



# Development and Application of a Simple Model Assessment R Tool (SMART)

Daiwen Kang

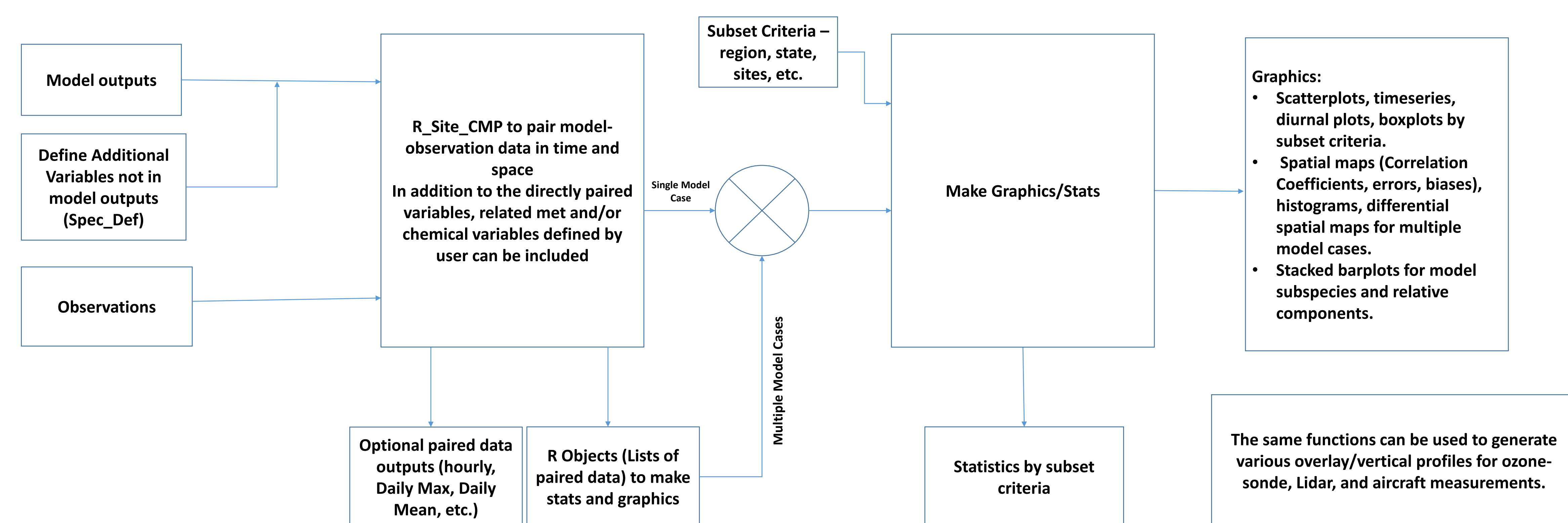
National Exposure Research Laboratory, U.S. Environmental Protection Agency, Research Triangle Park, North Carolina

Daiwen Kang | kang.daiwen@epa.gov | 919-541-4587

## Introduction

As an environment for statistical computing and graphics, the programming language R is widely used and freely available. The EPA's Atmospheric Model Evaluation Tool (AMET) is developed using R to produce statistical analysis and generate various graphics. Complementary to AMET, SMART provides an alternative tool to perform model evaluation for simple projects independent of databases and other languages for pre-processing. With the design of unique data structure and common interfaces to process model outputs and observations, SMART can easily pair any model outputs with observations from common observation networks and special field studies. The large memory requirement and slow processing speed associated with the R platform (the language and running environment) when dealing with large datasets are resolved by using SQLDF and its dependent packages and parallel computing packages. Saved R objects serve as pseudo database tables to perform statistical analysis and generate graphics for various model scenarios. With functionality to subset and merge for the designed data structure, it is easy and fast to generate statistics and graphics for any subsets of the model-observation pairs in time and space. Since R is the only language used for pairing model output with observations and performing statistical analysis, it is completely portable to any computing environment with R installed. To run SMART, except for the shell script that is used to define all the environmental variables (such as input and output directories and variable names), users are not required to have any knowledge of R.

## Design



## Data Structure and Functionality

The primary data structure used in this tool is multiple-level list that saves memory and storage, an example is shown below:

The first element of List of 1990 (1<sup>st</sup> level)

\$ :List of 6 (2<sup>nd</sup> level)

..\$ sid : chr "010030010"

..\$ vars : chr [1:2] "O3" "NOX"

..\$ units : chr [1:2] "ppb" "ppb"

..\$ datetime : chr [1:744] "2011-07-01 00:00:00 GMT" ...

..\$ var.names: chr [1:6] "O3" "O3.obs" "O3.mod" "NOX" ...

..\$ cd :List of 4 (3<sup>rd</sup> level)

...\$ pair.datetime: chr [1:744] "2011-06-30 18:00:00"...

...\$ O3.obs : num [1:744] 64 58 52 52 48 46 NA 18 20 ...

...\$ O3.mod : num [1:744] 68.1 64.7 62.2 59.5 55.2 ...

...\$ NOX.mod : num [1:744] 1.76 2.04 2.17 2.41 2.81 ...

**Functions** to convert list to data frame (list2df), subset a list with criteria (list.select), merge and combine lists (cdata.merge) are implemented in addition to the functions associated with the "rlist" package.

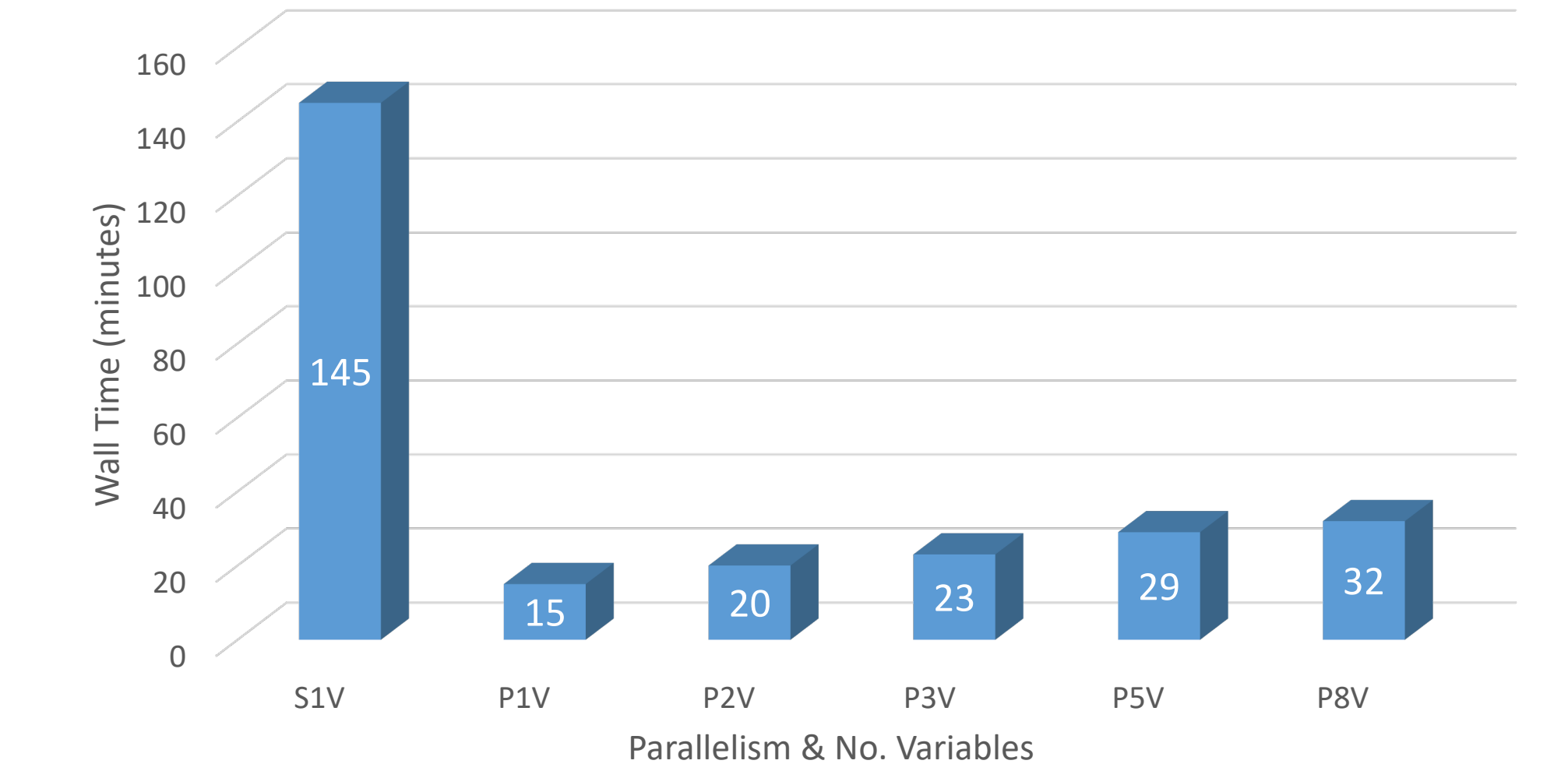
**Common interfaces (polymorphism)** to read in model data (WRF, CMAQ, MPAS, etc) and observational data distinguished by network names/formats (AQS, AIRNOW, CASTNET, CSN, ICARTT, IMPROVE, NADP, SEARCH, Sonde, TOLNet, etc.):

```
getCelldata <- function(site.list, mod.file,
mod.vars, in_units=NULL, llay=1, ulay=1,
tz="GMT", isParallel=F, cfactor=1, ...)
read.OBS <- function(network="AQS", ...)
```

## Memory and Speed

The major drawbacks associated with the base R packages and the computing environment are the limitation of memory and computing speed when dealing with large datasets. For example, the annual AQS hourly dataset can easily exceeds 2 GB that is beyond R's capacity to handle. But with the adoption of SQLDF package, that transparently sets up a database and imports the data frames into that database, performs SQL-like queries to the database using a heuristic method, thus bridges the dataset with the R environment avoiding loading the entire file into memory, R can access any large files efficiently.

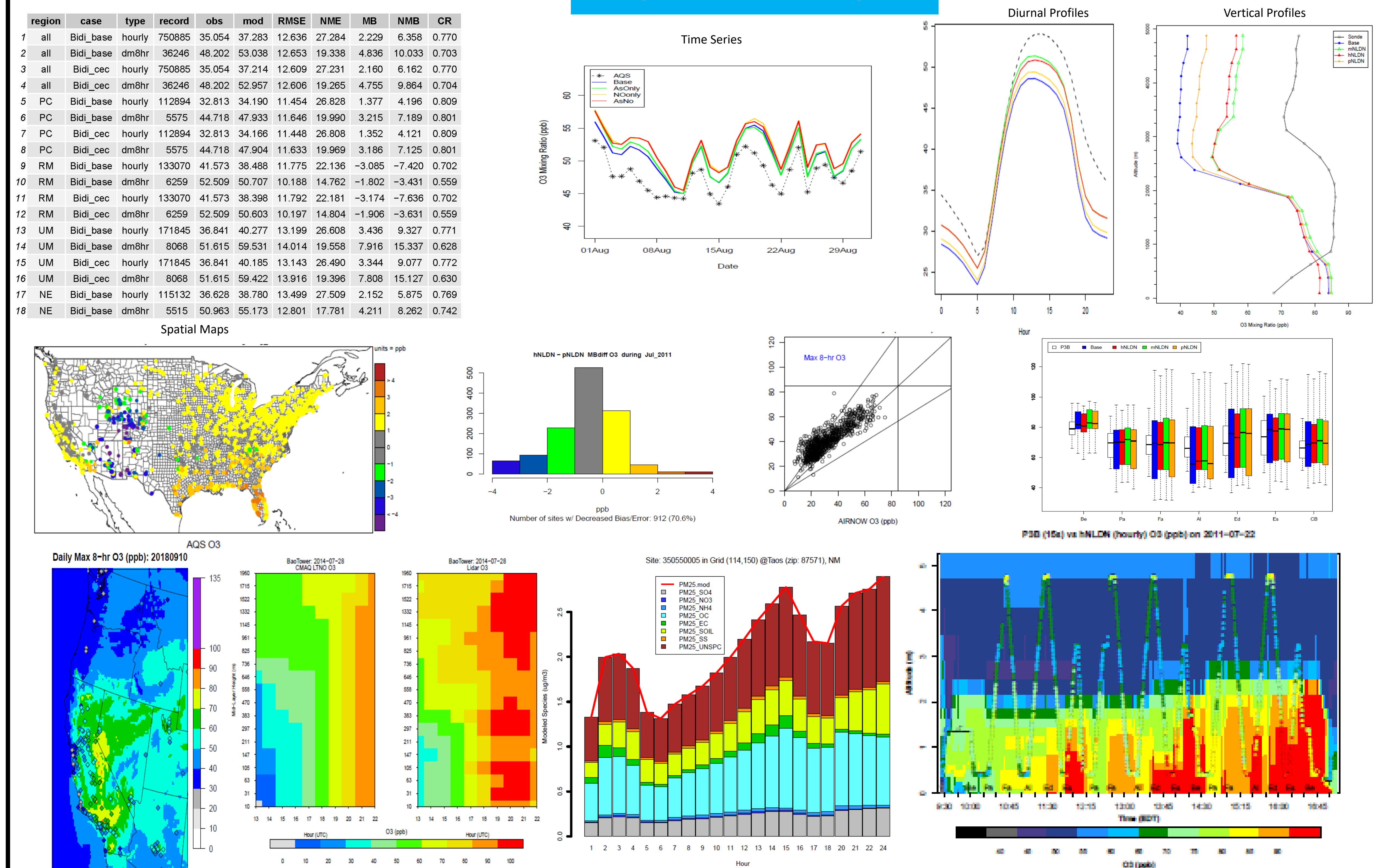
The slow processing speed when dealing with large dataset and complex computations is overcome by using the parallel computing packages in R. As shown in the right figure, during the site compare process to process the entire monthly data for 12 km US CMAQ simulation and AQS observations, when the parallel processing is implemented, dramatic speedup is achieved even with multiple variables.



The running wall time using serialized and parallel processing (16 processors) for different number of variables. S1V: single processor with 1 variable, P1V: multiple processors (16) with 1 variable, P2V: multiple processors with 2 variables, etc.

It takes about 15 to 30 minutes to generate the statistics and the generic graphics for one variable with 2-4 model cases. The time to make the vertical profiles varies with model cases and graphic types, but it ranges from a few minutes to about 1 hour.

## Sample Stats and Graphics



## Disclaimer

The views expressed in this presentation are those of the authors and do not necessarily represent the views or policies of the U.S. Environmental Protection Agency.