# HTCondor and TORQUE/Maui
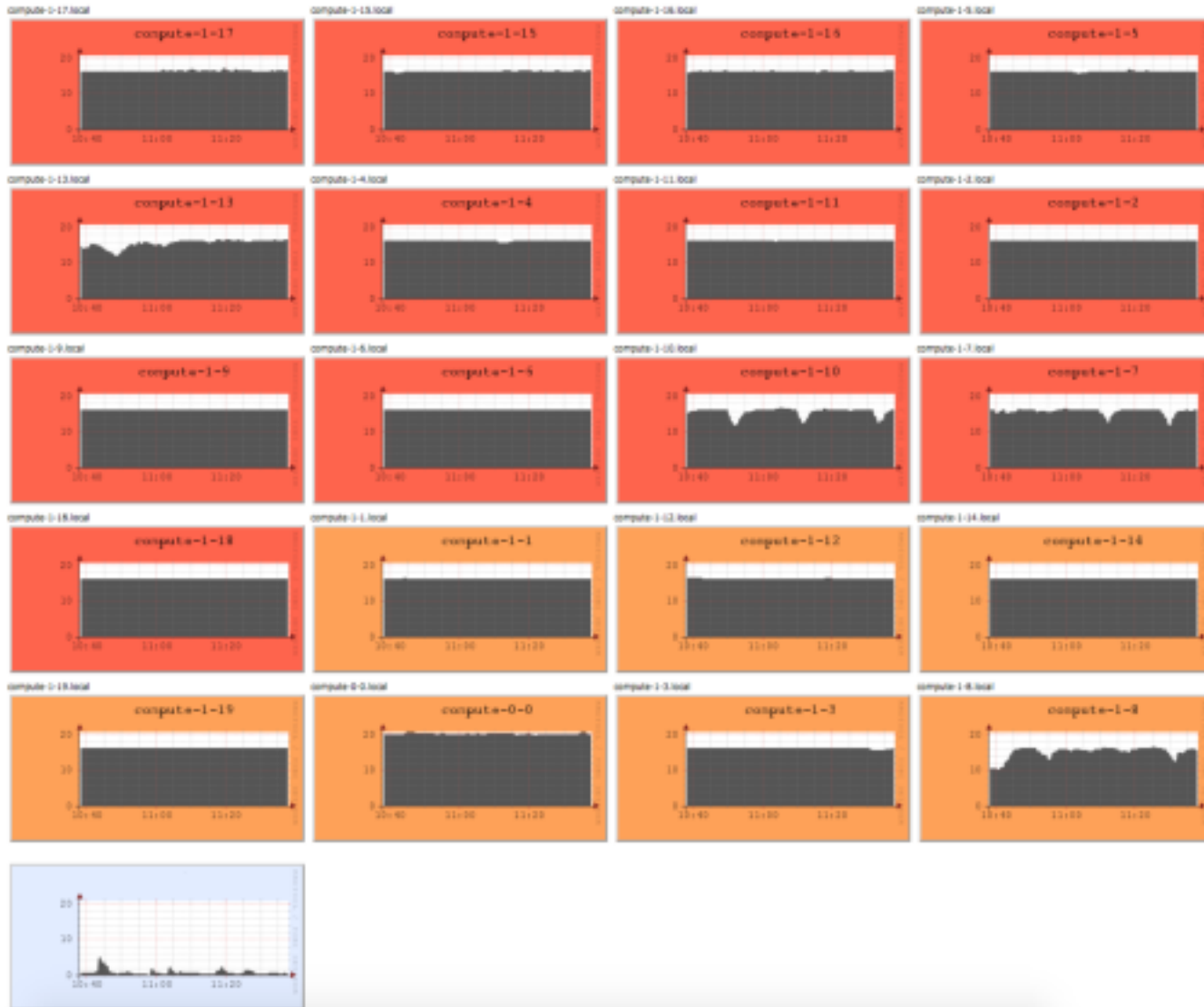## Resource Management  Integration
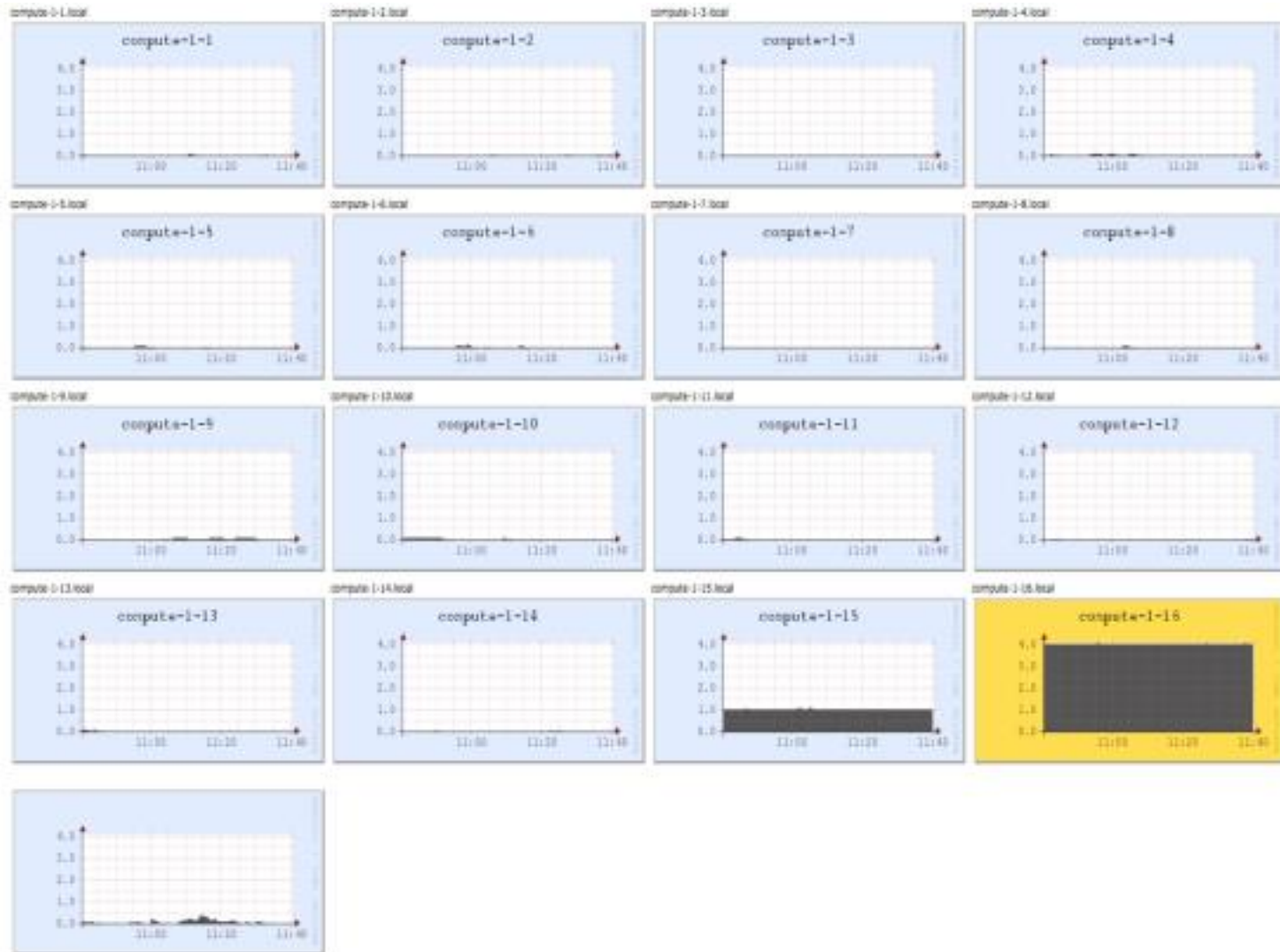
**Grigoriy Abramov**

Systems Administrator,  Research Data Center

Chicago, IL

# Problems:

- Low efficiency of available computational resources

- Disproportionate utilization of each individual cluster

- Absence of sharing computational resources

- Absence of a unified, resource-management platform

- Variations on Linux OS releases

- Distributed ownership of computational resources

# Advantage:

Availability of TORQUE/Maui resource management
and Portable Batch System (PBS)
on all HPC computational resources

# Goals:

- Cost-efficient optimization of computational resources

- Greater access to computational resources for faculty members and research groups who do not have the resources to obtain their own clusters

# Working Principles:

- Departments and research groups retain full ownership of their clusters and have priority in performing their computations

  o All running computational jobs submitted by "guests" on shared resources should be removed from the queue when needed

# Challenges:

- Finding an optimal platform that does not require system reinstallation or significant, configuration updates that would interrupt an already-running computation

# OS and Applications:

- ROCKS cluster OS

- TORQUE/Maui (Adaptive Computing)
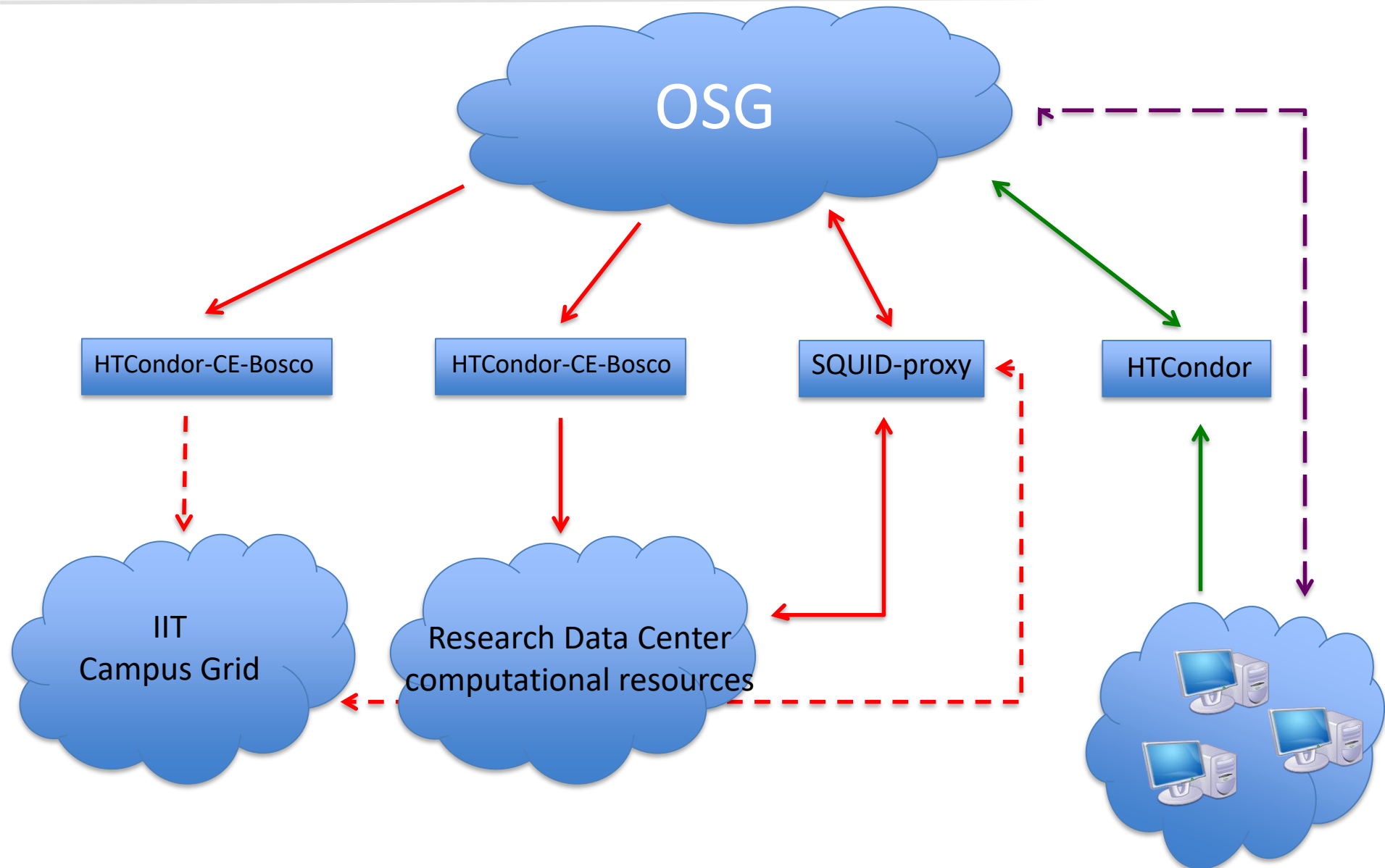
- HTCondor, HTCondor-CE-BOSCO, or both

Open Science Grid

IIT Compute Resource

# GridIIT/OSG Computing Grid Diagram

# Implementation:

- Single-user account and GID on all HPC clusters:
  <span style="color:red">guestuser/guestuser</span>
- PBS Queue-Manager (qmgr) configuration
- Maui configuration
- Installation of HTCondor-CE-Bosco or HTCondor
- Access to computational resources via Secure Shell (SSH)
- Testing computing grid resources for incoming/outgoing traffic

# Implementation | TORQUE configuration

Create and define queue <span style="color:red">grid</span> at qmgr prompt:

**<span style="color:red">create queue grid</span>**
**<span style="color:red">set queue grid queue_type = Execution</span>**
set queue grid max_user_queuable = 14
set queue grid resources_default.walltime = 48:00:00
set queue grid resources_default.ncpus = 1
**<span style="color:red">set queue grid acl_group_enable = True</span>**
**<span style="color:red">set queue grid acl_groups = guestuser</span>**
set queue grid kill_delay = 120
set queue grid keep_completed = 120
**<span style="color:red">set queue grid enabled = True</span>**
**<span style="color:red">set queue grid started = True</span>**

# Implementation | Maui configuration

- Priority
- Preemption
- Preemption policy
- Partitioning
- QOS – Quality of Services

# Implementation | Maui configuration (cont.)

RMCFG[base]  TYPE=PBS            SUSPENDSIG=15

PREEMPTPOLICY                    SUSPEND

NODEALLOCATIONPOLICY             PRIORITY

QOSCFG[hi]                       QFLAGS=PREEMPTOR

QOSCFG[low]                      QFLAGS=NOBF:PREEMPTEE

CLASSWEIGHT                      10

CLASSCFG[batch]                  QDEF=hi    PRIORITY=1000

CLASSCFG[grid]                   QDEF=low PRIORITY=1

# Implementation | MAUI configuration (cont.)

GROUPCFG[users]      PRIORITY=1000      QLIST=hi      QDEF=hi      QFLAGS=PREEMPTOR
GROUPCFG[guestuser]  PRIORITY=1         QLIST=low     QDEF=low     QFLAGS=PREEMPTEE
USERCFG[guestuser]   PRIORITY=1         QLIST=low     QDEF=low     QFLAGS=PREEMPTEE


PARTITIONMODE ON
NODECFG[compute-1-1]      PARTITION=grid
NODECFG[compute-1-2]      PARTITION=grid


SYSCFG[base]            PLIST=default, grid&
USERCFG[DEFAULT]        PLIST=default
GROUPCFG[guestuser]     PLIST=default:grid   PDEF=default


* Maui service needs to be restarted

- Test job submission via PBS as guestuser on compute cluster

- In submit script, the below-listed options should be presented:

  <span style="color:red">#PBS -q grid</span>
  <span style="color:red">#PBS -W x="PARTITION:grid"</span>

- Reliability of PREEMPTION needs to be verified

- Install and configure HTCondor or HTCondor-CE-BOSCO

- Add on computational cluster's head node following lines to file

  <span style="color:blue">*../bosco/glite/bin/pbs_local_submit_attributes.sh*</span>
  <span style="color:red">#!/bin/sh</span>
  <span style="color:red">echo "#PBS -q grid"</span>
  <span style="color:red">echo '#PBS -W x="PARTITION:grid"'</span>

- Submit test job from remote server via command:

  <span style="color:red">bosco_cluster –t guestuser@your_cluster_name.edu</span>

# GridIIT and OSG Shared Computational Resources Over a 6-Month Period
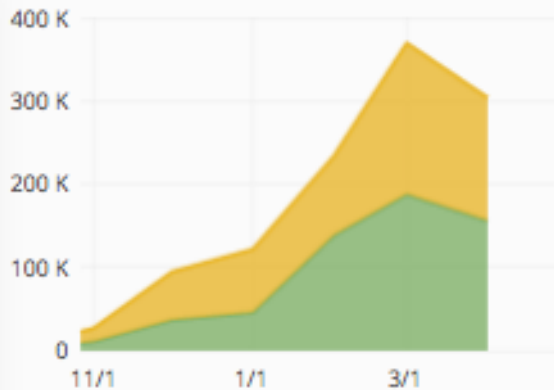


*Opportunistic    *Dedicated
Source:  http://gracc.opensciencegrid.org