# Pacific Research Platform and The Pacific Rim

**Presenter: Tom DeFanti, Research Scientist, QI, Co-PI**

**Larry Smarr, Calit2 Director and PI**

**Ilkay Altintas, Camille Crittenden, Ken Kreutz-Delgado, Phil Papadopoulos, Tajana Rosing, Frank Wuerthwein, co-PIs**

**John Graham, Senior Development Engineer**

**Dima Mishin, Isaac Nealey, Joel Polizzi, Mark Yashar, Programmers**

**UC San Diego and UC Berkeley**

# 2015 Vision: The Pacific Research Platform will Connect Science DMZs Creating a Regional End-to-End Science-Driven Community Cyberinfrastructure



Source: John Hess, CENIC

**NSF CC\*DNI Grant**
**$6.3M 10/2015-10/2020**
**Year 5 Starts in 3 Weeks!**

PI: Larry Smarr, UC San Diego Calit2
Co-PIs:

- **Camille Crittenden, UC Berkeley CITRIS,**
- **Tom DeFanti, UC San Diego Calit2/QI,**
- **Philip Papadopoulos, UCI**
- **Frank Wuerthwein, UCSD Physics and SDSC**

**Letters of Commitment from:**
- **50 Researchers from 15 Campuses**
- **32 IT/Network Organization Leaders**

**ESnet: Given Fast Networks, Need DMZs and Fast/Tuned DTNs**

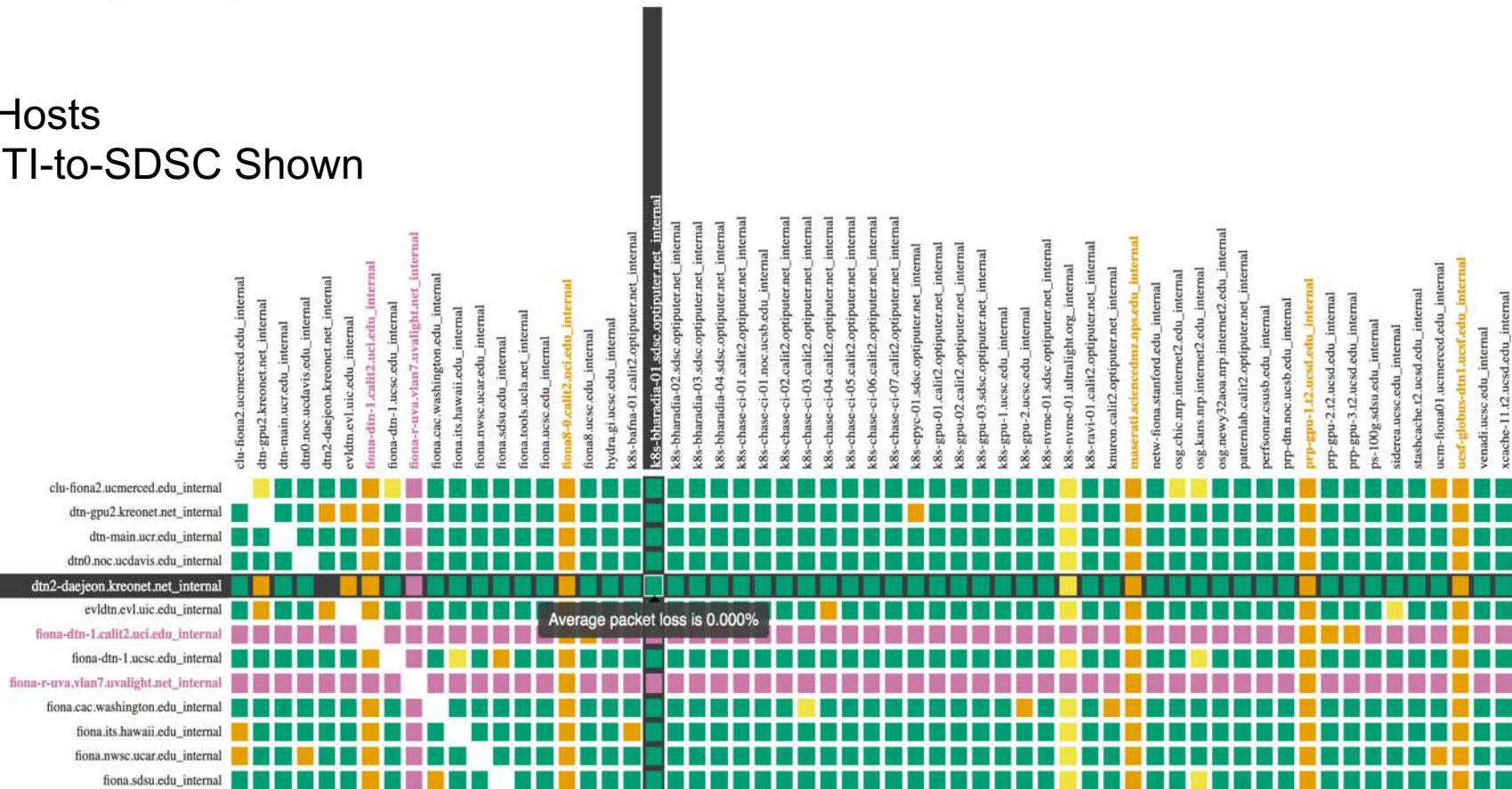# PRP Technical Deliverables
# 2015 - 2017

- **Phase 0: Tested Layer2 CENIC Networks and FIONAs—early 2015**
- **Phase 1: A Scalable Network for Optimizing Data Transfer—2015-2017**
  - **Layer 3 Data Transfer & Measurement Network**
  - **Tested, Debugged, Measured, Optimized, and MaDDash'ed Layer 3**
  - **Supported Rates up to 9.7/10, 37/40 Gb/s in 10GB Bursts**
  - **Included UvAmsterdam and Korea (KISTI)**
  - **Showed Full Bandwidth Utilization**
  - **Essentially No TCP Backoff on Long Distance Best-Efforts Networks**
- **This is What Most Other Research Platforms are Focusing on—Big Data Transfer**

# Nautilus Mesh - Latency - Loss

Loss rate is <= 0.001%    Loss rate is > 0.001%    Loss rate is >= 0.1%    Unable to find test data    Check has not run yet

⚠ Found a total of 8 problems involving 6 hosts in the grid

60 Hosts
KISTI-to-SDSC Shown

Average packet loss is 0.000%

# PRPv1 to PRPv2:
# The Transition from Network Diagnosing to Application Support

- **PRPv1 Designed, Built, and Installed ~40 Purpose-built FIONAs, Tuned to Measure and Diagnose End-to-End 10G, 40G and 100G**

- **But, Our PRP NSF Funding Requires Showing Use of the PRP by Scientists and Engineers—It's a Data Grant, not a Networking Grant**

- **Note: Our Scientists Clearly Need More than Bandwidth Tests**
  - **Teams of Scientists Want to Share Their Data at High Speed and Compute on It**
  - **They Need to Interoperate with Commercial and University Clouds**

- **So PRPv2 Added DMZ-Distributed Temporary Storage**
  - **1.7PB total in 14 ~200GB previous PRPv1 FIONAs in Campus DMZs**

# Detailed Real-Time Monitoring of PRP Nautilus:
## UCD, UCSD, UCI, UCSB, UCLA, UCR, Stanford, UCAR, UCM, UCSC, UHM Ceph
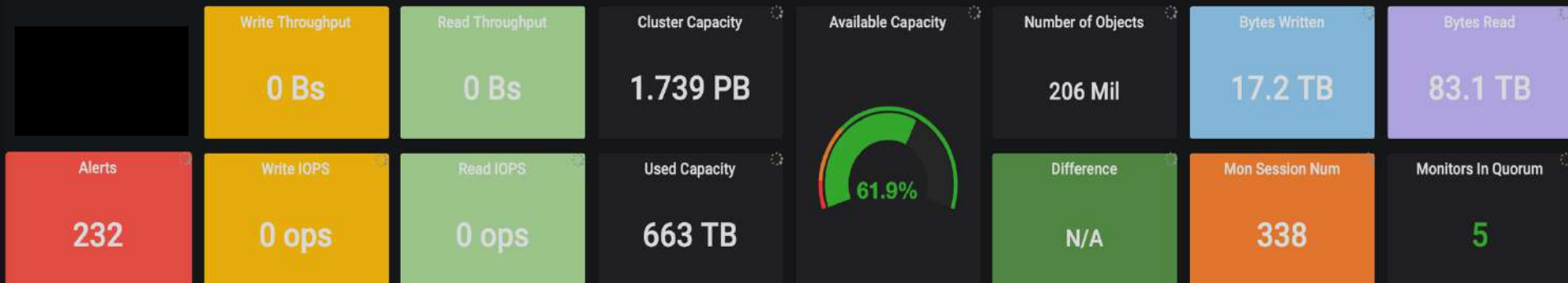


**Ceph - Cluster New**

Last 90 days | Refresh every 1m

Interval: auto

**This is Working Scratch Space for Data Transfer, Not Archival Research Storage**

### CLUSTER STATE

| Write Throughput | Read Throughput | Cluster Capacity | Available Capacity | Number of Objects | Bytes Written | Bytes Read |
|---|---|---|---|---|---|---|
| 0 Bs | 0 Bs | 1.739 PB | 61.9% | 206 Mil | 17.2 TB | 83.1 TB |

| Alerts | Write IOPS | Read IOPS | Used Capacity | | Difference | Mon Session Num | Monitors In Quorum |
|---|---|---|---|---|---|---|---|
| 232 | 0 ops | 0 ops | 663 TB | | N/A | 338 | 5 |

### OSD STATE

3/11/2019

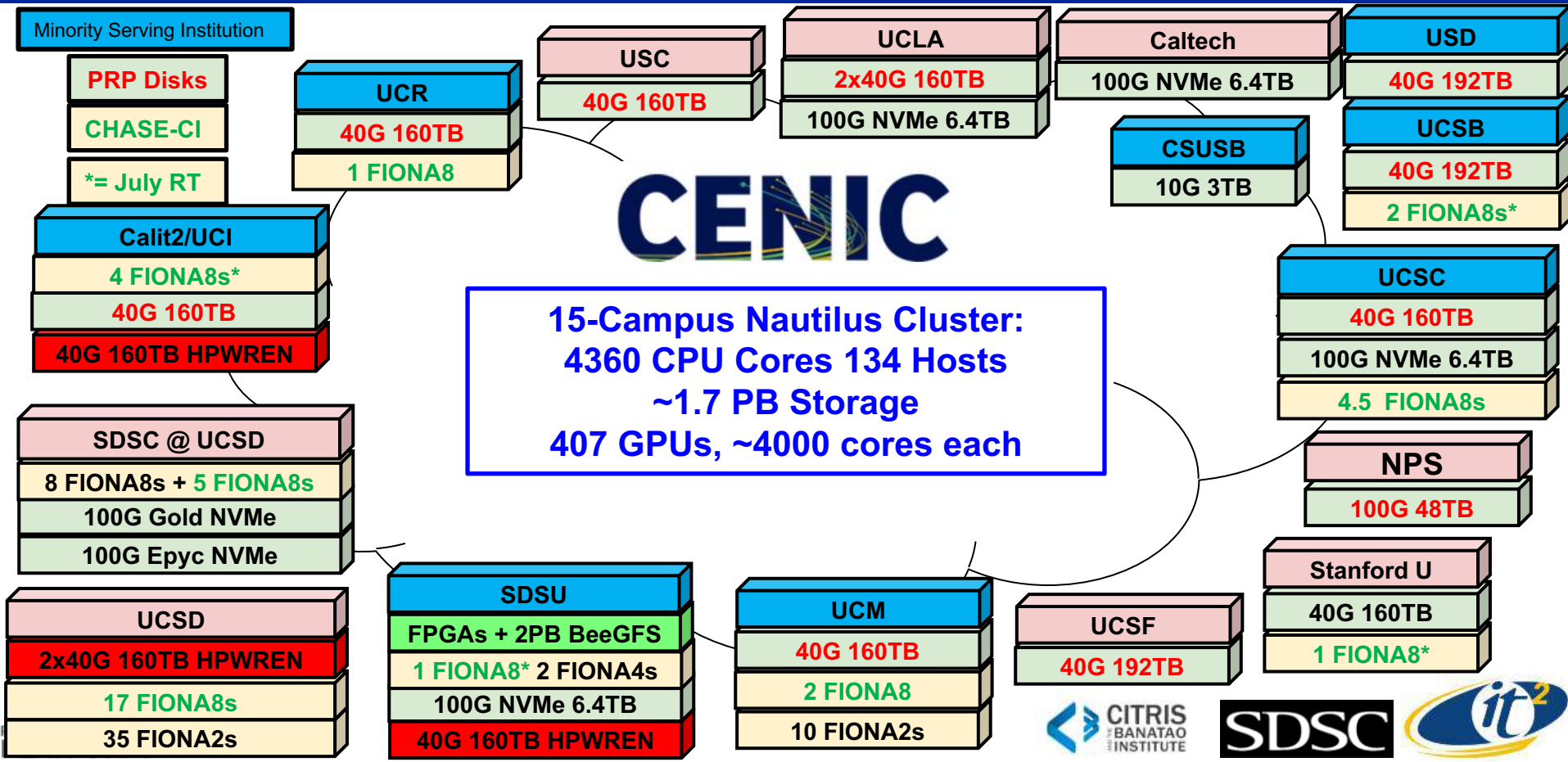| OSDs OUT | OSDs DO... | OSDs UP | OSDs IN | Avg PGs | Avg Apply Latency | Avg Commit Latency | Avg Op Write Latency | Avg Op Read Latency |
|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 153 | 153 | 214 | 17 ms | 17 ms | 0.14 ms | 0.03 ms |

> Alerts (3 panels)

> Ceph Versions (4 panels)

https://grafana.nautilus.optiputer.net/d/r6lloPJmz/ceph-cluster-new?orgId=1&from=now-90d&to=now

# Regional Scale Cluster: Connected by PRP's Use of CENIC 100G Network
## PRP's Nautilus Hypercluster

Nautilus

**Minority Serving Institution**

**PRP Disks**

**CHASE-CI**

**\*= July RT**

**USC**
40G 160TB

**UCLA**
2x40G 160TB
100G NVMe 6.4TB

**Caltech**
100G NVMe 6.4TB

**USD**
40G 192TB

**UCSB**
40G 192TB
2 FIONA8s*

**UCR**
40G 160TB
1 FIONA8

**CSUSB**
10G 3TB

**Calit2/UCI**
4 FIONA8s*
40G 160TB
40G 160TB HPWREN

## CENIC

**UCSC**
40G 160TB
100G NVMe 6.4TB
4.5 FIONA8s

**15-Campus Nautilus Cluster:**
**4360 CPU Cores 134 Hosts**
**~1.7 PB Storage**
**407 GPUs, ~4000 cores each**

**NPS**
100G 48TB

**SDSC @ UCSD**
8 FIONA8s + 5 FIONA8s
100G Gold NVMe
100G Epyc NVMe

**Stanford U**
40G 160TB
1 FIONA8*

**UCSD**
2x40G 160TB HPWREN
17 FIONA8s
35 FIONA2s

**SDSU**
FPGAs + 2PB BeeGFS
1 FIONA8* 2 FIONA4s
100G NVMe 6.4TB
40G 160TB HPWREN

**UCM**
40G 160TB
2 FIONA8
10 FIONA2s

**UCSF**
40G 192TB

CITRIS BANATAO INSTITUTE

SDSC

it²

# Grafana Showing State of Nautilus 9-10-19

MSU

UCB

UCM

Stanford

UCSC

Caltech    UCI    UCR

UCSD    SDSU

**CI-New: Cognitive Hardware and Software Ecosystem Community Infrastructure (CHASE-CI)**

**For the Period September 1, 2017 – August 31, 2020**

**SUBMITTED – January 18, 2017**

PI: Larry Smarr, Professor of Computer Science and Engineering, Director Calit2, UCSD
Co-PI: Tajana Rosing, Professor of Computer Science and Engineering, UCSD
Co-PI: Ken Kreutz-Delgado, Professor of Electrical and Computer Engineering, UCSD
Co-PI: Ilkay Altintas, Chief Data Science Officer, San Diego Supercomputer Center, UCSD
Co-PI: Tom DeFanti, Research Scientist, Calit2, UCSD

NSF Grant for 256 High Speed "Cloud" GPUs
For 32 ML Faculty & Their Students at 10 Campuses
To Train AI Algorithms on Big Data

PRP
PACIFIC RESEARCH PLATFORM

SDSC    UC San Diego    it²

# Road Trip! Installing Community Shared Storage and GPUs in June, December & January at UC Merced, UC Santa Cruz, UC Riverside, and Stanford



160-192TB added to 14 Existing PRPv1 FIONAs

New FIONA8 at UCSC

PACIFIC RESEARCH
PLATFORM

# PRP Engineers Designed and Built Several Generations of Optical-Fiber Big-Data Flash I/O Network Appliances (FIONAs)

## UCSD-Designed FIONAs Solved the Disk-to-Disk Data Transfer Problem *at Near Full Speed* on Best-Effort 10G, 40G and 100G Networks



Two FIONA DTNs at UC Santa Cruz: 40G & 100G
Up to 192 TB Rotating Storage

Add Up to 8 Nvidia GPUs Per 2U FIONA
To Add Machine Learning Capability

FIONAs Designed by UCSD's Phil Papadopoulos, John Graham, Joe Keefe, and Tom DeFanti

| Top Nautilus GPU users August 2019 | | | | | |
|---|---|---|---|---|---|
| PI | Campus | August 2019 GPU SU | FIONA8 Equivalent | August 2019 CPU SU | August 2019 Mem SU |
| Frank Wuerthwein | UCSD | 80084 | 13.90 | 398124.41 | 8.13864E+14 |
| Mark Alber | UCR | 40761 | 7.08 | 37131.21 | 6.60061E+13 |
| Hao Su | UCSD | 16396 | 2.85 | 42547.91 | 2.78718E+14 |
| Nuno Vasconcelos | UCSD | 10991 | 1.91 | 11218.07 | 9.11693E+13 |
| Jeff Krichmar | UCI | 6587 | 1.14 | 6997.06 | 2.20582E+13 |
| Falko Kuester | UCSD | 6211 | 1.08 | 35404.91 | 5.68019E+14 |
| Anshul Kundaje | Stanford | 6063 | 1.05 | 1481.62 | 5.38638E+13 |
| Ravi Ramamoorthi | UCSD | 4822 | 0.84 | 6767.49 | 3.83436E+13 |
| Larry Smarr | UCSD | 4359 | 0.76 | 3171.25 | 2.20892E+13 |
| Manmohan Chandraker | UCSD | 3788 | 0.66 | 3304.47 | 1.02188E+14 |
| Tom DeFanti | UCSD | 3203 | 0.56 | 2040.4 | 8.82778E+12 |
| Nuno Vasconcelos | UCSD | 2293 | 0.40 | 3797.22 | 3.37342E+13 |
| Kurt Schoenhoff | Australia | 1921 | 0.33 | 4910.91 | 1.79054E+13 |
| Nuno Vasconcelos | UCSD | 1888 | 0.33 | 1017.46 | 1.67571E+13 |
| Dinesh Bharadia | UCSD | 1771 | 0.31 | 5724.15 | 2.71821E+13 |
| Padhraic Smyth | UCI | 1387 | 0.24 | 647.53 | 1.09787E+13 |
| Jurgen Schulze | UCSD | 1330 | 0.23 | 10.88 | 3.9717E+12 |
| Larry Smarr | UCSD | 1314 | 0.23 | 0.57 | 2.34185E+12 |
| Jurgen Schulze | UCSD | 1306 | 0.23 | 0.7 | 1.92583E+12 |
| Nuno Vasconcelos | UCSD | 1209 | 0.21 | 5984.29 | 1.33191E+13 |
| Eric Shearer | UCI | 1131 | 0.20 | 1308.7 | 3.85832E+12 |

Top Nautilus GPU Users in August 2019

FIONA8 equivalent: running an 8-GPU machine 24x7x30

Top User is IceCube in OSG background mode

Others are ML

# 2017: PRP Connected 70 UCSD SunCAVE and 20 UCM WAVE 4K Screens to Share VR
# 2018: Added their 90 Game GPUs to PRP/OSG for Machine Learning Computations



UC Merced WAVE 20 Screens 20 GPUs
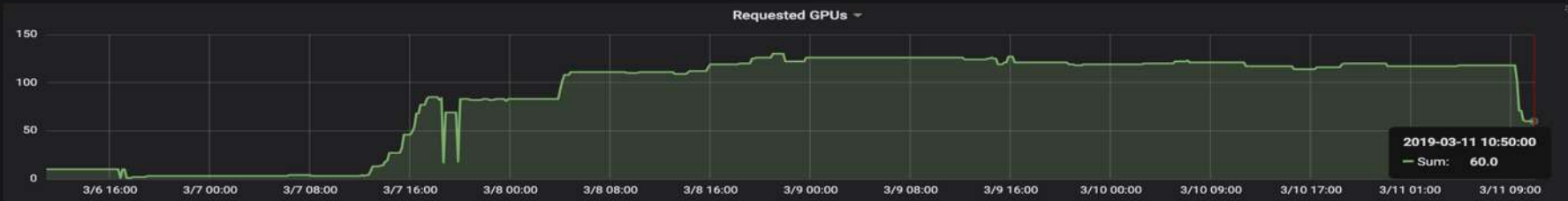


UCSD SunCAVE 70 Screens 70 GPUs

Leveraging UCM Campus Funds and NSF CNS-1456638 & CNS-1730158 at UCSD

By Amble - Own work, CC BY-SA 3.0, https://commons.wikimedia.org/w/index.php?curid=8773726

- IceCube Neutrino Observatory has been using 120 Nautilus GPUs since March 8

- This would cost $2,880/day in a commercial cloud (at $1/hr) or ~$20,000/week

- An 8-GPU FIONA8 for Nautilus costs $20,000 to buy

GPU Simulations are Needed *to Improve Ice Model*.
=> Results in Significant Improvement in Pointing Resolution
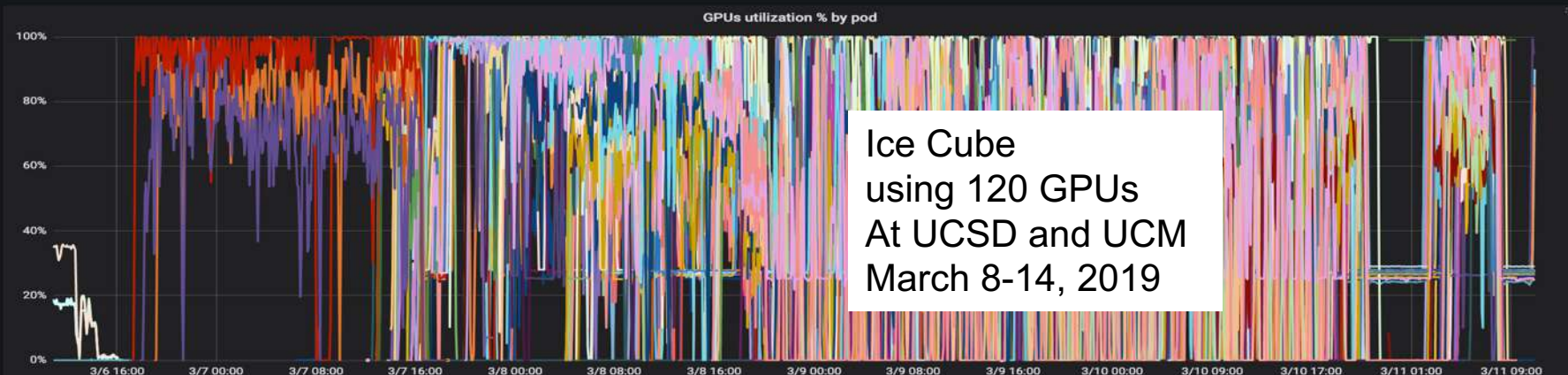for Multi-Messenger Astrophysics

K8s / Compute Resources / Namespace GPUs

namespace  osggpus

Requests

Requested GPUs

2019-03-11 10:50:00
— Sum:  60.0

Utilization

GPUs utilization % by pod

Ice Cube
using 120 GPUs
At UCSD and UCM
March 8-14, 2019

March 13, 2019: "This morning there was a big demo in the SunCAVE. The IceCube pods were kicked out automatically when the SunCAVE GPUs were in use, and restarted when the demo was over. **No admin intervention needed**."—Igor Sfiligoi, SDSC

# Very Cost-Effective for Academic Machine Learning and Data Sharing

- **Data science researchers need DTNs with lots of storage, encryption and lots of GPUS**
- **One UC spends $40,000 in cloud GPU per published grad student paper**
- **Another spends $20,000 for undergrad ML AWS access in just one course**
- **Instead, add to our Nautilus hypercluster (or clone it & federate):**
  - **UCSD ECE Department bought 4 FIONA8s, buying 4 more**
  - **UCSD Physics Department. bought 6 FIONA8s**
  - **UCSD CSE researchers bought 4 FIONA8s to add to Nautilus, buying 20 more**
  - **UCSD Instructional IT has 13 FIONA8s for Machine Learning/AI class labs**
- **Working Storage on Nautilus FIONAs is**
  - **very inexpensive (12TB drives are ~$430 each—16 per FIONA. FISA encrypted drives @ same cost)**
  - **and very high speed (most FIONAs are 40/100G and are located in ScienceDMZs)**

Clemson's Alex Feltus: "I cannot wait to add a node to the Nautilus compute fabric!" (He didn't wait)

PRP
PACIFIC RESEARCH
PLATFORM

it²

# UCSD's Information Technology Services Adapted PRP FIONA8s To Support Data Science Courses

Instructional Data Science
Machine Learning Platform:

Instead of Spending
~$20,000/Quarter/Course on
Commercial Clouds:
97 Courses over 6 Quarters →
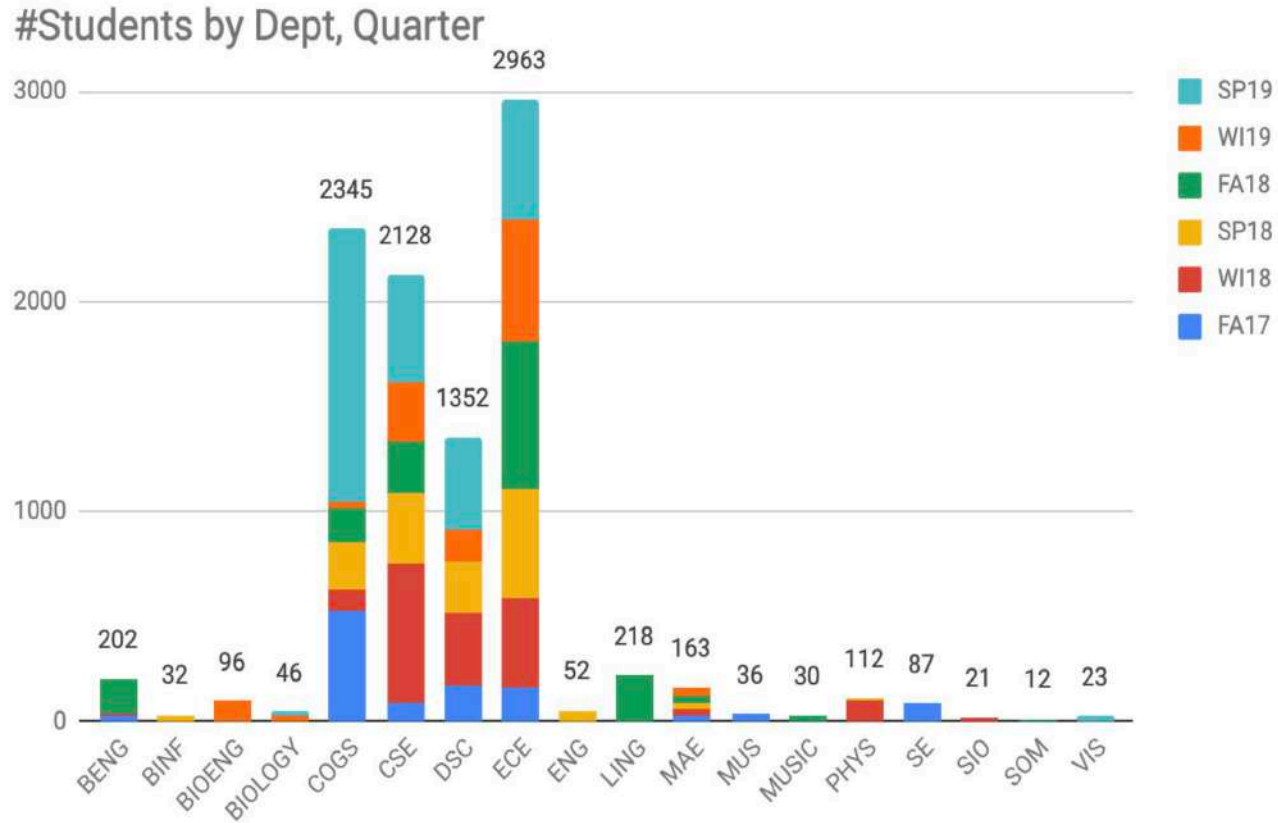$4M vs. $240K over 12 Quarters
At least 20,000 Students

Adam Tilghman, ITS

- 104 GPUs
  - 80 GTX 1080Ti
  - 16 RTX 2080Ti
  - 1.05 Petaflops
- 28 nodes
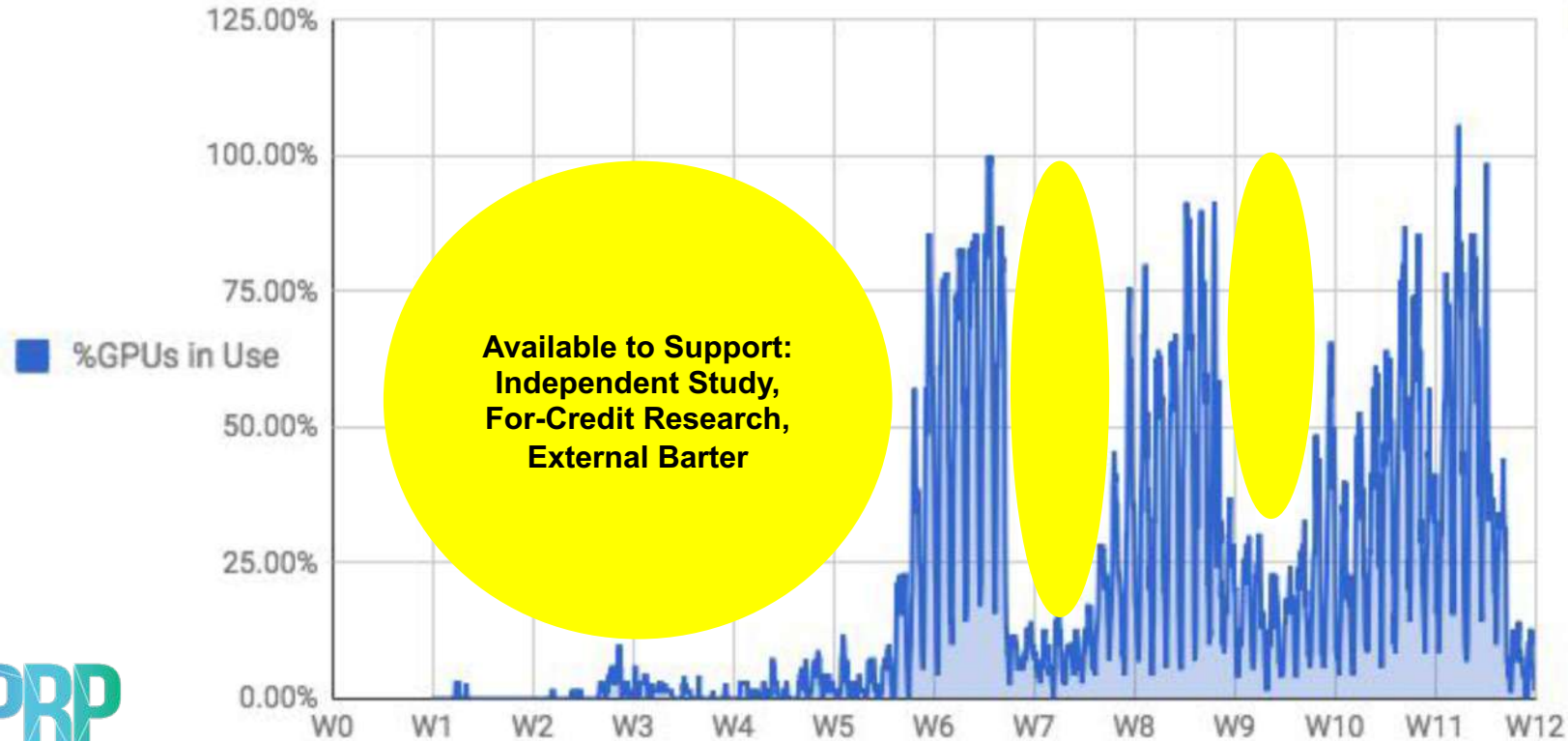  - 544 CPU cores
  - 6.5TB RAM
- 40TB Flash Storage

- 10G networking

Source: UCSD ITS

Source: UCSD ITS

# Student GPU Demand Is Variable
# Allowing for Other Student Uses



WI18 Instructional GPU Utilization

%GPUs in Use

Available to Support:
Independent Study,
For-Credit Research,
External Barter

Source: UCSD ITS

# 405 Research GPUs in Nautilus 9-10-2019

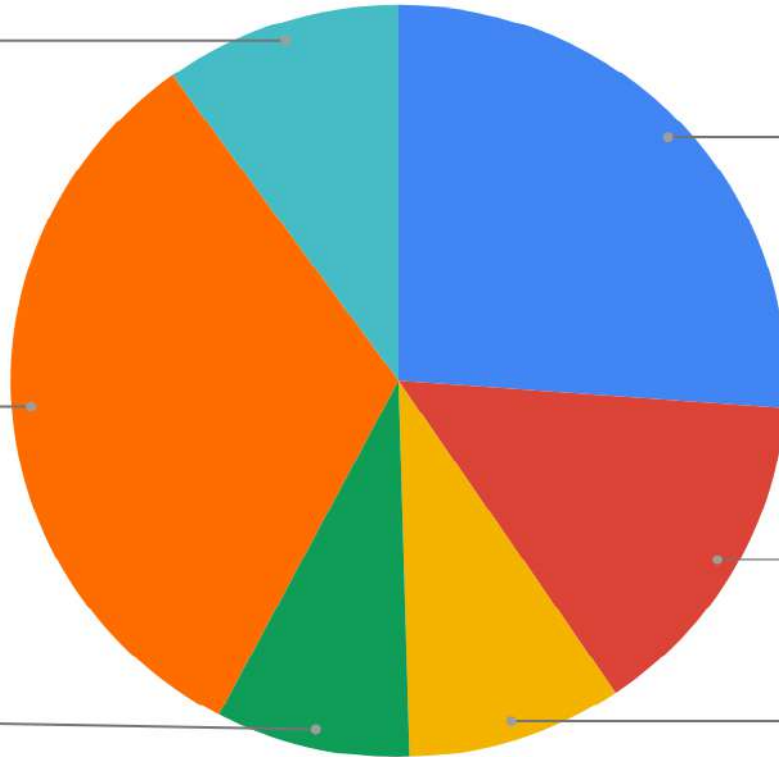134 Nautilus Nodes



UCM-WAVE
9.9%

PRP
26.1%

SunCAVE
32.4%

CHASE-CI
14.4%

OSG
8.1%

PRIVATE
9.0%

# Why PRPv2 Adopted Kubernetes

- **PRP FIONAs Are Coupled by Kubernetes Into the "Nautilus Hypercluster"**
  - **Kubernetes "Pods" Encapsulate Application Container(s), Storage Resources, and Execution Options**
  - **Implements PRP Cooperative Research Groups Support with Policy-Based Scheduling by Use of CILogon and Kubernetes Namespaces—704 Users in Namespaces as of 7/15/19**
  - **Allows Cloud Native Storage Integration (e.g., Rook/Ceph/EdgeFS)**
  - **Enables Us to Update Overnight, without local assistance, a RP Scaling Necessity**
  - **Emerging Solutions for Sophisticated SDN Overlay Network, Firewall, and Network Policy Controls**
- **Allows Easy User Job Scaling to Heterogeneous Platforms:**
  - **Deskside, Rack-Mounted, Supercomputers, even IOT Gizmos like ML on Remote Cameras**
  - **Amazon Elastic Container Service for Kubernetes (Amazon EKS)**
  - **Google Kubernetes Engine (GKE) (TensorFlow)**
  - **Microsoft Azure Kubernetes Service (AKS)**
  - **Also Comet and other XSEDE assets**

"Kubernetes with Rook/Ceph Allows Us to **Manage** Petabytes of Distributed Storage and GPUs for Data Science, While We Measure and Monitor Network Use."
--John Graham, Calit2/QI UC San Diego

"Towards The NRP" 3-Year Grant Funded by NSF $2.5M October 2018

PI Smarr
Co-PIs Altintas
Papadopoulos
Wuerthwein
Rosing
DeFanti

REGIONAL RESEARCH AND EDUCATION NETWORKS IN THE UNITED STATES

Original PRP

NRP Pilot

Announced May 8, 2018
Internet2 Global Summit

NSF CENIC Link

CENIC/PW Link

# NRP_GridFTP - Throughput

Throughput >= 7500Mbps    Throughput < 7500Mbps    Throughput <= 5000Mbps    Unable to retrie

⚠ Found a total of 4 problems involving 3 hosts in the grid

Now testing to AWS and Internet2



| | fiona01.tacc.utexas.edu | dtn1-v6.nysernet.org | hcc-fiona-v6.unl.edu | gpn-fiona-mizzou.scidmz.rnet.missouri.edu | dtn1.nysernet.org | dtn0.lsanca.pacificwave.net | osg.kans.nrp.internet2.edu | fiona.sce.pennren.net | ps-40g-gridftp.calit2.optiputer.net | osg.newy32aoa.nrp.internet2.edu | dtn0.uog.edu | hcc-fiona.unl.edu | osg.chic.nrp.internet2.edu |

Average throughput is 12.456Gbps
Average throughput is 9.727Gbps

Hosts (rows):
fiona01.tacc.utexas.edu
dtn1-v6.nysernet.org
hcc-fiona-v6.unl.edu
gpn-fiona-mizzou.scidmz.rnet.
dtn1.nysernet.org
dtn0.lsanca.pacificwave.net
osg.kans.nrp.internet2.edu
fiona.sce.pennren.net

**TNRP = PRP (CENIC, PNWGP, FRGP, HI, and MREN) + OSG + ESnet + Quilt + NRP Pilot (I2, KINBER, Learn, GPN, NYSERnet) + MCNC + NM Tribal + …**
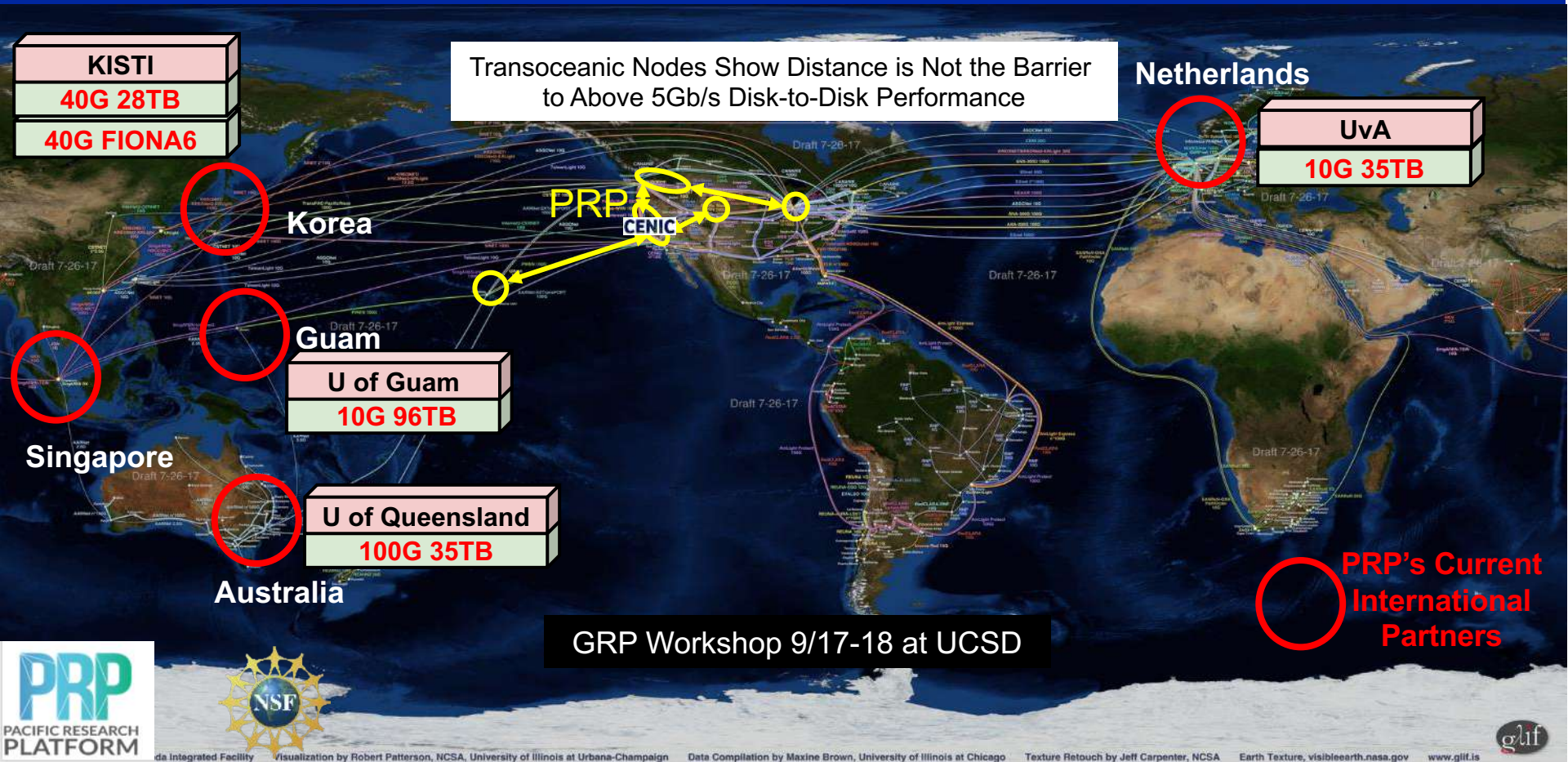
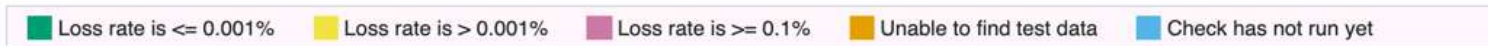Towards The NRP (TNRP)
3-Year $2.5M NSF Grant
OAC-1826967

Third NRP Workshop
September 24-25
Minneapolis

Original PRP

NRP Pilot

CENIC/PW Link

I2 CENIC Link

# Nautilus Has International Nodes
# The Global Research Platform is Emerging (1GRP Next Week Here!)

**KISTI**
**40G 28TB**
**40G FIONA6**

Transoceanic Nodes Show Distance is Not the Barrier to Above 5Gb/s Disk-to-Disk Performance

**Netherlands**

**UvA**
**10G 35TB**

PRP
CENIC

**Korea**

**Guam**
**U of Guam**
**10G 96TB**

**Singapore**

**U of Queensland**
**100G 35TB**

**Australia**

**PRP's Current International Partners**

GRP Workshop 9/17-18 at UCSD

PRP
PACIFIC RESEARCH PLATFORM

NSF

# Excellent Performance California to UQ (100G)



Nautilus Mesh - Latency ucsd - Loss

ElastiFlow: See Inter-cluster Campus-Level Traffic Flow Grouped by AS

# PRP Tech Coming

- **Support users with IoT/Robotics/Augmented Reality needs**
  - **Nvidia Jetson Xaviers and Nanos**
- **Also FPGA data-center boards (Xilinx U200s, Micron SB-852)**
  - **Compute: application acceleration (e.g.,TensorFlow)**
    - **Climate/weather segmentation**
    - **Inferencing**
    - **Satellite imagery ortho-rectification (align w/wildfire maps)**
  - **100G SDX P4 build out (SDSU, USC, NU, FIU, UCSD, Caltech)**
- **And Tensor Flow Cores and TPUs**
  - **Nvidia 2080-Ti cards: 544 Tensor Cores each, 4,352 per FIONA8**
  - **Our Nautilus Users can Access Google Cloud TPUs**
  - **Google Edge TPU Coral Development Boards and USB-C Edge TPU Accelerator/co-Processor**

Passive Option

![PRP Pacific Research Platform logo]

# P2PRP: Pacific to Pacific Rim Platform!

- **Top down Great Networking with 10-100Gbps Science DMZ Performance is a Necessary *but not Sufficient Condition* for Data-Driven Researchers**
  - **They need Science DMZs & DTNs with Lots of Low-Cost Storage, Encryption, Large RAM CPUs, GPUs, TPUs, FPGAs, and High-Availability Computing**

- **Measuring and Monitoring is Key to Better Usage and Security**

- **Compatibility with CloudBank, Google, Microsoft, and Amazon Clouds, and NSF/DOE Supercomputers Helps Ensure Scalability and Continuation**

- **Convergence with Open Science Grid/I2 Brings In Global Experience**

PRP
PACIFIC RESEARCH
PLATFORM

it²

# PRP/TNRP/CHASE-CI Support and Community:

- **US National Science Foundation (NSF) awards to UCSD, NU, and SDSC**
  - **CNS-1456638, CNS-1730158, ACI-1540112, ACI-1541349, & OAC-1826967**
  - **OAC 1450871 (NU) and OAC-1659169 (SDSU)**
- **UC Office of the President, Calit2 and Calit2's UCSD Qualcomm Institute**
- **San Diego Supercomputer Center and UCSD's Research IT and Instructional IT**
- **Partner Campuses: UCB, UCSC, UCI, UCR, UCLA, USC, UCD, UCSB, SDSU, Caltech, NU, UWash UChicago, UIC, UHM, CSUSB, HPWREN, UMo, MSU, NYU, UNeb, UNC,UIUC, UTA/Texas Advanced Computing Center, FIU, KISTI, UVA, AIST**
- **CENIC, Pacific Wave/PNWGP, StarLight/MREN, The Quilt, Kinber, Great Plains Network, NYSERNet, LEARN, Open Science Grid**
- **Internet2, DOE ESnet, NCAR/UCAR and Wyoming Supercomputing Center**