



Journal based Async Indexer

Chetan Mehrotra | @chetanmeh | Oakathon - August 2017



Current Async Indexer Design

- Diff based
- Diff cost \propto Content change between 2 async indexer runs
- Starts lagging if
 - Rate of content change is high for long time
 - Bulk change done in short time say via package import
- Suffer from same problem as "hitting the observation queue limit"
- Problem seen on both Segment and Document setups
- Indexable content is mostly small subset of all content changes

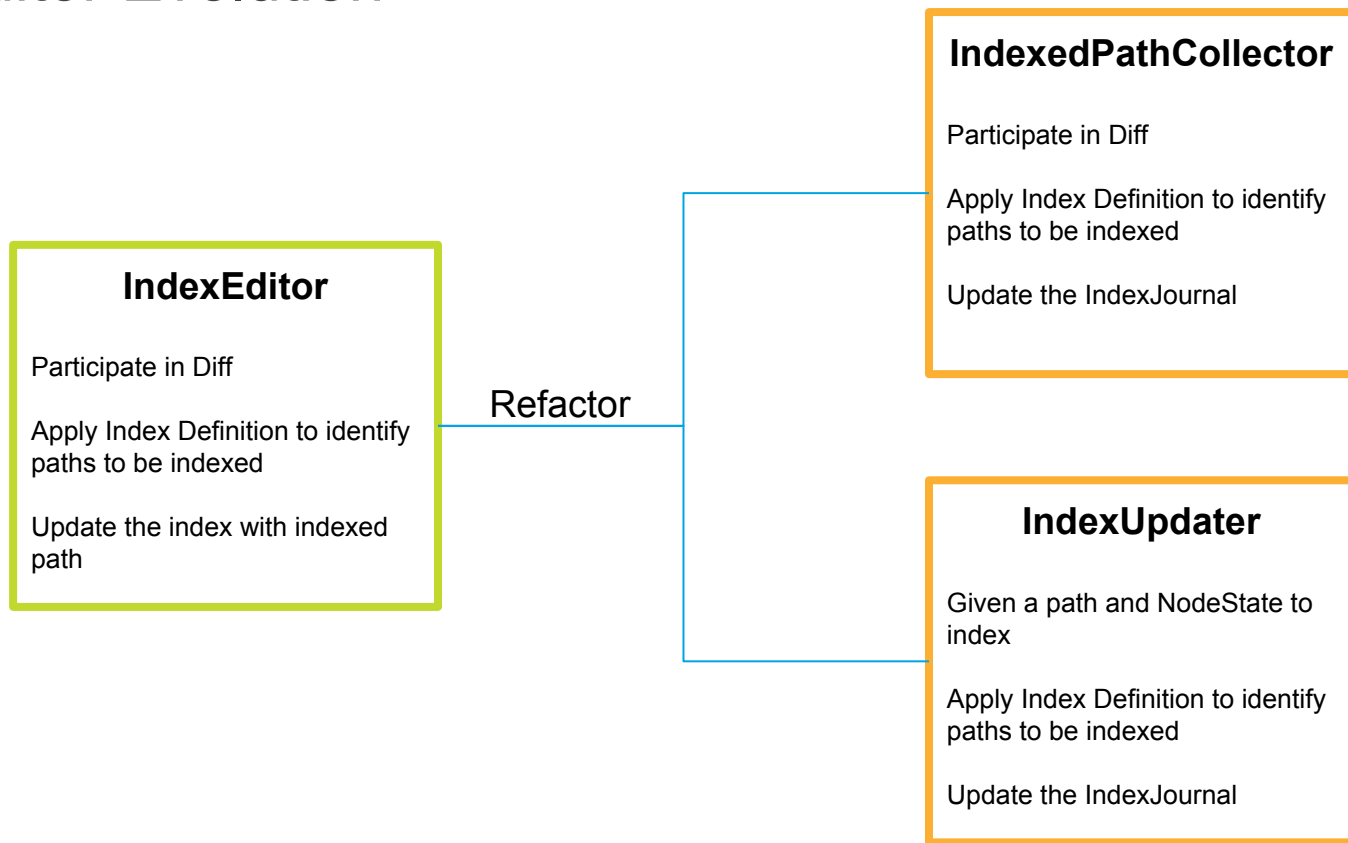
Update a journal of indexable paths as part of commit itself

- OAK-2683
- Proposal - Journal Based Async Indexer

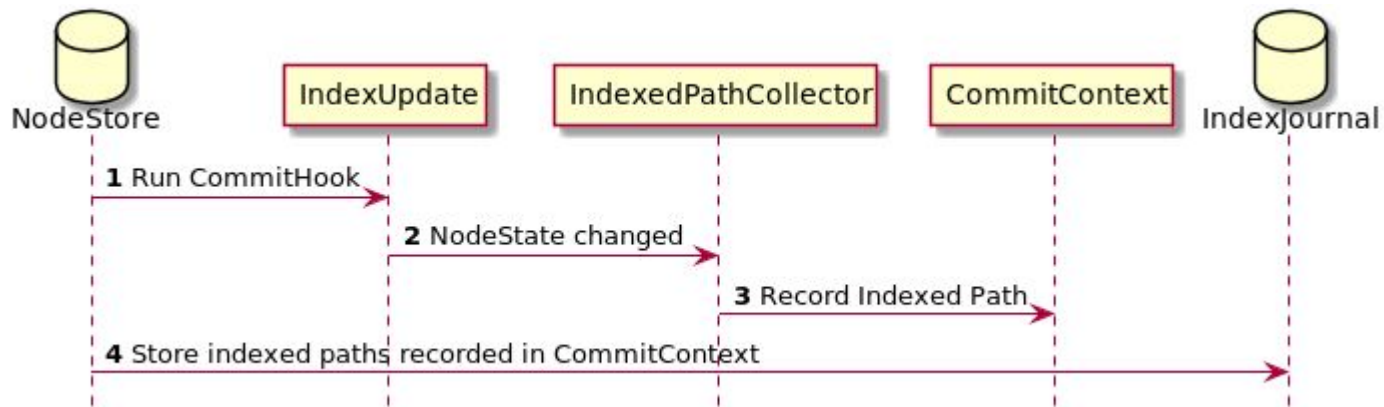
Benefits

- Avoids wasted effort on diff
- Effort spent in identifying indexable content is distributed
- Enables support for incremental indexing in smaller batches

IndexEditor Evolution



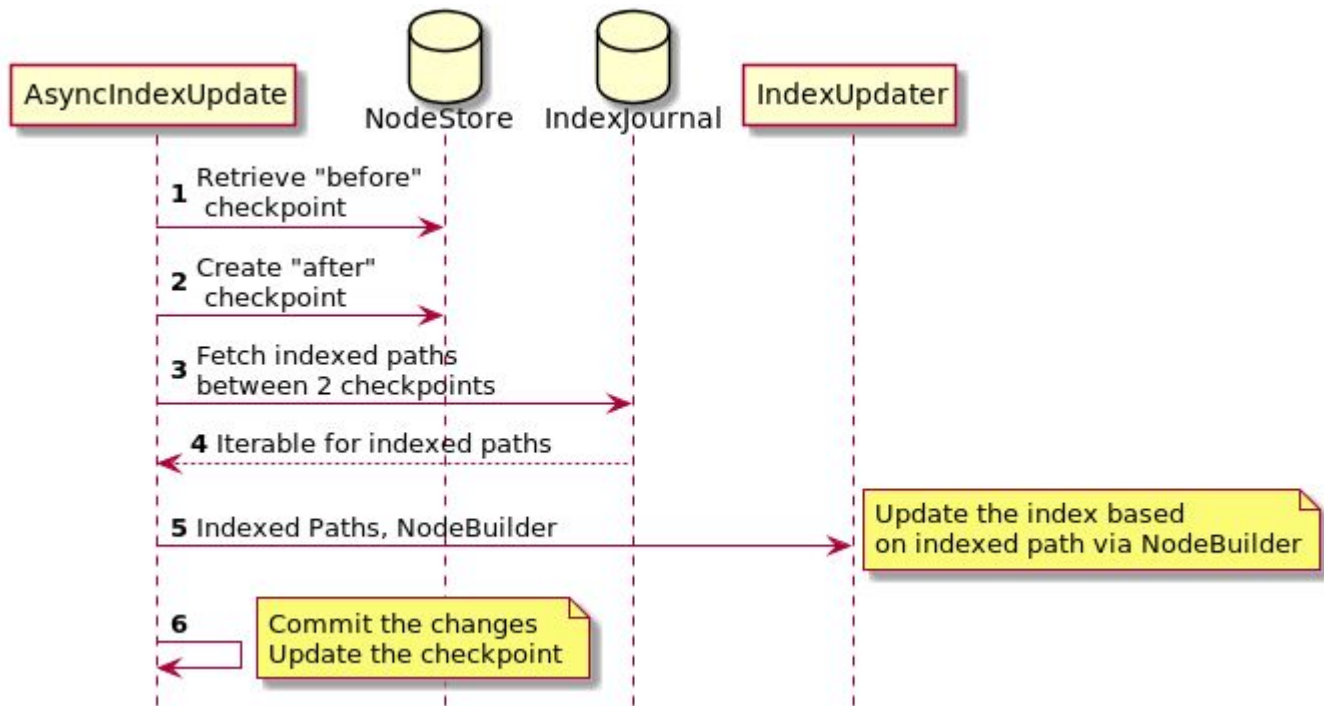
Commit Flow



Commit Flow

- Each index implementation participates in commit (prior to this only sync index editors participated in commit)
- IndexPathCollector would **identify** paths which are to be indexed based on their index definitions
- Collected paths would then be recorded in IndexJournal

IndexUpdate Flow



IndexUpdate Flow

- AsyncIndexUpdate runs as usual and gets NodeStore for before and after
- Instead of diff obtains a iterable for indexed paths from IndexJournal
- Each index provides an IndexUpdater which
 - Is given the index path and before and after state
 - Computes the changes to be done in the stored index

IndexJournal

- A new api in oak-store-spi
- Implementation provided by NodeStore
- A journal of
 - Paths which are to be indexed
 - Index paths in which the path is to be indexed
- Durable and consistent
- Iterable for changed paths
 - May have paths repeated
 - May have path entries even if not indexed

IndexJournal - DocumentNodeStore

- Built on top of existing Journal support
- An evolution of work done in
 - OAK-4808 Index external changes as part of NRT indexing
- Journal current records all changed paths
- Extend it also record which of those paths are indexed and under what index
- Provide a way to query for such indexed paths between 2 checkpoints (new api)

IndexJournal - SegmentNodeStore

- **Open Item**
- Currently does not exist
- To be discussed

Implementation - Done in phases

- Implement this in phased manner
- Phase 1
 - Supported only for DocumentNodeStore
 - Define and implement IndexJournal API
 - Refactor existing IndexEditor to IndexPathCollector and IndexUpdater
 - Support both diff based and journal based flow
- Phase 2
 - Determine if required for SegmentNodeStore setups
 - Implement IndexJournal for SegmentNodeStore



Adobe

MAKE IT AN EXPERIENCE