

# Optimizing Collective Communication for Petascale Supercomputers

<http://upc.lbl.gov>



## Introduction

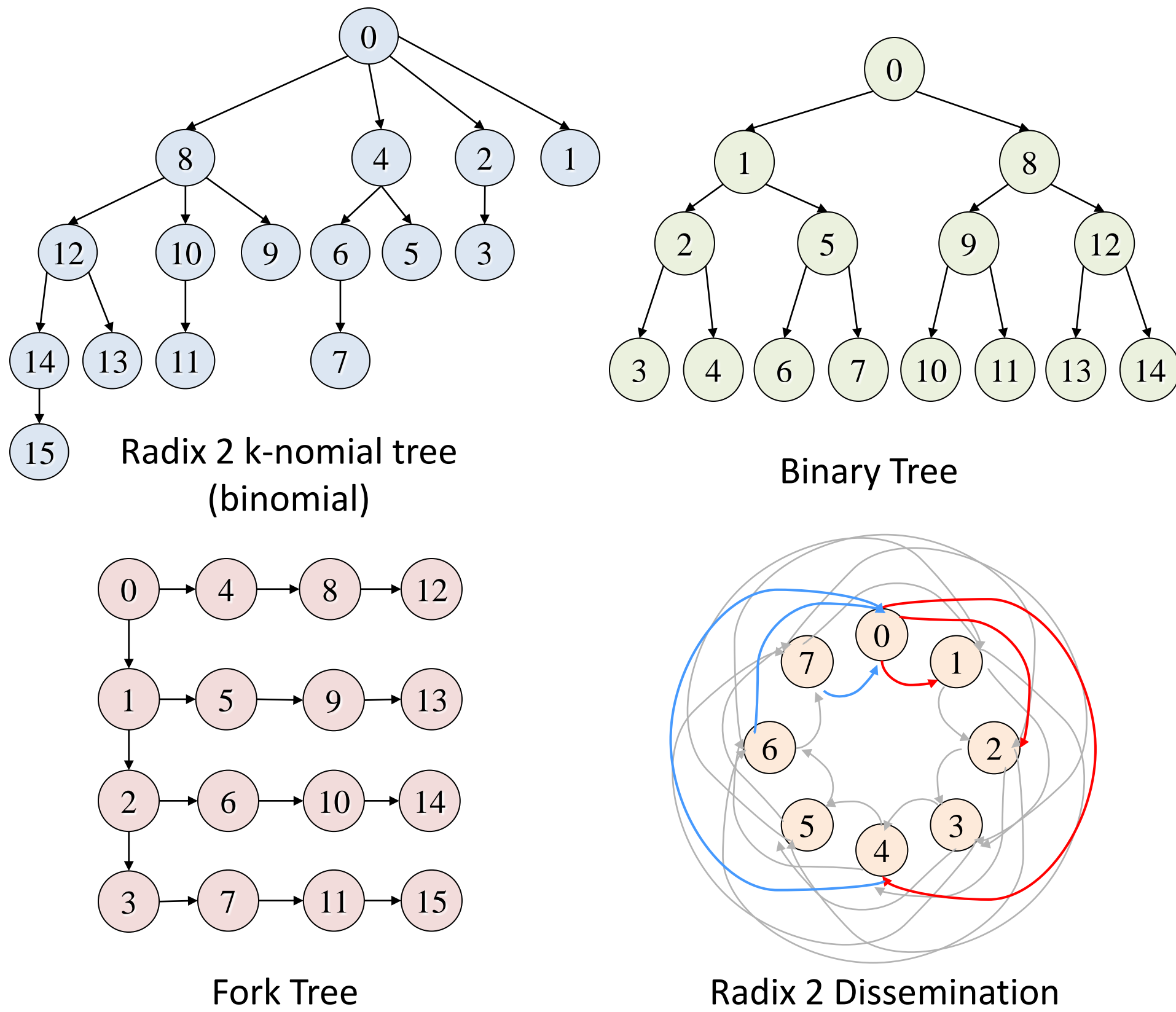
### Collective communication

- Cooperative data-movement beyond one-to-one communication
- Common building blocks for many applications
- Key bottleneck of performance scalability

### GASNet

- Portable high-performance communication primitives
- Used to implement partitioned global address space languages: e.g., UPC, Titanium, Co-array FORTRAN, and Chapel

## Example Communication Topologies



## Collectives for PGAS Languages

### Teams

- Thread-centric: Programmer explicitly specifies the threads that take part in the collective through a language level team construction API.
- Data-centric: Programmer only specifies the data for the collective. Runtime system then figures out where the data resides and performs the collective.

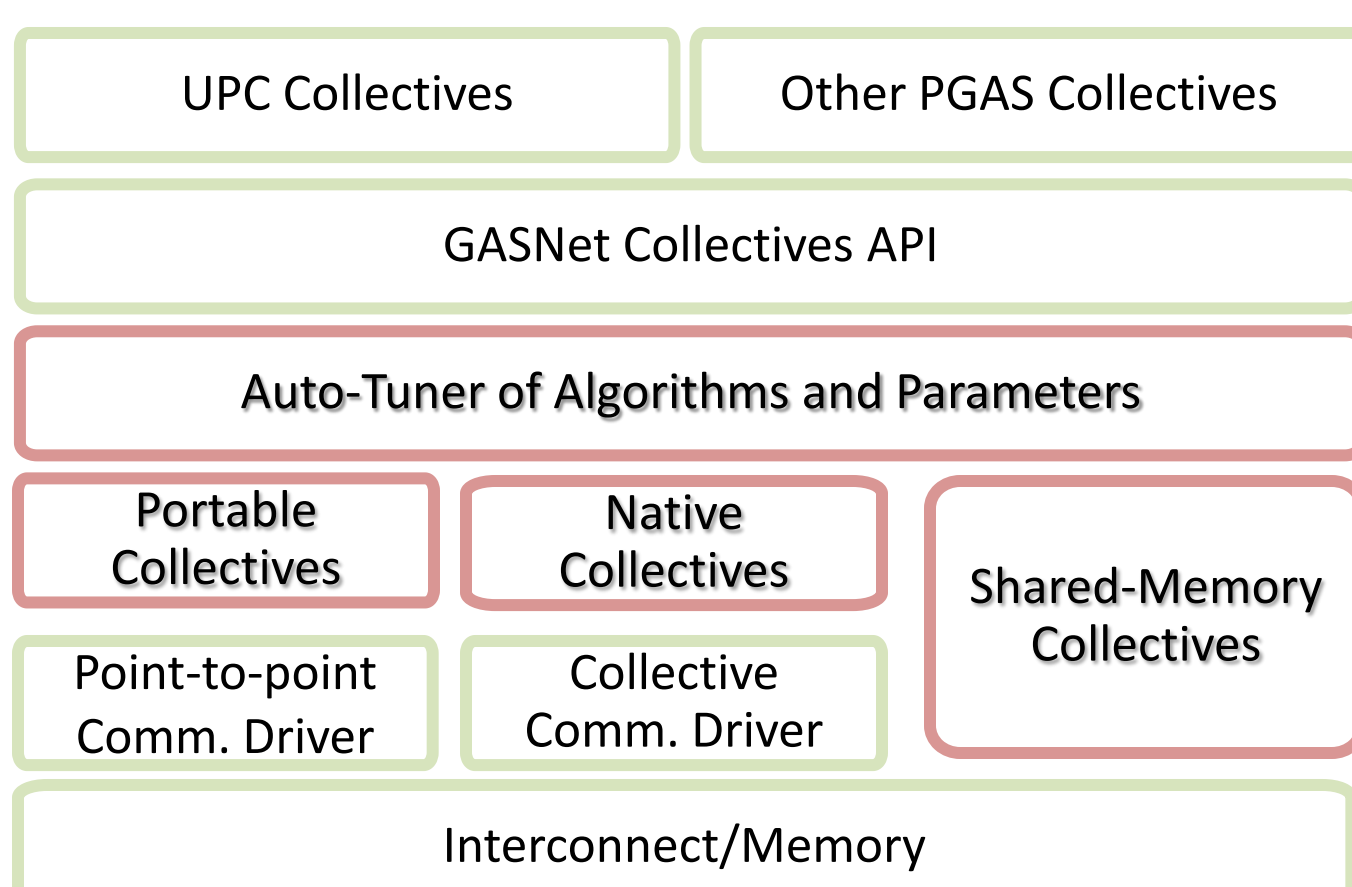
### Synchronization modes

- Loose: Data movement can start and as soon as first thread enters collective and continue until last thread leaves the collective.
- Middle: Data movement into and out local memory can occur only when the data-owner thread is in the collective operation.
- Strict: Data movement can start only after all threads have entered the collective and must finish before any thread leaves the collective.

### Optimizations

- Non-blocking collective operations that facilitate overlapping communication and computation
- Network-specific optimizations for leveraging hardware features
- Automated performance tuning for accommodating different application characteristics on multiple platforms

## Organization of GASNet Collectives



## Performance Auto-tuning

### Offline tuning

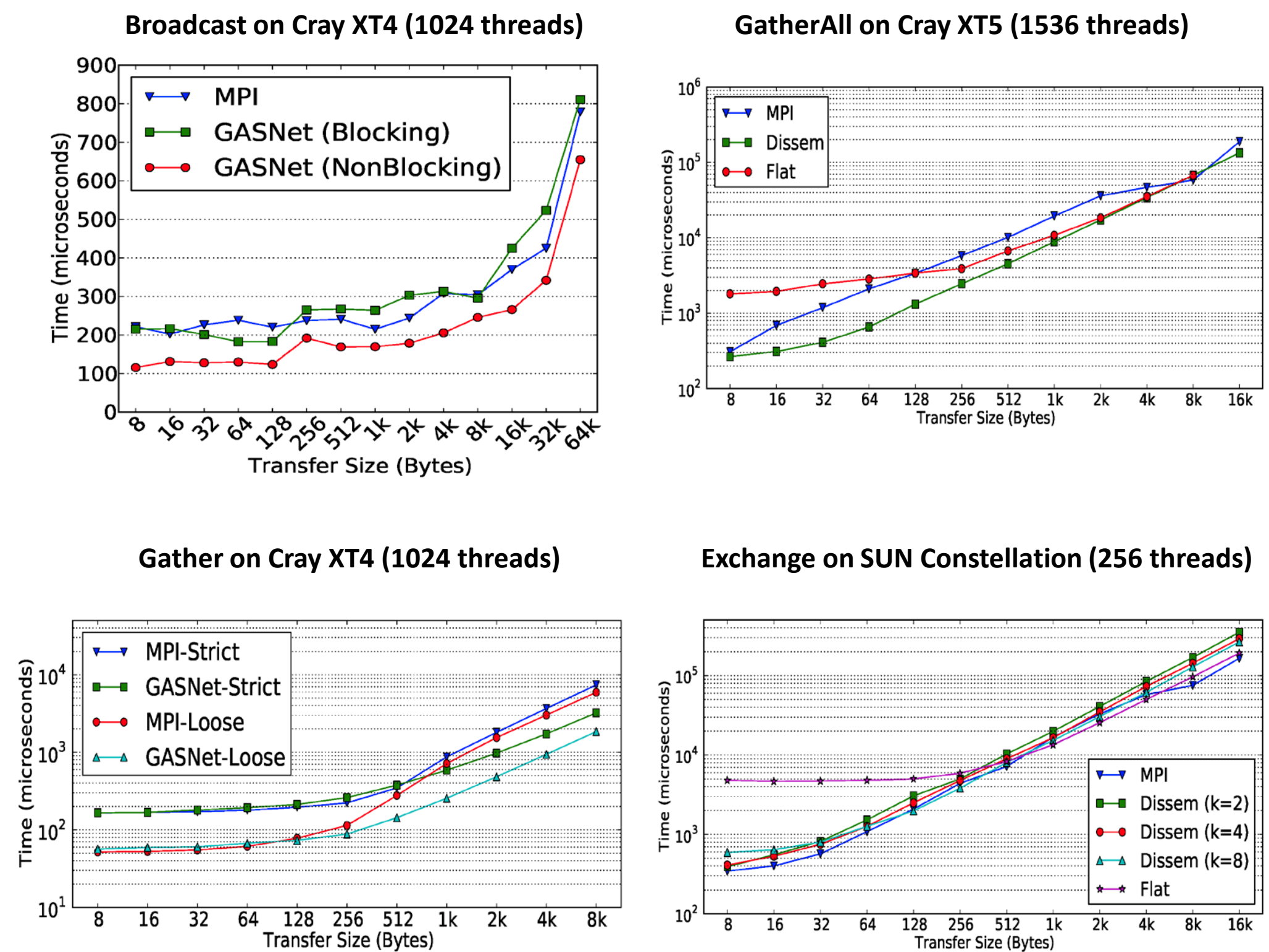
- Optimize for platform common characteristics
- Minimize runtime tuning overhead

### Online tuning

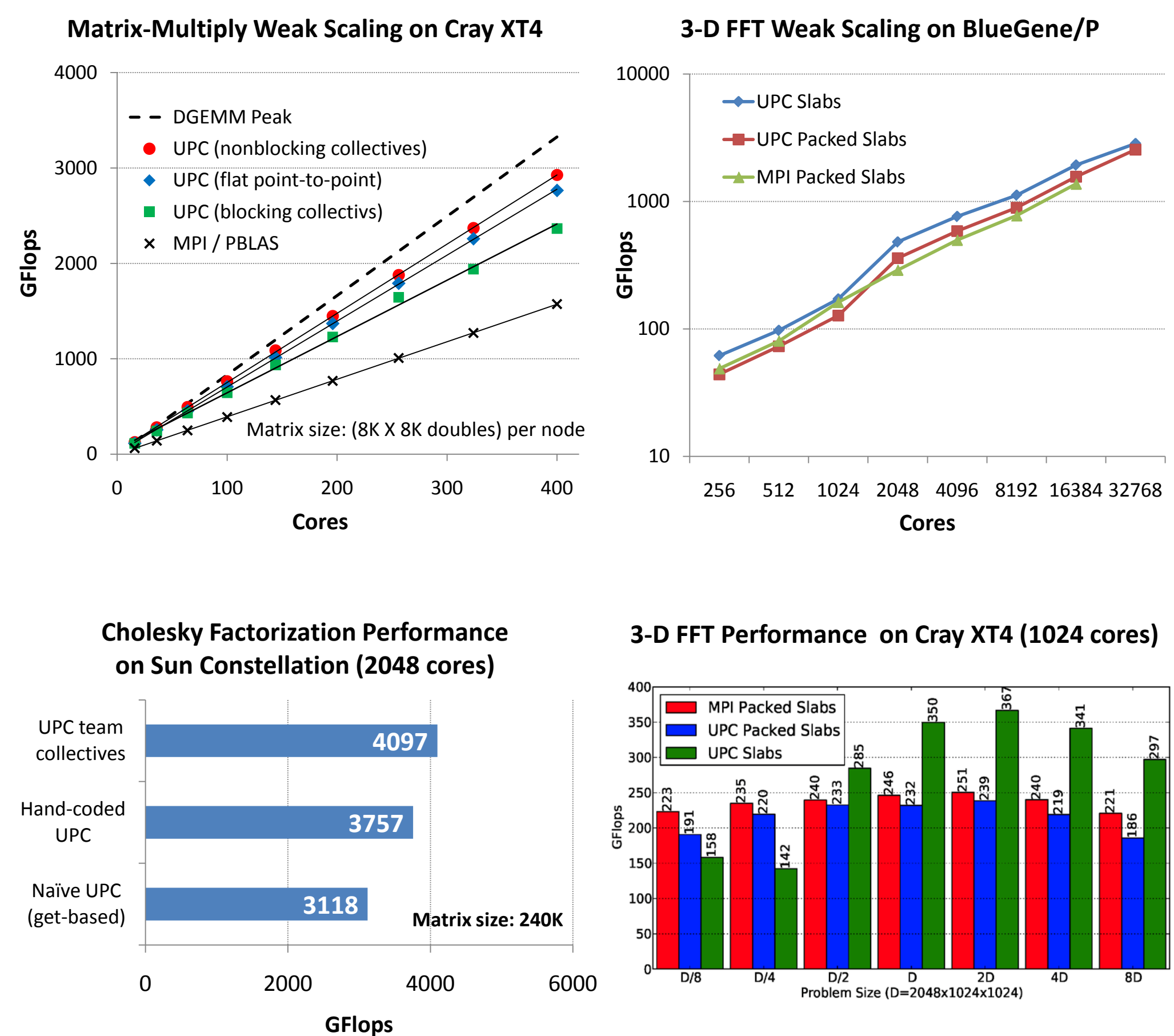
- Optimize for application runtime characteristics
- Refine offline tuning results

Performance Influencing Factors	Performance Tuning Space
<b>Hardware</b> <ul style="list-style-type: none"> <li>CPU</li> <li>Memory system</li> <li>Interconnect</li> </ul> <b>Software</b> <ul style="list-style-type: none"> <li>Application</li> <li>System software</li> </ul> <b>Execution</b> <ul style="list-style-type: none"> <li>Process/thread layout</li> <li>Input data set</li> <li>System workload</li> </ul>	<b>Algorithm selection</b> <ul style="list-style-type: none"> <li>Eager vs. rendezvous</li> <li>Put vs. get</li> <li>Collection of well-known algorithms</li> </ul> <b>Communication topology</b> <ul style="list-style-type: none"> <li>Tree type</li> <li>Tree fan-out</li> </ul> <b>Implementation-specific parameters</b> <ul style="list-style-type: none"> <li>Pipelining depth</li> <li>Dissemination radix</li> </ul>

## Micro-benchmarks



## Applications



## Conclusion

### High Productivity

- Portable performance from multi-core PCs to petascale supercomputers
- Compact and clean UPC code

### Scalable Performance

- 3-D FFT (communication intensive)
  - Weak scaling: **38%** over MPI (16K cores)
  - strong scaling: **20%** over MPI (16K cores)
- Numerical linear algebra: highly scalable performance up to **2X** MPI/PBLAS

