

# Ellexus Ltd: The I/O Profiling Company

*Dr Rosemary Francis, CEO, Good I/O evangelist*

## How to recognize I/O bottlenecks and what to do about them

Rosemary will be sharing industry perspectives on how to recognise I/O bottlenecks and what to do about them. The delicate and often dynamic balance between I/O, CPU and memory can hide some easy wins in terms of improving throughput on-prem and reducing costs in the cloud. Equally, improving I/O is also about reducing the load on shared storage and not just about the incremental improvements of individual applications.



The I/O Profiling Company - Protect. Balance. Optimise.

[www.ellexus.com](http://www.ellexus.com)

# Ellexus Ltd: The I/O Profiling Company

**Products:** We make system telemetry tools to help you

- improve application performance,
- protect shared storage, and
- manage application dependencies for migration.

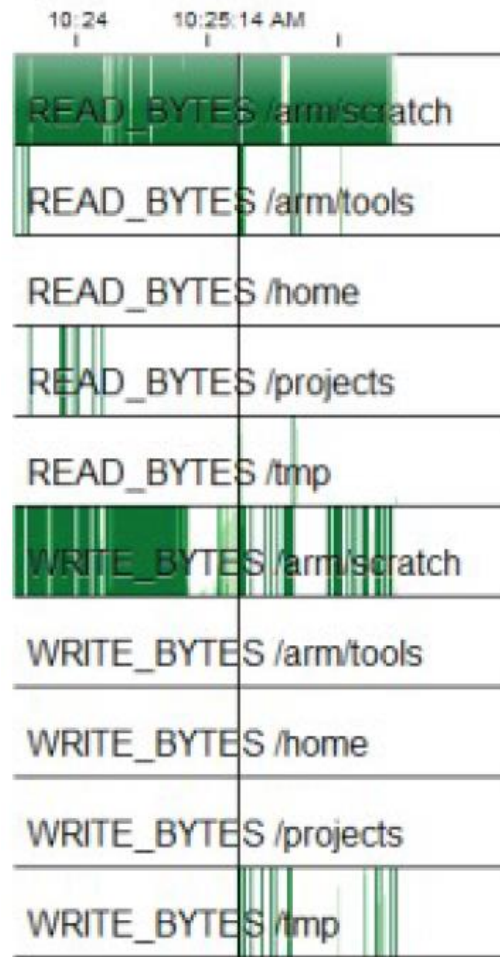
**Customers include:**



# Solving the noisy neighbour problem

How we worked with Arm to develop our technology

## Example of a rogue job from Arm:



This application is overloading shared storage by putting data in the wrong place.

Temporary data is written to shared storage

Local storage is unused

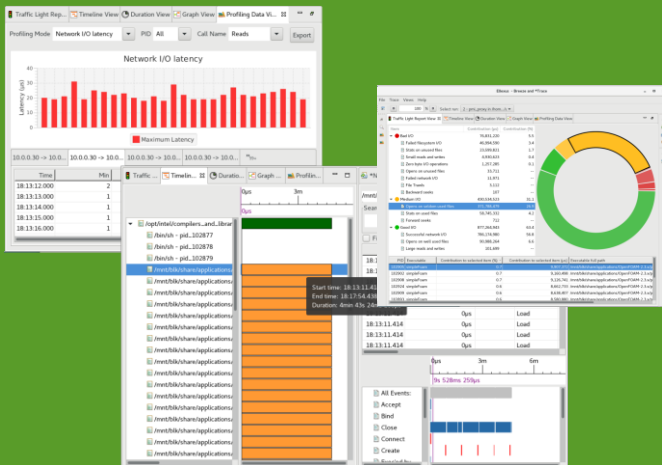


# Ellexus enterprise products

Take control of the way you access your data



Detailed I/O profiling  
Application discovery



## Dependencies

What do I need to include in my container?

How do I migrate this tool chain?

## I/O profiling

What resources do I need?

## Debug and triage

Why am I not getting the results I expect?



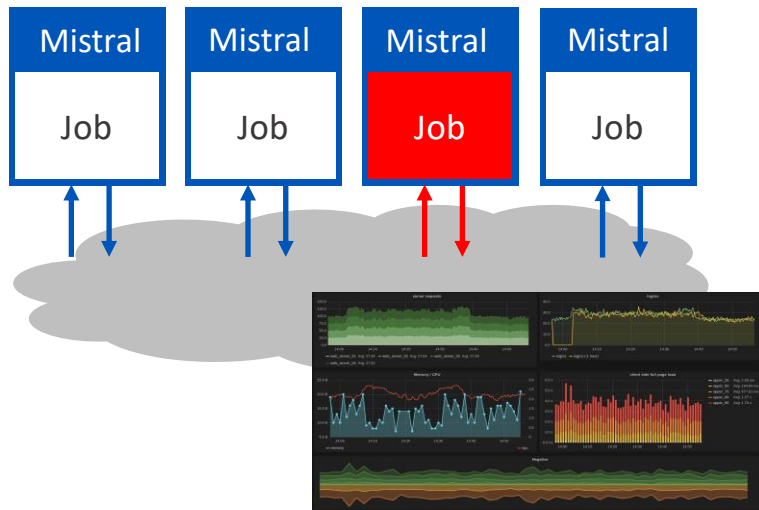
# Ellexus enterprise products

Take control of the way you access your data

Live telemetry for on-premises  
clusters and cloud

Protect storage and find bottlenecks

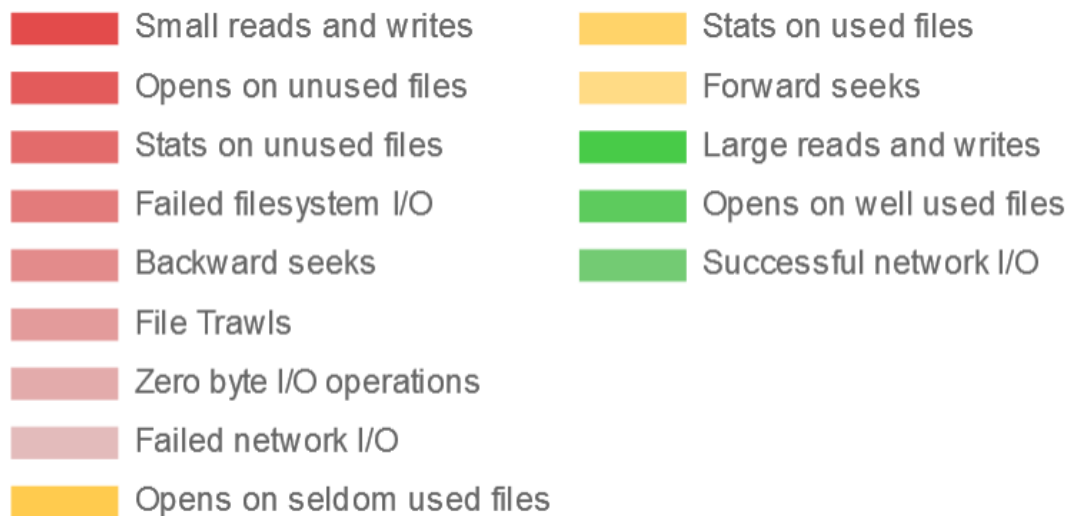
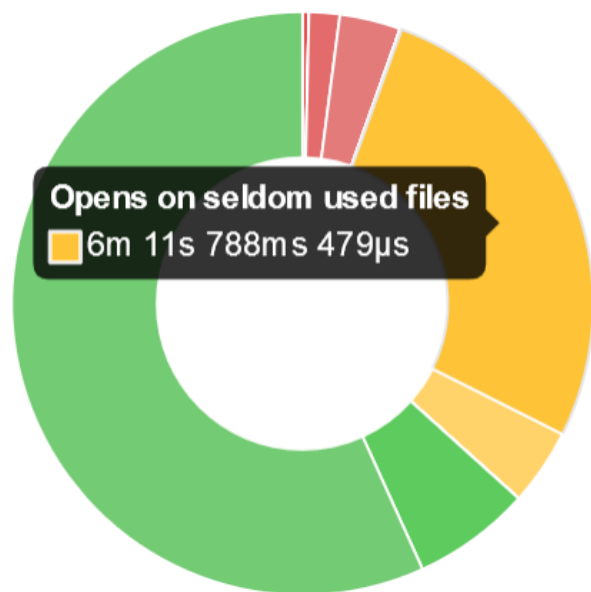
Cost management and forecasting



Live system telemetry:  
I/O monitoring in production

# Tuning and sizing:

How much time are you wasting doing bad I/O?



# Case study:

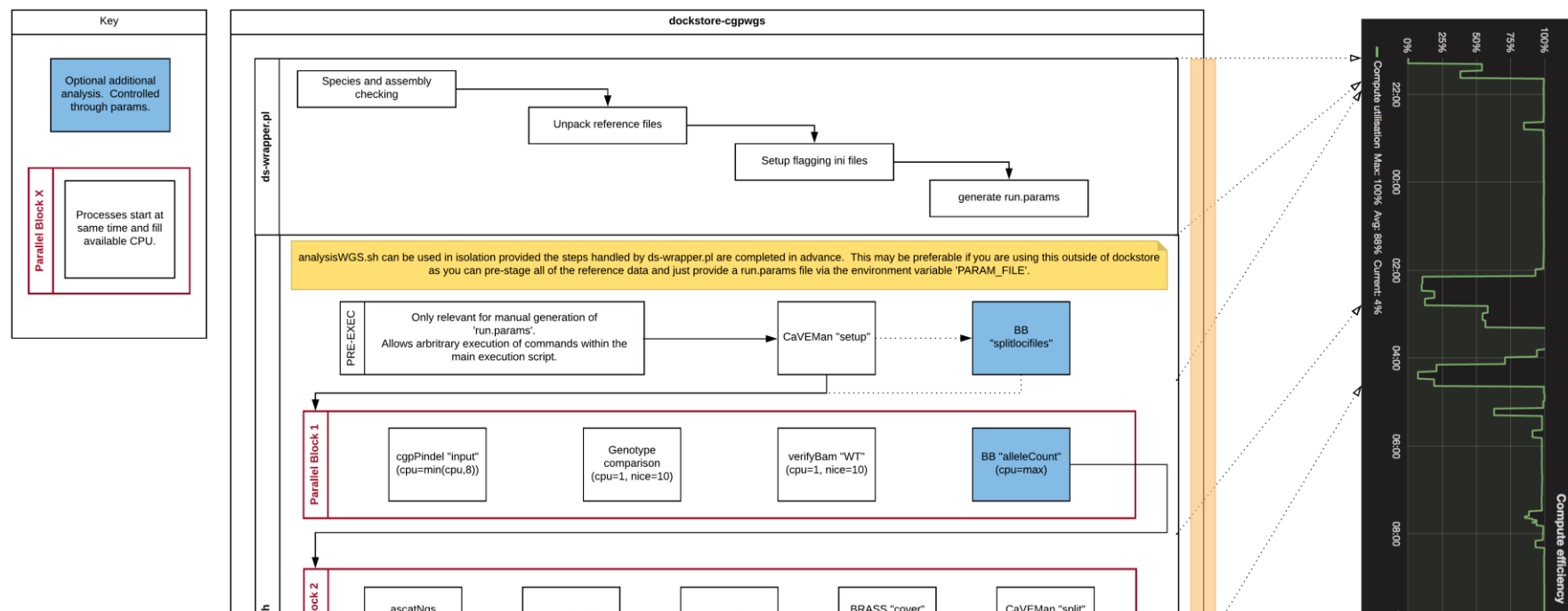
## Tuning cancer pipelines at the Sanger Institute

The Pancancer project: 2,000 whole genomes at multiple HPC sites

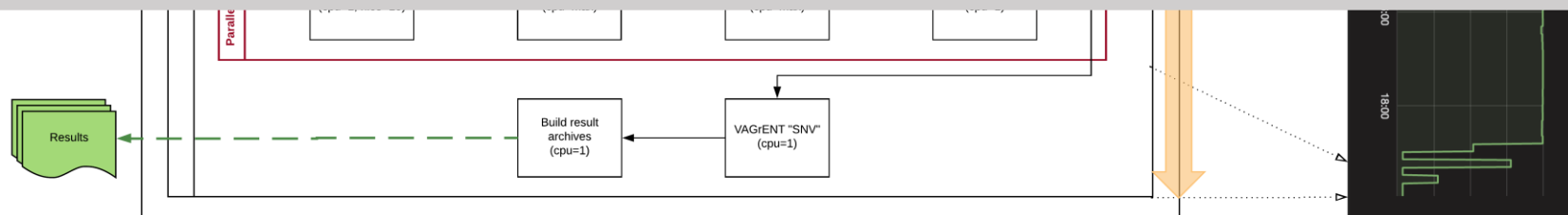
- Containerised pipelines for portability
- I/O tuned with Ellexus tools
- Storage now needs to be sized correctly



# Tuning cancer pipelines at the Sanger Institute



Runtime was reduced from 32hr to 18hr  
through profiling I/O and tuning deployment

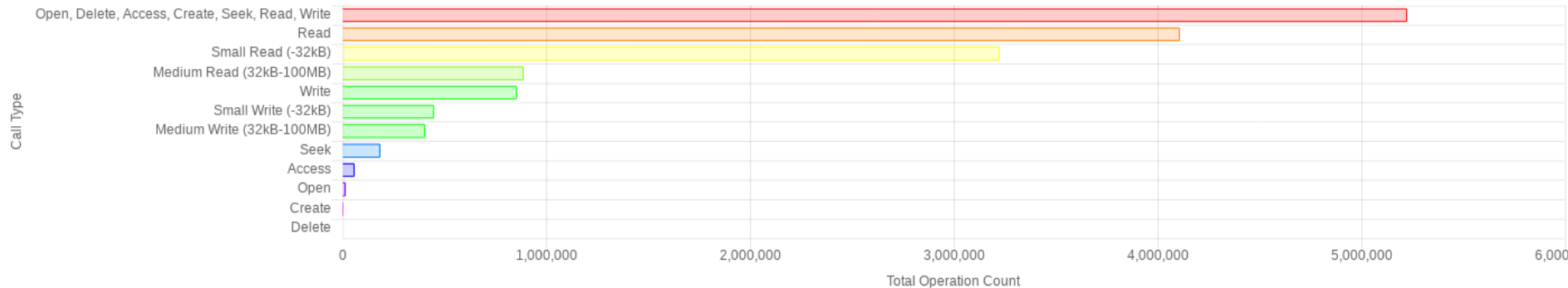




# Profiling the cancer pipeline

AWS m5.xlarge 4vCPU 16GB

Number of I/O operations() by type



Size of read and write operations()



# Storage comparison

	Time*		Cost per month (\$)	
GP2	52m 23s	100%	174.11	100%
Magnetic EBS	1h 01m 44s	118%	174.43	100%
Provisioned 100 IOPS	1h 42m 01s	195%	184.61	106%
Throughput optimised HDD	1h 19m 32s	152%	189.01	109%
150GB NVMe	51m 27s	98%	191.79	110%
Provisioned 500 IOPS	54m 22s	104%	215.01	123%

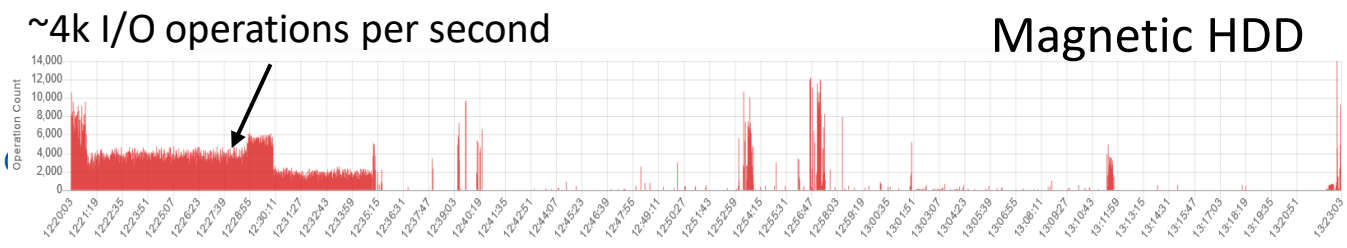
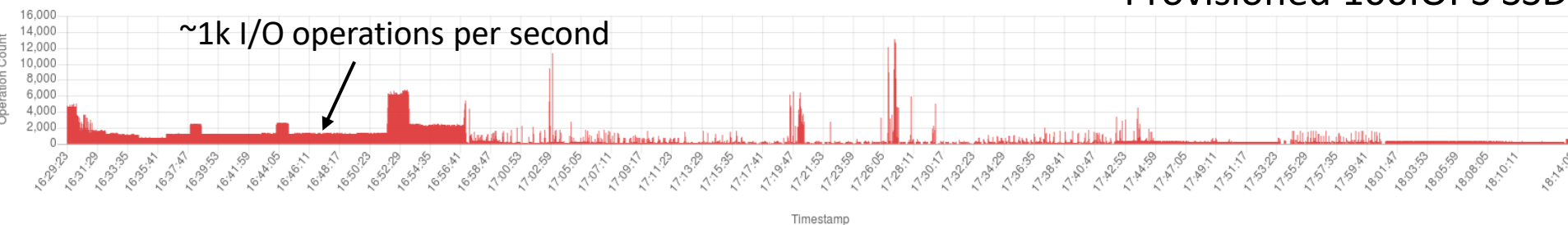
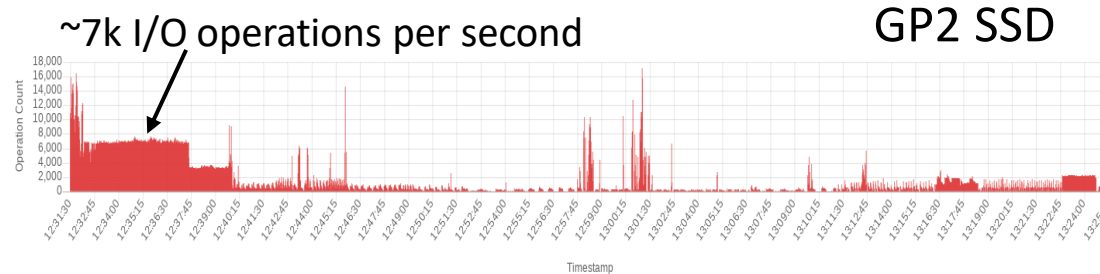
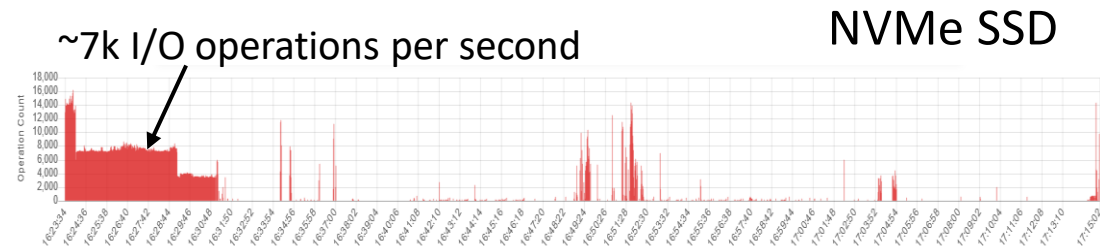
⇒ The Provisioned IOPS SSDs performed very badly

⇒ AWS default option, GP2 is the best

⇒ NVMe was only 2% faster for a 10% price increase

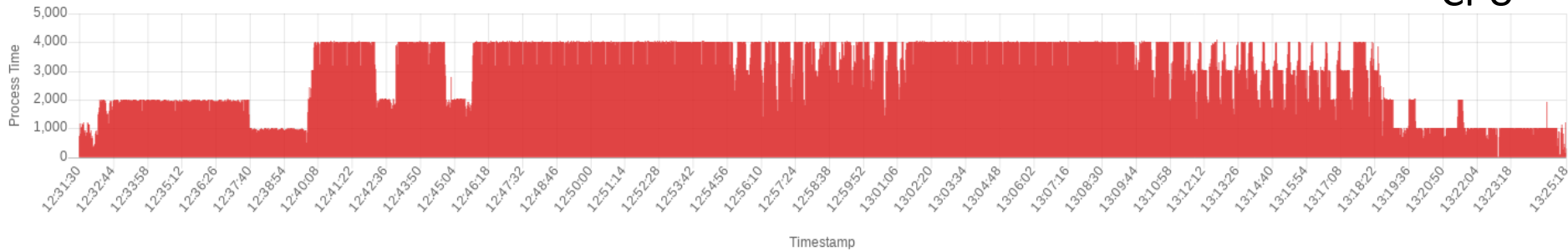


# I/O Operations() over time

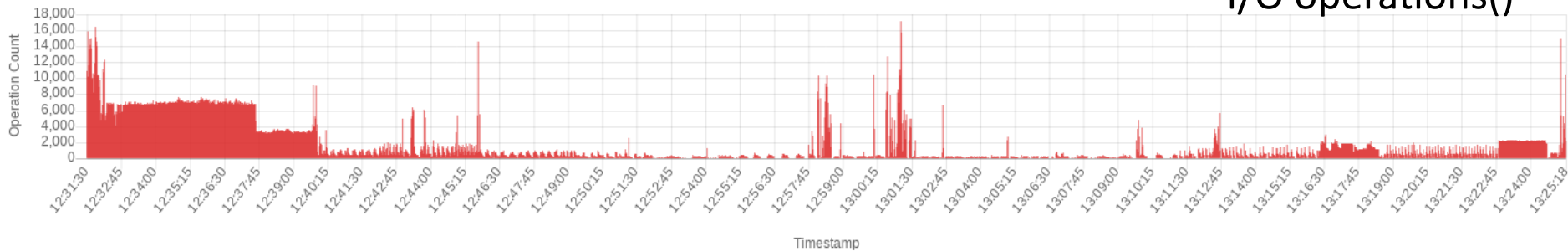


# CPU and I/O Profile (on GP2 SSD)

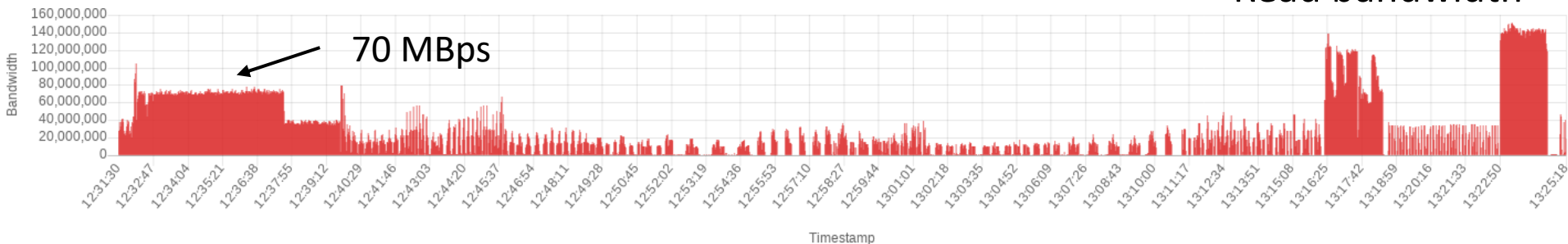
CPU



I/O operations()



Read bandwidth



# More CPU and less memory: m5.xlarge vs c5.xlarge

(still on GP2 SSD default storage)

## **M5.xlarge**

4 vCPU

16GB

Runtime: 53min

Cost: \$0.21

## **c5.xlarge**

4 vCPU

8GB

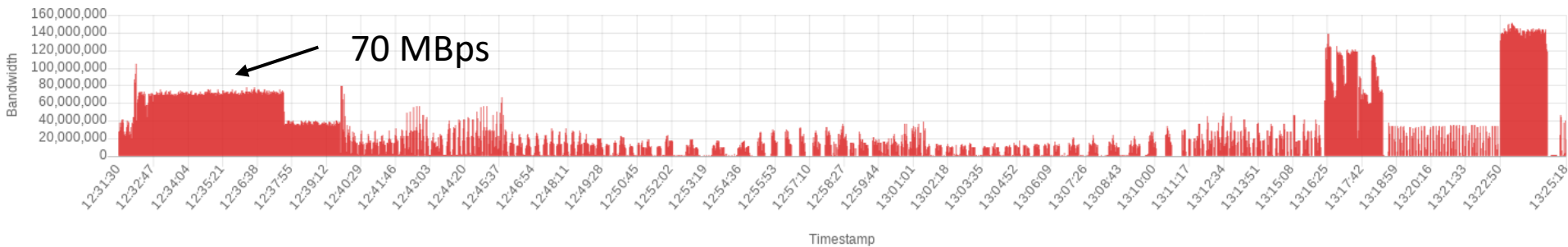
Runtime: 44min

Cost: \$0.16

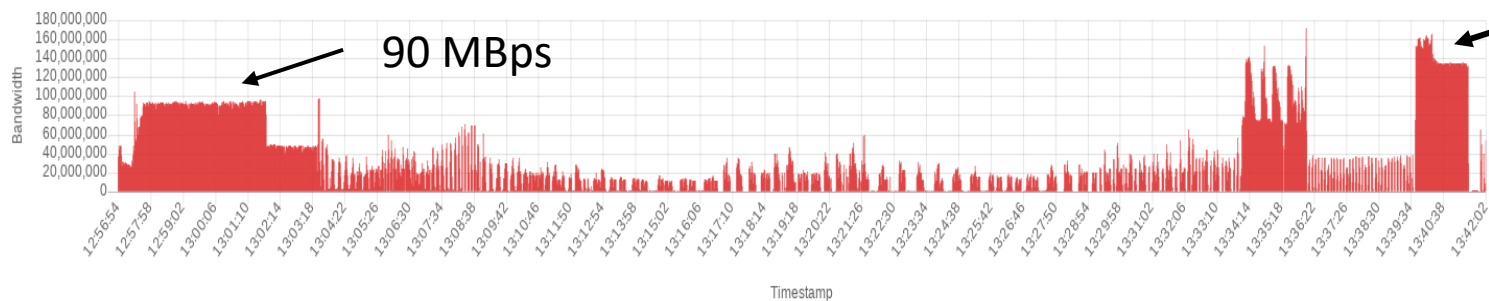


# Read bandwidth: m5.xlarge vs c5.xlarge

## Read bandwidth for mx.large 16GB



## Read bandwidth for cx.large 8GB



I/O limited  
by running  
out of AWS  
burst credits  
at the end





# How long did this work take?

Sizing the storage and compute correctly took three days  
... and we saved 10-40% of cloud costs for the project.

“Improving run time often doesn't require extensive rewrites.  
Knowing where to look is key.”

Keiran Raine, Cancer researcher, Sanger Institute



Ellexus: The I/O Profiling Company  
[www.ellexus.com](http://www.ellexus.com)



## Dependency hygiene flow at Qualcomm

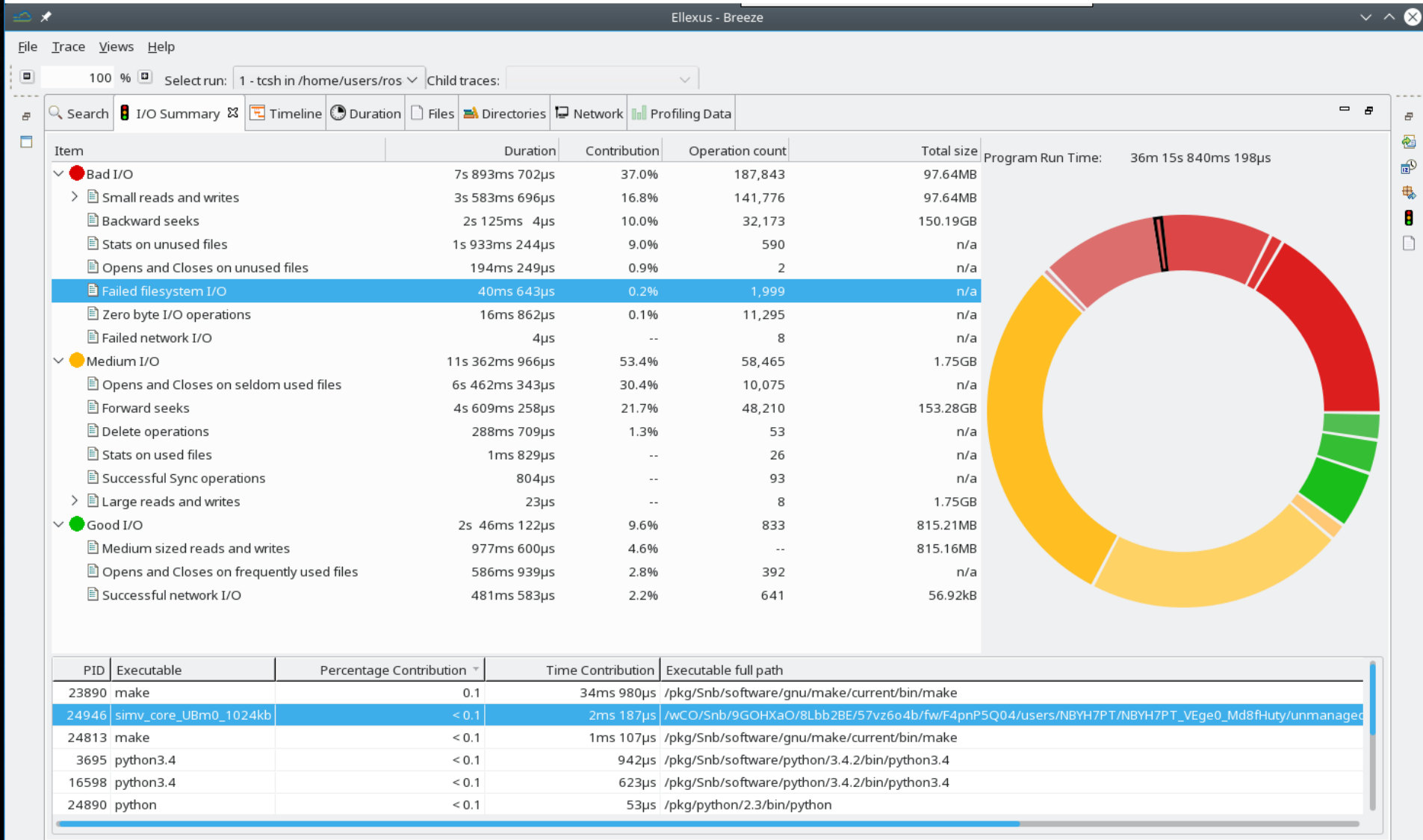
Breeze is used to trace thousands of workflows to automatically identify the mount points, file and network dependencies of each flow for migration.

### Disclaimer

The following trace was collected at Qualcomm, tracing a Synopsys VCS flow, but all identifiers and data have been modified or removed. No conclusions can be drawn from the following screenshots about the IT infrastructure, tools or usage at Qualcomm or Synopsys.



# Dependency hygiene flow at Qualcomm



# File dependencies, mounts points, packages for migration and containerization

Ellexus - Breeze

FileTraceViewsHelp

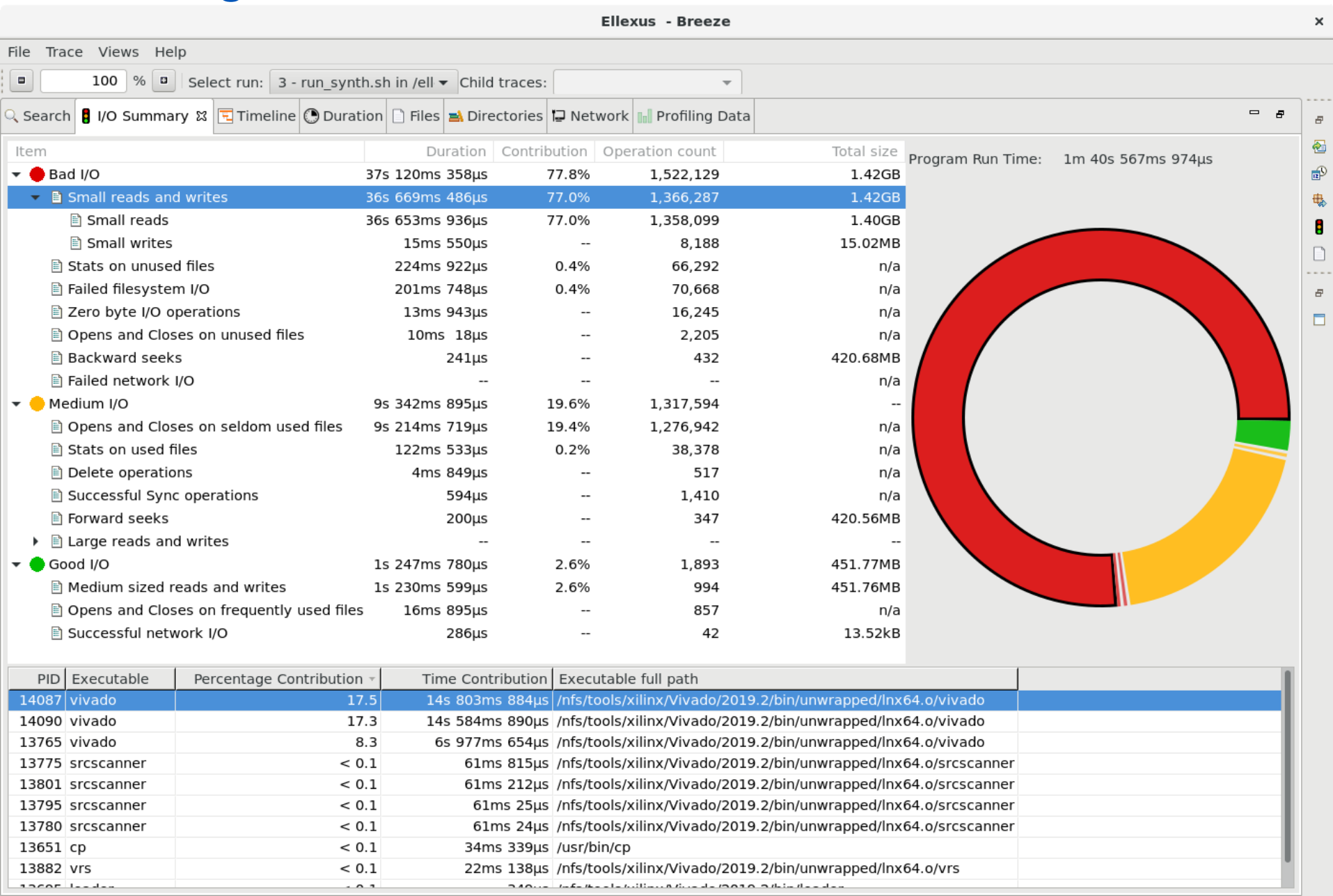
100 % Select run: 1 - tcsh in /home/users/ros Child traces:

Search I/O Summary Timeline Duration Files Directories Network Profiling Data

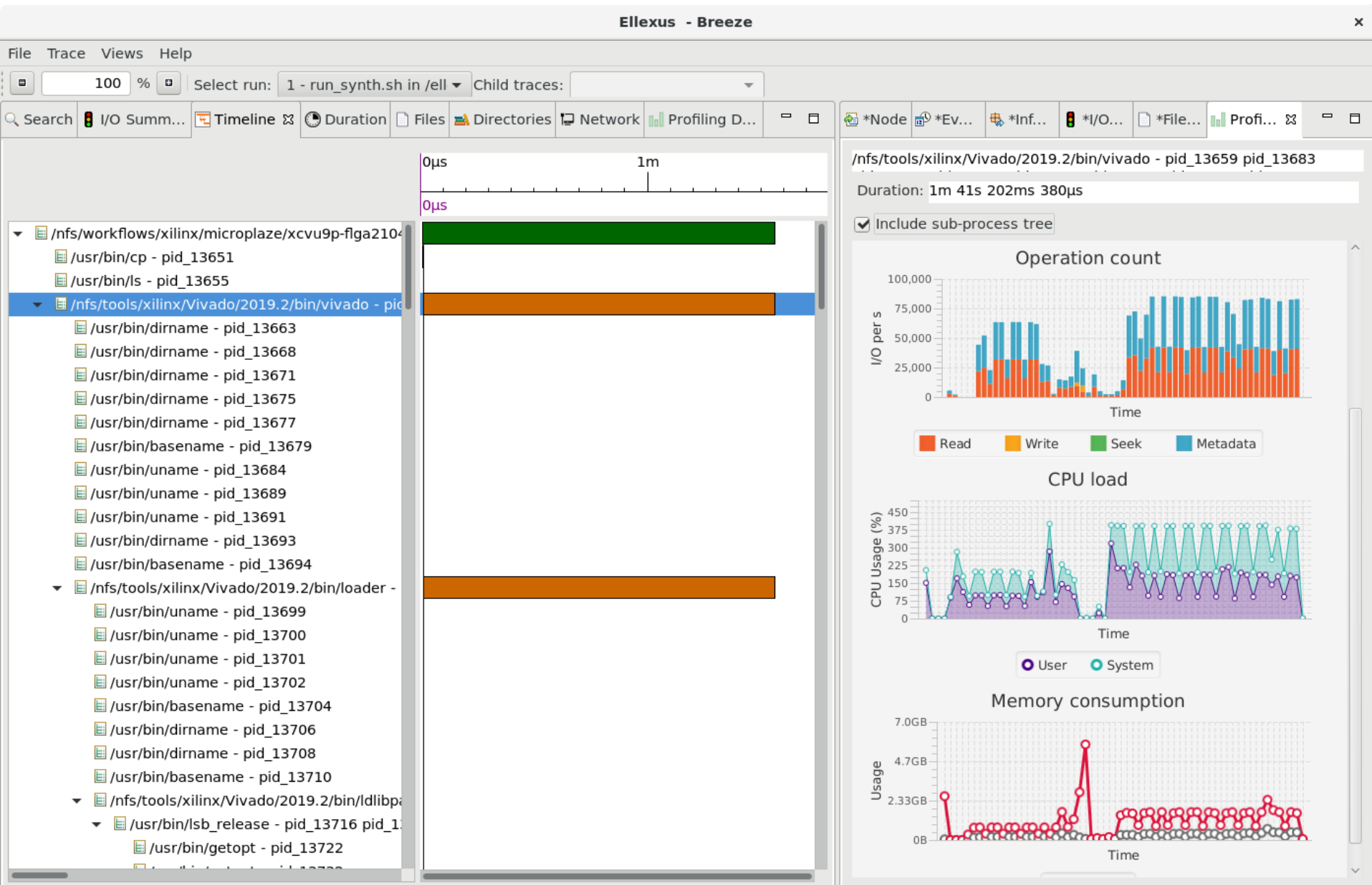
Mount Point	Location	Read	Small (< 32kB) Read	Large (≥ 100MB) Read	Write			
Mount Point	File Type	Filename	Full Path	Package	# Call	# Call	# Call	# Call
/	Data File	dumb	/usr/share/terminfo/d/dumb	terminfo-base-6.1-lp150.4.3...	7	7	0	0
/	Data File	xterm	/usr/share/terminfo/x/xterm	terminfo-base-6.1-lp150.4.3...	1	1	0	0
/	Data File	bindkey.tcsh	/etc/profile.d/bindkey.tcsh	tcsh-6.20.00-lp150.1.9.x86_64	5	4	0	0
/	Data File	complete.tcsh	/etc/profile.d/complete.tcsh	tcsh-6.20.00-lp150.1.9.x86_64	8	7	0	0
/	Data File	hosts.equiv	/etc/hosts.equiv	netcfg-11.6-lp150.1.1.noarch	2	1	0	0
/	Shared Library	libz.so.1	/lib64/libz.so.1	libz1-1.2.11-lp150.2.3.1.x86...	0	0	0	0
/	Shared Library	libxml2.so.2	/usr/lib64/libxml2.so.2	libxml2-2-2.9.7-lp150.2.6.1....	0	0	0	0
/	Shared Library	libselinux.so.1	/lib64/libselinux.so.1	libselinux1-2.6-lp150.2.14.x...	0	0	0	0
/	Shared Library	libgthread-2.0.so.0	/usr/lib64/libgthread-2.0.so.0	libgthread-2_0-0-2.54.3-lp1...	0	0	0	0
/	Shared Library	libgobject-2.0.so.0	/usr/lib64/libgobject-2.0.so.0	libgobject-2_0-0-2.54.3-lp15...	0	0	0	0
/	Shared Library	libglib-2.0.so.0	/usr/lib64/libglib-2.0.so.0	libglib-2_0-0-2.54.3-lp150.3...	0	0	0	0
/	Shared Library	libattr.so.1	/lib/libattr.so.1	libattr1-32bit-2.4.47-lp150.2...	0	0	0	0
/	Shared Library	libattr.so.1	/lib64/libattr.so.1	libattr1-2.4.47-lp150.2.16.x...	0	0	0	0
/	Shared Library	libacl.so.1	/lib/libacl.so.1	libacl1-32bit-2.2.52-lp150.3...	0	0	0	0
/	Shared Library	libacl.so.1	/lib64/libacl.so.1	libacl1-2.2.52-lp150.3.3.1.x8...	0	0	0	0
/	Shared Library	librt.so.1	/lib/librt.so.1	glibc-32bit-2.26-lp150.11.9...	0	0	0	0
/	Shared Library	libutil.so.1	/lib/libutil.so.1	glibc-32bit-2.26-lp150.11.9...	0	0	0	0
/	Shared Library	libcrypt.so.1	/lib/libcrypt.so.1	glibc-32bit-2.26-lp150.11.9...	0	0	0	0
/	Shared Library	ld-linux.so.2	/lib/ld-linux.so.2	glibc-32bit-2.26-lp150.11.9...	0	0	0	0
/	Shared Library	libnsl.so.1	/lib/libnsl.so.1	glibc-32bit-2.26-lp150.11.9...	0	0	0	0
/	Shared Library	libdl.so.2	/lib/libdl.so.2	glibc-32bit-2.26-lp150.11.9...	0	0	0	0
/	Shared Library	libpthread.so.0	/lib/libpthread.so.0	glibc-32bit-2.26-lp150.11.9...	0	0	0	0
/	Shared Library	libc.so.6	/lib/libc.so.6	glibc-32bit-2.26-lp150.11.9...	0	0	0	0
/	Shared Library	libm.so.6	/lib/libm.so.6	glibc-32bit-2.26-lp150.11.9...	0	0	0	0
/	Shared Library	libnss_dns.so.2	/lib64/libnss_dns.so.2	glibc-2.26-lp150.11.9.1.x86_...	0	0	0	0
/	Shared Library	ld-linux-x86-64.so.2	/lib64/ld-linux-x86-64.so.2	glibc-2.26-lp150.11.9.1.x86_...	0	0	0	0
/	Shared Library	libutil.so.1	/lib64/libutil.so.1	glibc-2.26-lp150.11.9.1.x86_...	0	0	0	0
/	Shared Library	libnss_files.so.2	/lib64/libnss_files.so.2	glibc-2.26-lp150.11.9.1.x86_...	0	0	0	0
/	Shared Library	librt.so.1	/lib64/librt.so.1	glibc-2.26-lp150.11.9.1.x86_...	0	0	0	0
/	Shared Library	libc.so.6	/lib64/libc.so.6	glibc-2.26-lp150.11.9.1.x86_...	0	0	0	0
/	Shared Library	libcrypt.so.1	/lib64/libcrypt.so.1	glibc-2.26-lp150.11.9.1.x86_...	0	0	0	0
/	Shared Library	libresolv.so.2	/lib64/libresolv.so.2	glibc-2.26-lp150.11.9.1.x86_...	0	0	0	0
/	Shared Library	libpthread.so.0	/lib64/libpthread.so.0	glibc-2.26-lp150.11.9.1.x86_...	0	0	0	0



# Profiling EDA tools: Xilinx Vivado



# Profiling EDA tools: Xilinx Vivado





# I/O profiling to improve DL\_POLY with the Hartree Centre

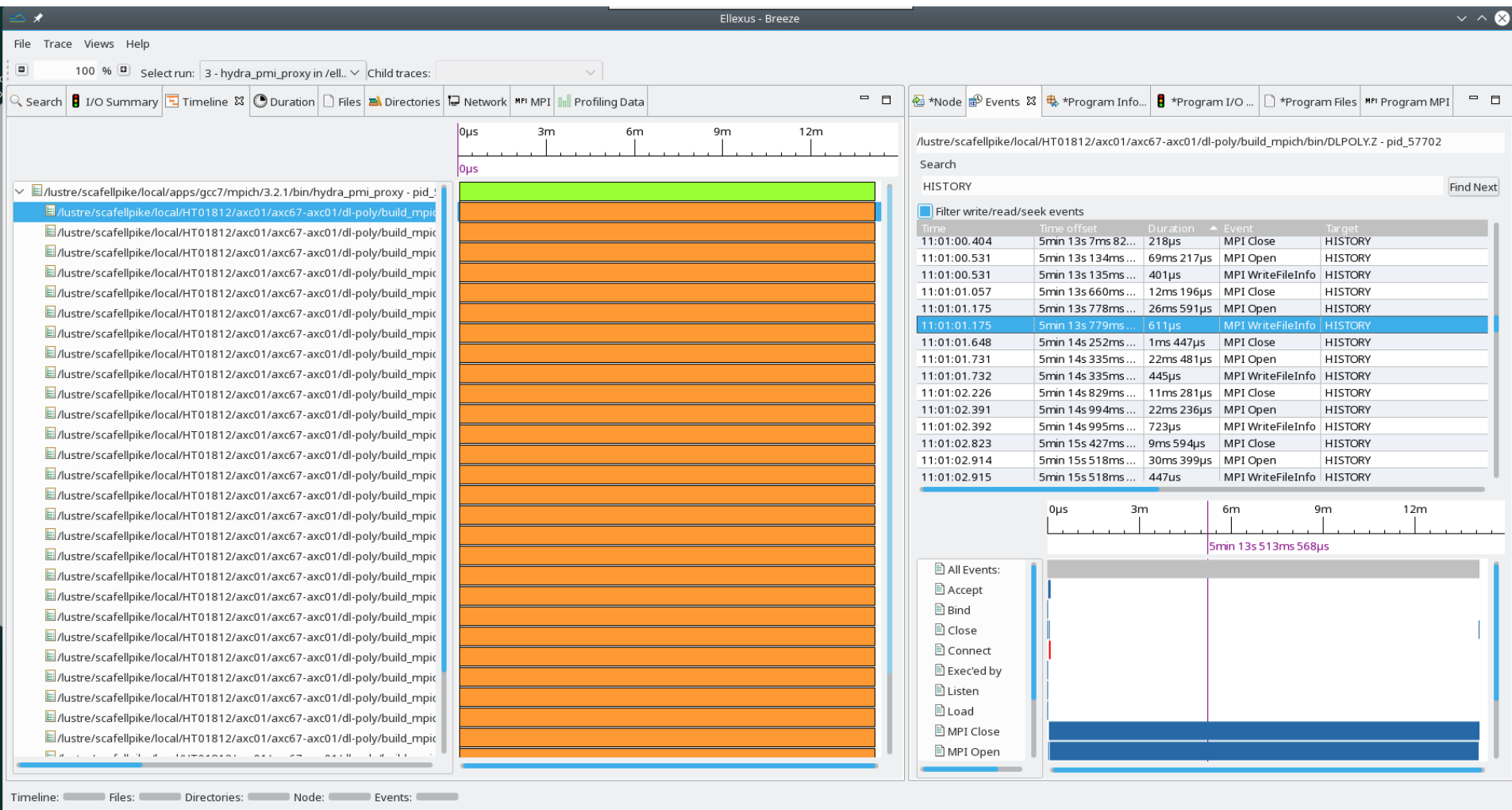
DL\_POLY is a general purpose classical molecular dynamics tool that uses MPI I/O.



Ellexus: The I/O Profiling Company  
[www.ellexus.com](http://www.ellexus.com)

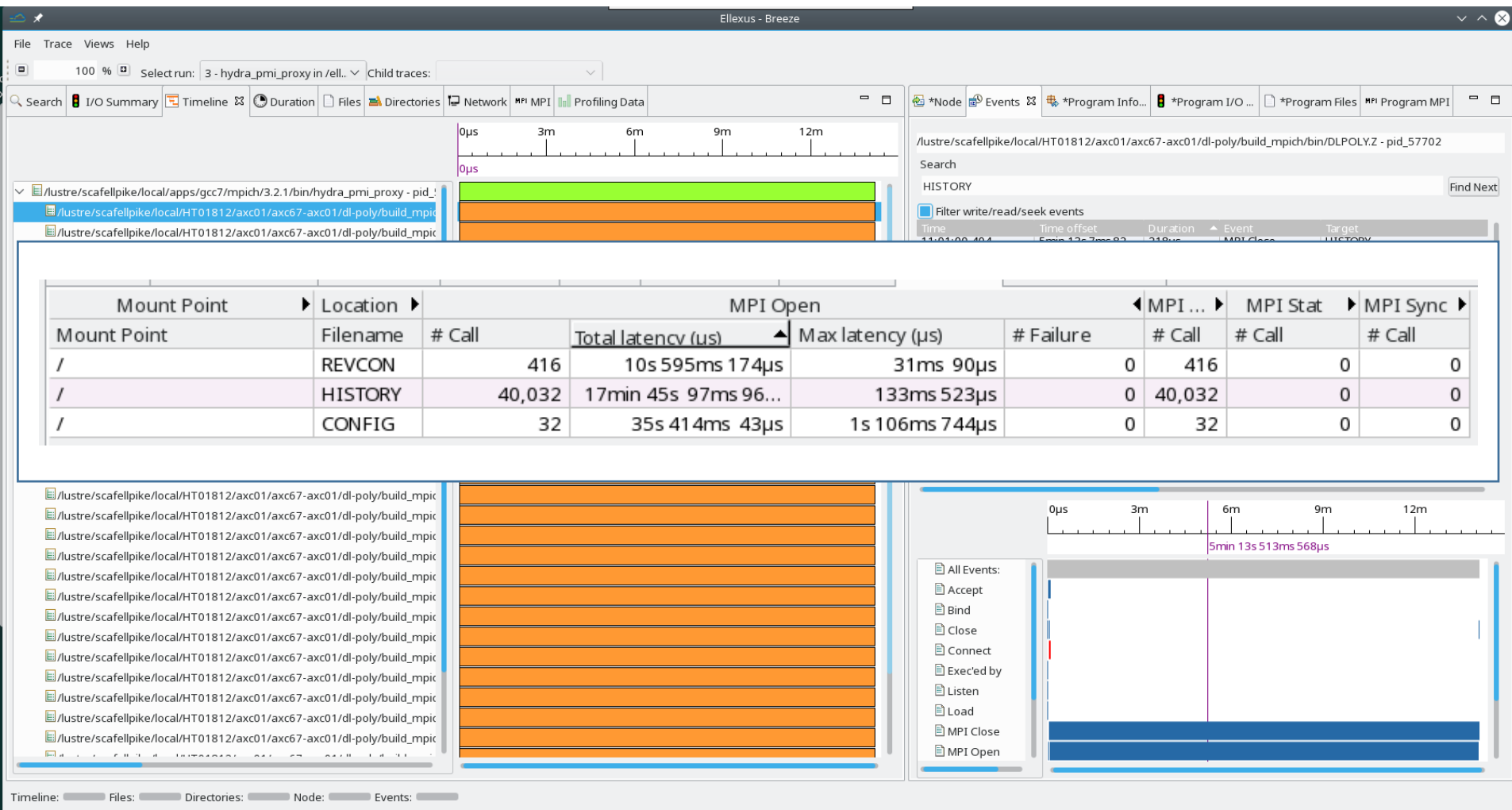
# I/O profiling to improve DL\_POLY with the Hartree Centre

DL\_POLY is a general purpose classical molecular dynamics tool that uses MPI I/O.  
It was opening the HISTORY file from every rank, but only writing from one.

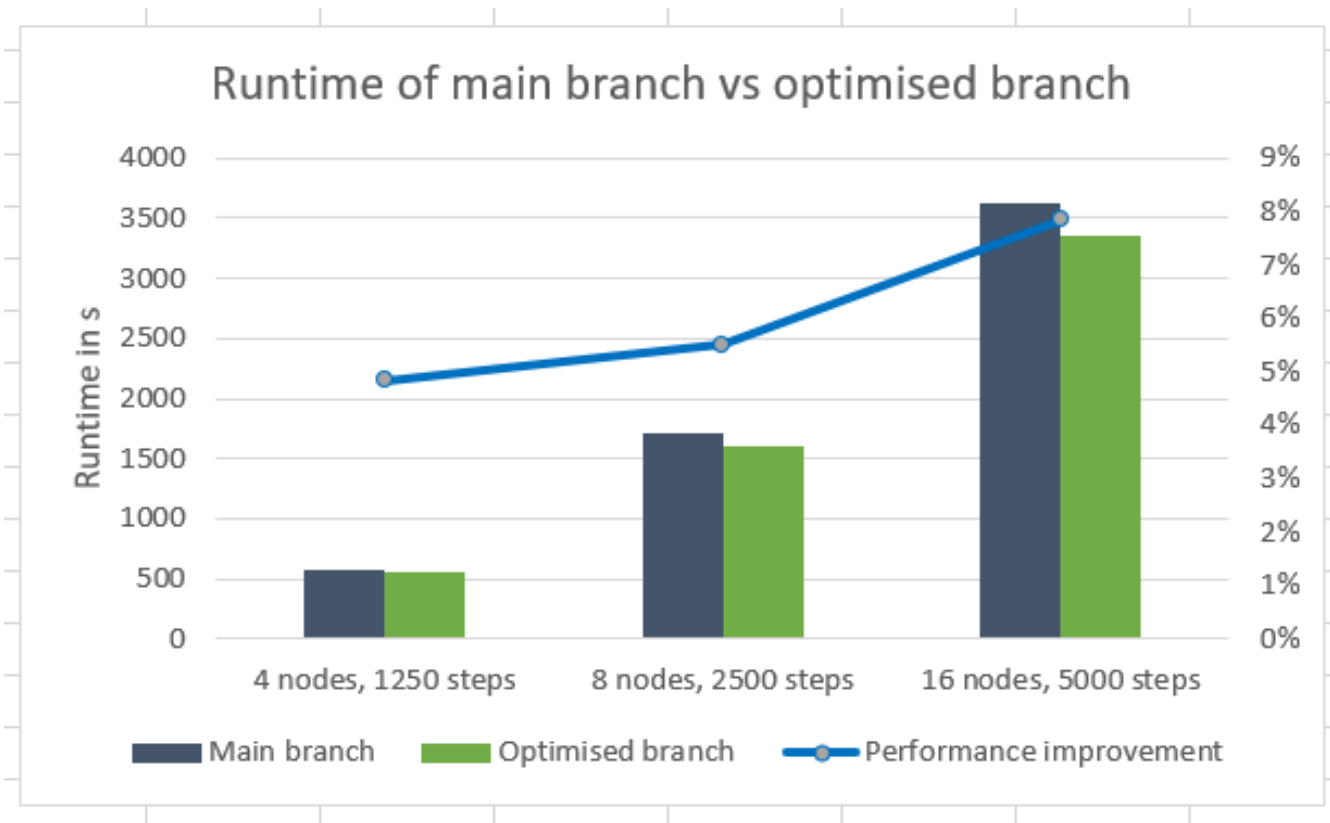


# I/O profiling to improve DL\_POLY with the Hartree Centre

DL\_POLY is a general purpose classical molecular dynamics tool that uses MPI I/O.  
It was opening the HISTORY file from ever rank, but only writing from one.



# I/O profiling to improve DL\_POLY with the Hartree Centre



Removing unnecessary opens gave significant performance improvements



# I/O profiling and what to do with the results

Profile in production

Optimization

Steering

- Data location, Scheduling, Filesystem, Burst buffers



The I/O Profiling Company - Protect. Balance. Optimise.

[www.ellexus.com](http://www.ellexus.com)

# Ellexus Ltd: The I/O Profiling Company

*Dr Rosemary Francis, CEO, Good I/O evangelist*

## Thanks for listening



The I/O Profiling Company - Protect. Balance. Optimise.

[www.ellexus.com](http://www.ellexus.com)