# LArSoft roadmap

Erica Snider
March 9, 2023
CSAID Roadmap Meeting

# LArSoft

- The Collaboration
  - Experiments, laboratories, software projects collaborating to produce shared, detector-independent software for LArTPC simulation, reconstruction and analysis

- LArSoft "Project"
  - Fermilab team (SciSoft) who support the software sharing paradigm, own the common infrastructure / architecture, provide user support / software expertise to experiments

- The experiments
  - Develop, contribute, validate and support the algorithm code

- Governance
  - Monthly meeting with offline leadership from all experiments
  - Quarterly meetings with spokes-level experiment representatives ("Steering Group")
    - Charged with oversight of project work
    - Approve an annual work plan. See https://LArSoft.org

# LArSoft work plan

- Describes the high-level plan of work for the project team.

  – Developed through process of one-on-one meetings with experiments followed by iterations on the draft until presentation to / approval by the Steering Group

  – Reflects experiment requirements and requests
  – Implements the strategic directions for the shared code of the collaboration

  LArSoft / SciSoft play a strong leadership role in defining direction, strategy

- [The 2023 LArSoft Work Plan](#)

**Fermilab**

# 2023 LArSoft Work Plan

Strategic directions of the 2023 work plan

1.  Support multi-threading to optimize running on grid resources

2.  Enable / facilitate optimized running on GPU and HPC resources

3.  Facilitate / simplify integration of machine learning workflows

4.  Support heterogeneous detector readouts in simulation and reconstruction

5.  Provide a multi-experiment capable event display framework

6.  Expand adoption of community / industry supported tools

# 2023 LArSoft Work Plan

Strategic directions of the 2023 work plan

1.  Support multi-threading to optimize running on grid resources

2.  Enable / facilitate optimized running on GPU and HPC resources

3.  Facilitate / simplify integration of machine learning paradigms

4.  Support heterogeneous detector readouts in simulation and reconstruction

5.  Provide a multi-experiment capable event display framework

6.  Expand adoption of community / industry supported tools

LArSoft held "LArTPC Multi-threading and Acceleration Workshop" Mar 2–3
(will come back to this at end…)

**≢ Fermilab**

# 1. Support multi-threading

- All experiments report using multiple grid slots to accommodate memory of jobs

    – Running single-threaded programs leads to significant underutilization of CPU
    – Memory use driven by event-level data
    – Need sub-event level multi-threading or more granular data management strategies


- Project work
    – Working with experiments to ensure thread-safety in common and experiment-specific code
    – Implementing multi-threading in common services
    – Past integration of contributions from SciDAC4 efforts

**Fermilab**

# 2. Enable / facilitate optimized running on GPU and HPC

- Many LArTPC computing problems highly parallelizable that can benefit from hardware acceleration
  - Low-level data and signal processing
  - Simulation
  - Machine learning

- Multi-threading, GPU acceleration create paths to optimized running on HPC
  - Several experiments / projects have experience with LArSoft code on HPC
  - Demonstrated in SciDAC4 work by Giuseppe Cerati, Sophie Berkman, et al.

- Project work
  - Focus on making low-level data structures suited to GPU processing,
  - Work with experiments on specific algorithms (to be identified)
  - Current Spack migration well suited to needs of HPC-enabled builds

🔷 **Fermilab**

# 3. Facilitate / simplify integration of machine learning workflows

- Most ML efforts within experiments are completely external to LArSoft
  - MicroBooNE experience:
    - ML-based analysis branch all but isolated to small group of analyzers.
    - Separate data production workflows required, which slowed data availability
  - Integration into LArSoft would alleviate all these issues

- Some ML algorithms benefit from GPU acceleration at inference stage
  - Highly dependent on the problem and solution
  - Some overlap with acceleration work previously noted

- Experiment groups in ICARUS and ND-LAr working on fully ML workflows

- Project work
  - Ensure configurations, inputs and outputs are available to ML interfaces
  - Assist experiment groups with interfacing to LArSoft
  - Past integration of Sonic-derived GPUaaS into LArSoft targeted ML inferencing

# 4. Support heterogeneous detector readouts in sim and reco

- Primarily aimed at accommodating pixelated readouts (ND-LAr)
  - Also intended to allow future detectors to have completely different readout schemes

- Project work

  - Adapt geometry and simulation systems
    - Portions of reconstruction code must differ
      - Will be provided by experiments

  - Geometry: requires re-factoring readout from volume geometry
    - Several wire-plane readout configurations already supported
    - Readout geometry currently tightly intertwined with more generic volume geometry

  - Past work adapted simulation via similar abstraction of anode simulation
    - The "artg4tk / LArG4" re-factoring completed several years ago

Fermilab

# 5. Provide a multi-experiment capable event display framework

- A persistent and vocal ask from many experiments

- Would add value in exactly the same way that common sim/reco do.


- Project work

  – Design, develop event display framework, or adapt an existing ED to requirements

  – Experiments provide customizing code

  Requires local ED / visualization expertise, which is currently lacking

  – Can view this as a request to build this expertise

# 6. Expand adoption of community / industry supported tools

- A good strategy wherever possible and cost effective

- Recent major examples
  – Migration to GitHub (last year)
  – Migration to Spack (continuing)

- Project work

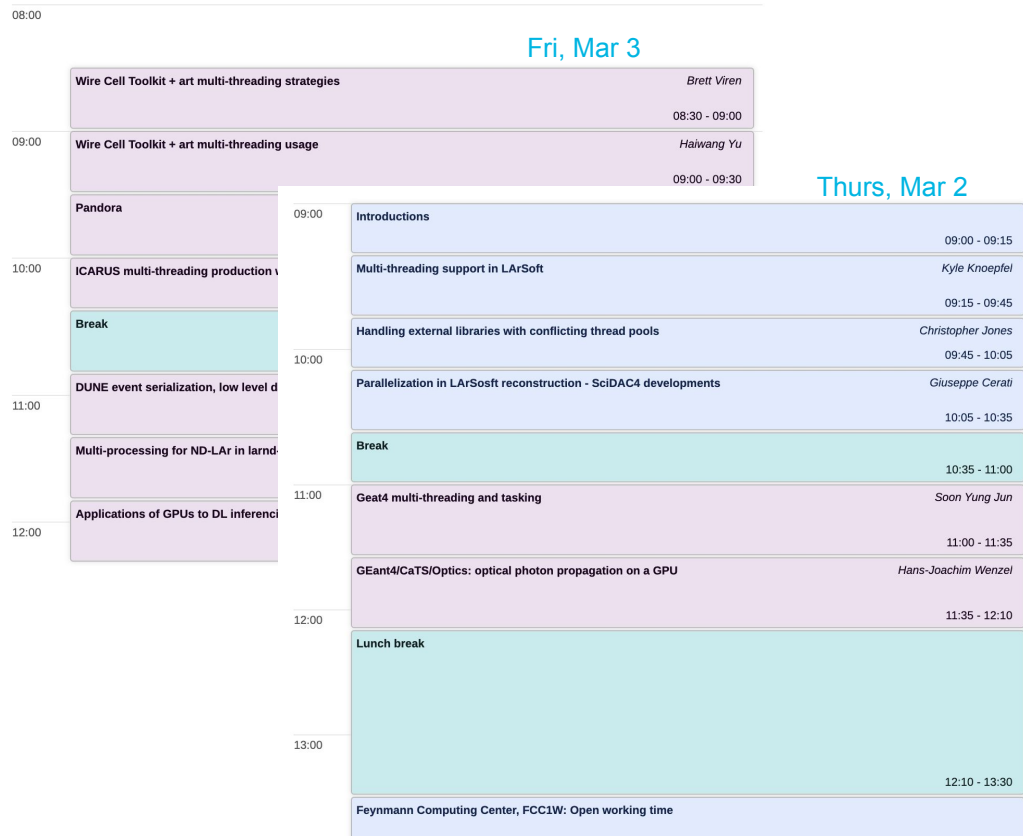  – Nothing beyond existing work currently in plan, but always seeking opportunities

# LArTPC multi-threading and acceleration workshop

Workshop goals:

- Invited developers and representatives from *art*, GEANT4, LArTPC SciDAC4, LArTPC neutrino experiments

- To learn the multi-threading and acceleration capabilities of frameworks and common toolkits used by LArTPC experiments;

- To share experiences across experiments about existing resource utilization and throughput problems that lend themselves to multi-threaded or acceleration solutions;

- To explore how multi-threading and acceleration is being used to address these problems and open avenues to the use HPC resources more broadly;

- To discuss the results of applying these techniques and capabilities
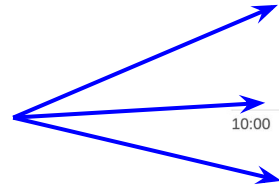
**Fermilab**

# Organization of the program

- Held last week, Mar 2–3

  - One day of presentations spread across two morning sessions

- Introduction by Adam Lyon

  - Framed broad context and directions within HEP

- Well attended:

  - 43 registered
  - 27 online + 10 in the room on Thursday at peak
  - At least 7-10 in the room + 20-24 online DC

- Engaged audience and robust discussion

  - Could see several adapting plans based on what they were learning

  - Some stated this explicitly

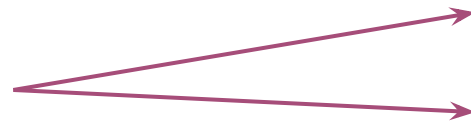- Consider it a success from this perspective

# Organization of the program
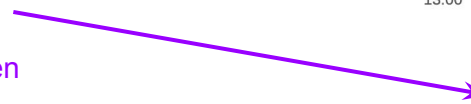
Common tools and support

Simulation tools

Open working time:  13:30 – 17:00
(per request)

Zoom will remain open
during this time

| 09:00 | **Introductions** | |
| | | 09:00 - 09:15 |

**Multi-threading support in LArSoft** — *Kyle Knoepfel*
09:15 - 09:45

**Handling external libraries with conflicting thread pools** — *Christopher Jones*
09:45 - 10:05

**Parallelization in LArSosft reconstruction - SciDAC4 developments** — *Giuseppe Cerati*
10:05 - 10:35

**Break**
10:35 - 11:00

**Geat4 multi-threading and tasking** — *Soon Yung Jun*
11:00 - 11:35

**GEant4/CaTS/Optics: optical photon propagation on a GPU** — *Hans-Joachim Wenzel*
11:35 - 12:10

**Lunch break**
12:10 - 13:30

**Feynmann Computing Center, FCC1W: Open working time**

# Organization of the program

Experiment tools and experience I

Experiment tools and experience II

| | |
|---|---|
| **Wire Cell Toolkit + art multi-threading strategies** | *Brett Viren* |
| | 08:30 - 09:00 |
| **Wire Cell Toolkit + art multi-threading usage** | *Haiwang Yu* |
| | 09:00 - 09:30 |
| **Pandora** | *Ryan Cross et al.* |
| | 09:30 - 10:00 |
| **ICARUS multi-threading production workflow** | *Tracy Usher et al.* |
| | 10:00 - 10:20 |
| **Break** | |
| | 10:20 - 10:50 |
| **DUNE event serialization, low level data processing and production** | *Thomas Junk* |
| | 10:50 - 11:20 |
| **Multi-processing for ND-LAr in larnd-sim and ndlar_flow** | *Matt Kramer* |
| | 11:20 - 11:50 |
| **Applications of GPUs to DL inferencing** | *Michael H L Wang* |
| | 11:50 - 12:20 |

08:00
09:00
10:00
11:00
12:00

15

# LArTPC multi-threading and acceleration workshop

- Some observations

  - DUNE + SBN experiments all working on multi-threading, HPC, GPUs

  - DUNE ND-LAr work outside LArSoft was particularly interesting
    - 3rd generation simulation / reconstruction
      - After Pandora/LArSoft, Wire-Cell
      - Heavy use of ML in workflows
    - Written in python
    - Vector data structures / operations replace loops
    - Framework provides native GPU support, data provenance

    Not clear how / if this will connect to LArSoft in future
    Not stated, but seems management and coders may not agree on strategies

  - DUNE working on alternatives to multi-threading to manage memory issues
    - Operate at APA level (defines TPC-level readout unit) for horizontal drift FD
    - Allows, but does not require multi-threading

LAr Soft

Fermilab

# LArTPC multi-threading and acceleration workshop

- Some observations

  - Wire Cell Toolkit
    - Full-featured sim/reco/R&D framework with configuration, plugins, component factories, interface classes, code aggregation methods, lib/package building
    - Notably not provenance, which they get from interfacing with LArSoft
    - Event loop / file IO can operate at APA-level for DUNE FD
      - Multiple APAs in flight simultaneously without reading entire event
      - Some limitations to multi-threading capabilities when run in art / LArSoft depending on structure of input files
    - Production simulation may not produce low level data objects
      - That level of data would then be exclusive to WCT

  - Discussions related to future of LArSoft
    - Questions about the function of common frameworks if always passing files between what are effectively independent, stand-alone applications.
    - What should LArSoft be doing to ensure we are adding the most value possible?

Fermilab

# SciSoft team

- Vito di Benedetto
- Patrick Gartung
- Chris Green
- Robert Hatcher
- Kyle Knoepfel (co-lead)
- Lynn Garren (ret.)
- Marc Paterno
- Saba Sehrish
- Erica Snider (co-lead)
- Mike Wang
- Hans Wenzel

# The end

# Backup

# Why a multi-threading and acceleration workshop?

1.  Resource optimization and throughput bottlenecks on existing resources
    *   All LArTPC neutrino experiments at the lab report significant fraction of jobs running on more than a single grid slot due to memory consumption
    *   Many LArTPC computing problems are parallelizable and would benefit from various types of acceleration

2.  HPC
    *   Funding agencies pushing lab / experiments to use more HPC
    *   Many experiments / groups have experience with this already
    *   Multi-threading / optimizing for GPU also help with this transition, or are already part of it

3.  Uniformity of LArTPC technology
    *   LArTPCs are well-suited for direct sharing of code, techniques, technologies

Erica Snider          LArSoft Roadmap

**Fermilab**