

Proposals for S2I2 focus areas from U.S. CMS

This document lists 4 high priority focus area proposals of U.S. CMS for S2I2 and suggestions for the backbone for sustainable software. We are not prioritizing amongst the 4 high priority focus area proposals and we don't list any medium priority focus area proposals for now. These proposals are presented in more detail to concretize the previous discussions.

1. Data Analysis Systems

- LHC: CMS uses as much compute resources for analysis as for central production, and most of the disk storage resources deployed by CMS are needed to host data for analysis. Moreover, the overall path from primary data to publication involves a lot of human effort by physicists (writing loopers, making ntuples, synchronizing results between multiple groups, ...) that has little to do with physics. → do something different that saves both hardware and human (=physicist) resources.
 - Reduce the event size for analysis => reduce total disk storage required
 - Increase the event rate per core of processing => reduce total CPU required, accelerate time to solution for physics.
 - Change paradigm towards “declarative programming” => reduce human effort to do analysis
 - Eliminate “copies” of data by being able to “add user defined columns” to existing data structures => improve functionality, thus accelerating time to solution for physics
 - Explore using industry standard big data solutions to achieve the above.
 - Build upon NSF funded DIANA project here.
- Rationale for focus area prioritization:
 - Physics impact: By greatly increasing the scale of analysis systems, physics measurements will be completed more quickly and with fewer resources, thus allowing either a greater number of measurements or measurements of greater sophistication and reach.
 - Resources impact: This effort could reduce the amount of resources needed for physics data analysis by orders of magnitude, which in itself could cut the overall resources needed for HL-LHC computing in half.
 - Sustainability impact:
 - Interest/Expertise: The university-based DIANA project has already taken the lead in this area, and would form the foundation of this new effort.
 - Leadership:
 - Value: All LHC experiments will benefit from improved data analysis systems.

- Research/Innovation: This area is a natural place to collaborate with computer scientists in the area of big-data tools and declarative programming paradigms.
 - Expected outcome of joint project with S2I2:
 - Scalable data analysis platform that accelerates time to solution by x100 to x1000 using novel infrastructures and systems
 - S2I2 work with OSG-LHC and Ops program on successive prototypes that are used to produce science in Run 3.
 - S2I2 does R&D on “software infrastructure to integrate into a platform”
 - OSG-LHC provides software lifecycle support
 - a. This can alternatively be done inside the S2I2 if the S2I2 is willing to provide software lifecycle support all the way from conception to operation to retirement of software infrastructure. This includes packaging, integration, deployment support, operations support (e.g. ticketing system etc.), and retirement planning.
 - Ops program provide T2 hardware to deploy prototypes on, and sysadmin support to operate the successive platform prototypes as actual facilities.

2. Machine Learning Applications

- Application specific problems pertaining to:
 - Analysis
 - reconstruction/trigger
 - Push down to physics objects and below (basic primitives)
- Expected outcome from project with S2I2:
 - Intellectual guidance on what ML approaches make sense for what types of problems.
 - CS involvement & training
 - Personnel with deep long term expertise may have to be cross funded between RP & S2I2 or Ops & S2I2 if application work is outside the scope of S2I2. As it's not possible to have deep expertise without doing the work, the individuals with the expertise can not be 100% funded via S2I2 if the work is out of scope of S2I2.
 - Software toolkits, e.g. adaptation of industry standard tools for HEP context.
 - Work with OSG-LHC and Ops program on making platforms available to physicists for R&D on applications of ML to LHC.
 - S2I2 work with OSG-LHC to provide deployable containers/software environments that LHC physicists can use to do application R&D.

- c. Fast turnaround processing with near-infinite elasticity:
how to provide access and store output
- Rationale for focus area prioritization:
 - Physics impact: The very fast turnaround of analysis results that could be possible with new approaches to data access and organization would lead to rapid turnaround for new science.
 - Resources impact: Optimized data access will lead to more efficient use of resources, thus holding down the overall costs of computing.
 - Sustainability impact: This effort would improve the reproducibility and provenance tracking for workflows (especially analysis workflows), making physics analyses more sustainable through the lifetime of the HL-LHC.
 - Interest/Expertise: University groups have already pioneered significant changes to the data access model for the LHC through the development of federated storage systems, and are prepared to take this further. Other groups are currently exploring the features of modern storage systems and their possible implementation in experiments.
 - Leadership:
 - Value: All LHC experiments will benefit from new methods of data access and organization, although the implementations may vary due to the different data formats and computing models of each experiment.
 - Research/Innovation: This effort would rely on partnerships with data storage and access experts in the CS community, some of whom are already providing consultation in this area.

4. Reconstruction, Trigger Algorithms

- Address the pileup/multiplicity induced exponential scaling issue of conventional HEP reconstruction algorithms
- Vectorization and advanced architectures (KNL, GPU, FPGA, but also future versions of XEON with ever increasing width of the vector units) of existing and new algorithms
 - How do we guarantee that CMS makes as much use as possible of all the silicon we buy ?
 - Example problem: CMS has a fully functional reconstruction for KNL that is x5 slower than on XEON, despite the fact that KNL has x10 more flops. So there is a relative factor of x50 or so in effectiveness of use of the silicon we buy for KNL vs XEON. If Intel were to merge features of the KNL into future XEON chips, we are likely getting worse use of the silicon we buy unless we do some serious R&D on our algorithms, and their implementations.
 - Class of reconstruction algorithms that are re-written and optimized for vectorized architectures
- Rationale for focus area prioritization:

- Physics impact: Pileup mitigation will be the fundamental technical issue of HL-LHC physics, and improvements to the reconstruction algorithms designed for modern architectures will be important for realizing the physics potential of the detectors.
- Resources impact: There are significant computing resources at HPC centers that could be made available to HL-LHC experiments at little cost, but many optimizations of existing code will be required to fully take advantage of them.
- Sustainability impact:
- Interest/Expertise: University groups are already making progress in the use of chipsets such as GPUs for specific HEP applications, such as track pattern recognition and fitting. New detector elements that are expected for HL-LHC upgrade could especially benefit from pattern recognition on new architectures, and groups that are building these detectors will likely get involved.
- Leadership: It is likely that there will be some overlap with work done at DOE HPC centers, but NSF HPC centers might require independent efforts.
- Value: All LHC experiments will benefit from these techniques, although many implementations will likely be experiment-specific given differing detector configurations.
- Research/Innovation: Much assistance will be required from the computing and software engineering communities to help prepare algorithms for new architectures.

5. Backbone for sustainable software

- Modernization of software development process for scientists
 - Individual experts can improve the software by orders of magnitude by just understanding the algorithms and intended optimizations and applying the appropriate optimizations. How can we improve the overall process that the quality of software and its optimization is better out of the box.
- Support for testbeds for validation and scaling
 - ops program has the hardware. S2I2 has people and R&D for capabilities. So scaling and performance verification vs scale should be a collaboration between S2I2 and Ops program.
- Tool support
 - Packaging, distribution
- Training
 - Best practices, training, workshops, ...