

**whamcloud**

The logo for Whamcloud features the word "whamcloud" in a bold, lowercase, sans-serif font. A thick blue horizontal line underlines the text. To the right of the text, a blue graphic element consists of a large, stylized letter 'D' that overlaps the end of the underline and extends upwards, with a smaller blue arc above it.

# Improvements for Traversing Large Directory

- Fan Yong  
yong.fan@whamcloud.com

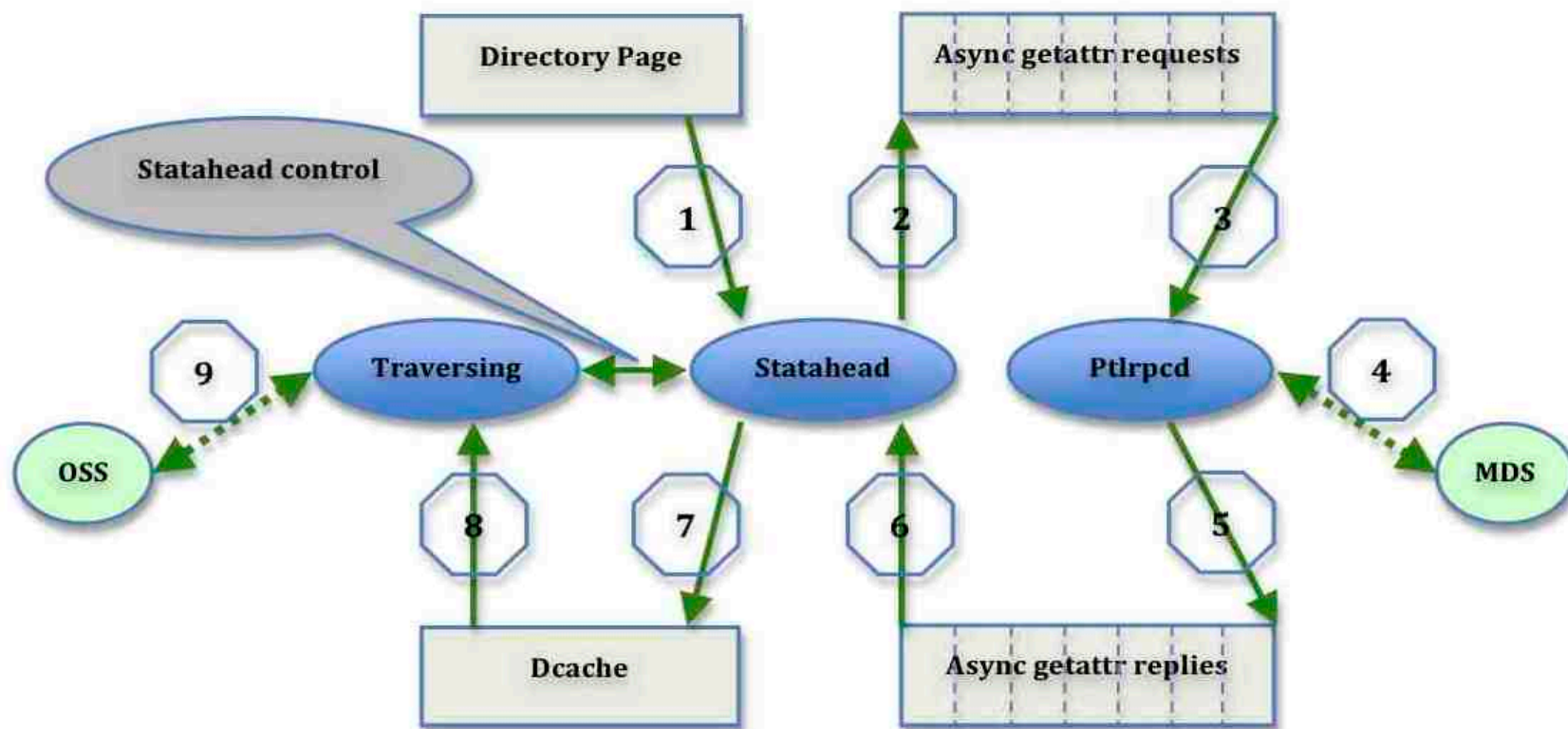
## Overview

Focus on sequential traversing large directory on single client: ls/du/find

- Current statahead
- Improvements
  - Asynchronous glimpse lock
  - Statahead load balance
  - Aggregate attributes to save RPCs
  - Efficient Idlm/object hash
- Test results
- Further works

# Existing statahead

- Only pre-fetch attributes from MDS
  - Single pipeline for async getattr



## Asynchronous glimpse lock

- Prefetch glimpse lock from OST
  - Async RPC without waiting
  - Cache file size/blocks if AGL granted
  - No glimpse callback if OST can't grant AGL
  - Re-enqueued by traversing thread if AGL non-granted/cancelled

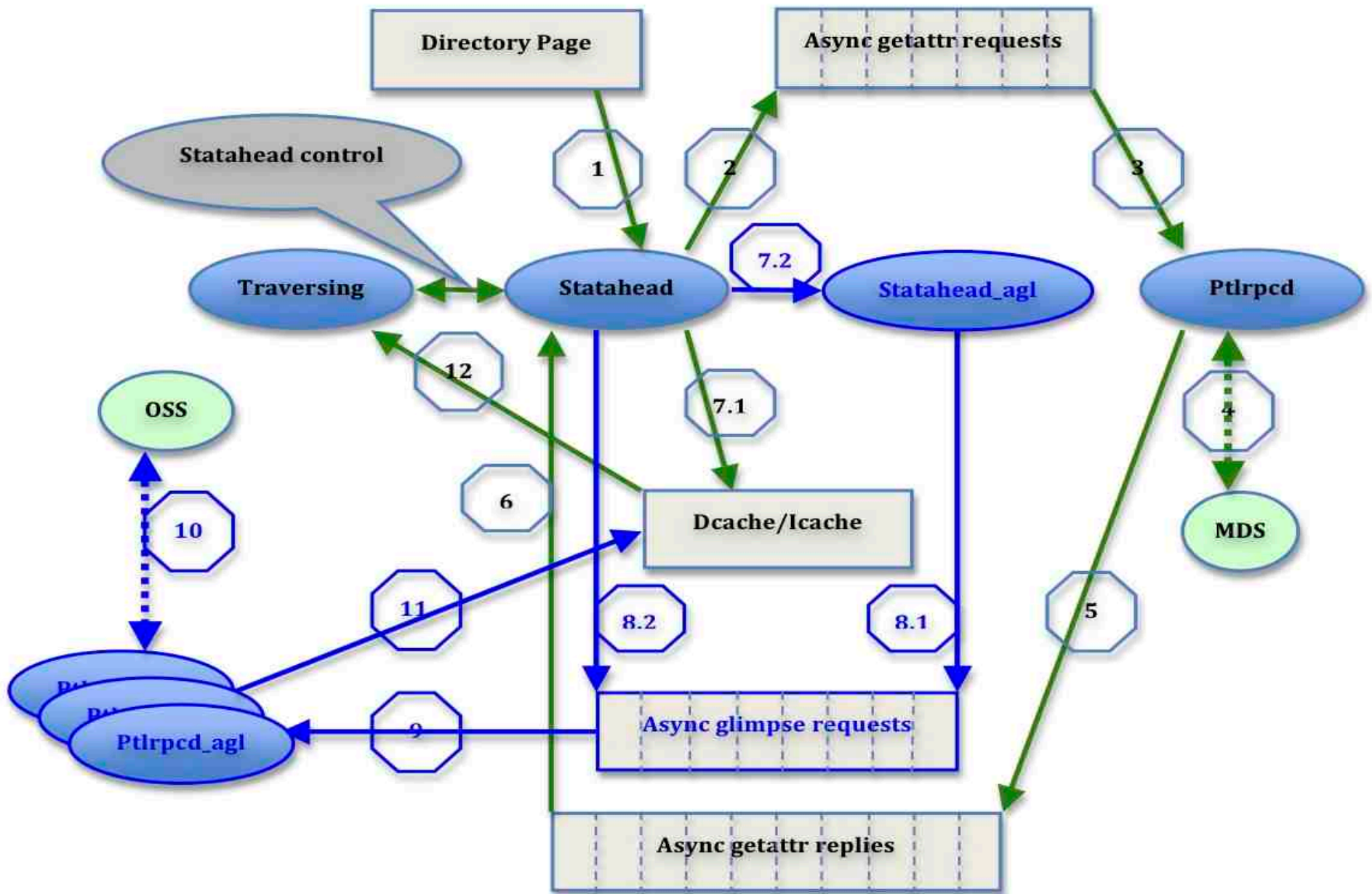
May cause additional RPCs if someone held conflicting locks on some items, but it works well for large directory, since stat/read cases are more common than writes.

## Statahead load balance (1/2)

- Independent portal RPC services (ptlrpcd\_agl)
  - Dedicated to processing AGL RPC
    - Single ptlrpcd on client may be overloaded to process all non-I/O async RPCs when traversing large directory
    - Unbalanced load, other CPUs maybe idle
  - Separate pipelines for async getattr and AGL, decrease unnecessary implied dependence between MDS-side RPC and OSS-side RPC
    - Async RPCs for getattr and glimpse can be processed in parallel (for different files), no contention between ptlrpcd and ptlrpcd\_agl
  - Multiple ptlrpcd\_agl services can scatter AGL loads among more CPUs, and more efficient for multiple-striped cases.

## Statahead load balance (2/2)

- Unbalanced load between statahead and traversing
  - Statahead thread
    - Async getattr pipeline Input & Output
    - AGL pipeline Input
  - Traversing thread
    - AGL pipeline Output
  - AGL is more time-consuming than async getattr for multiple-striped cases. Maybe cause statahead slow down as to RPC processing time can't be hidden.
- Independent LLITE level statahead\_agl thread
  - Statahead\_agl thread
    - AGL pipeline Input
  - Statahead thread
    - Async getattr pipeline Input & Output
    - Help to process AGL pipeline Input when idle
  - Traversing thread
    - AGL pipeline Output





## Aggregate attributes to save RPCs

- MDS-side attributes for traversing directory:
  - 1) Basic attributes, like mode/owner/flags
  - 2) Stripe information
  - 3) Access Control List
  - 4) Default ACL for directory

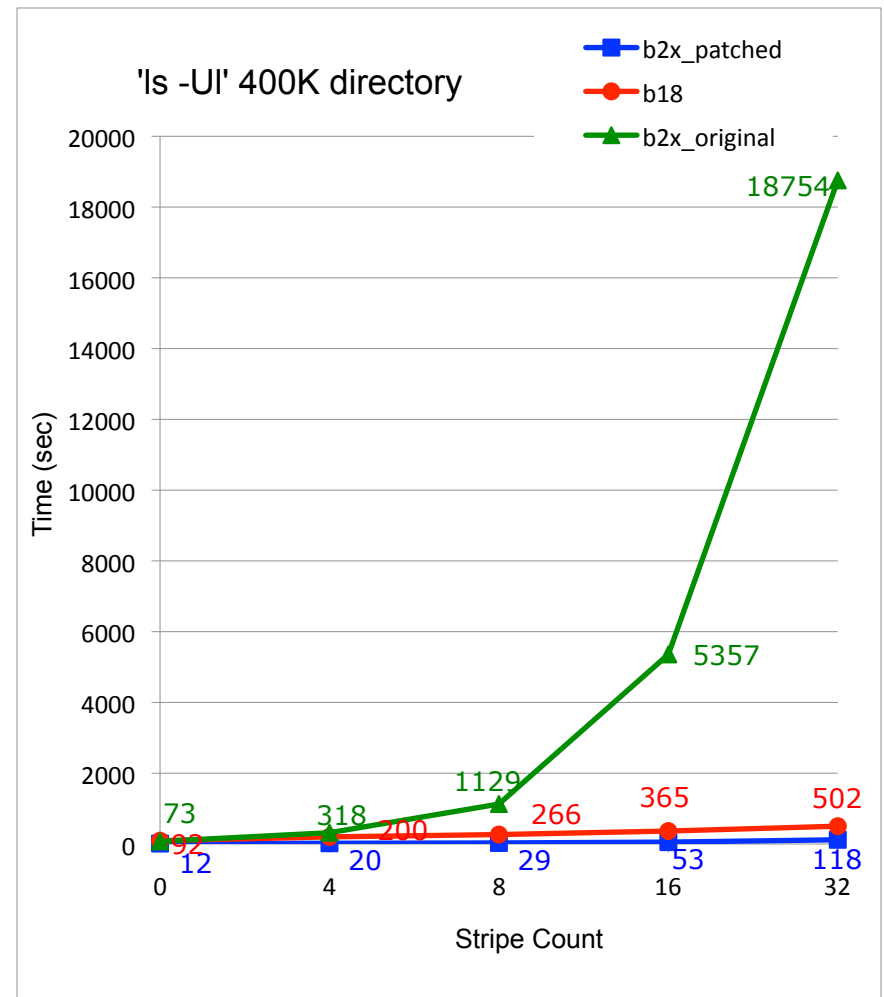
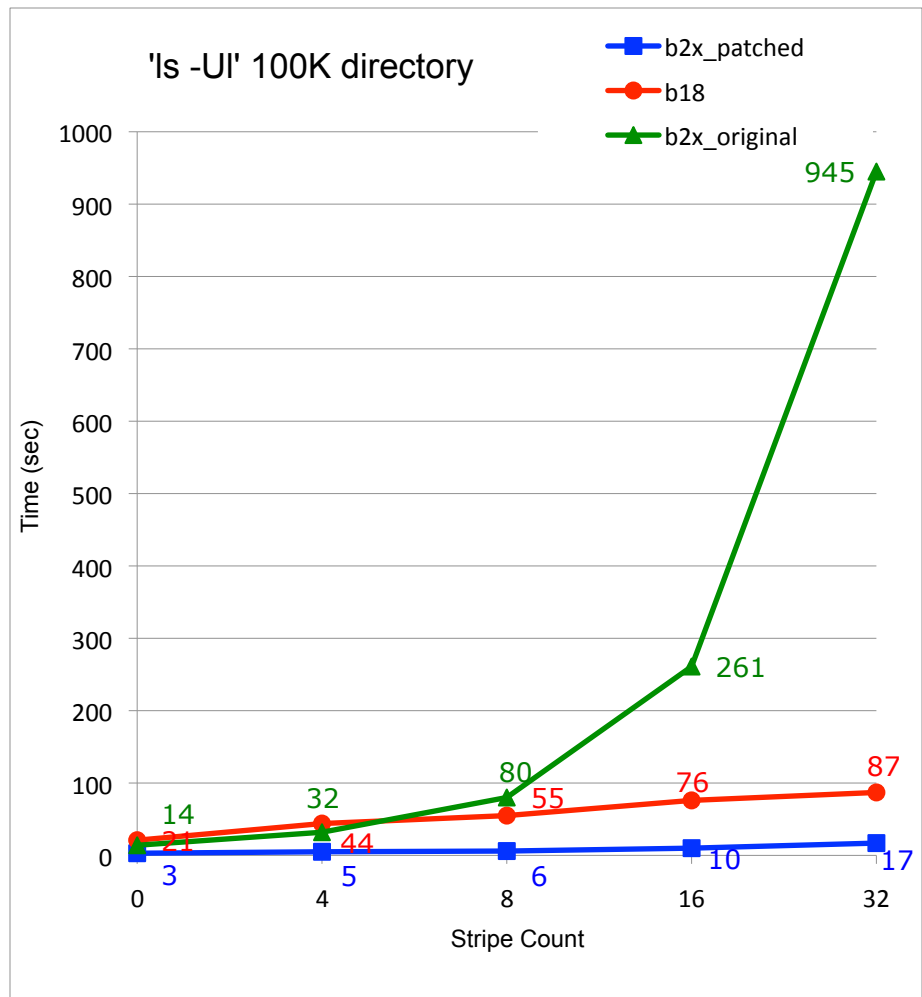
Currently, 1) & 2) & 3) can be obtained through single RPC; and further combining 4) can save about 50% RPC for the best cases.

- Maybe more in future

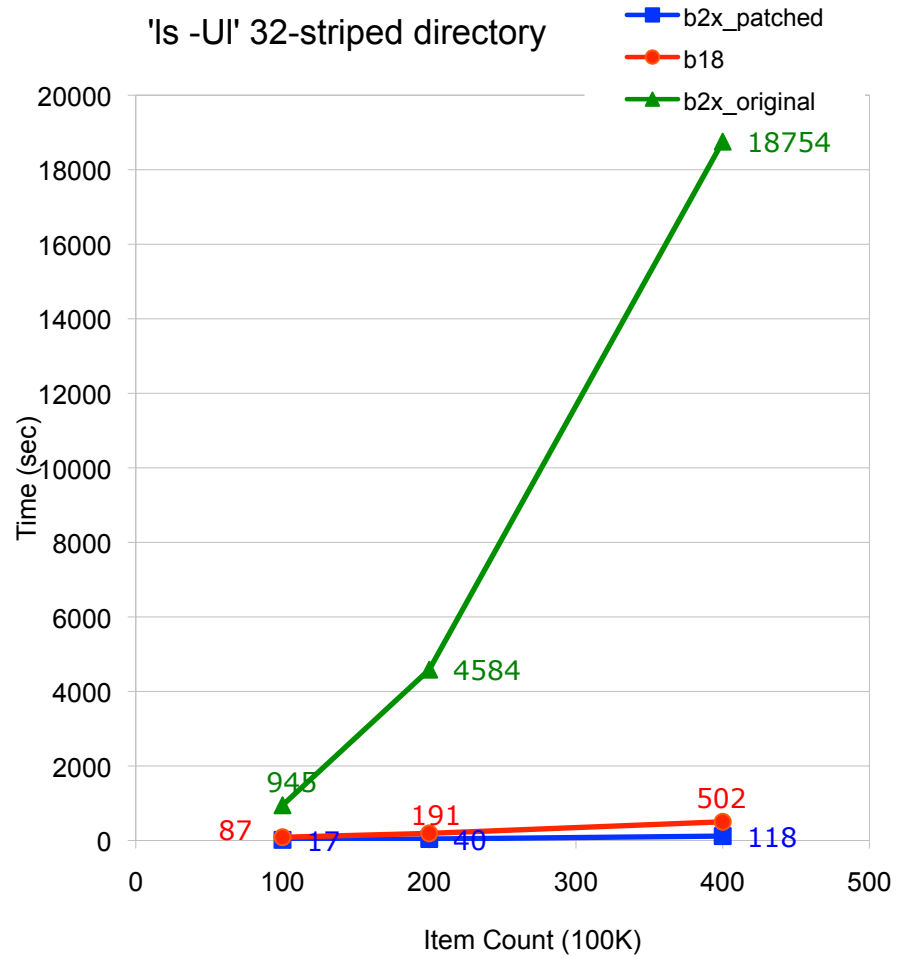
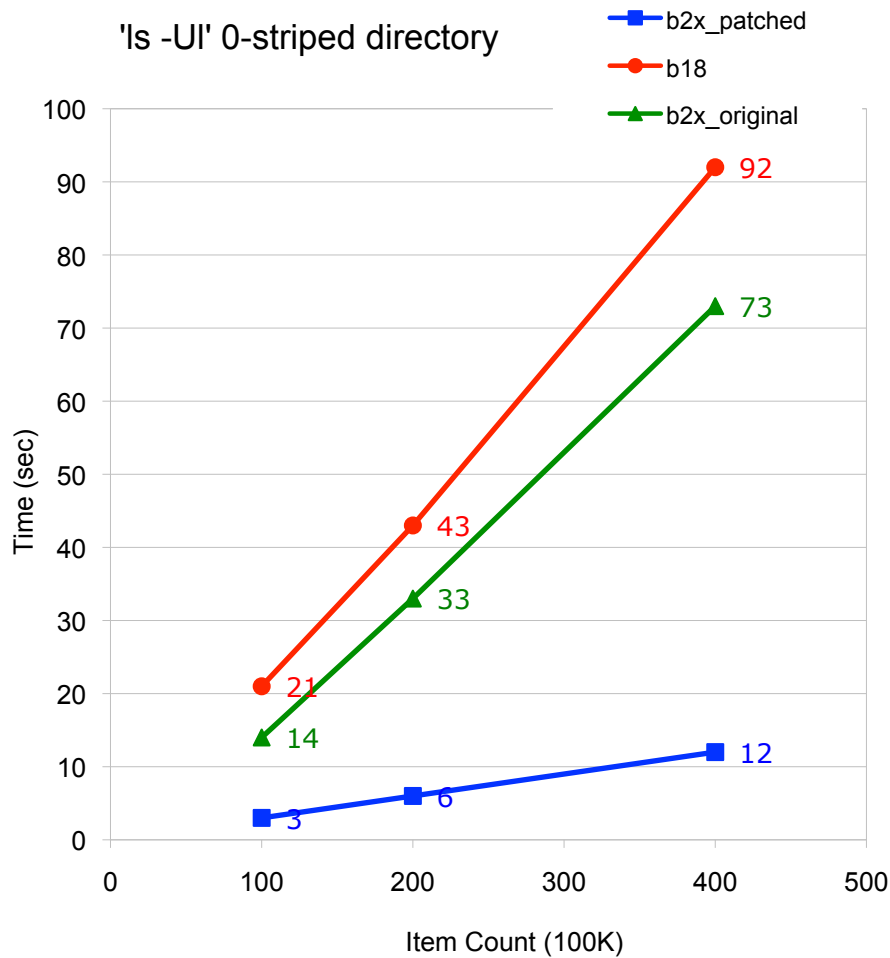
## Efficient Idlm/object hash

- Hash greatly affects traversing performance
  - MDS: object hash & Idlm hash
  - OSS: Idlm hash
  - Client: object hash & Idlm hash
- Hash bucket depth is quite important
  - Basically depends on hash function
    - Scatter Idlms/objects throughout buckets as evenly as possible
  - RAM size and filesystem size are associated
  - Experience: Max depth for millions of Idlms/objects hash
    - lustre-2.0 shows nearly 3K
    - no more than 50 is expected

# 'ls -UI' on single client (1/2)



# 'ls -UI' on single client (2/2)



## Further works

- Ptlrpcd threads pool
  - Common used for kinds of async operations: statahead, AGL, IO
- Aggregate RPC
  - Same type requests in single RPC: getattr/glimpse
  - Similar as readdir+ (mainly) for high latency network
  - Getattr-by-Fid & extent IBITS lock
- Size-on-MDS
  - Global eviction



**Thank You**

- Fan Yong  
yong.fan@whamcloud.com