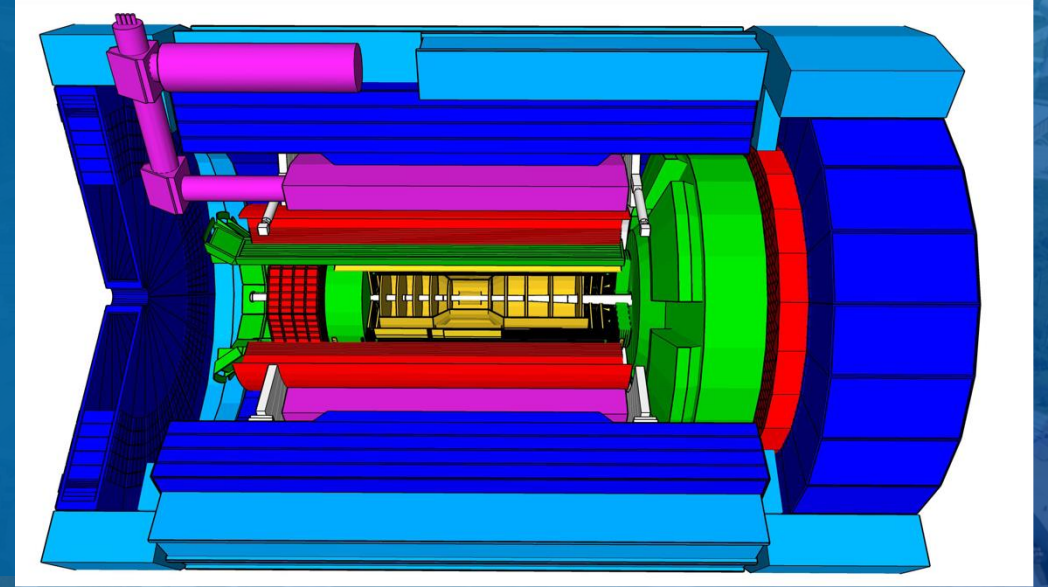


JUN 16, 2026

TOWARDS AI-BASED RECONSTRUCTION -- AI-READY DATA



Zhiwan Xu
Argonne National Lab

NP AMSC DATA PROVIDERS PROGRAM (DAPP)

Design principle: [10.5281/zenodo.19710980](https://zenodo.org/record/19710980)
BIC application + metadata: [arxiv:2606.07667](https://arxiv.org/abs/2606.07667)

AI-READY DATA DESIGN

Design principle

- Physical measurement (x, y, z, A, t) is naturally a "point-cloud" data structure: sparse, heterogeneous, uniquely distinctive from industrial LLM training data.
 - 4 Design principles: **unified representation, explicit uncertainty, metadata-defined semantics, minimal scope**
 - Available in Zenodo: [10.5281/zenodo.19710980](https://zenodo.org/record/19710980).
- Why not ROOT/HDF5? → **Nested tree structures require iterative lookup; no clear separation between trainable features and truth labels**
- Why not existing reconstruction formats? → **Designed for physicist-driven analysis, not AI training; tightly coupled to detector-specific software, hard to generalize across experiments.**

AI-READY DATA DESIGN

Measurement and Label array for supervised learning

TABLE I. Structure of the `measurements` array.

Field	Type	Description
<code>event</code>	<code>uint64</code>	Event identifier.
<code>detector</code>	<code>uint16</code>	Subdetector identifier.
<code>hit</code>	<code>uint16</code>	Hit index within a subdetector.
<code>sample</code>	<code>uint16</code>	Sample index within a hit.
<code>x, y, z</code>	<code>float32×3</code>	Spatial coordinates associated with the sample.
<code>pos_cov_xx, xy, xz, yy, yz, zz</code>	<code>float32×6</code>	Upper triangle of the position covariance matrix.
<code>amplitude</code>	<code>float32</code>	Measured signal amplitude.
<code>amplitude_uncertainty</code>	<code>float32</code>	Uncertainty on the amplitude.
<code>time</code>	<code>float32</code>	Time associated with the measurement.
<code>time_uncertainty</code>	<code>float32</code>	Uncertainty on the time measurement.

TABLE II. Structure of the `labels` array, with contributor-specific fields repeated for each retained contributor $n \in \{0, \dots, N - 1\}$.

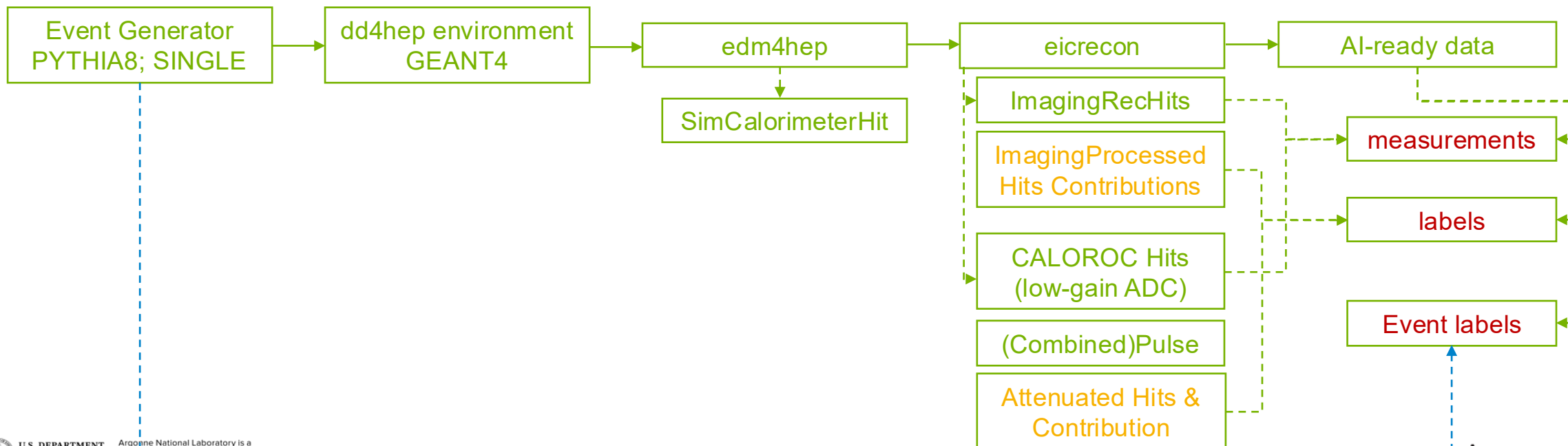
Field	Type	Description
<code>event</code>	<code>uint64</code>	Event identifier.
<code>detector</code>	<code>uint16</code>	Subdetector identifier.
<code>hit</code>	<code>uint16</code>	Hit index within a subdetector.
<code>sample</code>	<code>uint16</code>	Sample index within a hit.
<code>x, y, z</code>	<code>float32×3</code>	Truth spatial coordinates.
<code>deposit_energy</code>	<code>float32</code>	Truth total deposited energy.
<code>time</code>	<code>float32</code>	Truth time associated with the measurement.
<code>n_contributions</code>	<code>uint16</code>	Total number of contributing particles.
<code>con{n}_particle</code>	<code>uint32</code>	Index of contributor n in the event-level record.
<code>con{n}_x, y, z</code>	<code>float32×3</code>	Interaction spatial coordinate of contributor n .
<code>con{n}_deposit_energy</code>	<code>float32</code>	Deposited energy from contributor n .
<code>con{n}_time</code>	<code>float32</code>	Time of interaction of contributor n .
<code>con{n}_px, py, pz</code>	<code>float32×3</code>	Momentum vector of contributor n at production.
<code>con{n}_energy</code>	<code>float32</code>	Energy of contributor n at production.
<code>con{n}_pid</code>	<code>int32</code>	PDG identifier of contributor n .
<code>con{n}_parent_pid</code>	<code>int32</code>	PDG identifier of the immediate parent of contributor n .
<code>con{n}_parent_energy</code>	<code>float32</code>	Energy of the immediate parent of contributor n .

- Suitable for a variety of detector readouts, including tracking, calorimeter, spectator detectors. The uncertainty covariance matrix is key information for AI to learn detector resolution.
- Label array design principles: granularity, contributors, spatiotemporal, ancestry, minimality.

BIC CASE APPLICATION

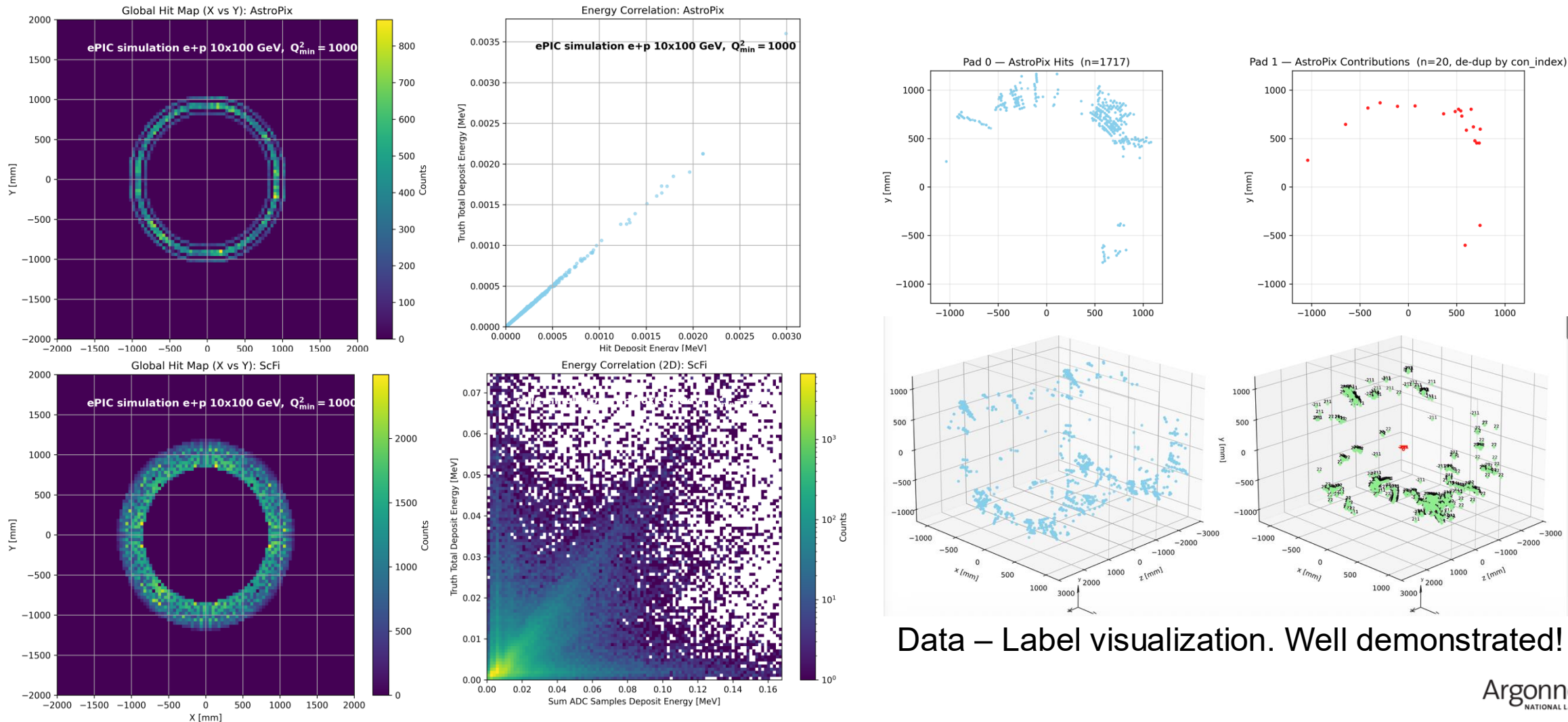
DAPP Argonne led detector data thrust

- Advantage: Hybrid dual readout: **AstroPix silicon pixels** (500 μm , single-sample) for spatial imaging + **Pb/ScFi waveform per SiPM** (25 ns multi-sample) for energy sampling.
- 3 arrays: measurements (readout), labels (per-hit truth + top-3 contributors), event_label (MC particles + decay chain). [arxiv:2606.07667](https://arxiv.org/abs/2606.07667)



AI-READY DATA VISUALIZATION

DIS data 10x100 GeV, min Q2=1000, npz format, 1k events ~ 5M

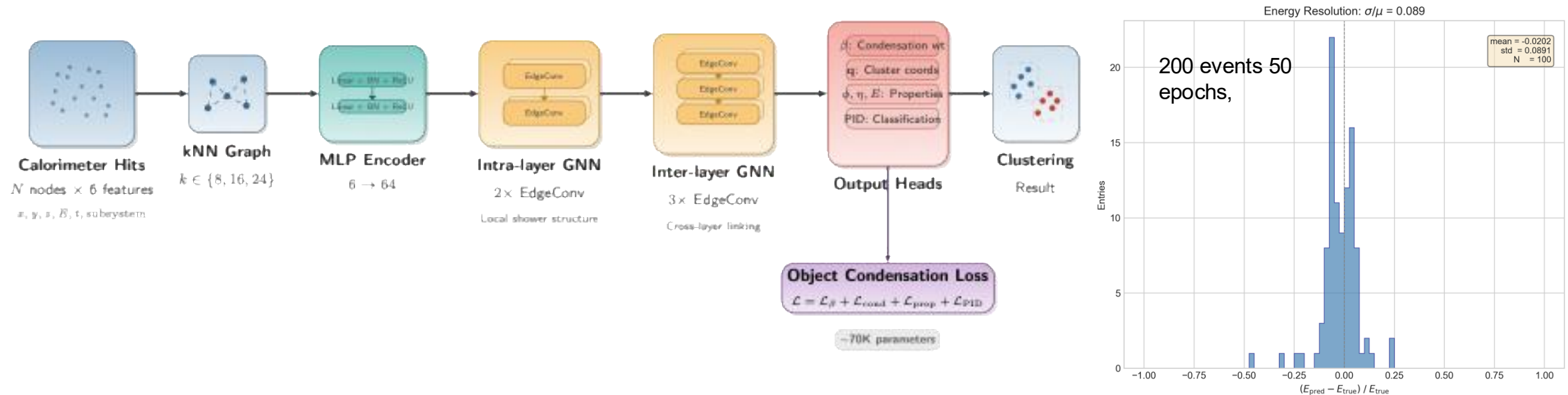


Data – Label visualization. Well demonstrated!

FIRST GLIMPSE OF AI-BASED RECONSTRUCTION

Build AI models based on the AI-ready BIC data

- AI applications: reconstruction (clustering, particle ID, energy regression), anomaly detection, fast simulation, jet tagging, and physics inference from raw detector data.
- (Ongoing) AI demonstration with hierarchical GNN with Object Condensation for shower cluster reconstruction in the BIC. Toy training on 10 GeV γ simulation: $E_{res} \sim 9\%$



Thank you!

Argonne 
NATIONAL LABORATORY



U.S. DEPARTMENT
of ENERGY