



Federated data services with abundant compute resources

Thomas C. Schulthess

Collaborators

@CSCS: Sadaf Alam, Lucas Benedicic, Miguel Gila, Cristian Mezzanotte, Colin McMurtrie, Marcel Schöngens, Maxime Montinasso, + team

@JSC: Dirk Pleiter, Thomas Lippert, Anna Lührs, Boris Orth + team



@CINECA: Giovanni Erbacci, Roberto Mucci + team



@CRAY: Jacob Balma, Jef Dawson, Mark Stavely, Rakhi Anand

@Microsoft: Chris Basoglu

@NERSC, UC Berkeley: Shane Canon, Doug Jacobsen, Fernando Perez

@EPFL: Nicola Marzari, Giovanni Pizz + team; Jeff Muller + team

@Forschungszentrum Jülich: Katrin Amunz, Timo Dickscheid, Sonja Grün, Alex Peyser + team

Deep Learning toolkits on Cray XC system @ CSCS

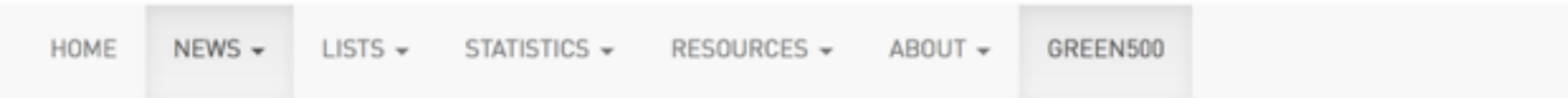
DL Toolkit	C++ & GPU backend	MPI	Working on Cray XC	Fully supported @ CSCS
TensorFlow	yes	no	yes	yes
TensorFlow+MPI	yes	yes	yes	in progress
MXNet	yes	no, ext. to use MPI	yes	in progress
Caffe-MPI	yes	yes	yes	in progress
CNTK	yes	yes	yes	in progress
Spark	no (Java + ext. to use GPUs)	no	yes	in progress

- Installing a DL toolkit on Cray XC is similar to installing any HPC application
 - few extra libraries are needed to satisfy dependencies
- Staging a toolkit can be done with SLURM (our resource manager at CSCS)
 - some toolkits (like Spark) require SSH to be available on compute nodes

DL toolkits on CSCS's Cray XC systems

- Many toolkits use familiar HPC technology
 - developed in C++
 - CUDA/GPU-aware libraries
- HPC filesystems like Luster provide high throughput and data resilience
- MPI enables access to high performance network such as Cray Aries
 - low latency, high bandwidth
- Running a toolkit at scale requires a large amount of GPU accelerated nodes
 - HPC platforms are providing this capacity
- At scale the toolkits face common issues of HPC applications:
 - network and synchronisation boundedness
 - I/O and data access issues

HPC systems are very well suited platforms for running Deep Learning workloads at scale



Home / News / Cray Takes on Fourth Paradigm

Cray Takes on Fourth Paradigm

Michael Feldman | December 19, 2016 18:09 CET



XC50 supercomputer plus Microsoft's Cognitive Toolkit was used to scale up training

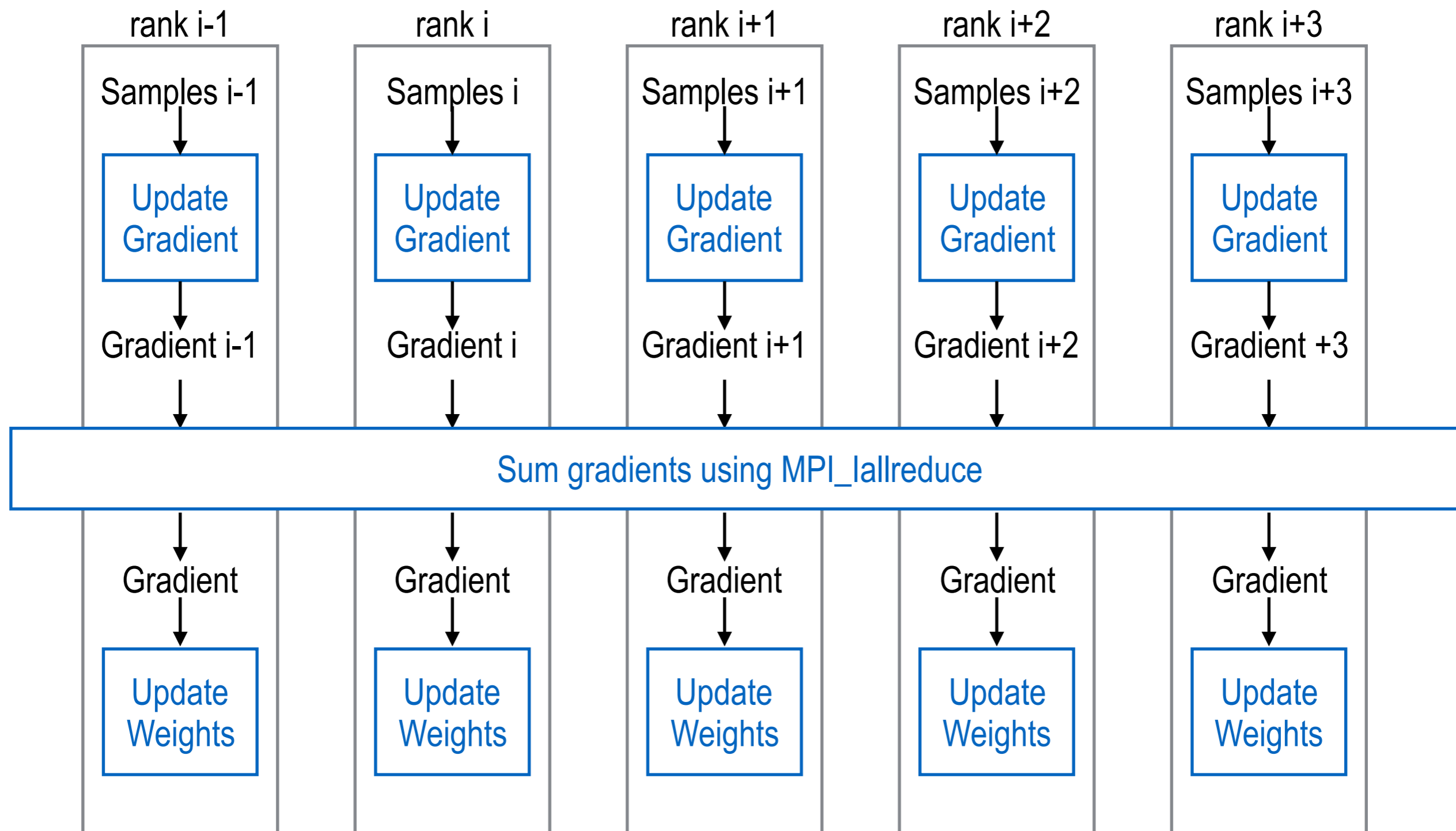
Cray's recent announcement about how its **XC50 supercomputer plus Microsoft's Cognitive Toolkit was used to scale up training of a neural network** serves as a proof point on how topflight HPC technologies can be used to push the boundaries of deep learning. But the company's long game is to bring supercomputing into the realm of deep learning and the broader category of data analytics in a more generalized fashion.

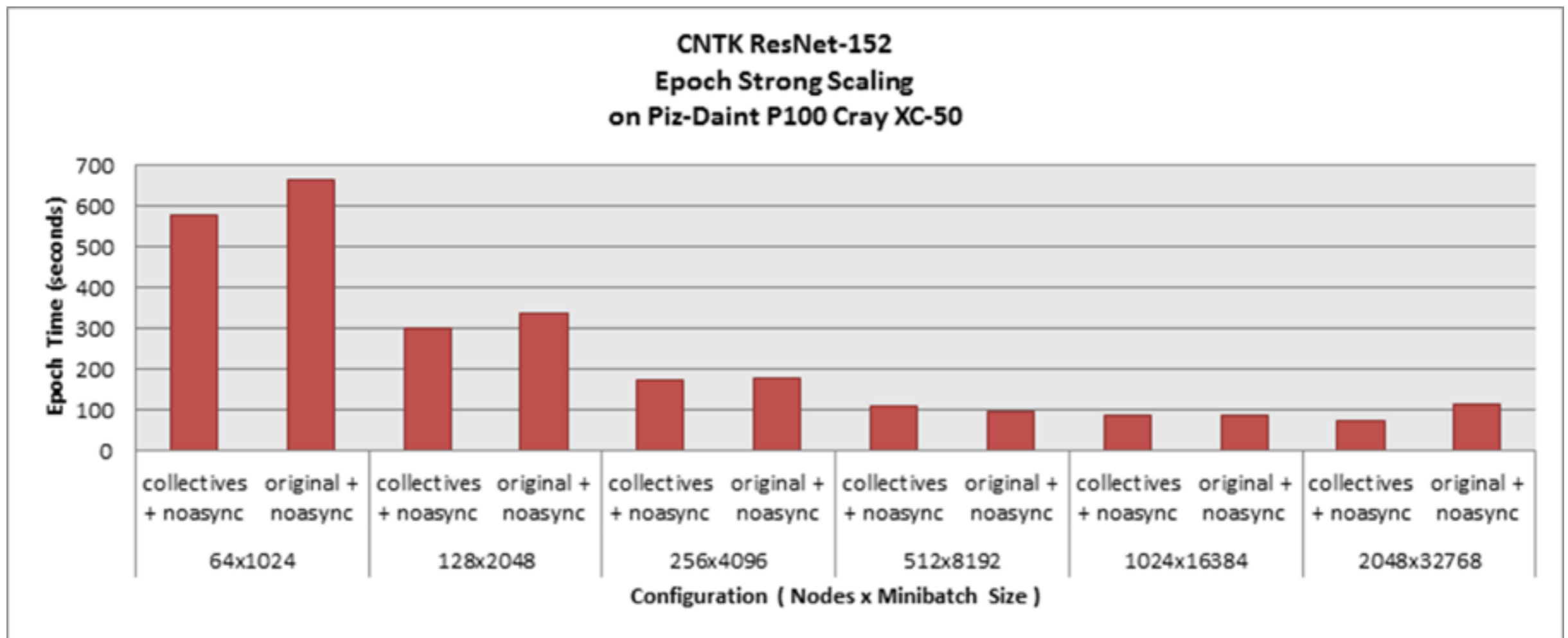
The XC50 in question is "Piz Daint," the 9.8-petaflop system housed at the Swiss National Supercomputing Centre (CSCS), which was recently upgraded with 4,500 of NVIDIA's new P100 GPUs. Each of those GPUs is able to deliver over 5.3 teraflops of double precision floating point performance, 10.6 teraflops of single precision performance, or 21.2 teraflops of half precision performance. Those lower precision flops make these devices especially useful for deep learning codes, where 64-bit float pointing operations are, for the most part, overkill.

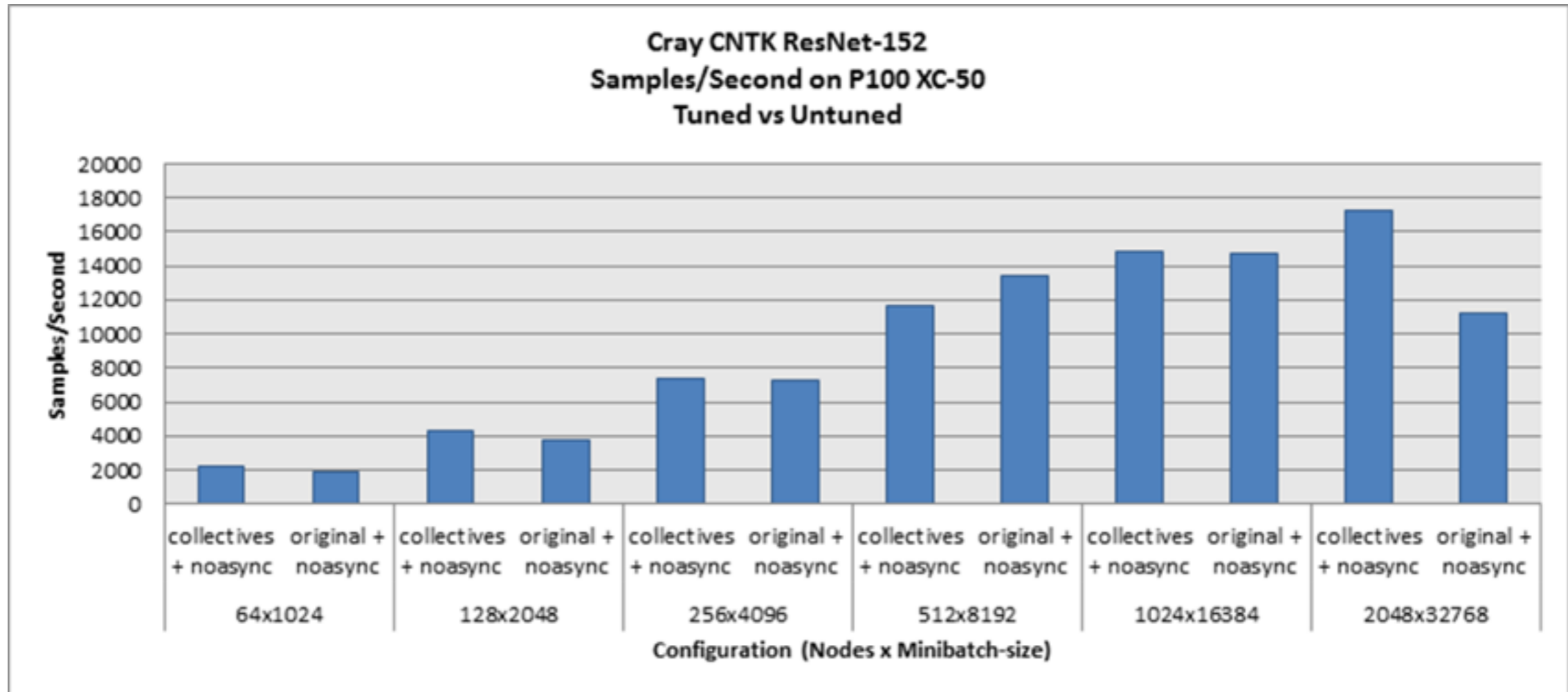
Despite all this concentrated performance, training neural networks is so computationally intense that many GPUs are often needed to in train these models in a reasonable amount of time – minutes or hours, rather than weeks or months. The problem is that with conventional compute clusters and first-generation deep learning frameworks, it's hard to scale these applications without running into some rather daunting bottlenecks.



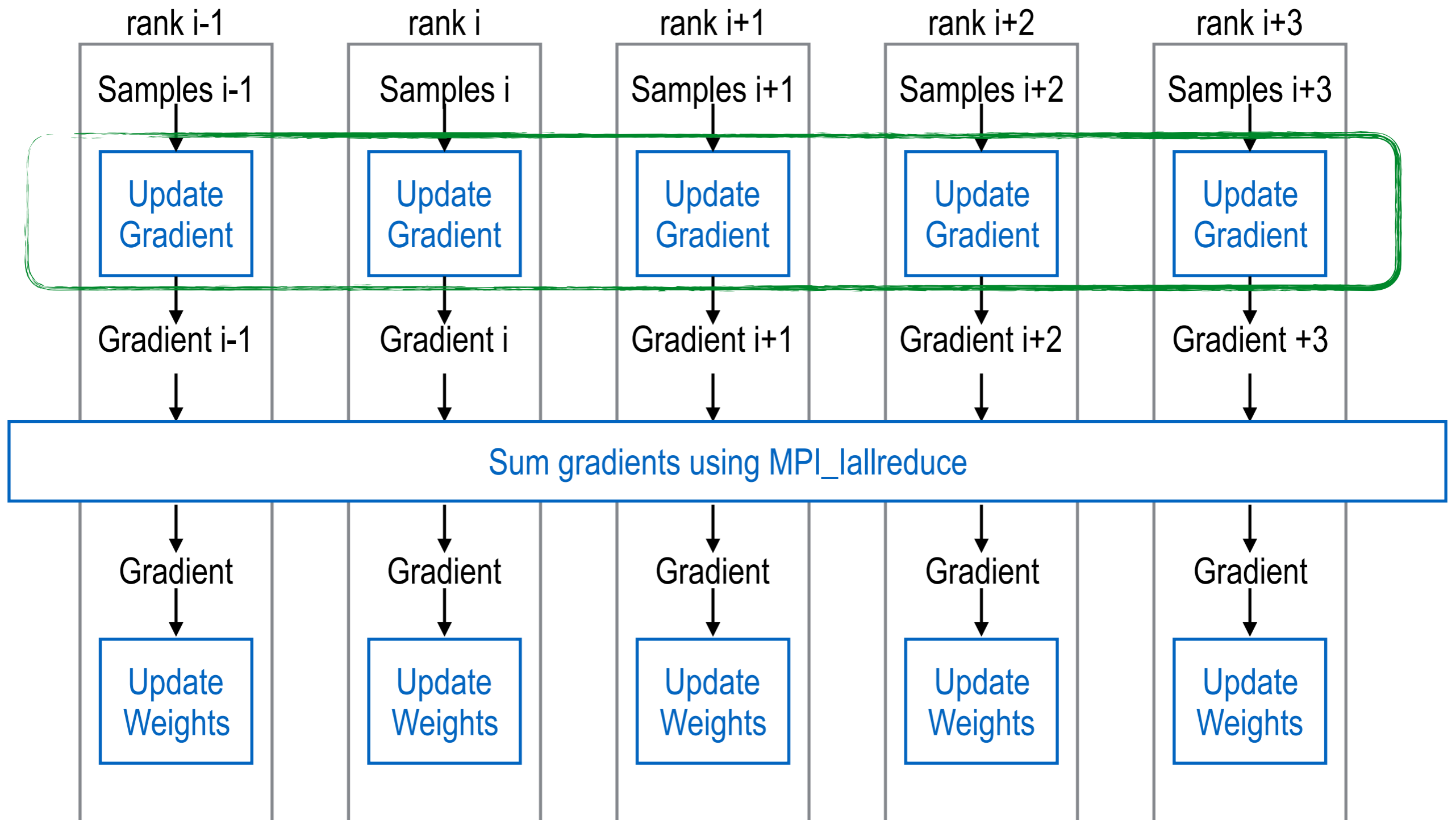
Scaling CNTK with MPI







Scaling CNTK with MPI



“We develop algorithms, we don’t have time to deal with C/C++ or MPI”

–a famous computer science colleague working in machine learning

... echoed by many scientists working with data

Interactive Notebook

Import TensorFlow and start an interactive session

```
In [1]: import tensorflow as tf  
sess = tf.InteractiveSession()
```

Build a computation graph

```
In [2]: matrix = tf.constant([[1., 2.]])  
negMatrix = tf.neg(matrix)
```

Evaluate the graph

```
In [3]: result = negMatrix.eval()  
print(result)  
[[-1. -2.]]
```

Nishant Shukla (2017)

Using Jupyter to get our heads around interactive supercomputing

- Jupyter allow users to
 - integrate development, execution of computation, pre- and post processing with visualisation into one “workflow”;
 - share these workflows in a team; and
 - document their work
- Need to provide an environment where web-based front-end of the notebook is separated from the computation backend
 - run some of the computation on a supercomputer
- Interactive simply means sub-second response time
 - this requires properly organising data in memory/storage sub-system

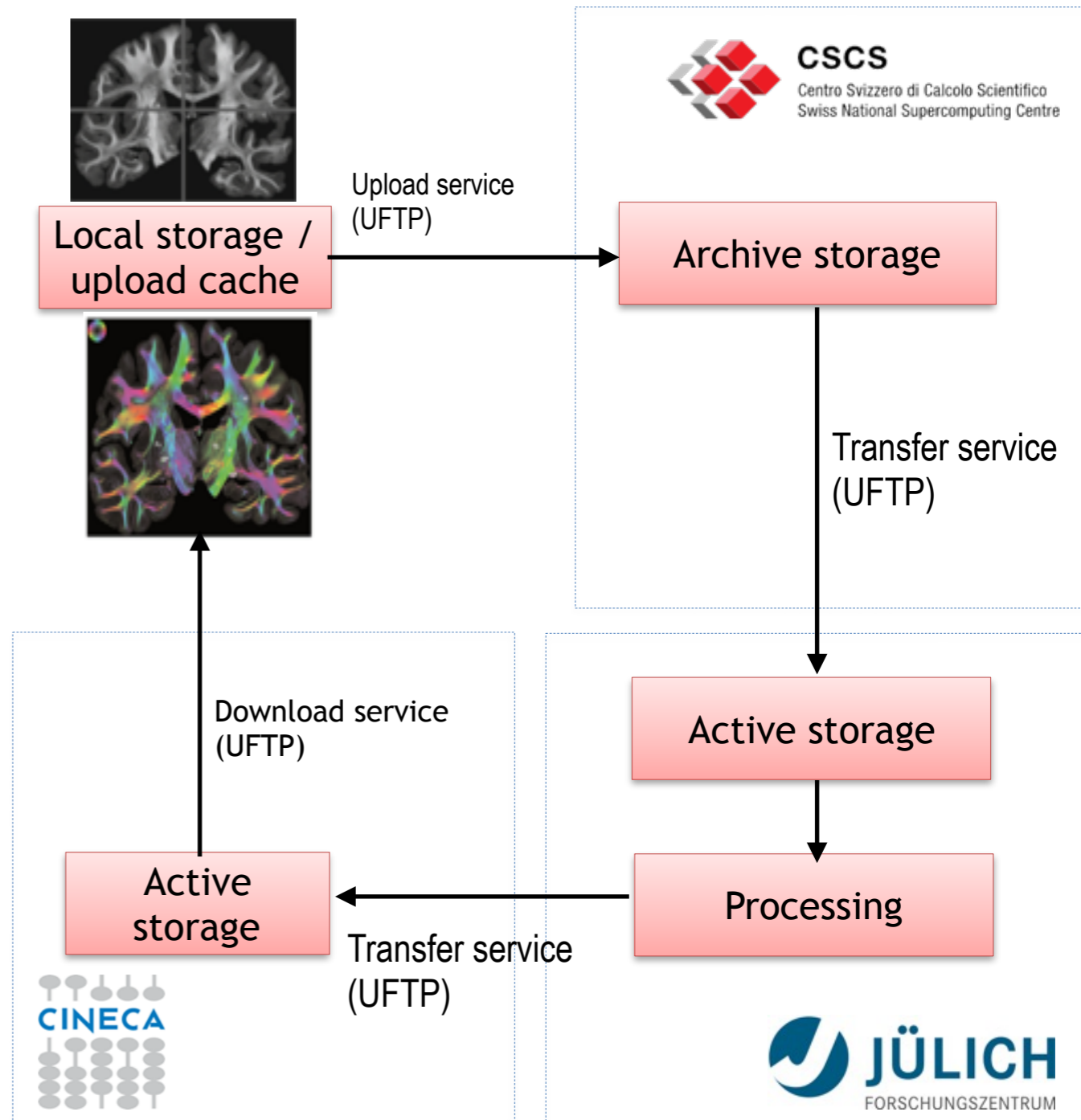
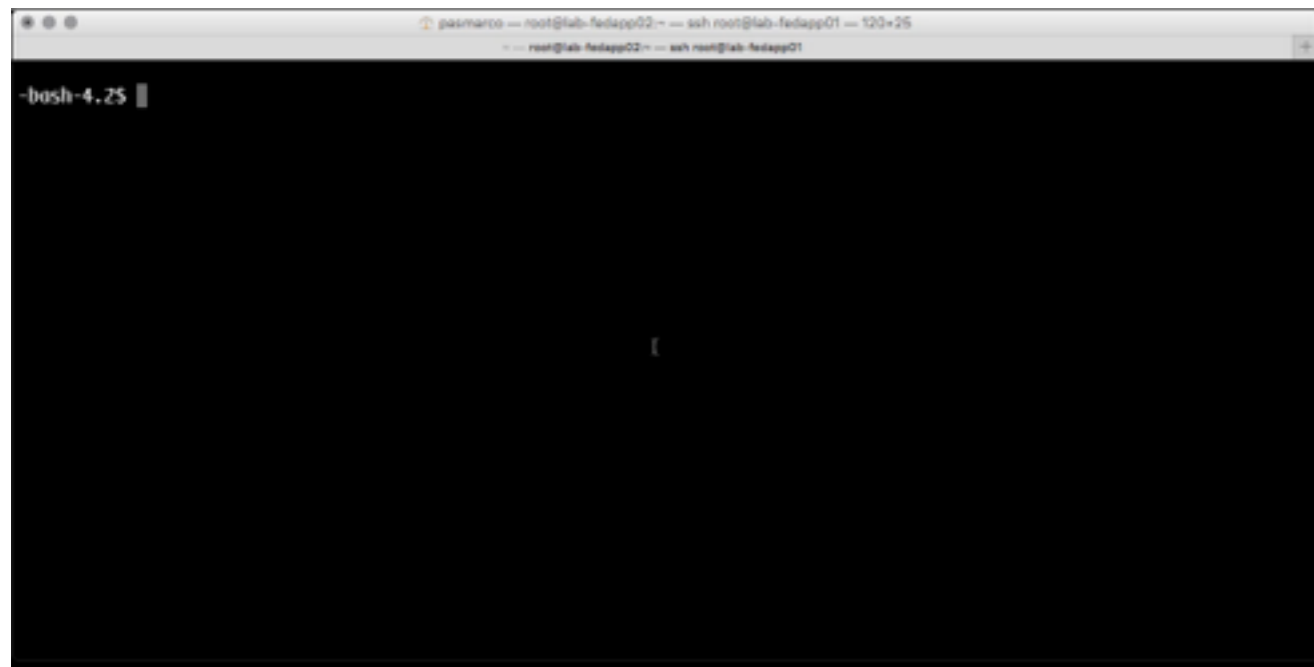
**... and the service should be up most
of the time (like 99+ %)**

Federated Data Pilot Project



www.worldatlas.com

Federated data infrastructure – proof of concept



A few more practical things ...

- Scientists, not system administrators will want to manage teams and access to data
 - this poses a fundamental challenge to POSIX style file systems
 - so we have to map between object-store approaches that give users control and POSIX file systems with appropriate performance and in use on supercomputers
- Software deployment needs to be easy
 - virtual machines; or
 - light-weight VMs: Docker / Shifter as a service

Architectural Developments – Traditional Architecture

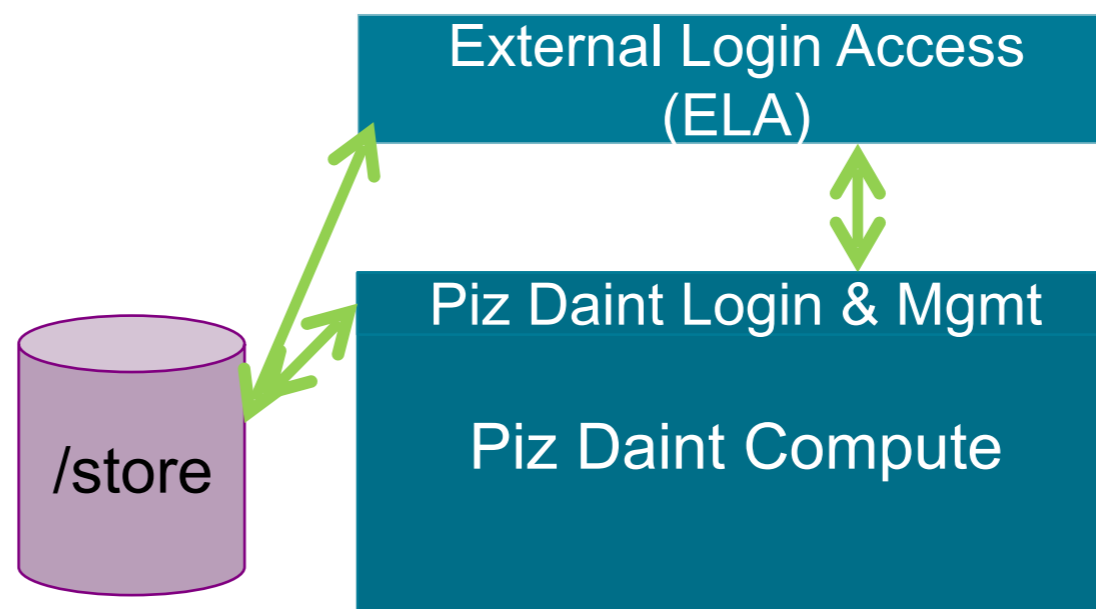
Research
Community



CSCS User



CSCS

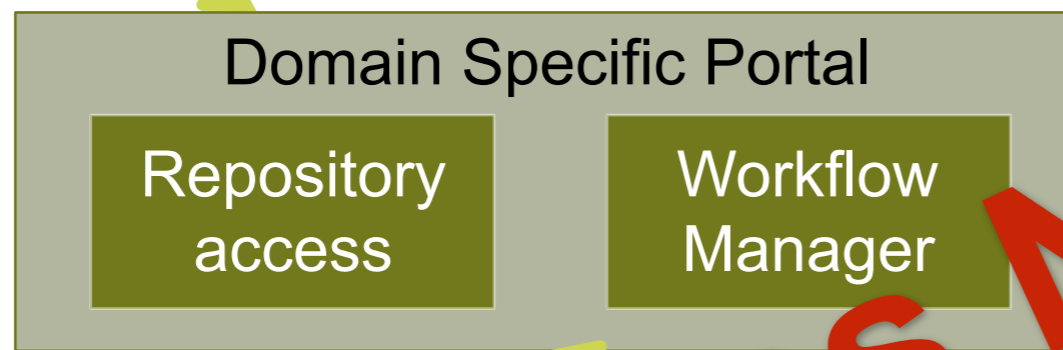


Architectural Developments – Improved Architecture Based on External Portal

Research Community

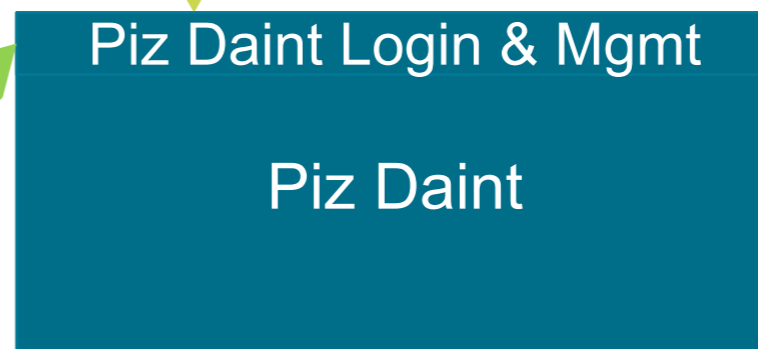


CSCS User



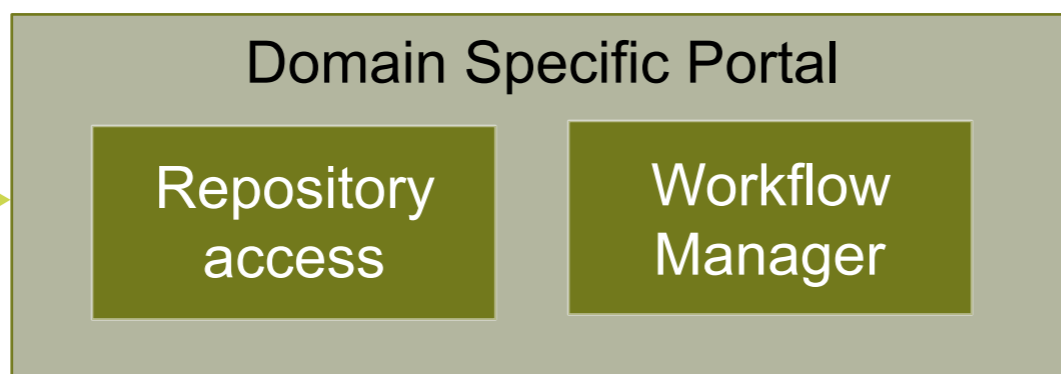
Does Not Scale

CSCS



Architectural Developments – Service Oriented Architecture (SOA)

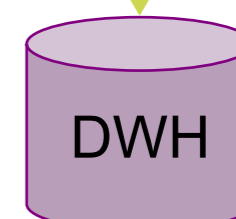
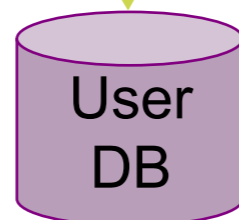
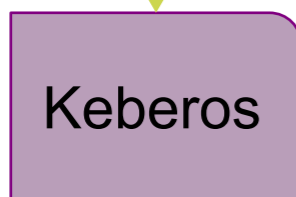
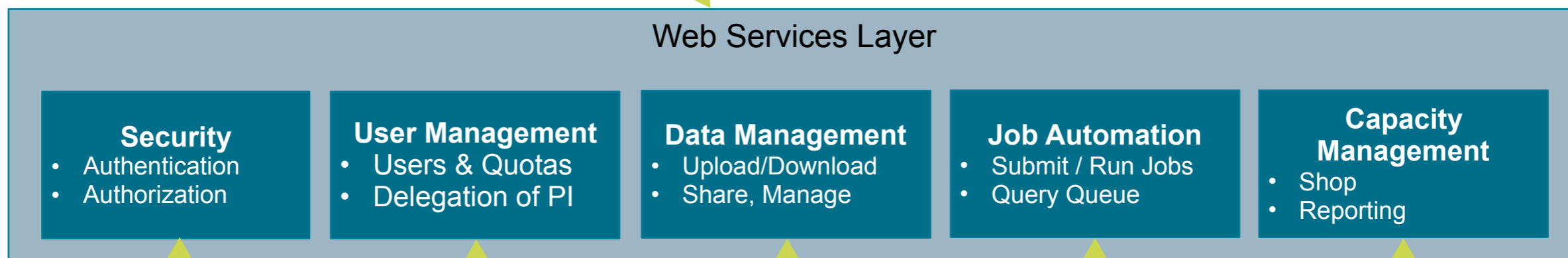
Research Community



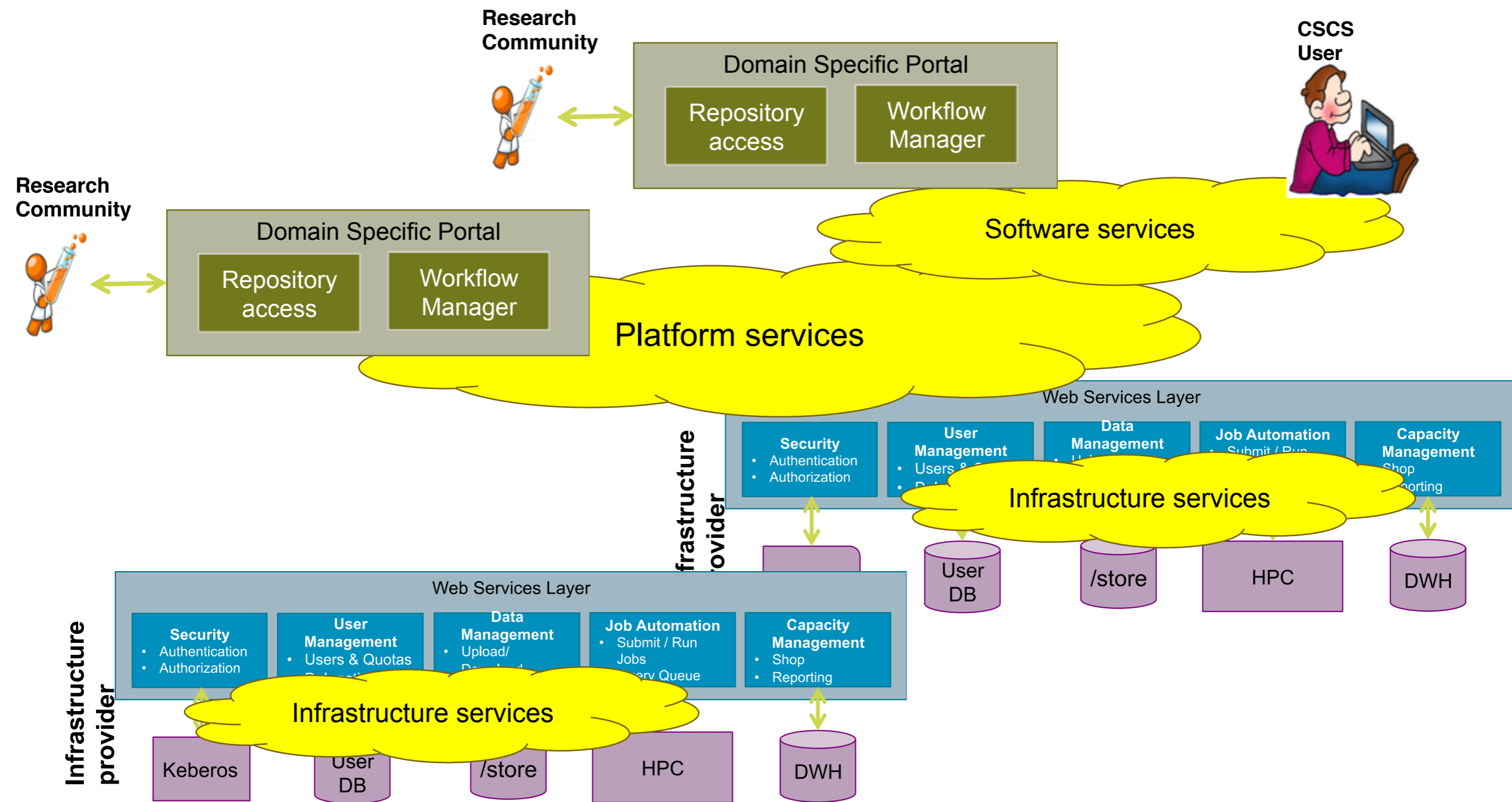
CSCS User



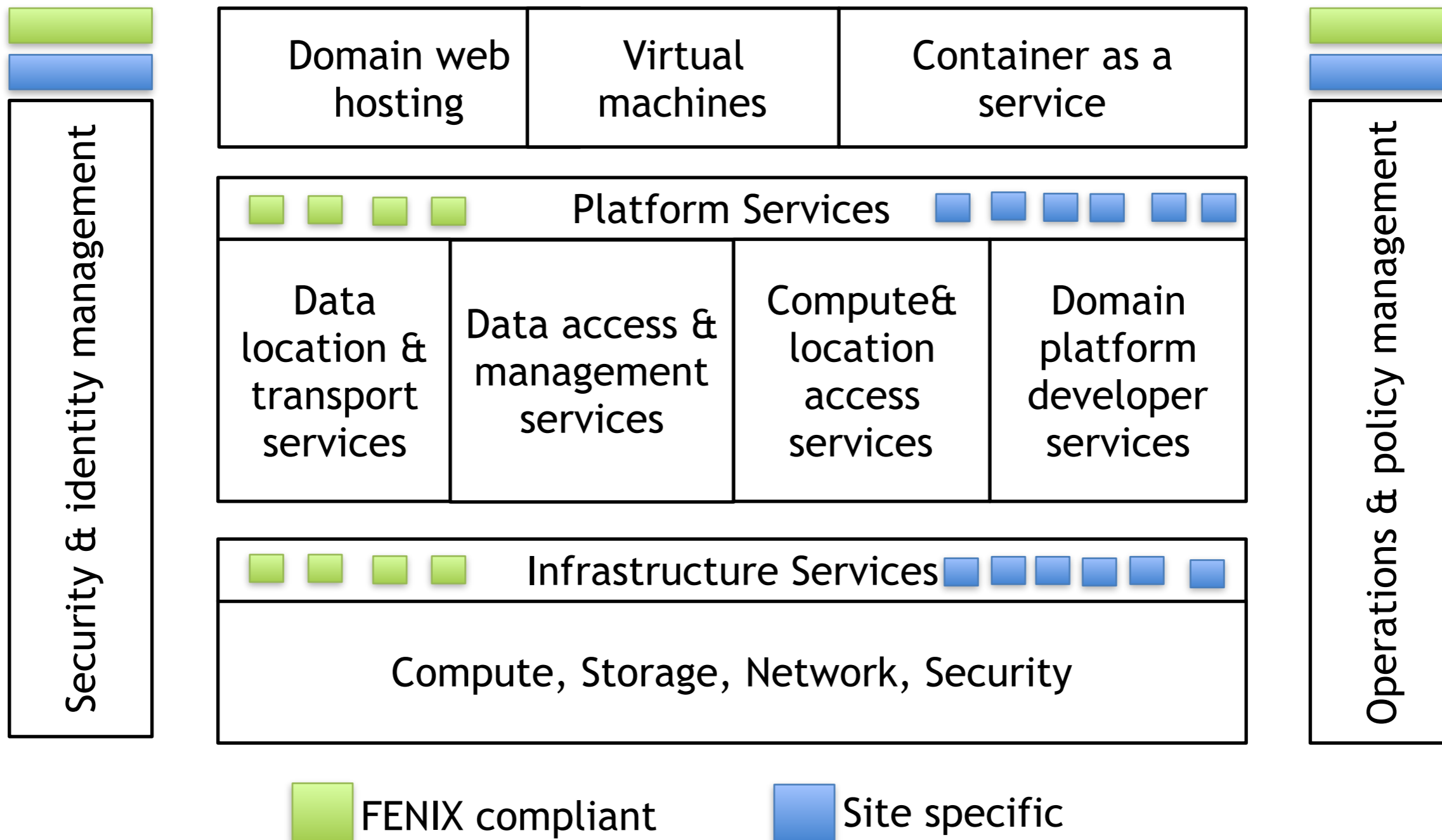
CSCS



Supporting Federation using SOA



The data center's concerns of SOA



2017 SC

Platform for Advanced Scientific Computing
Conference

Lugano
Switzerland

26-28 June 2017

