

OrangeFS

Hochleistungs-Ein-/Ausgabe

Michael Kuhn

Wissenschaftliches Rechnen
Fachbereich Informatik
Universität Hamburg

2016-05-06

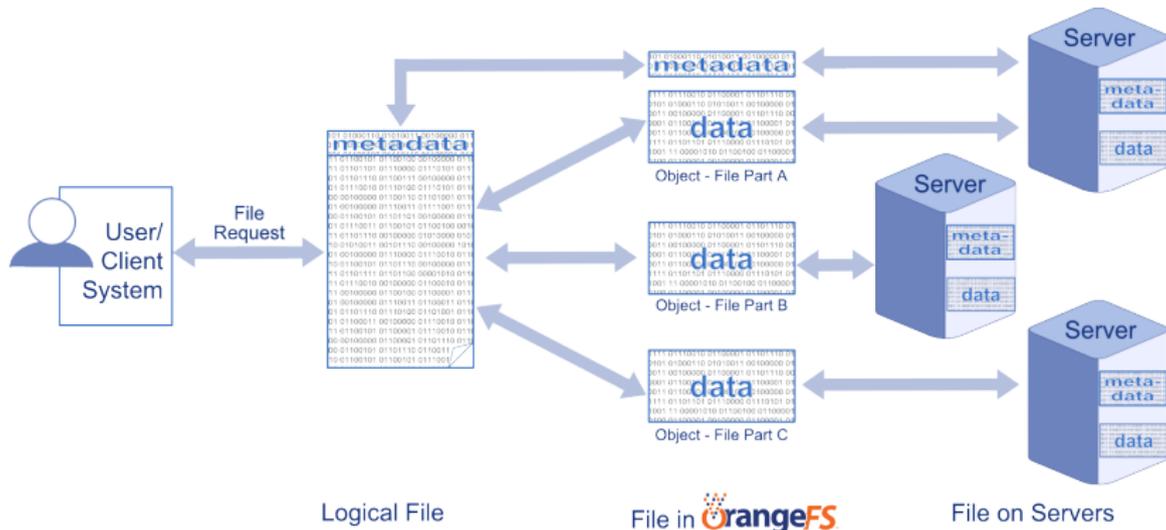
- Die SSD-Unterstützung wird dadurch realisiert, dass die Pfade für Daten und Metadaten separat konfiguriert werden können
- Bis Version 2.9 wurden alle Einträge eines Verzeichnisses von einem Metadatenserver verwaltet
 - Die eigentlichen Daten waren schon immer verteilt

- OrangeFS hat unter anderem Unterstützung für kollektive und nicht-zusammenhängende Ein-/Ausgabe

- OrangeFS wird mit folgenden Optionen konfiguriert:
 - Unterstützung für TCP und InfiniBand
 - Auf west[1-10] läuft jeweils ein Daten- und ein Metadatenserver
 - Daten und Metadaten werden im selben Verzeichnis gespeichert
 - Das Verzeichnis kann sich in einem beliebigen lokalen Dateisystem befinden
 - Es wird eine Protokolldatei geschrieben
- c_lush erlaubt das Ausführen von Befehlen auf mehreren Knoten
 - Mit dem mkfs-Parameter wird das Verzeichnis für Daten und Metadaten initialisiert
 - Mit dem rmfs-Parameter wird das Verzeichnis wieder gelöscht

- Die Metadaten einer Datei befinden sich auf einem einzelnen Metadatenserver
 - Unter Umständen müssen aber mehrere Metadatenserver kontaktiert werden, um die Datei zu finden

Funktionsweise... [2]



- Die logische Datei besteht aus Metadaten und Daten
- Die physikalische Datei besteht aus einem Metafile und drei Datafiles
 - Der erste Server enthält sowohl das Metafile als auch das erste Datafile
 - Die restlichen Server enthalten jeweils ein Datafile

- Um eine Datei zu finden, müssen möglicherweise mehrere Metadatenserver kontaktiert werden
 - Zuerst wird das Wurzelverzeichnis durchsucht, dann die nächste Pfadkomponente etc.

Schnittstellen

- Zugriff wird in drei Ebenen eingeteilt
 - C-Bibliothek: `fopen`, `fread` etc.
 - POSIX-Bibliothek: `open`, `read` etc.
 - Funktionieren mit OrangeFS- und POSIX-Dateisystemen
 - System-Call-Bibliothek: `pvfs_open`, `pvfs_read` etc.
 - Funktioniert nur mit OrangeFS-Dateisystemen
- Das Direct Interface kann explizit oder implizit genutzt werden
 - Dafür stehen vier Bibliotheken zur Verfügung:
`liborangefsposix`, `liborangefs`, `libofs` und `libpvfs2`

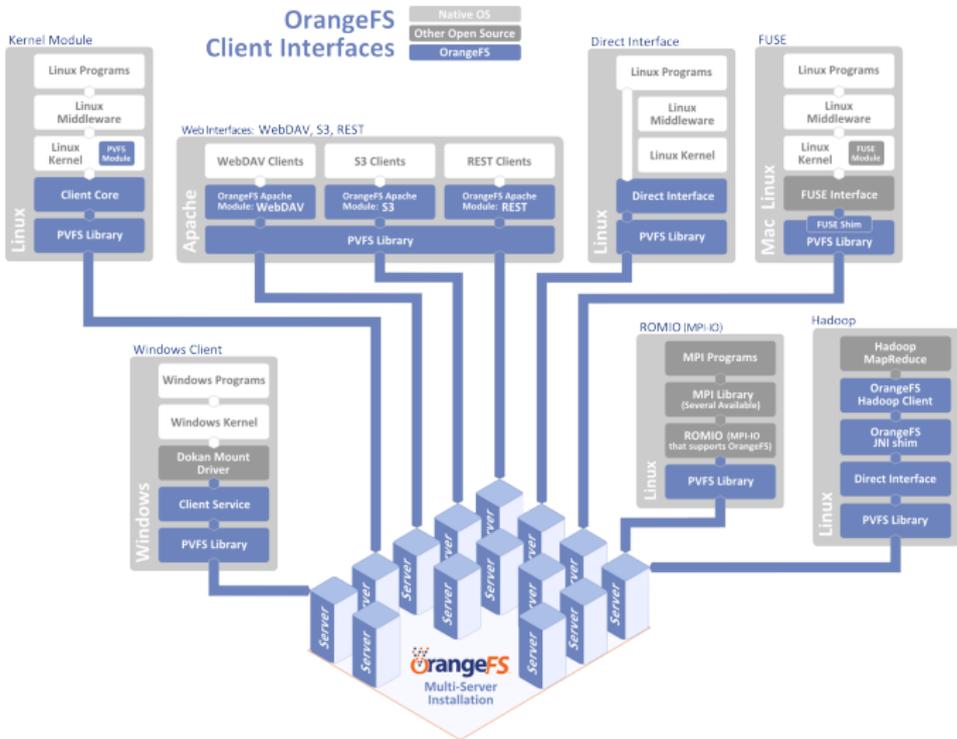
- Explizite Nutzung bedeutet, dass die Anwendungen gegen die entsprechenden Bibliotheken gelinkt werden
- Implizite Nutzung bedeutet, dass die entsprechenden Bibliotheken über den Preload-Mechanismus geladen werden

Schnittstellen...

- `liborangefsposix`
 - Deckt die C- und POSIX-Bibliotheken ab
 - Eine Zusammenfassung der anderen Bibliotheken
- `liborangefs`
 - Deckt die System-Call-Bibliothek ab
 - Enthält POSIX-PVFS-Funktionen
 - `pvfs_open`, `pvfs_read`, `pvfs_write`, ...
 - Verhalten sich wie die entsprechenden POSIX-Funktionen



Schnittstellen... [2]



Schnittstellen...

- POSIX-Schnittstellen erlauben auch Zugriff auf OrangeFS-Funktionalität
- Geregelt über spezielle Flags beim open-Aufruf
 - Für `pvfs_open` und `open` aber nicht für `fopen`
 - Beispiele:
 - Verteilungsfunktion (`PVFS_HINT_DISTRIBUTION_NAME`),
 - Anzahl der Datafiles (`PVFS_HINT_DFILE_COUNT_NAME`),
 - Caching (`PVFS_HINT_CACHE_NAME`) etc.
- Meistens nur Auswirkungen beim Erstellen

Verteilungsfunktionen

- Unterstützung für vier Verteilungsfunktionen
- `basic_dist`
 - Speicherung in einem Datafile
 - Eventuell nützlich für kleine Dateien
 - Keine Parameter einstellbar
- `simple_stripe` (Standard)
 - Round Robin mit fester Streifenbreite
 - Entspricht der Verteilungsfunktion in Lustre
 - Parameter:
 - `strip_size`: Streifenbreite (64 KiB)

Verteilungsfunktionen...

- twod_stripe
 - Verteilung nach 2D-Muster über Gruppen von Datafiles
 - Erlaubt Einteilung von Servern in Gruppen
 - Parameter:
 - num_groups: Anzahl Datafile-Gruppen
 - strip_size: Streifenbreite
 - group_strip_factor: Anzahl Streifen pro Server und Gruppe bevor nächste Gruppe benutzt wird
- varstrip_dist
 - Round Robin mit variabler Streifenbreite
 - Erlaubt Anpassung an spezielle Datenformate
 - Parameter:
 - strips: Streifenbreiten pro Server
 - Beispiel: 0:32K; 1:128K; 2:64K; 3:128K;

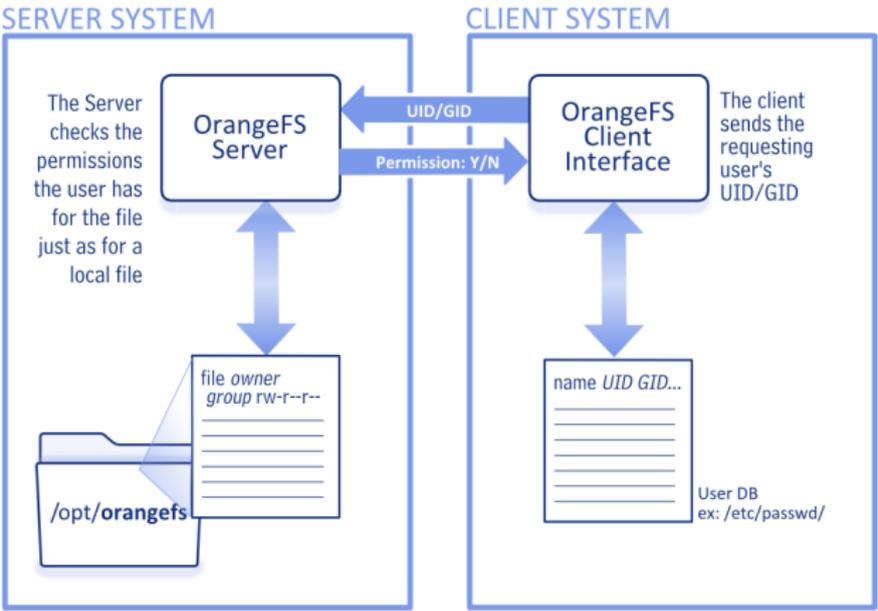
- Die `twod_stripe`-Verteilungsfunktion kann genutzt werden, um die Anzahl der zu kontaktierenden Server einzuschränken
 - Beispiel: n Clients greifen gemeinsam auf eine Datei der Größe m zu, dabei ist jeder Client für m/n zusammenhängende Bytes verantwortlich
 - Bei entsprechender Wahl der Parameter kann dafür gesorgt werden, dass jeder Client nur auf eine Untermenge der verfügbaren Server zugreift, insgesamt aber auf alle Server zugegriffen wird
- Die `varstrip_dist`-Verteilungsfunktion kann beispielsweise genutzt werden, um sicherzustellen, dass logisch zusammenhängende Teile einer Datei auf einem Datenserver gespeichert werden

Sicherheit

- Unterstützung für drei Sicherheitsmechanismen
 - Standard, schlüssel-basiert und zertifikat-basiert mit LDAP
 - Alle Mechanismen arbeiten mit Timeouts
- Zugriffskontrolle wird über Credentials geregelt
 - Standardmäßig senden Clients Credentials
 - Server überprüft Berechtigungen auf Basis der Credentials

Sicherheit... [2]

Default Security



Sicherheit...

- Benötigt keinerlei Konfiguration
 - Einfache Installation und Konfiguration
 - Gut geeignet für Evaluationen und Tests
- Hohe Leistung
 - Keine zusätzlichen Abfragen und Dienste notwendig
- Keine hohe Sicherheit
 - Clients können beliebige Credentials senden
- Auch bei anderen Dateisystemen häufig anzutreffen
 - Lustre arbeitet in der Standardkonfiguration ähnlich

Sicherheit...

- Benutzt private und öffentliche Schlüssel zur Authentifizierung
 - Jeder Server und Client hat eigenes Schlüsselpaar
 - Server nutzen einen Schlüsselspeicher, der alle öffentlichen Schlüssel enthält
- Credentials werden signiert
 - Server prüft zusätzlich Signatur
- Schlüsselspeicher ist schwierig zu handhaben
 - Muss bei Änderungen neu generiert und verteilt werden
 - Server müssen neugestartet werden

Sicherheit...

- Alle Server teilen sich eine Zertifizierungsstelle
 - Zertifizierungsstelle stellt weitere Zertifikate aus
- Jeder Benutzer hat eigenes Zertifikat
 - Wird im Home-Verzeichnis gespeichert
 - Zuordnung von Zertifikat zu Benutzer- und Gruppen-ID mit Hilfe von LDAP
- Höherer Aufwand als andere Mechanismen
 - LDAP-Installation falls noch nicht vorhanden
 - Komplexere und zusätzliche Schritte notwendig
 - Allerdings keine Server-Neustarts notwendig

Schichten

- Job Manager
 - Operationen bestehen aus mehreren Schritten (Jobs)
 - Koordination aller Jobs und Thread-Verwaltung
- Flows
 - Bezeichnet abstrakte Datenflüsse
 - Koordination mit anderen Schichten
- BMI
 - Abstraktion des Netzwerks (Buffered Message Interface)
 - Verwaltung der Netzwerkkommunikation
- Trove
 - Abstraktion der Speichers
 - Verwaltung von Key-Value-Paaren und Bytestreams

Operationen

- Zwei Varianten von allen Operationen
 - PVFS_sys_op und PVFS_i sys_op
 - Synchron und asynchron
- Eigentliche Funktionalität asynchron implementiert
 - Synchrone Variante ruft asynchrone auf
 - PVFS_sys_wait und PINT_sys_release
- Vielzahl an structs und unions
 - Generische Datenstrukturen für alle Operationen

- PINT_smc_b: State machine control block
- PINT_client_sm: Client state machine
- PINT_smc_b_alloc: Allokiert einen neuen Control Block
 - Dabei wird festgelegt welchen Typ die State Machine hat
- PINT_sm_frame: Ein Control Block kann mehrere Frames enthalten
- Die Client State Machine enthält eine Union, die die Informationen für die eigentliche Operation enthält
 - Kann somit für alle Operationen benutzt werden
- Abschließend wird die Client State Machine gepostet, d.h. zur Ausführung gebracht

State Machines

- Operationen werden als State Machines abgearbeitet
 - .sm-Dateien enthalten Beschreibung
 - statecomp-Werkzeug generiert .c-Dateien aus .sm-Dateien
 - Eigener Parser und Scanner
- State Machines bestehen aus Zuständen und Übergängen
 - Immer in genau einem Zustand
 - In jedem Zustand wird eine Funktion oder geschachtelte State Machine ausgeführt
 - Rückgabewert bestimmt Übergang
 - Beginn mit erstem Zustand und Ende mit terminate- bzw. return-Zustand

- Der `terminate`-Zustand beendet die Abarbeitung der State Machine
- Der `return`-Zustand kehrt aus einer geschachtelten State Machine zurück

- Die Client State Machine für die create-Operation startet im `init`-Zustand
 - Darin wird die `create_init`-Funktion ausgeführt
 - Danach wird ein Übergang in den `parent_getattr`-Zustand vollzogen
- Im `parent_getattr`-Zustand wird die geschachtelte State Machine `pvfs2_client_getattr_sm` ausgeführt
 - Tritt kein Fehler auf, erfolgt ein Übergang in den `parent_getattr_inspect`-Zustand
 - Ansonsten wird in den `cleanup`-Zustand übergegangen

- Die Client State Machine enthält eine dedizierte Variable für geschachtelte State Machines vom Typ `getAttr`

State Machines... [1]

```
1 machine pvfs2_client_create_sm
2 {
3     ...
4     state cleanup
5     {
6         run create_cleanup;
7         CREATE_RETRY => init;
8         default => terminate;
9     }
10 }
```


- Über `error_code` wird der Rückgabewert des Zustands gesetzt
 - Dieser Rückgabewert wird zur Bestimmung des Übergangs genutzt
- Der eigentliche Rückgabewert der Funktion gibt an, ob die Ausführung zurückgestellt (`SM_ACTION_DEFERRED`), vollständig ist (`SM_ACTION_COMPLETE`) oder terminiert werden soll (`SM_ACTION_TERMINATE`)

Message Pairs

- Kommunikation über Message Pairs
 - Anfrage durch Client, Antwort durch Server
- Abarbeitung durch spezielle State Machine
 - Name: pvfs2_msgpairarray_sm
 - Kümmert sich um Senden und Empfangen
 - Fehlgeschlagene Operationen werden wiederholt
 - Zusätzlich En-/Decoding der Nachrichten
 - Aufruf eines definierbaren Callbacks

- Dem Callback wird die Antwort des Servers übergeben
 - Dazu wird auf die Eltern-State-Machine zugegriffen und die Daten dorthin kopiert

- 1** OrangeFS
 - Orientierung
 - Einleitung
 - Installation und Konfiguration
 - Funktionsweise
 - Schnittstellen
 - Verteilungsfunktionen
 - Sicherheit
 - Interne Funktionsweise
 - Zusammenfassung

2 Quellen

Quellen I

- [1] OrangeFS Development Team. OrangeFS.
<http://www.orangefs.org/>.
- [2] OrangeFS Development Team. OrangeFS Documentation.
<http://docs.orangefs.com/>.
- [3] OrangeFS Development Team. OrangeFS Wiki.
<http://www.orangefs.org/trac/orangefs>.
- [4] Shuangyang Yang, Walter B Ligon III, and Elaine C Quarles. Scalable Distributed Directory Implementation on Orange File System. *7th IEEE International Workshop on Storage Network Architecture and Parallel I/O*, 2011.