# SBN production techniques and outlook

Giuseppe Cerati & Steven Gardiner

IF Production Data Processing Mini-Workshop

17 January 2024

# Outline
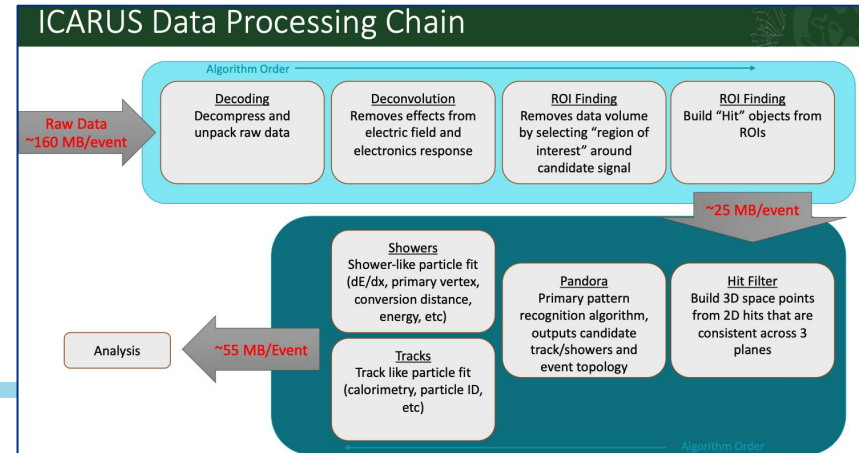
- Production overview
  - what data we get, what MC
- Keep up processing
- Campaign definitions
  - List of fcls, what output to store, where
- Disk management
  - Use FTS to transfer file to/from tape
- Data transfers with Rucio
  - CNAF and SLAC experience
- Metadata and its usage
  - What do we store and what is essential to keep in the future
- Monitoring capabilities

🟦 **Fermilab**

# Overview

- The SBN program includes the ICARUS and SBND experiments
  - Joint management of computing resources
- Avoid in-house tools, use as much as possible what provided by CSAID
  - POMS, fife_utils, etc.
- ICARUS and SBND collect data from different streams:
  - On-beam triggered data: BNB (both), NuMI off-axis (ICARUS only)
  - Off-beam triggered data: model of in-time cosmics background
  - Off-beam min bias: model of non-triggering cosmic background
- Data processing includes several steps:
  - typically grouped in 2 stages
  - MC processing adds GEN, G4, detsim
  - "Overlay" workflow under development
  - Alternative workflows also used (e.g. ML)

I. Caro Terrazas

## ICARUS Data Processing Chain

Algorithm Order →

Raw Data
~160 MB/event

**Decoding**
Decompress and unpack raw data

**Deconvolution**
Removes effects from electric field and electronics response

**ROI Finding**
Removes data volume by selecting "region of interest" around candidate signal

**ROI Finding**
Build "Hit" objects from ROIs

~25 MB/event

**Showers**
Shower-like particle fit (dE/dx, primary vertex, conversion distance, energy, etc)

**Tracks**
Track like particle fit (calorimetry, particle ID, etc)

**Pandora**
Primary pattern recognition algorithm, outputs candidate track/showers and event topology

**Hit Filter**
Build 3D space points from 2D hits that are consistent across 3 planes

Analysis

~55 MB/Event

Algorithm Order

# Keepup processing

- Not yet operational for SBND, design expected to be similar to ICARUS

- ICARUS uses POMS for handling its keepup processing

- Keepup processing is broken up into two main stages:
  - stage0 and stage1+caf+larcv

- Each data stream is handled by an individual workflow (8 total)

- Keepup processing automated via POMS automatic job submission

**🟦 Fermilab**

# Campaigns and definitions

- The campaign parameters are defined in a .cfg file and processed with POMS
  - example /exp/icarus/app/poms_test/cfg/icarus_run2_keepup_production.cfg
  - parameters can also be modified via POMS, such as the number of jobs, how many events, and destination of output using POMS too
  - use POMS slicing feature to handle large datasets; can cause delays if desired % not reached and next slice not triggered. More dynamic handling would be valuable.
- Handling of duplicate files:
  - for "first-gen" files (e.g. output from processing raw data through stage0) using sam_dataset_duplicate_kids is good enough to get rid of these duplicates but will not work for "children" files of duplicates. This tool is available to in fife_utils.
  - for the "second-gen" (say the stage1 output files) we created a script that looks at the parent file and if the file is retired from samweb, then deletes the child file. This script is icarus-specific and can be improved upon
- Handling of failures:
  - use POMS recovery feature based on SAM process "*consumed_status*"
  - Jobsub AutoRelese feature to deal with jobs exceeding resource usage

# Monte Carlo Production Workflow (1)

- SBN MC Production organized in "Production pushes"
  - Different MC samples produced under the same "push" for consistent studies across the collaboration(s)
  - Statistics and physics defined by physics groups and the same across samples of the same "push"

- MC Sample Request
  - Physics Groups (and/or individual analysers) fill a Google form with their MC request
  - Basic information needed: group contact, statistics, code version, list of fcl files, files to be saved
  - Only workflows tested locally are run by production. We also request for expected usage of resources
  - Form populates a spreadsheet that is used by the production group to organize and document each sample configuration

# Monte Carlo Production Workflow (2)

- POMS campaign setup
    - Production group creates .cfg configuration files with basic information for each "push"
        - Defines Samweb dataset and file names, code version, setup workflow stages, calls metadata injector, and define files to be saved
        - Each .cfg launch scrypt defines a new job_type in POMS
    - SBND creates a dataset of fcl files, defining input files and basic metadata like run/subrun. ICARUS uses POMS to define input and basic metadata.
    - Base campaign created for each workflow.
        - Minimal information need to be changed per each MC sample: list of .fcl files, campaign name, and number of slices (in case of large statistics)
    - Small test run for every campaign, requester is contacted to validate physics output while production evaluates resources usage, properly created datasets, final files location, etc

- Running the full sample
    - Each sample can be split into slices (priorities defined by physics groups) w/ all submissions handled by POMS
    - Once done the sample is checked offline for duplicate files, that are deleted and retired from SAM using a set of scripts based on metadata information.
    - Final sample is documented and advertised to the SBN collaboration
        - MC sample Request form used for sample documentation and bookkeeping

🔶 **Fermilab**

# Disk management

- FTS used to retrieve raw data files from DAQ, assign SAM metadata

  – Modest adjustments expected to adapt to metacat instead

- Data pool disk space (/pnfs/sbn/data and /pnfs/sbn/data_add) held in common between SBND and ICARUS

  – No technical enforcement of per-experiment quotas at present

- Tape read-back expensive, we avoid it where we can in workflow design

  – May become an issue for "overlay" workflow in the future

# Data transfers with Rucio

- ICARUS has been using Rucio for data transfers to offsite locations
- Transfer of RAW data to CNAF (duplicate archival copy)
  - triggered asynchronously by CNAF experts on a weekly basis
  - files at CNAF also declared on SAM and are available to users through grid processing
  - xrootd not supported at the moment
- Transfer of files for ML workflow at SLAC
  - automatically triggered during production campaigns (both data and MC)
  - still work in progress, issues with non-deterministic vs deterministic remote storage element

# Metadata and its usage

- ICARUS and SBND share metadata format and overall technical handling of it
  - metadata injector called in the launch script defined in the POMS job_type
  - Current version include critical fields e.g.: run number, release version, stage and file format, fcl file used

- Metadata variables planning still in the early stages for SBND
  - Final design can be adjusted in light of infrastructure changes

🎇 Fermilab

# Monitoring

- POMS and fifemon provide useful information on status of production

- As tools continue to evolve, maintaining similar monitoring capabilities will be important

🔹 **Fermilab**

# Summary

- SBN has been using as much as possible production tools provided by CSAID
- We welcome migration to more modern tools, and hope migration will be straightforward
- Production is critical at this time for our experiments, so it will be important to share the migration plan with enough contingency and give the opportunity for careful tests before retiring the current tools

‡ Fermilab