

Network File System Version 4
Internet-Draft
Intended status: Standards Track
Expires: 23 October 2026

T. Haynes
Hammerspace
21 April 2026

Deviceid-Scoped Layout Recall for NFSv4.2
draft-haynes-nfsv4-recalldevice-01

Abstract

The Parallel Network File System (pNFS) allows for the metadata server to use CB_LAYOUTRECALL to recall a layout from a client by file id or file system id or all. It also allows the server to use CB_NOTIFY_DEVICEID to delete a deviceid. It does not provide a mechanism for the metadata server to recall all layouts that have a data file on a specific deviceid. This document presents an extension to RFC7862 to allow the server to recall layouts from clients based on deviceid.

Note to Readers

Discussion of this draft takes place on the NFSv4 working group mailing list (nfsv4@ietf.org), which is archived at <https://mailarchive.ietf.org/arch/browse/nfsv4/>. Working Group information can be found at <https://datatracker.ietf.org/wg/nfsv4/about/>.

Status of This Memo

This Internet-Draft is submitted in full conformance with the provisions of BCP 78 and BCP 79.

Internet-Drafts are working documents of the Internet Engineering Task Force (IETF). Note that other groups may also distribute working documents as Internet-Drafts. The list of current Internet-Drafts is at <https://datatracker.ietf.org/drafts/current/>.

Internet-Drafts are draft documents valid for a maximum of six months and may be updated, replaced, or obsoleted by other documents at any time. It is inappropriate to use Internet-Drafts as reference material or to cite them other than as "work in progress."

This Internet-Draft will expire on 23 October 2026.

Copyright Notice

Copyright (c) 2026 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to BCP 78 and the IETF Trust's Legal Provisions Relating to IETF Documents (<https://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect to this document. Code Components extracted from this document must include Revised BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Revised BSD License.

Table of Contents

1. Introduction	2
2. Requirements Language	3
3. Client Capability Advertisement	3
4. Extension to CB_LAYOUTRECALL - Recall Layout from Client . .	4
5. Extraction of XDR	6
6. Security Considerations	7
7. IANA Considerations	7
8. References	7
8.1. Normative References	7
8.2. Informative References	8
Acknowledgments	8
Author's Address	9

1. Introduction

In the Network File System version 4 (NFSv4) with a Parallel NFS (pNFS) metadata server ([RFC8881]), there is no mechanism for the metadata server to recall layouts from the client when a particular deviceid (see Section 3.3.14 of [RFC8881]) either temporarily or permanently is no longer available.

The Flex Files layout type ([RFC8435]) is a primary motivating use case. A Flex Files layout describes data spread across multiple data servers, each identified by a distinct deviceid, which may reside in separate power fault domains.

One use case is when the deviceids in a layout are separated by power fault domains. Each layout might describe 3 different devices, each contained in a different power fault domain. In such a scenario, a single fault domain can have the power removed and not cause the loss of access to the data. However, client I/O will be impacted as the client still has to perform WRITES (see Section 18.32 of [RFC8881])

to the unavailable device, send LAYOUTERRORs (see Section 15.6 of [RFC7862]) to inform the metadata server of NFS4ERR_NXIO (see Section 15.1.16.3 of [RFC8881]).

If the metadata server had the means to recall layouts by deviceid, a lot of this unnecessary traffic could be eliminated. While the metadata server could use LAYOUTRECALL4_ALL to recall all of a client's layouts, that would evict layouts pointing at healthy devices unnecessarily. The deviceid-specific recall is surgical: only the layouts referencing the unavailable device are returned, leaving unaffected layouts in place. Similarly, while the metadata server could recall affected layouts one by one via LAYOUTRECALL4_FILE, that generates one RPC per file and the work can instead be offloaded to the client with a single recall.

Besides the use case above, consider if the metadata server wants to set the NOTIFY4_DEVICEID_DELETE in the CB_NOTIFY_DEVICEID callback (see Section 20.12 of [RFC8881]). This flag cannot be set if a layout is outstanding for a deviceid. While the metadata server can revoke all such layouts, there is no way to know that the client has acknowledged that revocation and hence is still not doing I/O to other data files in the layout. The metadata server could fence those layouts as well (see Section 12.5.5 of [RFC8881]), but that can be an expensive operation.

Using the process detailed in [RFC8178], the revisions in this document become an extension of NFSv4.2 [RFC7862]. They are built on top of the external data representation (XDR) [RFC4506] generated from [RFC7863].

2. Requirements Language

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "NOT RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in BCP 14 [RFC2119] [RFC8174] when, and only when, they appear in all capitals, as shown here.

3. Client Capability Advertisement

Before the server may send a CB_LAYOUTRECALL with LAYOUTRECALL4_DEVICEID, it MUST know that the client supports the new union arm. Per [RFC8178] Section 6, a server MUST NOT send a new callback operation or new union arm to a client that has not indicated support for it.

A client that supports LAYOUTRECALL4_DEVICEID signals this by setting a new flag in the eia_flags field of the EXCHANGE_ID operation (see Section 18.35 of [RFC8881]):

```
///      const EXCHGID4_FLAG_SUPP_RECALL_DEVICEID = 0x02000000;
///
```

The specific bit value 0x02000000 is a placeholder. EXCHGID4 flag bit values are assigned from the reserved flag space defined in Section 18.35 of [RFC8881] and are coordinated across in-flight NFSv4 Working Group extensions. The final value will be confirmed during Working Group Last Call.

A client that sets EXCHGID4_FLAG_SUPP_RECALL_DEVICEID in its EXCHANGE_ID request advertises that it supports handling CB_LAYOUTRECALL with a LAYOUTRECALL4_DEVICEID recall type.

A server MUST NOT send a CB_LAYOUTRECALL with LAYOUTRECALL4_DEVICEID to a client that did not set EXCHGID4_FLAG_SUPP_RECALL_DEVICEID in its EXCHANGE_ID request. A server that does not wish to support this capability MAY ignore the flag.

If the client does not set EXCHGID4_FLAG_SUPP_RECALL_DEVICEID, the server MAY fall back to recalling individual layouts via LAYOUTRECALL4_FILE (one CB_LAYOUTRECALL per layout file that references the unavailable deviceid), as described in Section 20.3 of [RFC8881]. This is less efficient but correct, and preserves interoperability with clients that predate this extension.

4. Extension to CB_LAYOUTRECALL - Recall Layout from Client

The original union layoutrecall4 (see Section 20.3.1 of [RFC8881]) is:

```
enum layoutrecall_type4 {
    LAYOUTRECALL4_FILE = LAYOUT4_RET_REC_FILE,
    LAYOUTRECALL4_FSID = LAYOUT4_RET_REC_FSID,
    LAYOUTRECALL4_ALL   = LAYOUT4_RET_REC_ALL
};

union layoutrecall4 switch(layoutrecall_type4 lor_recalltype) {
    case LAYOUTRECALL4_FILE:
        layoutrecall_file4 lor_layout;
    case LAYOUTRECALL4_FSID:
        fsid4                lor_fsid;
    case LAYOUTRECALL4_ALL:
        void;
};
```

The proposed extension is:

```
///      const LAYOUT4_RET_REC_DEVICEID  = 4;
///
///      enum layoutrecall_type4 {
///          LAYOUTRECALL4_FILE      = LAYOUT4_RET_REC_FILE,
///          LAYOUTRECALL4_FSID      = LAYOUT4_RET_REC_FSID,
///          LAYOUTRECALL4_ALL       = LAYOUT4_RET_REC_ALL,
///          LAYOUTRECALL4_DEVICEID  = LAYOUT4_RET_REC_DEVICEID
///      };
///
/// union layoutrecall4 switch(layoutrecall_type4 lor_recalltype) {
///     case LAYOUTRECALL4_FILE:
///         layoutrecall_file4 lor_layout;
///     case LAYOUTRECALL4_FSID:
///         fsid4                lor_fsid;
///     case LAYOUTRECALL4_DEVICEID:
///         deviceid4            lor_deviceid;
///     case LAYOUTRECALL4_ALL:
///         void;
///     };
```

Note that LAYOUT4_RET_REC_* constants are shared between layoutrecall_type4 (used in CB_LAYOUTRECALL) and layoutreturn_type4 (used in LAYOUTRETURN, see Section 18.44.1 of [RFC8881]). Adding LAYOUT4_RET_REC_DEVICEID = 4 therefore also extends layoutreturn_type4 with a LAYOUTRETURN4_DEVICEID value. A client that receives a LAYOUTRETURN with lor_recalltype set to LAYOUTRETURN4_DEVICEID and does not recognise it MUST return NFS4ERR_UNION_NOTSUPP per [RFC8178].

The server MUST NOT send CB_LAYOUTRECALL with LAYOUTRECALL4_DEVICEID to a client that has not set EXCHGID4_FLAG_SUPP_RECALL_DEVICEID (see Section 3). This satisfies the requirement in Section 6 of [RFC8178] that a server establish client awareness before sending new callback extensions.

The existing clora_iomode field in CB_LAYOUTRECALL4args (see Section 20.3.1 of [RFC8881]) applies normally: the client MUST return all layouts matching both the given deviceid and the given iomode. The server can determine that the client no longer has any layouts with the given deviceid and iomode once the client replies with NFS4ERR_NOMATCHING_LAYOUT.

5. Extraction of XDR

This document contains the external data representation (XDR) [RFC4506] description of the extension to CB_LAYOUTRECALL. The XDR description is presented in a manner that facilitates easy extraction into a ready-to-compile format. To extract the machine-readable XDR description, use the following shell script:

```
<CODE BEGINS>
#!/bin/sh
grep '^ *///' $* | sed 's?^ */// ??' | sed 's?^ *///$??'
<CODE ENDS>
```

For example, if the script is named 'extract.sh' and this document is named 'spec.txt', execute the following command:

```
<CODE BEGINS>
sh extract.sh < spec.txt > recalldevice.x
<CODE ENDS>
```

This script removes leading blank spaces and the sentinel sequence '///' from each line. XDR descriptions with the sentinel sequence are embedded throughout the document.

Note that the XDR code contained in this document depends on types from the NFSv4.2 `nfs4_prot.x` file (generated from [RFC7863]). This includes both nfs types that end with a 4, such as `offset4`, `length4`, etc., as well as more generic types such as `uint32_t` and `uint64_t`.

The extracted XDR extends types that are already defined in [RFC7863] rather than introducing new types. An implementer MUST integrate the extracted block into the base XDR by *replacing* the existing `layoutrecall_type4` enum declaration and `layoutrecall4` union declaration with the extended versions in this document, and by adding the `LAYOUT4_RET_REC_DEVICEID` constant in the same file area as the existing `LAYOUT4_RET_REC_*` constants. Appending the extracted block verbatim to the RFC 7863 XDR produces duplicate- type errors at compile time.

The `EXCHGID4_FLAG_SUPP_RECALL_DEVICEID` constant is purely additive and is added to the `EXCHGID4` flag constant block without replacing anything.

6. Security Considerations

The extension introduced by this document inherits the Security Considerations of CB_LAYOUTRECALL in Section 20.3 of [RFC8881] and of NFSv4.2 in [RFC7862]. Three aspects of the deviceid-scoped recall warrant explicit mention.

Blast radius: A single CB_LAYOUTRECALL with LAYOUTRECALL4_DEVICEID can cause a client to return many layouts in one callback exchange. This is not a new risk category -- LAYOUTRECALL4_FSID and LAYOUTRECALL4_ALL already allow wider-scoped recalls -- but administrators sizing callback-channel timeouts and client-side recall-handler resources SHOULD account for the possibility that a deviceid-scoped recall triggers per-layout cleanup work proportional to the number of layouts that reference the deviceid.

Cross-client scope: In deployments where a single deviceid is shared across clients serving different tenants or applications, a deviceid-scoped recall is cross-client by construction: every client that holds a layout referencing the deviceid is recalled. This is the same property as LAYOUTRECALL4_FSID and does not introduce a new information-exposure surface. Deployments that prefer per-tenant failure-domain isolation MAY treat this extension as an operational benefit: by assigning distinct deviceids per tenant, deviceid-scoped recall becomes a per-tenant operation with no cross-tenant effect.

Capability negotiation as defence in depth: The capability flag defined in Section 3 ensures that a server cannot silently change the callback wire format for a client that did not opt in. A client that does not advertise EXCHGID4_FLAG_SUPP_RECALL_DEVICEID continues to receive only the layoutrecall_type4 values defined in [RFC8881]. A server that sends LAYOUTRECALL4_DEVICEID to a client that did not advertise support is in violation of Section 6 of [RFC8178]; clients that encounter such a violation MUST return NFS4ERR_UNION_NOTSUPP and SHOULD log the server for operator attention.

7. IANA Considerations

This document has no IANA actions.

8. References

8.1. Normative References

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, DOI 10.17487/RFC2119, March 1997, <<https://www.rfc-editor.org/rfc/rfc2119>>.
- [RFC4506] Eisler, M., Ed., "XDR: External Data Representation Standard", STD 67, RFC 4506, DOI 10.17487/RFC4506, May 2006, <<https://www.rfc-editor.org/rfc/rfc4506>>.
- [RFC7862] Haynes, T., "Network File System (NFS) Version 4 Minor Version 2 Protocol", RFC 7862, DOI 10.17487/RFC7862, November 2016, <<https://www.rfc-editor.org/rfc/rfc7862>>.
- [RFC7863] Haynes, T., "Network File System (NFS) Version 4 Minor Version 2 External Data Representation Standard (XDR) Description", RFC 7863, DOI 10.17487/RFC7863, November 2016, <<https://www.rfc-editor.org/rfc/rfc7863>>.
- [RFC8174] Leiba, B., "Ambiguity of Uppercase vs Lowercase in RFC 2119 Key Words", BCP 14, RFC 8174, DOI 10.17487/RFC8174, May 2017, <<https://www.rfc-editor.org/rfc/rfc8174>>.
- [RFC8178] Noveck, D., "Rules for NFSv4 Extensions and Minor Versions", RFC 8178, DOI 10.17487/RFC8178, July 2017, <<https://www.rfc-editor.org/rfc/rfc8178>>.
- [RFC8881] Noveck, D., Ed. and C. Lever, "Network File System (NFS) Version 4 Minor Version 1 Protocol", RFC 8881, DOI 10.17487/RFC8881, August 2020, <<https://www.rfc-editor.org/rfc/rfc8881>>.

8.2. Informative References

- [RFC1813] Callaghan, B., Pawlowski, B., and P. Staubach, "NFS Version 3 Protocol Specification", RFC 1813, DOI 10.17487/RFC1813, June 1995, <<https://www.rfc-editor.org/rfc/rfc1813>>.
- [RFC8435] Halevy, B. and T. Haynes, "Parallel NFS (pNFS) Flexible File Layout", RFC 8435, DOI 10.17487/RFC8435, August 2018, <<https://www.rfc-editor.org/rfc/rfc8435>>.

Acknowledgments

Trond Myklebust and Paul Saab were involved in the initial requirements for this functionality.

Brian Pawlowski and Gorrry Fairhurst helped guide this process.

Internet-Draft

RECALL_DEVICE

April 2026

Author's Address

Thomas Haynes
Hammerspace
Email: loghyr@gmail.com