

Overview of the Data Assimilation Research Testbed

Draft: 25 April, 2002

Goals:

Data assimilation is the term used in atmospheric and oceanic sciences for the process of merging observations with a model. Data assimilation makes observations more useful by converting diverse and heterogeneous observations to regularly spaced and uniform quantities that can be interpreted more easily. At the same time, observational error can be reduced and information about model errors can be generated. In more advanced applications of data assimilation, models can be improved by confronting them with data. Data assimilation can also be used to perform Observing System Simulation Experiments (OSSE's) which evaluate the impact of existing or proposed observations for particular applications.

The data assimilation problem requires the coordination of expertise in many diverse areas. Model developers, observational specialists, and statisticians trained to do the 'filtering' that is at the core of assimilation algorithms, must all combine expertise to do this problem in an efficient fashion. However, software engineering and organizational practices have made it extremely difficult for experts in these different areas to interact. The result has been that data assimilation development efforts have generally been linked to a single model or observational set. Different assimilation methodologies have not been readily comparable because applying them to a different model has involved too much effort. At the same time, model developers and observationalists have often been tied to a single assimilation methodology without the ability to evaluate the abilities of other methodologies. The result is an inefficient use of resources where every research group is forced to have significant in-house expertise in all three aspects of the problem and is unable to interact well with external groups.

Both the software and organizational barriers to improving this situation can be overcome. Given ongoing community activities to build coordinated software frameworks, it is now possible to design a test-bed facility that would allow the naive combination of a numerical model, a set of observations, and a data assimilation methodology to produce assimilations. The Data Assimilation Research Testbed (DART) is a prototype for an assimilation testbed facility with two fundamental purposes: allowing assimilation algorithm developers to compare their methodologies in a fair way; allowing model developers and observationalists to explore the efficacy of various assimilation algorithms for their problems of interest. DART is designed to demonstrate that a mature test-bed could greatly increase the rate of development of improved data assimilation methodologies, in turn leading to improved datasets, better predictions, and a more efficient design of observational systems.

Implementation:

DART is a software facility that allows a hierarchy of different classes of users to experiment with data assimilation algorithms. At the core of DART is a software infrastructure which is designed to implement standardized interfaces to a wide array of models and observational sets. DART also includes an assortment of data assimilation algorithms built on top of this infrastructure, and a

wide assortment of models, ranging from highly idealized dynamical systems with a handful of variables to General Circulation Models (GCMs) which may be configured with more than a million state variables. A variety of both simulated and actual observation sets are also part of DART. Actual observations are obtained from a variety of sources and converted to a format consistent with the DART software infrastructure. The generation of simulated observation sets (a key component of OSSEs) using models and specified error characteristics for simulated instruments is one of the central capabilities of DART. DART also maintains a variety of auxiliary tools and user interfaces that allow users to experiment, combining assimilation methodologies, models, and observations in ways that can shed new light on their problems of interest.

User views:

An overview of the capabilities of DART can be given by discussing a hierarchy of ‘user views’ of the facility.

View 1: Analysis of previously executed data assimilation experiments:

DART provides an assortment of analysis tools that can give insight into the performance of a particular data assimilation experiment. These tools accept data from three different types of files which provide information about quantities related to observations, model state variables, and the behavior of the assimilation algorithm itself.

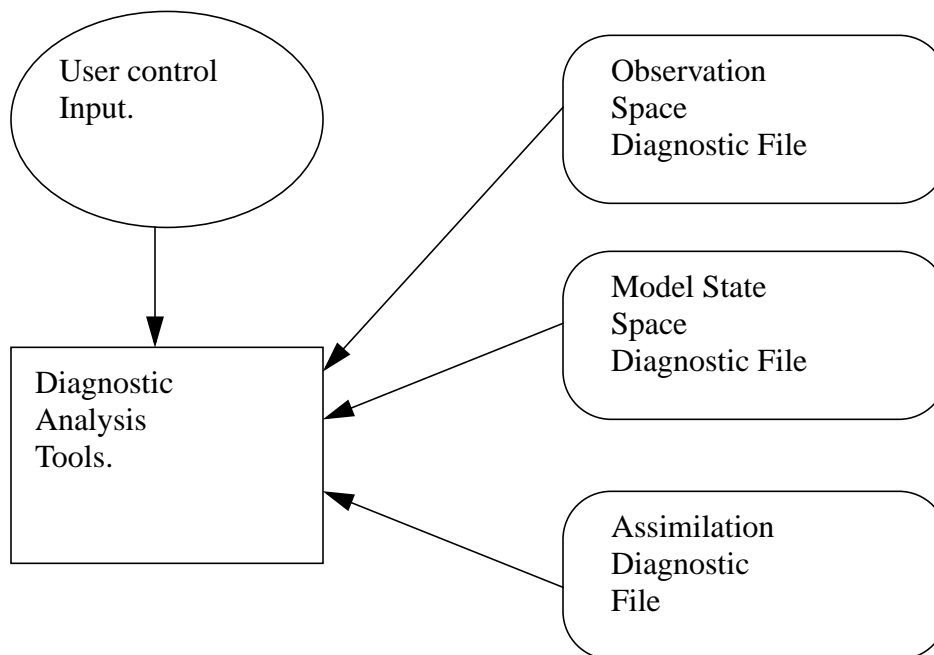
Observation space diagnostic files can provide information about the observing system itself such as the values of observations assimilated, meta-data defining the location and kind of the observations, error covariances of the observations, and, in synthetic observation assimilation experiments, the ‘true’ value of the observed quantity (what would have been measured with a hypothetical instrument with no error). In addition, information from the assimilation algorithm related to observations is also available in these files. These can include information about the (prior and posterior) estimates of the observations generated during an assimilation experiment. Since these (prior and posterior) estimates are formally probability distributions, information could include an estimate of the mean, estimates of higher order moments, a set of samples from the distribution, or a variety of other information. However, all data in observation space diagnostics files must have associated meta-data that describes the observations themselves. Diagnostic tools are available to produce such things as plots of error as a function of time for a repeated observation, overall assessments of error from observations as a function of time or space, etc.

State space diagnostics files provide information about the model state variables (or extended state variables that are functions of the state variables). These files contain meta-data describing the location and kind of model state variables and associated data. For many models and assimilation experiments, state space diagnostic data can be (quasi-)regularly distributed in space and time, but this is not a requirement. Often, it is natural to provide state space output of the (prior and posterior) model state estimate at the times for which observations are available. For instance, one might have the value of the (prior and posterior) model state estimates, some estimates of the error associated with these estimates, the ‘true’ value in synthetic observation experiments, etc. Diagnostic tools are available to produce cross section plots of the space-time state estimate. In a GCM experiment these could include time sections of state estimates at a particular point, plots of

the (error of the) state estimate on a particular vertical or horizontal surface, time series of globally integrated error of a particular field, etc.

Assimilation diagnostic files provide information about the behavior of the assimilation algorithm itself...

View 1 of DART: Analysis of previously existing assimilation experiments



When available should have some examples of diagnostic output, initially from one-dimensional models, here.

View 2: Exploration of Assimilation Experiments

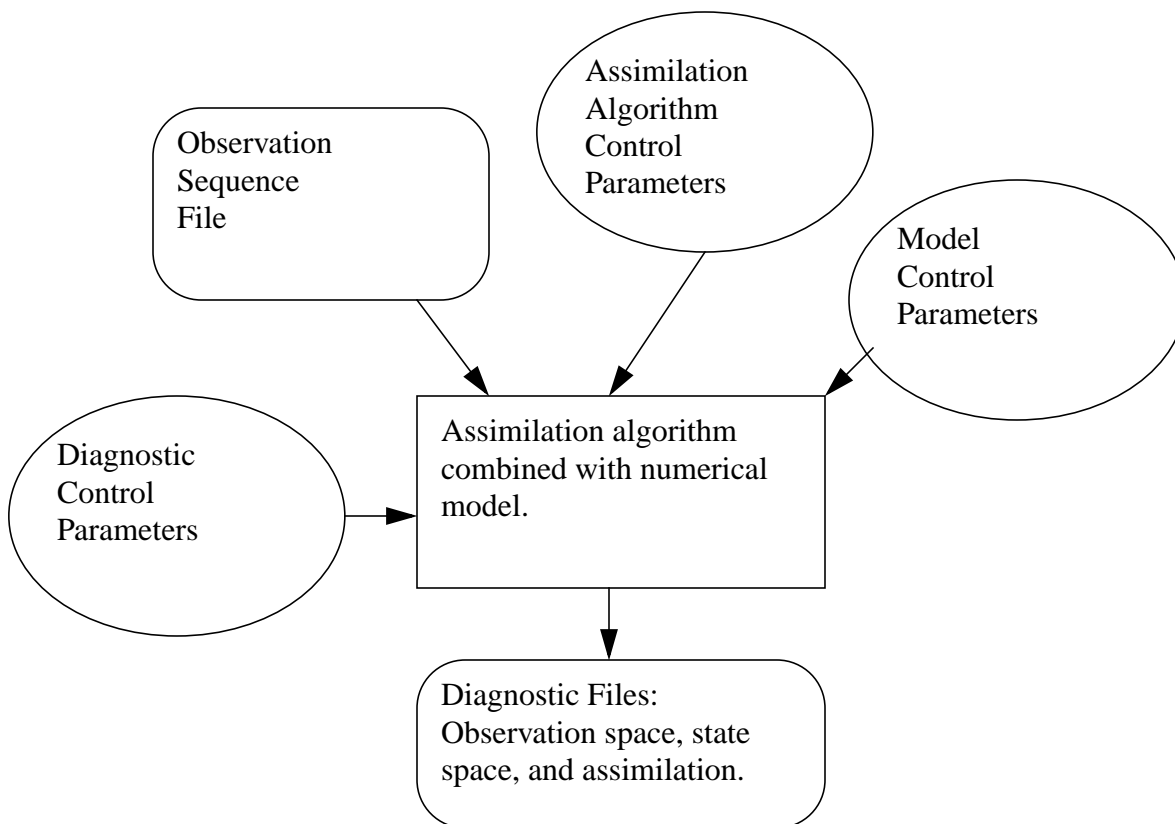
DART allows users to run a previously configured assimilation experiment, a program combining an assimilation algorithm and a numerical model, with modified parameters or different sets of previously defined observations. DART has been designed to make many options of both numerical models and assimilation algorithms available for run-time modification. Users can modify details of the numerical model (for instance the value of a diffusion coefficient in a GCM), details of the assimilation algorithm (for instance, the number of ensemble members used in an ensemble

Kalman filter), or details of the output diagnostics (for instance, selecting some specific subset of model state variables for more intensive diagnostic output).

DART assimilation experiments ingest an observation sequence file which provides meta-data defining the available observations as well as the observations themselves. A variety of pre-defined observation sequence files are available in DART to allow users to explore the impact of different observing systems on an assimilation experiment. The observation sequence files are in the same format as observation space diagnostic files and can be analyzed using the diagnostic tools described in View 1 above.

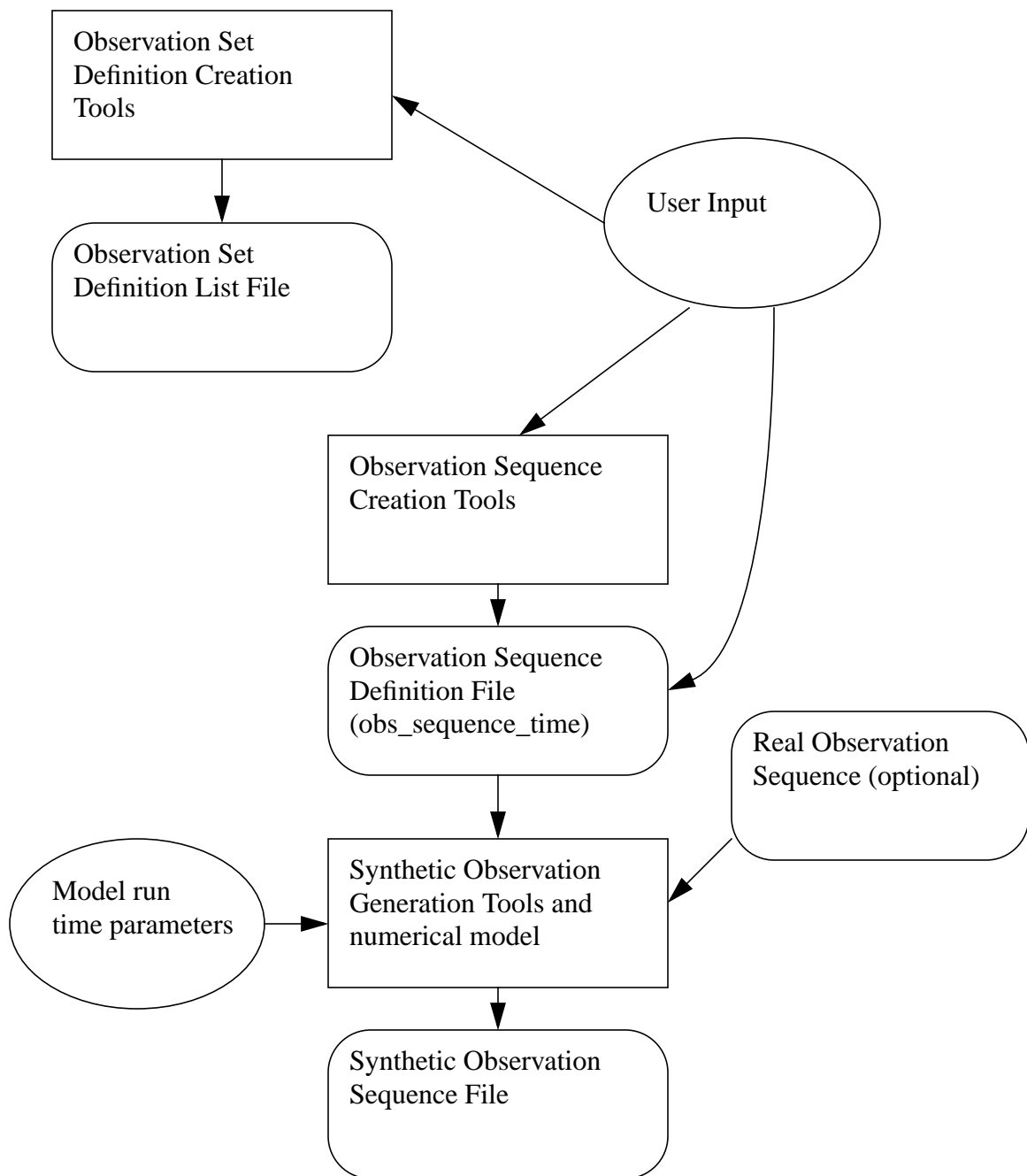
The programs that execute the assimilation experiments are currently implemented as Fortran-90 programs. Users can obtain access to particular configurations by either running previously created executables on DART computing facilities, or by obtaining source code along with configuration management software that allows them to create executables on their own platform.

View 2: Assimilation Experiment Exploration



View 3: Synthetic observation / Observation System Simulation Experiments

Many interesting applications of data assimilation involve the use of synthetic observations, where a model and a specified observing system (including both observation locations and times as well as error characteristics) are used to generate observations that can then be assimilated by the model. The same software that is used to perform assimilations can be used to generate synthetic observations.



DART allows users to specify sets of related observation definitions (a definition here refers to all information about an observation except its time and the actual value of the observation). These observation definitions can then be combined with information about time to generate sequences of observation definitions. A sequence of observation definitions can be in concert with a numerical model to generate synthetic observations using the model integration as input. The results is an obs_sequence with synthetic values which can be used as input to View 2 experiments. Synthetic observations can also be combined with real observations to explore the potential impact of enhancing an existing observing system.

View 4: Adding new models

DART provides a great deal of software infrastructure to assist in modifying existing numerical models for use with data assimilation methods. Nevertheless, adding a model to DART will require Fortran coding by someone who is intimately familiar with the details of the numerical model. Once these interfaces are added, the model should be able to be tested with a variety of assimilation methodologies and to make use of existing observation sequences.

View 5: Adding new observation sets

The tools used to generate synthetic observations are readily adapted to the introduction of new observational data sets into DART. Real observation sets are notoriously difficult to work with, so this process is likely to require effort by someone with expertise in a particular data set. Once the set is available in the format of a DART observational sequence, it can be used with a variety of assimilation techniques and models.

View 6: Adding new assimilation techniques: Although the DART infrastructure is designed to ease this process, it is likely that, at least during the early stages of the project, adding new assimilation algorithms may lead to requirements on the observation or model portions of DART that are unanticipated. This suggests that adding assimilation techniques may require some assistance from model experts in order to allow all models to be fully compliant. However, the promise of being able to make broad and fair comparisons of different assimilation algorithms is likely to encourage assimilation experts to try to place new algorithms in DART.