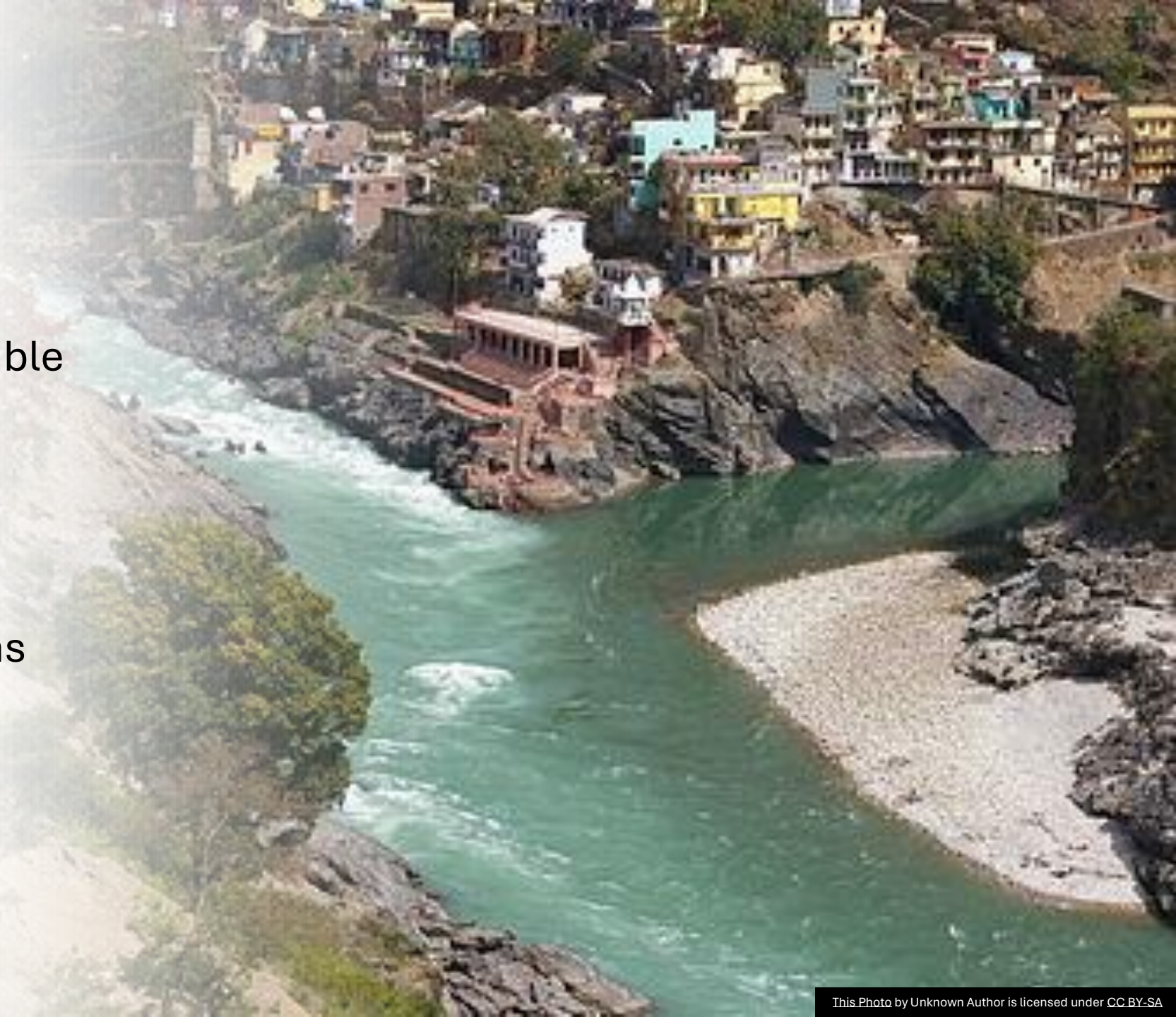


In-network transformations for AI/ML and other workloads

Rishabh Tewari, Deepak Bansal,
Gerald DeGrace, Srikanth Kandula

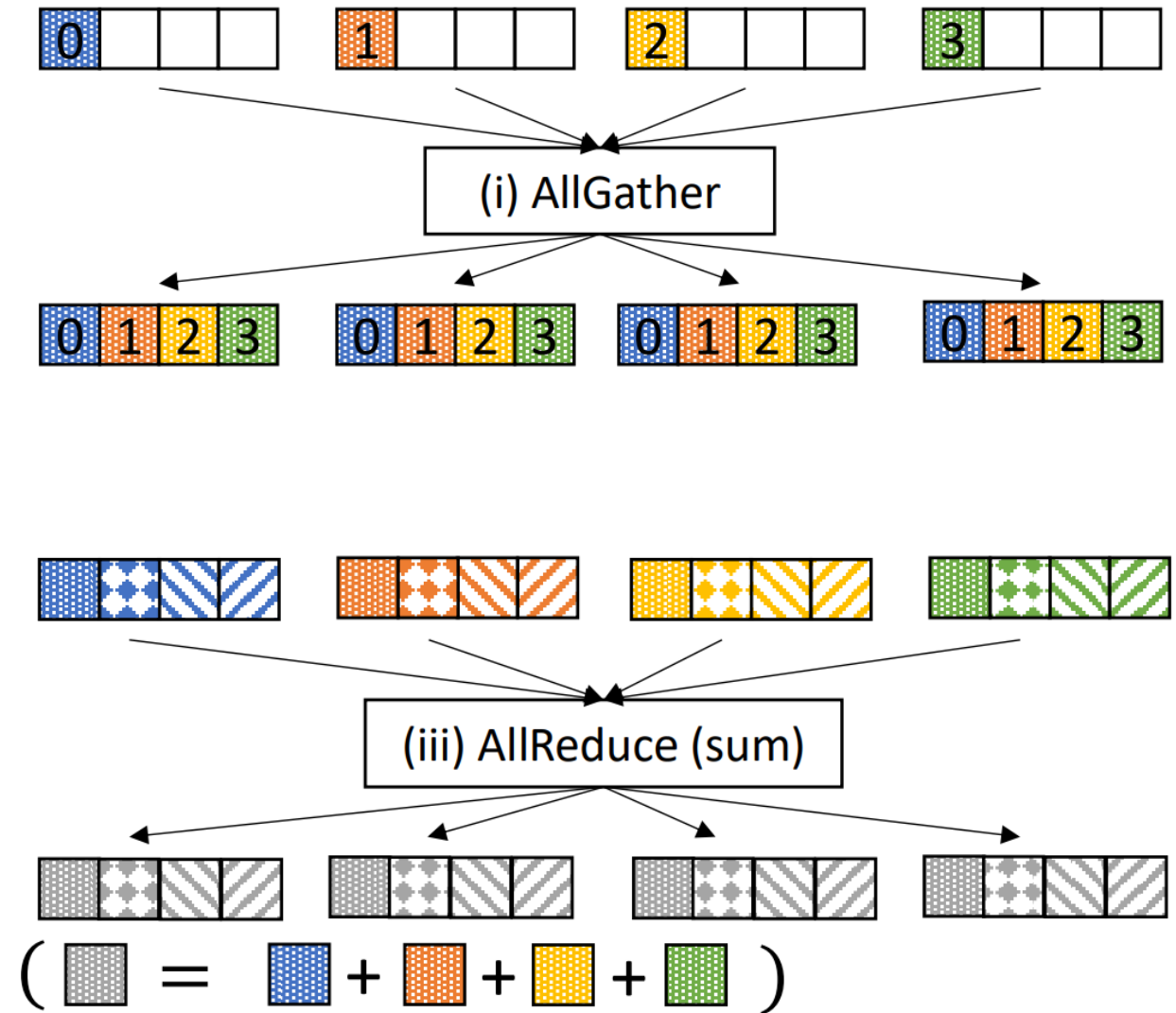
Why?

- Confluence between increasing share of new workloads and programmable network hardware
- Need Open API that allows interoperability between
 - Device implementations
 - Cloud providers
 - Application libraries



Which transformations?

- Multicast
- In-network aggregations
- Gradient compression
 - Sparse transforms
 - Floating point changes
- Future proof



Transform Abstractions

- One-to-one (e.g., gradient compression)
- One-to-many (e.g., multicast)
- Many-to-one (e.g., in-network aggregation)

Only many-to-one is stateful (requires payload state)

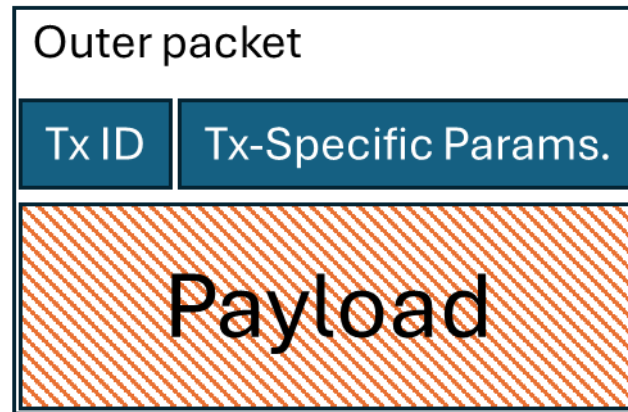
Can chain multiple transforms

* e.g., many-to-many

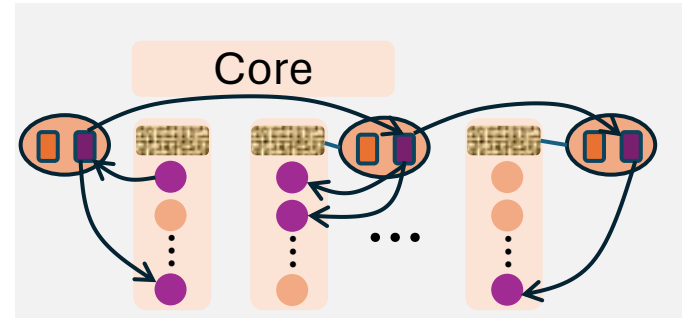
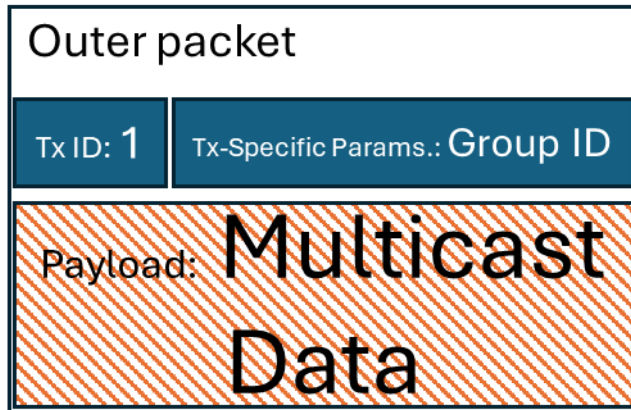
Design questions

- Carefully orchestrated paths, aggregation points vs. Disaggregated and opportunistic
 - sHarp vs. ATP
- Coexist with multiple L2, L3, L4 protocols vs. specialization

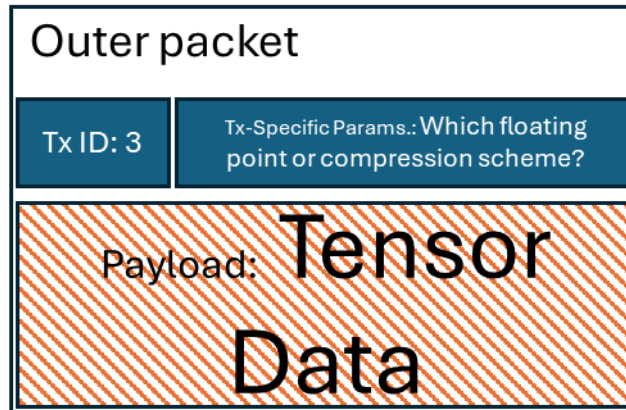
An initial stab at protocol



One-to-many Example



One-to-one example



Many-to-one

