



The BSC Tools: Extrace and Paraver

Based on slides by Judit Gimenez, BSC

EU H2020 Centre of Excellence (CoE)



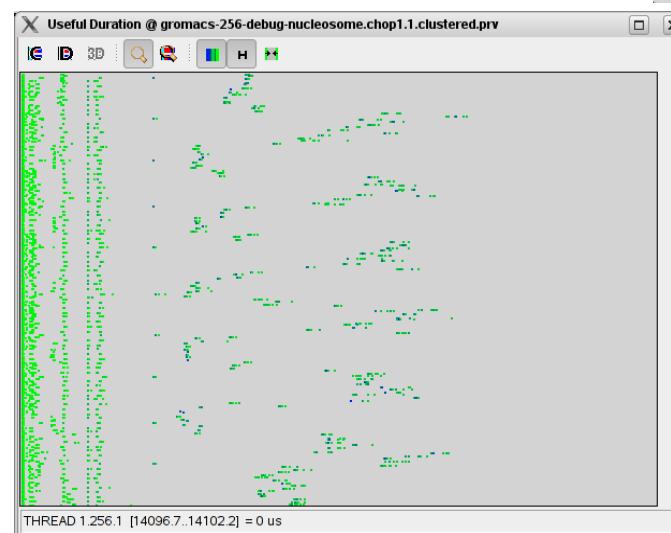
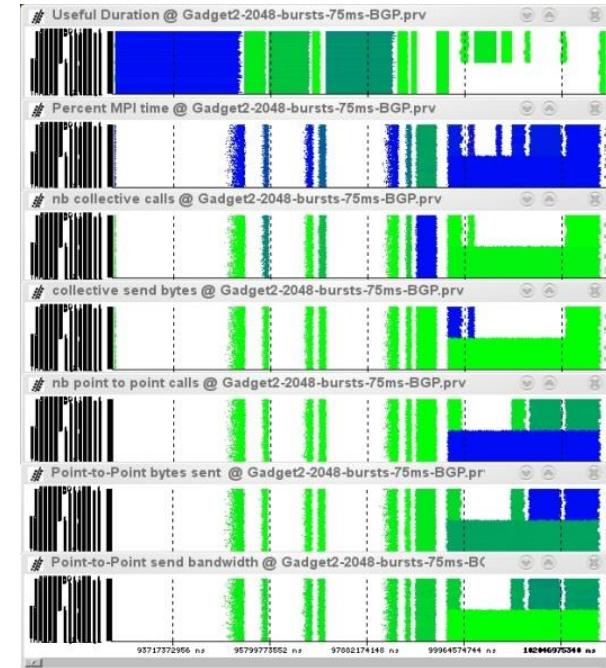
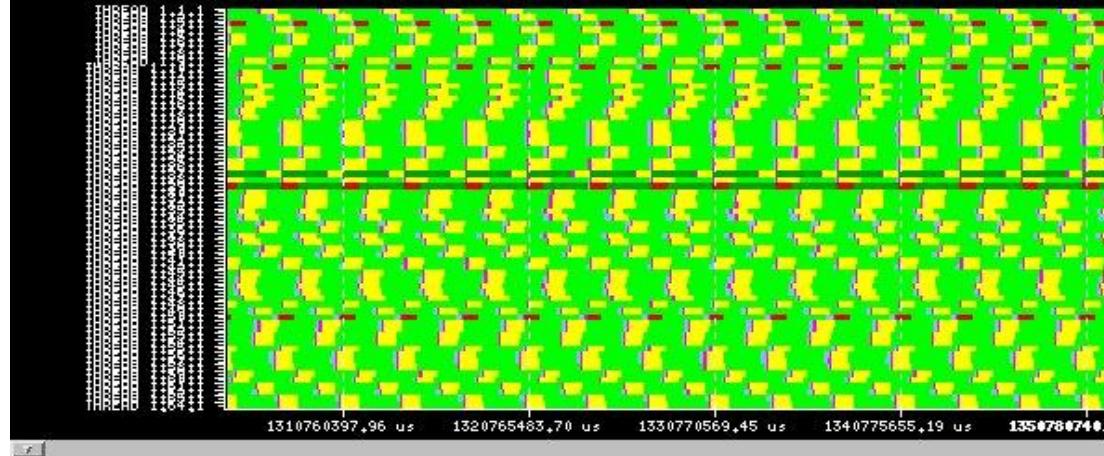
1 October 2015 – 31 March 2018

Grant Agreement No 676553

The BSC Tools



- Since 1991
- Based on traces
- Open Source
 - <http://tools.bsc.es>
- Core tools:
 - Extrae – instrumentation
 - Paraver (paramedir) – offline trace analysis
 - Dimemas – message passing simulator
- Focus
 - Detail, variability, flexibility
 - Behavioral structure vs. syntactic structure
 - Intelligence: Performance Analytics





Extrae



Extrae Features



- Parallel programming models
 - MPI, OpenMP, pthreads, OmpSs, CUDA, OpenCL, Java, Python...
- Platforms
 - Intel, Cray, BlueGene, MIC, ARM, Android, Fujitsu, Sparc...
- Performance Counters
 - Using PAPI interface
- Link to source code
 - Callstack at MPI routines
 - OpenMP outlined routines
 - Selected user functions (Dyninst)
- Periodic sampling
- User events (Extrae API)

No need
to
recompile
/ relink!



How does Extrae work?



- Symbol substitution through **LD_PRELOAD**
 - Specific libraries for each combination of runtimes
 - MPI
 - OpenMP
 - OpenMP+MPI
 - ...
- Dynamic instrumentation
 - Based on Dyninst (developed by U.Wisconsin/U.Maryland)
 - Instrumentation in memory
 - Binary rewriting
- Static link (i.e., PMPI, Extrae API)

Recommended



Using Extrae in 3 steps



1. Adapt the job submission script
 2. (Optional) Tune the Extrae XML configuration file
 - Examples distributed with Extrae at \$EXTRAE_HOME/share/example
 3. Run it!
-
- For further reference check the **Extrae User Guide**:
 - Also distributed with Extrae at \$EXTRAE_HOME/share/doc
 - https://tools.bsc.es/tools_manuels



Step 1: Adapt the job script to load Extrae (LD_PRELOAD)



```
(...)

### Application path
PISVM=../../bin/pisvm-train

### Input path
TRAINDATA=../../input/sdap_area_all_training.el

### Extrae path
export EXTRAE_HOME=/homec/deep/deep83/tools/extrاء
export EXTRAE_WORK_DIR=/work/$USER/pisvm/romeraw

### Run the application
srun ./trace.sh $PISVM -D -o 1024 -q 512 -c 10000
         g 16 -t 2 -m 1024 -s 0 $TRAINDATA

### Generate the trace
export TRACE_NAME=pisvm-romeraw-train.prv
${EXTRAE_HOME}/bin/mpi2prv
-f ${EXTRAE_WORK_DIR}/TRACE.mpits
-o ${TRACE_NAME}
```

train.sh

```
#!/bin/bash

### Load Extrae
source ${EXTRAE_HOME}/etc/extrاء.sh
export EXTRAE_CONFIG_FILE=../../../../config/extrاء.xml

### Load the tracing library (choose C/Fortran)
export LD_PRELOAD=$EXTRAE_HOME/lib/libmpitrace.so
#export LD_PRELOAD=$EXTRAE_HOME/lib/libmpitracef.so

# Run the program
$*
```

trace.sh

Select tracing
library



Step 1: Which tracing library?



- Choose depending on the application type

Library	Serial	MPI	OpenMP	pthread	CUDA
libseqtrace	✓				
libmpitrace[f] ¹		✓			
libomptrace			✓		
libpttrace				✓	
libcudatrace					✓
libompitrace[f] ¹		✓	✓		
libptmpitrace[f] ¹		✓		✓	
libcudampitrace[f] ¹		✓			✓

¹ include suffix “f” in Fortran codes



Step 2: Extract XML configuration



```
<mpi enabled="yes">
  <counters enabled="yes" /> ← Trace MPI calls + HW counters
</mpi>

<openmp enabled="yes">
  <locks enabled="no" />
  <counters enabled="yes" />
</openmp>

<pthread enabled="no">
  <locks enabled="no" />
  <counters enabled="yes" />
</pthread>

<callers enabled="yes">
  <mpi enabled="yes">1-3</mpi> ← Trace call-stack events @ MPI calls
  <sampling enabled="no">1-5</sampling>
</callers>
```



Step 2: Extract XML configuration (II)



```
<counters enabled="yes">
  <cpu enabled="yes" starting-set-distribution="cyclic">
    <set enabled="yes" domain="all" changeat-time="0">
      PAPI_TOT_INS, PAPI_TOT_CYC, PAPI_L1_DCM, PAPI_L2_DCM
    </set>
    <set enabled="yes" domain="all" changeat-time="500000us">
      ...
    </set>
    <set enabled="yes" domain="all" changeat-time="500000us">
      ...
    </set>
    <set enabled="yes" domain="all" changeat-time="500000us">
      ...
    </set>
  </cpu>
  <network enabled="no" />
  <resource-usage enabled="no" />
  <memory-usage enabled="no" />
</counters>
```

Select which HW
counters are
measured



Step 2: Extract XML configuration (III)



```
<buffer enabled="yes">
    <size enabled="yes">5000000</size> ← Trace buffer size
    <circular enabled="no" />
</buffer>

<sampling enabled="no" type="default" period="50m" variability="10m" /> ← Collect more measurements with sampling

<merge enabled="yes"
    synchronization="default"
    tree-fan-out="16"
    max-memory="512"
    joint-states="yes"
    keep-mpits="yes"
    sort-addresses="yes"
    overwrite="yes"
>
    $TRACE_NAME$
</merge>
```

Merge intermediate files into Paraver trace

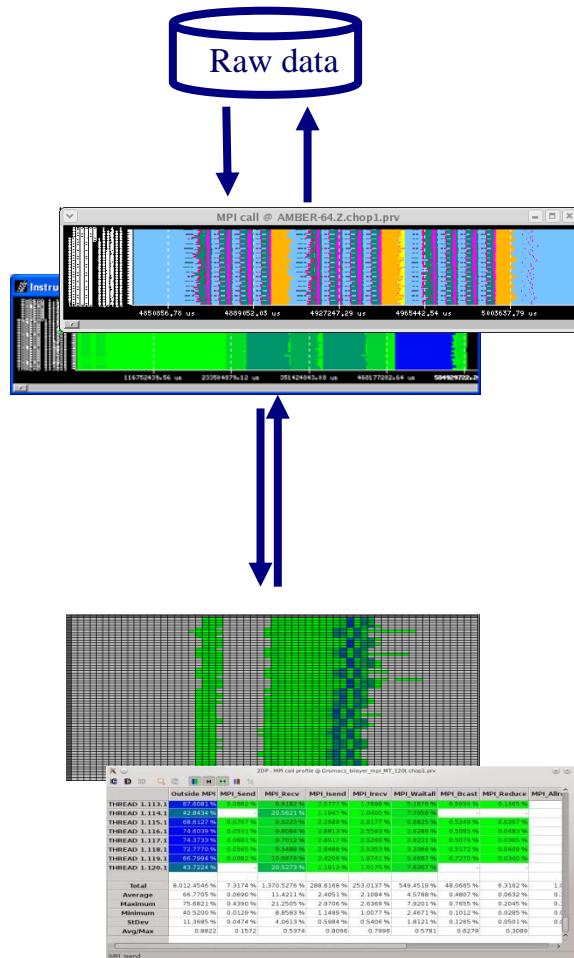




Paraver



Paraver – Performance data browser



Trace visualization/analysis
+ trace manipulation

Timelines

Goal = Flexibility
No semantics
Programmable

2/3D tables (Statistics)

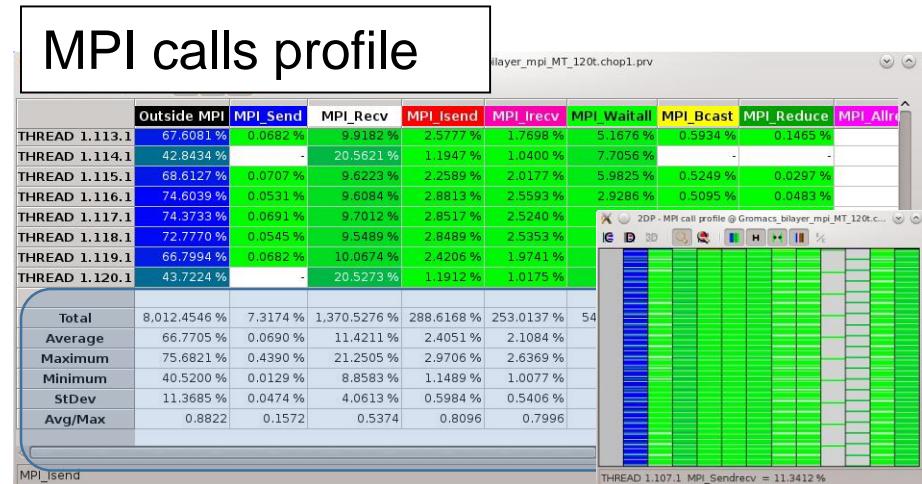
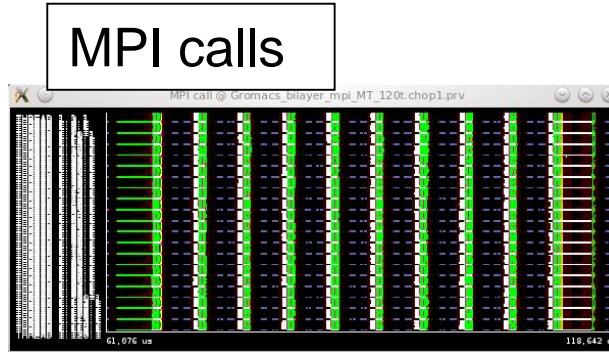
Comparative analyses
Multiple traces
Synchronize scales



Tables: Profiles, histograms, correlations



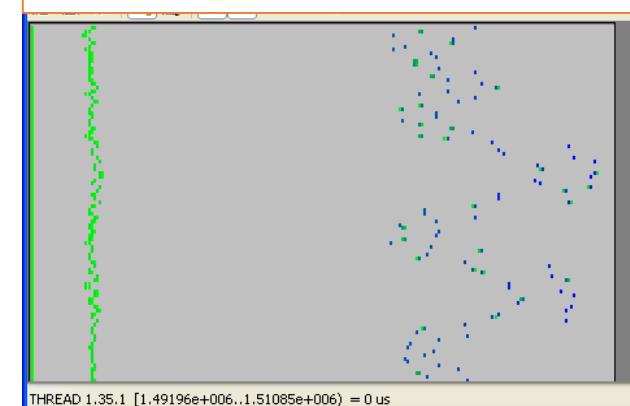
- From timelines to tables



Useful Duration



Histogram Useful Duration

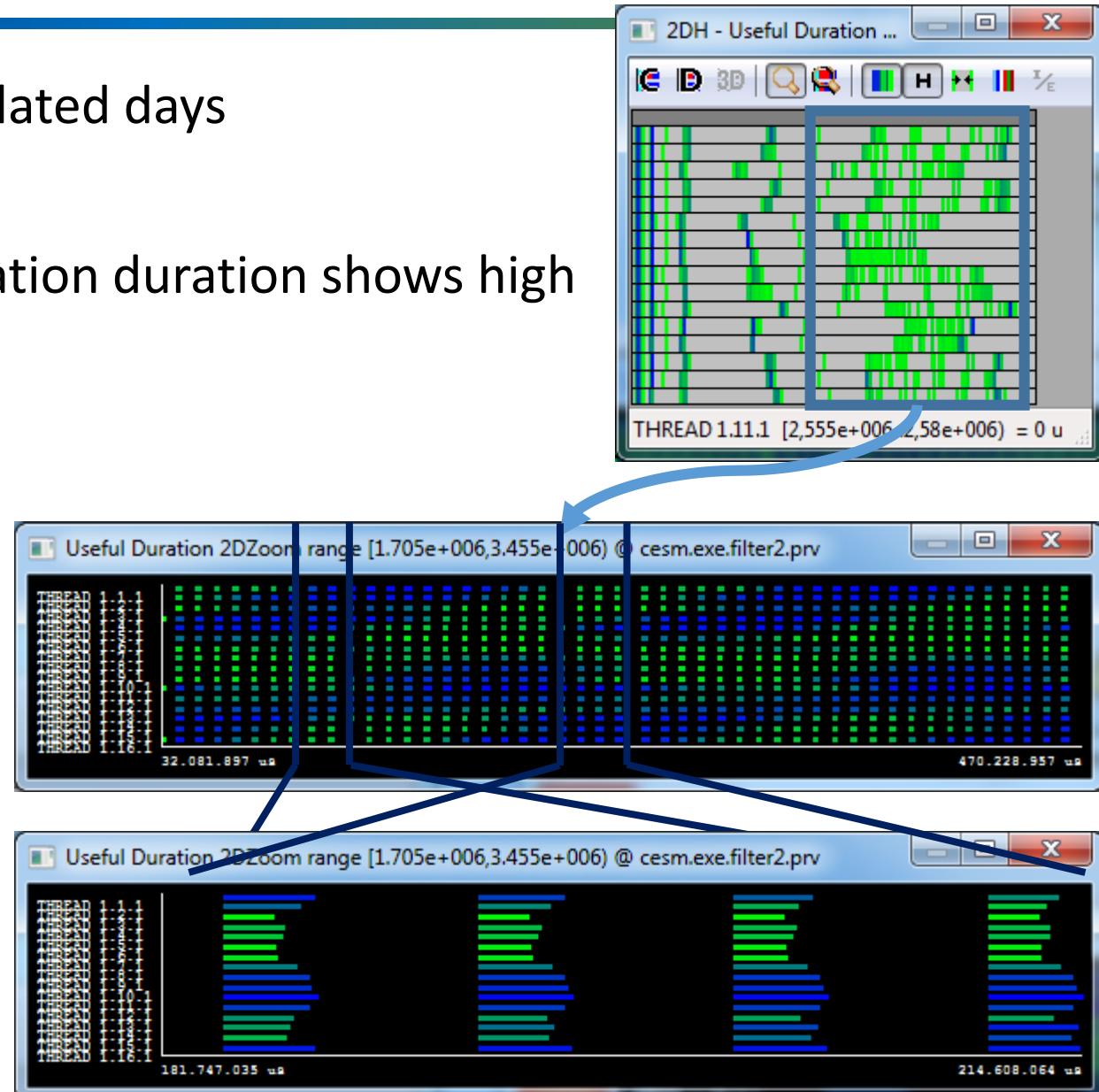


From tables to timelines



CESM: 16 processes, 2 simulated days

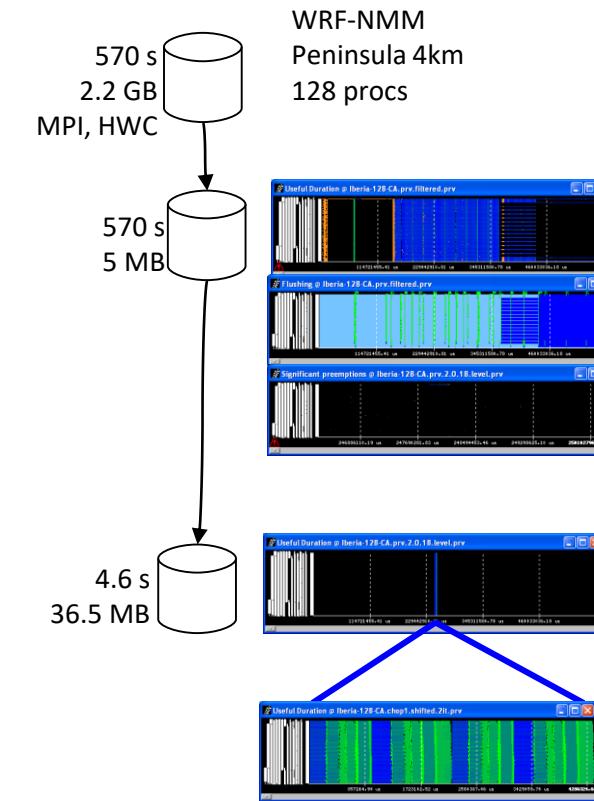
- Histogram useful computation duration shows high variability
- How is it distributed?
- Dynamic imbalance
 - In space and time
 - Day and night
 - Season ? ☺



Trace manipulation



- Data handling/summarization capability
 - Filtering
 - Subset of records in original trace
 - By duration, type, value,...
 - Filtered trace is still a Paraver trace and can be analysed with the same cfgs (as long as the data required has been kept)
 - Cutting
 - All records in a given time interval
 - Only some processes
 - Software counters
 - Summarized values computed from those in the original trace emitted as new even types
 - #MPI calls, total hardware count,...





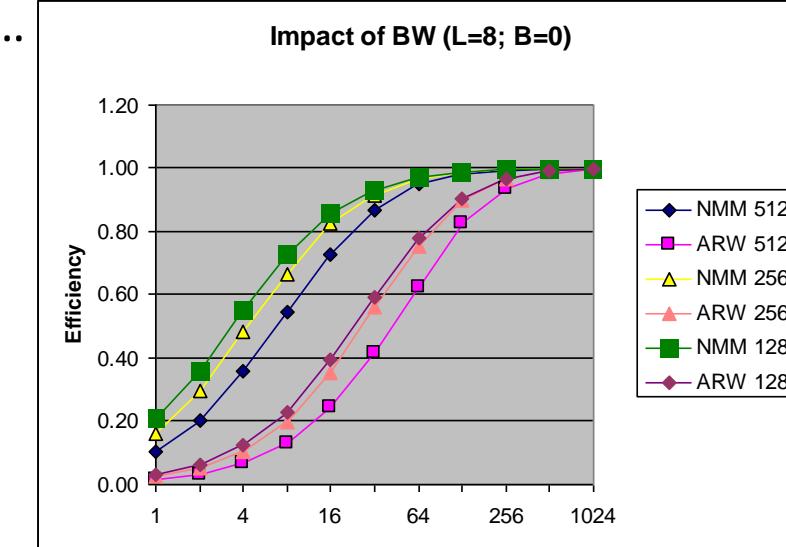
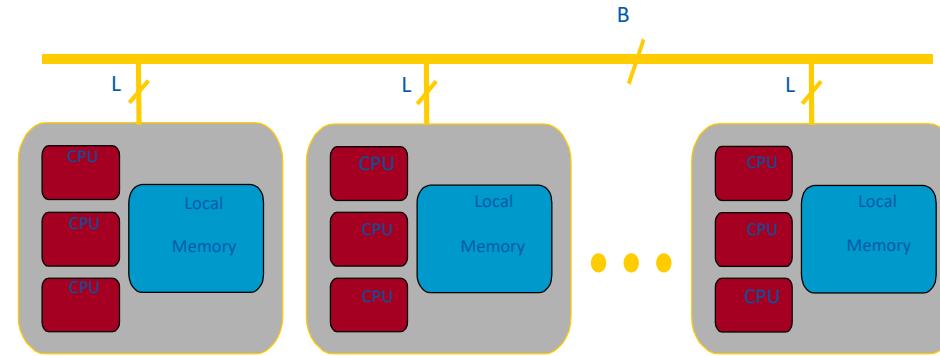
Dimemas



Dimemas: Coarse grain, Trace driven simulation



- Simulation: Highly non-linear model
 - MPI protocols, resource contention...
- Parametric sweeps
 - On abstract architectures
 - On application computational regions
- What-if analysis
 - Ideal machine (instantaneous network)
 - Estimating impact of ports to MPI+OpenMP/CUDA/...
 - Should I use asynchronous communications?
 - Are all parts equally sensitive to network?
- MPI sanity check
 - Modeling nominal
- Paraver – Dimemas tandem
 - Analysis and prediction
 - What-if from selected time window



Detailed feedback on simulation (trace)

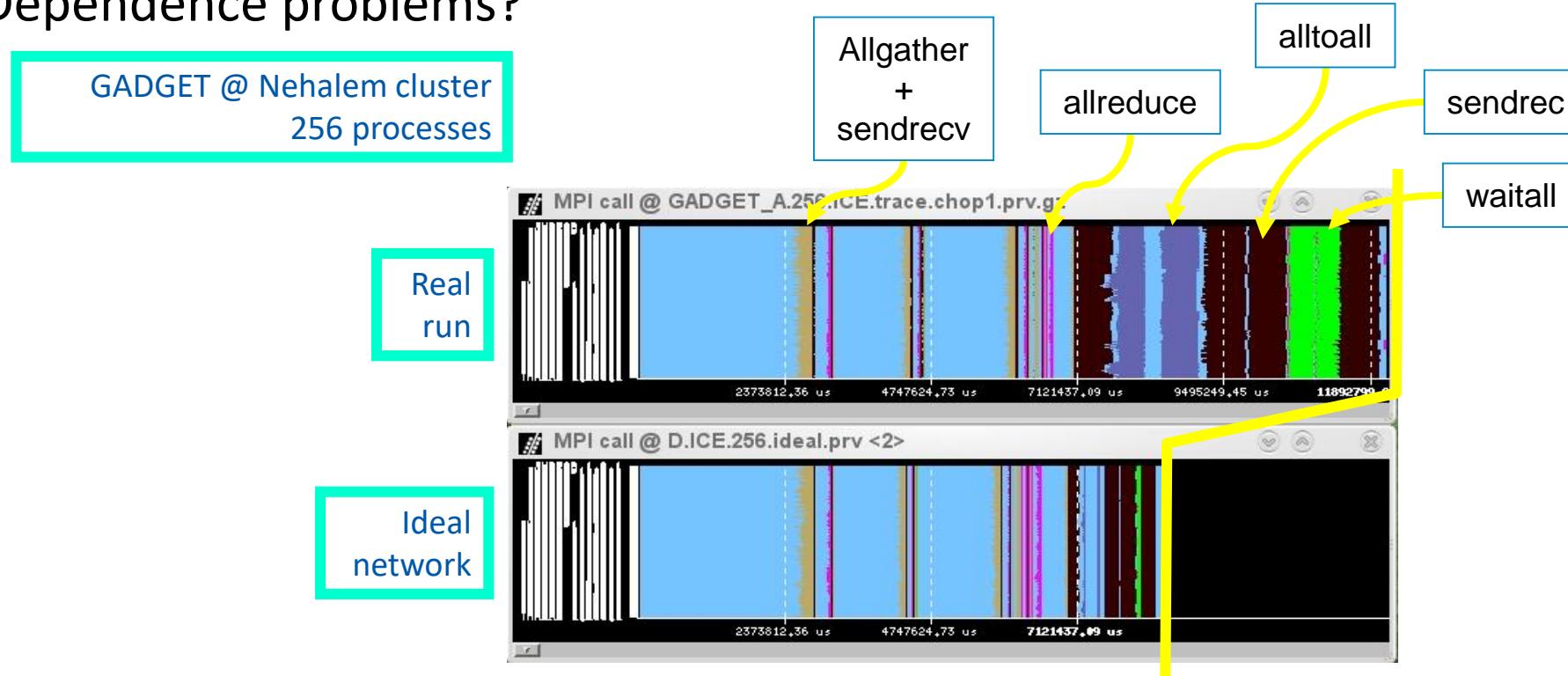


Ideal machine



The impossible machine: $BW = \infty$, $L = 0$

- Actually describes/characterizes Intrinsic application behavior
 - Load balance problems?
 - Dependence problems?





Performance Optimisation and Productivity

A Centre of Excellence in Computing Applications

Contact:

<https://www.pop-coe.eu>
<mailto:pop@bsc.es>

