# Background Fit to Satellite Observations

William F. Campbell

U.S. Naval Research Laboratory, Monterey, CA

## Background and Methodology

### What is it?

- Background fit to observations is a measure of how close, on average, a short forecast (background) is to the entire set of observations valid at the forecast time
- The measure of closeness typically chosen is the **innovation** (observation minus background) **standard deviation**
- Unlike most verification metrics, it is computed in **observation space** rather than model space
- Observations used to measure the background fit *need not be assimilated*, but the observation operator must be computed

### How to Compute

- For each observation (e.g. NPP ATMS channel 4) and for each assimilation window, *measure* how well the **background fit**s all of the observations by computing the *innovation standard deviation*
- Form the **ratio** of innovation standard deviations of the **experiment** to the **control**
- If the **experiment** fits the background **better** than the control, the **ratio** will be **less than 1**; if the **control** fits **better**, the ratio will be **greater than 1**
- Form a **time series** of these **ratios**, and evaluate whether the mean is **statistically significant**ly different from **1** (Student's t-test on the log of the time series with lag-1 autocorrelation correction)
- For that particular observation and time period:
  - If the confidence interval around the mean **includes 1**, the null hypothesis that **experiment** and **control** are indistinguishable cannot be rejected; else
  - If the mean is **less than 1**, the **experiment** is **better** than the **control**
  - If the mean is **greater than 1**, the **control** is **better** than the **experiment**
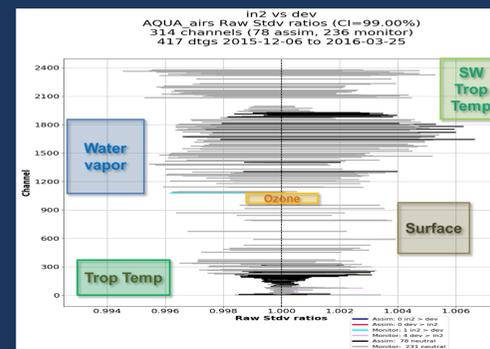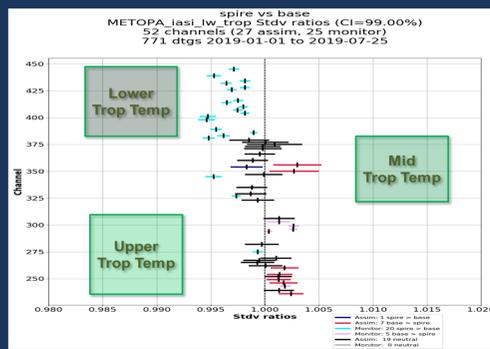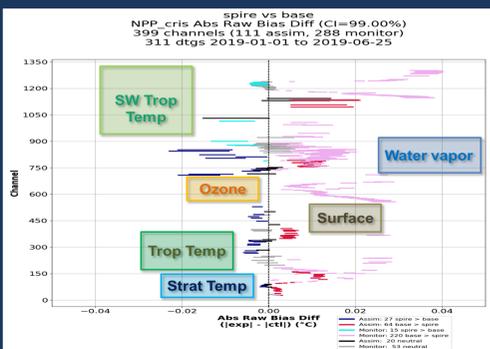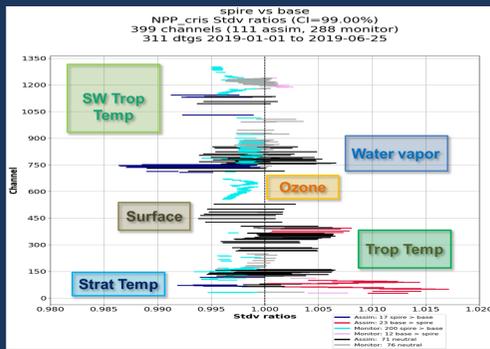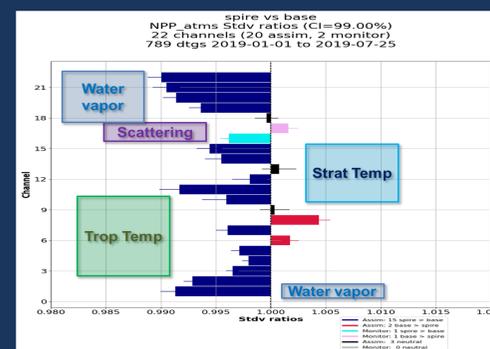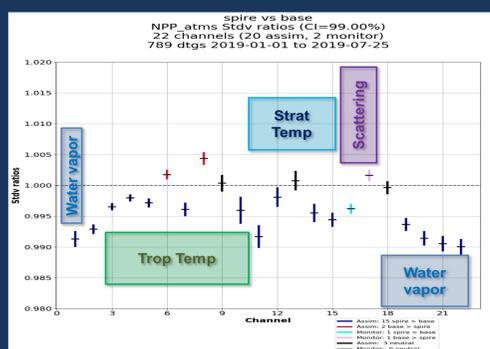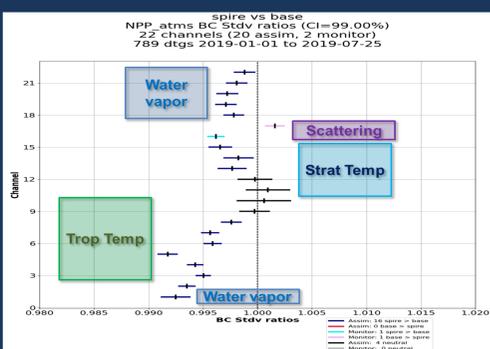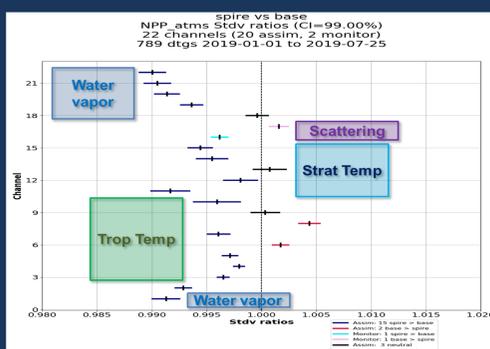
### Verification

- ECMWF uses background fit to observations as their **primary metric** in evaluating short-range (up to 5 days) forecasts[*]
- Decisions on the improvement or degradation of numerical experiments (relative to a control run) can be made **much faster** and (potentially) **more reliably** than with traditional forecast metrics
- **Any type of change** between experiment and control can be evaluated using this method: dynamics, numerics, observations, error covariances, physics, data assimilation, …
- The savings in **computer time** and, more importantly, **researcher time** are substantial

[*]Alan Geer, personal communication

### Diagnostic

- Fit to observations is also a **diagnostic tool**
- Millions of monitored and assimilated **observations** can be **used an additional time** to evaluate the prior forecast in observation space
- We can **aggregate similar observations** in plots and scorecards to facilitate **causal inferences**
  - For example, an experiment improves the background fit to radiances sensitive to moisture in the boundary layer, but degrades the background fit to radiances sensitive to stratospheric temperature
  - Diagnostic information like this provides guidance on how to resolve issues, resulting in an acceleration of the software development cycle

## Graphics Gallery



## Input and Output

### GUI



### Python Snippets



- The input for fit2obs.py consists of a set of ASCII files with the structure and naming convention as follows: {instrument}_fit2obs_{date}.{region}
- In each file are the counts and raw and bias-corrected innovation statistics for each satellite and channel of that instrument for that date.
- The user must compute these statistics and produce a set of input files in a directory structure that looks like {expname}/g{date}
- In each folder will be files for all of the instruments and all of the regions for that date
- Fit2obs.py is called separately for each named region



amsua_fit2obs_2019041900.global

### Detailed Description

By comparing raw innovations to 6-hour forecasts and computing the spatial standard deviation four times a day, we can construct a time series for the control run and each experiment, for every radiance observation type and channel. Each observation type and channel, including unassimilated channels that were monitored, will show either statistical improvement or degradation, or be neutral. The aggregation of these results are in accord with traditional forecast diagnostics that require multiple months of simulated global NWP; however, the fit to observations diagnostic requires only a few weeks of simulation. In addition, it has also been shown to be more sensitive than traditional diagnostics. Up until now, we have lacked this vital tool.

The tool consists of six routines written in Python 3. The main program is in fit2obs.py, and is where the bulk of the data handling is done. Fit2obs.py reads in the radiance statistics from either the *_ar_1.* files or the *_fit2obs_* files for a set of date-time groups for two runs, an experiment and a control. For each sensor and channel, a paired difference t-test is run on the log-transformed raw standard deviations, accounting for serial correlation in time. (Statistics routines are in fit2obs_stats.py). This information is passed to fit2obs_plots.py, which produces a set of graphics that are part of the final output. Fit2obs_scorecard.py also produces a scorecard (described below). Input arguments are handled by fit2obs_arguments.py, which includes an optional gui interface (requires the gooey package).

Fit2obs_plots.py takes a dictionary created by fit2obs.py, either a comparison of an experiment and a control, which is its primary use, or just the control. Graphics options include plot type (bar or line), orientation (vertical or horizontal), image format (png or pdf). Plots can be generated for all channels, or only for assimilated channels. The confidence level defaults to 95%, but can be specified on the command line. When given an experiment and a control dictionary and graphical options, it produces a set of images containing the ratios of raw innovation standard deviations with confidence intervals for each channel. If the confidence interval does not include 1.0, then it is plotted in red (magenta) if the mean ratio is greater than 1.0, indicating that the fit to the assimilated (monitored) channel is better in the experiment than the control. If the mean ratio is less than 1.0 and the confidence interval does not include 1.0, the plot is navy (cyan) for assimilated (monitored) channels. Confidence intervals that include 1.0 are plotted in black (gray) for assimilated (monitored) channels, indicating that the experiment and control fits to this channel are not statistically significant. Fit2obs_plots also generates a set of images containing the mean absolute bias differences with confidence intervals and the same color scheme as the standard deviation ratio plots. Options for plotting bias-corrected standard deviations and absolute differences are also available.

Fit2obs_scorecard.py produces an Excel spreadsheet with summary statistics of wins, losses, and ties for all channels, and for certain subsets (e.g. microwave temperature channels in the stratosphere for all instruments). For groups of channels, we chose a conservative criterion for declaring a group win (win = experiment fits the observations in the group better than the control): If number of individual channel wins exceeds the sum of losses and ties, we declare a group win; if number of individual losses exceed the sum of wins and ties, we declare a group loss; otherwise a group tie.

## Scorecard

- A **scorecard** can provide **summary guidance** as to whether enough evidence has been gathered to favor the experiment over the control or v.v., or whether not enough data has yet been evaluated to decide between them
- As the **time series** of ratios becomes **longer**, confidence intervals get **shorter**, and observations that were neutral will eventually[*] show either a **win** or a **loss**
- We can do this for individual observations, physically sensible groups of observations, and all observations aggregated together
- For groups of observations, we chose a conservative criterion for declaring a group win: If number of individual **wins** exceeds the sum of losses and ties, we declare a **group win**; if number of individual losses exceed the sum of wins and ties, we declare a **group loss**; otherwise a **group tie**

[*]Assumes statistical stationarity, which may not be true. Some side effects may occur. Consult your doctor(ate) before using the scorecard.



## Results, Conclusions, and Future Work

- Evaluating the background fit to observations is a useful **verification** tool in **observation space** that achieves statistical **significance much faster** than traditional verification methods in model space
- Substantial savings in **computer time** and **developer time** can be realized
- As it uses innovations that are already produced in the normal data assimilation process, there is **little overhead needed to compute** and evaluate their statistics
- A lightweight Python implementation **runs in minutes on a laptop for months of data** (satellite radiances only at present)
- Output plots for each satellite, instrument, and user-defined channel group yield useful **diagnostic** information to the researcher
- Output **scorecard** is a useful **summary and decision aid**

- Evaluation of reliability and fidelity to traditional forecast scores
- Fits to bias-corrected innovations, absolute bias
- Add other space-based observations such as AMVs, scatterometers, GPS-RO, Lidar, future sensors
- Regional statistics and scorecards to match traditional verification, or for limited area models
- Longer lead times (e.g. fit to 5-day forecasts), which will project onto model error
- OSE and OSSE applications, e.g. if new observations fit poorly, there may be a problem with the new observations or the forward model rather than the forecast model