



### NRP aims to address 2 challenges

- Provide an AI Education & Research Infrastructure to all of academia that academia can afford.
- Provide a "technology playground" where CS and Domain Science researchers can meet and work together to accelerate technology adoptions in light of the end of Moore's Law











# NRP brings CS R&D and Domain R&D onto the same platform

NRP blurs the lines between "testbed" and "production" CI

Create social cohesion to accelerate domain science adoption of new programming paradigms & architectures





### Algorithm x Hardware = Science

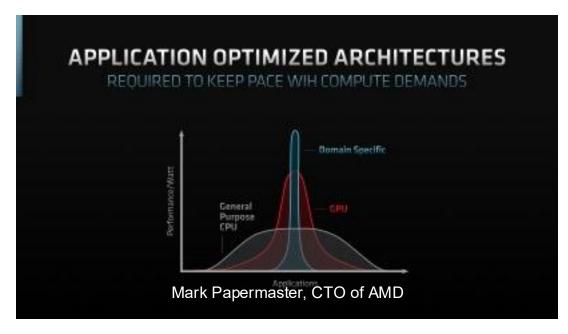
For decades, scientific progress was exponential in part because hardware performance per \$ increased exponentially due to Moore's Law.

Algorithmic performance did not scale nearly as fast, in most cases.

The end of Moore's Law threatens scientific progress.

#### "end of Moore's law" motivates new architectures





Performance improvements vs time slowed down by O(100)



NRP supports FPGAs, P4 switches, NVIDIA DPUs & DGXs

Committed to be a "Playground" of technologies, easily deployed & operated via BYOR and BYOD.



### **Advanced Technology Laboratory on NRP**

- Programmable computational capabilities emerged in devices of all kinds
  - Storage devices with embedded FPGAs => "Computational Storage"
  - GPUs on Network Interface Cards => "Data Flow Programming"
  - Programmable switches, down to individual ports => "Programmable Networks"
- We innovate nextGen systems in NRP to solve grand challenges of science
- Innovations made available to all of open science via our Open Infrastructure

Strategic Objective is to bring CS Research closer to Domain Research in the hope of decreasing time to adoption of new technologies & ideas

**NVIDIA BlueField DPU** 



P4 programmable switches

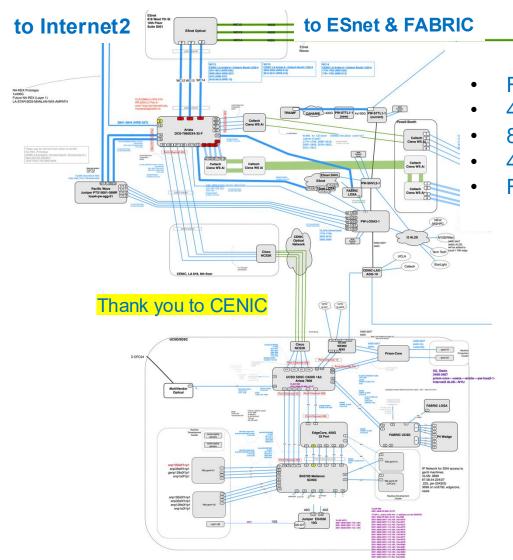








### **400G WAN Infrastructure**



#### **Infrastructure at SDSC:**

FPGAs: 32 U55C, 24 Bitware 520

400G P4 programmable switches

8 NVIDIA HGX w 8xA100 80G each

400TB of NVMe

FABRIC node

We own 400G capable nodes at MGHPCC, CERN, SDSC

We peer at 400G in LA with multiple networks via our 400G Arista switch

#### Real-Time In-Network Machine Learning on P4-Programmable FPGA SmartNICs with Fixed-Point Arithmetic and Taylor

Mohammad Firas Sada, John J. Graham, Mahidhar Tatineni, Dmitry Mishin, Thomas A. DeFanti, Frank Würthwein

As machine learning (ML) applications become integral to modern network operations, there is an increasing demand for network programmability that enables low-latency ML inference for tasks such as Quality of Service (QoS) prediction and anomaly detection in cybersecurity. ML models provide adaptability through dynamic weight adjustments, making Programming Protocol-independent Packet Processors (P4)-programmable FPGA SmartNICs an ideal platform for investigating In-Network Machine Learning (INML). These devices offer high-throughput, low-latency packet processing and can be dynamically reconfigured via the control plane, allowing for flexible integration of ML models directly at the network edge. This paper explores the application of the P4 programming paradigm to neural networks and regression models, where weights and biases are stored in control plane table lookups. This approach enables flexible programmability and efficient deployment of retrainable ML models at the network edge, independent of core infrastructure at the switch level.

Comments: To appear in Proceedings of the Practice and Experience in Advanced Research Computing

(PEARC25)

Subjects: **Distributed, Parallel, and Cluster Computing (cs.DC)**; Networking and Internet Architecture (cs.NI)

Cite as: arXiv:2507.00428 [cs.DC]

(or arXiv:2507.00428v1 [cs.DC] for this version) https://doi.org/10.48550/arXiv.2507.00428

Just one example use of our FPGAs, P4
Programming, and Al/ML.

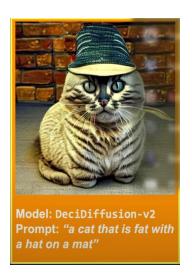


### Qualcomm Cloud AI 100 Ultra

- Multi-SoC PCle Cards
- 100B GenAl models on a single card
- Software tools
- 8 cards per server (2x4, w/ 4 on PCle switch)
- Fully programmable + C++ SDK
- Finetuning capabilities

- Models tested:
  - LLMs (varying sizes)
  - OpenAl Whisper (speech recognition model)
  - Stable Diffusion









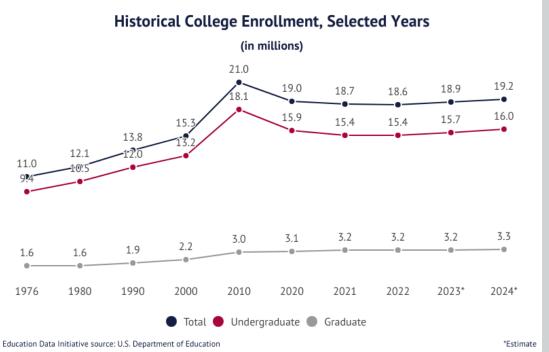


Work in progress: Understand the tokens/\$ and tokens/(sec\*power) performance for different models.

We are looking for more cost-effective inference hardware with which we can provide LLM service for a million students.

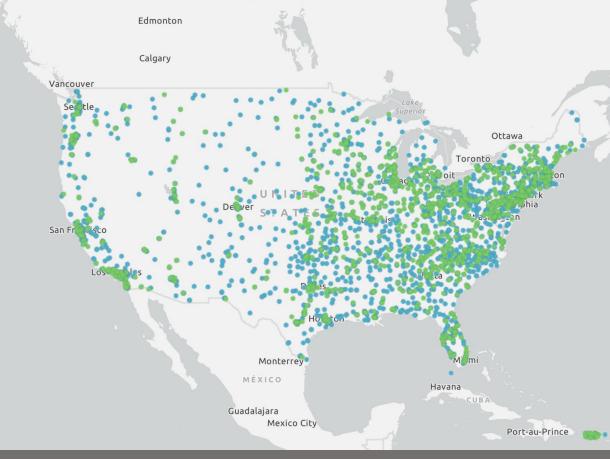






~20 Million students
 40% in community colleges
 43% in 4-year colleges
 17% in graduate/professional school

In the 2023-24 academic year, there were 1,632 public (blue) & 1,938 private non-profit (green) accredited degree granting colleges in the USA 3,580 non-profit colleges total





### Let's play with some numbers ...

- 20M students attend college in the USA
- 50% of students at San Diego State University want some sort of formal training in AI (survey result based on 8,000 student responses)
- 24% of UC San Diego undergraduates have used data & compute in their classroom in AY24

=> There are 5-10M college students in the USA that need data & compute resources for their classroom education



## Detailed look at UC San Diego



### Student Enrollment by Quarter

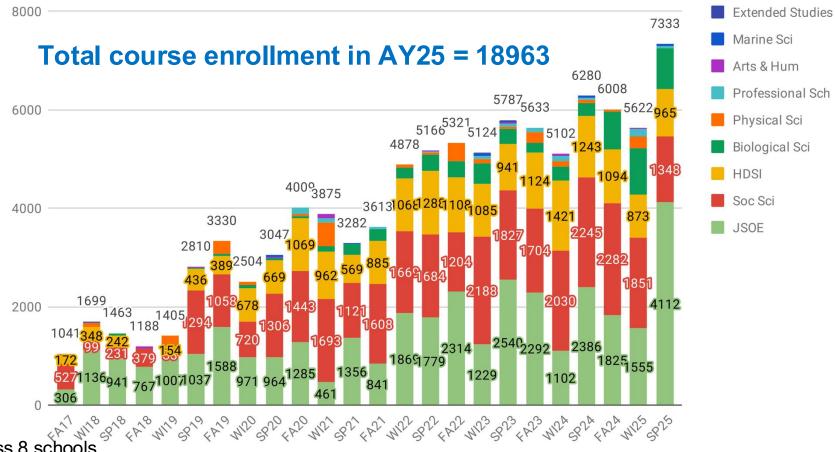
24% of all UCSD undergraduates used compute & data infrastructure in AY25 as part of their classroom education

Total in AY24	11092
UG Students	8538
Grad/Prof Students	2554

Some students take multiple courses/AY

34,955 undergraduates enrolled in Fall 2024





In AY25: 132 courses, 81 instructors, across 8 schools



### The beginning of a revolution in education

- Al redefines what's possible to teach in a classroom
  - Programming at the fingertips of every student
  - Simulations, data analytics, and visualizations on any topic
- ⇒ Students can experience vastly more realistic examples in the classroom to understand concepts & acquire knowledge
- AI makes personalized tutoring before, during, and after the classroom session possible.
- ⇒ More students understand more as a result of personalized attention, even in classes with hundreds or thousands of students

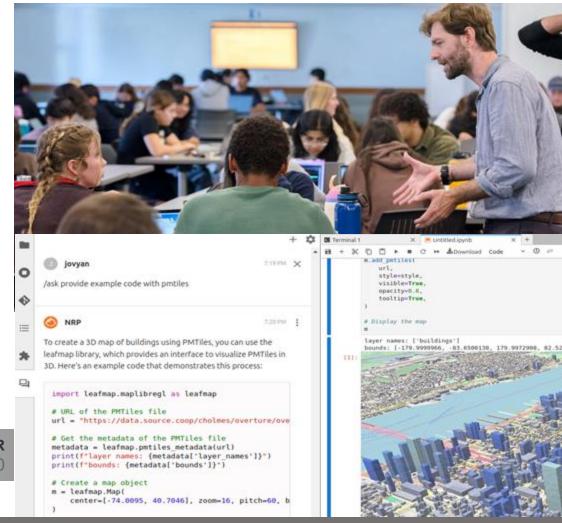


### Redefining what's possible

#### Coding with large language model

- 122 Students
- Active learning classroom
- Students visualize the Manhattan Skyline in 3D with the help of Al during class

CA 30x30 Planning & Assessment Prototype using LLMs on CENIC AIR https://huggingface.co/spaces/boettiger-lab/ca-30x30



# Operations is expensive ... unless it is amortized across large scales

- Of the 3,580 accredited, degree granting higher education non-profit institutions less than 200 are research intensive (R1), and have the scale to afford a group of sysadmins, cybersecurity, user support, ... professionals.
- ⇒ Less than 6% of colleges nationwide can afford AI infrastructure ops
- What if we aggregated the system administration, cybersecurity, and user support across 1,000 colleges?
- Colleges own their AI hardware, but we centralize the operations to benefit from economies of scale, just like the cloud providers.

National Research Platform (NRP), a prototype to achieve this vision



### NRP offers the community

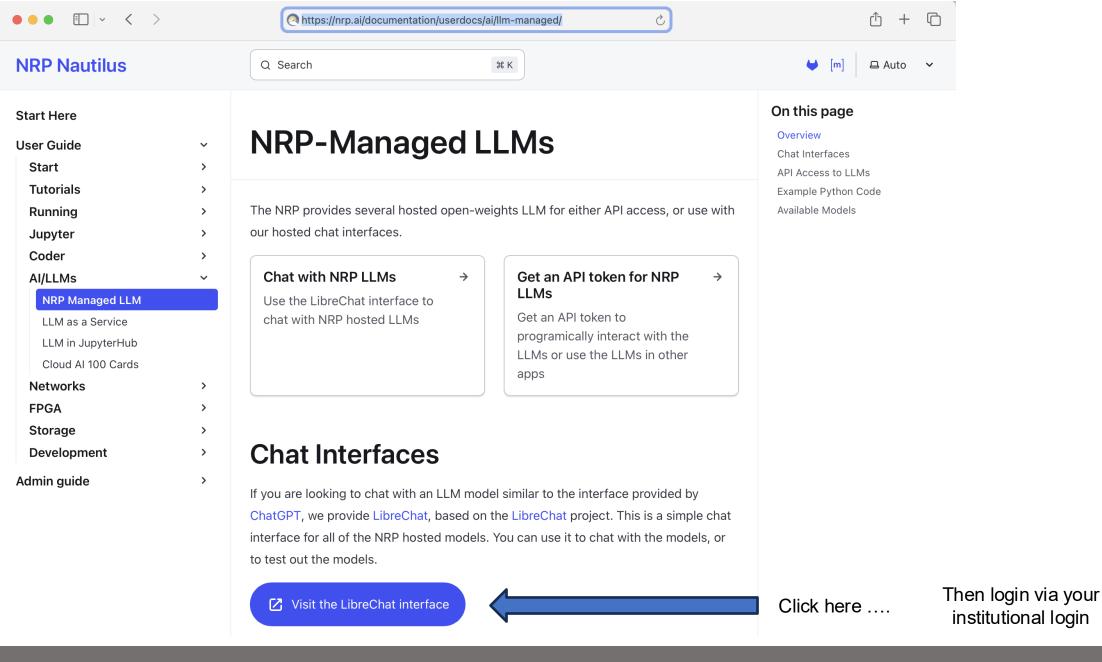
- To run your hardware from IPMI up
  - OS maintenance, security monitoring, ...
- Researchers, Educators, and students see a global scale Kubernetes cluster
  - While the cluster is shared, we restrict use of individual hardware to owners when necessary
  - Documentation and training for the community
- Lot's of software to reuse
  - E.g. we operate JupyterHub for the community & show people how to deploy & customize their own
    - 46 community-run JupyterHubs now on NRP
- Maintain lots of topical Chat channels for the community to learn from each other and interact with each other
  - Chat channels monitored by professionals, community is encouraged to interact & help each other
  - Starting to explore AI chatbots trained on the chats
- LLM as a service
  - We run most popular LLMs as services for all
  - ... and provide popular APIs to upload any AI model from Hugging Face onto our resources.



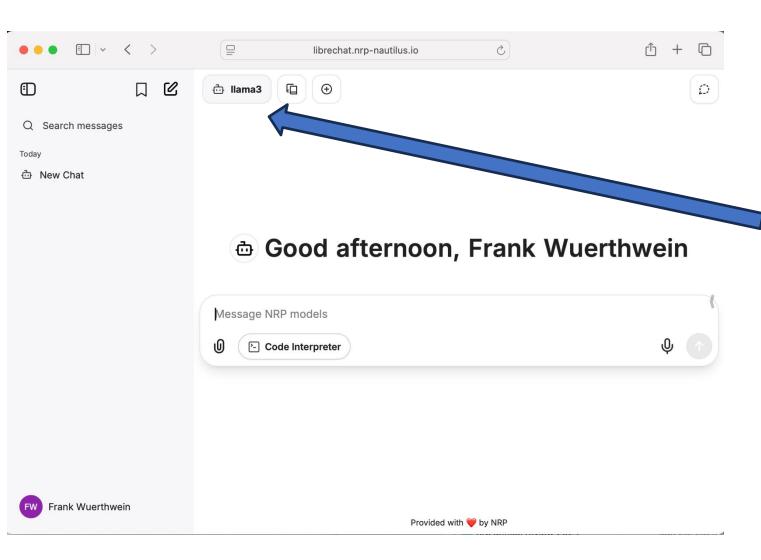
### Mental Picture of 3 types of "Al" Services

- Compute, Storage, Jupyter, Matrix
  - The initial set we started with
- LLM as a Service
  - Chat and API Access to popular LLMs that we run as service for the community
- Al workflows as a service => leverage National Data Platform
  - Imagine curated data, curated AI agents, curated RAG workflows, ... all available via a simple intuitive interface to compose custom AI services that fundamentally require workflow execution to instantiate the service.
  - Makes developing, instantiating, and operating your custom LLM easy

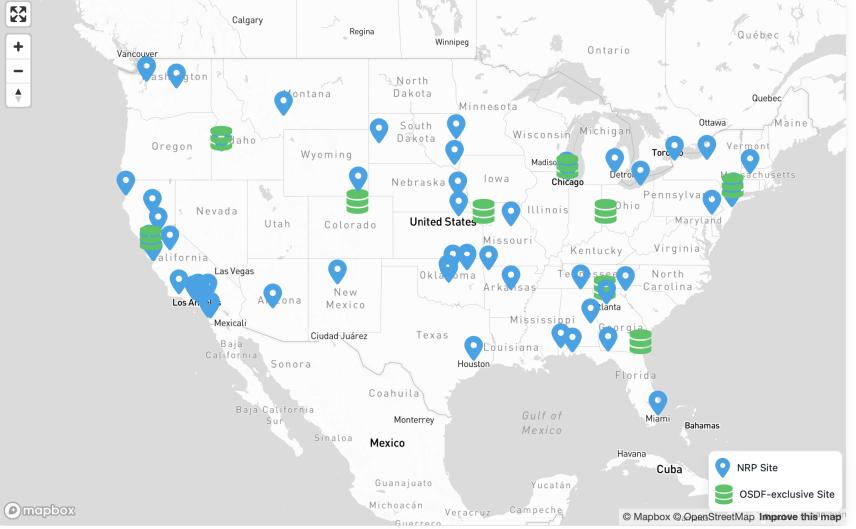








You can pick different models, and explore.



## NRPNATIONAL RESEARCH PLATFORM

The National Research Platform is a partnership of more than 50 institutions, led by researchers at UC San Diego, University of Nebraska-Lincoln, and Massachusetts Green High Performance Computing Center and includes contributions by the National Science Foundation, the Department of Energy, the Department of Defense, and many research universities and R&E networking organizations in the US and around the world.

Select a site or click on a site in the map

Sites

98
Sites hosting NRP nodes

CPU Cores

1,471
Total GPUs across all nodes

Nodes

CPU Cores

29,878
Total CPU cores across all nodes

In January we were at 72 locations and 22 PB of disk space

https://dash.nrp-nautilus.io

### NRP as of last week ...







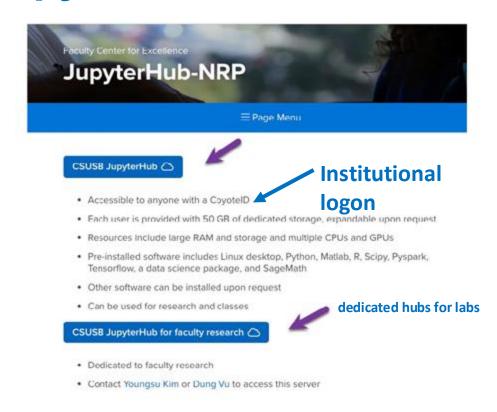
### **Example: CSU San Bernardino**

- < 1,100 faculty members
- < 19,000 students
- Serves 2 of CA's largest counties
- < Hispanic Serving Institution
- < 57% Pell Grant recipients
- < Many student oriented projects

**Non-R1 Institution** 



### JupyterHub User Interface at CSU, San Bernardino



#### Wide range of dedicated hubs:

#### Server Options

#### Advanced Options

Image	
0	Stack Minimal
0	Stack Datascience
0	Stack R-Studio, Vs-code for Dr. Becerra's class
0	Stack Desktop Apps - VS Code
O	Stack Desktop Apps - Pgadmin4
0	Stack Desktop Apps - Blender
0	Stack PySpark
0	Stack PyTorch2
0	Stack R-Studio
0	Stack R-Studio for BIOL-5050
0	Stack SageMath

https://csusb-metashape.nrp-nautilus.io: 3D modeling

https://csusb-vasp1.nrp-nautilus.io Viena Ab initio Simulation package ( VASP)

https://csusb-cousins-lab.nrp-nautilus.io: VASP simulation

https://csusb-becerra.nrp-nautilus.io AI/ML project

https://csusb-biol-5050.nrp-nautilus.io: Biology course

https://csusb-cse-salloum.nrp-nautilus.io Summer Research

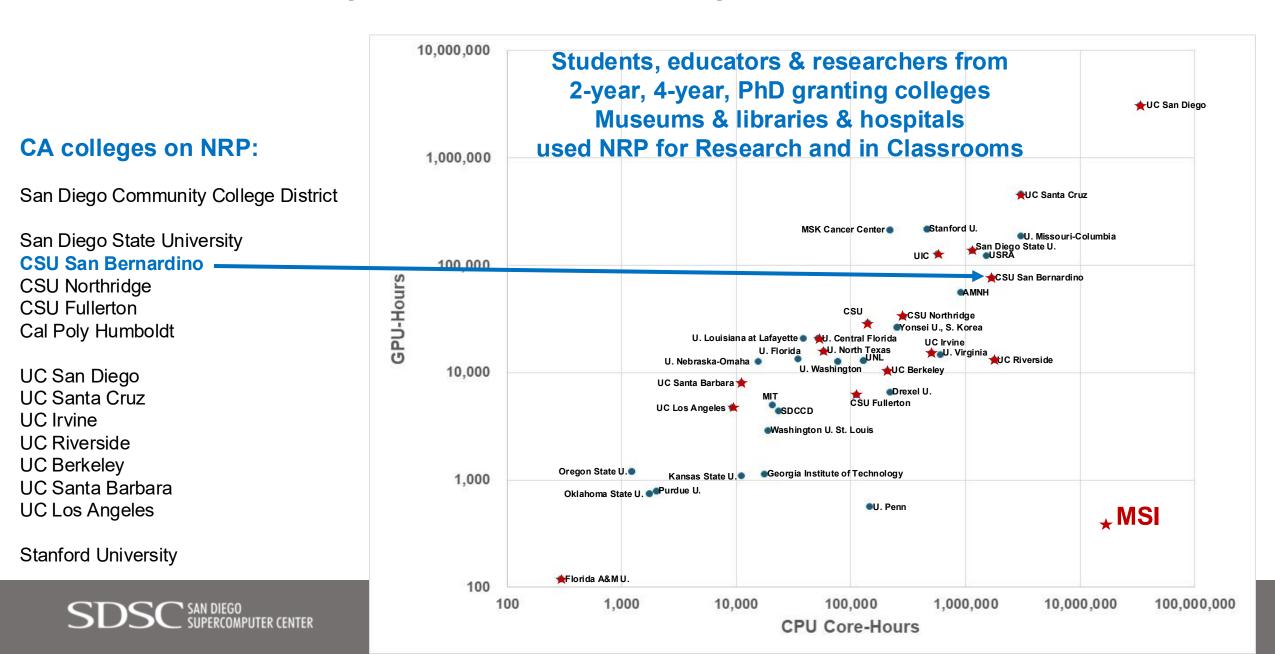
https://csusb-drhamoudahub.nrp-nautilus.io Data Analytics

https://csusb-ratnasingam.nrp-nautilus.io Data Analytics

https://csusb-zhang.nrp-nautilus.io AI/ML project



#### 63 Campuses had active namespaces in 2024 on NRP



#### 63 Campuses had active namespaces in 2024 on NRP

#### **CA** colleges on NRP:

#### San Diego Community College District

San Diego State University CSU San Bernardino

**CSU** Northridge

**CSU Fullerton** 

Cal Poly Humboldt

UC San Diego

**UC Santa Cruz** 

**UC** Irvine

**UC** Riverside

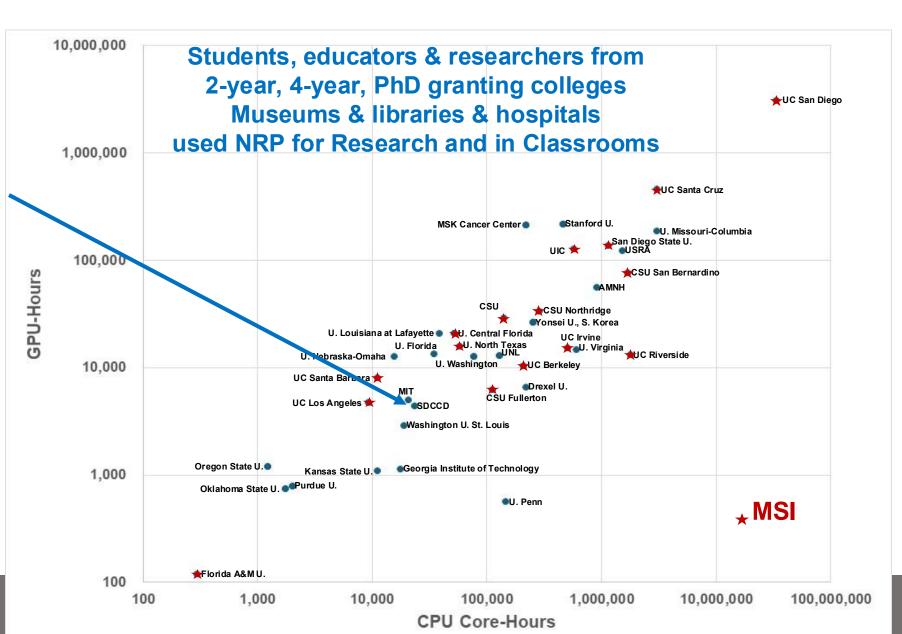
**UC Berkeley** 

UC Santa Barbara

**UC Los Angeles** 

Stanford University





# Community Colleges ... a crucial audience for NRP

- CCC system dominates the total college student population in CA
  - 70% of all of California's college students
- CCC transfer students to CSU & UC drive social mobility in CA
  - 30% of UC San Diego's incoming class are graduates from CCC
  - 50% of San Diego State Universities incoming class are graduates from CCC
- California has strong path from K-12 to CCC to UC/CSU
  - Public high schools in CA often have course offerings with their local CCC
  - High School students receive computing accounts from community colleges

NRP presently working with San Bernardino, San Diego, and Los Angeles Community College Districts



# Strategy for Scaling out

Engage with colleges from the network via the RENs and

from the curriculum via ADSA (Academic Data Science Alliance)



CENIC is a 501(c)(3) with the mission to advance education and research statewide by providing the world-class network essential for innovation, collaboration, and economic growth.

#### **Charter Associates:**

- California K-12 System
- California Community Colleges
- California State University System
- Stanford, Caltech, USC
- University of California System
- California Public Libraries
- Naval Postgraduate School

#### **Other Members Include:**

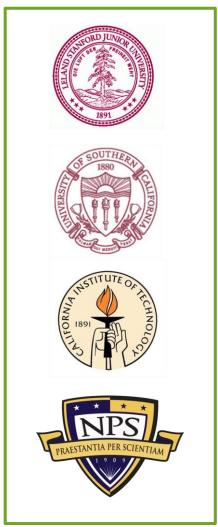
- Scientific and Cultural Institutions
- Private Colleges and Universities
- Hospitals and Specialized Medical Institutions
- Biomedical, Space, Environmental Research Organizations
- Tribal Nations

2M students enrolled ~10% of US college students









#### The Majority of GPUs in NRP Reside in the CENIC AI Resource (CENIC AIR)





### Research & Engagement w Industry

- Community Colleges and CSUs have a strong local focus
  - Serve local industries
- We aspire to create research capacity & engagement with local industries via their ownership of AI Infrastructure
- Two Examples:
  - Agriculture Technology => Instrumenting a Sonoma County Winery
  - Fusion Energy => facilitating collaboration between CSU San Bernardino,
     San Diego State University, and General Atomics



### Insta360 Camera & Tom on E-bike



## Path towards Sustainability



### We are in phase 1 of 3 phases

#### Phase 1:

- Establish Concept and scale out to 100 colleges and 10,000 students/year
- Today we operate hardware at 53 colleges
- Expecting to complete phase 1 in 2-3 years.

#### Phase 2:

- Scale out to 1,000 colleges and 1 Million students/year
- Scale out to 24x7 support for system & cybersecurity, but only 9-5 for user support & training
- Expecting this to be a 5-10 year program
- Estimating 10,000 GPUs needed to serve this population
- Phase 3: Sustainability through 501(c)(3) and membership fees
  - At 1 Million students per year, even a modest fee of \$10/(student\*year) sustains the program

Longer term: include high schools and public libraries



### **Summary & Conclusions**

- Opportunity is a requirement for Social Mobility
  - Education is the most effective guarantor of social mobility
  - Data & Compute is increasingly required for Education

We provide a Data & Compute Platform for Higher Ed that is cost-effective and aspire to make it available to all colleges in the USA

 We aspire to engineer transformative change by creating a long term sustainable AI Infrastructure that colleges nationwide can own together



### Acknowledgements

The NRP is supported by NSF grants OAC-1541349, OAC-1826967, OAC-2030508, OAC-1841530, OAC-2005369, OAC-21121167, CISE-1713149, CISE-2100237, CISE-2120019, & OAC-2112167

And by CENIC, Pacific Wave, MREN, GPN, NYSERNet, FLR, NEREN, SunCorridor, OARnet, SCLR, the Albuquerque GigaPoP, and Internet2

As well as a long list of Universities and colleges that host hardware, and share their hardware with the community.

