

# CUP-ECS Center Overview

PSAAP-III Annual Review

Prof. Patrick Bridges

September 29, 2022



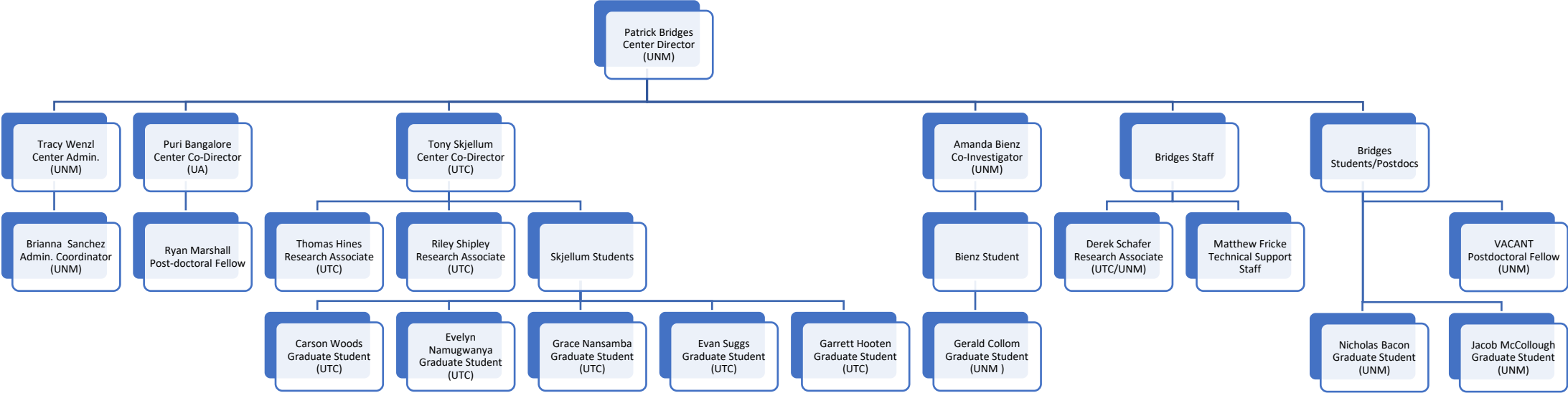
Center for Understandable, Performant Exascale Communication Systems



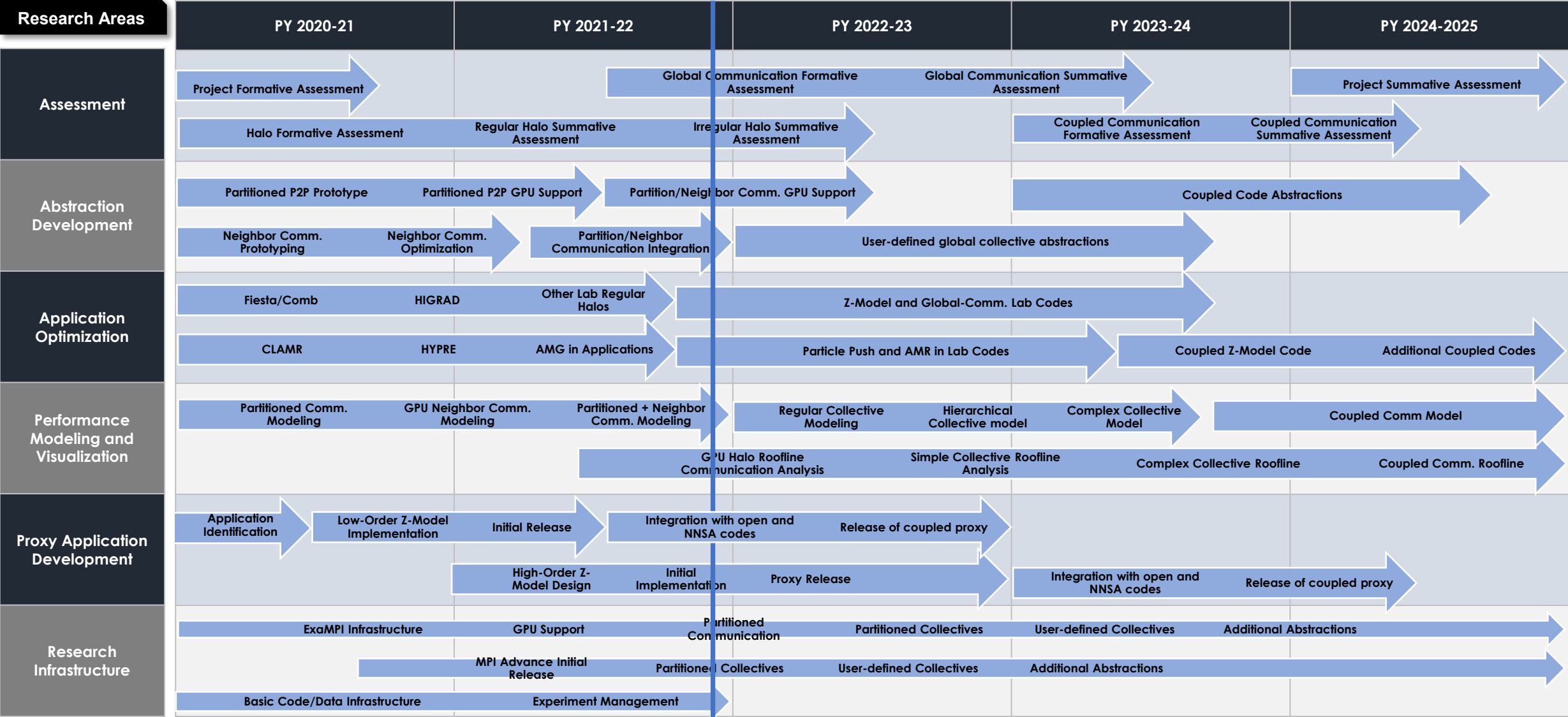
# Center Goals

- Mission: “Provide optimized, performance-transparent communication systems for NNSA exascale applications.”
- **Goal: Research, demonstrate and deploy better communication abstractions that make NNSA mission applications faster, more predictable, and easier to write**
- Approach
  1. Revisit and re-architect the relationship between exascale communication systems, applications, and hardware to support transformative scientific insights
  2. Research communication system innovations that accurately quantify, predict, abstract, and optimize exascale communication systems
  3. Develop and integrate enabling technologies and leverage these fundamental research advances in support of NNSA applications and systems
  4. Continuously refine research, development, and system integration based on feedback from NNSA collaborators and stakeholders.

# Center Personnel and Organizational Structure



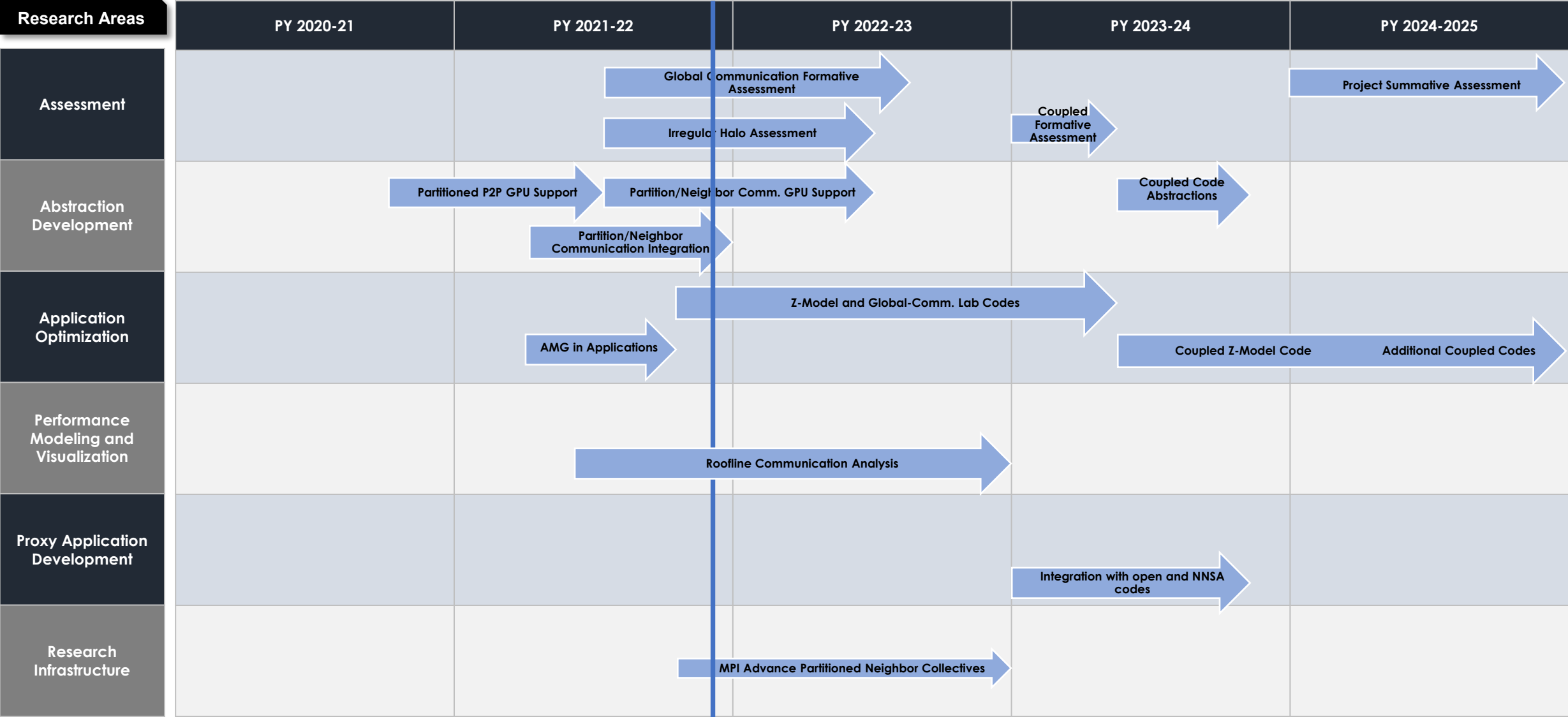
# 5-year Project Roadmap



# Project Risks and Mitigation

- Final Goal: Integration and assessment of optimized communication abstractions in benchmarks/applications, including coupled codes
- Non-Technical Risks: Student recruitment and retention
  - Student churn has slowed progress in some areas
  - Industry demand makes recruiting and retention difficult
  - Mitigation 1: Working on multiple recruitment paths to mitigate
  - Mitigation 2: Internships and staff placements to retain and builds center/lab ties
  - Mitigation 3: Educational materials to spin up new students quickly
- Technical Risks: Identify a critical path in the project plan to identify future risks and key challenges (from Year 1 Review) – next slides

# Technical Project Critical Path



# Technical Project Risks from Critical Path

- Assessment: Incomplete/inaccurate assessment of production codes due to lack of availability to all center personnel
- Abstractions: Robust, usable kernel and stream-triggered primitives
- Applications/Proxies: Lack of multi-physics benchmarks and codes
- Modeling: Identifying applications worth optimizing
- Infrastructure: Implementations that can be integrated in applications

# Addressing/Mitigating Technical Risk

- Assessment: Incomplete/inaccurate assessment
  - Engaging key staff and postdocs with production codes and frameworks (Currently xrage, HOSS, and Parthenon at LANL)
  - Creating benchmarks that can be calibrated from statistics gathered by these staff from production codes
  - Will need to broaden to include SNL/LLNL applications with help from lab staff and interns
- Abstractions: Current GPU triggering libraries (libmp) very fragile
  - Working directly to understand their limitations/tradeoffs
  - Engaging with NVIDIA for help running, optimizing, and debugging, some newer libraries may mitigate libmp problems
  - Need access/engagement with Cray and/or AMD in this area
- Application/Proxies: Lack of open multi-physics benchmarks and codes
  - Continued development of Z-Model proxy (Beatnik)
  - Searching for postdoc to work directly on optimization of a lab code and integration with Beatnik
- Modeling: Identifying applications worth optimizing
  - Leveraging other funding (NSF OAC Core) to fund a student on GPU Communication Roofline modeling and performance prediction
- Infrastructure: Implementations that can be integrated with production codes/systems: MPIAdvance



# For Reference: Challenges and Issues from Spring TST Meeting

- Continued student recruiting and retention issues - addressing as discussed on previous slides
- Hardware testbed acquisition delays – completed acquisition/standup
- Will need help with vendors on NDAs, hardware access - addressing as discussed on previous slides
- Continuation application delays – Year 3 continuation went off without a hitch

# Other Project Changes

- Prof. Abi Arabshahi (UTC) left project, David Walker from Cardiff joining project Nov. 1
- Derek Shafer moved to UNM from UTC
- Ryan Marshall being hired as postdoc at LANL, recruiting new postdoc to join project
- UNM using carry-forward to hire a postdoc, ad currently posted

# Center Meetings and Management

- Active slack channel for informal cross-team discussions
- Weekly Meetings (Zoom)
  - One center leadership meeting (PIs, key technical staff) on Zoom
  - Two agenda-run weekly technical planning meetings
  - Two open working weekly hackathon meetings
- Three to four in-person hackathons and/or leadership meetings
  - Two dedicated (summer, spring), one collocated with conference or other meeting (SC this year, or EuroMPI most years)

# Vendor Collaborations

- Working closely with NVIDIA/Mellanox on libmp GPU triggering issues
- Successfully using nsys for GPU/NIC communication measurement (but only have NIC metrics on IB systems)
- Need to work with AMD/Intel/Cray on GPU triggering on their systems
- Need to identify contacts for access to data plane/infrastructure processor units (emerging data center NIC offload platforms)

# Lab Interactions

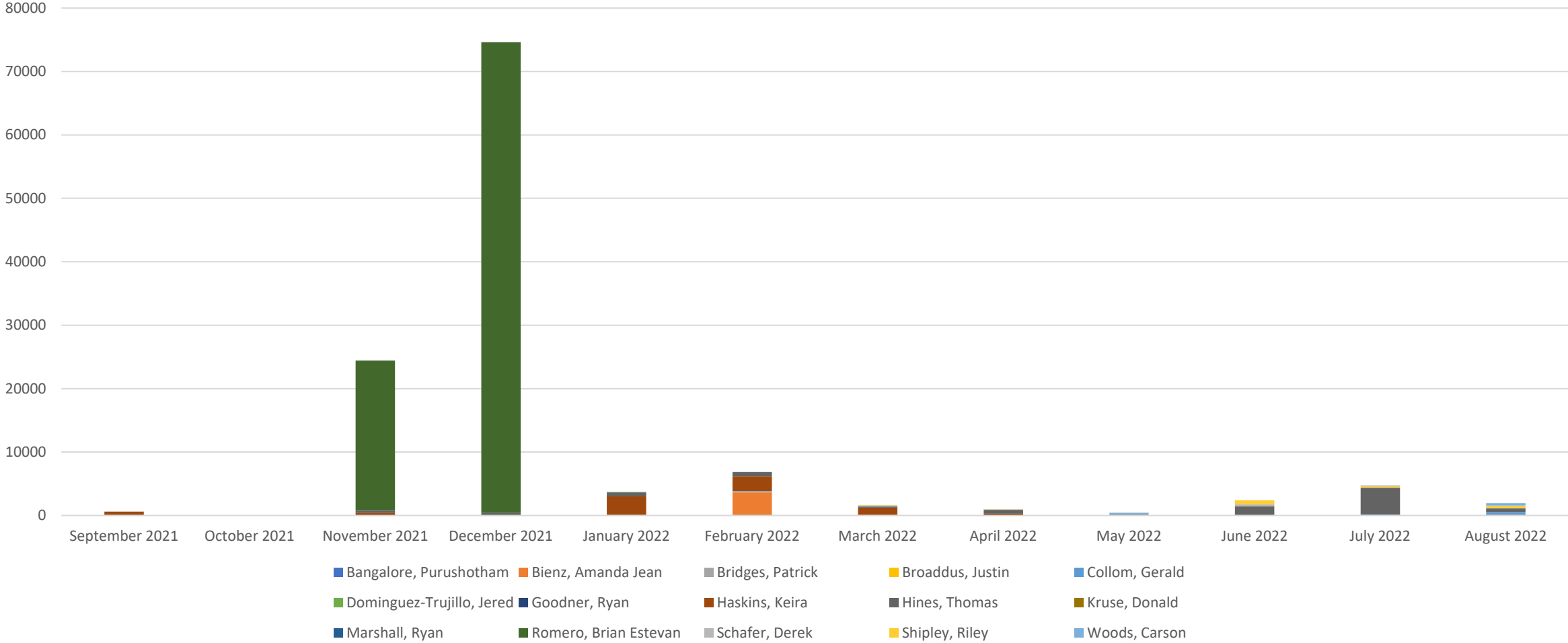
- Lab staff member participation in ongoing colloquium series
- Regular (i.e. weekly or bi-weekly) meetings with LANL, LLNL, and Sandia staff members by center leadership and students
- Lab internships of both PSAAP students and other students at participating institutions
  - Gerald Collom (UNM)
  - Keira Haskins (UNM)
  - Garrett Hooten (UTC)
  - Carson Woods (UTC)
- Center personnel joining NNSA labs as staff or postdoctoral fellows
  - Jared Dominguez-Trujillo (UNM), Technical Staff, LANL CCS
  - Ryan Marshall (UA), Postdoc, LANL XCP
  - Garrett Hooten (UTC), LLNL Staff position Jan. 2023
- Additional collaborations by center personnel with laboratory personnel on other awards/contracts



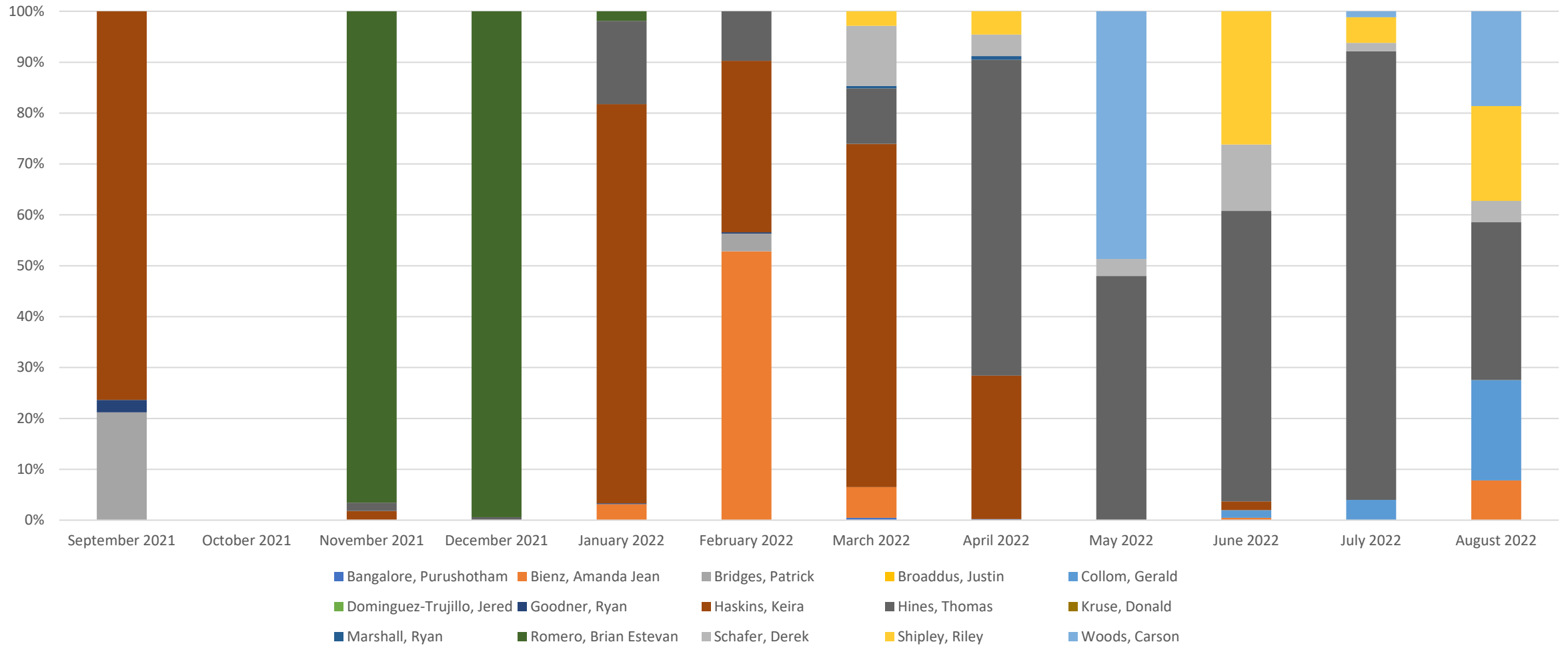
# Testbed Usage

- NVIDIA Testing
  - Frequent performance testing on Lassen, particularly for stream and kernel-triggered communication
  - Lassen environment has made non-vendor MPI testing challenging
  - Using local PSAAP-funded A100 test systems for fast turnaround R&D
  - Transitioning to LANL Chicoma for larger-scale NVIDIA testing
  - Dove-tails with separately-funded/allocated LANL production test runs
- AMD Testing
  - Looking forward to running on Tioga for AMD testing
  - Considering acquisition of local AMD test systems
  - Software support for kernel/stream triggering on AMD systems will be key

### Lassen Usage by Researcher/Month



### Lassen Usage Breakdown by Researcher





# Key Annual Review Recommendations

- Continue the Colloquium series, including inviting Tri-Lab presenters
- Quarterly virtual hackathons even when in-person hackathons are not possible
- Creation of a class on HPC networking and other outreach
- Arrange stronger vendor collaborations similar to ECP academic participants – already discussed
- Look broader than irregular AMR, for example unstructured meshes in HOSS
- Identify a critical path in the project plan to identify future risks and key challenges – already discussed

# Colloquium Series

- This year's speakers
  - [The challenge of sparsity in HPC codes](#) - Galen Shipman, Los Alamos National Laboratory
  - [MPI from the Ground Up: From Operations to Implementations, Part I](#) - Derek Schafer & Tony Skjellum, University of Tennessee at Chattanooga
  - [Revisiting a Classic: Bruck algorithm for non-uniform all-to-all communication](#) - Dr. Sidharth Kumar, University of Alabama at Birmingham
  - [GPU Integrated Communication](#) - Jim Dinan, NVIDIA
  - [The Kokkos Performance Portability Programming Model](#) - Christian Trott, Sandia National Laboratories
  - [PETSc and its communication module PetscSF](#) - Junchao Zhang, Argonne National Laboratory
- Most spring speakers deferred due to scheduling, plan to restart mid-fall

# Hackathon Status

- Two in-person hackathons in the past year:
  - February 8/9 - Focused on sharing knowledge and expertise monitoring application communication costs
  - July 19/20/21
    - Training new students in use and measurement of neighbor communication
    - Analysis/discussion of GPU communication options
    - Planning session on education materials/contents
- Weekly meeting schedule includes two regular short (2-hour) group hacking sessions with student and faculty mentoring

# Education outreach

- HPC Networking Class Development
  - Initial presentation as part of fall colloquia on HPC networking (background for center students and other personnel)
  - Assembling reading list and relevant articles/topics for same reason
  - Starting to plan for initial HPC networking special topics class – discussed at July hackathon
- Planning tutorials targeting lab and other HPC professionals on specific comm. topics e.g. partitioned communication
- Developing benchmarks/examples to use as materials in presentations, tutorials, and classes
- One focus on David Walker on project starting Nov. 1

# Look more broadly that irregular AMR

- Postdoc Ryan Marshall began work at LANL on HOSS in spring 2022, currently transitioning to LANL postdoc position
- Working on a survey paper on irregular communication abstractions on DOE applications
  - Starting with linear solver frameworks (HYPRE/Trilinos/PETSc)
  - Identifying additional applications, frameworks, and algorithms to study
- Need better benchmarks – ECP miniapps that are designed to support irregular meshes (MiniAero, LULESH) actually use regular meshes
  - Staff member Derek Schafer (joint funding with PSAAP and a LANL contract) working to irregular communication statistics from xRage and Parthenon

# Year 2 Research Areas and Directions

- Abstractions and Summative Assessment
  - Complete integration and evaluation of partitioned communication primitives in Comb proxy and Fiesta/HIGRAD applications
  - Begin integration and evaluation of optimized irregular neighbor collectives in HYPRE
  - Global communication formative – particle codes (EMPIRE, etc.), fast Fourier and fast multipole methods (Z-Model, FlecSPH, etc.)
- New Abstraction Development
  - Evaluation and optimization of prototype GPU partitioned communication in Comb proxy application
  - Design of partitioned neighbor collective abstraction as a general optimized halo exchange communication mechanism
- Fluid Interface Proxy
  - Design and begin implementation of parallel version of higher-order fluid interface model for use as a stand-alone proxy
- Research Infrastructure
  - Initial release of MPI Advance library with example usage in DOE applications
  - Design of general communication performance assessment experiment management system

# Center-supported publications since prior meeting

1. W. P. Marts, M. G. F. Dosanjh, S. Levy, W. Schonbein, R. E. Grant and P. G. Bridges, "MiniMod: A Modular Miniapplication Benchmarking Framework for HPC," *2021 IEEE International Conference on Cluster Computing (CLUSTER)*, 2021, pp. 12-22, doi: 10.1109/Cluster48925.2021.00028.
2. A. Bienz, L. N. Olson, W. D. Gropp and S. Lockhart, "Modeling Data Movement Performance on Heterogeneous Architectures," *2021 IEEE High Performance Extreme Computing Conference (HPEC)*, 2021, pp. 1-7, doi: 10.1109/HPEC49654.2021.9622742.
3. S. Ghosh, et al., "Towards Modern C++ Language Support for MPI," in *2021 Workshop on Exascale MPI (ExaMPI)*, St. Louis, MO, USA, 2021 pp. 27-35. doi: 10.1109/ExaMPI54564.2021.00009
4. D. Holmes, et al., "Partitioned Collective Communication," in *2021 Workshop on Exascale MPI (ExaMPI)*, St. Louis, MO, USA, 2021 pp. 9-17. doi: 10.1109/ExaMPI54564.2021.00007
5. D. Schafer, T. Hines, E. D. Suggs, M. Rufenacht and A. Skjellum, "Overlapping Communication and Computation with ExaMPI's Strong Progress and Modern C++ Design," *2021 Workshop on Exascale MPI (ExaMPI)*, 2021, pp. 18-26, doi: 10.1109/ExaMPI54564.2021.00008.
6. M. G.F. Dosanjh, A. Worley, D. Schafer, P. Soundararajan, S. Ghafoor, A. Skjellum, P. V. Bangalore, R. E. Grant, Implementation and evaluation of MPI 4.0 partitioned communication libraries, *Parallel Computing*, Volume 108, 2021, <https://doi.org/10.1016/j.parco.2021.102827>.
7. B. E. Romero, S. Poroseva, P. Vorobieff, and J. Reisner. "Comparison of 2D and 3D Simulations of a Shock Accelerated Inclined Gas Column." *APS Division of Fluid Dynamics Meeting Abstracts*, pp. P10-012. November, 2021.
8. B. E. Romero, S. Poroseva, P. Vorobieff, and J. Reisner. "Three-Dimensional Simulations of a Shock-Gas Column Interaction." *AIAA SCITECH 2022 Forum*, p. 1072. 2022. <https://doi.org/10.2514/6.2022-1072>
9. P. Haghi, Guo, A, Xiong, Q, et al. Reconfigurable switches for high performance and flexible MPI collectives. *Concurrency Computat Pract Exper*. 2022; 34( 6):e6769. doi:10.1002/cpe.6769
10. B. E. Romero, "FIESTA and Shock-Driven Flows." PhD diss. University of New Mexico, 2022
11. A. Bienz, S. Gautam and A. Kharel. "A Locality-Aware Bruck Allgather." *Proceedings of EuroMPI/USA 2022*. 2022.

By institution: UNM: 6 (3 primarily funded elsewhere); UTC: 3, UA: 2



Center for Understandable, Performant Exascale Communication Systems



# Year 3 Milestones

1. Formative assessment of irregular communication demands in DOE application, including but not limited to the LANL HOSS application
2. Submission of partitioned collective abstraction specification to MPI forum for future inclusion in MPI standard and revision based on community feedback.
3. Design and initial implementation of GPU-triggered neighbor collective abstractions in MPI Advance
4. Release of higher-order fluid interface model benchmark specification, implementation, and initial performance results.
5. Summative assessment of optimized performance of different GPU halo communication approaches in DOE benchmarks and applications.



# Research Activities and Accomplishments

PSAAP-III Annual Review

Prof. Patrick Bridges

September 29, 2022



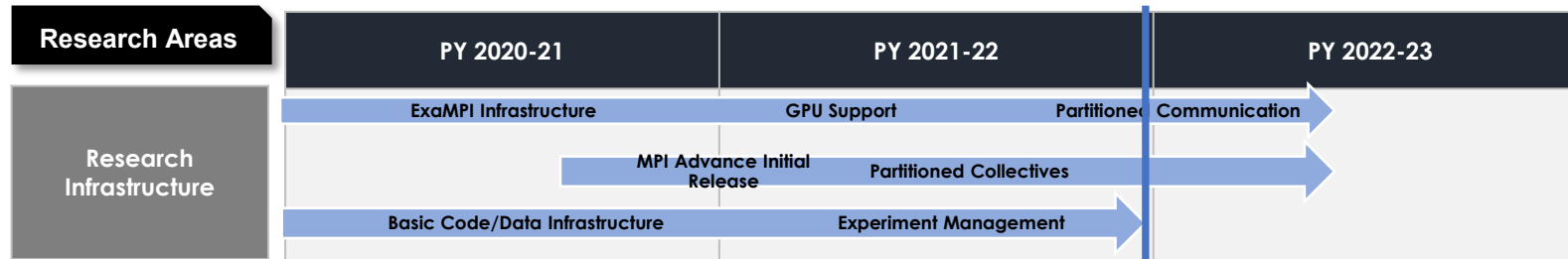
Center for Understandable, Performant Exascale Communication Systems



# High-Level Research Themes

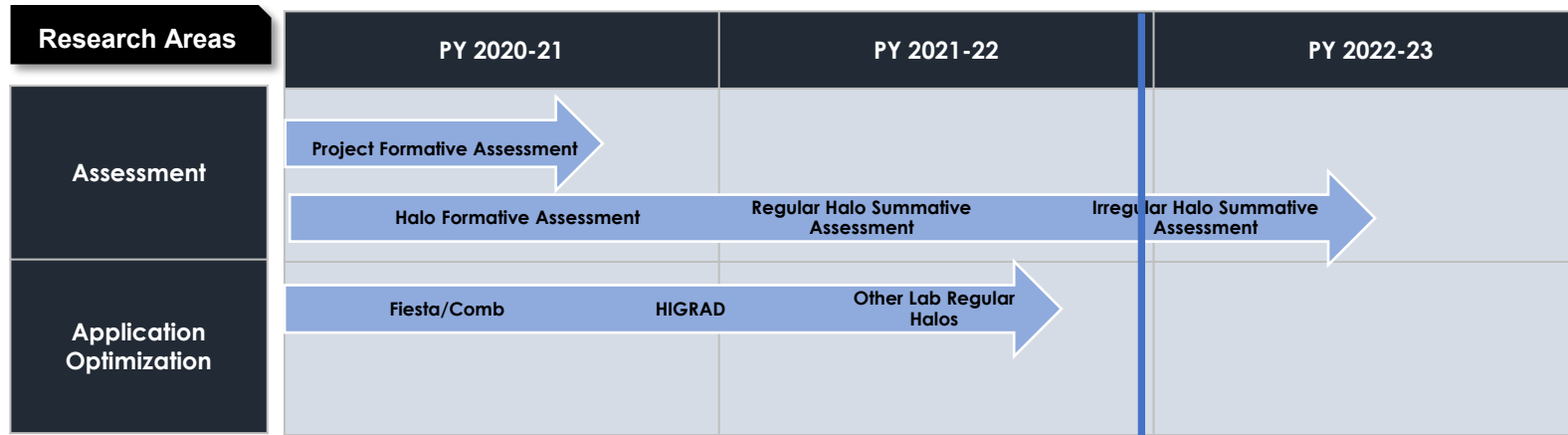
- Building research infrastructure to test, measure, and deploy new abstractions and implementations – e.g., MPI Advance and ExaMPI
- Exploring the many different fine-grain communication tradeoffs on GPU hardware and why current communication optimizations (e.g., in MPI) fail
- Quantifying the opportunities for regular halo optimization, using both point-to-point and collectives, versus a realistic baseline
- Examining same opportunities in different irregular communication benchmarks and applications
- Creating benchmarks to capture the communication behavior of more realistic applications
- Planning and initial outlining of educational materials for teaching students and professionals (1) modern HPC programming techniques and (2) HPC communication system basics
- Designing new communication abstractions to address the problems identified above for deployment through MPI Advance and ExaMPI

# Research Infrastructure



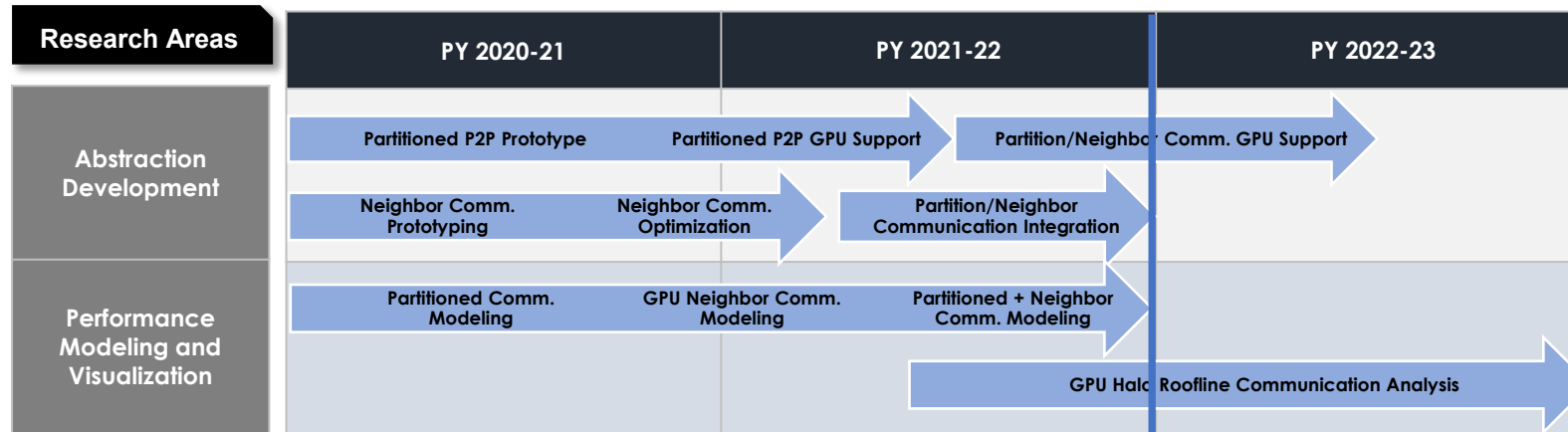
- MPI Advance
  - First release publicly available on GitHub – <https://github.com/mpi-advance>
  - Includes partitioned communication, locality-optimized collectives and neighbor collectives
  - Publicized and used in tutorial at EuroMPI, will be publicizing at SC and other upcoming venues
- ExaMPI
  - Compatibility features – Initial GPU-aware communication support added, ROMIO support in progress
  - Added Caliper tracing into core ExaMPI code paths for fine-grain measurement down to comm. fabric
- Data Management system
  - Overall design finalized and paper submitted for publication.
  - Additional related research on reproducibility approaches accepted to SC'22 workshop (CANOPIE-HPC)

# Regular Halo Assessment



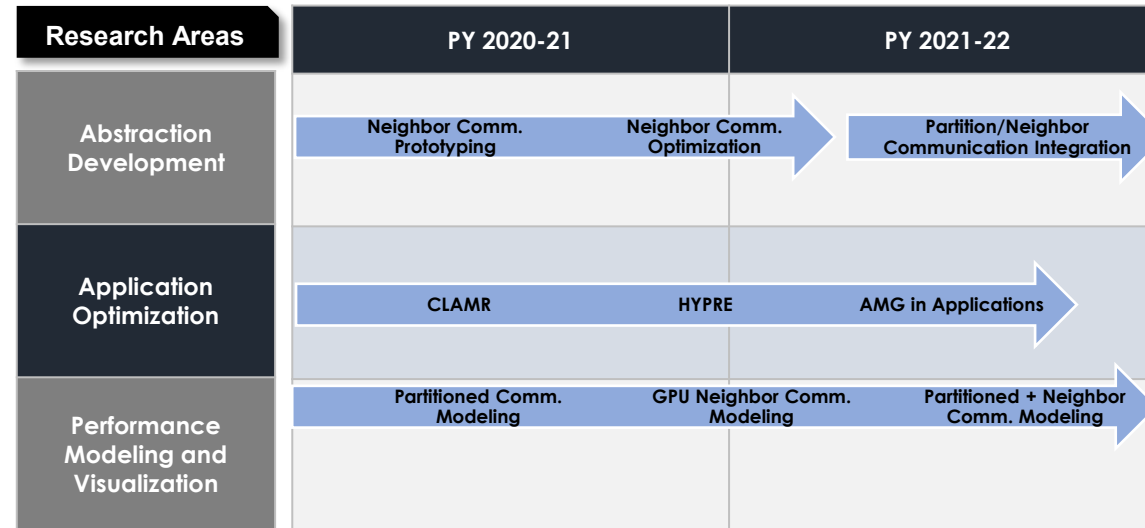
- Research question: how much room is there for aggressive optimization?
  - Baseline: Optimized HIGRAD to use GPU-aware communication, hand fused packing loops, supplemented with regular halo ping pong benchmark to study fine-grain costs
  - Optimizations: Considered partitioned pack/communicate and neighbor aggregation
  - Result: HIGRAD communication overheads reduced from 20% to 4%
  - Implications on partitioned communication for overlapping packing and sending, overheads of MPI datatypes, and message aggregation for regular halo exchanges

# GPUs and Point-to-Point Communication



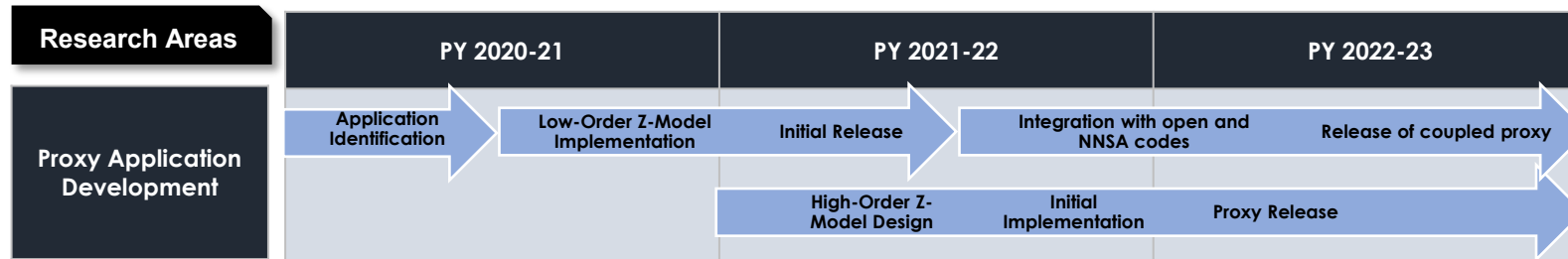
- Exploring the many different fine-grain communication tradeoffs on GPU hardware and why current communication optimizations (e.g. in MPI) fail
  - Partitioned P2P GPU abstractions near-finalized in MPI forum, partitioned collectives in progress
  - Fine-grain measurement of GPU communication performance characteristics using NVIDIA libmp.
  - Full implementation of GPU stream and kernel-triggered partitioned point-to-point in progress
  - Roadblock: Available stream/kernel triggering library brittle and fails in many cases (where NVIDIA says it shouldn't).
  - Mitigation: Working with NVIDIA to debug, examining alternative communication abstractions due to library limitations

# Irregular Communication



- Message aggregation powerful with irregular data layouts and communication
  - Tradeoffs with irregular halos due very different due to significant, harder-to-optimize packing costs
  - Implemented and tested optimized persistent neighbor collectives in HYPRE
  - Beginning evaluation of optimizing MPI\_Alltoallv in Fast Fourier Transforms (e.g. HeFFTe)
- Early-access to MPI and MPIX aggregating interfaces available via MPI Advance
- Characterizing irregular communication in multiple applications and benchmarks

# Assessment and Benchmark Development



- Developing novel approach for predicting application performance based on changing communication performance
- Implemented Beatnik, a fluid interface solver in Cabana/Kokkos/MPI
  - Full distributed low-order solve, brute-force single node high-order solve
  - Source release coming this fall via github (<https://github.com/CUP-ECS/beatnik>)
  - Successive versions will push on FFT, particle sort, and tree-code approaches
  - Looking at scalable high-order solve options in collaboration with LANL (e.g., Cabana, Parthenon, or FLECSi) and Stanford PSAAP (Legion/Regent)

# Education Efforts

- Developing training materials with a first focus of bringing new PSAAP students up to speed
  - Videos and online materials on basic workflows
  - Technical paper reading lists
  - Occasional center colloquium on core MPI topics
- Outlined a basic set of course topics on MPI semantics, layers, and basic implementation issues (e.g., key parameters)
- Potentially funding UTC research faculty for material development
- Goal: UNM/UTC course in the spring, tutorials at conferences next fall.