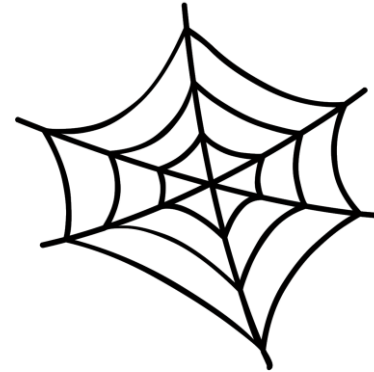


MOCHI BOOT CAMP



SESSION 6: NETWORK FABRICS & RPC TUNING



PHILIP CARNS
Argonne National Laboratory

NETWORK FABRIC SUPPORT IN MOCHI



Argonne National Laboratory is a
U.S. Department of Energy laboratory
managed by UChicago Argonne, LLC.



MERCURY NETWORK SUPPORT

“NA” plugins interface with communication libraries

▪ Libfabric

- Most important library for remote communication
- Actively supported and improving
- Supports many networks



Thank you to Intel
for contributing the
libfabric driver in
Mercury.

▪ NA SM

- Mercury’s own transport for shared memory, on-node communication
- Automatically used when “remote” peers are actually reachable on local node

▪ Self RPCs

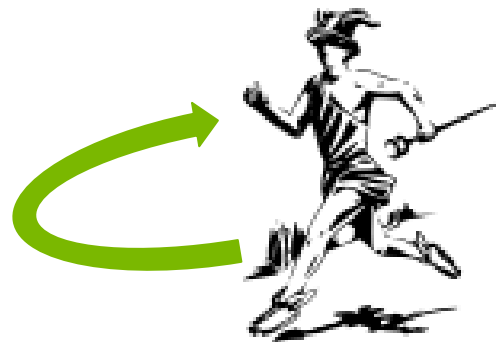
- Not really a transport, but a fast execution mechanism for in-process RPCs

▪ MPI: strictly for prototyping, not performant

▪ BMI: legacy code, not recommended

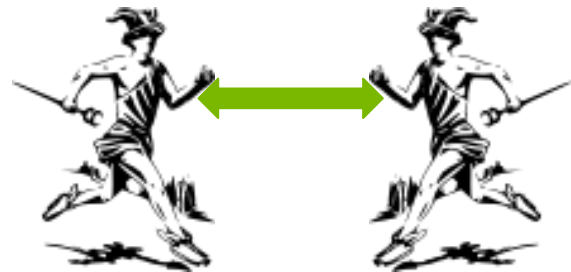
MERCURY “SELF” RPCS

Not really a transport, but good to know about!



- What if a process sends an RPC to itself?
- Mercury detects this and skips all NA plugins
 - RPC handler functions are invoked directly
 - Still delegated to the same Margo thread pool they would have been
 - Bulk transfers become memory copies
- Why is this important?
 - Flexible composition!
 - Components can communicate with other components using RPCs, regardless of whether peer component is in-process, on-node, or remote
- This functionality is enabled by default

MERCURY NA SM PLUGIN



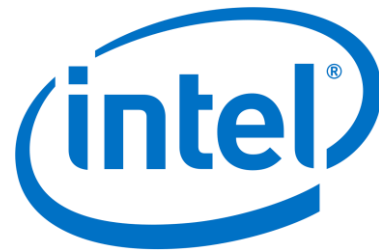
- Can be selected explicitly using the “na+sm://” network address prefix, as we have done in some of the hands-on exercises
- Will also be implicitly used when a remote process is detected locally
- On Linux:
 - Uses Cross-Memory Attach capability (single copy)
 - If bulk transfers fail, on some system you may need (as root):
`“echo 0 > /proc/sys/kernel/yama/ptrace_scope”`
- Others:
 - Uses conventional POSIX shared memory (2 copies)
- Enabled by default in spack or with NA_USE_SM variable in cmake

MERCURY LIBFABRIC: IB/VERBS



- Can be selected with a “verbs://” Mercury address prefix
- Uses the “verbs;ofi_rxm” provider combination in libfabric
- If using 1.8 or newer: “`export FI_MR_CACHE_MAX_COUNT=0`”
 - Disables new memory registration cache, which has some performance gotchas for Mochi
- If using anything older: “`export FI_FORK_UNSAFE=1`”
 - Mochi doesn’t care about fork; this is a compatibility issue with MPI
- Aside from above caveats, works great! Fast and stable.
- Performance tuning: be aware that memory registration can be relatively expensive (`margo_bulk_create()`), especially for small transfers

MERCURY LIBFABRIC: OMNIPATH



- Can be selected with a “psm2://” Mercury address prefix
- Uses the “gni” provider in libfabric
- “`export FI_MR_CACHE_MAX_COUNT=0`”
 - Allows Mochi and MPI to share interface if needed
- Performance is excellent
- Stability can be a problem, depends on workload, libpsm2 version, etc.
- All OmniPath systems also support verbs through emulation, at a performance penalty
- Memory registration (`margo_bulk_create()`) is free (noop), but not really because you pay performance cost at communication time on cold memory access

MERCURY LIBFABRIC: ARIES/GNI



- Can be selected with a “gni://” Mercury address prefix
- Uses the “gni” provider in libfabric
- `export FI_MR_CACHE_MAX_COUNT=0`
 - Tells Cray-MPICH to be less aggressive with resource consumption
- Stability is great, performance is (mostly) great but there are gotchas:
 - Communication across separately launched executables requires DRC, which in turn requires some explicit setup outside of Mercury
 - The above only works with libfabric 1.8.1rc or newer
 - Latency is poor when not busy-spinning (more on how to do this later)
 - Latency is poor (any polling mode) on KNL cores
 - Memory registration / copy tradeoffs are unclear, possibly similar to IB/Verbs

MERCURY LIBFABRIC: TCP

TODO: tin
Can string logo

- Can be selected with a “tcp://” Mercury address prefix
- Uses the “tcp;ofi_rxm” provider combination in libfabric
- Stability is good (libfabric 1.8.0 or newer)
- Performance is poor (when other alternatives are possible)
 - Part of the problem is emulation of RDMA over sockets
- There are other options that you may see in previous examples:
 - “sockets://” activates the legacy/reference sockets implementation in libfabric
 - “bmi+tcp://” activates TCP support in the BMI library (if enabled)
 - All have similar performance characteristics; “tcp;ofi_rxm” is now the most actively maintained, though

POLLING MODES

Hey network card: Are you done? Are you done? Are you done?

- How about now?
- Generally speaking, any of the network transports can be driven in “polling mode”, where a CPU core constantly checks for progress.
 - This is fast - no interrupt or context switch
 - This also eats a CPU core. The value of this tradeoff depends on the use case
 - Benchmarks do it all the time!
- Margo defaults to an adaptive model: it idles gracefully when there is nothing to do, at times will switch to busy mode
- You can force it to always busy poll by setting `progress_mode` to `NA_NO_BLOCK`

<https://xgitlab.cels.anl.gov/sds/sds-tests/blob/master/perf-regression/margo-p2p-latency.c#L91>

THE STATE OF MOCHI TRANSPORTS

September 2019

- Things are continually improving, particularly in the interaction between Mochi, libfabric, and MPI
- But as you can see, there are some quirks
- Want to keep up over time? Your best option right now is to monitor the scripts that we use for regression testing on different transports:

<https://xgitlab.cels.anl.gov/sds/sds-tests/tree/master/perf-regression>

- We update these over time with best practice for platforms at ANL

RPC TUNING



Argonne National Laboratory is a
U.S. Department of Energy laboratory
managed by UChicago Argonne, LLC.

BUILT-IN RPC PROFILING

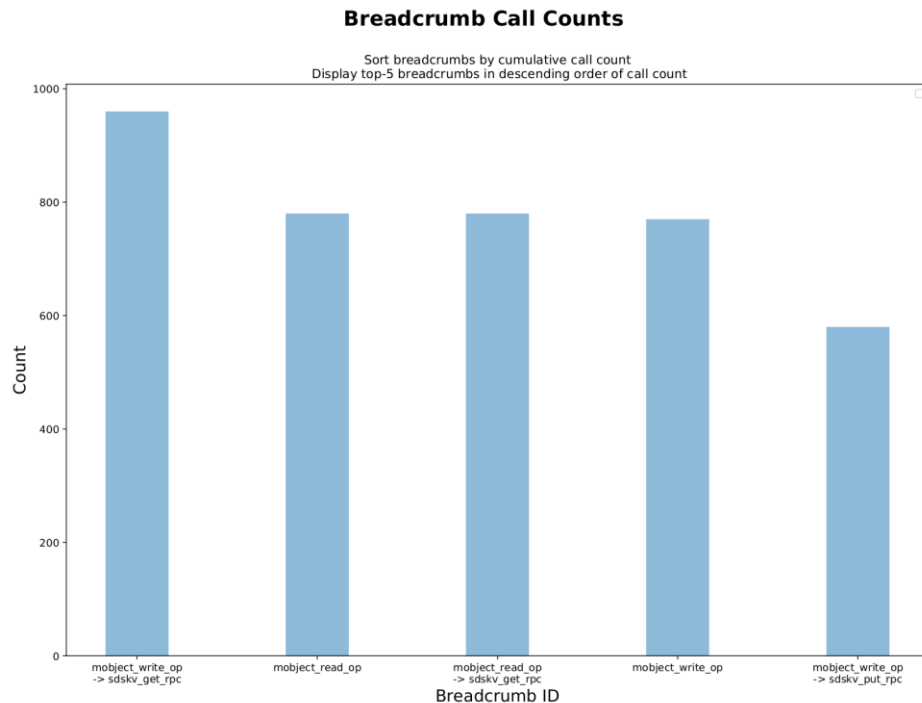
Coming soon!

- How do you tune a Mochi service, or at least understand it's performance?
- Almost every major operation in Mochi is an RPC
 - RPCs can be chained, local, remote, etc.
- If we can automatically instrument all RPCs in a service, then any Mochi service could be quickly characterized
- This is an upcoming feature, thanks in large part to work by Srinivasan Ramesh of the University of Oregon
- Usage is modeled after Darshan in many ways:
 - Characterization capability is compiled in, no code modification needed
 - Simple command line tools to produce starting point PDF analysis

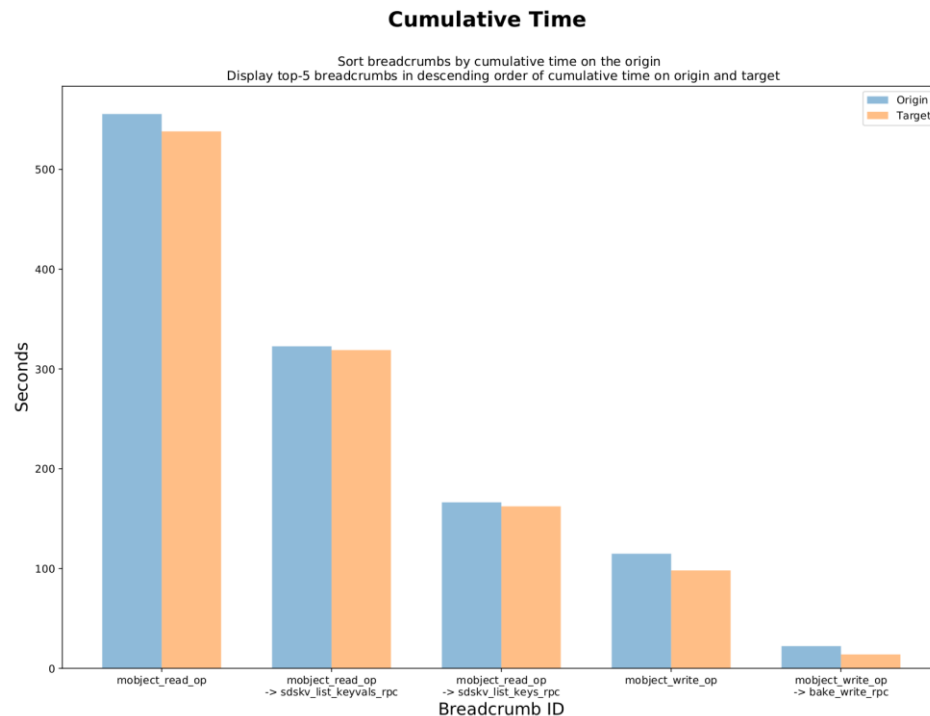
MOCHI RPC BREADCRUMBS

What RPCs are most prolific?

- This graph shows ranked list of 5 most frequently called RPCs
- Each includes lineage (1st, 3rd, and 5th RPC types in this case)



- This graph shows ranked list of 5 RPCs that cumulatively consumed the most service time
- Split by client and server side elapsed time
- Each includes lineage (2nd, 3rd, and 5th RPC types in this case)



RPC BREADCRUMBS

Status

- Mochi already includes mechanisms for tracking RPC time and tracking the lineage on chains of RPCs
 - The examples you run today include this capability already
- Remaining work is productizing the tools that can dump summary data and plot it
- There are many other possibilities for deeper analysis from that starting point

THANK YOU!



Argonne National Laboratory is a
U.S. Department of Energy laboratory
managed by UChicago Argonne, LLC.

