# Guidelines on debugging and introduction to Totalview

**Dominique Lucas**
**George Mozdzynski**

ECMWF

# Outline

- Different approaches to debugging

    - Printing

    - Compiling environment aiding tools

    - Signal trapping

- Debugging with totalview

**ECMWF**

# Don't Panic

- "Programming is an art" … and so is debugging
  - "You don't need a sledgehammer to crack a nut"
- Most problems are trivial and easy to fix.
- Look at stack trace/point of failure.
- Intuition, experience, luck all play a part.
- Check your code "again".
- Explain/show your code to somebody else.
- Use make to build your libraries/models, to guarantee that the same options are used everywhere.

ECMWF

# Some suggestions

- Is the problem reproducible on a rerun?

  - Does it fail in exactly the same way and place
- Try OMP_NUM_THREADS=1
- No stack trace?

  - suspect 'you' have trashed memory
- Try to reproduce problem at a lower resolution

ECMWF

# The Universal Debug Tool
## (the print/write statement)

- Will always be there for you!

- What to write out.

- Must be selective to keep file size(s) manageable.

- Can be activated with an DEBUG variable or Namelist entry.

- For "production" runs, try and switch off verbose output mode.

**ECMWF**

# Memory Constraints

- Task stack limit 4 Gbytes

- Thread stack limit

  - Master thread 4 Gbytes

  - Other threads 256 Mbytes

- Large arrays (>1 Mbyte) should be declared allocatable

  - allocatable arrays use the (per task) heap

  - no practical limit on heap size

- Fortran 90/95 (xlf90/xlf95) put local data on stack.

- Fortran 77 (xlf) puts statically allocated local data on heap.

ECMWF

# Eoj – check memory and CPU utilisation

```
eoj vers1.4 run at Wed Jun 11 17:35:46 GMT 2003 on hpca2501 for jobstep hpca2301.347796.0
Queued             : Wed Jun 11 17:11:09 GMT 2003 for      898 seconds
Dispatched         : Wed Jun 11 17:26:07 GMT 2003 for      579 seconds
Job Name           : edhm_model_fcgroup1
Step Name          : 0
Owner              : rdx
Unix Group         : rd
Account            : ecpreops
STDIN              : /dev/null
STDOUT             : /hpca/rdx_dir/log/mpm/edhm/fc/fcgroup1/model.1
STDERR             : /hpca/rdx_dir/log/mpm/edhm/fc/fcgroup1/model.1
Class              : np
Step Type          : General Parallel
Node Usage         : shared
Step Cpus          : 8
Total Tasks        :
Blocking           :
Node actual        : 1
Adapter Req.       : (csss,MPI,shared,US)
Resources          : ConsumableCpus(1) ConsumableMemory(900.000 mb)
#*#* Next 3 times NOT up-to-date (TOTAL CPU TIME given later IS accurate)
Step User   Time : 00:13:41.940000
Step System Time : 00:00:33.760000
Step Total  Time : 00:14:15.700000 (855.7 secs)
#*#* Last 3 times NOT up-to-date (TOTAL CPU TIME given later IS accurate)
Context switches : involuntary =        28637, voluntary    =        12378
                     per second =          49                          21
Page faults        : with I/O    =        13056, without I/O =     1203613
                     per second =          22                        2078
                   <--------- CPU --------> <------------ MEM ------------>
Node      ? #T #t secs/CPU    (Eff%) (Now%) max/TSK mb (Eff%) (Now% - mb   ) Task list
-------- - -- -- --------- ------ ------ ---------- ------ --------------- ---------
hpca1503 M  8  1    120.33 ( 20%) ( 52%)     764.85 ( 84%) ( 88% -   7680) 0:1:2:3:4:5:6:7:
-------- - -- -- --------- ------ ------ ---------- ------ --------------- ---------
 Elapsed =          579 secs                 900 mb = ConsumableMemory
 CPU Tot =       962.64 ( 0+00:16:02) Average:       963 s/node,       120 s/task
System Billing Units used by this jobstep = 1.037
```

# Debugging – compiler options

- checking:
    - argument checking: -qextchk
      $ xlf –qextchk prog.f –o prog
      Note that checking is done at compilation/linking

    - array bounds checking: -C
      $ xlf –C prog.f –o prog
      $ ./prog
      Note that checking is done at runtime

    - undefined reference checking
      $ xlf –qinitauto=FF –qsigtrap  \
          -qflttrap=inv:over:nanq:zero:en prog.f –o prog
      Note that checking is done at runtime

ECMWF

# Floating point exception

- In IEEE, a floating point exception sets status flag

- By default execution continues

- Trapping the exception requires software checking of status flags

- Utilities for enabling and checking

  - CALLs to fpgets and fpsets within program

- Automatic trapping by compiler (option –qflttrap)

  - high overhead to do precisely for whole program

  - "IMPrecise" option checks only at subprogram entry and exit

  - can apply flttrap to routines selectively

**ECMWF**

# Floating point exception

- IEEE exception types

    - OVerflow, UNDerflow, ZEROdivide, INValid, INEXact

- Other -qflttrap options

    - IMPrecise: check at routine exit and entry only

    - ENable

        - must specify in main program

        - may as well specify it everywhere

ECMWF

# Floating point exception

- Examples

    $ xlf –qflttrap=overflow:invalid:zerodivide:enable \
      –qsigtrap prog.f –o prog
    $ ./prog


- One can also use other exception handlers
- Relatively expensive … up to 20%

ECMWF

# Debugging – core files

- Core files – how to get a traceback
  - $ dbx ./prog core <<eof
    where
    eof

ECMWF

# ECMWF local signal trap - ECLIB

```fortran
INTEGER*4  CORE_DUMP_FLAG, IRETURN, SIGNALS(1), SIGNAL_TRAP
REAL A
CORE_DUMP_FLAG = 0
SIGNALS(1)    = 0
IRETURN = SIGNAL_TRAP(CORE_DUMP_FLAG, SIGNALS)
IF (IRETURN .LT. 0) THEN
  PRINT *, 'ERROR'
ELSE IF (IRETURN .EQ. 0) THEN
  PRINT *, 'FPE TRAPPING IS NOT SET'
ELSE
  PRINT *, 'FPE TRAPPING MODE =', IRETURN
ENDIF
call b(-2.)
end

subroutine b(a)
real a
write(*,*)sqrt(a)
return
END
```

- Link using $ECLIB, e.g.
  - $ xlf -c prog.f
  - $ xlf prog.o -o prog $ECLIB

ECMWF

# Signal_trap - arguments

- CORE_DUMP_FLAG:
  - = 0: no core dumped
  - Not = 0: core dumped
- SIGNALS – integer array with signals to trap
  - Signals(1)=0 => SIGFPE, SIGILL, SIGBUS, SIGSEGV, SIGXCPU.
  - See  "kill –l " for list of signals.
- Sample traceback …

```
Signal received: SIGTRAP - Trace trap
        Signal generated for floating-point exception:
                FP division by zero
Instruction that generated the exception:
        fdivs fr01,fr01,fr02
        Source Operand values:
                fr01 =   1.00000000000000e+00
                fr02 =   0.00000000000000e+00
Traceback:
        Offset 0x00000040 in procedure sub_
        Offset 0x00000048 in procedure ifs_model
        --- End of call chain ---
```

# Totalview

- Recompile your application without optimization (and –g):
  - -qnooptimize for all routines
  - -qsmp=noopt for routines with OpenMP
  - Beware –qsmp=omp implies optimization
  - Best choice: -g [–qoptdebug] -qfullpath
- Command Line interface exists
  - $ totalviewcli
- Recommended use of totalview in batch mode, i.e. with the GUI version.

ECMWF

# Totalview – GUI interface

- Before submitting batch job to launch totalview, request an X11 proxy at login time via ECaccess (ssh –X or NX).

- Include in your batch job the display:

  export DISPLAY=<your_ecaccess_display>
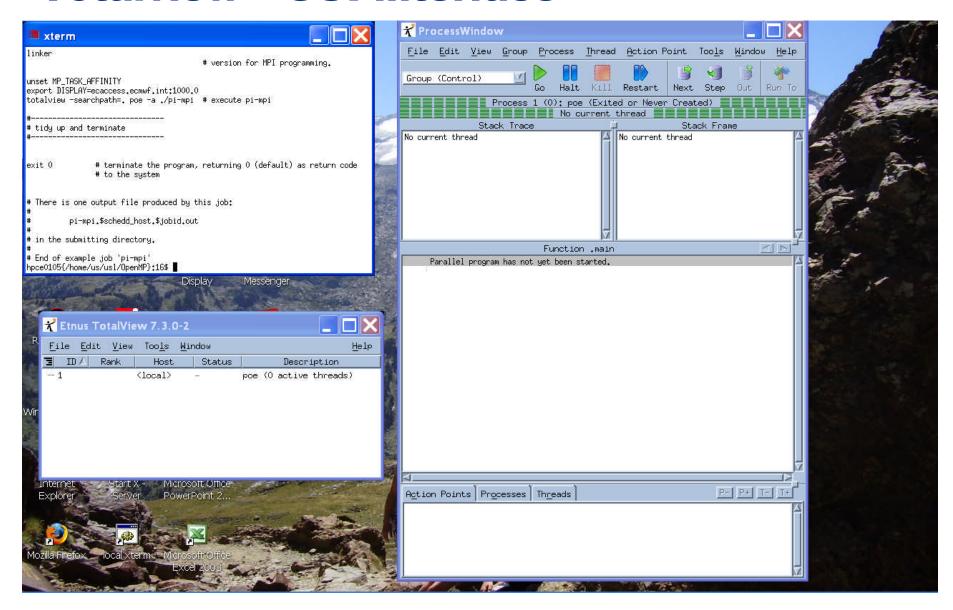
- If needed, include source code searchpath, e.g.

  searchpath='dir_1/,dir_2/,...,dir_n/'

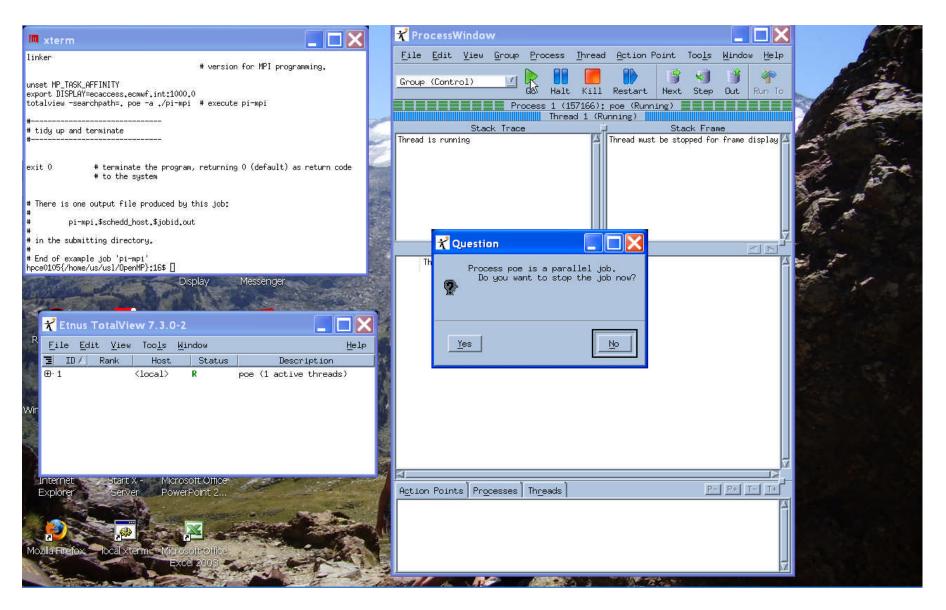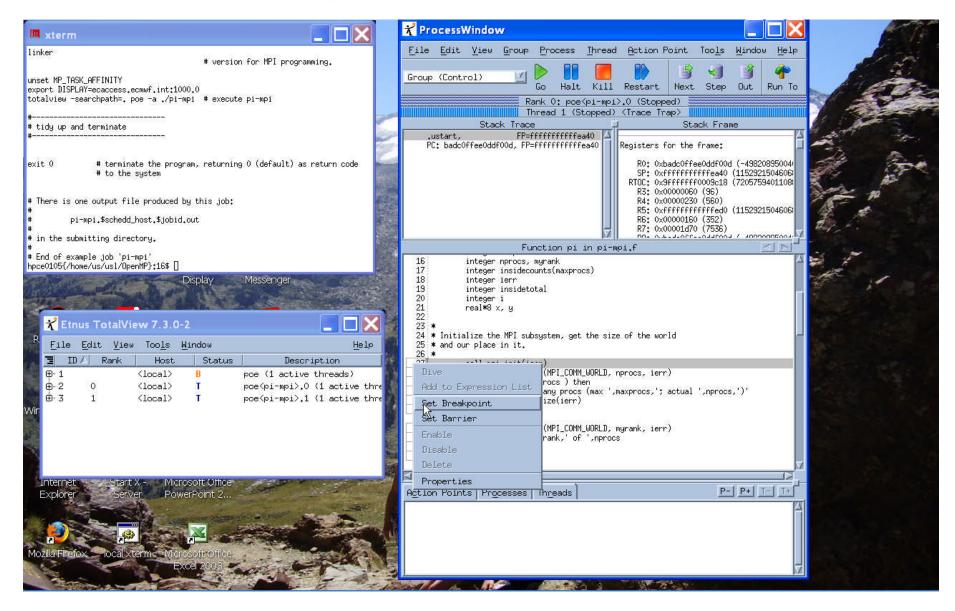- For MPI-parallel (load-leveller jobs):
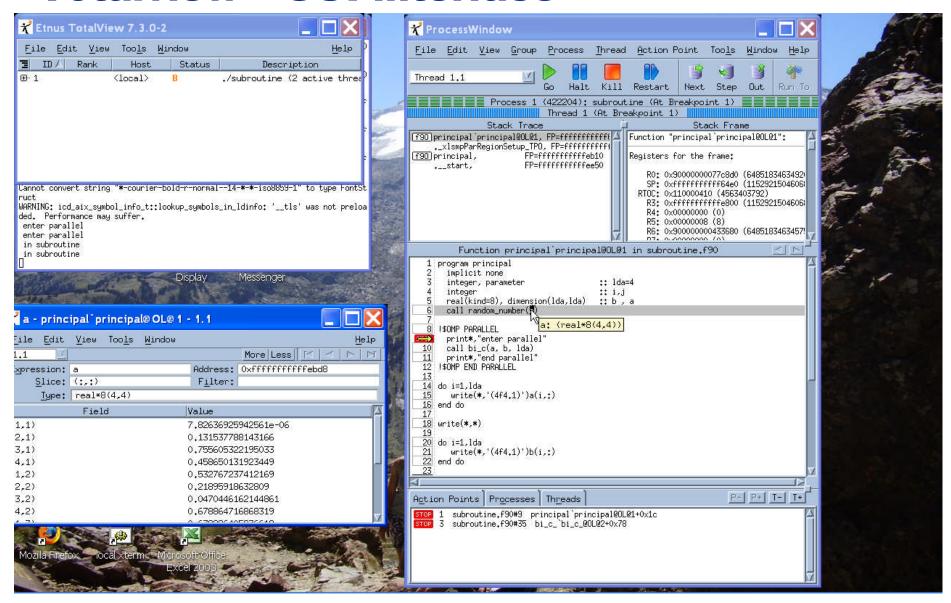
  totalview -searchpath=$searchpath poe -a <executable> <args>

- For serial/OpenMP only (interactive)

  totalview -searchpath=$searchpath <executable> -a <args>

ECMWF

# Totalview – GUI interface

# Totalview – GUI interface

# Totalview – GUI interface

# Totalview – GUI interface

# Totalview – GUI interface
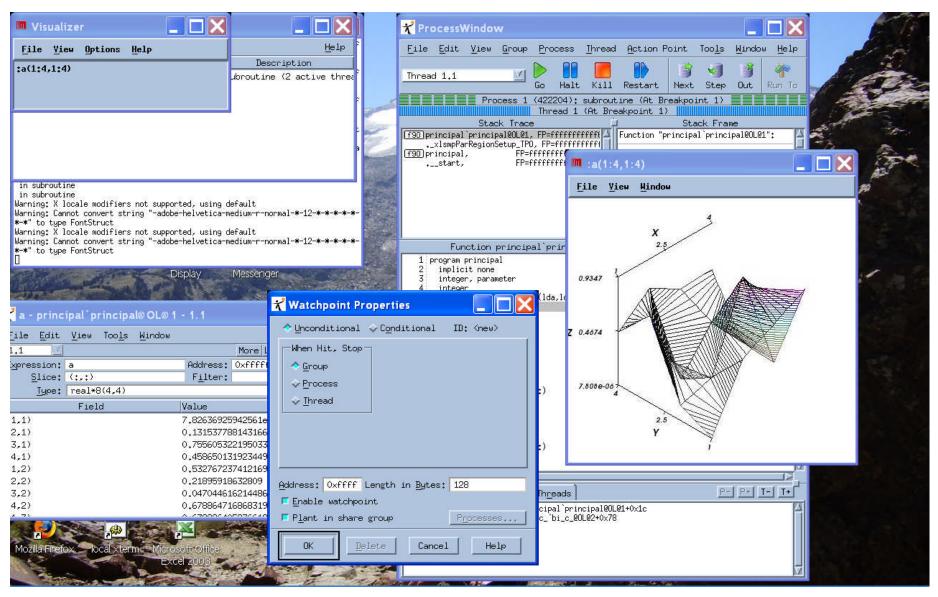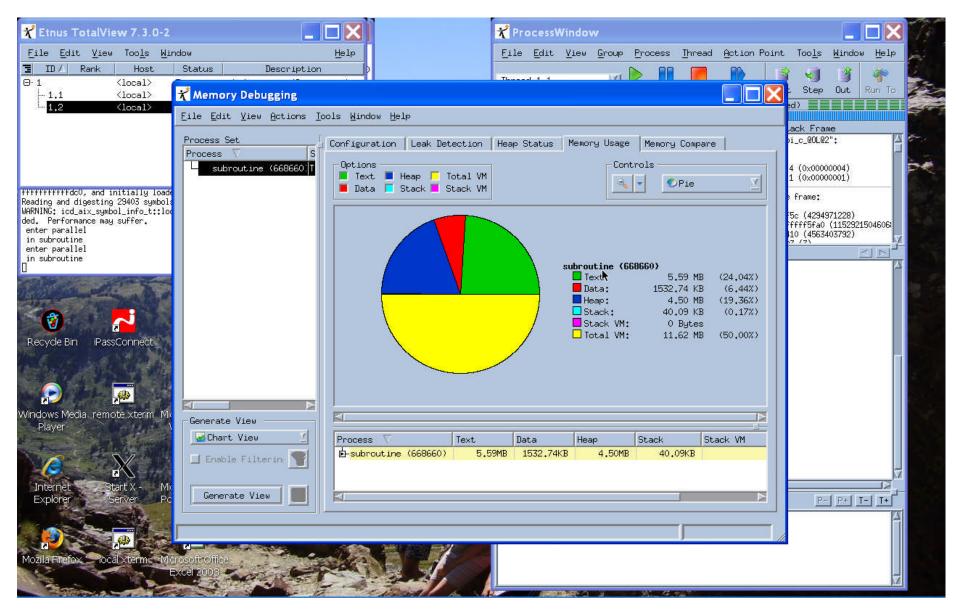
# Totalview – GUI interface

# Case Study – "bug" in EC-EARTH

- EC-EARTH – international project for earth system modelling.

- Experimental runs on our HPC systems.

- Coupled model runs, using MPMD parallel modelling approach:

  - poe -pgmmodel mpmd -cmdfile ${cmdfile} …

  with cmdfile like:

  oasis3.MPI1.x
  ifsMASTER -v ecmwf -e ewhx
  opa_exe.ORCA2_OASIS3.1.1
  appl-tm5.x

ECMWF

# Case Study – initial error

- Oasis compiled with –O3, TM5 with –g, IFS with –g and signal trap:

  1:signal_drhook(SIGABRT=6): New handler installed at 0x100ac5b8; old preserved at 0xa062ea70

  ERROR: 0031-250  task 3: Segmentation fault

  1:  Traceback:

  1:    Offset 0x0000086c in procedure pm_async_thread

  1:    Offset 0x000000dc in procedure _pthread_body

  1:    --- End of call chain ---

- Error (segmentation fault) in TM5 and IFS reports where it failed …

ECMWF

# Case Study – totalview (1)

- The job has been adapted to use totalview:

  export DISPLAY=galahad:0.0

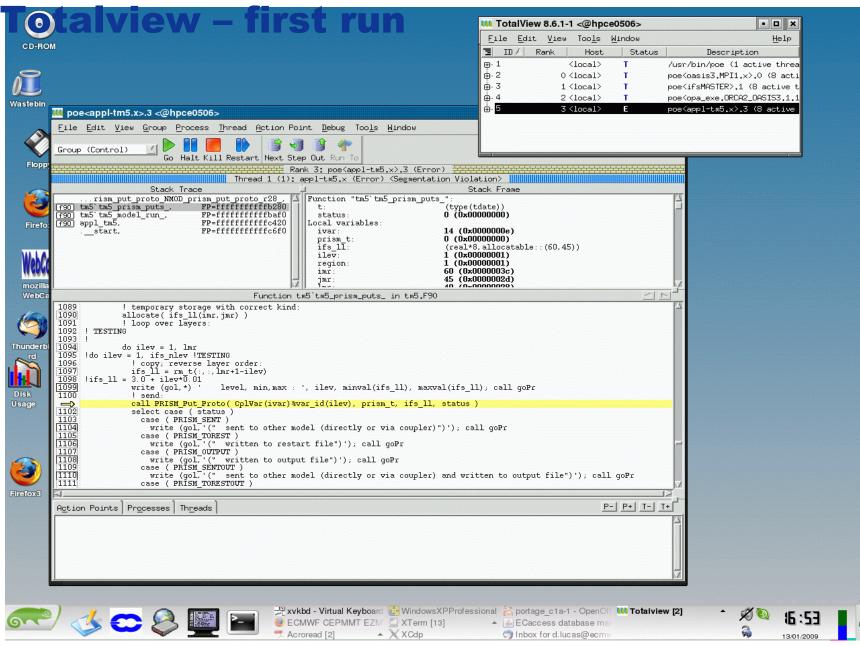  export MP_PGMMODEL='mpmd'

  export MP_CMDFILE="${cmdfile}"

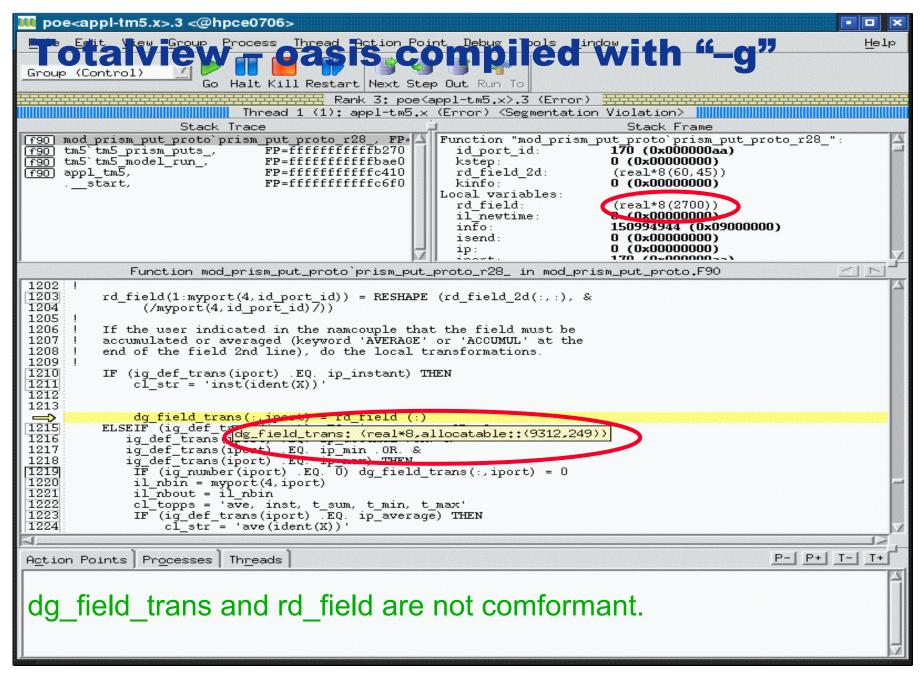  export MP_HOSTFILE="${hostfile}"

  export MP_LABELIO='yes'

  /usr/local/apps/toolworks/totalview.8.6.1-1/bin/totalview /usr/bin/poe

- Error in TM5 in a routine PRISM_Put_Proto (see following page).

ECMWF

# Totalview – first run

# Case Study – signal trap

- Activate signal trap in TM5:

  3: Signal received: SIGSEGV - Segmentation violation

  3:

  3: Traceback:

  3:   Location 0x000000010072c7fc

  3:   Offset 0x00001d08 in procedure __tm5_NMOD_tm5_prism_puts_, near line 1101 in file tm5.F90

  3:   Offset 0x00001b50 in procedure __tm5_NMOD_tm5_model_run_, near line 842 in file tm5.F90

  3:   Offset 0x000004d0 in procedure appl_tm5, near line 189 in file appl-tm5.F90

  3:   --- End of call chain ---

- Not much more information yet, because oasis not compiled with "-g".

- Recompiling oasis with "-g" and running totalview produces …

**ECMWF**

# Totalview – oasis compiled with "–g"



dg_field_trans and rd_field are not comformant.

# References

- Xlf Compiler Reference:

http://www.ecmwf.int/publications/manuals/hpcf_power6/xlf12.1_cr.pdf

- Xlf Optimisation and Programming Guide:

http://www.ecmwf.int/publications/manuals/hpcf_power6/proguide.pdf

- Totalview:

http://www.roguewave.com/support/product-documentation/totalview-family.aspx#totalview

ECMWF