

SADedupe:

Skew Area Inline Deduplication for Distributed Storage

Binqi Zhang, Bing Bing Zhou, Chen Wang*, Dong
Yuan, Albert Y. Zomaya

The University of Sydney, Sydney, Australia

*CSIRO, Sydney, Australia



THE UNIVERSITY OF
SYDNEY



Introduction – Deduplication

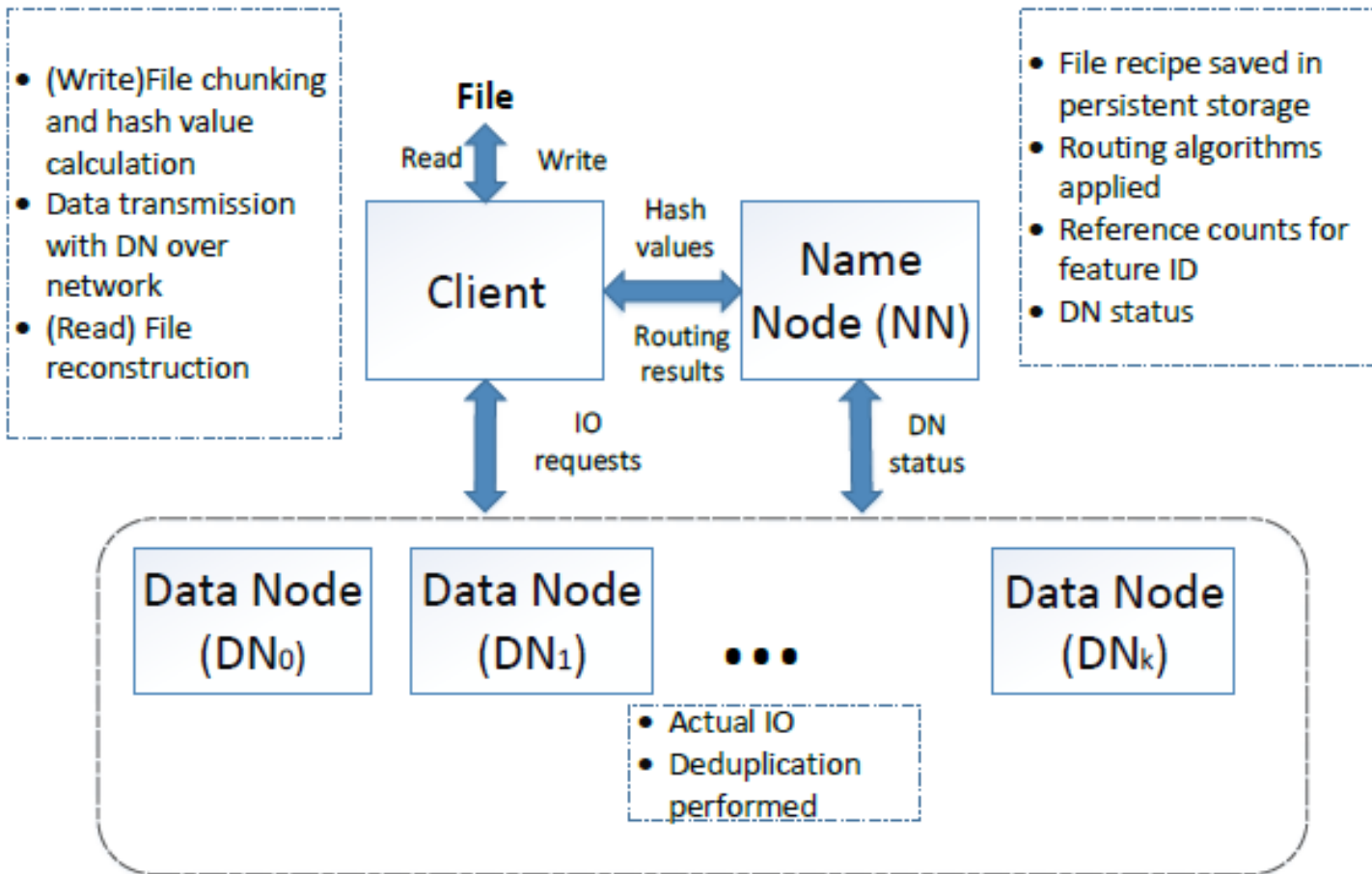
Routing

- Files -> Chunks
- Chunks -> Blocks & Hash calculation
- Extract the feature ID
- Use the feature ID to route the chunk to node

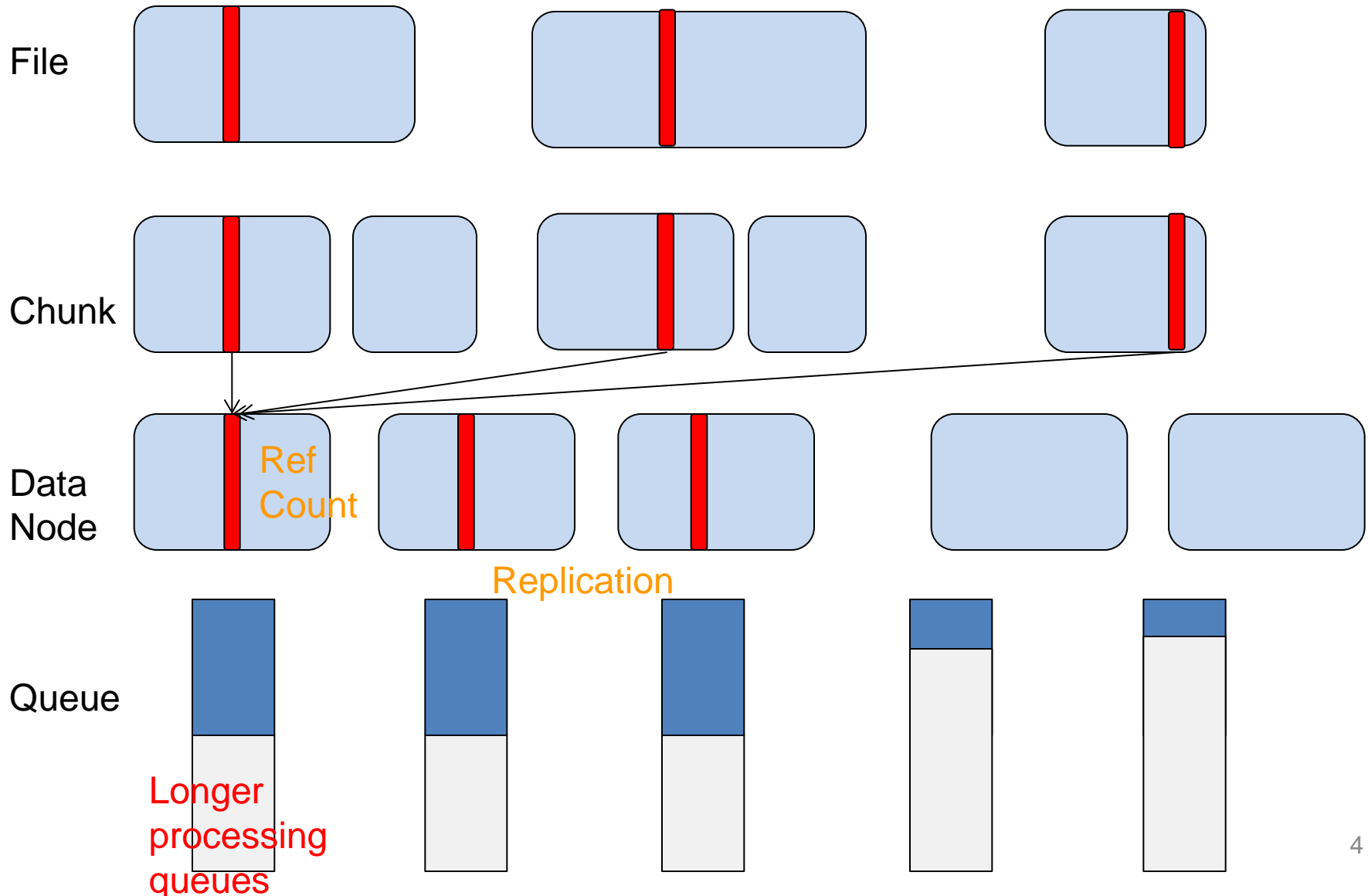
Deduplication

- Check all hash values of blocks
- If exist, then add reference
- If not, store the block

System architecture

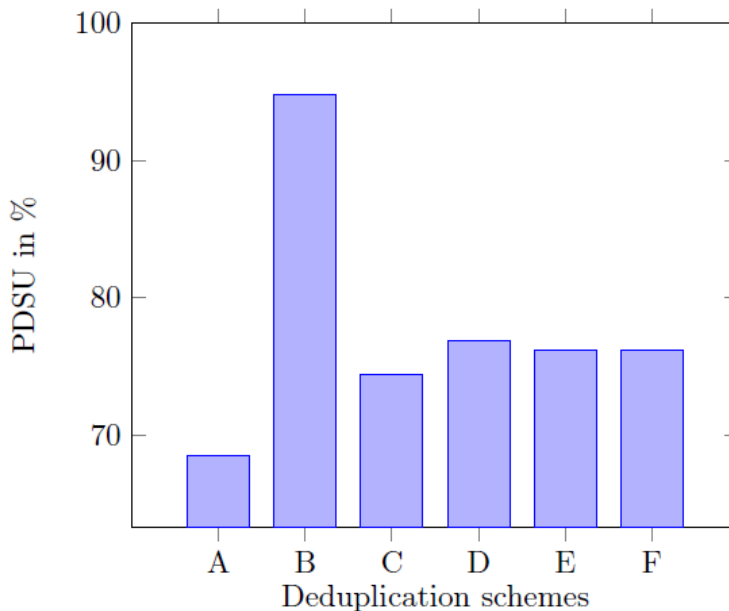


Problem



Algorithm & results

- We check the feature ID used for routing for its reference count
- Currently we use “capping” approach
- Standard deviation of post dedupe storage usage (PDSU) is examined. RT = reference count threshold



- A - Global deduplication on single node.
- B - Local deduplication on DNs with random routing.
- C - Local deduplication on DNs with MinHash routing.
- D - Local deduplication on DNs with MinHash routing, RT = 5.
- E - Local deduplication on DNs with MinHash routing, RT = 10.
- F - Local deduplication on DNs with MinHash routing, RT = 20.

Future work

- To find a better and bigger data set to illustrate the severity of the skew issue and impact to read performance
- To find a few more routing algorithms that optimize the load balancing
- Consider the replication

Thank you