



Modern Dataflow in Experimental Nuclear Science (and Tcl).

Ron Fox, Giordano Cerizza
Sean Liddick, Aaron Chester

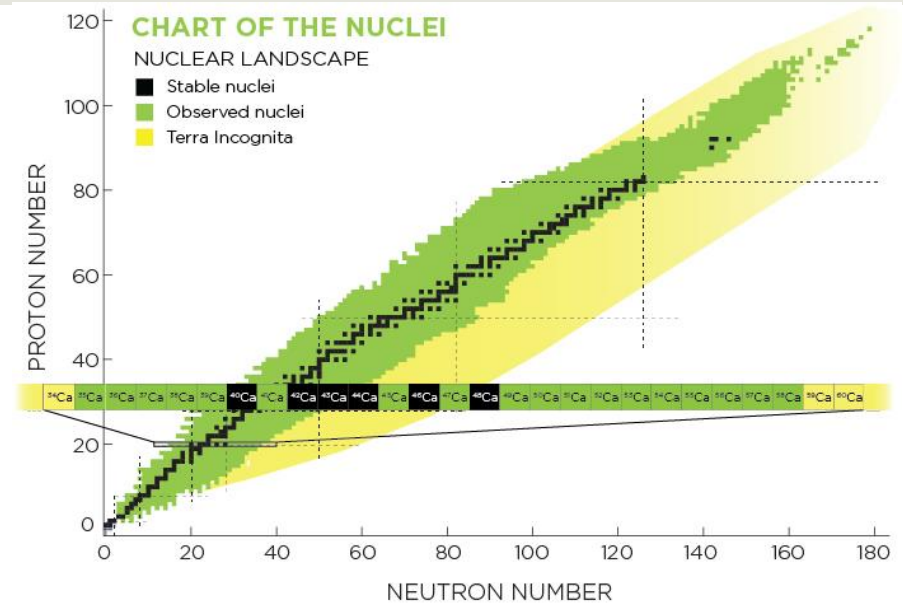
Talk Outline

- A bit about me and my Tcl history
- What is the National Superconducting Cyclotron Laboratory (NSCL)
- How data taking has evolved in experimental nuclear science
- E17011 an experiment with modern electronics – why it's computationally demanding
- Parallel resources available to us
- Message Passing Interface (MPI) and Tcl
 - Intro to MPI
 - Existing Tcl support
 - Tcl-Ish support we did.
- Applying MPITcl to an existing application
- What this means for experimental nuclear science at the NSCL

Tcl and me.

- Introduced Tcl/Tk at the National Superconducting Cyclotron Lab (NSCL) back in the 4.x days.
- Plugged into the community with a talk in New Orleans (Tcl 2004)
 - <https://www.tcl.tk/community/tcl2004/Papers/RonFox/>
 - NSCLSpecTcl – Histogramming package for experimental nuclear science.
- Tcl/Tk conference proceedings editor from Tcl2005 and on if memory serves.
- Tcl plays an important role in the NSCL experimental program.

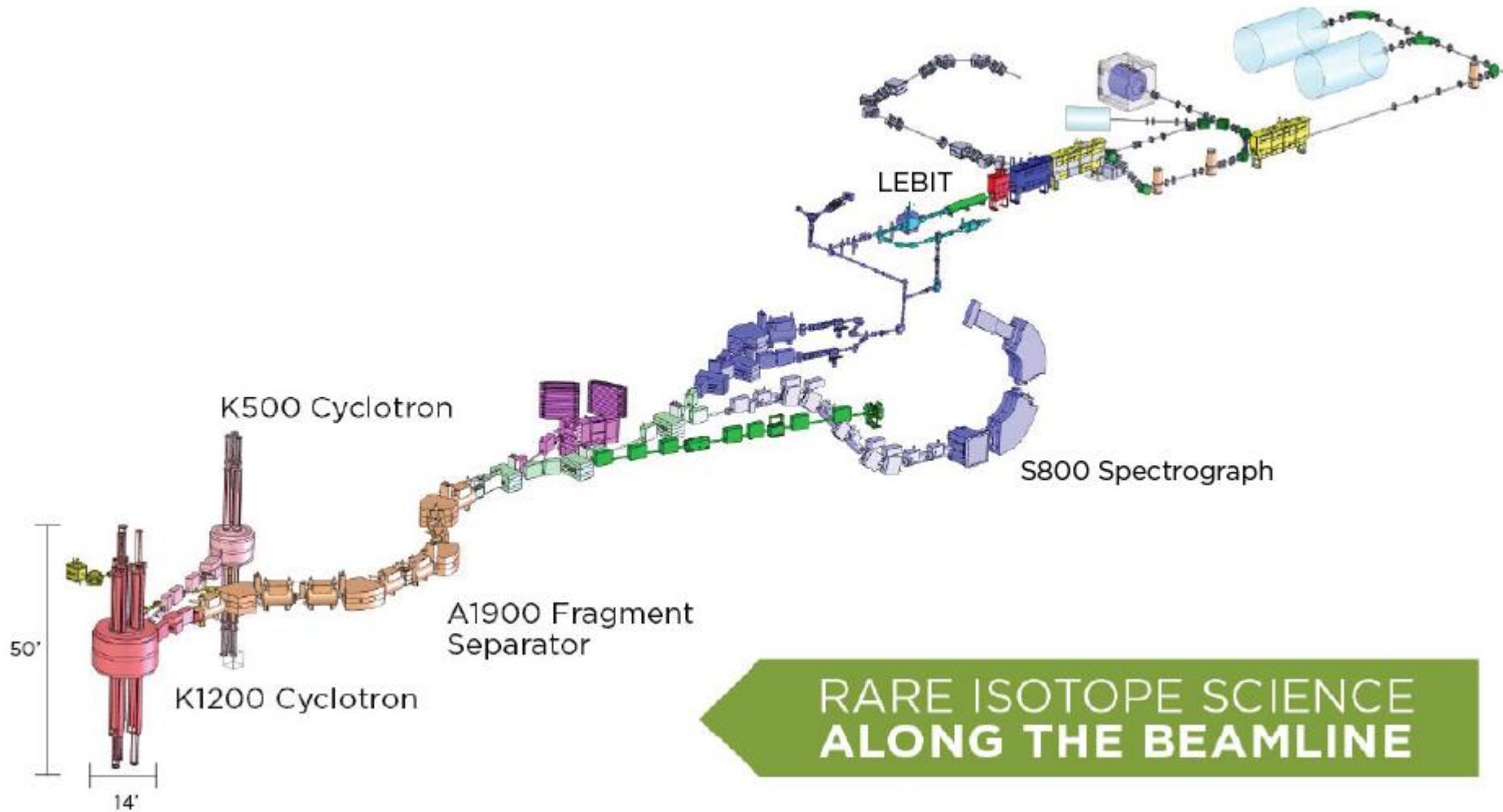
The National Superconducting Cyclotron Lab.



- Located at Michigan State University
- Funded by the National Science Foundation as a user facility

- Explore the properties of nuclear unstable nuclei
- Why and how do certain isotopes form.
- Where do the heavy elements come from?
- <http://www.nscl.msu.edu>

NSCL Block Diagram



**RARE ISOTOPE SCIENCE
ALONG THE BEAMLINE**

Science drivers for Rare Isotope Research

Science drivers (thrusts) from NRC RISAC 2007

Nuclear Structure	Nuclear Astrophysics	Tests of Fundamental Symmetries	Applications of Isotopes
-------------------	----------------------	---------------------------------	--------------------------

Intellectual challenges from NRC Decadal Study 2013

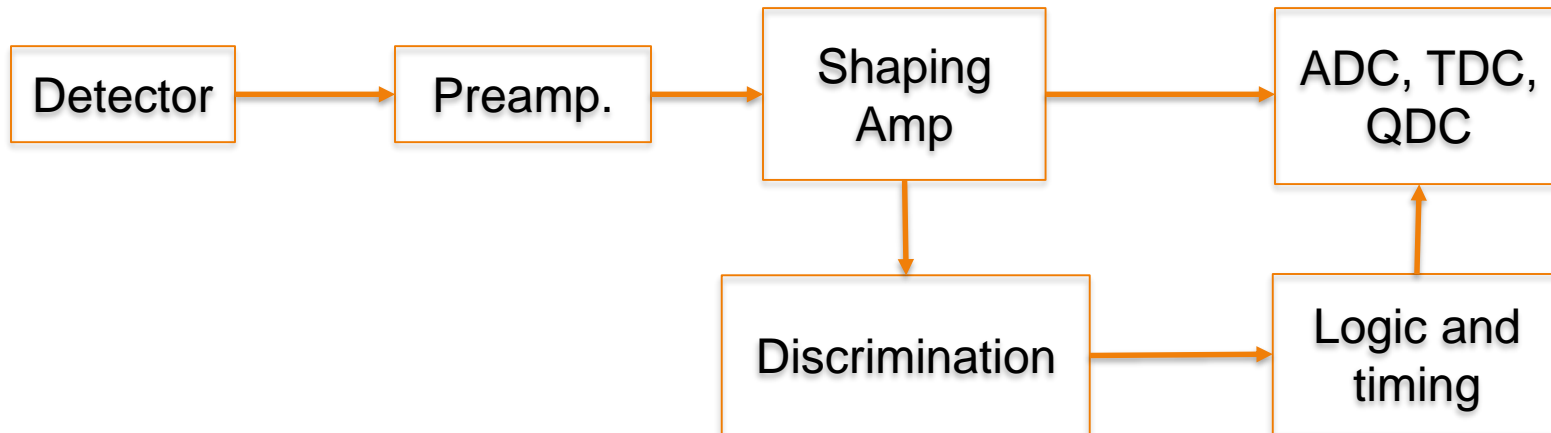
How does subatomic matter organize itself and what phenomena emerge?	How did visible matter come into being and how does it evolve?	Are fundamental interactions that are basic to the structure of matter fully understood?	How can the knowledge and technological progress provided by nuclear physics best be used to benefit society?
--	--	--	---

Overarching questions from NSAC Long Range Plan 2015

How are nuclei made and organized? What is the nature of dense nuclear matter?	Where do nuclei and elements come from? What combinations of neutrons and protons can form a bound atomic nucleus? How do neutrinos affect element synthesis?	Are neutrinos their own antiparticles? Why is there more matter than antimatter in the present universe?	What are practical and scientific uses of nuclei?
---	---	---	---

Overarching questions are answered by rare isotope research

Data Acquisition – old school (analog)



Important point – dead-times for a conversion are microseconds

Data Acquisition – old school (analog)

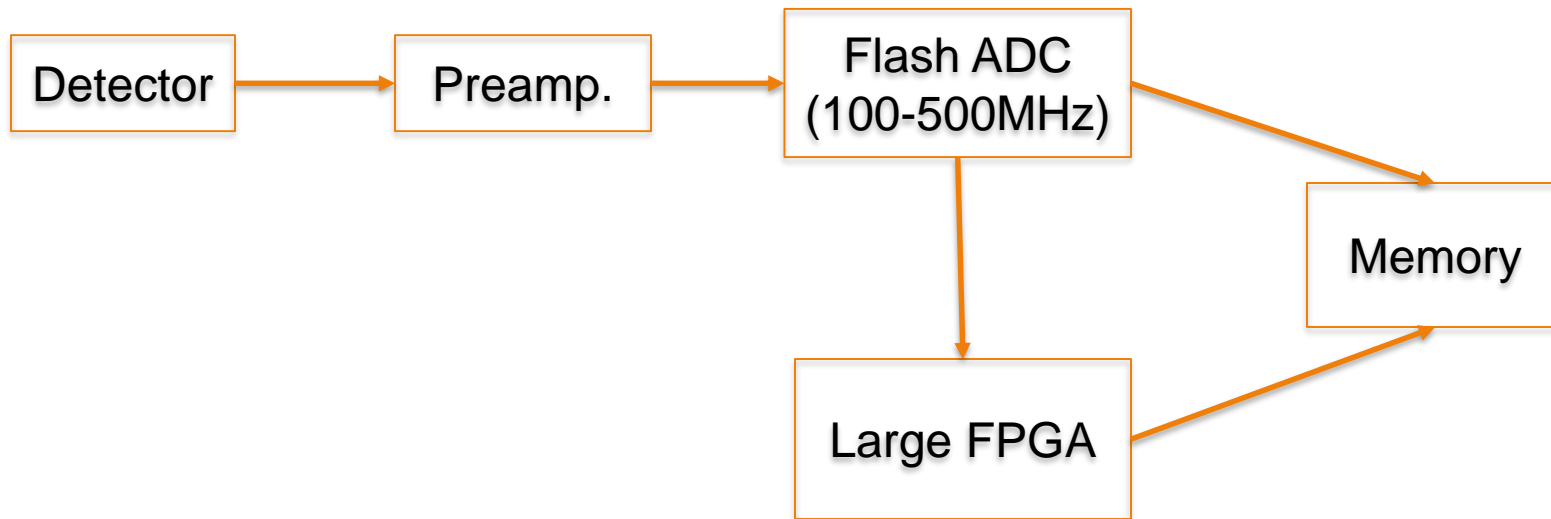


- Detector signals
- Pre-amplification
- Shaping/amplification
- Timing/triggering
- Digitizing modules

Each digitizing module
Gives one value per input:

- Pulse height
- Pulse charge integration
- Pulse timing relative to some reference time.

Modern Data Acquisition (digital)



Modern data acquisition (100MHz – 500MHz)



- Detector Signals
- Preamplification
- Digitization

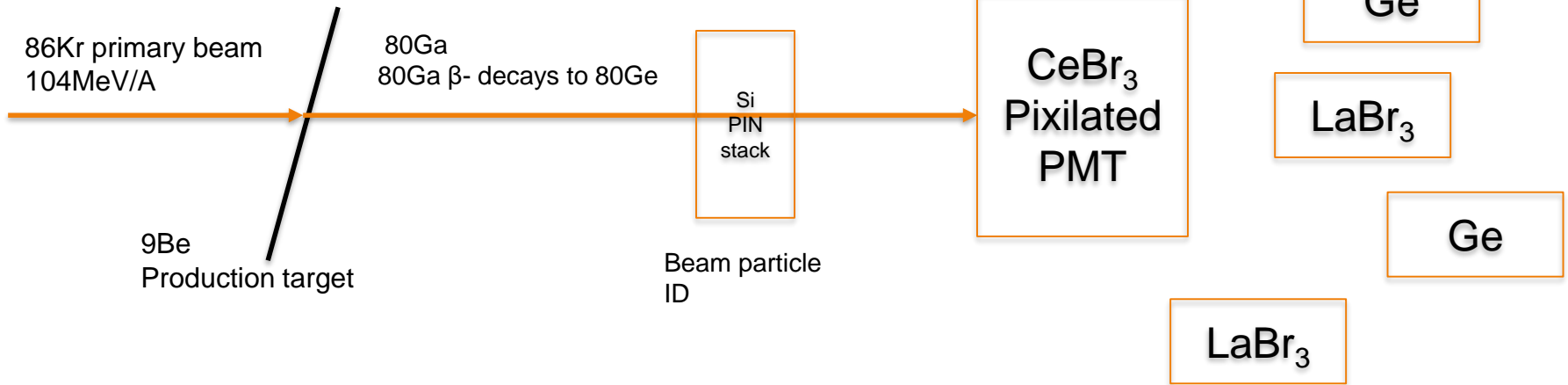
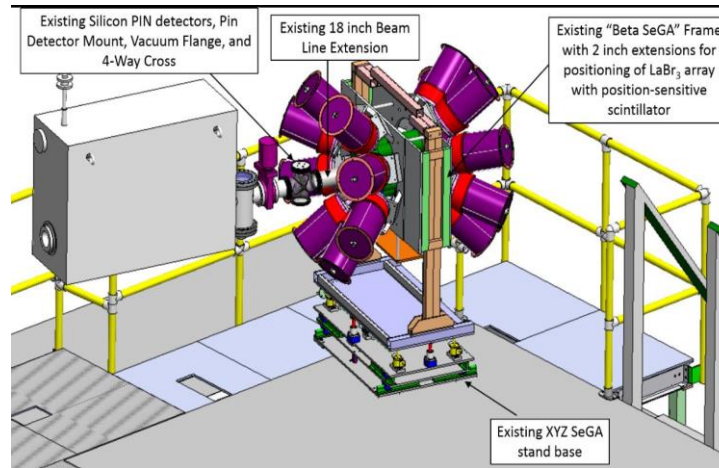
- Firmware can extract
 - Pulse ht.
 - Charge integral
 - Timing
- Keeping waveforms allows experiments that can't be done with analog electronics.
- Wave form analysis is computationally demanding
 - Wave forms bloat the data

E17011

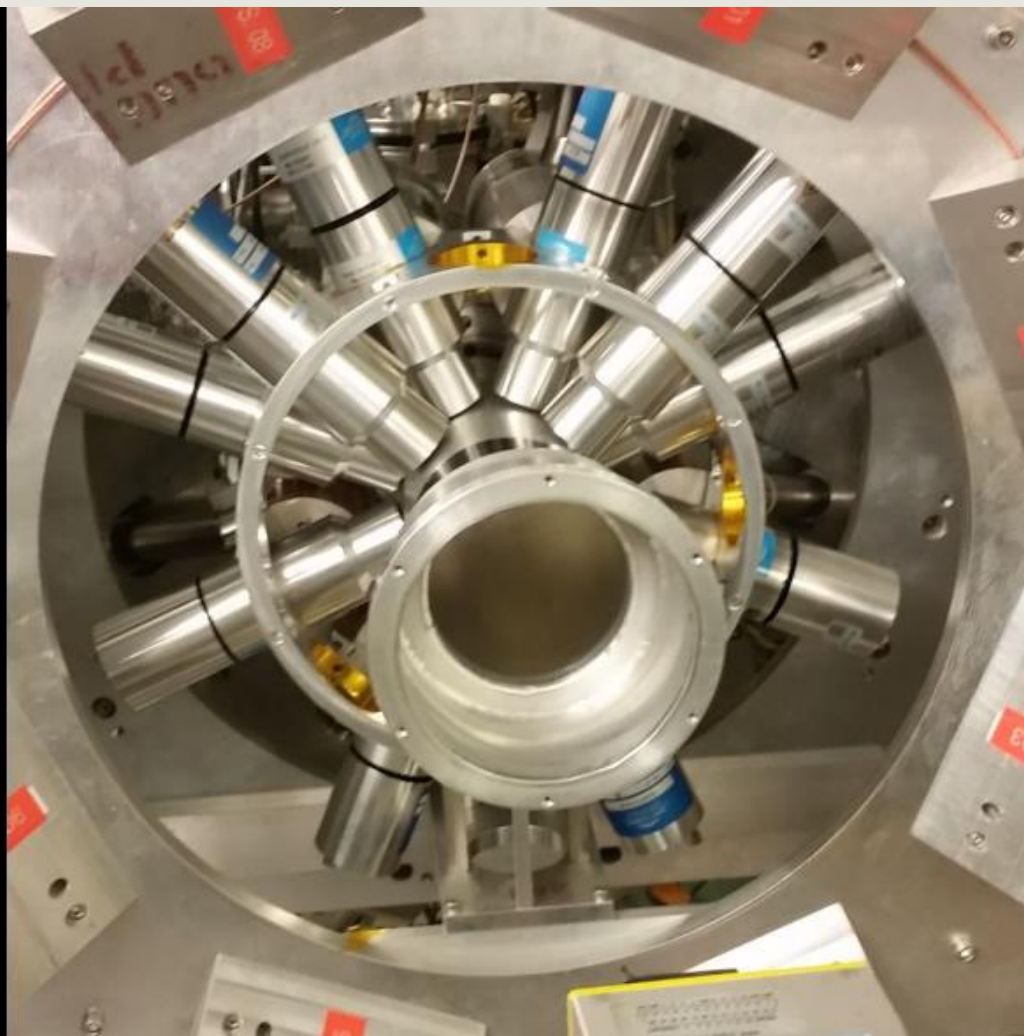
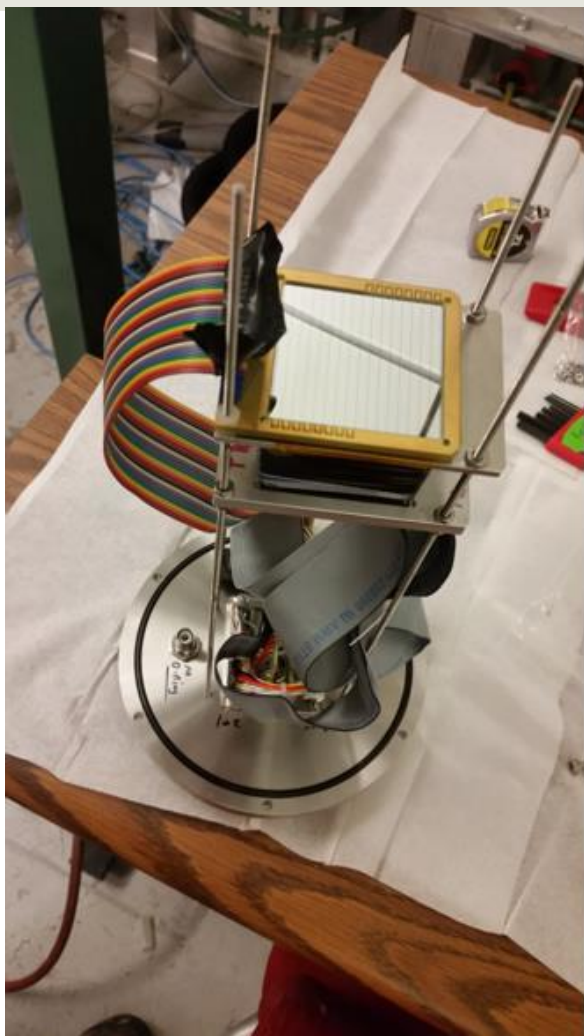
- Scheduled to run in January.
 - Look at beta decay of $80\text{Ga} \rightarrow 80\text{Ge}$
 - Look at the lifetime of the $0_2^+ \rightarrow 0_1^+$
 - Lifetime tell us something about the difference in the radius of the charge distribution of the two states.
- 200MB/second sustained – though modest trigger rate ($\sim 3\text{KHz}$).
- Will take 100TB+ of data
- Need good online and nearline analysis:
 - Are the detectors working.
 - Are we seeing what we think we should be seeing.
 - Should we ask for additional (discretionary time).

E17011 – block diagram

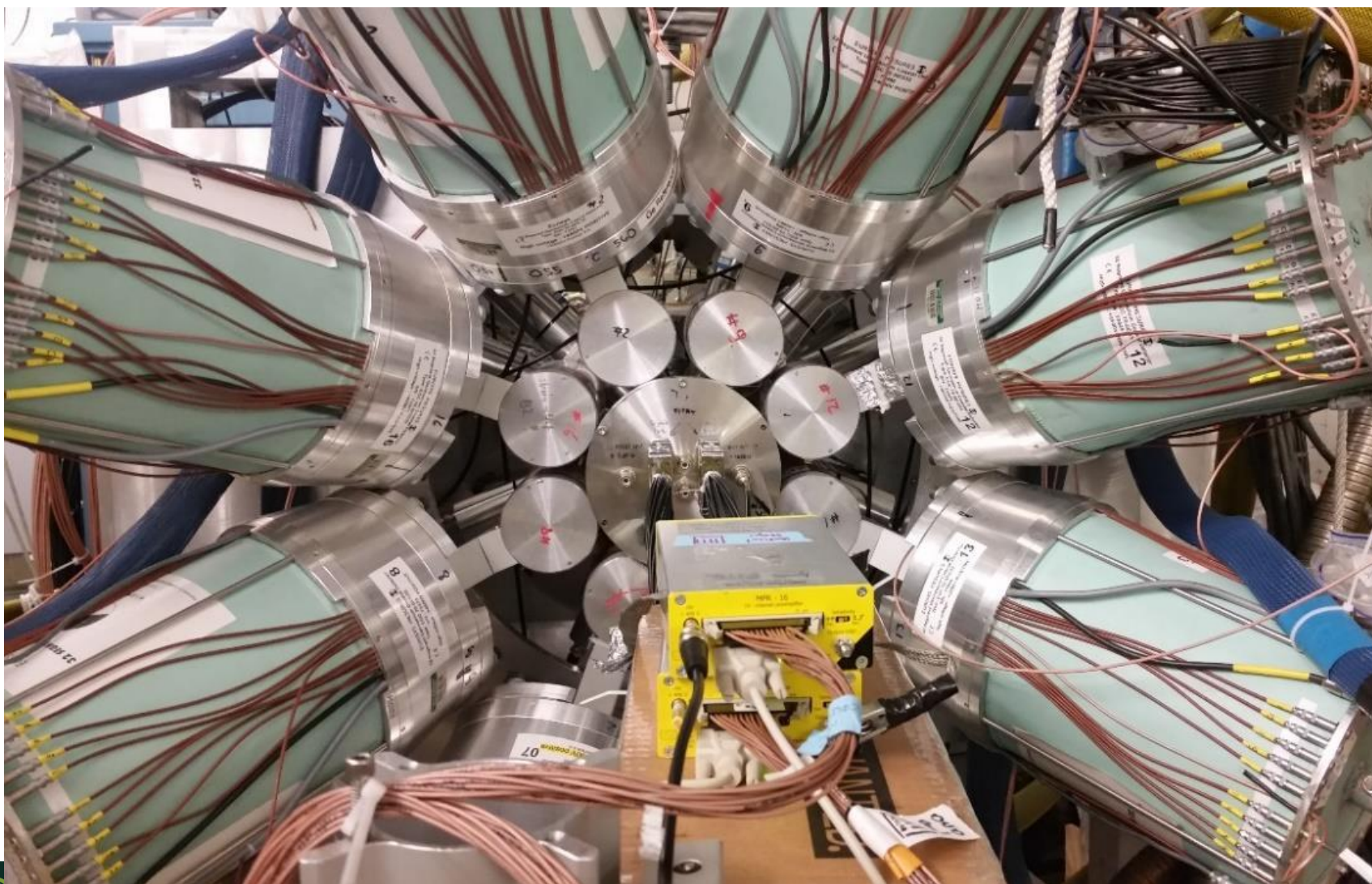
Sketch of experiment



Pictures pictures (CeBr_3 and LaBr_3 array)



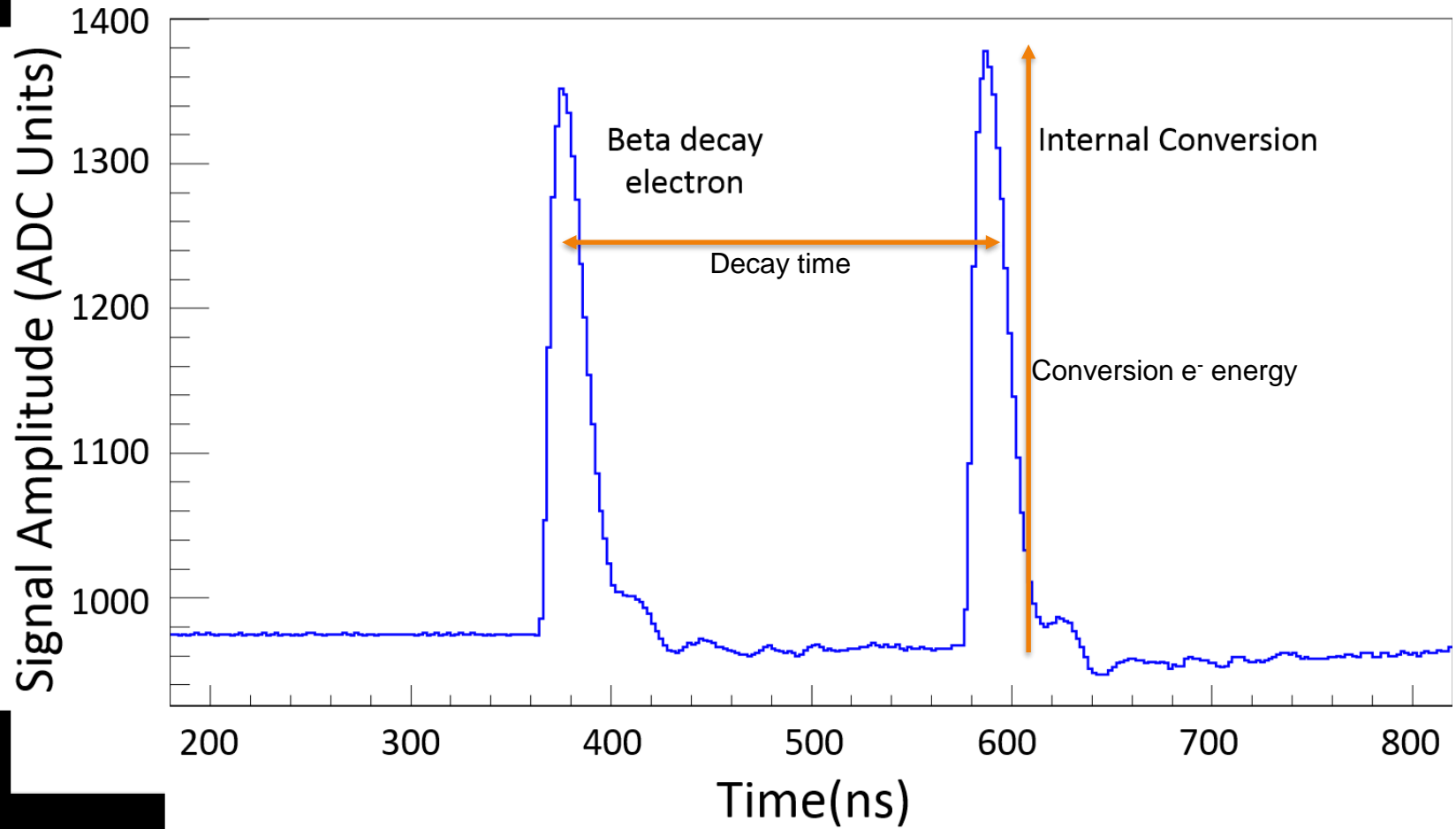
More pictures Ge Array (SeGA)



What happens to the implanted ions.

- ^{80}Ga decays to ^{80}Ge by β^- decay.
 - This decay is also detected in the CeBr_3 detector
 - This decay populates several energy levels of ^{80}Ge
- Of interest are the decays that populate the 0_2^+ state.
 - This eventually de-excites to the 0_1^+ state emitting a γ -ray (detected by the LaBr_3 array and/or SeGA) and a conversion electron.
 - The conversion electron produced by that decay is sensed by the CeBr_3
- Well it's not actually eventually.
 - Similar de-excitations have half lives of about 50ns.
 - We want the actual $\frac{1}{2}$ life.
- This is a short $\frac{1}{2}$ life. How to measure it.
 - Digitize the pulses in the CeBr_3
 - » Sum signal at 500MHz
 - » pixels at 250MHz
 - » Trace lengths of a few microseconds (on order 100 samples).

Sample trace from a similar experiment

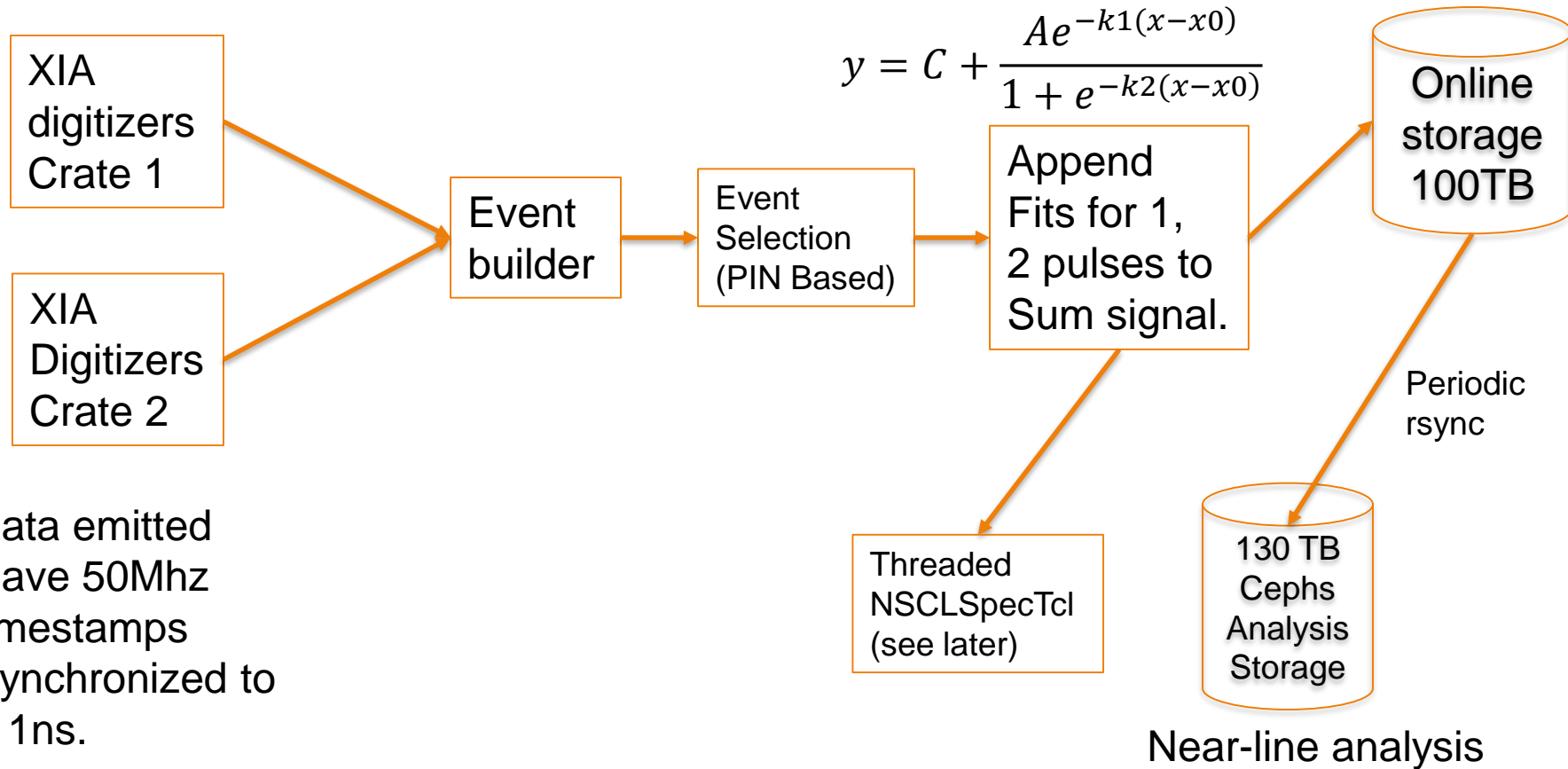


Where does that 200MB/sec come from?

- Since most of the CeBr3 detector lights up for a hit we about 200traces/event (maximal pixel is 'where' the event occurred).
- The data rate is dominated by traces from the CeBr3.
- Trigger rates may be 3KHz (modest)
- Data transfer rates will be a sustained 200MB/seconds.
- To see if the experiment is “working” we need to do some processing on all this stuff.
 - Determine if traces are single or double pulses.
 - Determine the characteristics of the pulse(s) – time and height.
- Good news though: Taking traces meas we can do the experiment.

This experiment is *really hard* to do
with old school electronics.

Data Flow:



Data emitted
Have 50Mhz
timestamps
Synchronized to
< 1ns.

Online analysis

- Fit the sum traces from the CeBr3.
 - Fit for both single and double pulses.
 - Use a heuristic to determine if the pulses are single or double.
- Make a pile of histograms (NSCLSpecTcl) and look at them online
- Keep up with the incoming data rate.

NOTE: Each fit costs 3.5ms to do using GSL's Levenberg-Marquardt.

Serial code isn't going to cut it.

Near-line Analysis – want to keep up with incoming data rate or better

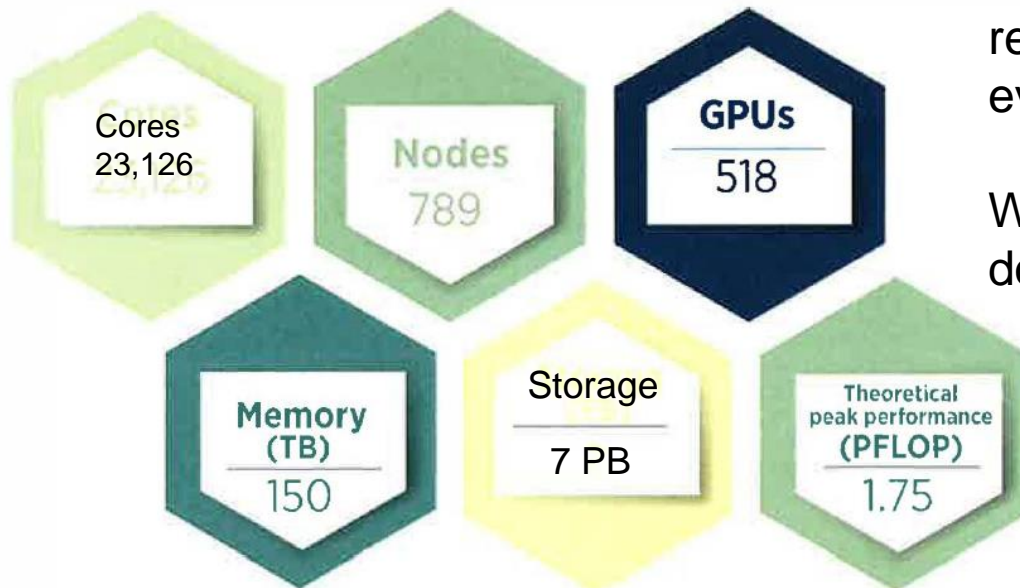
- Fit the remaining traces in the CeBr_3
 - Are they single or double pulses (heuristic)?
 - If double pulses extract the time difference as a parameter for histogramming.
- Correlate implantation events with decay events.
 - Using position and particle ID information
 - Timing between implantation and decay.
- These are computationally intensive (e.g. the fit is about 3.5ms/event). To make decisions about the experiment we need to analyze the data already taken faster than acquisition.
- Serial code isn't going to cut it ~2500 cores just for fitting all traces.

Parallel resources at the NSCL available to E17011

- Three high core count systems:
 - 1 26 core system. (Xeon E5-2690 v4 @ 2.60GHz)
 - 2 40 core systems (Xeon Gold 6148 @ 2.4GHz) – bought for this experiment
 - Used for online data flow and interactive ‘near-line’ analysis.
- Modest Linux cluster
 - 360 cores of various ages
 - Used for non-interactive ‘near-line’ partial analysis.
- That’s not going to be enough (to do the fitting of all signals at data rates needs about 2500cores).
- no GPU coprocessors ☹️

MSU Institute for Cyber Enabled Research (ICER)

HARDWARE

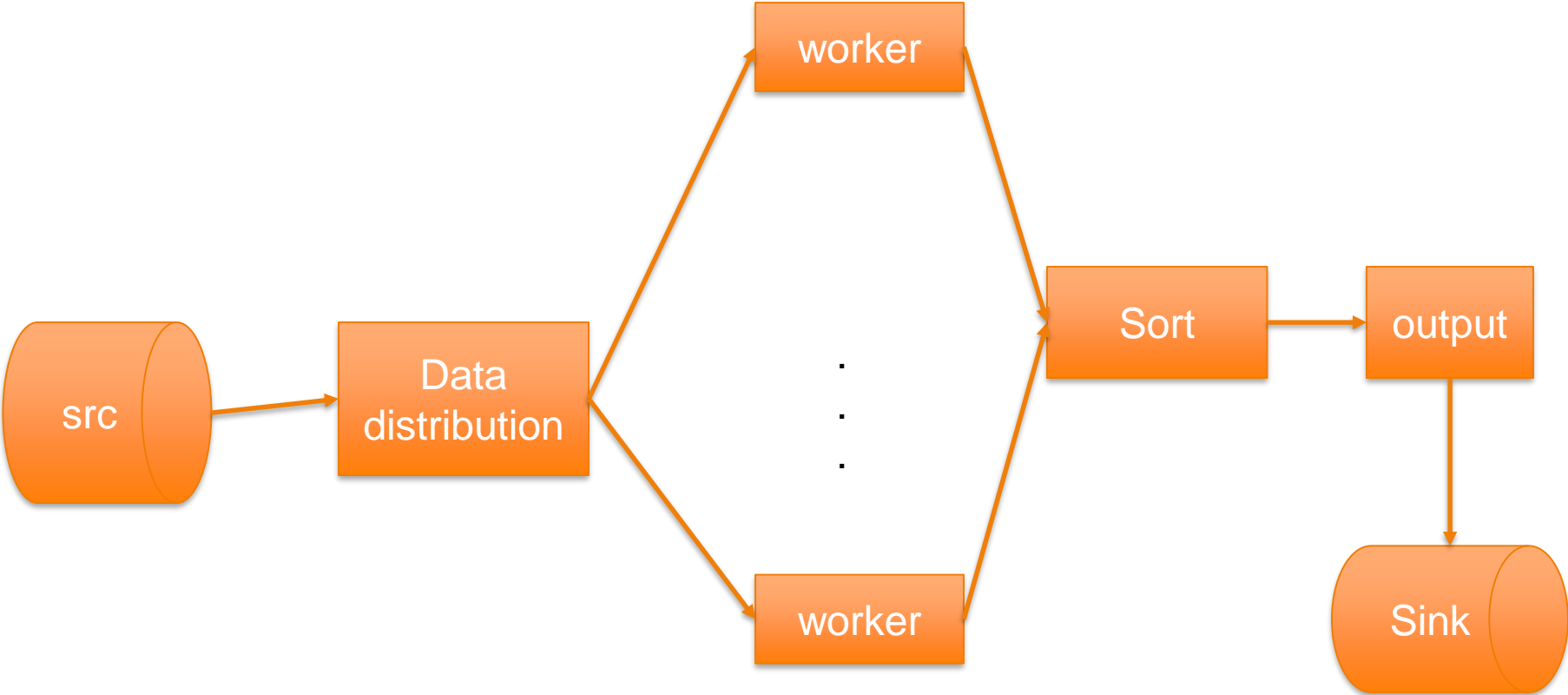


Naturally we've ~~listed after~~ sought ways to leverage this resource for near-line and maybe even online analysis.

Work to containerize our apps is done (thank you singularity)

Scheduling, however can be an issue: NSCL resources can be dedicated to E17011, ICER is shared across all university users.

Structure of event analysis parallel programs

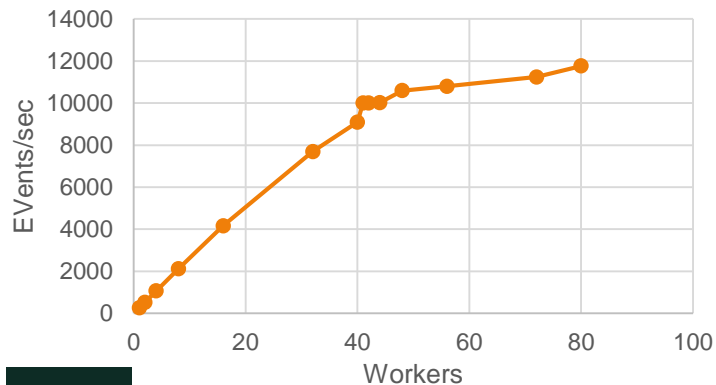


Meeting these needs.

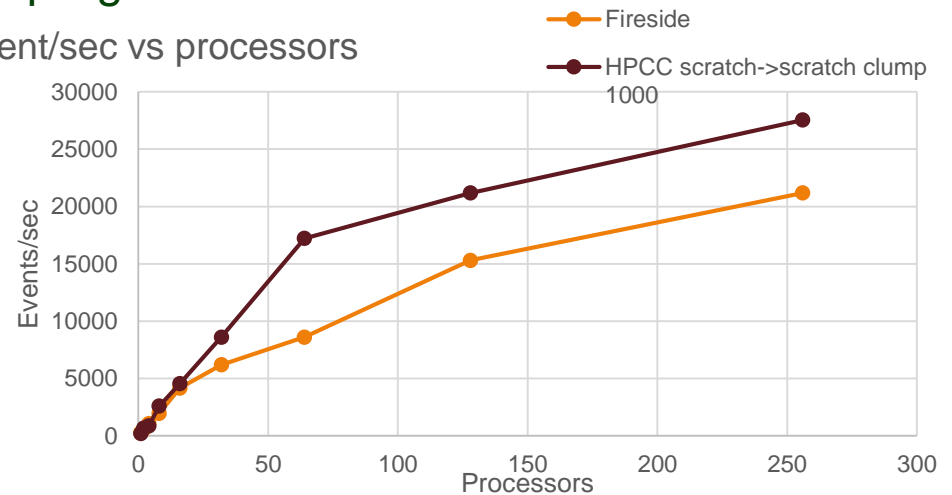
- Different types of parallelism
 - Threaded parallelism for the online/interactive stuff.
 - Distributed parallelism for near-line non-interactive stuff.
- Tools to make parallelization simpler
- Fitting:
 - Support for GPU 'accelerated' fitting residual and Jacobian computation ☹️
 - Machine learning for single/double pulse determination – most traces are single pulses

Example trace fitting the sum signal: same program threaded/cluster

Events/sec vs workers



Event/sec vs processors



MPI – cluster distributed parallel computing

- MPI - Message Passing Interface standard for writing distributed parallel programs.
 - OpenMPI <https://www.open-mpi.org/>
 - MPICH <https://www.mpich.org/>
- Multiple instances of the same process run in parallel.
 - Each process as a *rank* identifying it.
 - Processes can target messages to specific ranks.
 - Communicators can be formed to link groups of processes together.
- Messages require:
 - Rank – who we're sending to but:
 - Communicator – defines the process group in which the rank has meaning (MPI_COMM_WORLD – the entire application is pre-defined).
 - A Tag (integer)
 - Message data.
 - Type of data in the message (message data are strongly typed)
 - Number of items of that type being sent.

The MPI API is large and complex:

MPI	MPI_File_call_errhandler	MPI_Ineighbor_allgather	MPI_T_init_thread	MPI_Bcast	MPI_File_set_errhandler	MPI_Neighbor_alltoall
MPIX_Alloather_init	MPI_File_close	MPI_Ineighbor_allgather	MPI_T_pvar_get_info	MPI_Bsend	MPI_File_set_info	MPI_Neighbor_alltoallv
MPIX_Allgather_init	MPI_File_create_errhandler	MPI_Ineighbor_alltoall	MPI_T_pvar_get_num	MPI_Bsend_init	MPI_File_set_size	MPI_Neighbor_alltoallv
MPIX_Allreduce_init	MPI_File_delete	MPI_Ineighbor_alltoallv	MPI_T_pvar_handle_alloc	MPI_Buffer_attach	MPI_File_set_view	MPI_Op_c2f
MPIX_Alltoall_init	MPI_File_f2c	MPI_Ineighbor_alltoallv	MPI_T_pvar_handle_free	MPI_Buffer_detach	MPI_File_sync	MPI_Op_commutative
MPIX_Alltoallv_init	MPI_File_get_amode	MPI_Info_c2f	MPI_T_pvar_read	MPI_Cancel	MPI_File_write	MPI_Op_create
MPIX_Alltoallv_init	MPI_File_get_atomicsv	MPI_Info_create	MPI_T_pvar_readreset	MPI_Cart_coords	MPI_File_write_all	MPI_Op_f2c
MPIX_Bcast_init	MPI_File_get_byte_offset	MPI_Info_delete	MPI_T_pvar_reset	MPI_Cart_create	MPI_File_write_all_begin	MPI_Op_free
MPIX_Exscan_init	MPI_File_get_errhandler	MPI_Info_dup	MPI_T_pvar_session_create	MPI_Cart_get	MPI_File_write_all_end	MPI_Open_port
MPIX_Gather_init	MPI_File_get_group	MPI_Info_env	MPI_T_pvar_session_free	MPI_Cart_map	MPI_File_write_at	MPI_Pack
MPIX_Gatherv_init	MPI_File_get_info	MPI_Info_f2c	MPI_T_pvar_start	MPI_Cart_rank	MPI_File_write_at_all	MPI_Pack_external
MPIX_Neighbor_allgather_init	MPI_File_get_position	MPI_Info_free	MPI_T_pvar_stop	MPI_Cart_shift	MPI_File_write_at_all_begin	MPI_Pack_external_size
MPIX_Neighbor_alltoall_init	MPI_File_get_position_shared	MPI_Info_get	MPI_T_pvar_write	MPI_Cart_sub	MPI_File_write_at_all_end	MPI_Pack_size
MPIX_Neighbor_alltoallv_init	MPI_File_get_size	MPI_Info_get_nkeys	MPI_Test	MPI_Cartdim_get	MPI_File_write_ordered	MPI_Pcontrol
MPIX_Neighbor_alltoallv_init	MPI_File_get_type_extent	MPI_Info_get_nthkey	MPI_Test_cancelled	MPI_Close_port	MPI_File_write_ordered_begin	MPI_Probe
MPIX_Neighbor_alltoallv_init	MPI_File_get_view	MPI_Info_get_valuelen	MPI_Testall	MPI_Comm_accept	MPI_File_write_ordered_end	MPI_Publish_name
MPIX_Query_cuda_support	MPI_File_iread	MPI_Init	MPI_Testany	MPI_Comm_c2f	MPI_File_write_shared	MPI_Put
MPIX_Reduce_init	MPI_File_iread_all	MPI_Init_thread	MPI_Testsome	MPI_Comm_call_errhandler	MPI_Finalize	MPI_Query_thread
MPIX_Reduce_scatter_block_init	MPI_File_iread_at	MPI_Initialized	MPI_Topo_test	MPI_Comm_compare	MPI_Finalized	MPI_Raccumulate
MPIX_Reduce_scatter_init	MPI_File_iread_at_all	MPI_Intercomm_create	MPI_Type_c2f	MPI_Comm_connect	MPI_Free_mem	MPI_Recv
MPIX_Scan_init	MPI_File_iread_at_all	MPI_Intercomm_merge	MPI_Type_commit	MPI_Comm_create	MPI_Gather	MPI_Recv_init
MPIX_Scatter_init	MPI_File_iread_shared	MPI_Iprobe	MPI_Type_contiguous	MPI_Comm_create_errhandler	MPI_Gatherv	MPI_Reduce
MPIX_Scatterv_init	MPI_File_iwrite	MPI_Irecv	MPI_Type_create_darray	MPI_Comm_create_group	MPI_Get	MPI_Reduce_local
MPI_Abort	MPI_File_iwrite_all	MPI_Ireduce	MPI_Type_create_f90_comp	MPI_Comm_create_keyval	MPI_Get_accumulate	MPI_Reduce_scatter
MPI_Accumulate	MPI_File_iwrite_at	MPI_Ireduce_scatter	MPI_Type_create_f90_integ	MPI_Comm_delete_attr	MPI_Get_address	MPI_Reduce_scatter_block
MPI_Add_error_class	MPI_File_iwrite_at_all	MPI_Ireduce_scatter_block	MPI_Type_create_f90_real	MPI_Comm_disconnect	MPI_Get_count	MPI_Register_datarep
MPI_Add_error_code	MPI_File_iread	MPI_Irsend	MPI_Type_create_hindexed	MPI_Comm_dup	MPI_Get_elements	MPI_Request_c2f
MPI_Add_error_string	MPI_File_iread_all	MPI_Is_thread_main	MPI_Type_create_hindexed	MPI_Comm_dup_with_info	MPI_Get_elements_x	MPI_Request_f2c
MPI_Address	MPI_File_iread_at	MPI_Iscan	MPI_Type_create_hvector	MPI_Comm_f2c	MPI_Get_library_version	MPI_Request_free
MPI_Aint_add	MPI_File_iread_at_all	MPI_Iscatter	MPI_Type_create_indexed_b	MPI_Comm_free	MPI_Get_processor_name	MPI_Request_get_status
MPI_Aint_diff	MPI_File_iread_at_all_end	MPI_Iscatterv	MPI_Type_create_indexed	MPI_Comm_free_keyval	MPI_Get_version	MPI_Rget
MPI_Allgather	MPI_File_iread_at_all_end	MPI_Issend	MPI_Type_create_keyval	MPI_Comm_get_attr	MPI_Graph_create	MPI_Rget_accumulate
MPI_Alloc_mem	MPI_File_iread_at_all_end	MPI_Keyval_create	MPI_Type_create_keyval	MPI_Comm_get_errhandler	MPI_Graph_get	MPI_Rput
MPI_Allreduce	MPI_File_iread_at_all_end	MPI_Keyval_free	MPI_Type_create_keyval	MPI_Comm_get_info	MPI_Graph_map	MPI_Rsend
MPI_Alltoall	MPI_File_iread_at_all_end	MPI_Lookup_name	MPI_Type_create_keyval	MPI_Comm_get_name	MPI_Graph_neighbors	MPI_Rsend_init
MPI_Alltoallv	MPI_File_iread_at_all_end	MPI_Message_c2f	MPI_Type_create_keyval	MPI_Comm_get_parent	MPI_Graph_neighbors_count	MPI_Scan
MPI_Alltoallv	MPI_File_iread_at_all_end	MPI_Message_f2c	MPI_Type_create_keyval	MPI_Comm_group	MPI_Graphdims_get	MPI_Scatter
MPI_Attr_delete	MPI_File_iread_at_all_end	MPI_Mprobe	MPI_Type_create_keyval	MPI_Comm_idup	MPI_Grequest_complete	MPI_Scatterv
MPI_Attr_get	MPI_File_iread_at_all_end	MPI_Mrecv	MPI_Type_create_keyval	MPI_Comm_join	MPI_Grequest_start	MPI_Send
MPI_Attr_put	MPI_File_iread_at_all_end	MPI_Neighbor_allgather	MPI_Type_create_keyval	MPI_Comm_rank	MPI_Group_c2f	MPI_Send_init
			MPI_Type_create_keyval	MPI_Comm_remote_group	MPI_Group_compare	MPI_Sendrecv
			MPI_Type_create_keyval	MPI_Comm_remote_size	MPI_Group_difference	MPI_Sendrecv_replace
			MPI_Type_create_keyval	MPI_Comm_set_attr	MPI_Group_excl	MPI_Sizeof
			MPI_Type_create_keyval	MPI_Comm_set_errhandler	MPI_Group_f2c	MPI_Ssend
			MPI_Type_create_keyval	MPI_Comm_set_info	MPI_Group_free	MPI_Ssend_init
			MPI_Type_create_keyval	MPI_Comm_set_name	MPI_Group_incl	MPI_Start
			MPI_Type_create_keyval	MPI_Comm_size	MPI_Group_intersection	MPI_Startall
			MPI_Type_create_keyval	MPI_Comm_spawn	MPI_Group_range_excl	MPI_Status_c2f
			MPI_Type_create_keyval	MPI_Comm_spawn_multiple	MPI_Group_range_incl	MPI_Status_f2c
			MPI_Type_create_keyval	MPI_Comm_split	MPI_Group_rank	MPI_Status_set_cancelled
			MPI_Type_create_keyval	MPI_Comm_split_type	MPI_Group_size	MPI_Status_set_elements
			MPI_Type_create_keyval	MPI_Comm_test_inter	MPI_Group_translate_ranks	MPI_Status_set_elements_v

... and there's more...much more.

Approaches for encapsulating MPI

Subset wrapping

- Straight encapsulation of the MPI function interface.
 - Approach taken by Axel Kohlmeyer for the mpi package
 - See e.g. <https://core.tcl-lang.org/jenglish/gutter/packages/mpi.html>

```
5 Completeness of the implementation
5.1 Data types
5.2 Communicators
6 TclMPI Command Reference
6.1 tclmpi::init
6.2 tclmpi::finalize
6.3 tclmpi::abort <comm> <errorcode>
6.4 tclmpi::comm_size <comm>
6.5 tclmpi::comm_rank <comm>
6.6 tclmpi::comm_split <comm> <color> <key>
6.7 tclmpi::comm_free <comm>
6.8 tclmpi::barrier <comm>
6.9 tclmpi::bcast <data> <type> <root> <comm>
6.10 tclmpi::scatter <data> <type> <root> <comm>
6.11 tclmpi::allgather <data> <type> <comm>
6.12 tclmpi::gather <data> <type> <root> <comm>
6.13 tclmpi::allreduce <data> <type> <op> <comm>
6.14 tclmpi::reduce <data> <type> <op> <root> <comm>
6.15 tclmpi::send <data> <type> <dest> <tag> <comm>
6.16 tclmpi::isend <data> <type> <dest> <tag> <comm>
6.17 tclmpi::recv <type> <source> <tag> <comm> ?status?
6.18 tclmpi::irecv <type> <source> <tag> <comm>
6.19 tclmpi::probe <source> <tag> <comm> ?status?
6.20 tclmpi::iprobe <source> <tag> <comm> ?status?
6.21 tclmpi::wait <request> ?status?
7 Examples
```

- Still captures the flavor of the MPI API
- Still exposes explicitly the MPI API subset
- Still exposes the complexity.

What do Tcl MPI applications want to do:

- Send scripts executed in other ranks.
 - Special case of send to all or all others.
- Send data that can be handled by other ranks via callbacks.
 - Again special case of send to all or to others.
- Tcl – We know in advance: Everything has a string representation.
- Binary data – may be sent around by the C/C++ part of the application for C/C++ parts of the application to work on needs a way for that code to shove the TclMPI event handling stuff aside and take over.

mpitcl - MPI aware tcl shell.

- Provides MPI aware tclsh.
 - Must be run from mpirun.
- All ranks run this.
- Rank 0 is special – the ‘master’ interpreter it takes input from stdin (normally a file for cluster batch jobs).
- Provides all processes with the mpi namespace in which the mpi command ensemble lives.
- Rank 0 runs a thread to hoist MPI messages received to the event loop (vwait).
- Ranks other than zero run a main loop that accepts MPI messages and act on them under the assumption they come from mpitcl.
- MPI Tags are used to transparently dispatch messages to appropriate handlers.

mpi::mpi command ensemble subcommands:

- size - How many processes are in the application.
- rank – What is my rank in MPI_COMM_WORLD
- execute *where script* – Executes a *script* in the rank(s) defined by *where* where is a rank number “all” or “others”
- send *where data* - sends the *data* to *where*
- handle *script* - specifies *script* to handle data received. The script receives two parameters: sender rank and the data.
- stopnotifier - only legal in rank 0 – stops the event notifier thread.
- startnotifier - only legal in rank 0 – starts the notifier thread again.

Sample mpitcl scripts:

```
mpi::mpi stop notifier  
mpi::mpi send all exit
```

Mimimal MPI
script

```
proc receiver {rank data} {  
    puts "Received from $rank '$data'"  
    incr ::slaves -1  
}  
  
set slaves [mpi::mpi size]  
incr slaves -1;      # number of slave processes.  
  
mpi::mpi handle receiver  
  
mpi::mpi execute others {  
    mpi::mpi send 0 "Rank [mpi::mpi rank] is alive"  
}  
  
while {$slaves} {  
    vwait slaves  
}  
  
mpi::mpi stopnotifier  
mpi::mpi execute all exit
```

Soliciting/getting data from
workers

Applying tclmpi – The Circle is complete.

- NSCLSpecTcl
- Structure
- Parallelization of the interactive version (threaded).
- Creating a batch NSCLSpecTcl and using it with tclmpi.

NSCLSpecTcl is highly interactive

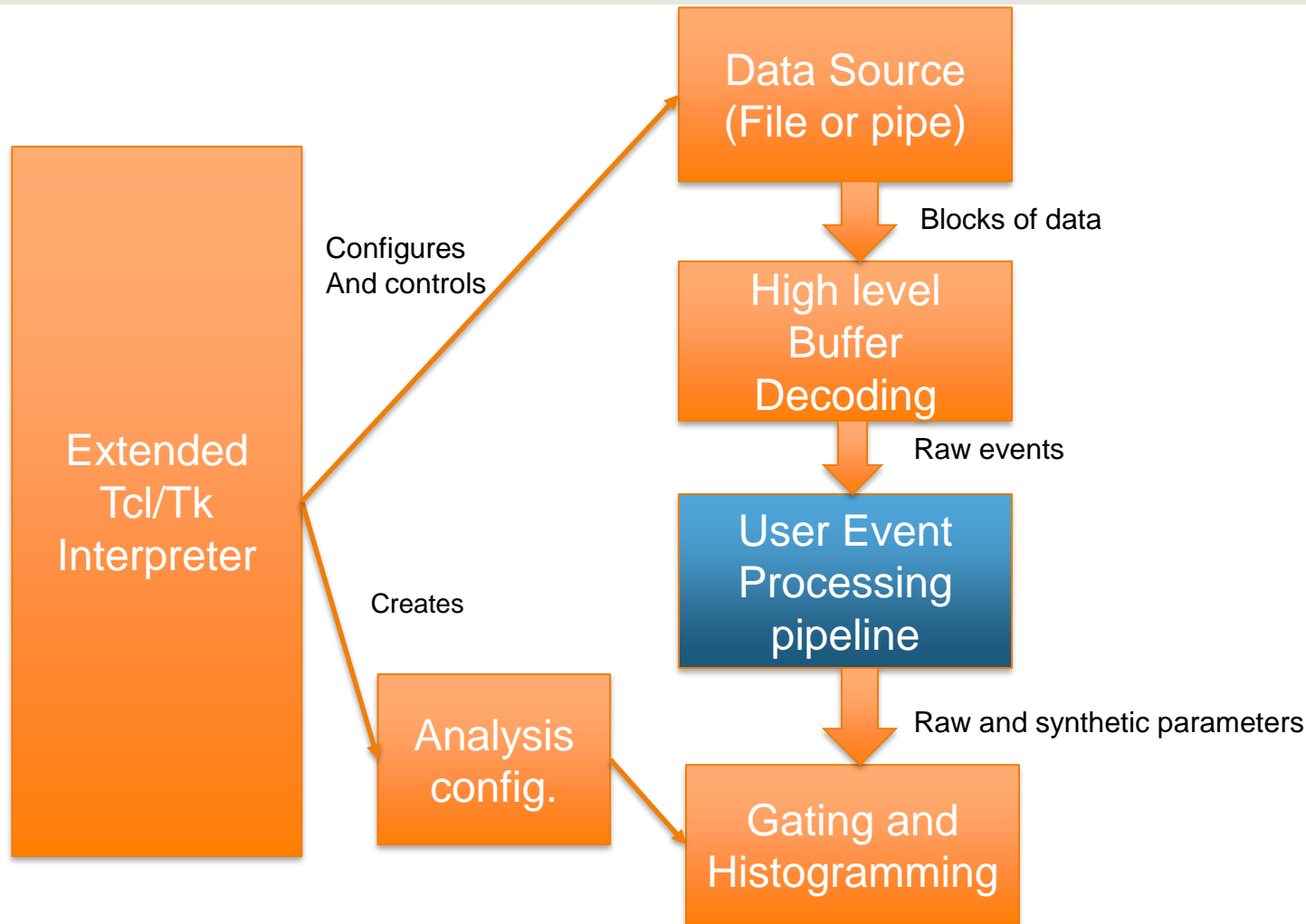
The image displays the NSCLSpecTcl software interface, which is highly interactive. It consists of several windows and panels:

- Left Panel:** A vertical menu with buttons for "Load spectra", "Save spectra", "Clear spectra", "Load configuration", "Save configuration", "Attach online", "Attach to file", "Attach list of files", "Attach filter file", "Detach", "Run title: >>> Unknown <<<", "Data Source: Test Test (Inactive)", "Run number: 0", "Analyzed buffers: 2116", and "Exit".
- TreeGUI Window:** The main configuration window with a menu bar (File, Edit, Data Source, Filters, Spectra, Gate, Help). It has tabs for Spectra, Parameters, Variables, Gates, and Folders. The "Parameters" tab is active, showing "Spectrum Type" (1D, 2D, Summary, Stripchart) and "Data Type" (Word (16 bits), Long (32 bits), Byte (8 bits)). It also includes a "Definition file" section with "Load" and "Save" buttons, and a "Cumulative" checkbox. Below this are buttons for "Create/Replace", "Clear", "Delete", "Gate", "Apply", "Duplicate", and "Ungate". A table lists parameters with columns for Name, Type, X Parameter, Low, High, Bins, Y Parameter, and Gate.
- Table in TreeGUI:**

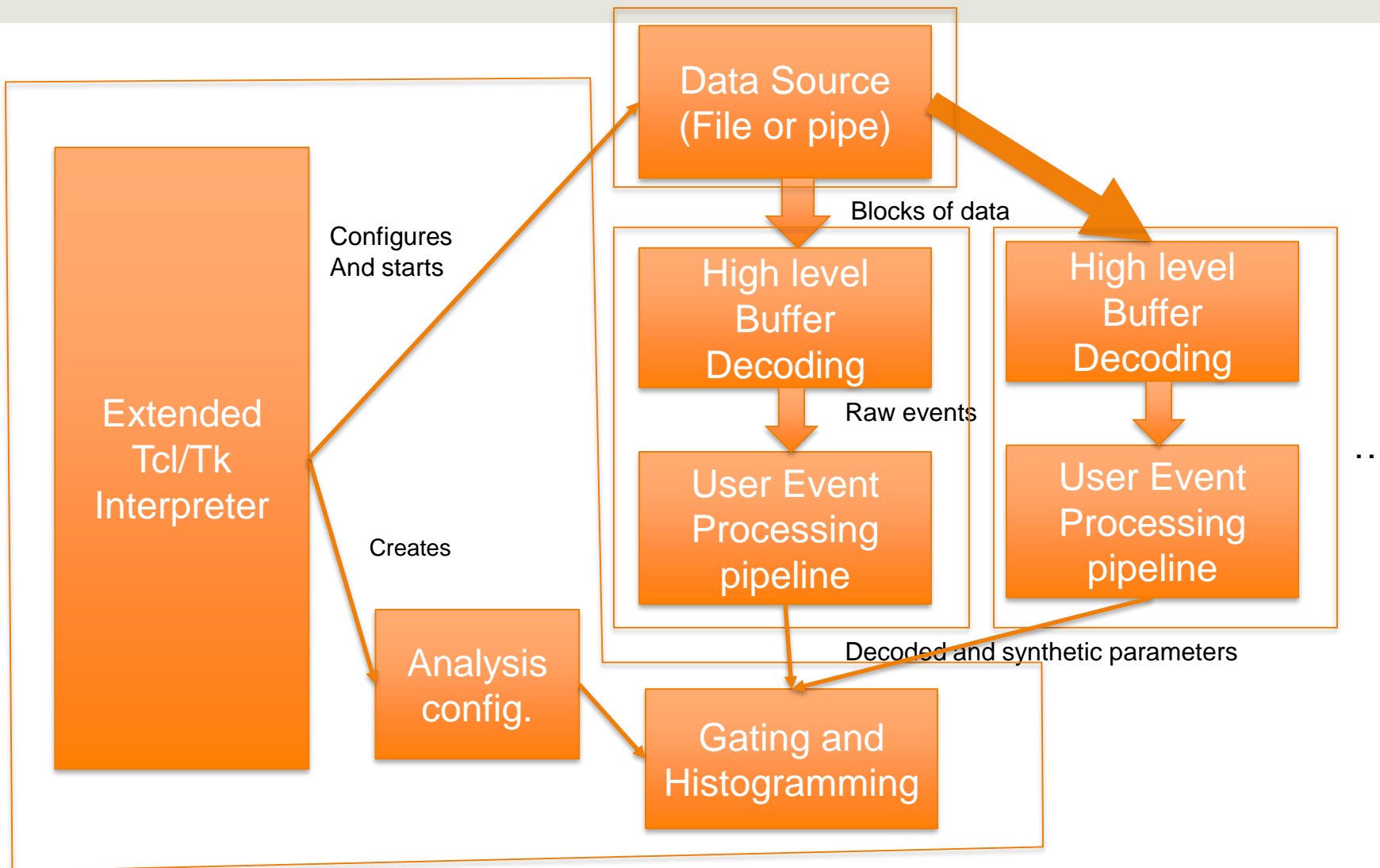
Name	Type	X Parameter	Low	High	Bins	Y Parameter	Low	High	Bins	Gate
2	1	event.raw.00	0	1023	512	event.raw.01	0	1023	512	
raw.00	1	event.raw.00	0	1023	1024					cut
raw.01	1	event.raw.01	0	1023	1024					
raw.02	1	event.raw.02	0	1023	1024					
raw.03	1	event.raw.03	0	1023	1024					
raw.04	1	event.raw.04	0	1023	1024					
raw.05	1	event.raw.05	0	1023	1024					
raw.06	1	event.raw.06	0	1023	1024					
raw.07	1	event.raw.07	0	1023	1024					
raw.08	1	event.raw.08	0	1023	1024					
raw.09	1	event.raw.09	0	1023	1024					
- Xamine Window:** A multi-panel plot window showing a 2D heatmap on the left and three 1D spectra on the right. The spectra are labeled "channel 1", "channel 2", and "channel 3". The x-axis is labeled "X -1.00" and the y-axis is labeled "Y 736.28". The z-axis is labeled "Counts 0". Below the plots are control buttons for "Geometry", "Zoom", "Update All", "Expand", "Marker", "Cut", "Display", "Update Selected", "UnExpand", "Summing Region", "Band", "Display +", "Info", "Log", "Map", "Integrate", and "Contour".
- SpecTcl Console Window:** A window titled "SpecTcl" with a menu bar (File, Console, Edit, Interp, Prefs, History, Help). It displays the following text:

```
loading history file ... 48 events added
buffer line limit: 512 max line length: unlimited
Main console display active (Tel8.6.2 / Tr8.6.2)
Done.
Starting treeparangui... Done
Version: SpecTcl-5.2-001 build on charlie Wed Aug 14 16:47:37 EDT 2019 by fox
(5.2) 49 % start
(5.2) 50 % stop
(5.2) 51 %
```

Simplified NSCLSpecTcl structure and parallelization approaches.

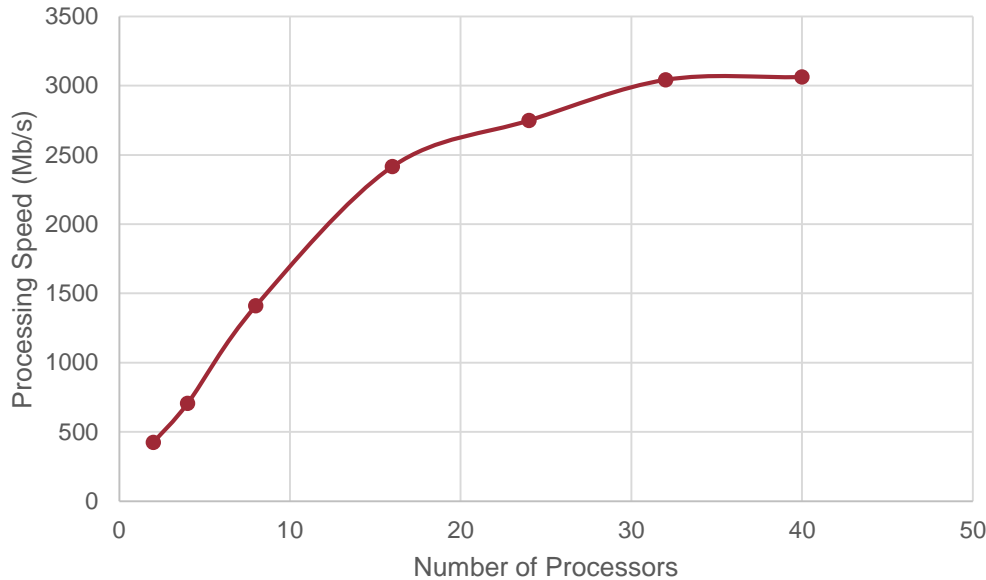


Threaded Parallelism (Giordano Cerizza)



User code must be thread-safe

Threaded Spectcl performance



- Roll off is at the performance limit of the SSD that contained the data.
- Analysis pipeline for this case is simple compared with E17011's.
- Good scaling up until SSD transfer limits.

MPI parallelism?

- We can recruit more cores if the application scales.
- We don't have to worry about thread safety since it's process parallelism.
- With an assist from container technology (e.g. singularity) we can get outside the NSCL to supercomputer centers (or ICER e.g.) with even more cores.
- BUT In almost all cases cluster computing doesn't allow dynamic interactivity.
 - Needed to turn NSCLSpecTcl into a batch program.
 - Needed to figure out how to easily parallelize it.
 - This is the original target of mpitcl.

MPI NSCLSpecTcl

- Each process is a complete batch NSCLSpecTcl
 - Batch NSCLSpecTcl is three packages
 - » spectcl – the base application code.
 - » mpispectcl MPI data sources and sinks.
 - » A user supplied package implementing the processing pipeline.
 - Batch/MPI NSCLSpecTcl has generalized data sources and sinks. **analyze** command sends blocks of events from source to sink. I've implemented:
 - » Source –file.
 - » Source – MPI (for workers – requests block of data from rank 0).
 - » Sink -- Analysis
 - » Sink – MPI (Distributes blocks of data to workers using MPI source).
- Rank 0 :
 - Tells each process (including itself) to read in the configuration scripts.
 - Tells other process to use an MPI Source and Analysis sink
 - Tells itself to use a file data source and MPI sink.
 - Tells everyone to start analyzing data.
- When analysis is complete Rank 0
 - Tells all other processes to send it spectrum data.
 - Sums the spectra into total spectra
 - Writes them out for visualization.

What this looks like:

```
mpi::mpi execute all {  
    package require spectcl  
    package require mpispectcl  
    package require MyPipeline;    # User event processing code is here.  
    source defs.tcl  
}
```

```
mpi::mpi execute others {  
    mpisource  
    analysissink  
}
```

```
filesource run-0003-00.evt  
mpisink  
mpi::mpi stopnotifier
```

```
mpi::mpi execute others analyze  
analyze
```

```
mpi::mpi startnotifier
```

Getting the data back:

```
set l [spectrum -list]
set f [open spectra.dat w]

foreach spectrum $l {
    set name [lindex $spectrum 1]
    getSpectrumFromWorkers $name
    swrite -format ascii $f $name;    # Writes a histo to file.
}

close $f
```


Summing a spectrum from worker nodes.

```
proc addData {name src data} {;           # sums the histo from 1 worker into the local histo.
    foreach datum $data {
        set value [lindex $datum end]
        set coords [lrange $datum 0 end-1]
        set current [channel -get $name $coords]

        incr current $value
        channel -set $name $coords $current
    }
    incr ::expected -1
}

proc getSpectrumFromWorkers name {
    mpi::mpi handle [list addData $name]
    set ::expected [mpi::mpi size]
    incr ::expected -1;

    set script "mpi::mpi send 0;      # Care must be taken to ensure substitutions are done
    append script "\["                ; # in the right process.
    append script "scontents $name]"
    mpi::mpi execute others $script
    while {$::expected > 0} {
        vwait ::expected;
    }
}
```

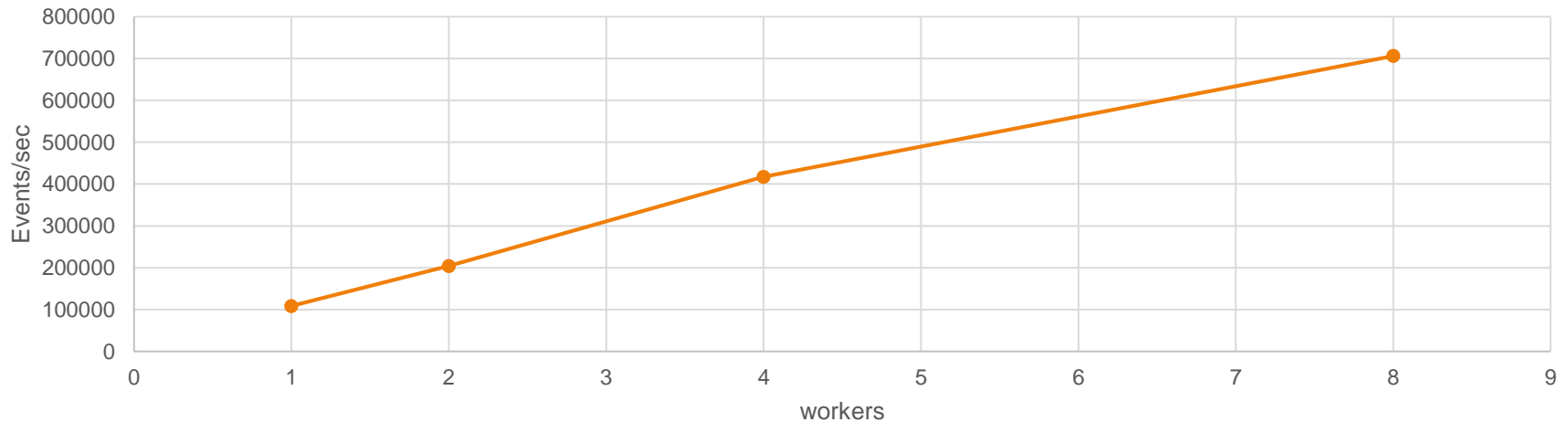
Does this work?

10Gbyte file with 2,754,450 events.

Workers	Time	MB/sec	Events/sec		
1	25.4	403.1496	108442.91		
2	13.5	758.5185	204033.33		
4	6.6	1551.515	417340.91		
8	3.9	2625.641	706269.23		

- 2.6Gbytes/sec is interconnect saturation (dual 10Gbit/sec ethernet)
- Event processing in the actual experiment will be more complex but we can scale to more workers.

Events/sec



Conclusions

- NSCL's Transition to modern, digital nuclear electronics poses problems for online and near-line data analysis.
- Experiments will increasingly require parallelism in online and near-line data handling:
 - Software tools to make it easy for naïve users to make use of parallelism by plugging in their event analysis code.
 - Large core count systems for interactive, online analysis (threaded parallel).
 - Clusters dedicated to the running experiment for near-line analysis (distributed parallel).
 - High speed interconnects to support the data flow bandwidth.
 - Large, multi-petabyte storage that's fast with fast interconnects.
- E17011 provides a laboratory to explore the techniques we'll need to apply to modern experiments.
- mpitcl is one technique to simplify parallel programming for “the masses”
 - Easily retrofitted a highly interactive, complex serial analysis program (2 days work)
 - Got scaling up to the interconnect bandwidth.