

IBM at ISC19

IBM Spectrum LSF & Scale User Group

IBM Spectrum LSF Update

Bill.McMillan@uk.ibm.com

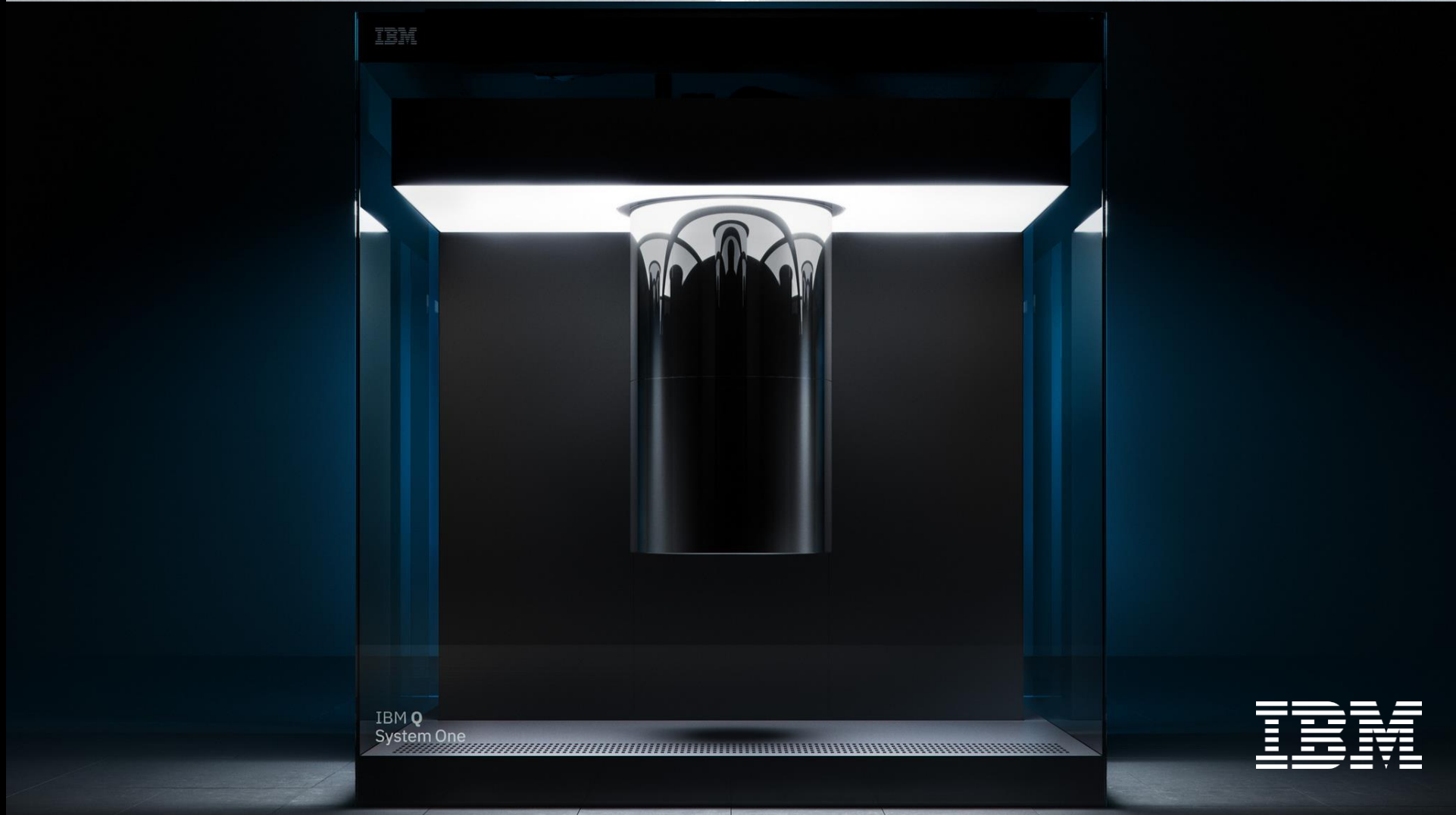
Global Offering Leader

IBM Spectrum LSF

IBM Cognitive Systems



LSF Update, June 2019 / © 2019 IBM Corporation



Notices and disclaimers

© 2019 International Business Machines Corporation. No part of this document may be reproduced or transmitted in any form without written permission from IBM.

U.S. Government Users Restricted Rights — use, duplication or disclosure restricted by GSA ADP Schedule Contract with IBM.

Information in these presentations (including information relating to products that have not yet been announced by IBM) has been reviewed for accuracy as of the date of initial publication and could include unintentional technical or typographical errors. IBM shall have no responsibility to update this information.

This document is distributed “as is” without any warranty, either express or implied. In no event, shall IBM be liable for any damage arising from the use of this information, including but not limited to, loss of data, business interruption, loss of profit or loss of opportunity. IBM products and services are warranted per the terms and conditions of the agreements under which they are provided.

IBM products are manufactured from new parts or new and used parts. In some cases, a product may not be new and may have been previously installed. Regardless, our warranty terms apply.”

Any statements regarding IBM's future direction, intent or product plans are subject to change or withdrawal without notice.

Performance data contained herein was generally obtained in a controlled, isolated environments. Customer examples are presented as illustrations of how those customers have used IBM products and the results they may have achieved. Actual performance, cost, savings or other results in other operating environments may vary.

References in this document to IBM products, programs, or services does not imply that IBM intends to make such products, programs or services available in all countries in which IBM operates or does business.

Workshops, sessions and associated materials may have been prepared by independent session speakers, and do not necessarily reflect the views of IBM. All materials and discussions are provided for informational purposes only, and are neither intended to, nor shall constitute legal or other guidance or advice to any individual participant or their specific situation.

It is the customer’s responsibility to insure its own compliance with legal requirements and to obtain advice of competent legal counsel as to the identification and interpretation of any relevant laws and regulatory requirements that may affect the customer’s business and any actions the customer may need to take to comply with such laws. IBM does not provide legal advice or represent or warrant that its services or products will ensure that the customer follows any law.

Notices and disclaimers continued

Information concerning non-IBM products was obtained from the suppliers of those products, their published announcements or other publicly available sources. IBM has not tested those products about this publication and cannot confirm the accuracy of performance, compatibility or any other claims related to non-IBM products. Questions on the capabilities of non-IBM products should be addressed to the suppliers of those products. IBM does not warrant the quality of any third-party products, or the ability of any such third-party products to interoperate with IBM's products.

IBM expressly disclaims all warranties, expressed or implied, including but not limited to, the implied warranties of merchantability and fitness for a purpose.

The provision of the information contained herein is not intended to, and does not, grant any right or license under any IBM patents, copyrights, trademarks or other intellectual property right.

IBM, the IBM logo, ibm.com and [names of other referenced IBM products and services used in the presentation] are trademarks of International Business Machines Corporation, registered in many jurisdictions worldwide. Other product and service names might be trademarks of IBM or other companies. A current list of IBM trademarks is available on the Web at "Copyright and trademark information" at: www.ibm.com/legal/copytrade.shtml.

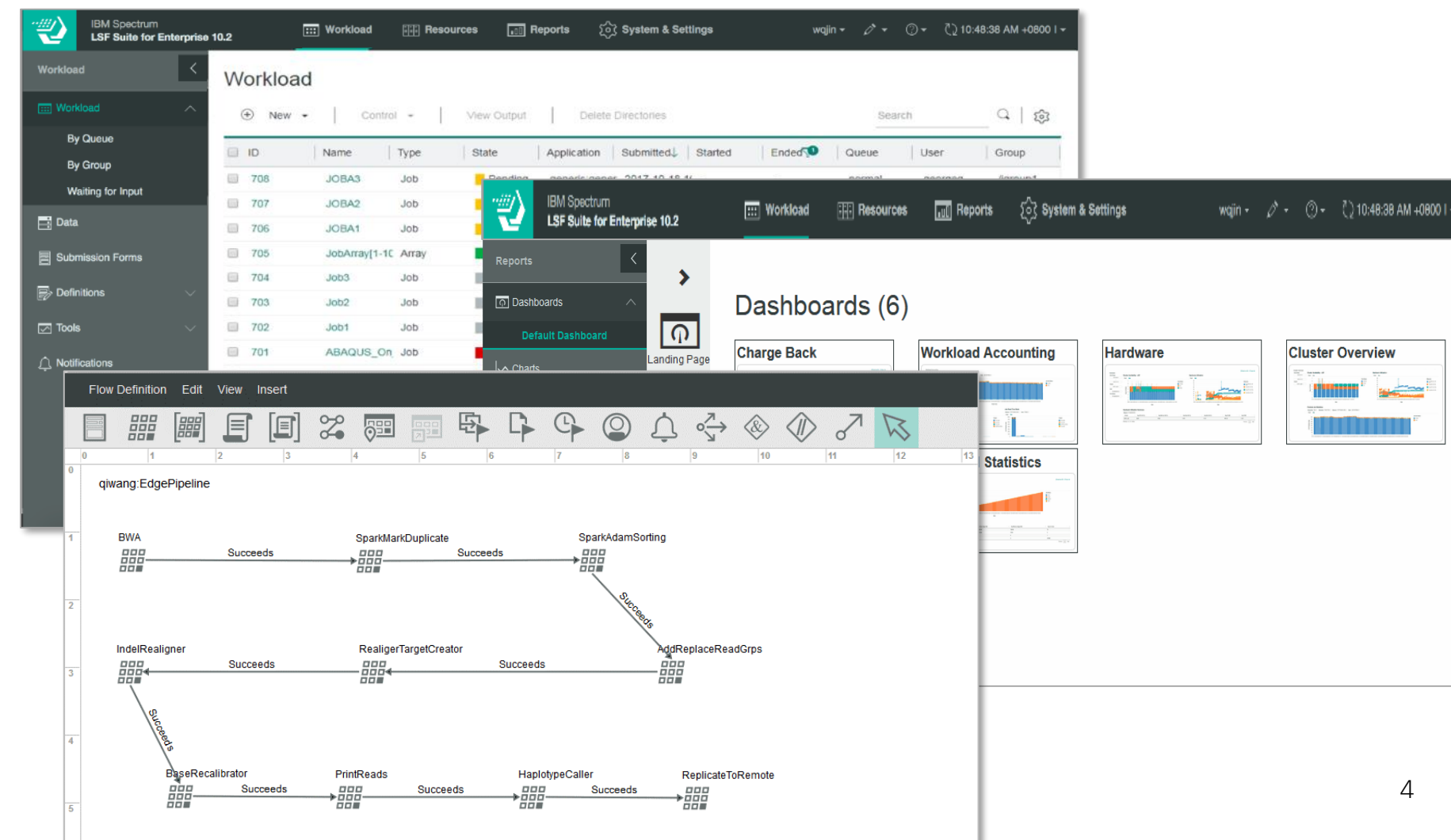
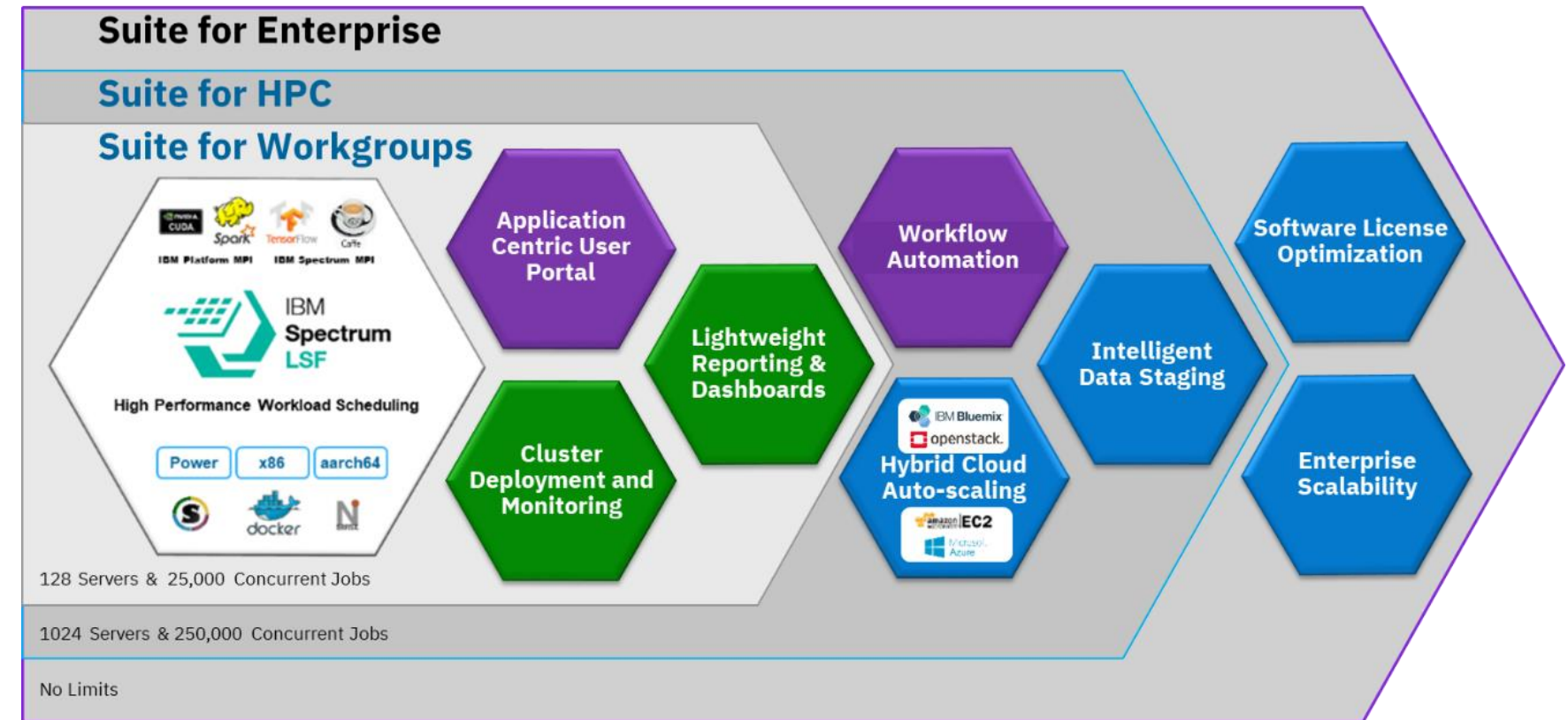
IBM Spectrum LSF Suite

Everything you need for Workgroup, HPC or Enterprise scale.

- **Enhanced Utilization of assets through effective scheduling and sharing policies**
- **Enhancing User Productivity through ease of use, accessibility and simplification**
- **Operational Efficiency through insight of how the HPC environment is being used**

The LSF Suite for HPC is available at no charge via the IBM Academic Initiative

Hourly Pricing now available.



Simplify the User Experience

Workload

+ New
Control
View Output
Delete Directories
Search

<input type="checkbox"/>	ID	Name	Type	State	Application	Submitted↓	Started	Ended	Queue	User	Group
<input type="checkbox"/>	708	JOBA3	Job	<div><div></div>Pending</div>	generic:gener	2017-10-18 16:00	-	-	normal	georgeg	/jgroup1
<input type="checkbox"/>	707	JOBA2	Job	<div><div></div>Pending</div>	generic:gener	2017-10-18 16:00	-	-	normal	georgeg	/jgroup1
<input type="checkbox"/>	706	JOBA1	Job	<div><div></div>Pending</div>	generic:gener	2017-10-18 16:00	-	-	normal	georgeg	/jgroup1
<input type="checkbox"/>	705	JobArray[1-10]	Array	<div><div></div></div>	generic:gener	2017-10-18 16:00	2017-10-18 16:00	-	normal	georgeg	-
<input type="checkbox"/>	704	Job3	Job	<div><div></div>Done</div>	generic:gener	2017-10-18 16:00	2017-10-18 16:00	2017-10-18 16:00	normal	georgeg	-
<input type="checkbox"/>	703	Job2	Job	<div><div></div>Done</div>	generic:gener	2017-10-18 16:00	2017-10-18 16:00	2017-10-18 16:00	normal	georgeg	-
<input type="checkbox"/>	702	Job1	Job	<div><div></div>Done</div>	generic:gener	2017-10-18 16:00	2017-10-18 16:00	2017-10-18 16:00	normal	georgeg	-
<input type="checkbox"/>	701	ABAQUS_On	Job	<div><div></div>Exited</div>	ABAQUS:AB/	2017-10-18 16:00	2017-10-18 16:00	2017-10-18 16:00	normal	georgeg	-
<input type="checkbox"/>	700	tt	Job	<div><div></div>Exited</div>	-	2017-10-17 16:00	2017-10-17 16:00	2017-10-17 16:00	normal	georgeg	-
<input type="checkbox"/>	699	ABAQUS_150	Job	<div><div></div>Done</div>	ABAQUS:AB/	2017-10-16 16:00	2017-10-16 16:00	2017-10-16 16:00	normal	georgeg	-

<<
<
Page 2 of 3
>
>>
10

Viewing 11 - 20 of 27

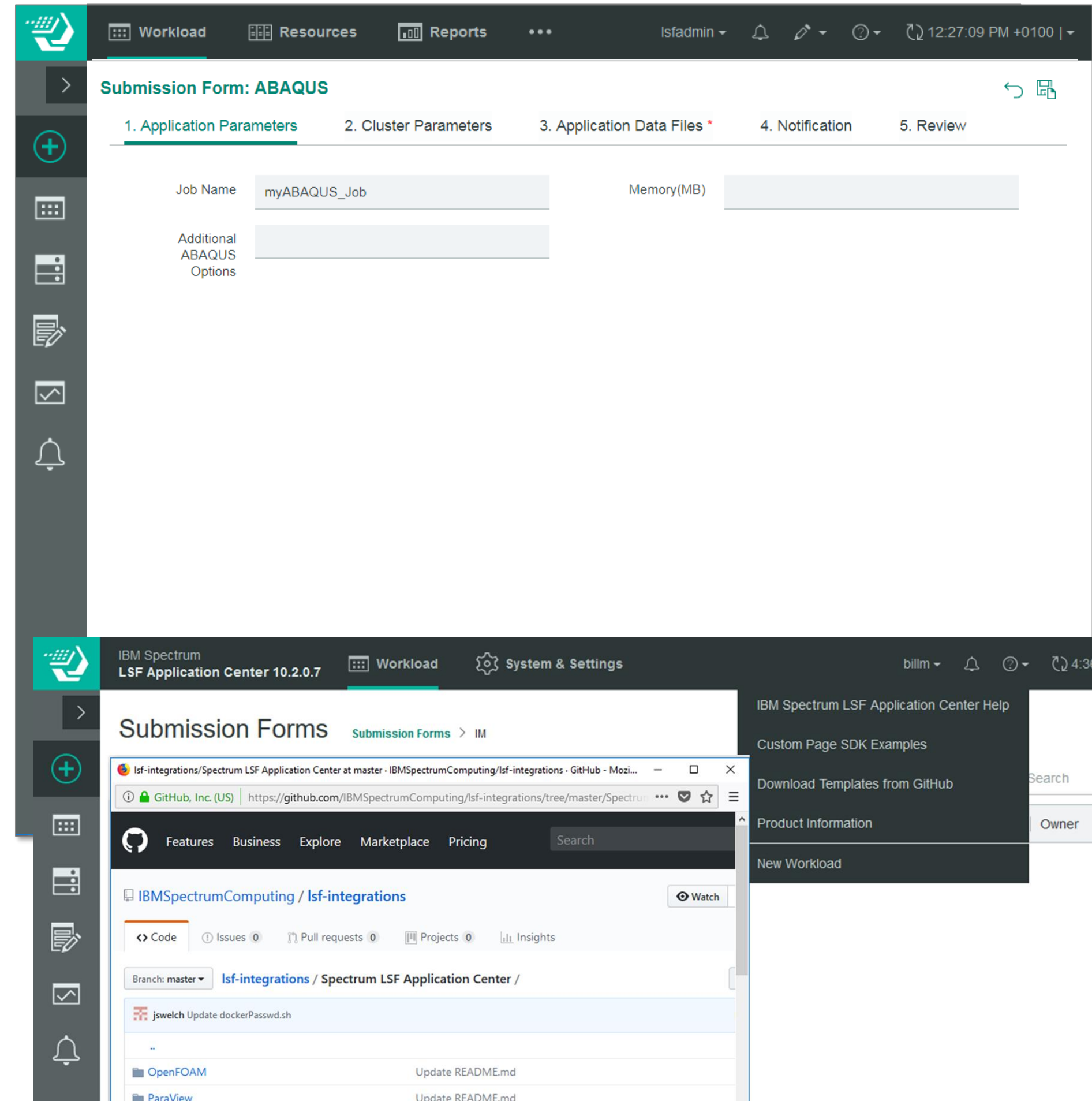
Guided Submission

The administrator or power user can define templates for specific applications.

- Display as a form or wizard (new in SPK8)

Templates:

- Simplify the use of the application exposing specific fields with (optional) predefined options/values.
- Reduces user training requirements, errors and support user
- Makes HPC more accessible – point and click submission
- Templates can handle dependencies and define how input/output files for an application should be viewed/graphed.
- Download additional/updated templates from Github



Example – Containerised Tensorflow + Tensorboard

LSF Suite for HPC 10.2.0.6

Workload

Resources

System & Settings

Reports

student3

10:26:29 AM -0700

Workload

New

Control

View Output

Delete Directories

Search

User = student3

Ended = Past Hour

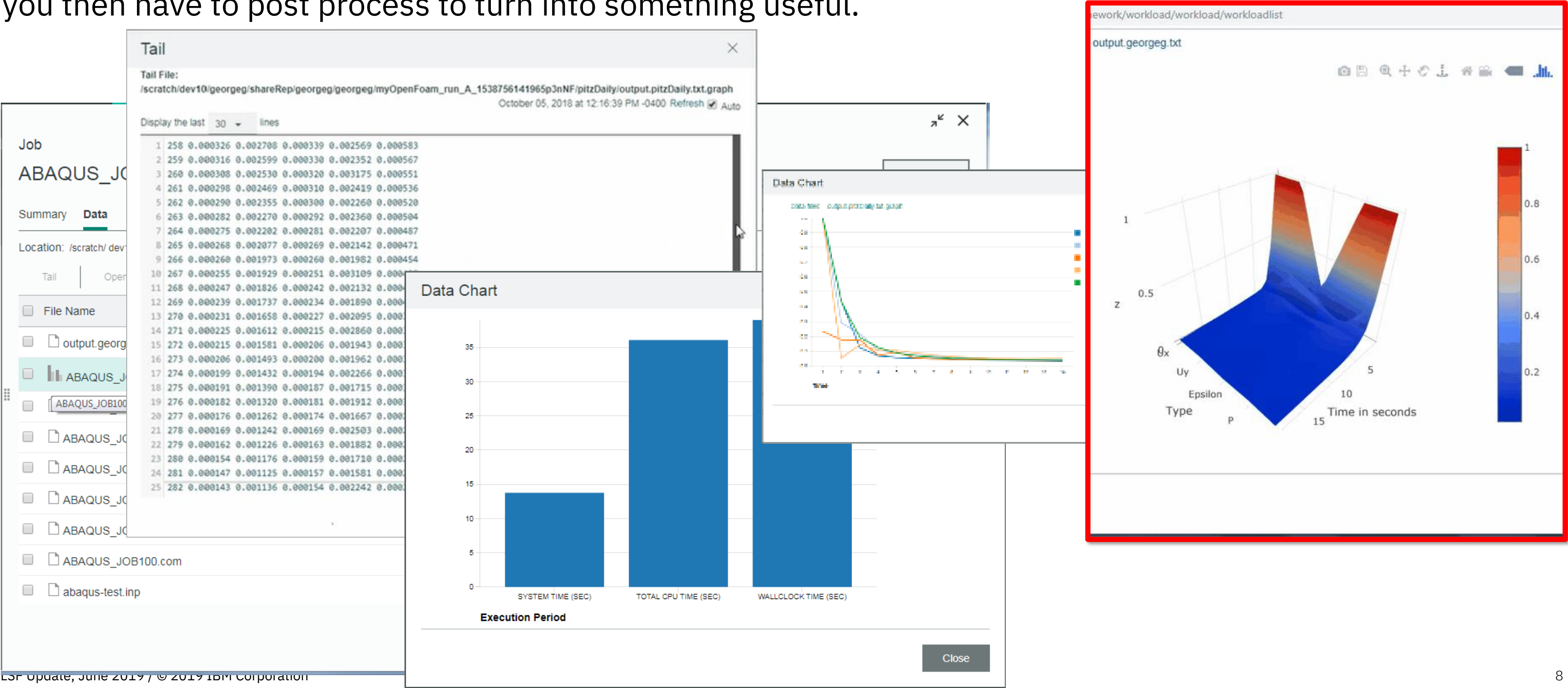
ID	Type	Name	State	Submitted	Started	Ended	User
----	------	------	-------	-----------	---------	-------	------

Page 0 of 0

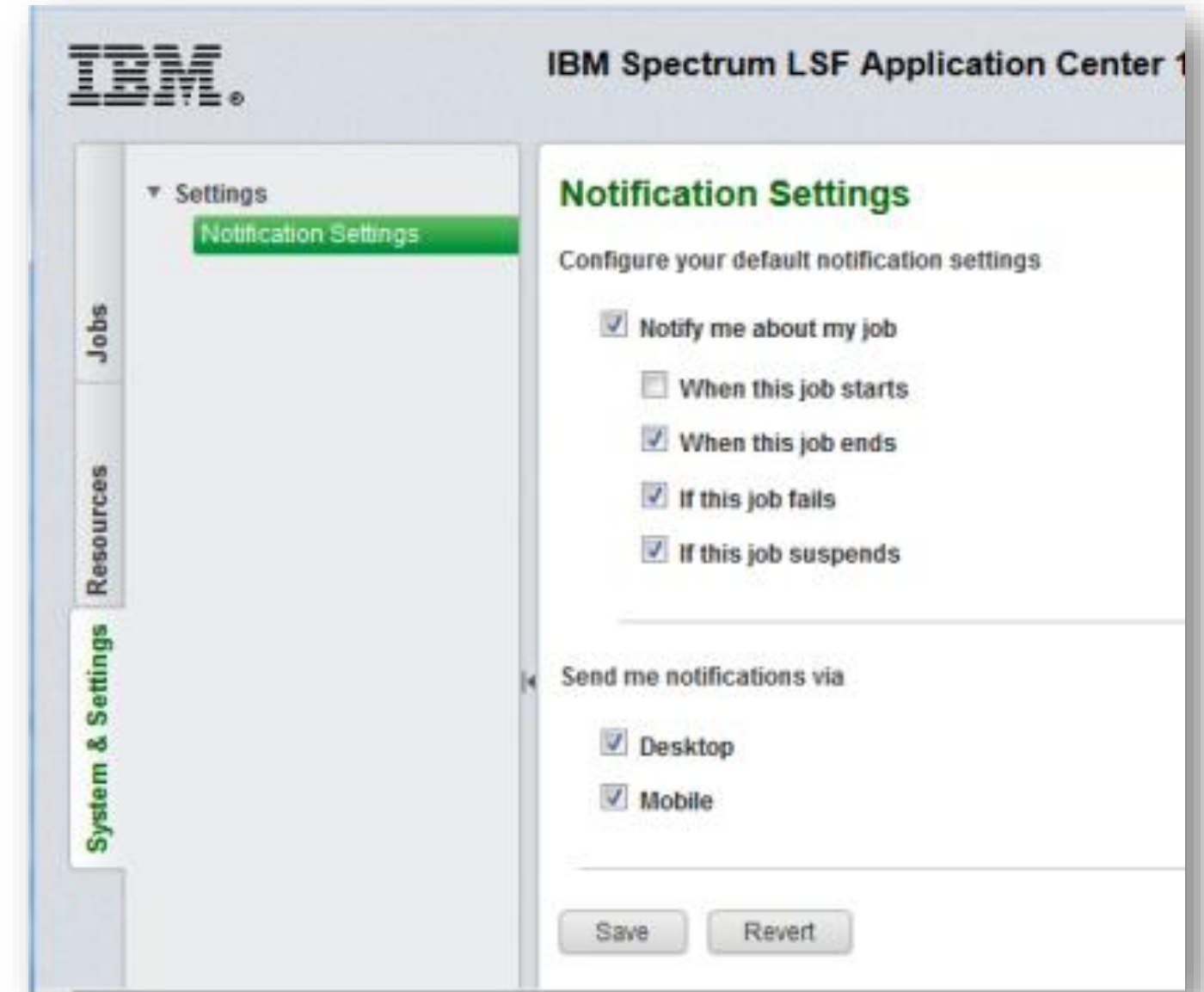
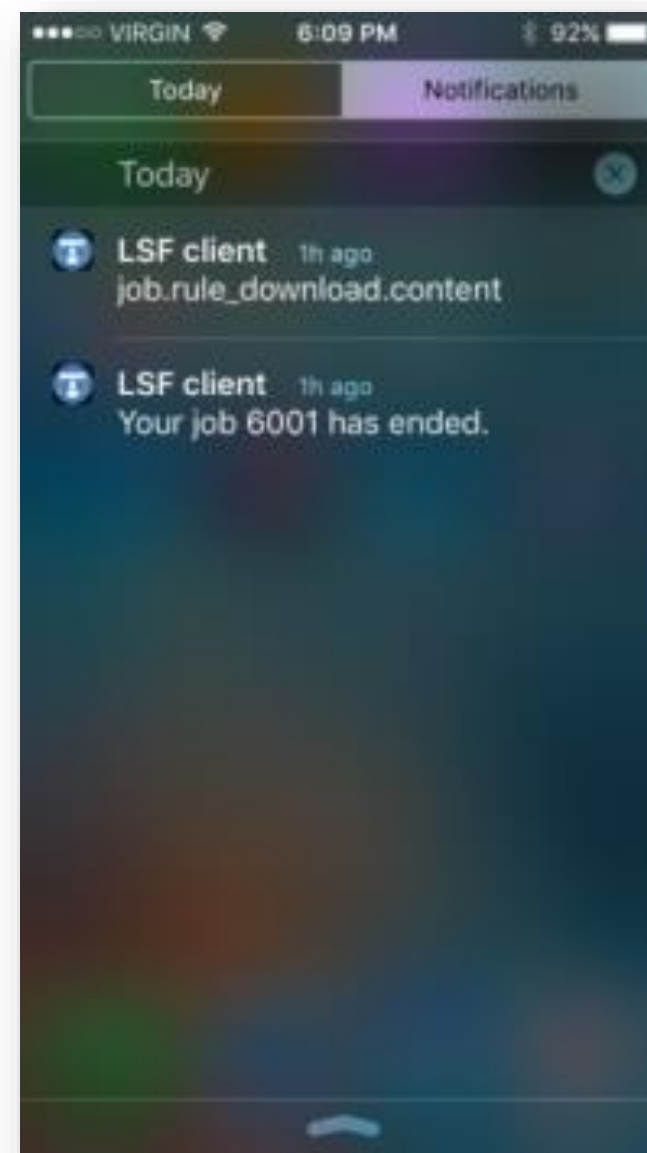
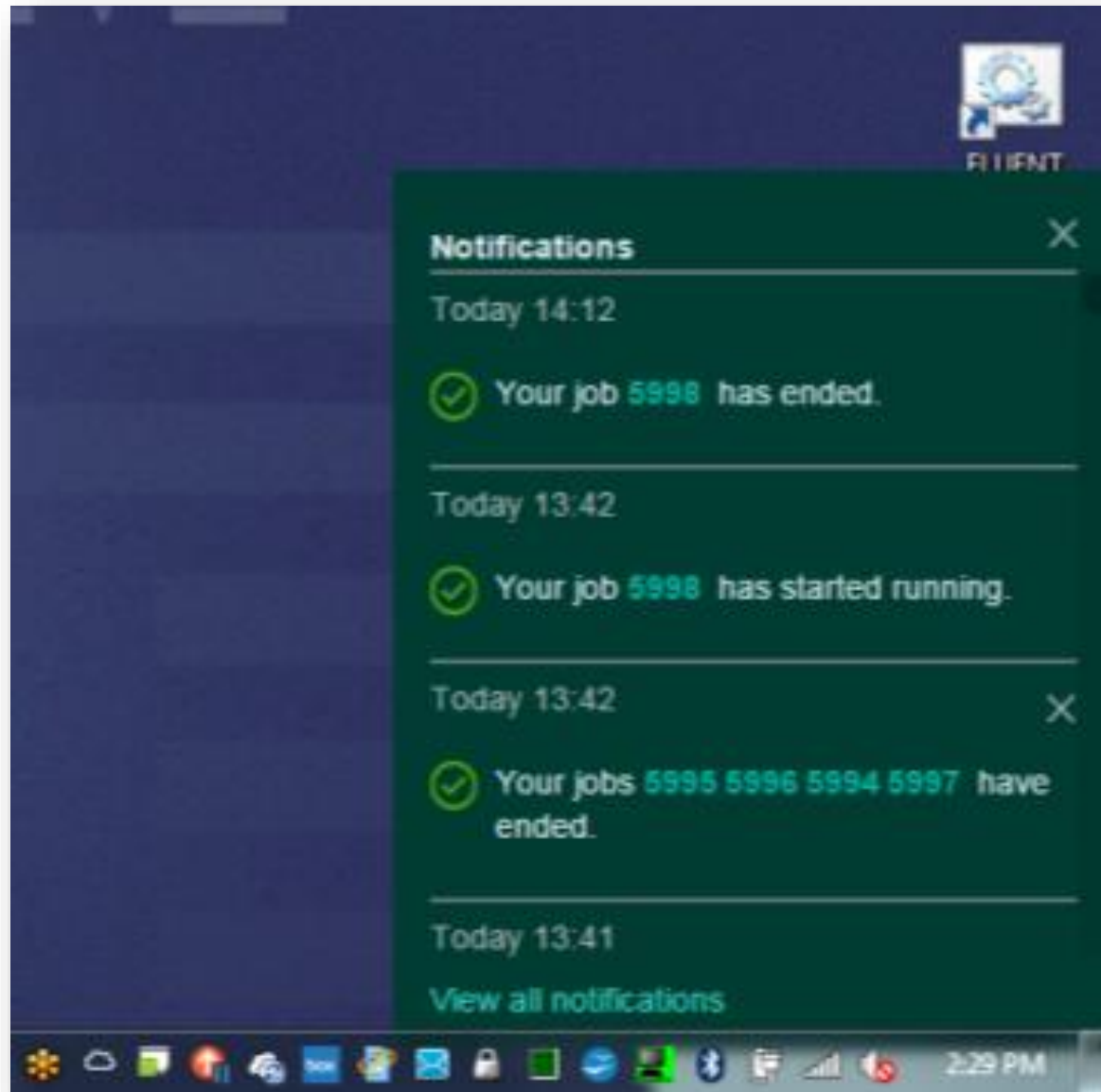
No data to display

Simplifying Viewing of Results

The output of most applications are one or more text files which you then have to post process to turn into something useful.



Inform the user of the progress, completion or failure of their work



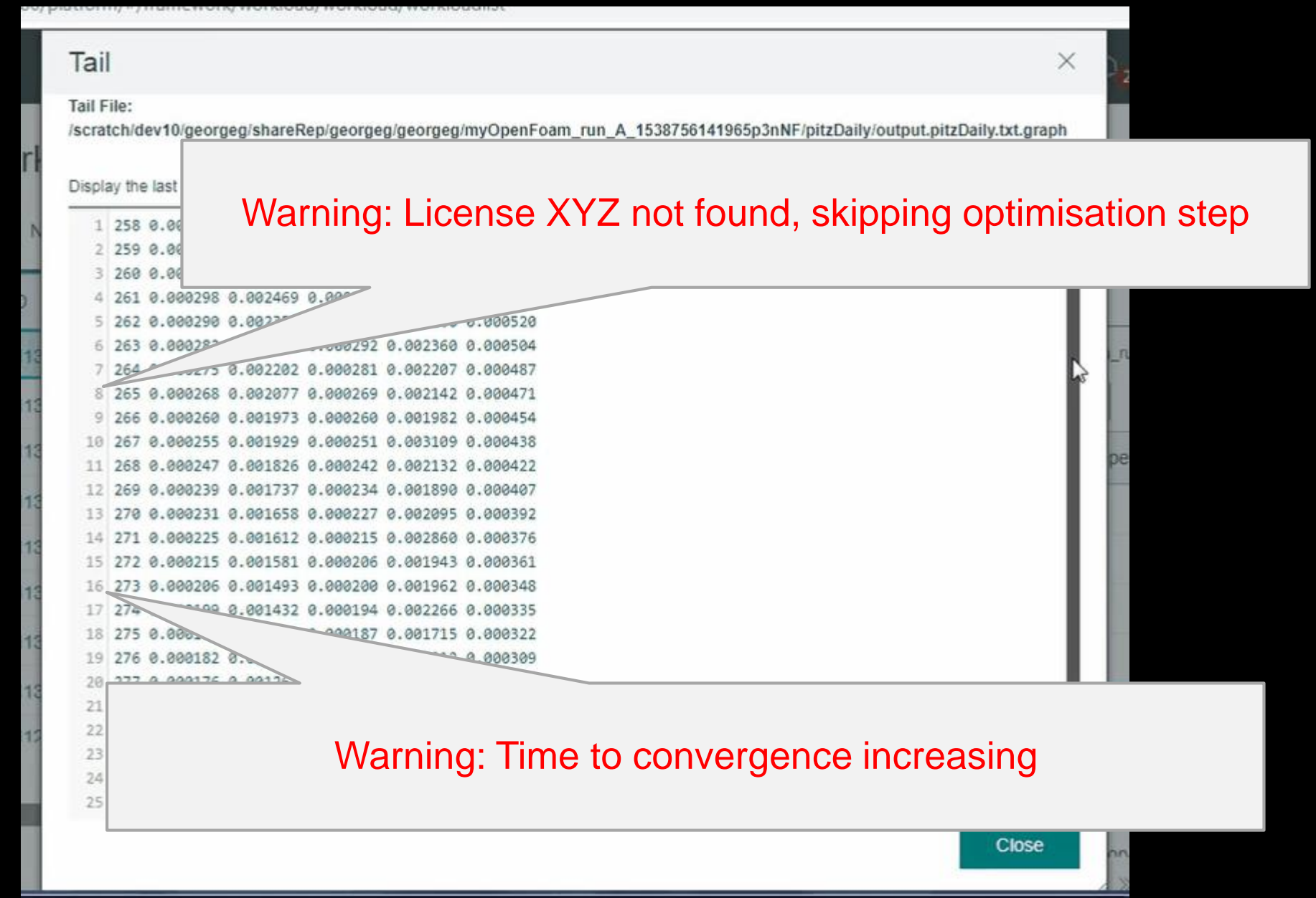
Email, desktop, browser or mobile

But what about jobs that are still running, but are not quite “right”?

Users repeatedly check the output of applications to check that the application is doing the right thing.

There may be something wrong, but it hasn't caused the application to fail – such as a solution not converging.

If the user doesn't notice this, the job may run for many hours and produce no useful output.



New Application Specific Watchdog

WorkloadResourcesSystem & Settings

georgeg

Workload

New

Control

View Output

Delete Directories

ID	Name	State	External Status	Contai
17776	Longjob	Running	High Memory Usage (500M)	-
18079	sleep 600	Exited	-	-
18080	sleep 111	Done	-	-
18081	sleep 500	Exited	Job hang, killed by system.	-
18082	sleep 200	Exited	Job hang, killed by system.	-
18083	with_Watchdog	Exited	High memory useage 106 M	-
18084	sleep 101			

Notifications

5 minutes ago

Job <17776>: High Memory Usage (500M)

6 minutes ago

Job <17776>: High Memory Usage (500M)

Show All Notifications

Search

efficiency

Type

Application

Job	generic:generic
Job	-
Job	-
Job	-
Job	-
Job	ABAQUS:ABAQU
Job	-

Executive Organizer

Job warning for <17776:Longjob>

noreply to: George Gao

From:

noreply@ib22b11.localdomain

To:

George Gao/Ontario/IBM@IBMCA

High Memory Usage (500M)

The following is the workload details:

ID:

17776

Name:

Longjob

Application:

generic:generic

Started:

2019-05-09 11:41:33

Today

Notifications

Today

LSF client 1h ago
job.rule_download.content

LSF client 1h ago
Your job 6001 has ended.

LSF Client
Job 11319
Fatal "Optimisation failed. Aborted"

LSF Client
Job 11321
Warning "Solution time diverging"

IBM & ASTON MARTIN RED BULL RACING



Focus Areas 2019-2021

Core Scheduling

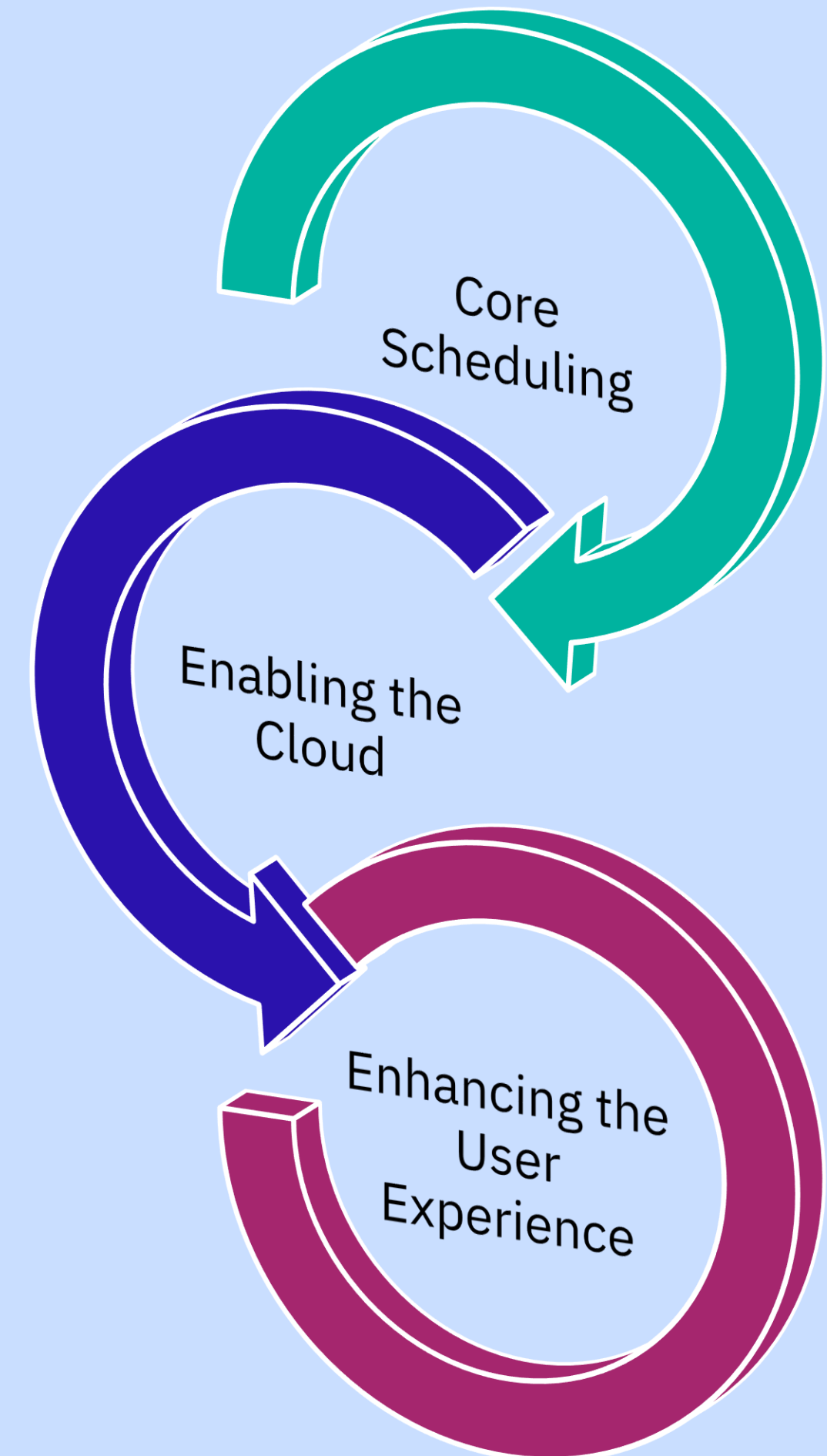
- **Performance & Scalability**
- **Workload Policies**
- **Technology: GPU's & Containers**

Enabling Multi-Cloud

- **When to forward work**
- **Ensuring the right data is available**
- **Intelligent autoscaling the Cloud**

Enhancing the User Experience

- **Simplifying HPC**
- **Computational Workflows**
- **Operational Visibility**



Cutting Edge Performance and Scalability

Continuous Performance Improvement:

LSF 10 delivers ~3.8x improvement in scheduling performance.

Every update contains new performance optimizations – 6 monthly updates.

Scalability:

The largest single cluster today is in excess of 12,000 hosts running large scale parallel simulations.

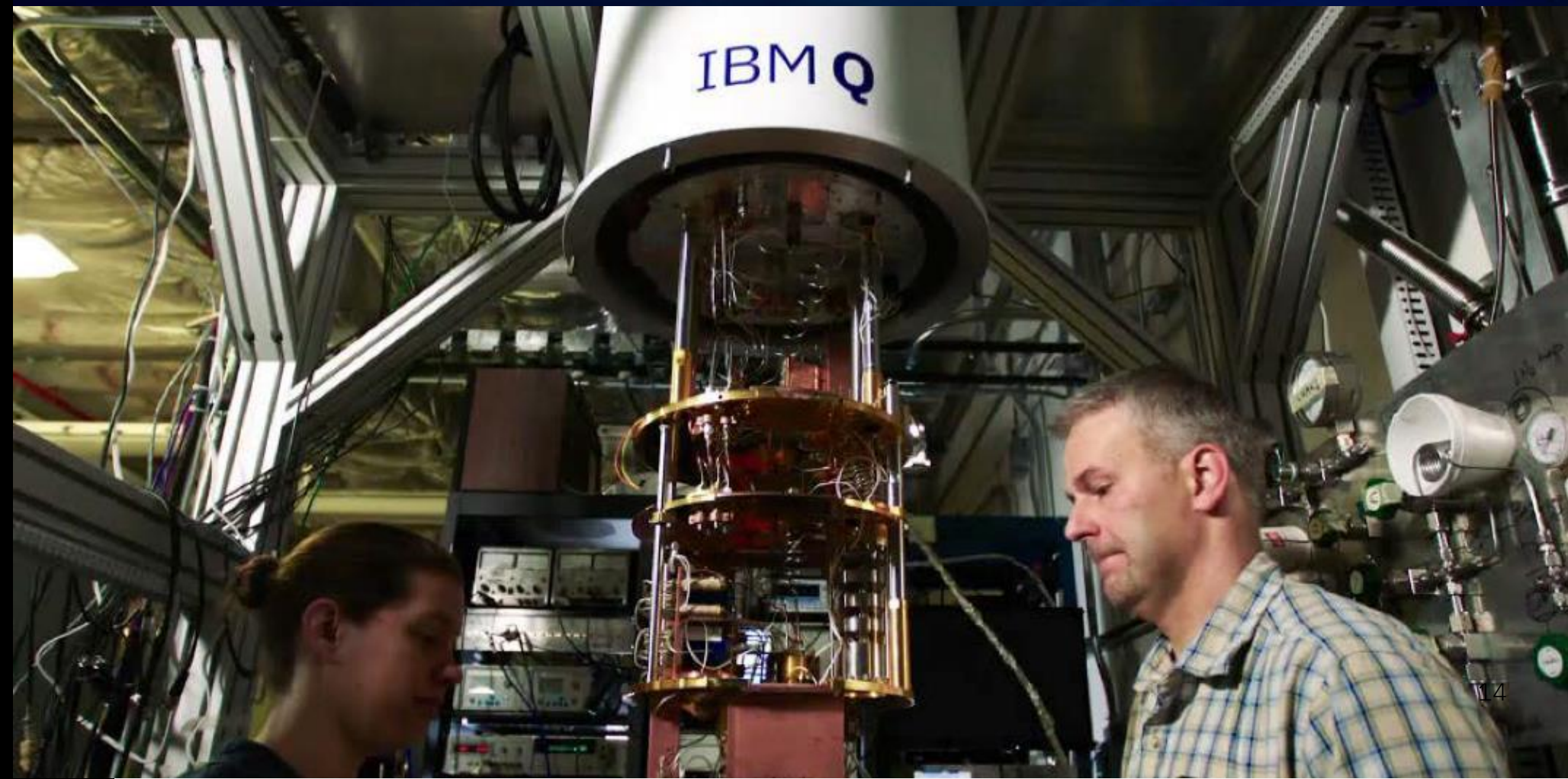
The largest high throughput cluster is around 6,500 hosts, running 7M+ jobs per day, and bursting 50,000+ cores to the cloud

CORAL and IBM Q

Sampling of OpenPOWER Members Contributing to Sierra & Summit



2/3rds of these use LSF

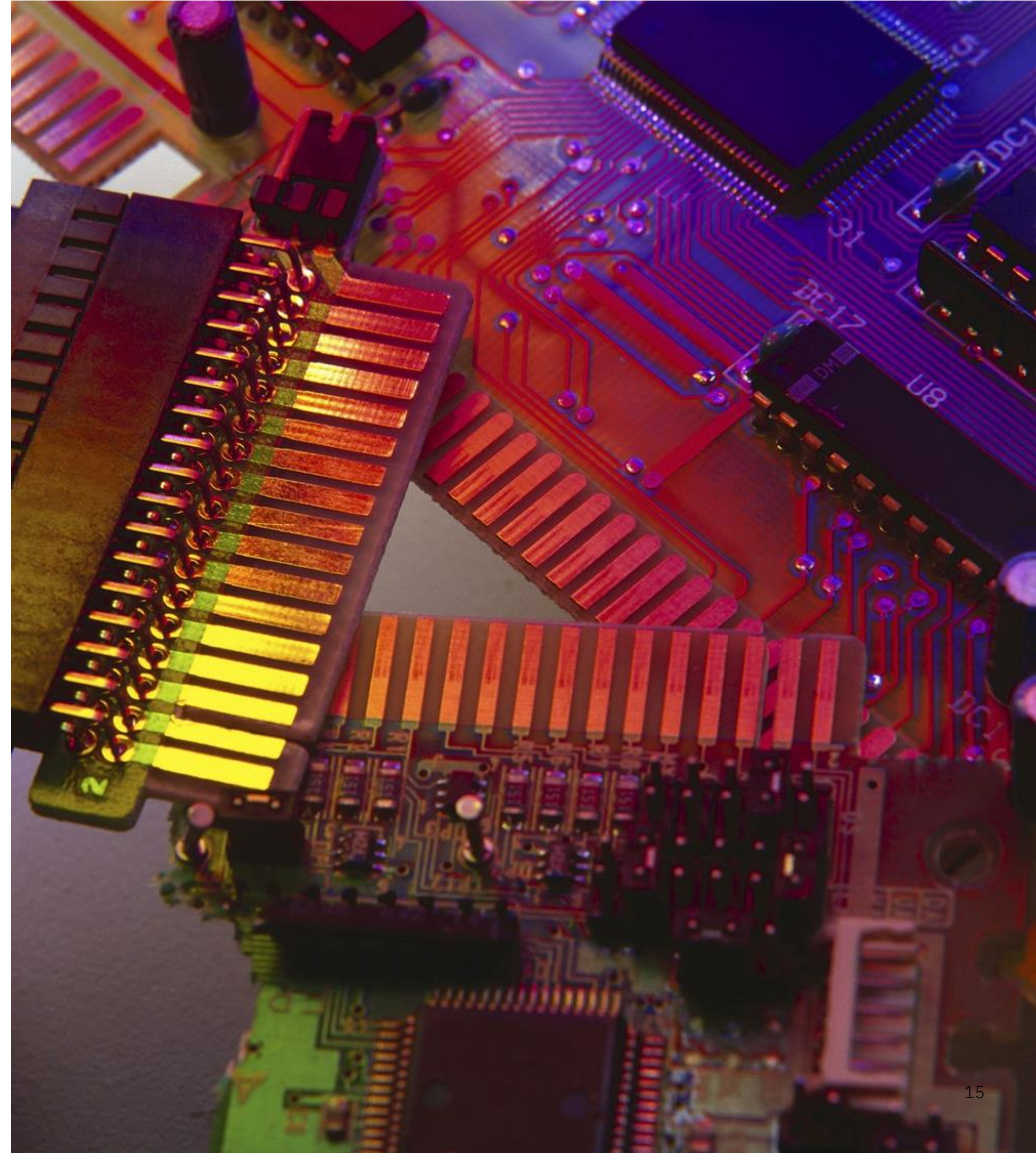


Market Leading GPU Support

Over 10 years of GPU functionality

Recent enhancements:

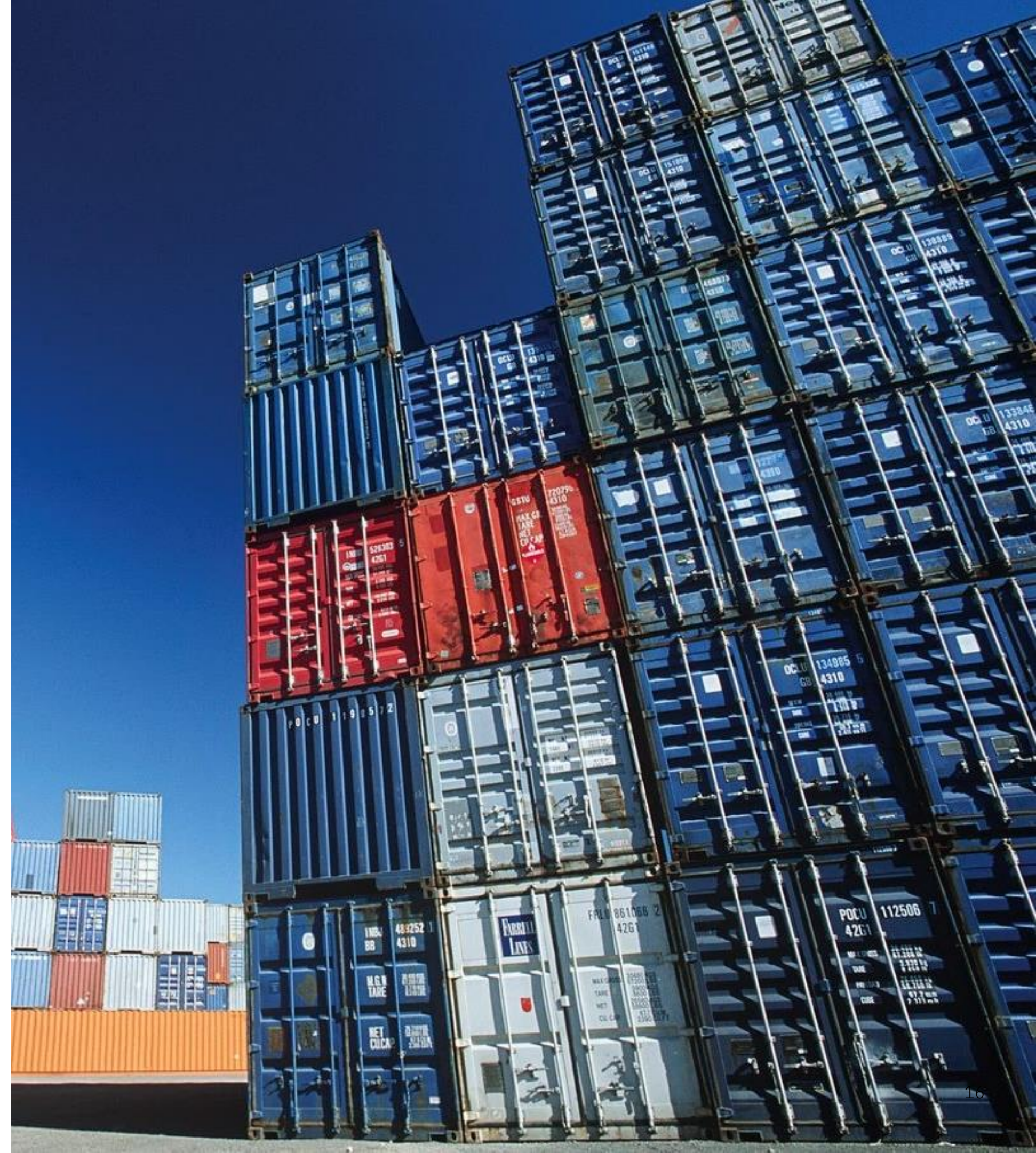
- “Zero Config” - LSF will now automatically detect and configure GPU support. This means that users can take advantage of GPU's as soon as Spectrum LSF installed.
- Simplified (-gpu) syntax
- GPU Fairshare & GPU Pre-emption
- Multi-MPS Support – Multiple MPS daemons per job and/or multiple jobs per MPS daemon (spk8)
- Additional affinity options (spk8)



Running Containerized Workloads with LSF

LSF 10 provides:

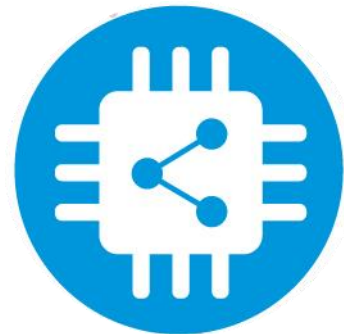
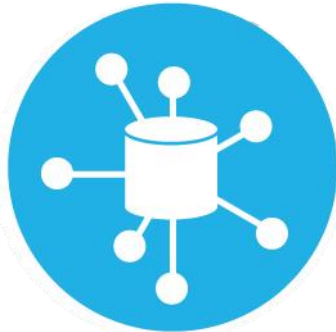
- Integrated support for running containerized applications with Docker, Shifter and Singularity.
- Transparent container access – users don't need to learn complex container syntax.
- All container startup and filesystem mounting is performed by LSF.
- The User never gains elevated privileges.
- Administrator Visibility of container use:
 - host, container name, tags, source repository, file path, size, install time, age, last used, last used by



HPC Administration in the Cloud Era

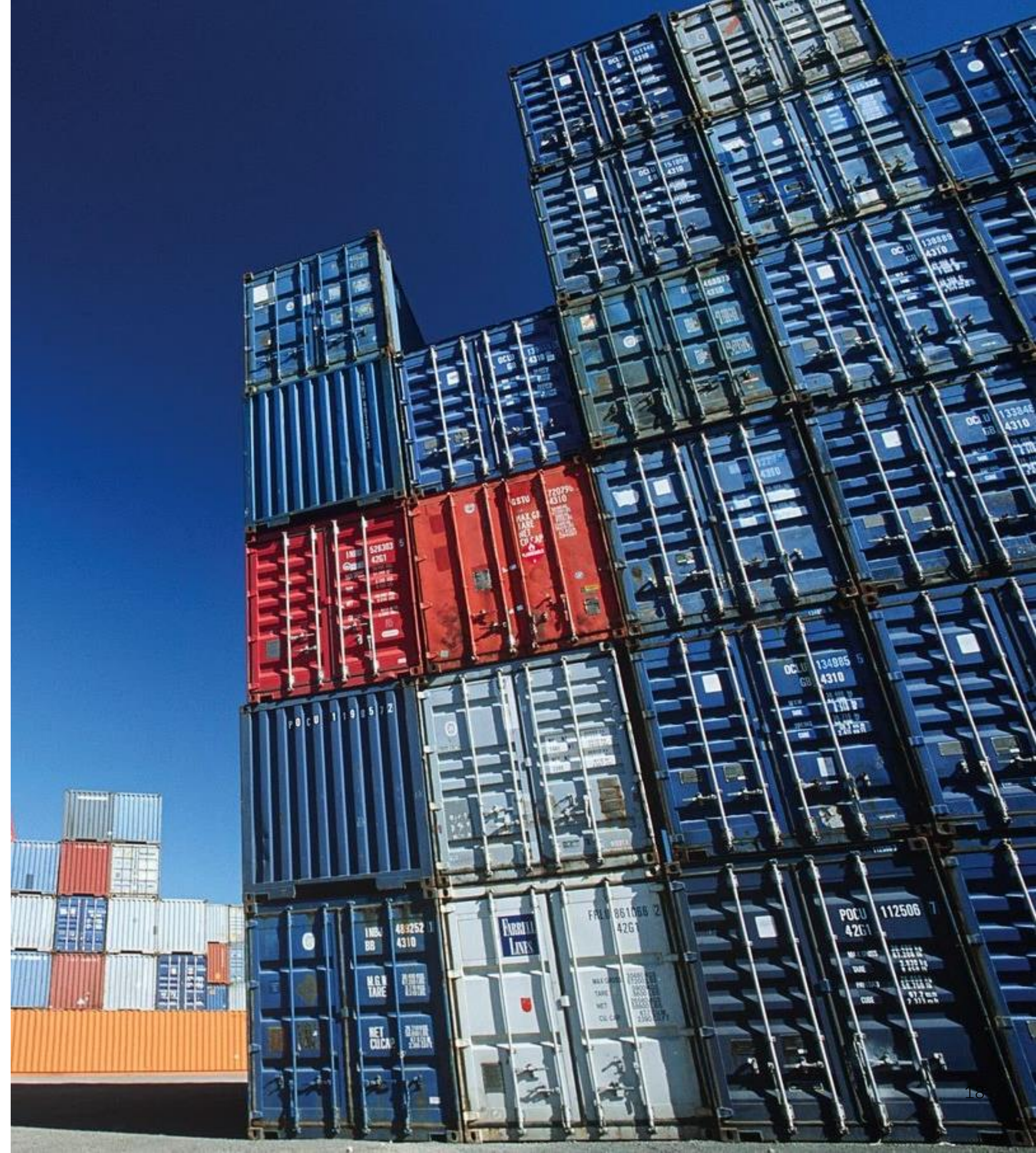
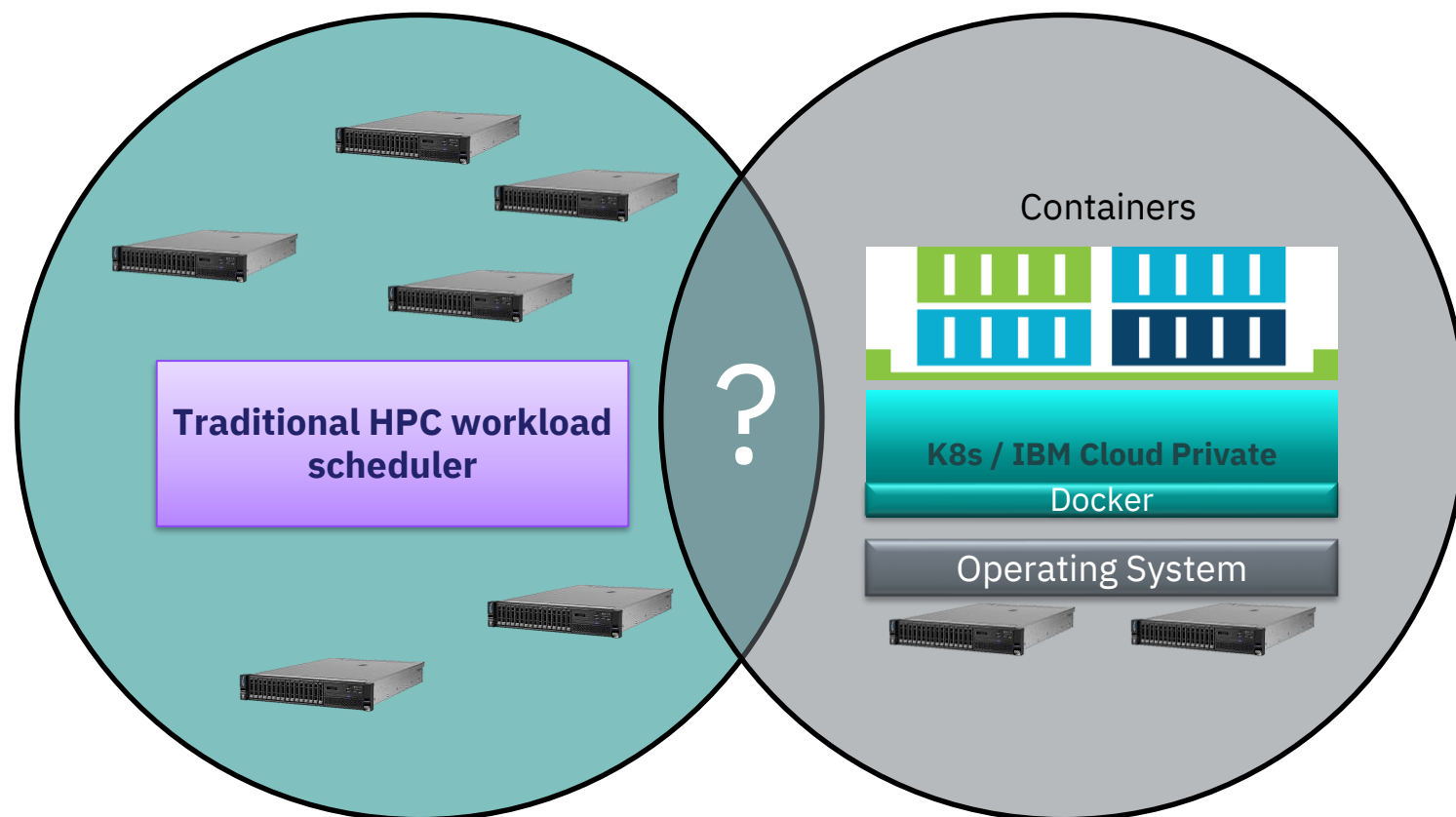
Numerous considerations and requirements to satisfy

- I need to on-board cloud native workloads without disrupting my existing HPC cluster
- I need enable DevOps tools and processes for my users
- I need to provide a secure multi-user environment that doesn't sacrifice performance



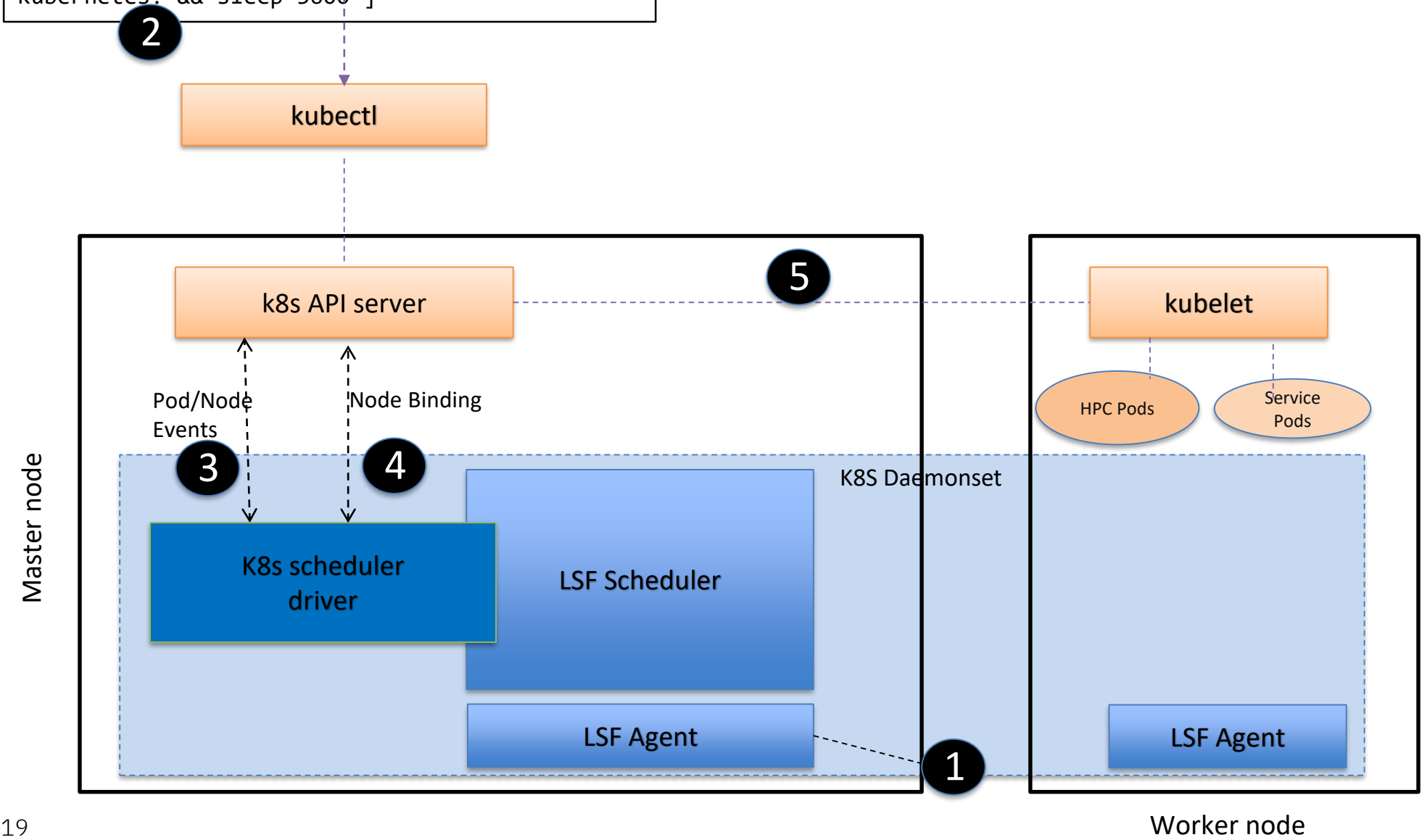
HPC and Kubernetes

- Can k8s address all HPC use cases and negate the need for a traditional workload scheduler?
- While many high performance computing workloads can be containerized, they expect various services to be available in the environment – containerizing these can be challenging.
- How to integrate k8s with a traditional workload scheduler to get the best of both worlds?



Deployment Option #1: LSF as a pod scheduler in a pure K8S environment

```
apiVersion: v1
kind: Pod
metadata:
  name: myapp-pod
  labels:
    app: myapp
  annotations:
    lsf.ibm.com/queue: "night"
    lsf.ibm.com/fairshareGroup: "project-1"
    lsf.ibm.com/gpu: "gpu=4:mode=shared"
spec:
  schedulerName: lsf
  containers:
    - name: myapp-container
      image: busybox
      command: ['sh', '-c', 'echo Hello
Kubernetes! && sleep 3600']
```



K8S native user experience

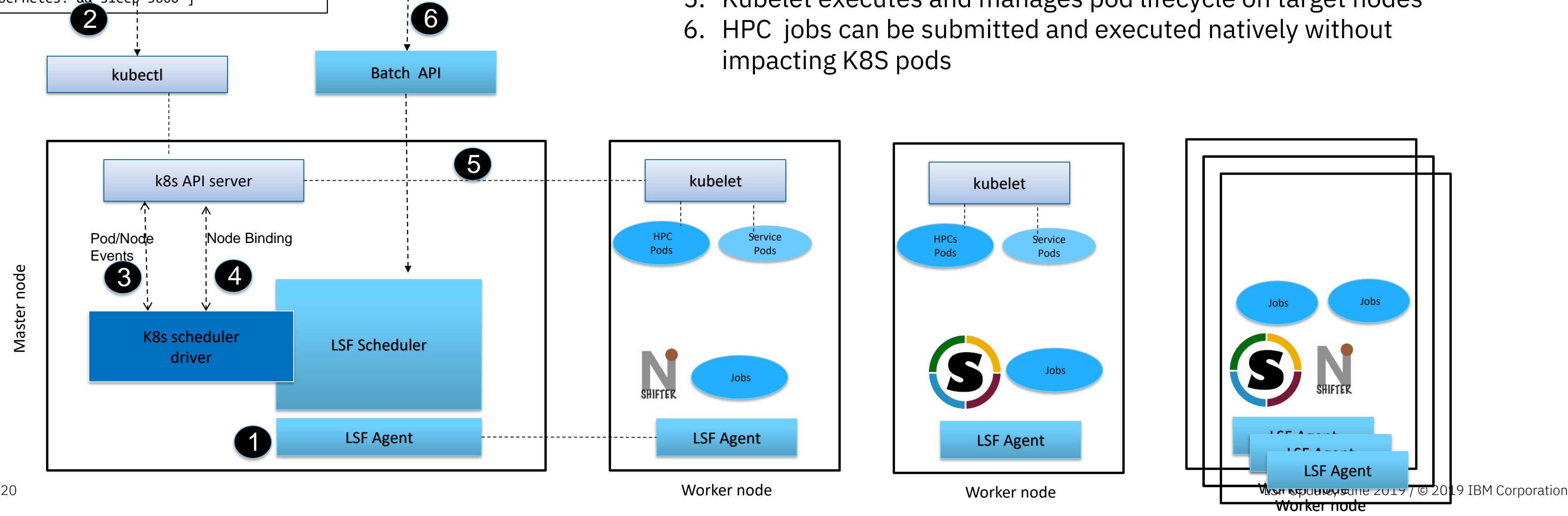
1. HPC Scheduler components are deployed as DaemonSet in K8S via Helm chart
2. Users submit workload into K8S API annotating pods with scheduler directives
3. Driver listens to API servers and translates pod requests into jobs in HPC Scheduler
4. HPC Scheduler makes decisions to bind pod to specific node based on policy
5. Kubelet executes and manages pod lifecycle on target nodes

Deployment Option #2: Existing Spectrum LSF augmented with Kubernetes

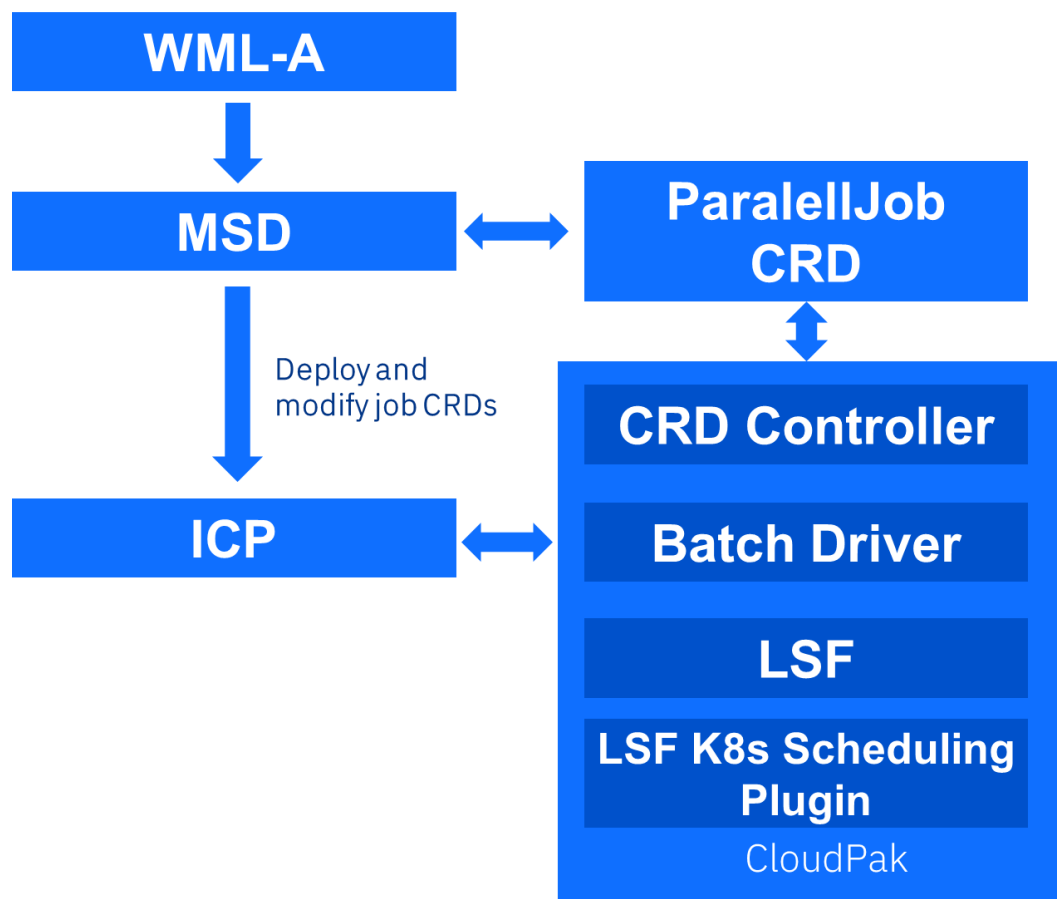
```
apiVersion: v1
kind: Pod
metadata:
  name: myapp-pod
  labels:
    app: myapp
  annotations:
    lsf.queue: "night"
    lsf.fairshareGroup: "project-1"
    lsf.ibm.com/gpu: "gpu=4:mode=shared"
spec:
  schedulerName: lsf
  containers:
    - name: myapp-container
      image: busybox
      command: ['sh', '-c', 'echo Hello
Kubernetes! && sleep 3600']
```

```
bsub -app singularity -n 5,10
mpijob.sh
```

- 1. LSF runs on bare OS, and connects to a k8s API server. LSF cluster can be larger than K8S nodes (>5000 nodes)
- 2. Users submit workload into K8S API annotating pods with scheduler directives
- 3. Driver listens to API servers and translates pod requests into jobs in LSF Scheduler
- 4. LSF Scheduler makes decisions to bind pod to specific node based on policy
- 5. Kubelet executes and manages pod lifecycle on target nodes
- 6. HPC jobs can be submitted and executed natively without impacting K8S pods



Watson Machine Learning Accelerator Use Case



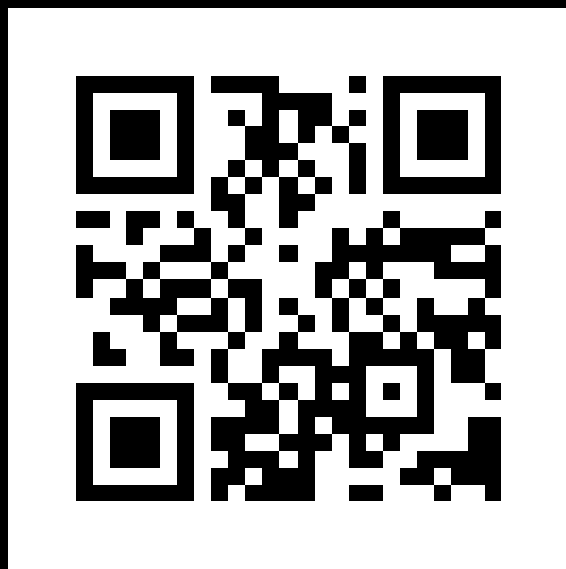
```
apiVersion: scheduling.batch/v1alpha1
kind: ParallelJob
metadata:
  name: large_model_train
  namespace: default
  annotations:
    lsf.ibm.com/queue: myQueue
spec:
  name: large_model_train #same with metadata/name
  description: This is a parallel job to run model training across hosts.
  priority: 100
  headerTask: group0
  placement:
    sameTerm: Zone | Rack | Host
  taskGroups:
    - metadata:
        name: parameter_server
      spec:
        replica: 1
        template:
          spec:
            containers:
              image: ubuntu
            resources:
              request:
                cpu: 1
                memory: 4096Mi
```

LSF specific annotations

Job level specification

Task level specification

```
- metadata:
  name: worker_nodes
  annotations:
    lsf.ibm.com/gpu: nvlink=yes
spec:
  placement:
    spanTerm:
      - topologyKey: node
      taskTile: 2
      - topologyKey: node
      taskTile: 4
  replica: 4
  template:
    spec:
      containers:
        image: nvida/cuda9.2
        name: task1
      resources:
        request:
          nvidia.com/gpu: 2
          memory: 16000Mi
```



IBM
Spectrum
LSF

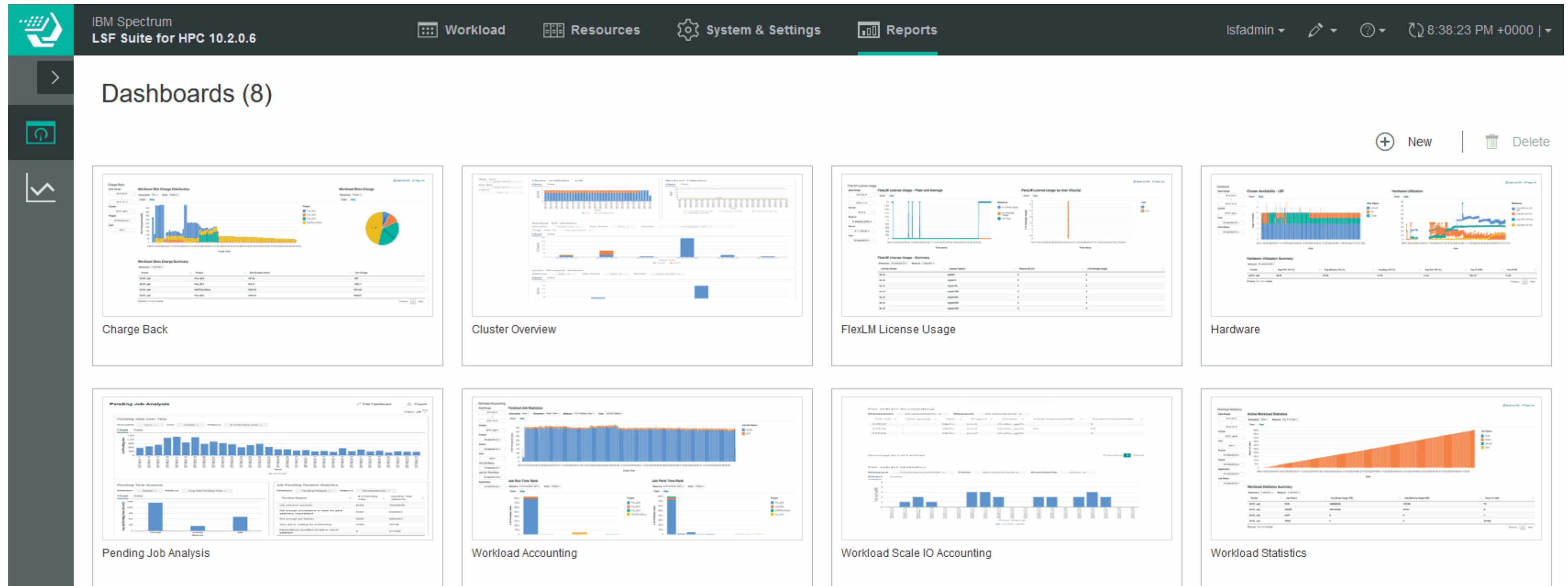
http://ibm.biz/SpectrumComputing_CloudPak_preview

IBM



IBM
Spectrum
Computing

Management visibility LSF Explorer & LSF RTM



Explorer enhancements for SPK8/SPK9

- Update to ElasticSearch 6.x
- Performance Optimizations
- Data consistency between LSF CLI and ES

RTM 10.2 enhancements:

- New responsive user interface with an updated look and feel that is consistent with the rest of the LSF family GUI's.
- Simplified installation
- Enhanced performance with Cacti 1.2