



IRIS-HEP Blueprint Meeting: *Analysis Systems on Scalable Platforms*

Live notes

Agenda: <https://indico.cern.ch/event/820946/>

Livenotes from AS topical meeting:

<https://docs.google.com/document/d/10TVHqzfBUacd8pL89loDUTqTSTw8KKfaffqVRx-YYaY/edjt>

Registered Attendees

In person (not complete)

Mark Neubauer (U. Illinois), Kyle Cranmer (NYU), Ben Galewsky (NCSA), Nils Krumnack (Iowa State University), Matthew Feickert (SMU), Lukas Heinrich (CERN), Mason Proffitt (U of Washington)...

See [participant list in agenda](#).

Remote (Vidyo)

Gordon Watts (U. Washington)
(Not complete)

Agenda & Notes

Welcome (Kyle)

Blueprint overview (Neubauer)

SSL (Rob)

Program of work:

Provide the Institute and the HL-LHC experiments with scalable platforms needed for development in context. SSL is on the path to production

- Provides access to infrastructure and environments
- Organizes software and resources for scalability testing
- Does foundational systems R&D on accelerated services
- Provides the integration path to the OSG-LHC production infrastructure

Challenges:

- community platform
- support groups and projects
- bespoke resources & configurations
- declarative & reproducible deployments
- services to build & manage artifacts
- scalable up and back down

Opportunities:

- spark ad hoc collaborations across traditional organizational boundaries (e.g. experiments)
- contributions and involvement from diverse resource providers
- new modes of infrastructure development, supporting more rapid innovation
- redeployable artifacts and reproducible patterns

Ben: Storage and networking becomes a big issue for the large data sets we're dealing with

Rob: Right now cloud ~30% (spot pricing) more expensive than on premises

Google provided several million cores to an MIT researcher over a weekend for ~\$20k

Nils: Memory/core and i/o bandwidth *could* be a problem on the cloud

AS (Kyle)

Just finished a 2-day workshop to go over the milestones and deliverables for the rest of Y1 and Y2

n.b. AS scalability milestones Y2Q1 --> SSL readiness; moving testing on whatever resources now onto a more "SSL-managed" infrastructure.

- prototypes of analysis systems components -> SSL (Aug '20)
-

Recent Ad-doc AS demos:

- REANA/RECAST demo'd at Kubecon 2018 : real analysis reproducibility, also shows that HEP can engage with modern tools and communities
- Scalability demo of Higgs rediscovery on Kubernetes (200 GB/s, 70 TB)

ServiceX

Coffea

Lindey: This week we were about to get a simple analysis through ServiceX + Coffea

MADMINER -- containerized workflows with mix of cpu and gpu/tpu (integration of sim, ml, stat inference); run on k8s w/ reana

AMPGEN, pyhf (fitting as a service)

- Have resources setup and available for users to upload information (e.g., pyhf json) and then the service performs the fit. Saves the user from being *required* to set things up on their own (services are simple, but not everyone has a nice GPU cluster ready to go)

Types of systems we'd like to test on:

- public cloud
- university resources
- grid/osg
- hpc systems
 - kubernetes as common substrate / infrastructure layer
 - non-priv k8s on HPC
 - move between grid, cloud, hpc

Containers on grid

- Hyperparameter tuning using Docker and the LHC WCG [[ACAT 2019 poster](#)]
- ATLAS machine learning Docker images for ML on LHC WCG [[talk](#)]

Adapting reana (user facing) for hpc w/ vc3

Reality - schedule

[Vast.ai](#) - this is a brokerage for GPU / ML software stack resources, a “narrow” cloud if you will (Andrew C.)

Lincoln (SSL architecture)

Desires for SSL:

- **a community platform**
 - Open to all working on software infrastructure in HEP
 - CILogon, Globus to provide single-sign on and federated identity
 - Lightweight user and group (project) management system
 - Infrastructure itself composable, reusable
 - CS cluster at UC good example
- **supports groups and projects**
 - Web and CLI interfaces for user management
 - Groups organized as a tree structure with arbitrary depth
 - Users can invite others, create sub groups, etc.
 - Smart clients inspect the tree and implement appropriate provisioning of resources
- **bespoke resources & configurations**
 - Mix of bespoke dedicated resources and capability for users to bring allocations on others
 - Container-based service orchestration on dedicated resources
 - VC3-like technology to connect to HPC/HTC resources for batch
 - Facilitate integration of commercial cloud resources when needed
- **declarative & reproducible deployments**
 - Want infrastructure built under the SSL to be easily reusable and deployable to other sites
 - Declarative nature of Kubernetes is a good fit and gets us a long way down that road.
 - SSL as an incubator for projects which then “graduate” to become full-fledged infrastructures that run on production resources.
- **services to build & manage artifacts**
- **scalable up and back down**
- **reduce cognitive load for developers and deployers**

Right now using GKE to test SSL deployment, want to soon pull that off into in-house resources

From university point of view would like to have simple version of what do they need to provide

- Eg. kubernetes
- Eg. a REANA instance

- From campus IT/research technology point of view ability to make a compelling case to the provost that they are providing resources that enable good science
-

Suggest some light-weight mechanism for discovery of resources

- eg. some machine with some special accelerator

Value of reporting / showing science that is happening on the contributed resources, incentivize resource providers at the universities. Capturing success stories, so university community understands how they can benefit from investments to shared campus resources, including staff.

Having a dashboard to help communicate the contributions.

Products that can be used outside HEP - giving higher visibility.

Will SSL provisioning, panel, etc. tools be open sourced and productized?

SSL should probably provide boilerplate text for facilities

Lukas (Rediscovering Higgs on google cloud)

https://docs.google.com/presentation/d/1zuNbJOcl4U1E5LHMqv0B5FDG-6ZVfAitlkjPEucQb4M/edit#slide=id.g52246c4d4c_0_4

CMS open data used (70 TB), ~25000 files

Legacy software from CMS scientific software stack, put on k8 and execute at large scale.

Last slide:

What did we learn from this?

- Higgs still exists. Good! :)
- Google Network can serve extreme data rates into compute nodes (2Gbps/core) and needs to be tamed:
 - incoming data staged in via gsutil.
cap bandwidth per code.
 - Disks that can reliably take in data at this rate
expensive and scarce (Local SSDs).
Write to Memory

- At highly parallel workloads, scheduling parameters becomes important
 - How many pods/s can you get into the API server
 - How quickly do they show up on the nodes
 - CERN: tweak qps in masters.
 - GCP: multi-cluster + threaded submission

Sanjay (Red Hat Openshift)

Google

- Researcher credits: https://edu.google.com/programs/credits/?modal_active=none

Afternoon discussion

Thinking that REANA is a first deployment target for SSL

[However, ServiceX is deployed on GKE and it would be interesting to use the HELM chart to deploy on k8]

REANA - what does a deployment look like on SSL?

Openshift (gives Kubernetes) cluster

Deploy 1st instance of REANA on SSL

Different set of people (Ben) than did the original (IRIS-HEP) deployment

Templated python scripts → HELM charts, YAML (goal to make it easily deployable and reproducible by anyone)

Things that need to be in the substrate:

- Storage
- Internet ingress
- Load balancing
- Object storage

Is there one SSL (as a service) or a standard?

We don't want site admin having to follow notes off a twiki to stand up kubernetes. We'd like the admin to connect their nodes to SSL console and everything is automatically deployed.

Process to join must be lightweight

I see a deployed pattern I should be able to deploy them to a local environment.

Dues: If you use SSL then you should at least publish deployment instructions so others can reproduce it.

An R&D Hybrid Cloud provider

Try to make the substrate as compatible as possible with CERN-IT and FermiLab

Mark: Q: If a “**federated**” **set of kubernetes clusters** is the substrated approach for SSL, do we lose anything there? Are there any blockers to this approach for the AS R&D plans on SSL?

A: does not seem so, but maybe HPC integration suffers a bit?

Rob: For scaling test, federating SSL clusters will be desirable.

Kyle: **Metric data** would be useful to have from SSL. Mark/Rob: yeah, this was discussed and dashboards and retention of results, mines ElasticSearch, [prometheus](#) (more standard tool)? Developing the complete suite for logging and metrics collection is out of scope of SSL, but maybe provide a few standard tools and a repository to collect notes on best practice.

Gordon: (later) - there were comments that not everyone in our group knows how to build a helm chart. While that is true, often the stuff we are all working on is going to fit into a helm chart eventually - so there is something to be said for being able to at least use it. docker/desktop can run kubernetes and I don't know how hard it is to make a flexible chart that can scale from a single node cluster to many, but this means that the person working on the component can basically run it in the environment they will eventually have to run in.

Burden on R&D areas to define the metrics and guide SSL on what to retain.

SSL: Service Level Agreement

- What is going to be monitored in the scalability tests (the y-axes)
 - Eg. things that should go into elastic search like throughput, time
- What configuration parameters are going to be scanned (the x-axis)
 - Eg.
- Agreement to publish deployment artifacts
- Request time / scheduling for scalability tests
- Probably make sure that this information is agreed upon before devoting these major resources

Opportunities & Planning for resources

Regional or not? Not necessary but could be useful

A set of tools that would allow one to re-create the SSL environment at another location.

NCSA

ISL (Integrated Systems Lab)

Openstack cluster

Blue Waters (Neubauer allocation - kind of short term)
Illinois Campus Cluster (opportunistic use of GPUs)
NSF MRI Deep Learning research platform

Redhat (Openshift)

[Massachusetts Open Cloud?](#)

CERN?

NYU (Kyle add this)

Fermilab (Rob will add this)

BNL?

SDSC (Edgar will add to this)

- PRP

Saturday, June 22

AS milestones

		Y2 planning								
	Analysis Systems	Type (M/D)	Y1Q1	Y1Q2	Y1Q3	Y1Q4	Y2Q1	Y2Q2	Y2Q3	Y2Q4
3										
4	Description									
5	Organize topical meetings, Analysis System group meetings, etc.									
6	List publicly-accessible repositories and other relevant documentation on the iris-hep.org website									
7	Collect and curate example analysis use cases with some existing reference implementation									
8	Survey of analysis systems efforts in the field to aid in planning for topical workshop									
9	Blueprint workshop coordinating resource needs for evaluating analysis systems coordinated by SSL with participation of operations program				Scheduled for June					
10	Develop initial specifications for user-facing interface to analysis system components									
11	Prototype awkward-array analyses in the scientific Python ecosystem									
12	Initial roadmap for ecosystem coherency									
13	Develop initial design for interface of analysis query system to the IDDS									
14	Translate analysis examples into new specifications, provide feedback, iterating as necessary									
15	Initial roadmap for high-level cyberinfrastructure components of analysis system									
16	Benchmarking and assessment of existing analysis systems									
17	Implement prototype query-based and cache-aware dispatch									
18	Establish analysis description database schema and integrate with archival tools like CAP, INSPIRE, HEPDATA, etc.									
19	GPU/accelerator-based implementation of statistical and other appropriate components									
20	Move prototypes of analysis system components to SSL									
21	Benchmarking and assessment of prototype analysis system components									

AS June Meeting Topic (also relevant blueprint WS)

SSL Blueprint

Implied SSL cyberinfrastructure for AS success:

- lowering barriers for participation in analysis systems
- single-PI, postdoc, students can contribute meaningfully to analysis
- connection to diversity, outreach
- Connection to onboarding & transfer (forward evolution) of analysis, especially for single-PI groups
- establishing how SSL can coordinate w/ campus IT efforts to replicate environments
 - Eg. NYU Tier3 and k8 research cluster can be used for testing this
 - Capture patterns
- Could be a pattern in the future for transitioning Tier2
- Edgar: UCI/Anyes - t3-box - needs to be k8s'd. SSL-box --> something like an SSL-distribution; data and caches and submission environment. Three objectives:
 - Making life easier for sys admins
 - Make infrastructure more flexible

- Increase capabilities for single PI
- Empowering the small groups - important for NSF funded institutions.
 - SSL side: easy deployment patterns
 - Why? If a small research group had an idea for an infrastructure component and they wanted to develop it in some isolated prototype environment
 - Examples:
 - RECAST was a new style of analysis that required some new **infrastructure components**. CERN cloud services facilitated that
 - Reproducible Open Benchmarks (ROB) product (eg. for benchmarking ML reconstruction and analysis components)
 - Alternative Event Processing
 - pyhf (being able to use GPUs and TPUs for benchmarking and scaling studies)
 - aligned w/ google - cloud credits
 - partnering with other groups (e.g. Stephen mentioned there is TPU research group)
 - Similarly, [ATLAS ML base Docker images](#) (designed to have a minimum useful base env, but has need for GPU resources for testing. Look to GRID, but elsewhere also)
 - if SSL can provide these quickly it would move these developments forward
 - Analysis Systems Onboarding / forward evolution User Story
- [kubvirt.io](#)
- SSL as a matchmaker / hub for relationships.
- BinderHub connected to SSL authentication as one of the capabilities in the SSL service catalogue
 - [Binder](#) (Example: [pyhf example stat analysis](#))

Rob: Openshift is popular because it makes Kubernetes easy to deploy and work
CVMFS works on Kubernetes (thank Igor and Dima)

1. simplifying life for site admins
2. empowering single PI's
3. increasing capabilities by making infrastructure more flexible

SSL milestones

Y1Q4

- roadmap of initial cyberinfrastructure components from AS.
- SSL substrate project (see below)
- REANA service deployed (Helm'd)

- ServiceX deployed (gke->ssl procedure defined)

Y2Q1

- A general app monitoring, benchmarking, metrics collection service
- Next gen Tier2 & Tier3
- HEPIX, Oct 14-18: <https://indico.cern.ch/event/810635/abstracts/>
- CHEP

Y2Q2

- Engagement with clients
 - reservation scheduling

Y2Q3

-

Y2Q4

-

SSL Substrate Project

- Evaluation of the base layers ([Openshift](#), k3s, ...)
- Write down the basic architecture, basic "SSL distribution"
- Build a thin k8s substrate layer from contributed resources
- Partners (identify point person)
 - Chicago (River - Lincoln/Andrew)
 - Illinois (Illinois Campus Cluster (ICC) - Tim Borner)
 - New York (legacy T3 & research tech's k8 cluster)
 - San Diego (campus k8s + PRP + UC Irvine - Edgar)
 - CERN (coordinate with existing patterns and successes)
 - Princeton (to be developed)
 - Fermilab (Glen Cooper, Liz, ... to be developed)
 - BNL (challenges)
- Documentation for administrators
 - Need to think about target audience not physicists
 - For example: at NYU we pointed IT team to <http://reanahub.io/>
 - Which points to <https://reana.readthedocs.io/en/latest/administratorguide.html>
- User and group management service
- Resource registry service}
- CI environment services
- Container registry service
- LHC software environment
- Deployment catalog service
- Need to articulate the vision to contacts at the partner institutions for planning purposes

SSL Service Level Agreement

- a loose contract - declarative
- what resources are there, needed?
- what do you want in terms of monitoring, etc.
- Agree to publishing deployment pattern and other artifacts
- given a user story, make sure the prerequisites are in place
- tooling for specific testing, metrics collection

SSL Scheduling

- Need to avoid everyone trying to do their scalability test at the same time on limited resources

SSL [CRM](#)

- Manage relationships for both
 - Resources: eg. partners in the SSL Substrate Project
 - Clients of SSL
- A serious CRM or just something light weight like slack / google groups etc.
- Do we need dedicated effort / role for managing this. Administrative / project management / community management

Connection to community infrastructure R&D

- WLCG light-weight sites
- Possible evolution of Tier-3 and Tier-2 resources
- [HEPIX](#)
- Computer science conference where accelerators are discussed
- GDB (too early?)
- Rob: This blueprint meeting has spawned an idea to deploy Tier-2 / Tier-3 and reach out to WLCG to engage with their planning (submit an abstract to HEPiX and plan to have attendance)

Outcomes from Infrastructure breakout

- Analysis service deployments imply re-thinking capabilities and flexibility of resources at existing and new LHC computing sites
 - simplifying life for site admins
 - empowering single PI's
 - **increasing capabilities by making infrastructure more flexible**

- Must identify a path to get there for these use cases:
 - One person site into an SSL
 - A "mixed T2" batch and SSL (re-blueprinting the Tier2 complex, solutions for Tier3s)
 - already deployed kubernetes? how to contribute? (e.g. PRP)
 - impedance circuitry to institutional campus research computing to contribute to SSL need to be developed (Tim) - establishing the pattern and process for campus research computing. (Irvine another use case)
- Converting a site, handling condor, CE, storage..
- Role of incentives - admins + PIs (already constrained) - value proposition has to be clear.
- Importance of engaging community wide - hepix, wlcg gdb, etc.
- Must make solid connections and provide conduits for industry contacts
 - Hugh Brock <hbrock@redhat.com> from Red Hat - for engagement
 - Cloud
 - Karan Batia(sp?) - Google
 - azure (contact gordon)
 - aws (sanjay), csu fresno contact
- "Before I thought, now I think" outcomes
 - K8 as common denominator
 - Rob: This blueprint meeting has spawned an idea to deploy Tier-2 / Tier-3 and reach out to WLCG to engage with their planning (submit an abstract to HEPiX and plan to have attendance)
 - Kyle: Fitting service / TPU engagement
- Identified a multi-site "substrate" project
- Service level agreement

Outcomes from Analysis breakout

Communicate AS goals

- Create greater functionality
- Reducing time to insight
- Lowering barriers for smaller teams, and
- Streamlining analysis preservation, reproducibility, and reuse

Communicate projects where they fit

Discussion of relevant milestones and deliverables for AS and how they translate into an ask for SSL. These plans and ask form an outcome for the analysis systems topical meeting. Communication of these plans and the request to SSL is an outcome of the blueprint workshop.

- Collect and curate example analysis use cases with some reference implementation
 - User stories:

- Event Processing (left blue box [on this page](#))
 - Template Fitting (center box)
 - Reinterpretation (right box)
 - Analysis evolution cycle: on-boarding / forward evolution / preservation
 - Cross-cutting “vertical slice on it’s side” :-)
 - Aligns with 3 AS goals
 - Less clear how it’s related to SSL
- Benchmark examples (concrete examples with reference implementation)
 - Z-peak
 - Reference: Coffea(awkward)+Service X (single-threaded from a single file with uproot: [zpeak.ipynb](#), [talk at HOW2019](#))
 - TSelector (exists and one-to-one with awkward capabilities: [ZPeak.h](#), [ZPeak.C](#))
 - RDataFrame ([this one](#) is more basic than the awkward and TSelector based ones, not a one-to-one comparison)
 - ATLAS Run3 EventLoop style (possible)
 - KubeCon CMS OpenData Higgs Demo
 - Reference: Containerized Root analysis shown at KubeCon
 - Awkward (to do)
 - A template fit with control regions and systematics (based on a ttbar example)
 - Need to check on the status of reference implementation for ttbar example. Alex did TRex fitter in top context, so should be good
 - Idea that evolved at the meeting: implement a fitting service (eg. send pyhf specification, run it on a GPU or TPU-enabled resource)
 - Analysis evolution cycle:
 - Examples in ATLAS using GitLab CI compatible with CERN Analysis Preservation, etc.
- Develop initial specifications for user-user-facing interfaces to analysis system components
 - Almost done
 - Starting with writing user stories for the 4 cases above
 - To each user story, description of a corresponding product
 - Event Processing:
 - Coffea, RDataFrame, ...
 - Template Fitting:
 - TRexFitter, HistFitter, (built from components like pyhf, HlstFactory, RooFit, RooStats inside)
 - Reinterpretation:
 - Recast, REANA
 - Analysis evolution cycle: on-boarding / forward evolution / preservation:
 - GitLab CI, CAP, etc.

- Based on experience with these existing products and the core user story, we will base a minimal initial specification for user-user-facing interface
- Translate analysis examples into new specifications, provide feedback, iterating as necessary
 - Partially done; target Dec 2019
 - See list above, some complete, some to do
- Benchmarking and assessment of existing analysis systems
 - **Target March 2020**; topic of this meeting
 - For each of the benchmark examples
 - Identify resources and capabilities needed for scalability testing of
 - **reference implementation**
 - **Action:** This should happen in the next ~month
 - each alternative implementation in a prototype system
 - **Action:** deadline for this information in ~March.
 - How should this be communicated, what does SSL need to know?
 - **Action:** schedule some meeting between SSL and corresponding product team in AS
 - Discuss: Analysis Evolution Life Cycle
 - The deployment patterns for the GitLab contribution are relevant
 - Parallel / reinforces goals in SSL for reducing barrier for small groups
 - Potential Integration test: can we pull out of CAP and run workflow on REANA instance deployed out of SSL
 - Need to decide if there is an ask to SSL for this example
 - Need to work out scheduling, communication between SSL and AS
- GPU/accelerator-based implementation of statistical / fitting tools (and other relevant components)
 - **Target March 2020**; topic of this meeting
 - Seems oddly specific in this context, but maybe good to call out this specific case
 - A fitting service would satisfy this. Not developed
 - Stephen from Google mentioned team that collaborates with researchers to optimize systems for TPUs.
 - **Action:** organize this coordination under SSL
 - Deadline to develop a prototype system ~Dec
- Move prototypes of analysis system components to SSL
 - **Target August 2020**; topic of this meeting
 - Scheduling
- Benchmark and assessment of prototype system components
 - **Target August 2020**; topic of this meeting
 - Scheduling

Communicate progress via drop in notes in EB + ad hoc meetings as necessary

Feedback on Blueprint

- Some preparatory documents
- Back to Ben's point about IRIS needing community development / management etc.
- Start planning earlier (obvs)
- AS contribution to accelerated inference blueprint