ECCO v4 development notes

Gaël Forget Department of Earth, Atmospheric and Planetary Sciences Massachusetts Institute of Technology

June 19, 2015

abstract

These notes pertain to the ECCO v4 state estimate, model setup, and associated codes (Forget et al., 2015). Section 1 points to the other elements of documentation that are available online, and associated download procedures. Section 2 provides guidance to ECCO v4 users interested in operating the ECCO v4 model set-up and/or reproducing the ECCO v4 solution. Section 3 documents the re-implemented estimation modules of MITgcm. Some of the included material in section 3 is expected to eventually move to the MITgcm manual. Throughout this document I try to rely on pre-existing documents rather than duplicating them. Links to pre-existing documents are indicated by blue colored font (e.g. 'manual' in the previous sentence).

Contents

1	downloads	3
	1.1 MITgcm	3
	1.2 ECCO v4 setup	4
	1.3 ECCO v4 solution	4
	1.4 Diagnostic Tools	5
2	MITgcm runs	7
	2.1 regression tests	7
	2.2 Iterative optimization test case	10
	2.3 full ECCO v4 runs	10
3	the generic pkg/ecco and pkg/ctrl	14
	3.1 usage: pkg/ecco	15
	3.2 usage: pkg/ctrl	17
	3.3 implementation: pkg/ecco and pkg/ctrl	19
	3.4 Legacy: pkg/ecco and pkg/ctrl	21

References

Forget, G., J.-M. Campin, P. Heimbach, C. N. Hill, R. M. Ponte, and C. Wunsch, 2015: Ecco version 4: an integrated framework for non-linear inverse modeling and global ocean state estimation. *Geoscientific Model Development Discussions*, 8 (5), 3653–3743, doi:10.5194/ gmdd-8-3653-2015, URL http://www.geosci-model-dev-discuss.net/8/3653/2015/.

1 downloads 1

This section documents locations and directions to download the MITgcm (section 1.1), the 2

ECCO v4 model setup (section 1.2), the ECCO v4 state estimate output (section 1.3), and 3 related diagnostic matlab tools (section 1.4). 4

1.1MITgcm

5

To install the MITgcm: 6

- Go to the MITgcm web-page @ mitgcm.org 7
- Install MITgcm using cvs as explained @ cvs 8
- Run MITgcm using testreport as explained @ manual, howto 9

Pre-requisites are cvs, gcc, gfortran (or alternatives), and mpi (only for parallel runs). For 10 example, my laptop setup, including mpi and netcdf, involved the following mac ports: 11

- cvs @1.11.23_1 (active) 12
- wget @1.14_5+ssl (active) 13
- gcc48 @4.8.2_0 (active) 14
- mpich-default @3.0.4_9+gcc48 (active) 15
- mpich-gcc48 @3.0.4_9+fortran (active) 16
- netcdf @4.3.0_2+dap+netcdf4 (active) 17
- netcdf-fortran @4.2_10+gcc48 (active) 18
- Overridding the default mac gcc and mpich with the above requires: 19
- sudo port select -set gcc mp-gcc48 20
- sudo port select -set mpich mpich-gcc48-fortran 21
- Using mpi and netcdf within MITgcm requires two environment variables: 22
- export MPI_INC_DIR=/opt/local/include 23
- export NETCDF_ROOT=/opt/local 24

²⁵ 1.2 ECCO v4 setup

Any MITgcm user can easily install the ECCO v4 setups using the setup_these_exps.csh shell script as explained @ README. It downloads global_oce_cs32/ (small setup), global_oce_llc90/ (bigger setup) and model inputs from global_oce_input_fields.tar.gz to a subdirectory called global_oce_tmp_download/. The user then wants to move its contents to MITgcm/verification/ (as shown in Fig.1) in order to allow for automated execution of the short benchmark runs via testreport using genmake2 (see section 2.1). Pre-requisites: having downloaded MITgcm (section 1.1) and mpi libraries (only if user wants to run the bigger global_oce_llc90/). The short benchmarks are ran on a daily basis to ensure continued compatibility with the

The short benchmarks are ran on a daily basis to ensure continued compatibility with the up to date MITgcm. While the short benchmarks only go for a few time steps, global_oce_llc90/

 $_{35}$ also is the basic setup that produces the 1992-2011 ECCO v4 ocean state estimate (Forget et al.,

³⁶ 2015) when configured accordingly (as explained in section 2.3). Thus running the short bench-

³⁷ marks (section 2.1) is a useful step towards re-producing the state estimate (section 2.3). It

³⁸ should also be noted that an adjoint version of the short benchmarks also exist that can readily

³⁹ be run by users who access to the TAF compiler.

Figure 1: MITgcm directory structure downloaded using cvs. The ECCO v4 directories indicated with "+" were downloaded separately using setup_these_exps.csh script and moved to MITgcm/verification/.

40 1.3 ECCO v4 solution

⁴¹ The state estimate output for ECCO v4-release 1 is available via this server which is linked to

42 ecco-group.org. The various subdirectories contain monthly fields, this documentation of the solution,

⁴³ in situ and model profiles, the grid specifications and ancillary data as explained in README.docx.

⁴⁴ For example a file (or a subdirectory) can be downloaded at the command line e.g. per

45 wget --recursive ftp://mit.ecco-group.org/ecco_for_las/version_4/release1/README.docx

46 1.4 Diagnostic Tools

To help ECCO v4 and MITgcm users analyze model output obtained either per section 1.3 or per section 2.3, two sets of Matlab tools are made freely available:

• download gcmfaces and MITprof using shell script (or see getting_started.m)

• download MITgcm/utils using cvs (basic functionalities only).

Any user can for example regenerate this documentation of the solution (the gcmfaces 'standard analysis') from the section 1.3 or section 2.3 output (expectedly organized according to Fig.2) simply by executing diags_driver.m¹ and diags_driver_tex.m² in the following sequence :

```
54 dirModel='release1_20150603_c651/';
```

- 55 dirMat='release1_20150603_c651/mat/';
- 56 dirTex='release1_20150603_c651/tex/';
- 57 nameTex='standardAnalysis';
- 58 %
- ⁵⁹ diags_driver(dirModel,dirMat,1992:2011);%requires gcmfaces and MITprof in path
- 60 diags_driver_tex(dirMat,{},dirTex,nameTex);%further requires m_map in path

¹This involves MITprof that also gets installed by this shell script.

²User needs to install m_map for mapping and plotting.

Figure 2: Directory structure as expected by gcmfaces and MITprof toolboxes. The toolboxes themselves can be relocated anywhere as long as their locations are included in the matlab path. Advanced analysis using diags_driver.m and diags_driver_tex.m will respectively generate the mat/ directory (for intermediate computational results) and the tex/ directory (for standard analysis). This diagnostic process relies on the depicted organization of GRID/ and solution/ for automation (user will otherwise be prompted to enter directory names) and depends on downloaded copies of fields to nctiles/ (local subdirectory).

```
gcmfaces/ (matlab toolbox)
 _sample_input/ (binary files)
 __ @gcmfaces/ (matlab codes)
 __gcmfaces_calc/ (matlab codes)
 . . .
MITprof/ (matlab toolbox)
 _profiles_samples/ (netcdf files)
  _profiles_process_main_v2/ (matlab codes)
  _profiles_stats/ (matlab codes)
   . . .
GRID/ (binary output)
release 1 solution/
_____diags/ (binary output)
 __nctiles/ (netcdf output)
 __MITprof/ (netcdf output)
 __mat/ (created by gcmfaces)
___tex/ (created by gcmfaces)
other solution/
 _diags/ (binary output)
 _ . . .
. . .
```

61 2 MITgcm runs

The following procedures, commands and submission scripts allow runs of the ECCO v4 MITgcm setup – either in short regression tests (section 2.1) or for multi-decadal simulations such as the full 20 year state estimate (section 2.3). Pre-requisite for sections 2.1 and 2.3: having downloaded the MITgcm (section 1.1) and the ECCO v4 setups (section 1.2). Pre-requisite for section 2.3: having downloaded forcing fields and a few other binary model inputs (listed below).

67 2.1 regression tests

Short benchmarks of the MITgcm and ECCO v4 setup are run using testreport command line
utility (see Fig.2; howto). Serial runs are executed simply at the command line e.g. per

```
70 ./testreport -t global_oce_cs32
```

71 OT

```
72 ./testreport -skipdir global_oce_llc90
```

The reader is referred to 'testreport –help' and how to for additional explanation about such 73 commands. If everything proceeds as expected then the result of the comparison with the 74 reference result is reported to screen as shown in abbreviated form in Fig. 3. Depending on 75 your machine environment the agreement with the reference result may be lower in which case 76 'testreport' may indicate 'FAIL' (e.g. see README). Despite the dramatic character of such 77 message, this is generally ok and does not prevent reproducing full model solutions accurately 78 (see section 2.3). If the testreport process gets interrupted then it is often safer to clean up 79 experiment directories (e.g., by executing ./testreport -clean -t global_oce_*) and start over. 80

```
----T-----
default 10
                          ----S-
GDM
          с
                    m
                       S
                                 m
                                    S
 раR
          g
                    е
                          m
             m
                m
                       .
                              m
                                 е
nnku
          2
             i
                       d
                          i
                а
                              a
                                 а
                                    d
                    а
          d
2 d e
      n
             n
                х
                    n
                          n
                              x
                                 n
Y Y Y Y>14<16 16 16 16 16 16 16 16 16
                                        pass
                                             global_oce_cs32
```

Figure 3: Abbreviated output of testreport to screen.

The above 'testreport' commands deserve a couple more specific comments. The first com-81 mand runs the global_oce_cs32/ benchmark solely. The second command will run all MITgcm 82 benchmarks including global_oce_cs32/ but not the global_oce_llc90/ benchmark that requires 83 at least 12 processors in forward (96 in adjoint) and therefore should not be run in serial mode 84 (doing so may crash your laptop). It is thus excluded by using the 'skipdir' option. It should 85 be stressed however that $global_oce_cs32/$ depends on the files in $global_oce_llc90/$ (which is the 86 main setup) rather than duplicating them. Therefore global_oce_llc90/ must not be removed 87 from MITgcm/verification for global_oce_cs32/ to work. 88

⁸⁹ Running the short benchmarks with mpi (assuming it has been installed) is equally simple:

90 ./testreport -of ../tools/build_options/linux_amd64_ifort+mpi_ice_nas \
91 -j 4 -MPI 96 -command 'mpiexec -np TR_NPROC ./mitgcmuv' \
92 -t global_oce_llc90

⁹³ for example will run the first forward benchmark of global_oce_llc90/ on 96 processors using an ⁹⁴ ifort compiler. Note that the specifics (number of processors and compiler choice) are to be ⁹⁵ determined by the user and are machine dependent.

Often in massively parallel computing environments, it is common that mpi jobs can only be run within a queuing system. The submission script in Fig.4 (that is also machine specific) provides an example on how to do it. It contains 3 hard-coded switches : fwdORad = 1 (2 for adjoint); numExp = 1 (2 for llc90); excludeMpi = 0 (1 for serial). This script should be located and submitted from MITgcm/verification. It is also common that compute nodes cannot access certain compilers, in which case the user may want to proceed in two steps:

102 1. compile outside of the queuing system using e.g. per

```
./testreport -of ../tools/build_options/linux_amd64_ifort+mpi_ice_nas \
    -j 4 -MPI 96 -command 'mpiexec -np TR_NPROC ./mitgcmuv' \
    -t global_oce_llc90 -norun
```

```
2. submit the Fig.4 script, after adding -q to the 'opt' variable to skip compilation.
```

Running adjoint benchmarks requires access to the TAF compiler. The calls to testreport 107 (see above) then only need to be slightly altered by appending the '-ad' option (for either 108 serial or mpi jobs) and replacing 'mitgcmuv' with 'mitgcmuv_ad' (only for mpi jobs). It should 109 also be noted that, unlike other MITgcm benchmarks, global_oce_cs32/ and global_oce_llc90/ 110 do not include any adjoint specific 'code_ad/' directory as they simply use the forward model 111 'code/' directory instead. Since testreport relies on the existence of 'code_ad/' for its adjoint 112 option though, it is necessary to soft link 'code/' to 'code_ad/' in both global_oce_cs32/ and 113 global_oce_llc90/ accordingly in order to to run their 'testreport -ad' versions. 114

Figure 4: Example script to run mpi testreport via a queueing system (machine dependent).

```
#PBS -S /bin/csh
#PBS -1 select=1:ncpus=16:model=ivy+4:ncpus=20:model=ivy
#PBS -1 walltime=02:00:00
#PBS -q devel
#PBS -m n
#environment variables and libraries
limit stacksize unlimited
module purge
module load modules comp-intel/2013.1.117 mpi-sgi/mpt.2.10r6 netcdf/4.0
#
setenv LD_LIBRARY_PATH ${LD_LIBRARY_PATH}:${HOME}/lib
setenv MPI_IB_TIMEOUT 20
setenv MPI_IB_RAILS 2
setenv MPI_IB_FAILOVER 1
setenv MPI_CONNECTIONS_THRESHOLD 2049
#local variables and commands
#-----
set fwdORad = 1
set numExp = 1
set excludeMpi = 0
#
if ( \{numExp\} == 1\} then
 set nameExp = global_oce_cs32
 set NBproc = 6
else
 set nameExp = global_oce_llc90
 set NBproc = 96
endif
#
if ( ${excludeMpi} == 1 ) then
 set opt = '-of ../tools/build_options/linux_amd64_ifort -j 4'
else
 set opt = '-of ../tools/build_options/linux_amd64_ifort+mpi_ice_nas -j 4'
endif
#
if ( ${fwdORad} == 1 && ${excludeMpi} == 0 ) then
  ./testreport ${opt} -MPI \
  ${NBproc} -command 'mpiexec -np TR_NPROC ./mitgcmuv' -t ${nameExp}
else if ( fwdORad == 2 && fexcludeMpi == 0 ) then
 ./testreport ${opt} -MPI \
 ${NBproc} -command 'mpiexec -np TR_NPROC ./mitgcmuv_ad' -ad -t ${nameExp}
else if ( ${fwdORad} == 1 && ${excludeMpi} == 1 ) then
  ./testreport ${opt} -t ${nameExp}
else if ( ${fwdORad} == 2 && ${excludeMpi} == 1 ) then
  ./testreport ${opt} -ad -t ${nameExp}
endif
exit
                                        9
```

¹¹⁵ 2.2 Iterative optimization test case

The global_oce_cs32/input_OI implements an iterative optimization test case. It case boils down to optimal interpolation (the model dynamics are not involved) solved by a variational method using the MITgcm adjoint (a diffusion equation in this test case). The pre-requisites are:

119 1. run the adjoint benchmark in global_oce_cs32/ via testreport (see section 2.1).

- 2. Go to MITgcm/lsopt and compile (see section 3.18 of manual).
- 3. Go to MITgcm/optim, replace 'natl_box_adjoint' with 'global_oce_cs32' in this Makefile,
 and compile as explained in section 3.18 of manual. An executable named 'optim.x' should
 get created in MITgcm/optim. If otherwise, please contact ecco-support@mit.edu
- 4. go to MITgcm/verification/global_oce_cs32/input_OI and type 'source ./prepare_run'
- ¹²⁵ Then the iterative optimization itself proceeds as follows
- 126 1. ./mitgcmuv_ad > output.txt
- 127 2. ./optim.x > op.txt
- 128 3. increment optimcycle by 1 in data.optim
- 4. go back to step #1, to run the next iteration
- ¹³⁰ 5. type 'grep fc costfunction000*' to display results (Fig. 5).

costfunction0000: fc =	4118.1987222194211	0.0000000
costfunction0001: fc =	1523.9310891186672	0.0000000
costfunction0002: fc =	1053.3611790049420	0.0000000
costfunction0003: fc =	790.10479375339185	0.0000000

Figure 5: Results of iterative optimization after 3 iterations carried out as explained in section 2.2.

131 2.3 full ECCO v4 runs

The 1992-2011 ECCO v4 ocean state estimate (Forget et al., 2015) is reproduced on a monhtly basis to ensure continued compatibility with the up to date MITgcm. Re-running the baseline 20 year solution (or any other 20 year of global_oce_llc90/) on 96 processors may take about 8 to 12 hours (depending on the computing environment). Reproducing the state estimate requires additional input to be downloaded (besides section 1.2; see below). Unlike for the short benchmarks of section 2.1, in the case of these longer model runs:

- the model is compiled and run outside of testreport.
- the model is compiled with compiler optimization.

• additional forcing and binary input is necessary.

• additional memory and/or disk space is necessary.

The reader is referred to how for a general explanation of such practice. The typical compi-142 lation sequence for the ECCO v4 forward model (i.e. the model setup in global_oce_llc90/) is 143 shown in Fig.6. The tamc.h_itXX and profiles.h_itXX headers (see Fig.6) allow for additional 144 time steps and in situ profiles input, respectively. Once done with compilation, the user typically 145 creates and enters a run directory, links the model executable and inputs into place (see Fig.7), 146 and submits a job to the queueing system (see Fig.8). The 'input_itXX/prepare_run' script 147 (Fig.7) makes a local copy of the model executable ('mitgcmuv'), of all namelists ('data*' from 148 the various 'input*/' directories downloaded in section 1.2) that set up the model at run time, 149 and links a few binary inputs such as the grid and bathymetry files. Assuming that the forcing, 150 observations, model parameter adjustments and initial condition directories were populated and 151 the script was edited (user need to specify forcingDir, obsDir, ctrlDir and pickDir in Fig.7) 152 accordingly then it will furthermore link their contents in the run directory. Users interested in 153 obtaining the necessary input files are advised to contact ecco-support@mit.edu regarding: 154

- 6 hourly forcing files over 1992-2011 (EIG*199? EIG*20??).
- insitu data sets (*feb2013*.nc) used in long benchmark (see below).
- model parameter adjustsments, a.k.a. the control vector (xx_*).
- the initial conditions (pickup*)

Once the model run has completed, one wants to verify that it accurately reproduces the 159 reference result – or detect that a mistake was made. To this end, a mechanism that is analogous 160 to testreport but is geared towards benchmarking long runs was introduced by Forget et al. 161 (2015). It is operated by testreport_ecco.m within Matlab. The pre-requisite is to add the 162 reference result directory 'MITgcm/verification/global_oce_llc90/results_itXX/' to the Matlab 163 path. As explained in Forget et al. (2015) testreport_ecco.m compares time series of global 164 mean variables, and other characteristics of the solution, to the reference state estimate values. 165 The array of tests can be extended to e.g. meridional transports by adding gcmfaces to the 166 Matlab path. The typical call sequence is indicated in the help of testreport_ecco.m and in 167 Fig.9 that also illustrates the typical display of the benchmarking results report to the user 168 screen. The expected level of accuracy for re-runs of the baseline 20 year solution (with an up 169 to date MITgcm code on any given computer) is reached when the displayed values are < -4170 (see Forget et al., 2015, for details). In cases when some of the tests were omitted (e.g. because 171 gcmfaces was not in the Matlab path) the display will show NaN for omitted tests. From the 172 generated model output, one may further easily compute and display many diagnostic quantities 173 using the gcmfaces standard analysis for example (see section 1.4). 174

Figure 6: Compilation directives, outside testreport, for intensive model runs. On a different machine (computer) another build option file such as linux_amd64_gfortran or linux_amd64_ifort11 should be used. To compile the adjoint, users need a TAF license and to replace 'make -j 4' with 'make adall -j 4'. Note : the '-mods=../code' specification can be omitted if the build directory contains the 'genmake_local' file).

```
cd verification/global_oce_llc90/build
../../../tools/genmake2 -optfile=\\
../../tools/build_options/linux_amd64_ifort+mpi_ice_nas -mpi -mods=../code
make depend
\rm tamc.h profiles.h
cp ../code/tamc.h_itXX tamc.h
cp ../code/profiles.h_itXX profiles.h
make -j 4
```

Figure 7: Example script to setup the 20 year ECCO v4 state estimate. It is implied that user has filled directories /bla, /blaa, /blaaa and /blaaa with appropriate forcing, observational, control vector, and pickup files.

```
#!/bin/csh -f
set forcingDir = ~/bla
                = ~/blaa
set obsDir
set ctrlDir
                = ~/blaaa
                = ~/blaaaa
set pickDir
source ../input_itXX/prepare_run
cp ../build/mitgcmuv .
\rm pick*ta EIG*
ln -s ${forcingDir}/EIG* .
ln -s ${obsDir}/* .
ln -s ${ctrlDir}/xx* .
ln -s ${pickDir}/pick* .
exit
```

Figure 8: Example script to run the 20 year ECCO v4 state estimate on 96 processors (machine dependent).

```
PBS -S /bin/csh
#PBS -1 select=1:ncpus=16:model=ivy+4:ncpus=20:model=ivy
#PBS -1 walltime=12:00:00
#PBS -q long
#environment variables and libraries
limit stacksize unlimited
module purge
module load modules comp-intel/2013.1.117 mpi-sgi/mpt.2.10r6 netcdf/4.0
#
setenv LD_LIBRARY_PATH ${LD_LIBRARY_PATH}:${HOME}/lib
setenv MPI_IB_TIMEOUT 20
setenv MPI_IB_RAILS 2
setenv MPI_IB_FAILOVER 1
setenv MPI_CONNECTIONS_THRESHOLD 2049
#run MITgcm
#-----
mpiexec -np 96 dplace -s1 ./mitgcmuv
exit
```

Figure 9: Calling sequence to be executed form within matlab to verify that their re-run of the 20 year ECO v4 state estimate is acceptably close to the released state estimate.

¹⁷⁵ 3 the generic pkg/ecco and pkg/ctrl

176 State estimation consists in minimizing a least squares distance, J(u), that is defined as

$$\mathbf{J}(\mathbf{\mathfrak{u}}) = \sum_{i} \alpha_{i} \times (\mathbf{d}_{i}^{T} \mathbf{R}_{i}^{-1} \mathbf{d}_{i}) + \sum_{j} \beta_{j} \times (\mathbf{\mathfrak{u}}_{j}^{T} \mathbf{\mathfrak{u}}_{j})$$
(1)

$$\mathbf{d}_i = \mathcal{P}(\mathbf{m}_i - \mathbf{o}_i) \tag{2}$$

$$\mathbf{m}_i = \mathcal{SDM}(\mathbf{v}) \tag{3}$$

$$\mathfrak{v} = \mathcal{Q}(\mathfrak{u}) \tag{4}$$

$$\mathfrak{u} = \mathcal{R}(\mathfrak{u}') \tag{5}$$

where d_i denotes a set of model-data differences, α_i the corresponding multiplier, $\mathbf{R_i}^{-1}$ the 177 corresponding weights, u_j a set of non-dimensional controls (of adjustable model parameters), 178 β_i the corresponding multiplier, and additional symbols appearing in Eqs. 2-5 are defined below. 179 The generic implementation of Eqs.1-5 and the adjoint interface within the MITgcm is 180 charted in Fig. 10. A basic presentation of Eqs. 1-5 and Fig. 10 can be found in Forget et al. 181 (2015). Details of the implementation within the MITgcm 'pkg/ecco' and 'pkg/ctrl' (a concern 182 for developers mainly) are provided later in sections 3.3 and 3.4. Most importantly, sections 3.1 183 and 3.2 document the generic features in 'pkg/ecco' and 'pkg/ctrl' and their practical application. 184 The presented features are tested daily via $global_oce_cs32/$ (section 2.1; adjoint experiment) 185 and tested monthly in real-life conditions via the full ECCO v4 run (section 2.3; forward run), 186 which will also serve for illustration in this section. 187



Figure 10: Chart of the organization and roles of MITgcm estimation modules. Additional details are reported in the MITgcm manual, in Forget et al. (2015), and in section 3 of this document.

188 3.1 usage: pkg/ecco

Model counterparts (m_i) to observational data (o_i) derive from adjustable model parameters (\mathfrak{v} ; see section 3.2) through the model dynamics (\mathcal{M} ; see Forget et al. 2015), diagnostic computations (\mathcal{D}), and averaging (or subsampling in 'pkg/profiles') in space and time (\mathcal{S}). For each cost function term the underlying uncertainty field ($\sqrt{\mathbf{R}_i}$) is specified by 'gencost_errfile'.³ The corresponding cost function multiplier (α_i) is specified by 'mult_gencost' (it is 1. by default).

The file name for the observational fields (o_i) is specified by 'gencost_datafile'. Normally 194 o_i (and m_i accordingly) is a time series of daily or monthly averages as specified by 'gen-195 cost_avgperiod'. In principle any periodicity should be possible but only 'month', 'day', 'step' 196 and 'const' are implemented. The observational time series may be split in yearly files finishing in 197 e.g. '_1992', '_1993', etc. Dense time series of model time steps can also be employed for testing 198 purposes (e.g. in this data.ecco). Climatologies of m_i can be formed from its time series to com-199 pare with observational o_i climatologies. This option is activated by the gencost_preproc='clim' 200 specification as illustrated in this data.ecco. Finally the gencost_avgperiod='const' option is 201 adequate when m_i is constant through time once the model initialization phase is complete.⁴ 202 Plain model-data misfits $(m_i - o_i)$ can be penalized directly (i.e. used in Eq. 1 in place of d_i). 203 More generally though penalized misfits (d_i in Eq. 1) derive from $m_i - o_i$ through a generic 204 post-processor (\mathcal{P} in Eq. 2). They can thus be smoothed in space at run time by setting 205 gencost_posproc='smooth' for example (see this data.ecco). 206

The physical variable in m_i is specified at run time via the first characters in 'gencost_barfile' 207 (to match the observed variable specified as o_i) as illustrated in this data.ecco and that data.ecco. 208 The list of implemented variables as of the MITgcm checkpoint c65m is reported in Tab. 1. In 209 cases when two different averages of the same variable may be needed in separate cost function 210 terms (e.g. daily and monthly) or simply for convenience then an extension starting with '-' 211 can be added to 'gencost_barfile' (such as '_day' and '_mon'). In cases when two cost function 212 terms may use the same m_i , the user may specify the same name (via 'gencost_barfile') in both 213 terms. In cases of three dimensional variables (see Tab. 1) the 'gencost_is3d' run-time option is 214 automatically set to .TRUE. (it is .FALSE. by default). The gencost_outputlevel=1 option will 215 output model-data misfit fields for offline analysis and visualization. 216

³The option for time varying error fields remains to be implemented in gencost.

⁴This feature remains to be added to daily benchmark.

Table 1: List of implemented 'gencost_barfile' options that can be specified at run-time in data.ecco (as of the MITgcm checkpoint c65m). An extension starting with '_' can be appended at the end of the variable names for convenience.

variable name	description	remarks
m_eta	sea surface height	free surface $+$ corrections
m_sst	sea surface temperature	first level temperature
m_sss	sea surface salinity	first level salinity
m_bp	bottom pressure	
m_ustress	zonal wind stress	
m_vstress	meridional wind stress	
m_uwind	zonal wind	
m_vwind	meridional wind	
m_atemp	atmospheric temperature	
m_aqh	atmospheric humidity	
m_{precip}	precipitation	
m_swdown	downward shortwave	
m_lwdown	downward longwave	
m_wspeed	wind speed	
m_siarea	sea-ice concentration	
m_siheff	sea-ice effective thickness	
$m_sihsnow$	snow effective thickness	
m_theta	temperature	three-dimensional
m_salt	salinity	three-dimensional
m_diffkr	diapycnal diffusion	three-dimensional, constant
m_kapgm	bolus velocity parameter	three-dimensional, constant
m_kapredi	isopycnal diffusion	three-dimensional, constant
$m_{geothermalflux}$	geothermal heating	'const'
$m_bottomdrag$	bottom drag	'const'

217 3.2 usage: pkg/ctrl

Three basic options are implemented for Eqs. 4-5: time variable two dimensional controls ('gentim2d'), time-invariant 2D controls ('genarr2d'), and time-invariant 3D controls ('genarr3d'). The 'gentim2d' run-time options are documented below as an example. Corresponding options exist in 'genarr2d' and 'genarr3d' except for the specifically time variable aspects (see below).

The control problem is non-dimensional by default, as reflected by the omission of weights

in control penalties $(\mathfrak{u}_j^T\mathfrak{u}_j, \text{Eq.1})$. Non-dimensional controls (\mathfrak{u}_j) are scaled to physical units (\mathfrak{v}_j)

through multiplication by their respective uncertainty fields $(\sigma_{\mathfrak{u}_j})$, as part of the generic preprocessor Q (Eq.4). An adjustable parameter are activated and specified by the first character

²²⁶ in 'xx_gentim2d_file' (as illustrated in this data.ctrl and that data.ctrl). The list of implemented

variables as of the MITgcm checkpoint c65m is reported in Tab. 2.

Table 2: List of implemented 'xx_gen????_file' (with ????? indicated under remarks) options that can be specified at run-time in data.ctrl (as of the MITgcm checkpoint c65m). An extension starting with '_' can be appended at the end of the variable names for convenience.

variable name	description	remarks
'xx_atemp'	atmospheric temperature	'gentim2d'
'xx_aqh'	atmospheric humidity	'gentim2d'
'xx_swdown'	downward shortwave	'gentim2d'
'xx_lwdown'	downward longwave	'gentim2d'
'xx_precip'	precipitation	'gentim2d'
'xx_uwind'	zonal wind	'gentim2d'
'xx_vwind'	meridional wind	'gentim2d'
'xx_tauu'	zonal wind stress	'gentim2d'
'xx_tauv'	meridional wind stress	'gentim2d'
'xx_etan'	initial free surface height	'genarr2d'
'xx_theta'	initial temperature	'genarr3d'
'xx_salt'	initial salinity	'genarr3d'
'xx_diffkr'	diapycnal diffusion	'genarr3d'
'xx_kapgm'	bolus velocity parameter	'genarr3d'
'xx_kapredi'	isopycnal diffusion	'genarr3d'
'xx_geothermal'	geothermal heating	'genarr2d'
'xx_bottomdrag'	bottom drag	'genarr2d'

The corresponding uncertainty must be provided (in the form of weights $1/\sigma_{u_j}^2$; via a file name specified by 'xx_gentim2d_weight') to scale u_j to physical units.⁵ Besides the scaling of u_j to physical units, generic pre-processor Q can include for example spatial correlation modeling (using an implementation of Weaver and Courtier, 2001). This feature is activated for e.g. the first set of controls by setting xx_gentim2d_preproc(1)='WC01'.⁶ As an alternative, one may set xx_gentim2d_preproc(1)='smooth' to apply the smoothing part of Weaver

⁵Options to specify an uncertainty field or constant instead remain to be implemented.

⁶The ctrlSmoothCorrel3D/2D switches and CPPs remain to be fully deprecated.

and Courtier, 2001 but omit the normalization part. Additional specification is possible in some cases (depending on 'xx_gentim2d_preproc') via sub-options 'xx_gentim2d_preproc_i' (integer), '..._r' (real), '..._c' (character string). For example, setting xx_gentim2d_preproc_i(1,1)=2 along with xx_gentim2d_preproc(1)='WC01' would apply the second correlation model defined in data.smooth to the first set of controls (the first correlation model would otherwise be used by default). The full list of implemented 'xx_gentim2d_preproc' options (as of the MITgcm checkpoint c65m) is reported in Tab. 3.

Table 3: List of implemented 'gentim2d' options (top) and associated 'xx_gentim2d_preproc' options (bottom) that can be specified at run-time in data.ctrl (as of the MITgcm checkpoint c65m).

	parmater name	type		role
	$xx_gentim2d_file$	character	activate	an adjustable parameter
	xx_gentim2d_weight	character	$^{\mathrm{sp}}$	ecify weight field(s)
	$xx_gentim2d_prepro$	c character	option	al features listed below
	xx_gentim2d_bound	s real (five values)	impose b	ounds (see that data.ctrl)
	mult_gentim2d	real	cost function	on multiplier (1. by default)
	gentim2dPrecond	real	precon	ditioner (1. by default)
xx_gentim2d_preproc		further specification	s (see text)	effect (in forward)
'WC01'		xx_gentim2d_preproc_i		activate correlation modeling
'smooth'		xx_gentim2d_preproc_i		activate plain smoothing
'docycle'		xx_gentim2d_preproc_i		average period replication
'rmcycle'		xx_gentim2d_preproc_i		periodic average subtraction
'variaweight'		(none)		time variable weights

In the case of time-variable parameter adjustments, the frequency is specificied by 'xx_gentim2d_period'.

Time variable weights can also be provided by specifying 'variaweight' as e.g. ' $xx_gentim2d_preproc(2)$ '.

In this case the 'xx_gentim2d_weight' file must contain as many records as the control parameter time series itself (\approx the duration of the run divided by 'xx_gentim2d_period'). In the case when several adjustments are sought in one model parameter (e.g. time mean and time variable forcing adjustments treated separetely) then an extension starting with '_' can be added to the 'xx_gentim2d_file' specification (e.g. '_mean' and '_anom'; see this data.ctrl).

Further time-variable ('gentim2d' only) options are available via xx_gentim2d_preproc= 248 'docycle' and 'rmcycle'. They can be combined with 'variaweight' (that occurs after 'docycle' 249 and 'rmcycle') to create many 'gentim2d' varieties. The example in this data.ctrl specifies that 250 adjustments to atmospheric temperatures are split in three terms: time mean ('xx_atempA'), sea-251 sonal cycle anomaly ('xx_atempB'; of zero time mean), and interannual anomaly ('xx_atempC'; 252 of zero time mean and *seasonal cycle*). In a real-life situation, a seasonal cycle would consist 253 of e.g. 12 monthly averages or 26 bi-weekly averages. In the short global_oce_cs32 benchmark 254 a seasonal cycle is represented by a cycle of just two time steps. The three corresponding 255 cycle durations would be specified as $xx_gentim2d_preproc_i=12$, 26 and 2 respectively. The 256 corresponding time mean would always be specified as xx_gentim2d_preproc_i=1. 257

²⁵⁸ With 'gentim2d', it can be imposed that adjustments stay bounded via 'xx_gentim2d_bounds'.

Within 'genarr2d' and 'genarr3d', the corresponding option rather imposes bounds on adjusted parameters. Another run-time parameter in data.ctrl is 'mult_gentim2d' (or the 'genarr2d', 'genarr3d' version) that sets the multiplier for the corresponding cost function penalty (β_j in Eq. 1; $\beta_j = 1$. by default). Pre-conditioner \mathcal{R} (Eq. 5) does not appear in the estimation problem itself (Eq.1), as it only serves to push an optimization process preferentially towards certain directions of the control space. It is specified by 'gentim2dPrecond' (which is 1. by default).⁷

²⁶⁵ 3.3 implementation: pkg/ecco and pkg/ctrl

The implementation of Eqs. 2 and 3 belongs in 'pkg/ecco' and 'pkg/profiles'⁸ whereas Eqs. 4 and 266 5 belong in 'pkg/ctrl'. This section depicts the generic features implementations in 'pkg/ecco' 267 (first paragraphs) and 'pkg/ctrl' (later paragraphs). The maximum numbers of generic cost 268 function terms (NGENCOST), 3D cost function terms (NGENCOST3D), and post-processing 269 options (NGENPPROC) are set at compile time in ecco.h. The maximum numbers of generic 270 time-variable 2D controls (maxCtrlTim2D), time-invariant 2D controls (maxCtrlArr2D), time-271 invariant 3D controls (maxCtrlArr3D) and pre- or post-processing options (maxCtrlProc) are set 272 at compile time in CTRL_SIZE.h. Other files involved in compiling 'pkg/ecco' and 'pkg/profiles' 273 are listed in section 3.4. Run-time options for 'pkg/ecco' and 'pkg/profiles' (that should be the 274 main aspect of interest for most users) are readily documented in sections 3.1-3.2. 275

The operations in \mathcal{D} and \mathcal{S} (see Eq.3) are mainly carried out as the forward model steps 276 through time, respectively by ecco_phys.F and cost_averagesfields.F. During cost_averagesfields.F, 277 cost_gencost_customize.F maps physical variables to generic arrays (according to 'gencost_barfile' 278 specified in data.ecco; see section 3.1) and cost_averagesgeneric. F then proceeds with time-279 averaging, and periodically outputs the time-averaged m_i to file. Climatologies of m_i can be 280 formed (as an optional feature) from its time series to compare with observational o_i climatolo-281 gies (see section 3.1). This part of the m_i processing is carried out within cost_generic.F after 282 the full time series has been written to file. Model-data misfits are then computed (Eq. 2) by 283 cost_generic.F that relies on ecco_toolbox.F for elementary operations and on cost_genread.F for 284 re-reading m_i from file. The calls to cost_generic.F are operated in a loop by cost_gencost_all.F. 285 The overall sequence of operations for one cost function term is charted in Fig.11. The 286 distinction between 'preproc' and 'posproc' matches that between Eqs. 3 and 2. Most con-287 cretely the pre-processing ends and post-processing starts at the computation of $m_i - o_i$ using 288 'ecco_diffmsk' in cost_generic.F. Besides the numerous possibilities offered by this generic code, 289 specific cost function terms that do not fit in the Fig.11 chart quite vet can be operated via 290 cost_gencost_all.F and freely take advantage of the rest of the generic capabilities (storage ar-291 rays, adjoint checkpoint storage, run-time parameters, etc.). Examples of how to do this include 292 cost_gencost_boxmean.F, cost_gencost_sshv4.F and cost_gencost_seaicev4.F. 293

In the implementation of Eqs. 4 and 5, generality and versatility is greatly improved by operating virtually all of the pre-processing during model initialization. This is done by ctrl_map_ini_genarr.F ('genarr2d', 'genarr3d') and ctrl_map_ini_gentim2d.F ('gentim2d'). By the end of the processing steps, the effective version of the parameter adjustments ('gentim2d',

⁷The 'genarr2d' cases is treated accordingly, but the 'genarr3dPrecond' implementation seems incomplete.

⁸Additional documentation of 'pkg/profiles' is available in the MITgcm manual and in Forget et al. (2015).

Algo	orithm 1 Generic cost function algorithm.	
1: f	unction $COST_GENERIC()$	\triangleright Argument list defines the cost function
2:	call ecco_zero	\triangleright Initialize local array to 0
3:	call ecco_cprsrl	\triangleright Copy mask to local array
4:	for irec = 1, nrecloop \mathbf{do}	\triangleright Loop over time steps, days or months
5:	call cost_gencal	\triangleright Get file names, pointers
6:	Begin cost_genread	\triangleright Read, process model field
7:	if no preproc then	
8:	call ecco_readbar	\triangleright Read one record
9:	else if preproc=clim then	
10:	$\mathbf{call} \ \mathbf{ecco_readbar} \ \mathbf{within} \ \mathbf{loop}$	\triangleright Average records
11:	end if	
12:	$\mathbf{End} \ \mathbf{cost_genread}$	
13:	call mdsreadfield	\triangleright Read observational field
14:	call ecco_diffmsk	\triangleright Compute masked model-data misfit
15:	$\mathbf{if} \text{ posproc}=\text{smooth } \mathbf{then}$	
16:	$call smooth_hetero2d$	\triangleright Smooth masked misfit
17:	end if	
18:	call ecco_addcost	\triangleright Add to cost function
19:	end for	
20: e	nd function	

Figure 11: Chart of the generic cost function routine in pkg/ecco.

²⁹⁸ 'genarr2d') or of the adjusted model parameters ('genarr3d') are written to disk (with '.effec-²⁹⁹ tive' in the file name). Adjusted time-invariant model parameters are generally set during model ³⁰⁰ initialization, so their adjustments are also operated before the model time-stepping starts by ³⁰¹ ctrl_map_ini_genarr.F (for 'genarr2d', 'genarr3d'). Adjusted forcing variables are however reset ³⁰² at each time-step during the model run, so their adjustments are operated at the same times ³⁰³ by e.g. exf_getsurfacefluxes.F or exf_getffields.F (for 'gentim2d'). The cost function term $\mathfrak{u}_j^T\mathfrak{u}_j$ ³⁰⁴ is computed once the model run is complete along with the rest of Eq.1 (see section 3.1).

It should be stressed that if a control parameter is activated at run-rime (see section 3.2) 305 but its uncertainty is not specified (or vice versa) then ctrl_readparms.F signals the inconsis-306 tency before stopping the model during its initialization (the user may otherwise overlook that 307 his specification would have no effect). It should also be noted that the effective version of 308 'gentim2d' adjustments always consists of a full time series of length \approx the duration of the run 309 divided by 'xx_gentim2d_period'. It may contain e.g. repeated seasonal mean adjustments or 310 a sequence of interannual adjustments. This choice allows for a uniform, simple and general 311 treatment of all varieties of 'gentim2d' adjustments needed during the model integration (that 312 boils down to the temporal interpolation carried out in ctrl_map_gentim2d.F). This approach 313 is particularly advantageous in the context of the checkpointed adjoint model development (see 314 MITgcm manual). It only implies a marginal overhead in disk storage as compared with the 315 adjoint checkpointing output itself or with e.g. six-hourly re-analysis forcing input files. 316

317 3.4 Legacy: pkg/ecco and pkg/ctrl

Much of the legacy code that has been distributed as part of 'pkg/ecco' and 'pkg/ctrl' in the past is now deprecated – it is superseeded by the generic cost function and control codes presented above. Most of the deprecated codes had not been tested or maintained for many years, and consist of variations of the same operations duplicated many times. Another issue was the lack of organization amongst the deprecated codes (unlike in Fig.10). The consensus was that there was no point in keeping them around much longer.

For the time being the deprecated codes still exist but they are not compiled anymore unless the 'ECCO_CTRL_DEPRECATED' compile option is added in e.g. 'ECCO_CPPOPTIONS.h' (see below for details). To further facilitate the transition from old to new setup, the ctrlUseGen run-time parameter allows a switch between the old and new (generic) treatment of control vectors (assuming that 'ECCO_CTRL_DEPRECATED' was defined at compile time). As a side note: there is one non-generic feature that ISN'T deprecated since it has not been re-implemented in generic fashion, which is the control of open boundary conditions.

The deprecation of the legacy codes leads to a vast reduction in the volume of estimation 331 codes (30% of the code treated by automatic differentiation, which includes the entire phys-332 ical model, was removed in the process), a vast addition of capabilities (new or pre-existing 333 functionalities are now available for any gridded data set), and a greatly improved flexibility 334 (virtually all options can now be switched on/off at run time). Furthermore, the ecco, ctrl and 335 autodiff packages were made independent of each other, and to follow common MITgcm coding 336 practices. For example they can now be switched on/off at run time, independently (by virtue 337 of useECCO, useCTRL, useAUTODIFF). 338

Compiling options are typically found in the 'code/' directory of any given setup of MITgcm (when customized) or in the corresponding MITgcm package (when using defaults). The most obvious difference between the new setup and an old setup is that CPP_OPTIONS.h now disregards ECCO_CPPOPTIONS.h and uses the following instead :

- AUTODIFF_OPTIONS.h contains the few compile directives of pkg/autodiff. The maximum numbers of time steps are set in tamc.h
- ECCO_OPTIONS.h contains compile directives of pkg/ecco. Very few remain necessary, since all generic cost function settings can now be chosen at run time. The maximum numbers of cost terms are set in ecco.h
- CTRL_OPTIONS.h contains compile directives of pkg/ctrl. Very few remain necessary, since all generic control settings can now be chosen at run time. The maximum numbers of controls are set in CTRL_SIZE.h
- along with MOM_COMMON_OPTIONS.h, GMREDI_OPTIONS.h, GGL90_OPTIONS.h,
 PROFILES_OPTIONS.h, EXF_OPTIONS.h, SEAICE_OPTIONS.h, DIAG_OPTIONS.h