalFunction>
    <!--defining if the function can be distributed to all the nodes-->
    <MultiInputFiles>false</MultiInputFiles>
    <!--defining if the function can be distributed to all the nodes-->
    <InputFileFormat>netCDF</InputFileFormat> <!--setting format of the input file-->
    <properties>
        <!--defining which datasets can use this function, type can be "include" or "exclude"-->
        <!--the contributor can define model,frequency,modelingRealm and variableName values for filtering-->
        <Model type="include"><!--refer to PCMDI documents to get acceptable values-->
            <value>xxx</value>
            ...<!--multiple values-->
        </Model>
        ... <!--frequency,modelingRealm,variableName-->
    </properties>
    <Controls>  <!--setting the parameters that input from the webpage-->
        <parameter>
            <!--defining name,description,tag,type and value range of the parameter-->
        </parameter>
        ... <!--multiple parameters-->
    </Controls>
    <InputFileParameters>  <!--setting the parameters about the input file(s)-->
        <parameter>
            <!--defining name,description,tag,type and value rangeof the parameter-->
        </parameter>
        ...
    </InputFileParameters>
    <OutputFileParameters>  <!--setting the parameters about the output file(s)-->
        <result>
            <filetype>xxx</filetype> <!--png/nc/txt...-->
            <filecount>1</filecount>
            <parameter>
                <!--defining name,description,tag,type and value range of the parameter-->
            </parameter>
            ...
        </result>
        ... <!--multiple kinds of results-->
    </OutputFileParameters>
    <pre-processing> <!--defining pre-processing type before invoking the script-->
        <pre-processing-type>Latlon</pre-processing-type> <!--Origin/Latlon/YearAvg/SeasAvg/LTM/Subset-->
        <AddToResults>false</AddToResults>
        <!--defining if the pre-processing result need to be added to the result files-->
    </pre-processing>
</function>
```
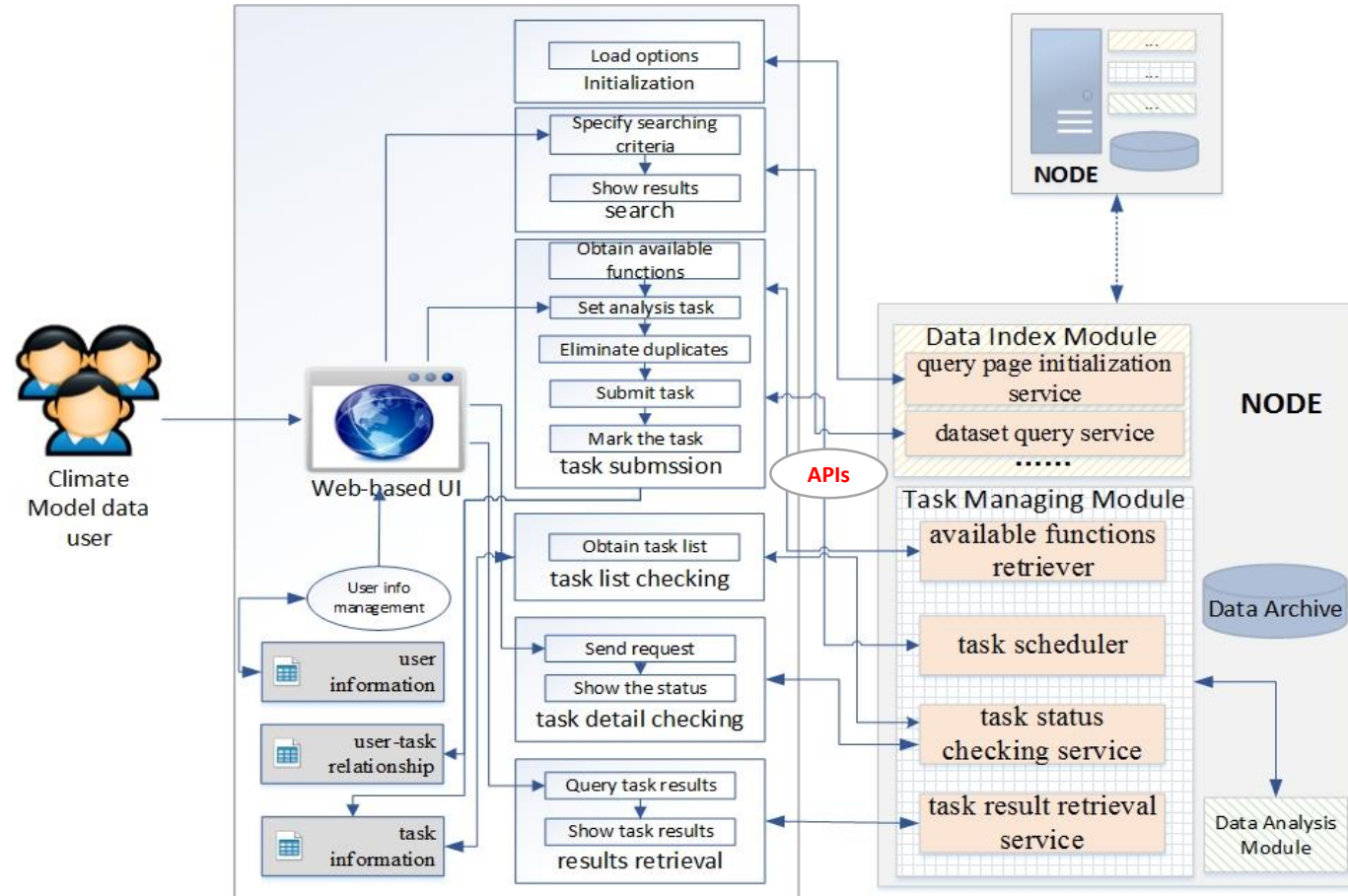
**properties**

**controls**

**Input/Output related**

**Pre-processing**

Easy way for extensions. Only the analytic script and its XML description are needed.
The script should have an I/O interface for command line and can be invoked by Java.

Shell:

```sh
#test.sh
#!/bin/sh
date=$1
filename=$2
echo "date:${date} filename:${filename}"
#command invoking: sh test.sh 19000101 test.nc
```

9

TSINGHUA UNIVERSITY
CENTER FOR EARTH SYSTEM SCIENCE
清华大学地球系统科学研究中心

```xml
<?xml version="1.0"?>
<function lang="NCL">  <!--defining the language of the analytic function-->
    <name>PolarNHEOF</name>  <!--the name of the analytic function-->
    <script>PolarNHEOF.ncl</script>  <!--the file name the script-->
    <description>xxx</description>  <!--description of the function-->
    <isGlobalFunction>true</isGlobalFunction>
    <!--defining if the function can be distributed to all the nodes-->
    <MultiInputFiles>false</MultiInputFiles>
    <!--defining if the function can be distributed to all the nodes-->
    <InputFileFormat>netCDF</InputFileFormat>  <!--setting format of the input file-->
    <properties>
        <!--defining which datasets can use this function, type can be "include" or "exclude"-->
        <!--the contributor can define model,frequency,modelingRealm and variableName values for filtering-->
        <Model type="include"><!--refer to PCMDI documents to get acceptable values-->
            <value>xxx</value>
            ...<!--multiple values-->
        </Model>
        ... <!--frequency,modelingRealm,variableName-->
    </properties>
    <Controls>  <!--setting the parameters that input from the webpage-->
        <parameter>
            <!--defining name,description,tag,type and value range of the parameter-->
        </parameter>
        ... <!--multiple parameters-->
    </Controls>
    <InputFileParameters>  <!--setting the parameters about the input file(s)-->
        <parameter>
            <!--defining name,description,tag,type and value rangeof the parameter-->
        </parameter>
        ...
    </InputFileParameters>
    <OutputFileParameters>  <!--setting the parameters about the output file(s)-->
        <result>
            <filetype>xxx</filetype> <!--png/nc/txt...-->
            <filecount>1</filecount>
            <parameter>
                <!--defining name,description,tag,type and value range of the parameter-->
            </parameter>
            ...
        </result>
        ... <!--multiple kinds of results-->
    </OutputFileParameters>
    <pre-processing> <!--defining pre-processing type before invoking the script-->
        <pre-processing-type>Latlon</pre-processing-type> <!--Origin/Latlon/YearAvg/SeasAvg/LTM/Subset-->
        <AddToResults>false</AddToResults>
        <!--defining if the pre-processing result need to be added to the result files-->
    </pre-processing>
</function>
```

**properties**

**controls**

**Input/Output related**

**Pre-processing**

Easy way for extensions. Only the analytic script and its XML description are needed.
The script should have an I/O interface for command line and can be invoked by Java.

Python:

```python
#test.py
    import sys
    date=sys.argv[1]
    filename=sys.argv[2]
    print date,filename
#command invoking: python test.py 19000101 test.nc
```

```xml
<?xml version="1.0"?>
<function lang="NCL">  <!--defining the language of the analytic function-->
    <name>PolarNHEOF</name>  <!--the name of the analytic function-->
    <script>PolarNHEOF.ncl</script>  <!--the file name the script-->
    <description>xxx</description> <!--description of the function-->
    <isGlobalFunction>true</isGlobalFunction>
    <!--defining if the function can be distributed to all the nodes-->
    <MultiInputFiles>false</MultiInputFiles>
    <!--defining if the function can be distributed to all the nodes-->
    <InputFileFormat>netCDF</InputFileFormat> <!--setting format of the input file-->
    <properties>
        <!--defining which datasets can use this function, type can be "include" or "exclude"-->
        <!--the contributor can define model,frequency,modelingRealm and variableName values for filtering-->
        <Model type="include"><!--refer to PCMDI documents to get acceptable values-->
            <value>xxx</value>
            ...<!--multiple values-->
        </Model>
        ... <!--frequency,modelingRealm,variableName-->
    </properties>
```

**properties**

```xml
    <Controls>  <!--setting the parameters that input from the webpage-->
        <parameter>
            <!--defining name,description,tag,type and value range of the parameter-->
        </parameter>
        ... <!--multiple parameters-->
    </Controls>
```

**controls**

```xml
    <InputFileParameters>  <!--setting the parameters about the input file(s)-->
        <parameter>
            <!--defining name,description,tag,type and value rangeof the parameter-->
        </parameter>
        ...
    </InputFileParameters>
    <OutputFileParameters>  <!--setting the parameters about the output file(s)-->
        <result>
            <filetype>xxx</filetype> <!--png/nc/txt...-->
            <filecount>1</filecount>
            <parameter>
                <!--defining name,description,tag,type and value range of the parameter-->
            </parameter>
            ...
        </result>
        ... <!--multiple kinds of results-->
    </OutputFileParameters>
```

**Input/Output related**

```xml
    <pre-processing> <!--defining pre-processing type before invoking the script-->
        <pre-processing-type>Latlon</pre-processing-type> <!--Origin/Latlon/YearAvg/SeasAvg/LTM/Subset-->
        <AddToResults>false</AddToResults>
        <!--defining if the pre-processing result need to be added to the result files-->
    </pre-processing>
</function>
```

**Pre-processing**

Easy way for extensions. Only the analytic script and its XML description are needed.
The script should have an I/O interface for command line and can be invoked by Java.

**R:**

```r
#test.R
    args <- commandArgs(trailingOnly = TRUE)
    date<- as.numeric(args[1])
    filename<- as.character(args[2])
    print(date)
    print(filename)
#command invoking: Rscript test.R 19000101 test.nc
```

- **CAFE nodes:**
  - Data layer:mybatis3.2.3+mysql
  - Service layer: Spring4.0
  - Interaction: Spring MVC4.0+REST API

- **Analytic scripts:**
  - NCL+NCO+CDO

- **Web-based UI:**
  - Service: pHp+yii
  - Database: mysql
  - Web Server: apache2

**Central server**

**P2P**

CAFE Nodes

**Web UI**

- Distributed dataset access and searching
- Parameter setting and task submitting
- Built-in functions
  - EOF analysis( region, NH, SH )
  - Long term mean( region, NH, SH)
  - Trend Analysis( region, NH, SH)
  - Seasonal cycle analysis(NH, SH)
  - Time series ( Annual, Seasonal)
- Task management, result display and downloading

- Compare and analyze sea surface temperature data(tos) among different model output and the observation data

- Model: GFDL-CM3 and IPSL-CM5B-LR

- Observation data:  AVHRR

- Time range: 198501~200512

- Spatial range: 0°N~90°N , 110°E~280°E

- Analysis method: EOF Analysis (Specified regions)

# Prototype and Use Case-use case

- At first, user can search for data by specifying querable attributes. Selected three datasets may be archived on several nodes.

● At first, user can search for data by specifying querable attributes. Selected three datasets may be archived on several nodes.

### Data Files

| institute | model | experiment | modelingRealm | variableName | ensembleMember | temporalStart | temporalEnd | |
|-----------|-------|------------|---------------|--------------|----------------|---------------|-------------|--------|
| IPSL | IPSL-CM5A-LR | historical | ocean | tos | r1i1p1 | 185001 | 200512 | select |
| NCEI | Obs-AVHRR | obs | ocean | tos | global | 198201 | 201412 | select |
| NOAA-GFDL | GFDL-CM3 | historical | ocean | tos | r1i1p1 | 197501 | 200512 | select |

**1**

### Selected

| institute | model | experiment | modelingRealm | variableName | ensembleMember | temporalStart | temporalEnd | |
|-----------|-------|------------|---------------|--------------|----------------|---------------|-------------|----------|
| IPSL | IPSL-CM5A-LR | historical | ocean | tos | r1i1p1 | 185001 | 200512 | unselect |
| NCEI | Obs-AVHRR | obs | ocean | tos | global | 198201 | 201412 | unselect |
| NOAA-GFDL | GFDL-CM3 | historical | ocean | tos | r1i1p1 | 197501 | 200512 | unselect |

Function: EOF Analysis(Specified Region) ▼

### *Temporal Range

Start: 1980 ▼        End: 2005 ▼

### *Spatial Range

North: 90
South: 0
West: 110
East: 280

Task Name: test1

Submit Task

● After submitting the task, user can check status of each task.

● User can obtain the results when the status is *finished*.

- Maps and time series graph are provided.

- Results as NetCDF files and Text files are provided for downloading.

# Conclusion

➢ CAFE can support web-based online batch analysis of distributed environmental data, as well as multi-node collaboration for data intercomparison.

➢ CAFE can dramatically decrease the amount of data need to be transferred from data centers to data users.

➢ CAFE can be easily extended to support more data analytics functions and more data formats.

➢ CAFE is very promising in facilitating overall research efficiency when dealing with large volume of environmental data that are distributedly maintained.

# Future Work

- Look forward to a close collaboration with the development of ESGF

- Support for OPeNDAP, WMS and WPS

- Support for user work space and management of tasks

- Autogeneration of XML descriptions of the analytic functions from user-input web page forms.

- Tracing data provenance and provide intermediate results

- Adding a third-party authentication mechanism and encryption for the APIs between the Web-based UI and CAFE nodes

# THANKS FOR YOUR ATTENTION

**CAFE working team**:

Hao Xu(CESS,xuhao13@mails.tsinghua.edu.cn )

Sha Li(CESS)

Wenhao Dong(CESS)

Shiming Xu(CESS)

Wenyu Huang(CESS)

Yanluan Lin(CESS)

Bin Wang(LASG,CESS)

Fanghua Wu(BCC)

Xiaoge Xin(BCC)

Li Zhang(BCC)

Zaizhi Wang(BCC)

Tongwen Wu(BCC)

Yuqi Bai (CESS, yuqibai@mail.tsinghua.edu.cn )