



Go further, faster®



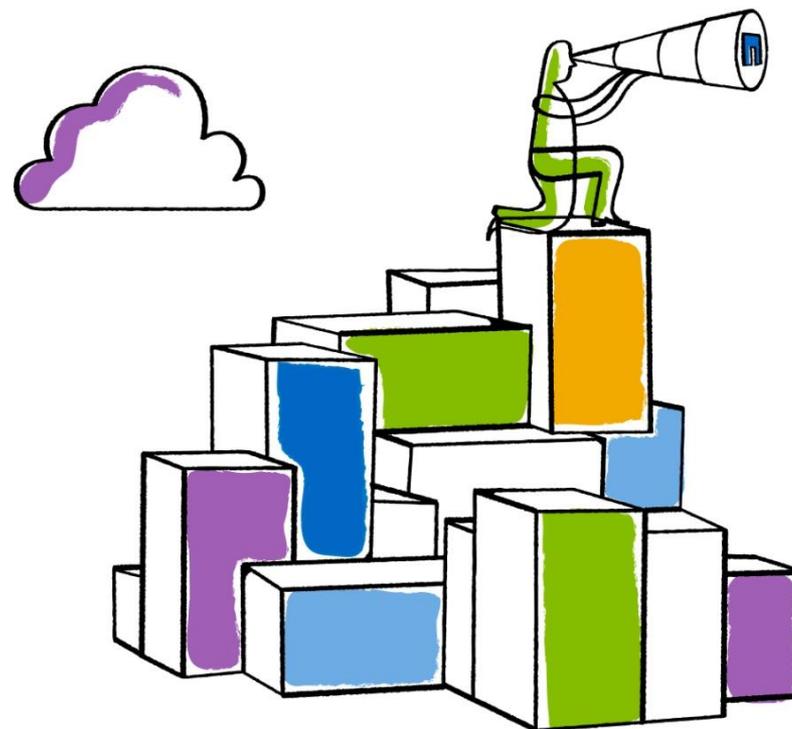
# I/O Hints Discussion

## IETF 81

Mike Eisler

mre dash ietf at eisler dot com

July 27, 2011





# Overview

- draft-eisler-nfsv4-enterprise-apps-01
- Proposes an IO\_ADVICE operation
  - Similar to fadvise()
- Also proposes new READ\_WITH\_ADVICE and WRITE\_WITH\_ADVICE operations
- Several controversies



## Controversies

- Overlaps with draft-hildebrand-nfsv4-fadvise-02.txt
- Need to provide stronger justification/use case for
  - IO\_ADVISE4\_PREFETCH\_OPPORTUNISTIC
  - IO\_ADVISE4\_RECENTLY\_USED
- The need for READ\_WITH\_ADVICE and WRITE\_WITH\_ADVICE



# Overlaps with draft-hildebrand-nfsv4-fadvise-02.txt

- Proposed merged list:
  - IO\_ADVICE4\_NORMAL - as in FADVISE i-d
  - IO\_ADVICE4\_SEQUENTIAL - as in FADVISE i-d
  - IO\_ADVICE4\_RANDOM - as in FADVISE i-d
  - IO\_ADVICE4\_WILLNEED\_TO\_READ - same as FADVISE\_WILLNEED and IO\_ADVICE4\_PREFETCH.
  - IO\_ADVICE4\_DONTNEED - as in FADVISE i-d
  - IO\_ADVICE4\_NOREUSE - as in FADVISE i-d
  - IO\_ADVICE4\_MIGHTNEED\_TO\_READ – same as IO\_ADVICE4\_PREFETCH\_OPPORTUNISTIC
  - IO\_ADVICE4\_WILLNEED\_TO\_WRITE - same as IO\_ADVICE4\_INTENT\_TO\_WRITE
  - IO\_ADVICE4\_RECENTLY\_USED - as in enterprise apps i-d
- To get the IO\_ADVICE4\_SEQUENTIAL\_CACHE behavior, include both IO\_ADVICE4\_SEQUENTIAL and IO\_ADVICE4\_WILLNEED\_TO\_READ in the IO\_ADVICE operation.
- To get the IO\_ADVICE4\_SEQUENTIAL\_DONTCACHE behavior, include both IO\_ADVICE4\_SEQUENTIAL and IO\_ADVICE4\_NOREUSE in the IO\_ADVICE operation.



# Justification for IO\_ADVISE4\_PREFETCH\_OPPORTUNISTIC

- Sometimes one is certain a prefetch is needed (e.g. sequential reads), and other times one speculates it is needed
- IO\_ADVISE4\_PREFETCH is for the certain case
- IO\_ADVISE4\_PREFETCH\_OPPORTUNISTIC is for the speculative case where it costs the server little to perform
  - E.g. an application reads data that contains a reference to data in another block (possibly in another file, possibly in another server)
- A server that is lean on free/cold cache space might prefetch block pointers instead of the block itself



## Justification for IO\_ADVISE4\_RECENTLY\_USED

- Data can go cold in the server's cache while it stays warm in the client's cache
- In order to meet service level objectives including in the face of client restart, the server needs to know which data is warm
- Data that gets LRUed out of server's primary cache (e.g. DRAM) can be placed in secondary cache (e.g. flash memory)



## The need for `READ_WITH_ADVICE` and `WRITE_WITH_ADVICE`

- The objective was to handle the case where the client is indicating advice that applies to just one I/O operation and leaves the `IO_ADVICE` hint intact
- E.g. Overall the file has a random workload, but the client knows when it reads a particular block that the block will be immediately written (e.g. database record update)
  - So server need not cache the block
  - And if the server's file system is log based, this provides advance notice to find free space
- This class of use cases can be handled by doing (for example)

```
IO_ADVICE IO_ADVICE4_WILLNEED_TO_WRITE ;  
  READ ; IO_ADVICE previous_hint
```

  - But this leads to some other issues ...



## New Issues

- How many hints does a server support per open-owner/file pair
  - Very relevant to whether proposal drops `READ/WRITE_WITH_ADVICE`
  - Client should know how many



## New Issues (continued)

- Proposal 1: IO\_ADVISE response to include count of number of hints the server has on the file
  - E.g., a client requests two hints on two non-overlapping byte ranges
    - The second IO\_ADVISE response indicates just one hint (the last hint) is in effect
    - If there a maximum of one hint, then since NORMAL is also a hint, then this means that despite the byte range, the hint always applies to the entire file
  - If the server supports between 2 and  $2^{64}$  hints then the specification needs to define which of the remaining hints apply to orphaned byte ranges
    - Hint with nearest offset?
    - Least recently sent hint?
    - Most recently sent hint?
    - ...



## New Issues (continued)

- Proposal 2: Drop byte range from IO\_ADVISE arguments
  - Unambiguously dictates that the protocol supports exactly one hint per open-owner/open file pair
  - Much simpler, if limiting
    - LAYOUTCOMMIT provides a lesson here
  - Clashes with POSIX standard for fadvise, but does anyone implement multiple byte ranges?



# New Issues (continued)

- pNFS issue: should IO\_ADVISE be allowed on requests to data server (DSes)
  - Very relevant to whether proposal drops READ/WRITE\_WITH\_ADVICE
  - In order to satisfy per I/O hint use case, protocol must allow this possibility
  - Then what does it mean to send IO\_ADVISE to both MDS and DS?
  - This is analogous to COMMIT: MDS decides whether pNFS client can COMMIT to DS.
  - Proposal is to solve it the same way
  - New flag NFL4\_UFLG\_IO\_ADVISE\_THRU\_MDS in the field nfl\_util of the file layout
    - if set, IO\_ADVISE MUST NOT be sent to DS
    - if not set,
      - IO\_ADVISE MAY be sent to DS but will not impact other DSes
      - Hint will not outlive the layout

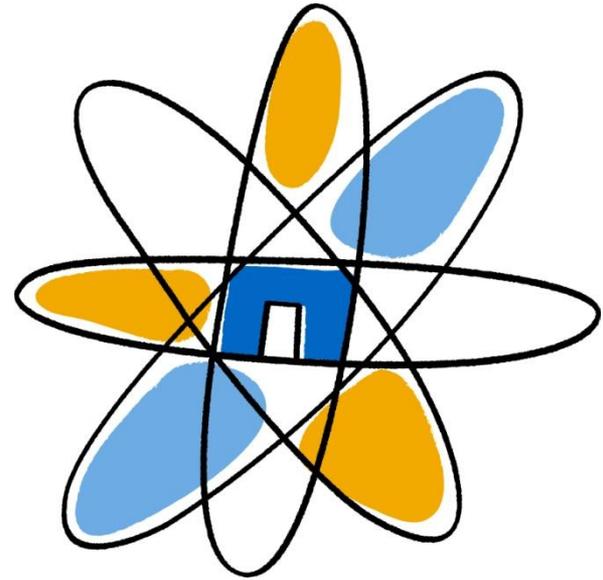


# Proposal for Moving forward

- Since READ/WRITE\_WITH\_ADVICE are contentious, drop those operations and address multiple hint and pNFS issues
- Combine the two I-Ds (just the hint stuff from the enterprise apps I-D), using merged hint list presented earlier
- New operation is called IO\_ADVICE since it supports both POSIX fadvise and non-fadvise requirements
- Incorporate into NFSv4.2



*Thank you*



# Q&A

