



Welcome to the second  
Mochi boot camp!



## Vision

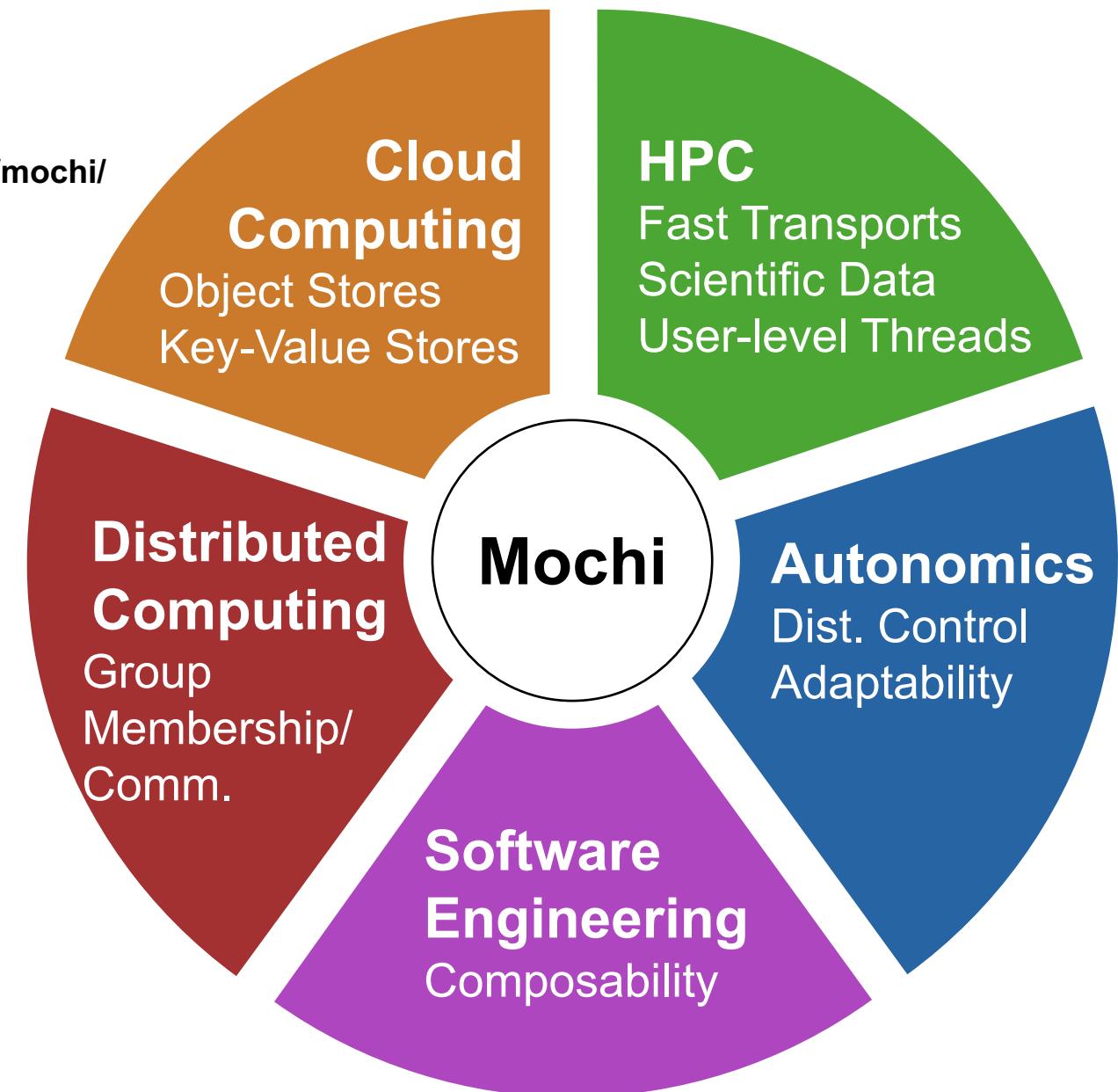
**Lowering the barriers to distributed services in computational science.**

## Approach

- Familiar models (key/value, object, file)
- Easy to build, adapt, and deploy
- Lightweight, user-space components
- Modern hardware support

## Impact

- Better, more capable services for specific use cases on high-end platforms
- Significant code reuse
- Ecosystem for service development



# Mochi: What are we trying to accomplish?

**We're trying to transform HPC data services from a monoculture to an ecosystem.**

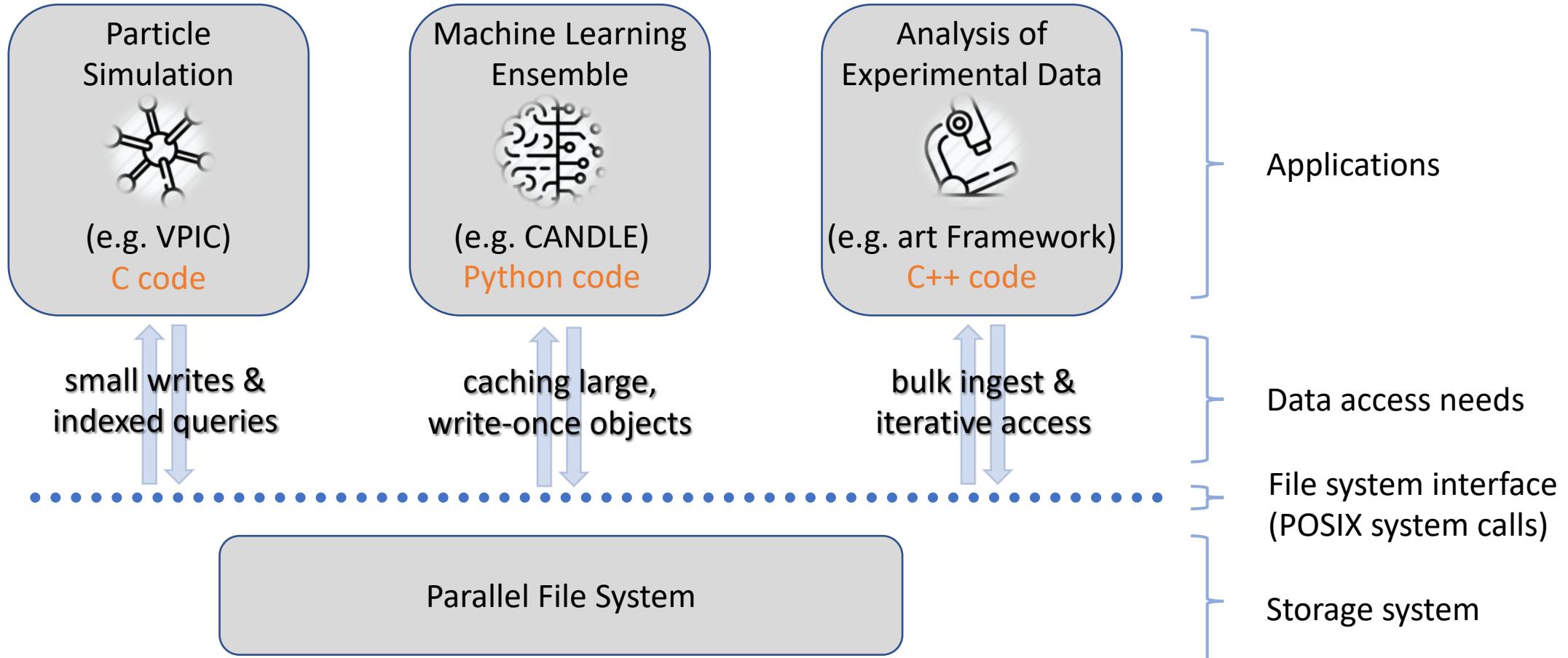
- Redefining how teams design and develop distributed services for use in HPC systems.
- Providing a portable "programming model" for these services.
- Providing a set of core building blocks.
- Demonstrating the methodology and tools with DOE science use cases.

**We're trying to foster a community of service developers.**

- Developing a set of training materials that will help others employ the tools.
- Making all these building blocks available to the larger community.

# How is this traditionally done in HPC?

## File system monoculture for data (dis)service



*All applications use the same “one size fits all” file system interfaces, semantics, and policies for data access.*

# How is this traditionally done in HPC?

The internals of a parallel file system (PVFS2 example)

1. All applications bend to the same service: single data model (string of bytes), decomposition (blocks), and consistency.

PVFS2 Client

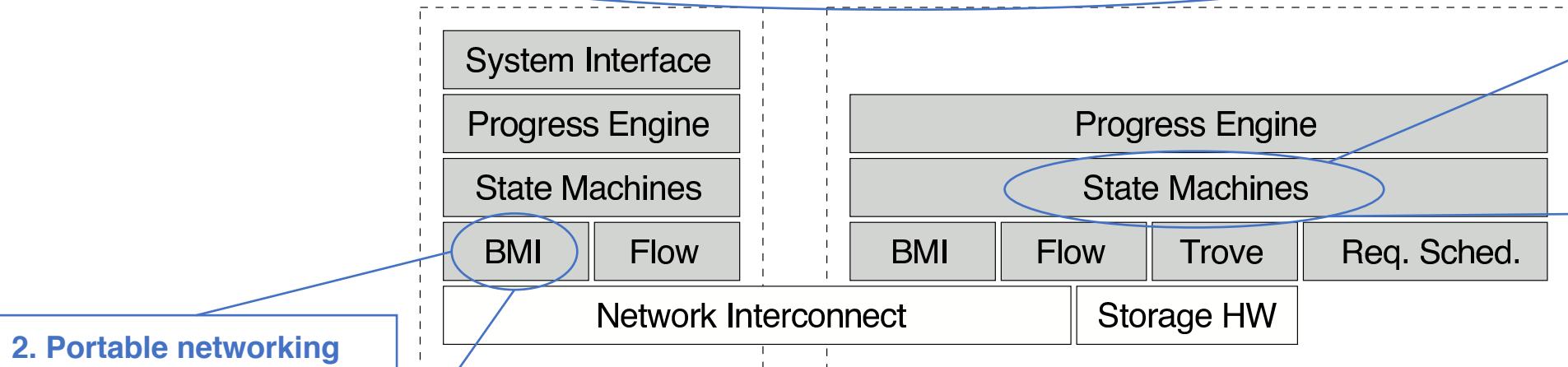
PVFS2 Server

2. Portable networking layer (BMI) is the only component usable separate from PVFS.

Service developers build RPC abstraction on top of it.

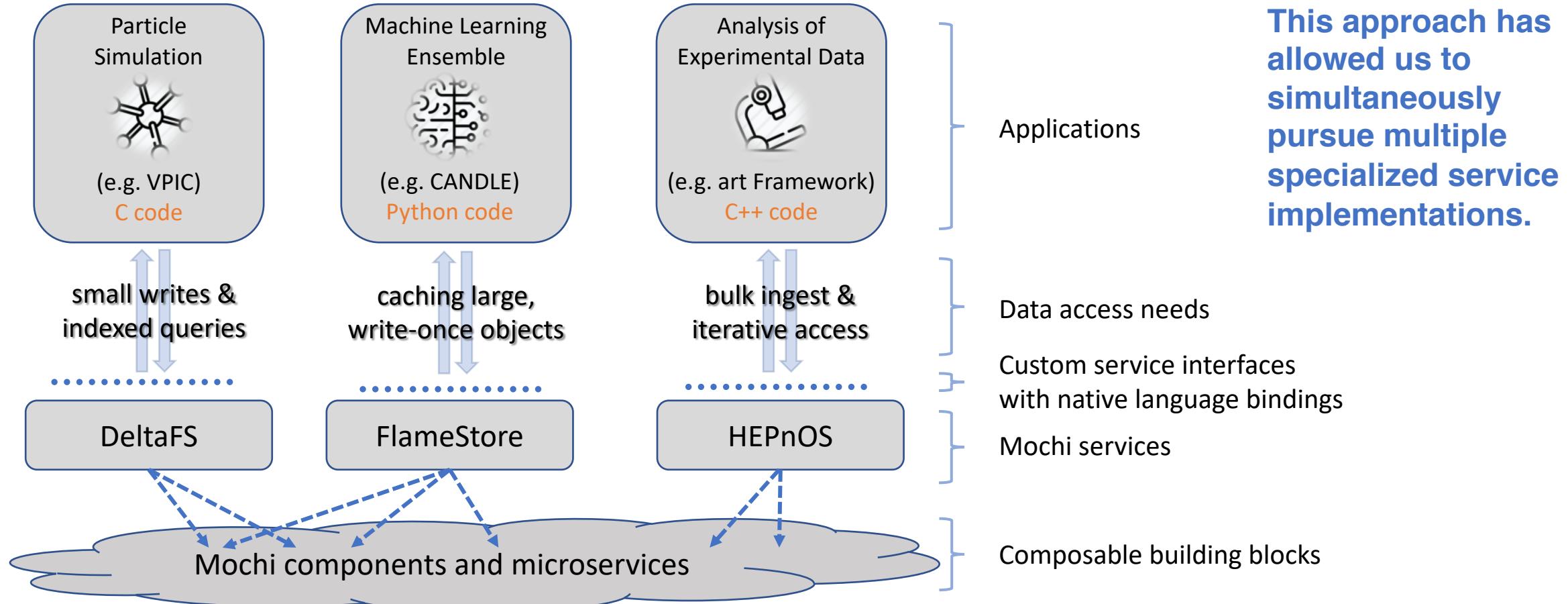
3. Programming is accomplished with proprietary state machine language. Effective but steep learning curve.

(Many) other researchers employ PVFS in their work, but it is difficult to build a new service from it.



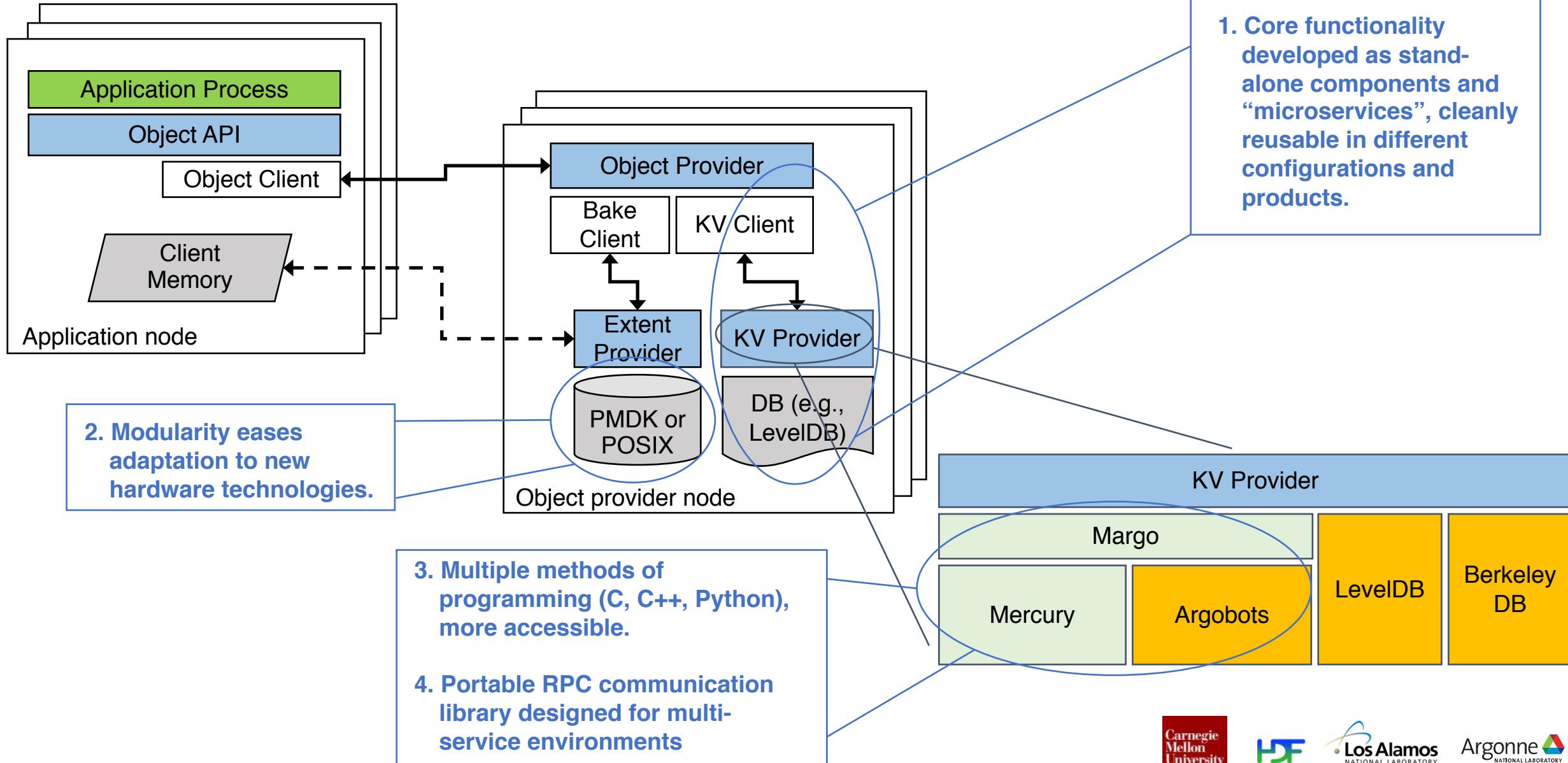
# What's new in the Mochi approach?

## An ecosystem of services co-existing and reusing functionality



*Instead of “one size fits all”, Mochi data services present tailored interfaces, semantics, and policies for data access while still leveraging robust building blocks.*

# What's new in the Mochi approach?



# Example Mochi Composed Services

## Particle Trajectory Assembly (DeltaFS)

### Goals

- Extreme scale file system metadata
- In-situ indexing with fast file retrieval

### Features

- Specialized directories to efficiently support trillions of files in a single directory
- Software routing to scalably manage connections
- POSIX(-like) API

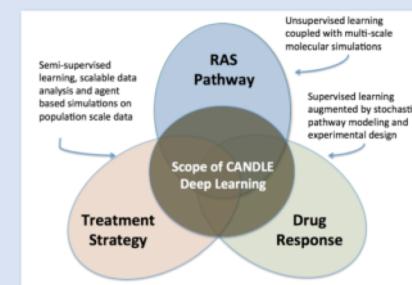
## Transient Storage for Deep Neural Networks (FlameStore)

### Goals

- Store deep neural network models during a deep learning workflow
- Retain most promising candidate models

### Features

- Flat namespace
- Python API (Keras models)



## Fast Event-Store for High-Energy Physics (HEPnOS)

### Goals

- Manage physics event data through multiple analysis phases
- Retain data in the system to accelerate analysis

### Features

- Write-once, read-many
- Hierarchical namespace (datasets, runs, subruns)
- C++ API



# Mochi: Core Components

A programming environment for distributed services

## Mercury

Portable RPC comm.

- Support for shared memory and multiple HPC transports: IB, GNU, PSM2, ...
- RDMA for bulk data movement
- Busy wait or wake on network event

<https://mercury-hpc.github.io/>



## Argobots

Lightweight runtime for concurrent execution.

- Utilize HW and OS constructs for performance
- hwloc-aware
- Custom scheduling and placement

<http://www.argobots.org/>

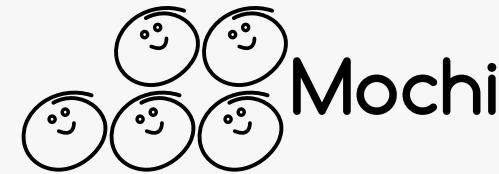


## Margo

Simple service development.

- Multi-threaded model
- Lightweight threads created to handle RPCs
- Express Mercury operations as blocking functions
- Uses Argobots to manage concurrency

<https://xgitlab.cels.anl.gov/sds/margo>

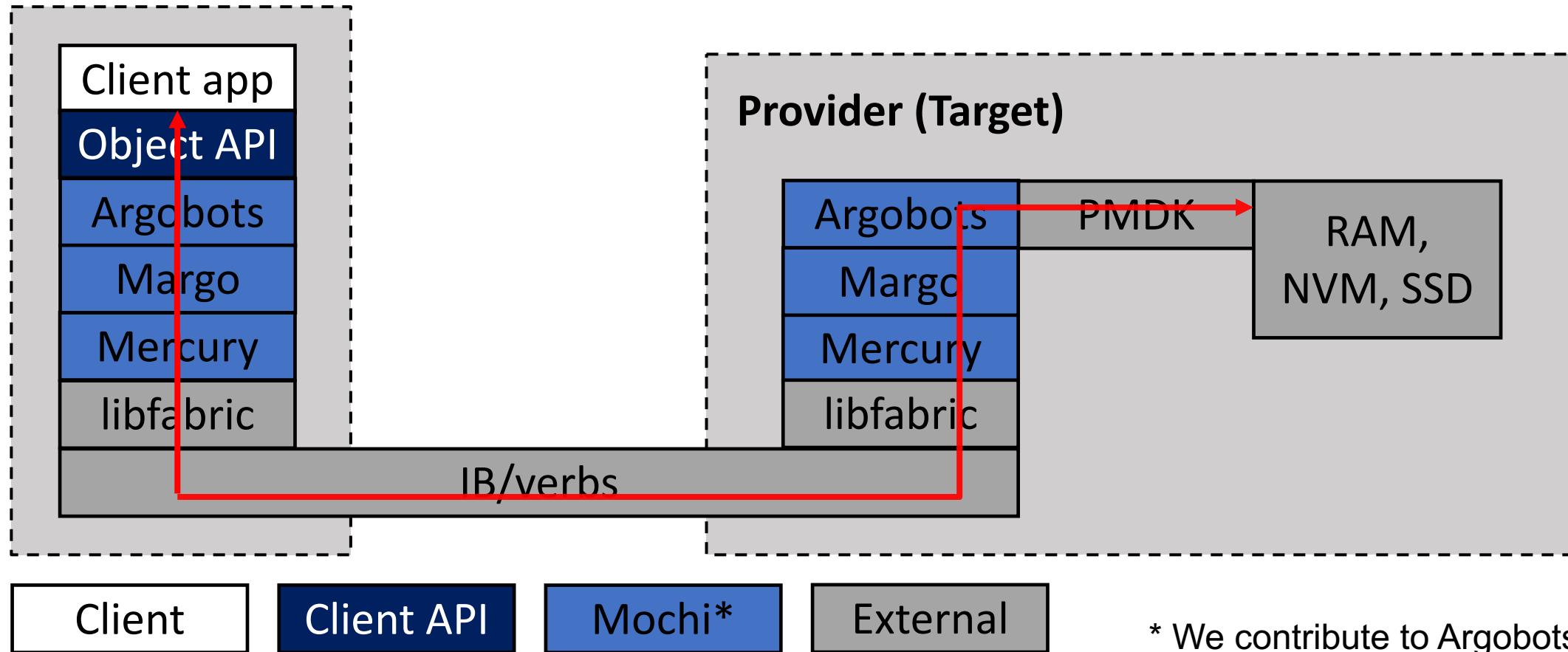
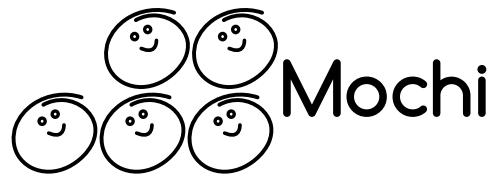


... and Thallium, if writing in C++!



# Microservice Example: BAKE

A Composed Service for Remotely Accessing Objects



\* We contribute to Argobots, but it's primarily supported by P. Balaji's team.

P. Carns et al. "Enabling NVM for Data-Intensive Scientific Services."  
INFLOW 2016, November 2016.

# Agenda

8:30 – 9:00	Welcome and Introductions	Rob Ross
9:00 – 9:20	Summit System Details	Rob Latham
9:20 – 10:00	Mochi Core: Margo and Thallium	Matthieu Dorier
10:00 – 10:30	<i>Break</i>	
10:30 – 11:00	Core: SSG	Shane Snyder
11:00 – 11:20	Microservices: Bake	Phil Carns
11:20 – 12:00	Microservices: SDSKV and SDSDKV	Rob Latham
12:00 – 1:00	<i>Lunch</i>	
1:00 – 1:30	Composing Mochi Services	Matthieu Dorier
1:30 – 1:45	Mochi Best Practices	Phil Carns
1:45 – 3:00	Mochi Service Hacking Hands-On	
3:00 – 3:30	<i>Break</i>	
3:30 – 4:15	Hands-On Continued	
4:15 – 4:30	Recap, Feedback, and Next Steps	

# Mochi Components: A Reference

	<b>Component</b>	<b>Summary</b>
<b>Core</b>		
	<b>Argobots</b>	Argobots provides user-level thread capabilities for managing concurrency.
	<b>Mercury</b>	Mercury is a library implementing remote procedure calls (RPCs).
	<b>Margo</b>	Margo is a C library using Argobots to simplify building RPC-based services.
	<b>Thallium</b>	Thallium allows development of Mochi services using modern C++.
	<b>SSG</b>	SSG provides tools for managing groups of providers in Mochi.
<b>Utilities</b>		
	<b>ABT-IO</b>	ABT-IO enables POSIX file access with the Mochi framework.
	<b>ch_placement</b>	ch-placement is a library implementing multiple hashing algorithms.
	<b>MDCS</b>	MDCS exposes remotely accessible counters for monitoring purposes.
	<b>Shuffle</b>	Shuffle provides a scalable all-to-all data shuffling service.
<b>Microservices</b>		
	<b>BAKE</b>	Bake enables remote storage and retrieval of named blobs of data.
	<b>POESIE</b>	Poesie embeds language interpreters in Mochi services.
	<b>REMI</b>	REMI is a microservice that handles migrating sets of files between nodes.
	<b>SDSKV</b>	SDSKV enables RPC-based access to multiple key-value backends.
	<b>SDSDKV</b>	SDSDKV provides a distributed key-value service using Mochi components.

# Mochi Components: Core

# Argobots

**Argobots provides user-level thread capabilities for managing concurrency.**

Argobots implements a lightweight runtime system that supports computation and data movement with massive concurrency. It leverages the lowest-level constructs in the hardware and OS: lightweight notification mechanisms, data movement engines, memory mapping, and data placement strategies.

Location	<a href="https://github.com/pmodels/argobots">https://github.com/pmodels/argobots</a>
Programming Language	C
Dependencies	
Status	Production
Python Wrappers	

# Mercury

**Mercury is a library implementing remote procedure calls (RPCs).**

Mercury supports the execution of remote procedures in a variety of scenarios, including over high-performance network fabrics. Mercury also provides facilities to support handling of large RPC arguments (i.e., bulk data), including taking advantage of RDMA when possible.

Location	<a href="https://github.com/mercury-hpc/mercury">https://github.com/mercury-hpc/mercury</a>
Programming Language	C
Dependencies	Libfabric (or other supported communication library)
Status	Production
Python Wrappers	

# Margo

**Margo is a C library using Argobots to simplify building RPC-based services.**

Margo includes a set of Argobots-aware wrappers that translate Mercury's event-driven communication model into a more intuitive sequential user-level thread model. Margo implements a communication progress loop that consolidates best practices for polling and communication, and manages threads to execute handler functions for each RPC request.

Location	<a href="https://xgitlab.cels.anl.gov/sds/margo">https://xgitlab.cels.anl.gov/sds/margo</a>
Programming Language	C
Dependencies	Argobots, Mercury
Status	Production
Python Wrappers	<a href="https://xgitlab.cels.anl.gov/sds/py-margo">https://xgitlab.cels.anl.gov/sds/py-margo</a>

# Thallium

**Thallium is a C++14 library wrapping Margo allowing development of RPC-based services using all the power of modern C++.**

Thallium wraps all Margo, Mercury, and Argobots concepts into C++ classes with automated memory management, replaces Mercury's serialization mechanism with a template-based mechanism, and uses modern C++'s variadic templates to turn RPCs into callable objects.

Location	<a href="https://xgitlab.cels.anl.gov/sds/thallium">https://xgitlab.cels.anl.gov/sds/thallium</a>
Programming Language	C++14
Dependencies	Argobots, Margo, Mercury
Status	Production
Python Wrappers	

# Scalable Service Groups (SSG)

**SSG provides tools for managing groups of providers in Mochi.**

SSG aids in the development of distributed services by allowing the tracking of groups of providers and notification of change in the group (e.g., losing a participant). The current SSG implementation leverages the SWIM weakly consistent group membership protocol to detect failures and evict noncommunicating members of the group.

Location	<a href="https://xgitlab.cels.anl.gov/sds/ssg">https://xgitlab.cels.anl.gov/sds/ssg</a>
Programming Language	C
Dependencies	Margo, PMIX (optional), MPI (optional)
Status	Prototype
Python Wrappers	<a href="https://xgitlab.cels.anl.gov/sds/py-ssg">https://xgitlab.cels.anl.gov/sds/py-ssg</a>

# Mochi Components: Utilities

# ABT-IO

**ABT-IO enables POSIX file access with the Mochi framework.**

ABT-IO works by managing blocking POSIX file operations with an Argobots execution stream so that file I/O can proceed concurrent with other Mochi operations.

Location	<a href="https://xgitlab.cels.anl.gov/sds/abt-io">https://xgitlab.cels.anl.gov/sds/abt-io</a>
Programming Language	C
Dependencies	Argobots
Status	Production
Python Wrappers	

# ch-placement

**ch-placement is a library implementing multiple hashing algorithms.**

ch-placement can be used by a distributed storage system to access multiple hashing algorithms, distance metrics, and virtual node settings using a consistent API. This is typically used for mapping data and/or work across multiple providers.

Location	<a href="https://xgitlab.cels.anl.gov/codes/ch-placement">https://xgitlab.cels.anl.gov/codes/ch-placement</a>
Programming Language	C
Dependencies	
Status	Production
Python Wrappers	

# Mochi Diagnostic Counter Service (MDCS)

**MDCS exposes remotely accessible counters for monitoring purposes.**

Performance understanding can be complex in distributed systems. MDCS provides a base capability for remote counter access helpful in the development of monitoring capabilities for Mochi services.

Location	<a href="https://xgitlab.cels.anl.gov/sds/mdcs">https://xgitlab.cels.anl.gov/sds/mdcs</a>
Programming Language	C
Dependencies	Margo
Status	Prototype
Python Wrappers	

# Three Hop All-to-All Shuffle

**Shuffle provides a scalable all-to-all data shuffling service.**

Shuffle reduces the memory/network resource load of an all-to-all communication pattern by adding additional hierarchy (3 hops: local in, remote, and local out). Shuffle provides an asynchronous “enqueue and return” API. It supports RPC batching and implements its own flow control mechanisms on top of Mercury.

Location	<a href="https://github.com/pdlfs/deltafs-shuffle">https://github.com/pdlfs/deltafs-shuffle</a>
Programming Language	C/C++ with a C API
Dependencies	Mercury, MPI
Status	Production
Note	Related libraries: Nexus and Mercury-progressor

# Mochi Components: Microservices

# Bake

**Bake enables remote storage and retrieval of named blobs of data.**

Bake is a microservice for storing blobs on nonvolatile memory or file-based storage backends, using the PMDK object API for efficient storage to nonvolatile memory backends and POSIX calls for file-based backends. Bake does not itself provide sharding across multiple backends/instances.

Location	<a href="https://xgitlab.cels.anl.gov/sds/bake">https://xgitlab.cels.anl.gov/sds/bake</a>
Programming Language	C
Dependencies	Margo, PMDK or POSIX
Status	Production for PMDK backend, prototype for POSIX file backend
Python Wrappers	<a href="https://xgitlab.cels.anl.gov/sds/py-bake">https://xgitlab.cels.anl.gov/sds/py-bake</a>

# Poesie

## **Poesie embeds language interpreters in Mochi services.**

Poesie clients can send code (e.g., Lua, Python) to Poesie providers for remote execution on their behalf, providing a flexible method of extending an RPC-based service. Currently we are evaluating Poesie as a method for enabling rapid reconfiguration of Mochi composed services.

Location	<a href="https://xgitlab.cels.anl.gov/sds/poesie">https://xgitlab.cels.anl.gov/sds/poesie</a>
Programming Language	C
Dependencies	Margo, Lua and/or Python
Status	Prototype
Python Wrappers	

# REsource Migration Interface (REMI)

**REMI is a microservice that handles migrating sets of files between nodes.**

REMI works in terms of *filessets*, collections of files in the directory hierarchy. REMI uses RDMA and memory mapping to efficiently transfer potentially large filesets at once, managing creation of the hierarchy at the destination and notification on completion. This provides a basic capability for data management useful for constructing adaptive or hierarchical data services.

Location	<a href="https://xgitlab.cels.anl.gov/sds/remi">https://xgitlab.cels.anl.gov/sds/remi</a>
Programming Language	C++
Dependencies	Thallium
Status	Prototype
Python Wrappers	<a href="https://xgitlab.cels.anl.gov/sds/py-remi">https://xgitlab.cels.anl.gov/sds/py-remi</a>

# sds-keyval (SDSKV)

**SDSKV enables RPC-based access to multiple key-value backends.**

SDSKV is a recognition of the importance of KV stores in modern data management. SDSKV enables RPC-based access to multiple KV backends, including LevelDB, Berkeley DB, and in-memory databases. SDSKV is typically helpful for metadata and small data management. SDSKV does not itself provide sharding across multiple backends/instances.

Location	<a href="https://xgitlab.cels.anl.gov/sds/sds-keyval">https://xgitlab.cels.anl.gov/sds/sds-keyval</a>
Programming Language	C/C++
Dependencies	LevelDB, Berkeley DB, or other backend KV store
Status	Production
Python Wrappers	<a href="https://xgitlab.cels.anl.gov/sds/py-sdskv">https://xgitlab.cels.anl.gov/sds/py-sdskv</a>

# SDSDKV

**SDSDKV provides a distributed key-value service using Mochi components.**

SDSDKV leverages SDSKV, SSG, and ch-placement to provide key-value services sharded over multiple instances on a variety of different backends. SSG allows for detection of failure of members of the service. Current implementation is not resilient to faults.

Location	<a href="https://xgitlab.cels.anl.gov/sds/sdsvkv">https://xgitlab.cels.anl.gov/sds/sdsvkv</a>
Programming Language	
Dependencies	ch-placement, SDSKV, SSG
Status	Prototype
Python Wrappers	



This work is in part supported by the Director, Office of Advanced Scientific Computing Research, Office of Science, of the U.S. Department of Energy under Contract No. DE-AC02-06CH11357; in part supported by the Exascale Computing Project (17-SC-20-SC), a joint project of the U.S. Department of Energy's Office of Science and National Nuclear Security Administration, responsible for delivering a capable exascale ecosystem, including software, applications, and hardware technology, to support the nation's exascale computing imperative; and in part supported by the U.S. Department of Energy, Office of Science, Office of Advanced Scientific Computing Research, Scientific Discovery through Advanced Computing (SciDAC) program.

<http://www.mcs.anl.gov/research/projects/mochi/>