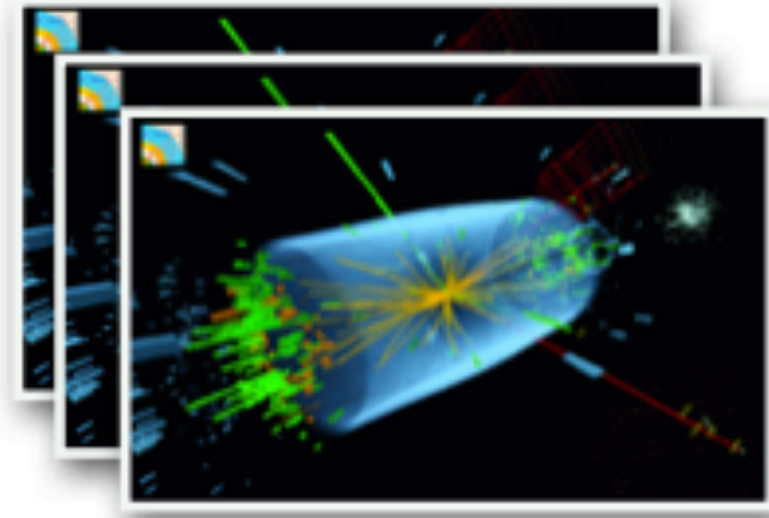




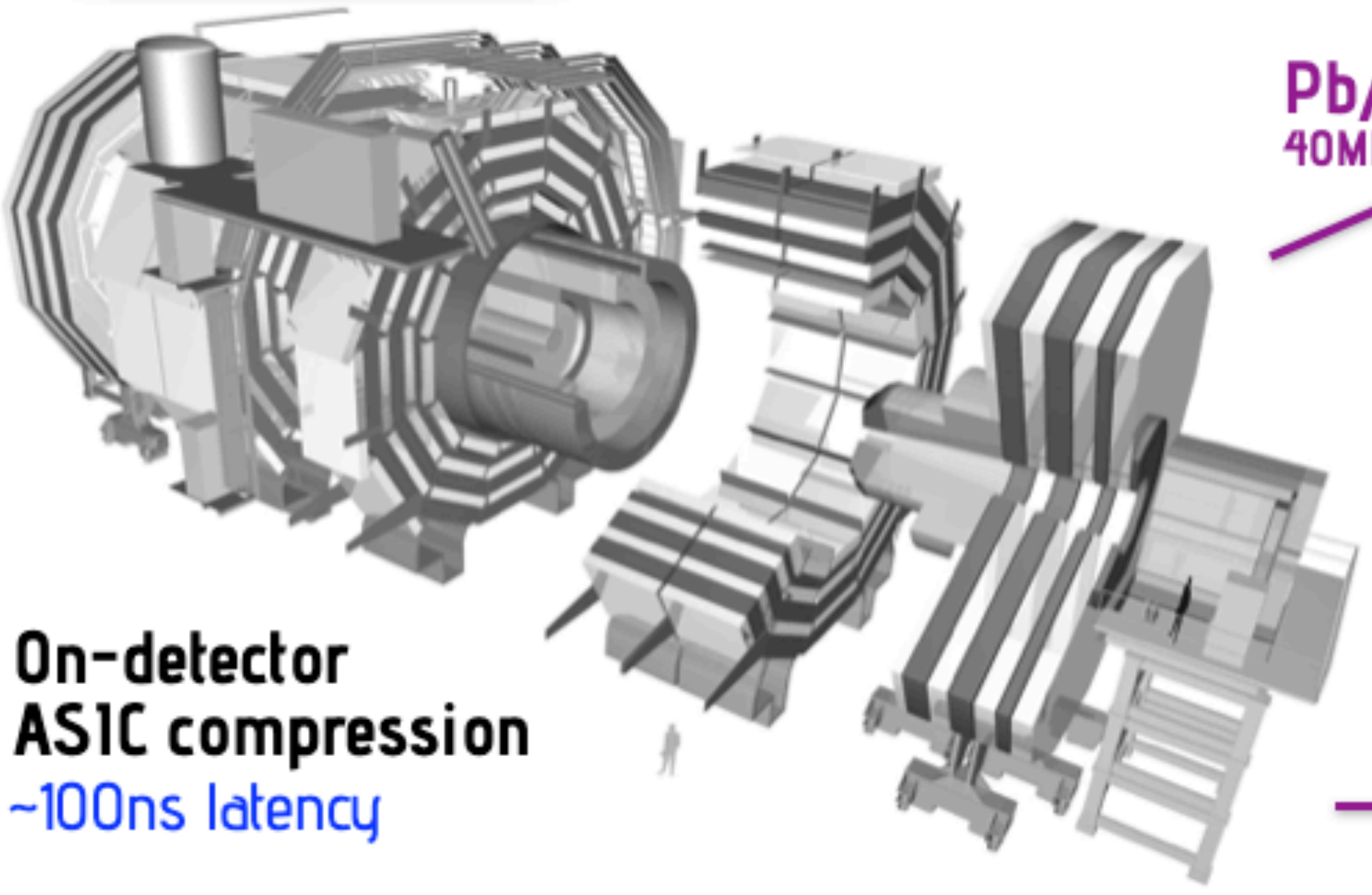
Real time AI

Ben Hawks, Sergo Jindariani, Jovan Mitrevski, Nhan Tran (Fermilab) for
the hls4ml team (fastmachinelearning.org)



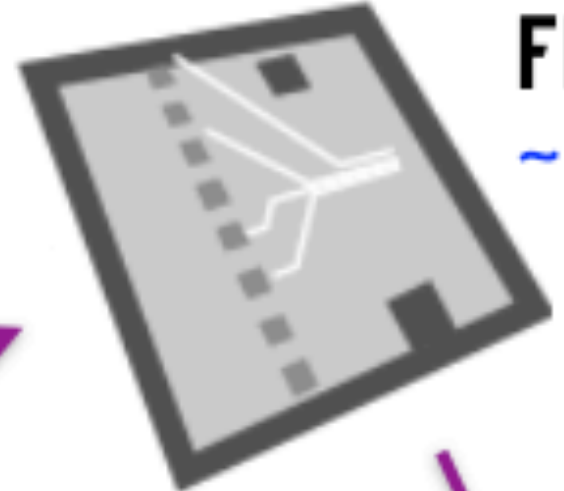
CMS Experiment

40MHz collision rate
~1B detector channels



On-detector ASIC compression
~100ns latency

Pb/s
40MHz

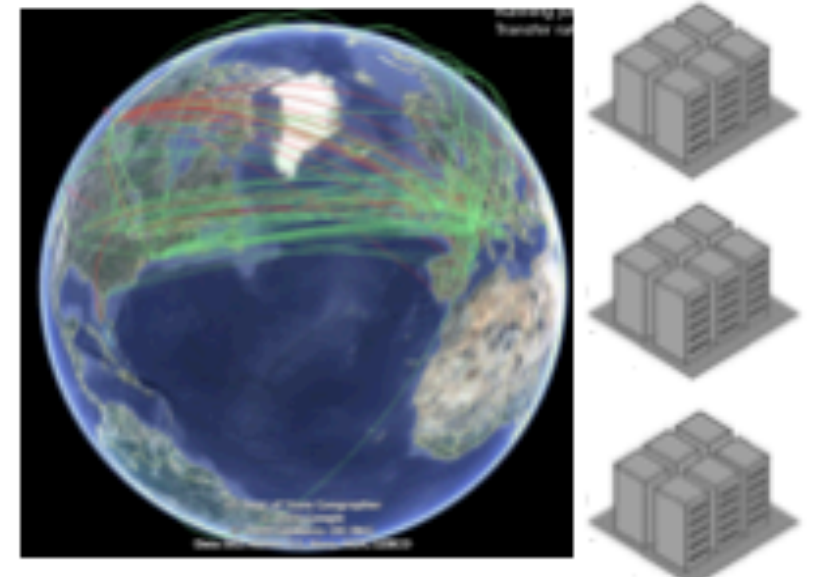
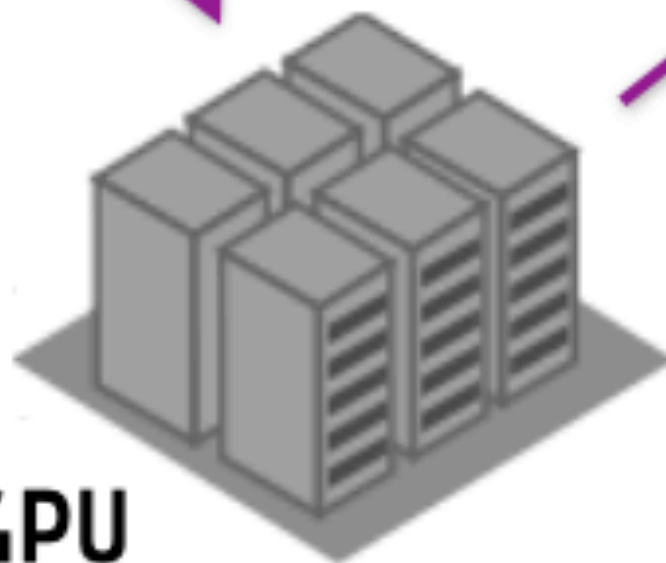


FPGA filter stack
~μs latency

10s Tb/s
100s kHz

On-prem CPU/GPU filter farm
~100 ms latency

10s Gb/s
~5 kHz



Worldwide computing grid
Exabyte-scale datasets

INTELLIGENT DATA REDUCTION

$\sim 1 \text{ PB/S}$

$\sim 1 \text{ PB/DAY}$

Compute Latency

1 ns

1 μs

1 ms

1 s



40 MHz



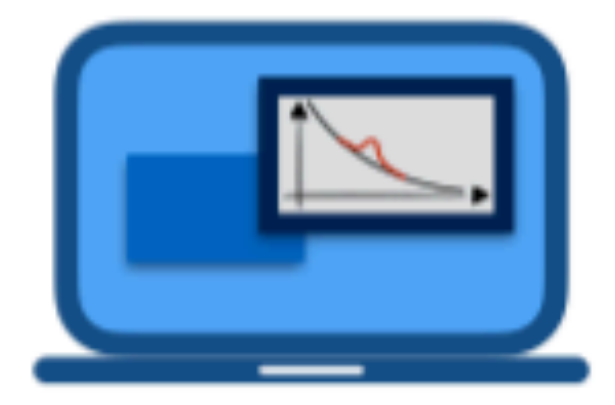
L1 Trigger

100 kHz



High-Level Trigger

1 kHz
1 MB/evt



Offline

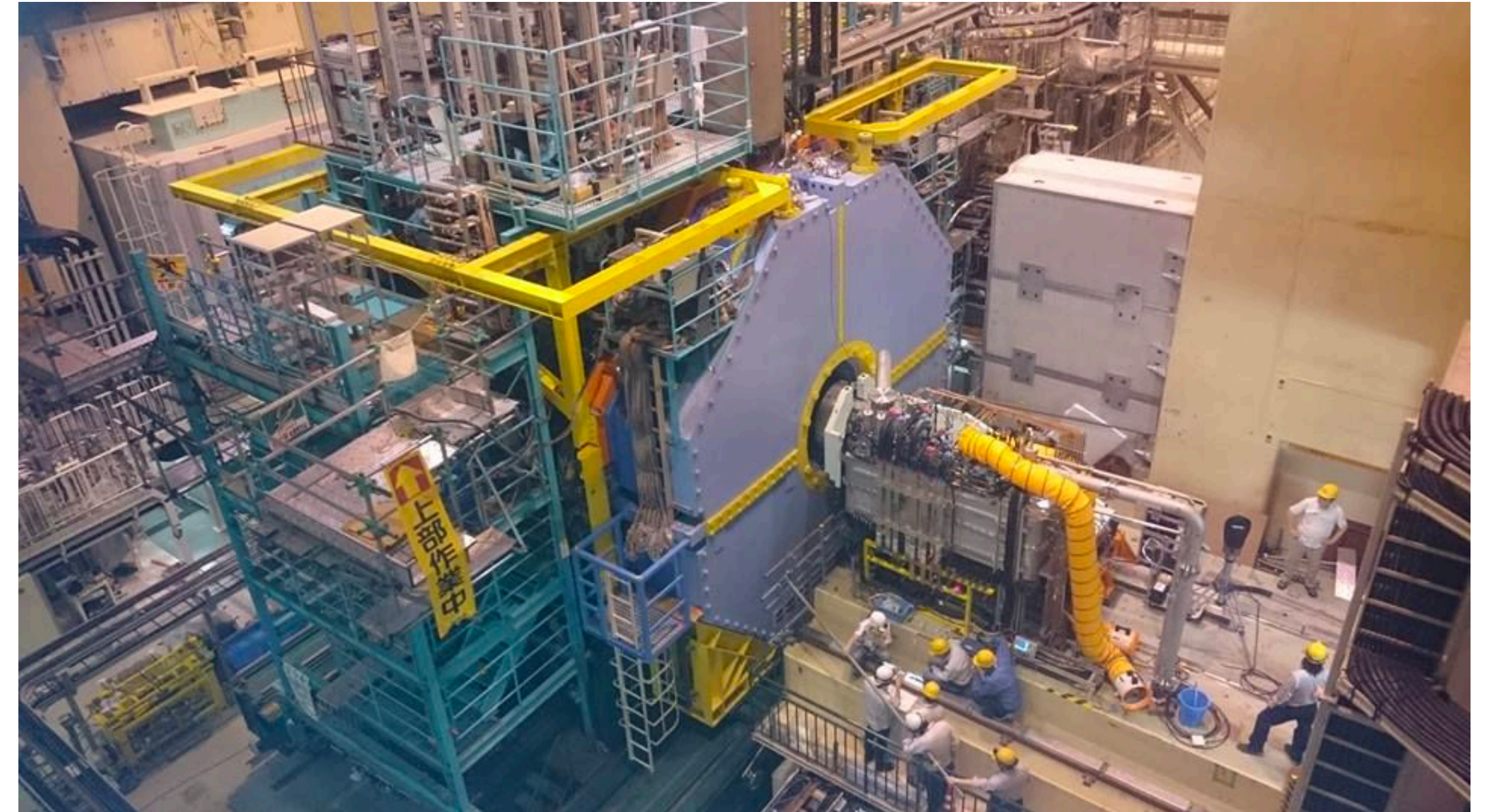


More real-time AI

NA62, CERN



Belle-II, Japan



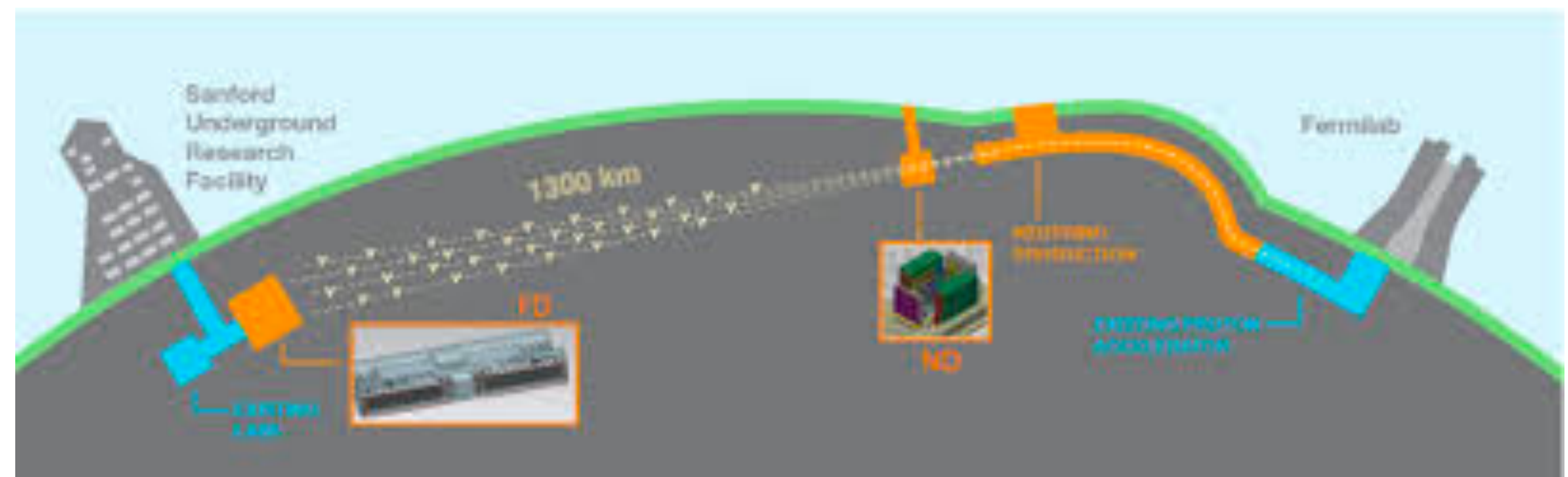
More real-time AI

5. Researchers expect DUNE's data system to catch about 10 neutrinos per day—but must be able to catch thousands in seconds if a star goes supernova nearby.



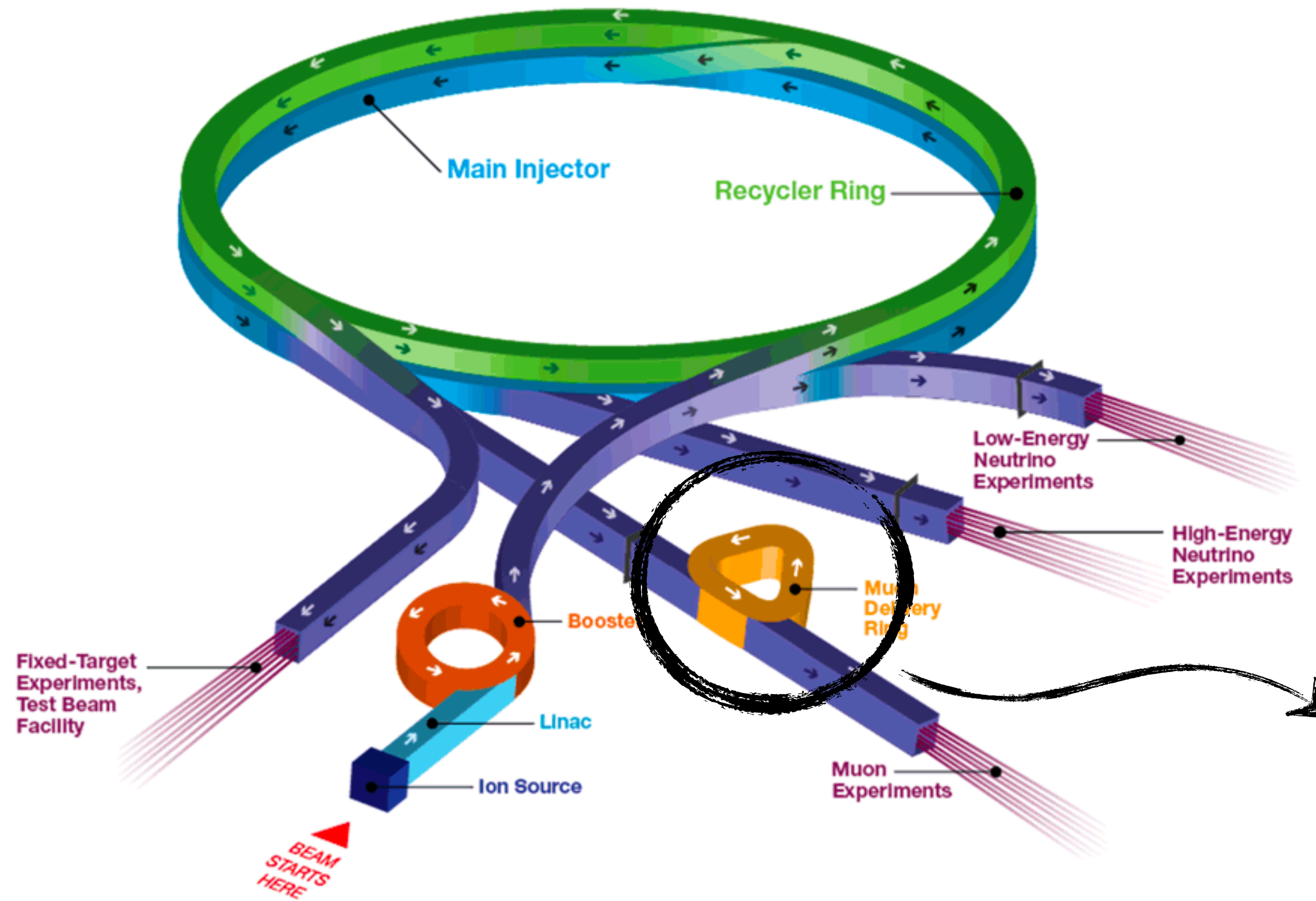
A supernova could produce thousands of neutrino events within seconds.

Real-time filtering at the sub-millisecond scale required for next generation neutrino detectors to identify neutrinos from supernovae, $P \sim (25 \text{ years})^{-1}$

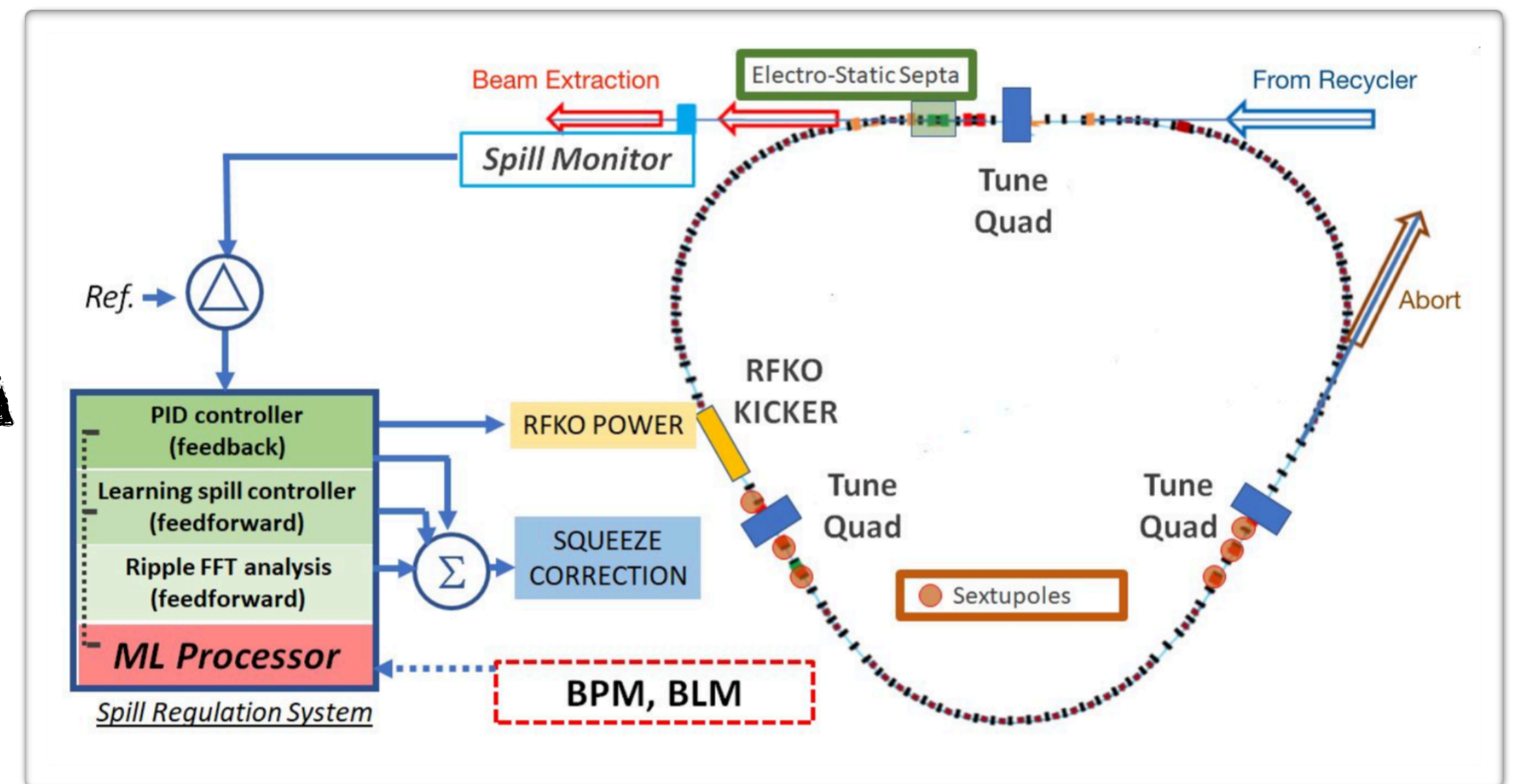


More real-time AI

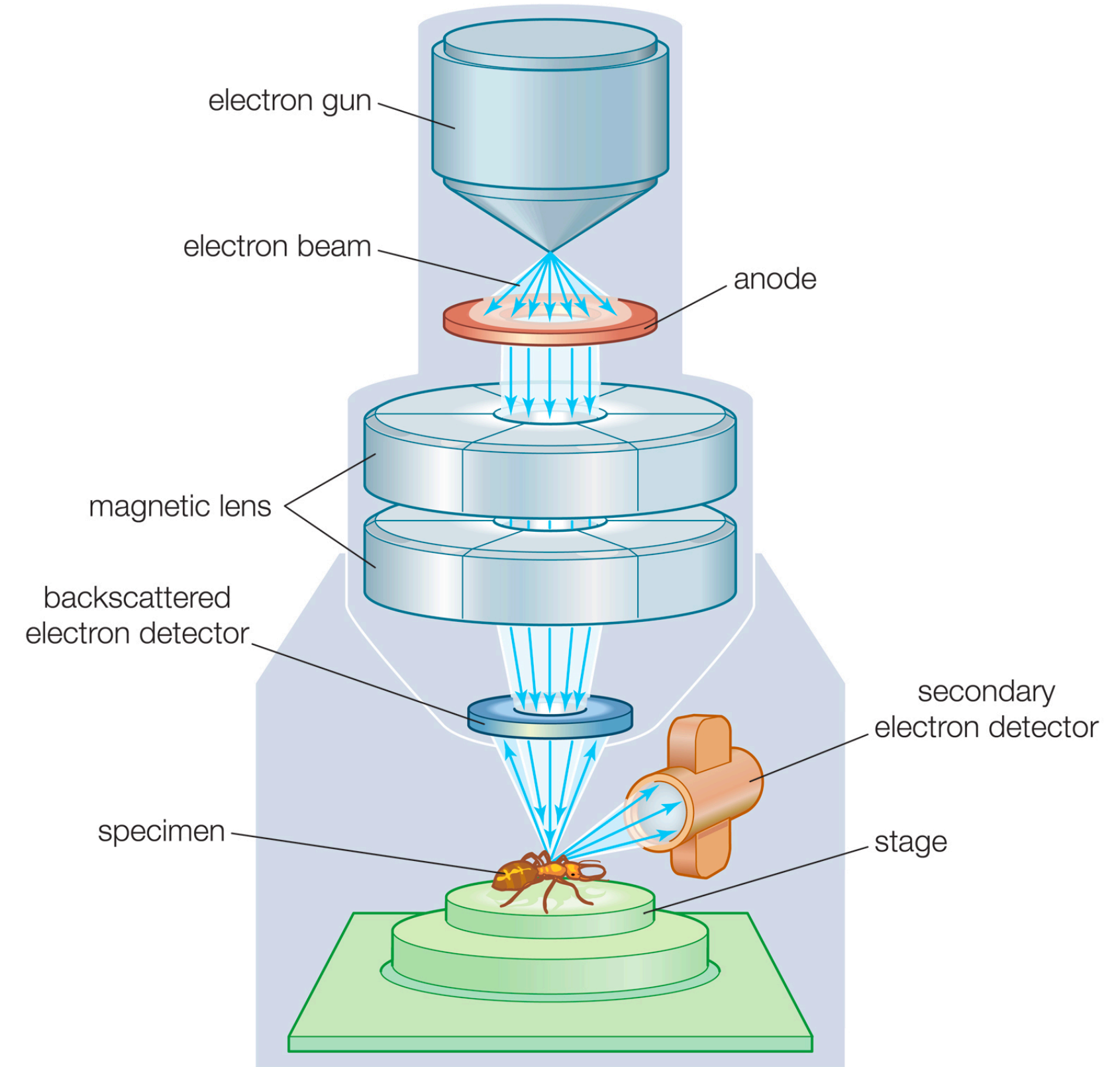
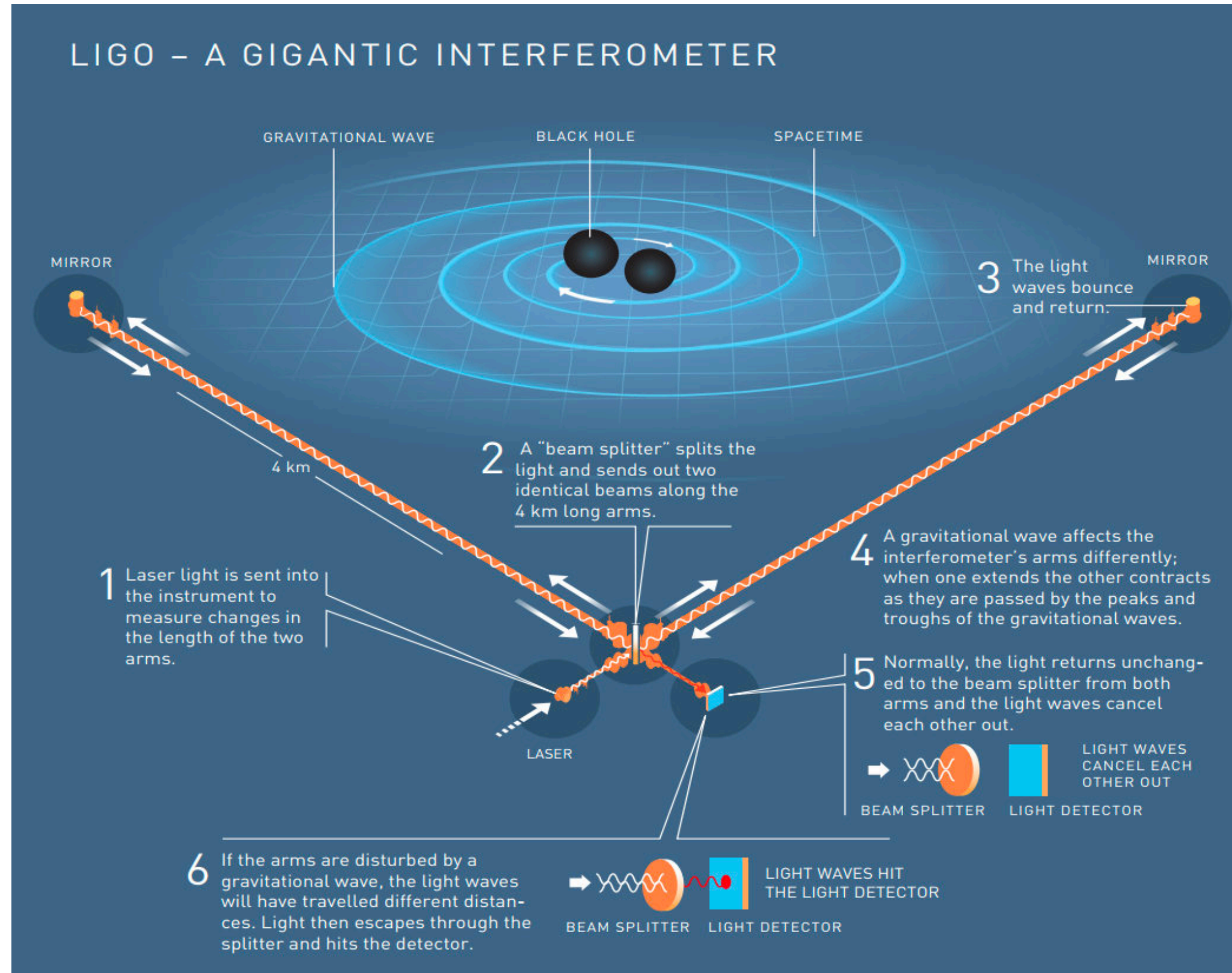
Fermilab Accelerator Complex

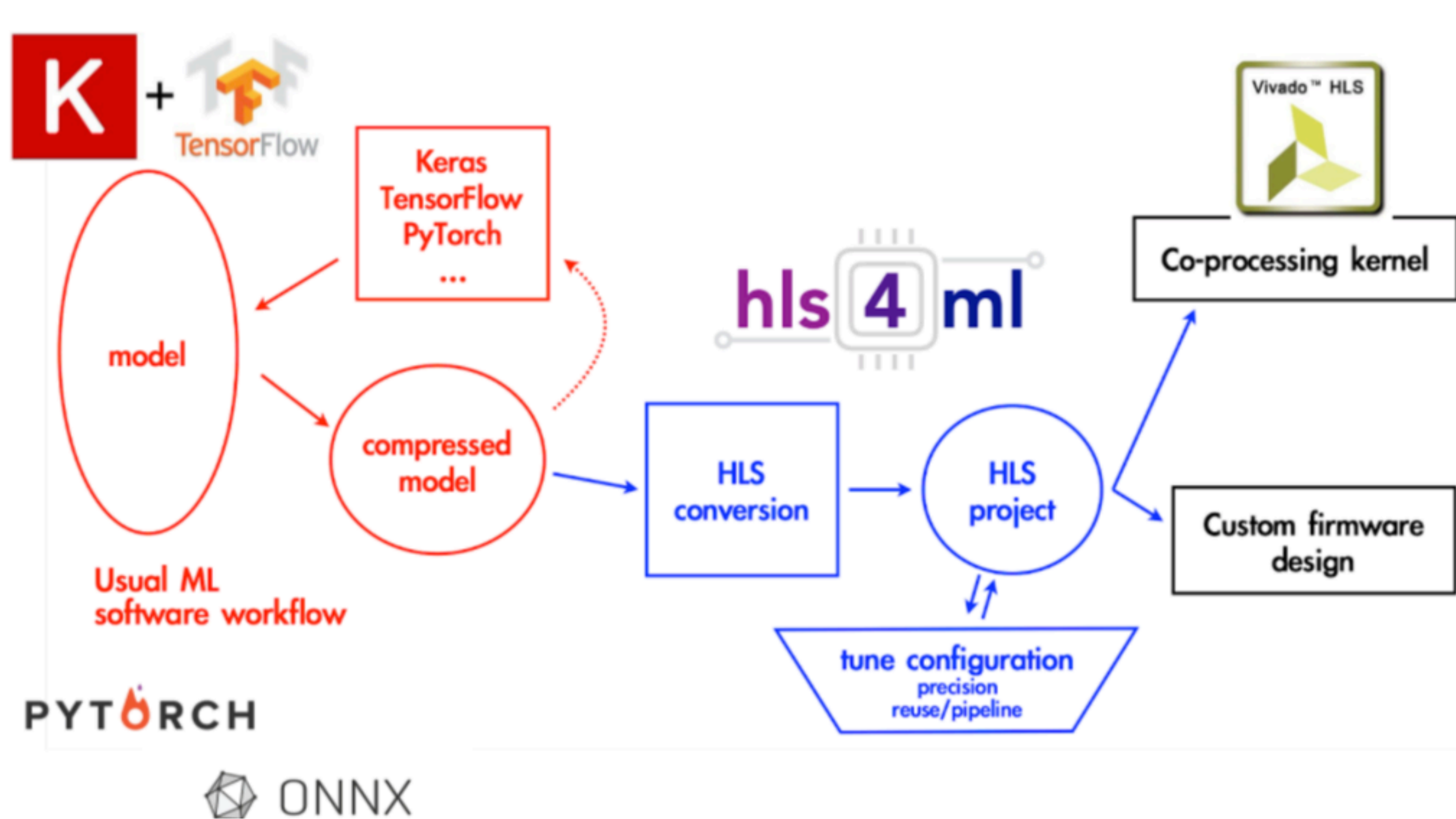


Reinforcement learning for millisecond-scale magnet current tuning required



And more...





Mentor
Catapult
 Coming Soon

Today's **hls4ml** hands on



Google Cloud



- Part 1:
 - Get started with hls4ml: train a basic model and run the conversion, simulation & c-synthesis steps
- Part 2:
 - Learn how to tune inference performance with quantization & ReuseFactor
- Part 3:
 - Perform model compression and observe its effect on the FPGA resources/latency
- Part 4:
 - Train using QKeras “quantization aware training” and study impact on FPGA metrics