

OrangeFS

Hochleistungs-Ein-/Ausgabe

Michael Kuhn

Wissenschaftliches Rechnen
Fachbereich Informatik
Universität Hamburg

2017-05-12



Universität Hamburg

DER FORSCHUNG | DER LEHRE | DER BILDUNG



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

2 Quellen

Historie

- 1993: PVFS Version 0 als NASA-Projekt
- 1994: PVFS Version 1 mit TCP-Unterstützung
- 1997: Veröffentlichung als Open Source
- 2003: Veröffentlichung von PVFS Version 2
- 2008: Start von OrangeFS als Entwicklungszweig
- 2010: OrangeFS ersetzt PVFS als Hauptversion
- 2011: Unterstützung für Windows-Clients
- 2016: Kernel-Client in Linux 4.6 hinzugefügt

- 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

Installation

■ Sehr einfache Benutzung und Installation

```

1 $ {apt,dnf,yum} install bison flex libdb-dev[el]
   ↳ libattr-dev[el]
2
3 $ wget .../orangefs-x.y.z.tar.gz
4 $ tar xf orangefs-x.y.z.tar.gz
5 $ cd orangefs-x.y.z
6
7 $ ./configure --prefix=${PREFIX} --enable-shared
8 $ make --jobs=$(nproc)
9 $ make install

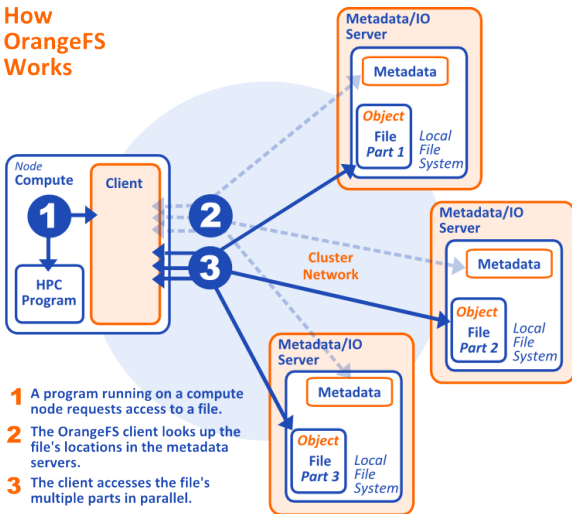
```


- 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



Funktionsweise [2]

How OrangeFS Works



- 1 A program running on a compute node requests access to a file.
- 2 The OrangeFS client looks up the file's locations in the metadata servers.
- 3 The client accesses the file's multiple parts in parallel.

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

Funktionsweise...

- OrangeFS ist objekt-basiert
 - Vier Hauptobjekttypen: Datafile, Metafile, Directory, Symlink
- Dateien/Verzeichnisse bestehen aus mindestens zwei Objekten
 - Eines für Metadaten und die restlichen für Daten
- Server arbeiten mit Objekten, Bytestreams und Key-Value-Paaren
 - Objekte haben eindeutige Handles
 - Bytestreams für Daten
 - Key-Value-Paare für Attribute und Metadaten
 - Objekte nutzen Bytestreams und/oder Key-Value-Paare



Funktionsweise...

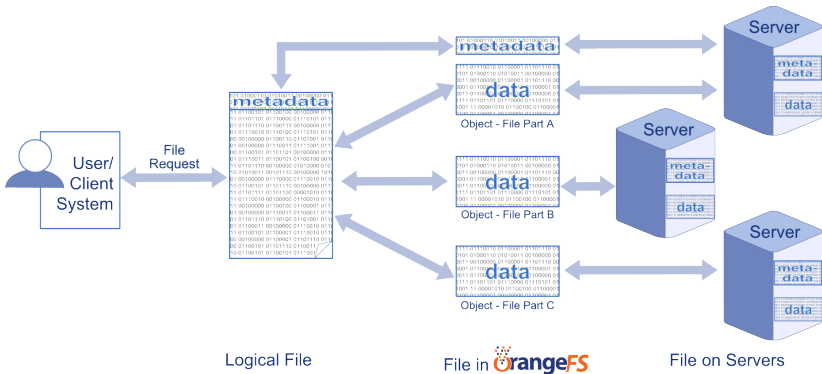
- Key-Value-Paare werden mit der Berkeley DB verwaltet
 - Alternativ mit LMDB
- Bytestreams werden im lokalen Dateisystem gespeichert

```
1 $ ls -lh
2 drwxr-xr-x 4 wr wr 4.0K Jun 25 20:45 3b3f08ea
3 -rw----- 1 wr wr 8.0K Jun 25 20:45 collections.db
4 -rw----- 1 wr wr 8.0K Jun 25 20:45 storage_attributes.db
5 $ ls -lh 3b3f08ea/
6 drwxr-xr-x 66 wr wr 4.0K Jun 25 20:45 bstreams
7 -rw----- 1 wr wr 8.0K Jun 25 20:45 collection_attributes.db
8 -rw----- 1 wr wr 992K Jun 28 19:42 dataspace_attributes.db
9 -rw----- 1 wr wr 16 Jun 25 20:45 __db.001
10 -rw----- 1 wr wr 228K Jun 27 18:38 keyval.db
11 drwxr-xr-x 2 wr wr 4.0K Jun 28 19:42 stranded-bstreams
12 $ ls -lh 3b3f08ea/bstreams/00000021/
13 -rw----- 1 wr wr 416M Jun 28 22:48 4666666666665515.bstream
```


Funktionsweise...

- Metafilés repräsentieren Dateien
 - Speichern Metadaten wie Eigentümer und Berechtigungen
 - Enthalten außerdem Handles mehrerer Datafiles (inklusive deren Verteilung)
 - Größe wird nicht explizit gespeichert
- Datafiles repräsentieren Daten
 - Verteilt über alle Server
 - Keine Metadaten, nur eigentliche Dateiinhalte
- Directories repräsentieren Verzeichnisse
 - Speichern Metadaten wie Eigentümer und Berechtigungen
 - Außerdem Handles von Directory-Data-Objekten
- Symlinks repräsentieren symbolische Verweise

Funktionsweise... [2]

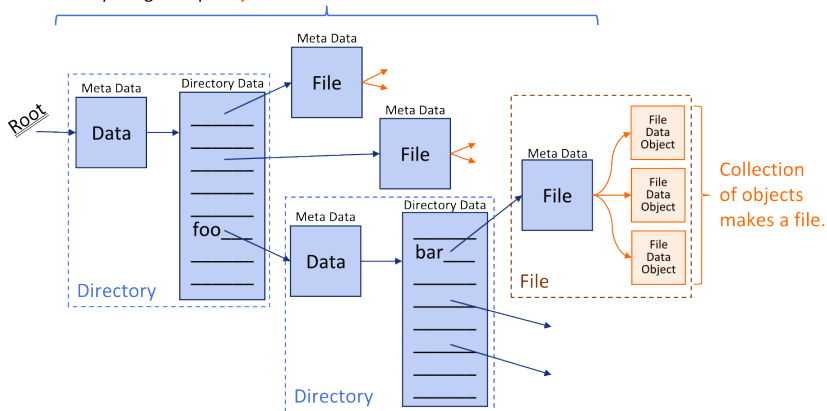


- Die logische Datei besteht aus Metadaten und Daten
- Die physikalische Datei besteht aus einem Metafile und drei Datafiles
 - Der erste Server enthält sowohl das Metafile als auch das erste Datafile
 - Die restlichen Server enthalten jeweils ein Datafile

Funktionsweise... [2]

Logical Flow of Distributed Metadata

Everything except **objects at the end of the chain** is **metadata**.



- 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...

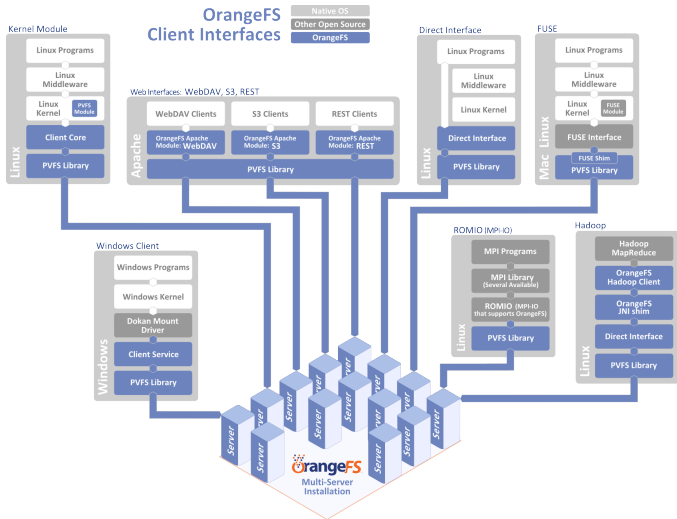
- Anwendungen können direkt gegen entsprechende Bibliotheken gelinkt werden
- Alternativ Preloading von Bibliotheken

```
1 $ gcc -o posix_app -L${PREFIX}/lib posix_app.c -lorangefsposix
2 $ ./posix_app
3 $ gcc -o orange_app -L${PREFIX}/lib orange_app.c -lorangefs
4 $ ./orange_app
```

```
1 $ export LD_LIBRARY_PATH=${PREFIX}/lib
2 $ export LD_PRELOAD=${PREFIX}/lib/libofs.so
3 $ ./legacy_app
```



Schnittstellen... [2]



Schnittstellen...

```
1 int fd;  
2 PVFS_hint myhint = NULL;  
3 int layout = PVFS_SYS_LAYOUT_RANDOM;  
4  
5 PVFS_hint_add(&myhint, PVFS_HINT_LAYOUT_NAME, sizeof(layout),  
6             ↪ &layout);  
7  
8 fd = pvfs_open("/path/to/file", O_CREAT | O_RDWR | O_TRUNC |  
9             ↪ O_HINTS, 0644, myhint);
```

Listing 1: OrangeFS-spezifische Hinweise [3]

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

Verteilungsfunktionen...

```

1 PVFS_sys_dist* new_dist;
2 PVFS_size strip_size = 131072;
3
4 new_dist = PVFS_sys_dist_lookup("simple_stripe");
5
6 PVFS_sys_dist_setparam(new_dist, "strip_size", &strip_size);
7 PVFS_sys_create(entry_name, parent_ref, attr, credentials,
   ↪ new_dist, &resp_create, NULL, hints);

```

Listing 2: Verteilungsfunktion setzen [3]

Verteilungsfunktionen...

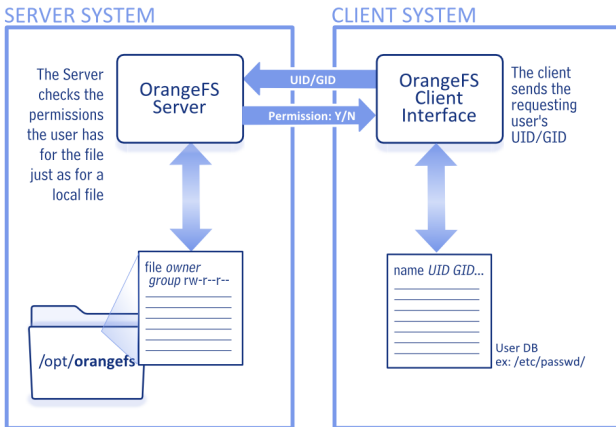
- Zusätzlich Unterstützung für vier Layouts
 - Legen die Reihenfolge der Server fest
- PVFS_SYS_LAYOUT_NONE
 - Reihenfolge wie in der Konfigurationsdatei angegeben
 - Start beim ersten Server
- PVFS_SYS_LAYOUT_ROUND_ROBIN (Standard)
 - Reihenfolge wie in der Konfigurationsdatei angegeben
 - Start bei einem zufälligen Server
- PVFS_SYS_LAYOUT_RANDOM
 - Reihenfolge der Server ist zufällig (ohne Wiederholung)
- PVFS_SYS_LAYOUT_LIST
 - Reihenfolge der Server ist wie explizit angegeben

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

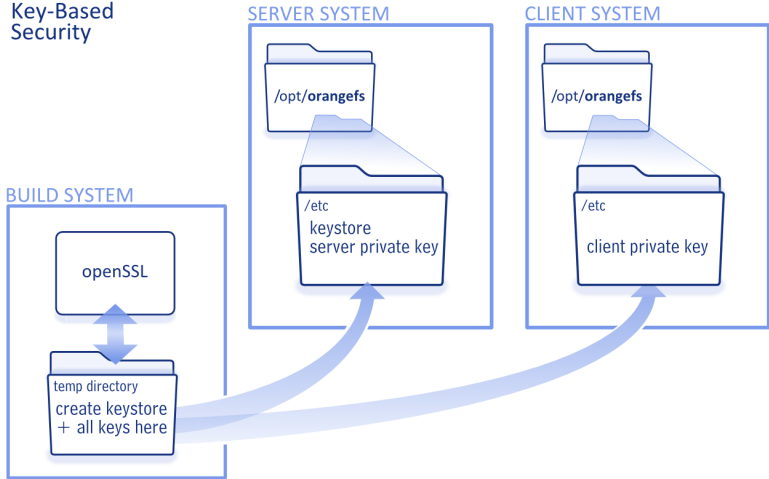


Sicherheit...

- 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... [2]

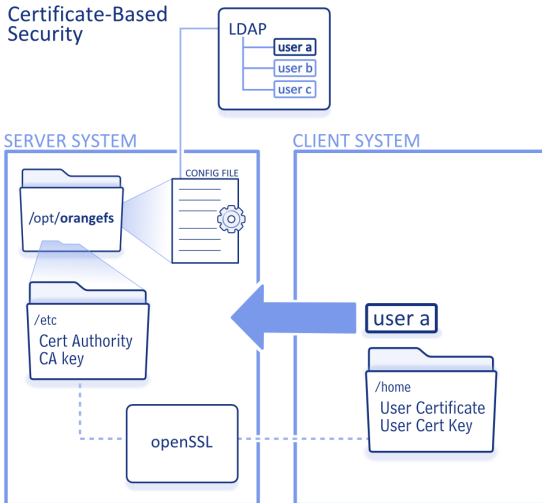
Key-Based Security



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... [2]



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

Operationen... [1]

```
1 PVFS_error PVFS_isys_create (...)
2 {
3     PINT_smc_b *smcb = NULL;
4     PINT_client_sm *sm_p = NULL;
5     ...
6     PINT_smc_b_alloc(&smcb, PVFS_SYS_CREATE, ...);
7     sm_p = PINT_sm_frame(smcb, PINT_FRAME_CURRENT);
8     ...
9     sm_p->u.create.object_name = object_name;
10    ...
11    sm_p->parent_ref = parent_ref;
12    sm_p->object_ref = parent_ref;
13    ...
14    return PINT_client_state_machine_post(smcb, ...);
15 }
```

- PINT_smcb: State machine control block
- PINT_client_sm: Client state machine
- PINT_smcb_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



State Machines... [1]

```
1 static PINT_sm_action create_init (struct PINT_smcb  
   ↪ *smcb, job_status_s *js_p)  
2 {  
3     struct PINT_client_sm *sm_p = PINT_sm_frame(smcb,  
   ↪ PINT_FRAME_CURRENT);  
4  
5     PINT_SM_GETATTR_STATE_FILL(  
6         sm_p->getattr,  
7         sm_p->object_ref,  
8         PVFS_ATTR_COMMON_ALL | PVFS_ATTR_DIR_HINT |  
9         PVFS_ATTR_CAPABILITY |  
   ↪ PVFS_ATTR_DISTDIR_ATTR,  
10        PVFS_TYPE_DIRECTORY,  
11        0);  
12  
13    return SM_ACTION_COMPLETE;  
14 }
```

- 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 }
```

State Machines... [1]

```
1  static PINT_sm_action create_cleanup (struct PINT_smbc
   ↪ *smcb, job_status_s *js_p)
2  {
3    ...
4    else if (...)
5    {
6        sm_p->u.create.stored_error_code = 0;
7        sm_p->u.create.retry_count++;
8        js_p->error_code = CREATE_RETRY;
9        return SM_ACTION_COMPLETE;
10   }
11   ...
12   PINT_SET_OP_COMPLETE;
13   return SM_ACTION_TERMINATE;
14 }
```

- 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



State Machines... [1]

```
1 machine pvfs2_client_create_sm
2 { ...
3     state create_setup_msgpair
4     {
5         run create_create_setup_msgpair;
6         success => create_xfer_msgpair;
7         default => cleanup;
8     }
9     state create_xfer_msgpair
10    {
11        jump pvfs2_msgpairarray_sm;
12        success => crdirent_setup_msgpair;
13        default => cleanup;
14    }
15    ...
16 }
```


Message Pairs... [1]

- Nachrichten werden für Versand kodiert
 - Bei Empfang wieder dekodiert
- Unterschiedliche Encodings für Anfragen und Antworten
 - Unterschiedliche Felder interessant
 - Beispielsweise für Create-Operation:
 - Anfrage: Dateisystem-ID, Credentials, Datafile-Parameter
 - Antwort: Metatile-Handle und -Attribute
- Unterstützung für mehrere Encoding-Typen
 - Standardmäßig “little endian bytefield encoding”
 - Geschachtelte Makros, die Daten umwandeln und kopieren



Message Pairs... [1]

```
1  decode_fields_6_struct(PVFS_servreq_create, ...);
2
3  #define decode_fields_6_struct(name, ...) \
4      decode_fields_6_generic(name, struct name,
5          ↪ ...)
6
7  #define decode_fields_6_generic(name, sname, ...) \
8      static inline void encode_##name (char **pptr,
9          ↪ const sname *x) { \
10         encode_##t1(pptr, &x->x1); \
11         encode_##t2(pptr, &x->x2); \
12         ...
13     } \
14     static inline void decode_##name (char **pptr,
15         ↪ sname *x) { \
16         ...
17     }
```


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

Zusammenfassung

- Unterstützung für mehrere Schnittstellen
 - MPI-IO, C, POSIX, System Call
- Einfache Installation und Konfiguration
 - Wenige Abhängigkeiten
 - Unproblematisch, da alles im User-Space
- State Machines zur Abarbeitung von Prozessen
 - Unterstützung für geschachtelte State Machines
- Unterstützung für mehrere Verteilungsfunktionen
 - Anpassung an Software-/Hardware-Umgebung
- Unterschiedliche Sicherheitsmechanismen
 - Einfach, schlüssel- und zertifikat-basiert

- 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.