

McKernel Specifications
Version 1.7.1-0.7

Masamichi Takagi, Balazs Gerofi, Tomoki Shirasawa, Gou Nakamura
and Yutaka Ishikawa

Monday 18th January, 2021

Contents

List of Figures

List of Tables

Chapter 1

インターフェイス

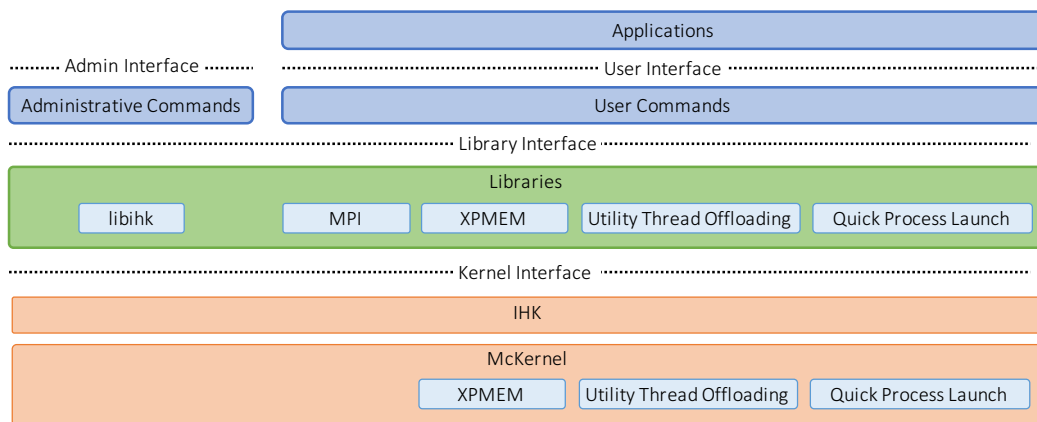


Figure 1.1: McKernel software stack

McKernel のソフトウェアスタックを図??に示す。本章では、ユーザ向けの User Interface と、アプリ向けの Library Interface と、コマンドやライブラリ向けの Kernel Interface を説明する。

本章の想定読者は、以下の 3 種類のユーザまたは開発者である。

- McKernel のコマンドを用いてアプリを実行するユーザ
- McKernel のライブラリインターフェイスを使用してアプリを開発する開発者
- McKernel のカーネルインターフェイスを使用してライブラリを開発する開発者

ユーザインターフェイス、ライブラリインターフェイスの関連ファイルは以下の通り。なお、インストールディレクトリを<install>とする。

インストール先	インターフェイス	説明
<install>/bin/mcexec	ユーザ	プロセス起動コマンド
<install>/bin/eclair	ダンプ解析ツール	
<install>/bin/vmcore2mckdump	ダンプ形式変換ツール	
<install>/rootfs/usr/lib64/libuti.so	ライブラリ	Utility Thread Offloading ライブラリ
<install>/include/uti.h	ライブラリ	Utility Thread Offloading ライブラリ ヘッダファイル
<install>/include/qlmpilib.h	ライブラリ	高速プロセス起動ヘッダファイル
<install>/lib/libqmpi.so	ライブラリ	高速プロセス起動ライブラリ
<install>/lib/libqfort.so	ライブラリ	高速プロセス起動ライブラリ (Fortran プログラム用)
<install>/lib/libxpmem.so	ライブラリ	XPMMEM ライブラリ
<install>/include/xpmem.h	ライブラリ	XPMMEM ライブラリヘッダファイル

12 以下、これら3種のインターフェイスを説明する。

13 1.1 プロセス起動コマンド

14 書式

```

15 mcexec [-c <cpu_id>] [-n <nr_partitions>] [-t <nr_threads>]
16 [-M (--mpol-threshold=<min>)] [-h (--extend-heap-by=<stride>)]
17 [-s (--stack-premap=<premap_size>)[,<max>]] [--mpol-no-heap] [--mpol-no-bss]
18 [--mpol-no-stack] [--mpol-shm-premap] [-m <numa_node>] [--disable-sched-yield]
19 [-O] [<os_index>]
20 <program> [<args>...]

```

21 オプション

22

-c <cpu_id>	mcexec を実行する CPU の番号を<cpu_id>に設定する。指定がない場合は 0 が用いられる。
-n <nr_partitions>	1 計算ノードの CPU 群を<nr_partitions>の区画に分割し、第 <i>i</i> 番目に起動された mcexec プロセスから起動される McKernel スレッドが第 <i>i</i> 番目の区画のみを利用するように設定する。分割は物理コア単位で行われる。こうすることで、1 ノード<nr_partitions>プロセスの MPI+OpenMP 実行において CPU を適切に使い分けることができる。
-t <nr_threads>	mcexec のスレッド数を<nr_threads>に設定する。このオプションが指定されない場合は、OMP_NUM_THREADS 環境変数が定義されている場合はその値+4 に設定し、存在しない場合は McKernel に割り当てられた CPU 数+4 に設定する。mcexec スレッドは McKernel からの要求を処理する。同時に多くの要求がなされる可能性があるため、この数は <McKernel のスレッド数 + α > に設定する必要がある。
-M (--mpol-threshold=<min>)	<min>以上のサイズのメモリを要求したときのみ、ユーザが設定したメモリ割り当てポリシーが適用されるようにする。<min>は K, M, G (k, m, g でもよい) の単位を付けた場合、それぞれ KiB, MiB, GiB の指定になる。指定がない場合はサイズに関係なくユーザが設定したメモリ割り当てポリシーが適用される。
-h (--extend-heap-by=<step>)	ヒープの拡大時にヒープサイズを少なくとも<step>バイト拡大する。また、ヒープの終了アドレスをラージページサイズにアラインする。<step>は K, M, G (k, m, g でもよい) の単位を付けた場合、それぞれ KiB, MiB, GiB の指定になる。指定がない場合は 4 KB が用いられる。
-s (--stack-premap=<premap_size>,<max>)	プロセス生成時にスタック領域のうち<premap_size>バイトをプリマップする。また、スタックの最大サイズを<max>に設定する。<premap_size>, <max>は K, M, G (k, m, g でもよい) の単位を付けた場合、それぞれ KiB, MiB, GiB の指定になる。指定がない場合、<premap_size>は 2 MB、<max>は ulimit コマンドまたは setrlimit() システムコールで設定された値が用いられる。
--mpol-no-heap	ヒープへのメモリ割り当て時にユーザの設定したメモリ割り当てポリシーに従わない。
--mpol-no-stack	スタックへのメモリ割り当て時にユーザの設定したメモリ割り当てポリシーに従わない。
--mpol-no-bss	bss へのメモリ割り当て時にユーザの設定したメモリ割り当てポリシーに従わない。
--mpol-shm-premap	/dev/shm を用いた共有メモリをプリマップする。
-m <numa_node>	メモリを<numa_node>番目の NUMA ノードから割り当てる。割り当てが不可能な場合は他の NUMA ノードから割り当てる。
--disable-sched-yield	sched_yield() 関数を何も行わない関数に置き換える。
-0	McKernel に割り当てられた CPU 数より大きい数のスレッドまたはプロセスの生成を許可する。指定がない場合は許可しない。許可されていない場合に、CPU 数より大きい数のスレッドまたはプロセスを clone(), fork(), vfork() などで生成しようとすると、当該システムコールが EINVAL エラーを返す。
<os_index>	プロセス起動先 OS インスタンスを<os_index>番に設定する。省略した場合は 0 番の OS インスタンスに起動する。

23 説明

24 <program>で指定された実行可能ファイルを args で指定された引数で、McKernel 上に起
25 動する。

26 mcexec の動作を変える環境変数は以下の通り。

書式	説明
MCEXEC_WL=<path1> [:<path2>...]	<path1>, <path2>, ... 以下に存在する McKernel 用実行ファイルについて、 <code>mcexec</code> の指定を省略する。なお、指定ディレクトリ以下に実行可能ファイルが存在しても、以下のケースでは Linux で実行される。 <ul style="list-style-type: none"> ● McKernel が動作していない場合 ● コマンドが 64 ビット ELF バイナリではない場合 ● コマンド名が <code>mcexec</code>, <code>ihkosctl</code>, <code>ihkconfig</code> である場合
MCEXEC_ALT_ROOT=<path>	<code>ld-linux.so</code> などのローダを探す際に、<path>と実行可能ファイルの <code>.interp</code> セクションに記載されたパスを結合したパスを探す。
MCKERNEL_RLIMIT_STACK= <premap_size>, <max>	(非推奨) プロセス生成時にスタック領域のうち<premap_size>バイトをプリマップする。また、スタックの最大サイズを<max>に設定する。<premap_size>, <max>は K, M, G (k, m, g でもよい) の単位を付けた場合、それぞれ KiB, MiB, GiB の指定になる。指定がない場合、<premap_size>は 2 MB、<max>は <code>ulimit -s</code> コマンドまたは <code>setrlimit()</code> システムコールで設定された値が用いられる。なお、本環境変数の代わりに <code>mcexec</code> の <code>--stack-premap</code> オプションを使用することを推奨する。

27 使用例は以下の通り。この例では `ls -ls` を McKernel 上で実行する。

28 \$ `mcexec ls -ls`

29 戻り値

30 <program>の exit status を返す。

31 1.2 ダンプ採取・解析

32 カーネルダンプの採取と解析のステップは以下の通り。

33 1. 以下のいずれかの方法でダンプファイルを作成する。

34 (a) IHK の関数 `ihk_os_makedumpfile()` または IHK のコマンド `ihkosctl` を用いて、
35 McKernel 形式のダンプファイルを作成する。

36 (b) Linux の `panic` を契機に `makedumpfile` 形式のダンプファイルを作成する。また、
37 コマンド `vmcore2mckdump` を用いて McKernel 形式に変換する。

38 2. `eclair` と呼ぶコマンドを用いてダンプファイルを解析する。

39 以下、関連コマンドのインターフェイスを説明する。

40 1.2.1 ダンプ解析コマンド

41 書式

42 `eclair [-ch] [-d <dump>] [-k <king>] [-o <os_index>] [-l] [-i]`

43 オプション

44

-c	NMI 受付時のコンテキストをスレッドとして扱う。それぞれのコンテキストは 1000000+ (CPU 番号) という TID を持つスレッドとして扱われる。スレッドとして扱うことで、割り込み処理のバックトレースを表示することができる。
-h	利用法を表示する。
-d <dump>	ダンプファイル名を指定する。指定がない場合は mcdump が用いられる。
-k <king>	カーネルイメージファイル名を指定する。指定がない場合は kernel.img が用いられる。
-o <os_index>	OS インスタンスのインデックスを指定する。指定がない場合は 0 が用いられる。
-l	run-queue にユーザスレッドが存在しない CPU について、idle() を実行しているスレッドが存在するように見せかける。
-i	Interactive mode と呼ぶ、デバッグ対象マシンに存在するメモリを直接参照した解析を行う。なお、ダンプ時に interactive mode を指定する必要がある。

45 説明

46 <dump>で指定された eclair 形式のダンプファイルを<os_index>で指定された OS インデックスを持つ OS として、<king>で指定されたカーネルイメージファイルを使って解析する。
47 ダンプ解析コマンド内では、gdb が動作しており、gdb と同じコマンドを利用できる。
48 McKernel は、マルチスレッドの単一プロセスに見える。まず、最初に、以下のコマンドを実行して、ダンプ解析コマンドにスレッド一覧を覚えさせる必要がある。

51 (eclair) info threads

52 quit コマンド実行時に、inferior の切り離し許可をユーザに求める。これには、y と応答すること。ダンプ解析コマンドは、gdb のコマンドの、bt コマンドと x コマンドをサポートする。

54 1.2.2 ダンプ形式変換コマンド

55 書式

56 vmcore2mckdump <vmcore> <file_name>

57 オプション

<vmcore>	makedumpfile 形式のダンプファイルのファイル名
<file_name>	変換先ダンプファイルのファイル名

58

59 説明

60 <vmcore>で指定された makedumpfile 形式のダンプファイルから McKernel に関連する部分を取り出し<file_name>で指定されたファイルに eclair 形式で出力する。

62 1.3 高速プロセス起動ライブラリインターフェイス

63 McKernel は、複数種の MPI プログラムを起動しさらにそれを繰り返すジョブにおいて MPI
64 プログラム起動時間を短縮する機能を提供する。利用例は以下の通り。

- 65 ● アンサンブルシミュレーションとデータ同化を繰り返す気象アプリケーション
66 このアプリではジョブスクリプトでそれぞれの MPI プログラムを交互に起動する。この起動時間を短縮する。
67

68 本機能を利用するためにはジョブスクリプトとアプリケーションを修正する必要がある。
69 ジョブスクリプトの修正方法を例を用いて説明する。

70 修正前

```
71 /* アンサンブルシミュレーションと同化を 10 回繰り返す */
72 for i in {1..10}; do
73
74     /* 100 ノードを用いるアンサンブルシミュレーションを 10 個並列に動作させる */
75     for j in {1..10}; do
76         mpiexec -n 100 -machinefile ./list1_$j p1.out a1 & pids[$i]=$!;
77     done
78
79     /* p1.out の終了を待つ */
80     for j in {1..10}; do wait ${pids[$j]}; done
81
82     /* アンサンブルシミュレーションで用いたのと同じ 1000 ノードを用いてデータ同化
83        を行う */
84     mpiexec -n 1000 -machinefile ./list2 p2.out a2
85 done
```

86 修正後

```
87 for i in {1..10}; do
88     for j in {1..10}; do
89         /* mpiexec を ql_mpiexec_start に置き換える */
90         ql_mpiexec_start -n 100 -machinefile ./list1_$j p1.out a1 & pids[$j]=$!;
91     done
92
93     for j in {1..10}; do wait ${pids[$j]}; done
94
95     ql_mpiexec_start -n 1000 -machinefile ./list2 p2.out a2
96 done
97
98 /* p1.out と p2.out は常駐しているため、ql_mpiexec_finalize で終了させる。
99     mpiexec への引数と実行可能ファイル名で MPI プログラムを識別しているため、
100     実行時と同じものを指定する。 */
101 for j in {1..10}; do
102     ql_mpiexec_finalize -machinefile ./list1_$i p1.out a1;
103 done
104 ql_mpiexec_finalize -machinefile ./list2 p2.out a2;
```

105 アプリケーションの修正方法を擬似コードを用いて説明する。計算を何度も行えるよう
106 なループ構造を持たせ、また ql_client() を計算完了後に呼び出すようにする。

```
107     MPI_Init();
108     先行・後続 MPI プログラムとの通信準備
109 loop:
110     foreach (Fortran の) モジュール
111         コマンドライン引数・パラメタファイル・環境変数を用いた初期化処理
112     先行 MPI プログラムからのデータ受信・スナップショット読み込み
113     計算
114     後続 MPI プログラムへのデータ送信・スナップショット書き出し
115     /* ループボディの終わりに ql_client() を挿入する */
```

```

116     if(ql_client() == QL_CONTINUE) { goto loop; }
117     MPI_Finalize();

```

118 以下、コマンドや関数のインターフェイスを説明する。

119 1.3.1 MPI プロセス開始再開コマンド

120 書式

```

121     ql_mpiexec_start -machinefile <hostfile> [<mpiopts>...] <exe> [<args>...]

```

122 オプション

オプション	内容
-machinefile <hostfile>	ホストファイル
<mpiopts>	mpiexec のオプション
<exe>	実行可能ファイル
<args>	実行可能ファイルの引数

123

124 説明

125 <exe>で指定された MPI プログラムを開始する。または再開指示待ちの状態にある MPI
126 プログラムに次の計算開始を指示する。本コマンドは MPI プログラムの一回の計算の完了
127 と共に終了する。また、MPI プログラムは<hostfile>の内容、<mpiopts>、<exe>とで識別
128 する。

129 ql_mpiexec_{start,finalize} コマンドから MPI プログラムに次の動作、引数、環境変
130 数を渡すために用いるファイルは、環境変数 QL_PARAM_PATH が定義されている場合はその下に、
131 そうでない場合はホームディレクトリ下に作成される。当該ディレクトリは ql_mpiexec_start
132 コマンドを実行するノード、各 MPI プロセスが実行される計算ノードからアクセスできる必
133 要がある。

134 また、環境変数 MPIEXEC_TIMEOUT によるタイムアウトおよび複数の実行可能ファイルの
135 指定はサポートしない。

136 戻り値

戻り値	説明
0	正常終了
0 以外	異常終了

137

138 エラー時出力

139 エラーメッセージは mpiexec が出力するエラーメッセージの他に ql_mpiexec_start 独自
140 に以下のメッセージを出力する。

メッセージ	意味
unknown option: <opt>	未知のオプション <opt>が指定された
bad option: <opt>	オプション <opt>の指定が誤っている
unsupported option: <opt>	オプション <opt>はサポートしていない
'.' is unsupported	'.' はサポートしていない
unable to read hostfile(<hostfile>): <reason>	<hostfile>を<reason>の理由により読み込めない
could not open hostfile(<hostfile>): <reason>	<hostfile><reason>の理由によりオープンできない
<hostfile> not exist	<hostfile>が存在しない
specify -machinefile option	-machinefile オプションが指定されていない
no user program	<exe>が指定されていない
socket directory not exist	ソケット通信用のディレクトリが存在しない
ql_server not execution <reason>	ql_server の起動に<reason>の理由により失敗した
ql_mpiexec_start: socket(<reason>)	ql_mpiexec_start コマンドの socket 操作で<reason>のエラーが発生した
ql_mpiexec_start: bind(<reason>)	ql_mpiexec_start コマンドの bind で<reason>のエラーが発生した
ql_mpiexec_start: listen(<reason>)	ql_mpiexec_start コマンドの listen で<reason>のエラーが発生した
ql_mpiexec_start: connect(<reason>)	ql_mpiexec_start コマンドの connect で<reason>のエラーが発生した

141 1.3.2 MPI プロセス終了指示コマンド

142 書式

143 ql_mpiexec_finalize -machinefile <hostfile> [<mpiopts>...] <exe>

144 オプション

オプション	説明
-machinefile <hostfile>	ホストファイル
<mpiopts>	mpiexec のオプション
<exe>	実行可能ファイル

145

146 説明

147 ql_mpiexec_start によって起動された MPI プログラムを終了させる。本コマンドは MPI
 148 プログラムの終了と共に終了する。また、MPI プログラムは<hostfile>の内容、<mpiopts>、
 149 <exe>とで識別する。

150 戻り値

戻り値	説明
0	正常終了
0 以外	異常終了

151

152 エラー時出力

153 エラーメッセージは `mpiexec` が出力するエラーメッセージの他に `ql_mpiexec_finalize` 独自に以下のメッセージを出力する。

メッセージ	意味
unknown option: <opt>	未知のオプション<opt>が指定された
bad option: <opt>	オプション<opt>の指定が誤っている
unsupported option: <opt>	オプション<opt>はサポートしていない
'.' is unsupported	'.' はサポートしていない
unable to read hostfile(<hostfile>): <reason>	<hostfile>を<reason>の理由で読み込めない
could not open hostfile(<hostfile>): <reason>	<hostfile>を<reason>の理由でオープンできない
<hostfile> not exist	ホストファイルが存在しない
specify -machinefile option	-machinefile オプションが指定されていない
no user program	<exe>が指定されていない
socket directory not exist	ソケット通信用のディレクトリが存在しない
not found mpi process	mpiexec プロセスが存在しない

154

155 1.3.3 計算の再開・終了関数 (C 言語)

156 書式

157 `ql_client(int *argc, char ***argv)`

158 引数

引数	説明
<code>argc</code>	引数の数へのポインタ
<code>argv</code>	引数文字列の配列へのポインタ

159

160 説明

161 `ql_mpiexec_{start,finalize}` コマンドによる指示を待ち、指示結果を返す。本関数は、
162 MPI プログラム内で、一回の計算の完了後に呼び出す。

163 戻り値

戻り値	説明
<code>QL_CONTINUE</code>	次の計算の開始が指示された
<code>QL_EXIT</code>	MPI プログラムの終了が指示された、あるいは当該プロセスが <code>ql_mpiexec_start</code> コマンドで起動されていない

164

165 1.3.4 計算の再開・終了関数 (Fortran)

166 書式

167 `subroutine QL_CLIENT(ierr)`

168 引数

引数	型	説明
ierr	INT	戻り値

169

170 説明

171 MPI プログラム内で一回の計算の完了後に呼び出され、`ql_mpiexec_{start,finalize}`
 172 コマンドによる指示を待ち、指示結果を返す。なお、本関数を使用するためには `libqlfort.so`
 173 を `LD_PRELOAD` でロードする必要がある。また、コンパイラは GNU Fortran Compiler また
 174 は Intel Fortran Compiler をサポートする。Intel Fortran Compiler 使用時は、コンパイルオ
 175 プションに `-shared-intel` を指定する必要がある。

176 戻り値

戻り値	説明
1	次の計算の開始が指示された
0	MPI プログラムの終了が指示された、あるいは当該プロセスが <code>ql_mpiexec_start</code> コマンドで起動されていない

177

178 1.4 高速プロセス起動カーネルインターフェイス

179 1.4.1 swapout システムコール

180 書式

181 `int swapout(char *filename, void *workarea, size_t size, int flag)`

182 引数

引数	説明
filename	スワップファイル名へのポインタ
workarea	作業領域へのポインタ
size	作業領域のサイズ
flag	swapout の動作制御用フラグ

183

184 説明

185 プロセスのメモリ領域のファイルへの待避（スワップアウトと呼ぶ）とファイルからの復
 186 元（スワップインと呼ぶ）を行う。

187 処理ステップは以下の通り。

188 1 filename が NULL または flag が 1 の場合はステップ 6 に進む。そうでない場合はス
 189 テップ 2 に進む。

- 190 2 スワップアウト処理を行う。
- 191 3 flag が2の場合は、ステップ5に進む。そうでない場合は、ステップ4に進む。
- 192 4 mcexec へ制御を移し、スワップアウト完了の同期と、`ql_mpiexec_{start,finalize}`
193 による指示を待った後、本システムコールに制御を戻す。
- 194 5 スワップイン処理を行う。さらに呼び出し元に戻る。
- 195 6 mcexec へ制御を移し、`ql_mpiexec_{start,finalize}` による指示を待った後、本シス
196 テムコールに制御を戻す。さらに呼び出し元に戻る。

197 戻り値

戻り値	説明
0	正常終了
-1	エラー
-ENOMEM	メモリが不足
-EINVAL	引数が不正

198

199 1.5 Utility Thread Offloading ライブラリインターフェイス

200 インターフェイスは「McKernel 仕様付録 (Utility thread offloading ライブラリ編)」に記載
201 する。

202 1.6 Utility Thread Offloading カーネルインターフェイス

203 McKernel は、スレッドを Linux の CPU にマイグレートする機能 (Utility Thread Offloading,
204 UTI と呼ぶ) を提供する。UTI のカーネルインターフェイスは、第??節で説明するライブラ
205 リによって用いられる。

206 以下、関連システムコールのインターフェイスを説明する。

207 1.6.1 McKernel スレッドの Linux へのマイグレートシステムコール

208 書式

```
209 int util_migrate_inter_kernel(uti_attr_t *attr)
```

210 説明

211 attr と環境変数 UTI_CPU_SET で指定された、CPU 位置とスレッドの振る舞いの記述に基
212 づき、呼び出し元スレッドを Linux CPU にマイグレートさせる。

213 環境変数 UTI_CPU_SET はビットマップ形式で CPU 位置を示す。また、uti_attr_t は以
214 下のように定義される。

```
215 #define UTI_MAX_NUMA_DOMAINS (1024)
216
217 typedef struct uti_attr {
218     uint64_t numa_set[(UTI_MAX_NUMA_DOMAINS + 63) / 64];
```

```

219     /* スレッド配置先 NUMA ノードを表すビットマップ */
220     uint64_t flags;
221     /* CPU 位置とスレッドの振る舞いを表すビットマップ */
222 } uti_attr_t;

```

223 `uti_attr_t` の `flags` はビットマップで、ビット 1 は対応する CPU 位置の指示または振
224 る舞いの記述が有効であることを示す。ビット位置と指示・記述の対応は以下の通り。

```

225 #define UTI_FLAG_NUMA_SET (1ULL<<1)
226 /* numa_set フィールドで指定した NUMA ノードへ配置する */
227 #define UTI_FLAG_SAME_NUMA_DOMAIN (1ULL<<2)
228 /* 呼び出し元と同一 NUMA ノードへ配置する */
229 #define UTI_FLAG_DIFFERENT_NUMA_DOMAIN (1ULL<<3)
230 /* 呼び出し元とは異なる NUMA ノードへ配置する */
231 #define UTI_FLAG_SAME_L1 (1ULL<<4)
232 #define UTI_FLAG_SAME_L2 (1ULL<<5)
233 #define UTI_FLAG_SAME_L3 (1ULL<<6)
234 /* 呼び出し元とそれぞれのレベルのキャッシュを共有する CPU へ配置する */
235 #define UTI_FLAG_DIFFERENT_L1 (1ULL<<7)
236 #define UTI_FLAG_DIFFERENT_L2 (1ULL<<8)
237 #define UTI_FLAG_DIFFERENT_L3 (1ULL<<9)
238 /* 呼び出し元とそれぞれのレベルのキャッシュを共有しない CPU へ配置する */
239 #define UTI_FLAG_EXCLUSIVE_CPU (1ULL<<10)
240 /* CPU を専有させると効率的に動作する。
241     例えば、mwait 命令を用いている。*/
242 #define UTI_FLAG_CPU_INTENSIVE (1ULL<<11)
243 /* CPU サイクルを多く使用する。例えば、ネットワーク
244     デバイスのイベントキューを繰り返しポーリングする。*/
245 #define UTI_FLAG_HIGH_PRIORITY (1ULL<<12)
246 /* スケジューラのプライオリティを上げると効率的に動作する。
247     例えば、ネットワークデバイスのイベント待ちをする。*/
248 #define UTI_FLAG_NON_COOPERATIVE (1ULL<<13)
249 /* co-operative スケジューリングを行っていない。例えば、
250     イベント待ちになった際に sched_yield() を呼ぶ、ということをしていない。*/

```

251 なお、McKernal から Linux への 1 度のマイグレートのみ可能である。

252 戻り値

0	正常終了
-1	エラー

253

254 エラー時の `errno` の値

-ENOSYS	<code>util_migrate_inter_kernel</code> がサポートされていない。
-EFAULT	<code>attr</code> にアクセスできない。

255

256 1.6.2 スレッド生成先 OS 指定システムコール

257 書式

```
258 int util_indicate_clone(int mod, uti_attr_t *attr)
```

259 説明

260 呼び出し元スレッドが発行する clone システムコールの動作を変え、スレッド生成後
261 直ちに mod に指定した OS へマイグレートさせる。CPU 位置と Linux のスケジューラ設
262 定は、attr と環境変数 UTI_CPU_SET で指定されたヒントに基づいて決定される。この関
263 数は、pthread_create() などで Linux へスレッドを生成させるために用いる。本関数も、
264 util_migrate_inter_kernel と同様、McKernel から Linux への 1 度のマイグレートのみ可
265 能である。

mod の取りうる値と意味は以下の通り。

SPAWN_TO_REMOTE	Linux へ生成
SPAWN_TO_LOCAL	McKernel へ生成

266

267 戻り値

0	正常終了
-1	エラー

268

269 エラー時の errno の値

ENOSYS	util_indicate_clone がサポートされていない。
EINVAL	mod に未定義の値を指定した。
EFAULT	attr にアクセスできない。

270

271 1.6.3 カーネル種別取得システムコール

272 書式

```
273 int get_system()
```

274 説明

275 呼び出し元スレッドが動作している OS の種別を返却する。なお、本関数の名称は次バー
276 ジョンにて is_mckernel() 等に変更される予定である。

277 戻り値

0	OS が McKernel
-1	エラー (OS が Linux)

278 エラー時の `errno` の値

ENOSYS	Linux で呼び出した
--------	--------------

279

280 1.7 XPMEM ライブラリインターフェイス

281 XPMEM を使うことによって、あるプロセスがマップしたメモリ領域を他のプロセスからマッ
282 プできるようになる。利用方法は以下の通り。第 1 のプロセスのメモリ領域を第 2 のプロセ
283 スがマップしようとしているとする。

- 284 1. 第 1 のプロセスがメモリ領域を `xpmem_make()` を用いて XPMEM segment として登録
285 する。また、segment id を第 2 のプロセスに渡す。
- 286 2. 第 2 のプロセスが `xpmem_get()` で当該 XPMEM segment に対するアクセス許可を得る。
- 287 3. 第 2 のプロセスが `xpmem_attach()` で当該 XPMEM segment を自身の仮想アドレス範
288 囲にマップする。
- 289 4. 第 2 のプロセスが当該メモリ領域に対する操作を行う。
- 290 5. 第 2 のプロセスが `xpmem_detach()` で当該メモリ領域をアンマップする。
- 291 6. 第 2 のプロセスが `xpmem_release()` で当該 XPMEM segment に対するアクセス許可が
292 不要になったことをドライバに伝える。
- 293 7. 第 1 のプロセスが `xpmem_remove()` を用いて当該 XPMEM segment を破棄する。

294 以下、関連関数のインターフェイスを説明する。

295 1.7.1 Get Version Number

296 書式

```
297 int xpmem_version (void)
```

298 説明

299 This function gets the XPMEM version.

300 戻り値

$\neq -1$	XPMEM version number
-1	Failure

301

302 1.7.2 Expose Memory Block

303 書式

304

```
305 xpmem_segid_t xpmem_make(  
306     void *vaddr,  
307     size_t size,  
308     int permit_type,  
309     void *permit_value)
```

310 説明

311 `xpmem_make()` shares a memory block specified by `vaddr` and `size` by invoking the
312 XPMEM driver. `permit_type` is for the future extension. Use `XPMEM_PERMIT_MODE` for this
313 version. `permit_value` specifies the permissions mode expressed as an octal value.

314 This function is expected to be called by the source process to obtain a segment ID
315 to share with other processes. It is common to call this function with `vaddr = NULL` and
316 `size = XPMEM_MAXADDR_SIZE`. This will share the entire address space of the calling process.

317 戻り値

$\neq -1$	64-bit segment ID (<code>xpmem_segid_t</code>)
-1	Failure

318

319 1.7.3 Un-Expose Memory Block

320 書式

321

```
322 static int xpmem_remove(xpmem_segid_t segid)
```

323 説明

324 The opposite of `xpmem_make()`, this function deletes the mapping specified by `segid`
325 that was created from a previous `xpmem_make()` call. All the attachments created by
326 `xpmem_attach()` are detached and all the permits obtained by `xpmem_get()` are revoked.

327 Optionally, this function is called by the source process, otherwise automatically called
328 by the driver when the source process exits.

329 戻り値

0	Success
-1	Failure

330

331 1.7.4 Get Access Permit

332 書式

333

```
334 xpmem_apid_t xpmem_get(  
335     xpmem_segid_t segid,  
336     int flags,  
337     int permit_type,  
338     void *permit_value)
```

339 説明

340 `xpmem_get()` attempts to get access to a shared memory block specified by `segid`.
341 `flags` specifies access mode, i.e. read-write (`XPMEM_RDWR`) or read-only (`XPMEM_RDONLY`).
342 `permit_type` is for the future extension. Use `XPMEM_PERMIT_MODE` for this version. `permit_value`
343 specifies the permissions mode expressed as an octal value.

344 This function is called by the consumer process to get permission to attach memory
345 from the source virtual address space associated with `segid`. If access is granted, an `apid`
346 will be returned to pass to `xpmem_attach()`.

347 戻り値

≠ -1	64-bit access permit ID (<code>xpmem_apid_t</code>)
-1	Failure

348

349 1.7.5 Release Access Permit

350 書式

351

```
352 int xpmem_release(xpmem_apid_t apid)
```

353 説明

354 The opposite of `xpmem_get()`, this function deletes any mappings associated with `apid`
355 in the consumer's address space. Optionally, this function is called by the consumer process,
356 otherwise automatically called by the driver when the consumer process exits.

357 戻り値

0	Success
-1	Failure

358

359 1.7.6 Attach to Memory Block

360 書式

361

```
362 static int xpmem_attach(  
363     struct xpmem_addr addr,  
364     size_t size,  
365     void *vaddr)
```

366 説明

367 This function attaches a virtual address space range from the source process.
368 struct xpmem_addr is defined as follows.

```
369 struct xpmem_addr {  
370     /** apid that represents memory */  
371     xpmem_apid_t apid;  
372     /** offset into apid's memory region */  
373     off_t offset;  
374 };
```

375 addr.apid is the access permit ID returned from a previous xpmem_get() call. addr.offset
376 is offset into the source memory to begin the mapping. The mapping is created at vaddr
377 with the size of size. Kernel chooses the mapping address if vaddr is NULL.

378 This function is called by the consumer to get a mapping between the shared source
379 address and an address in the consumer process' own address space. If the mapping is
380 successful, then the consumer process can now begin accessing the shared memory.

381 戻り値

≠ -1	Virtual address at which the mapping was created
-1	Failure

382

383 1.7.7 Detach from Memory Block

384 書式

385

```
386 int xpmem_detach(void *vaddr)
```

387 説明

388 This function detach from the virtual address space of the source process.

389 Optionally, this function is called by the consumer process, otherwise automatically
390 called by the driver when the consumer process exits.

391 戻り値

392

0	Success
-1	Failure

393 1.8 XPMEM カーネルインターフェイス

394 XPMEM は、あるプロセスがマップしたメモリ領域を他のプロセスからマップできるように
395 する。XPMEM のカーネルインターフェイスは、第??節で説明するライブラリによって用い
396 られる。

397 以下、関連する `ioctl()` のインターフェイスを説明する。

398 1.8.1 `ioctl` システムコール

399 書式

400 `int ioctl(int fd, int cmd, void* arg)`

401 説明

402 `cmd` で指定された操作を行う。`cmd` ごとの処理を表?? に示す。

Table 1.1: XPMEM デバイスに対する `ioctl` の各コマンドの処理

コマンド	説明
<code>XPMEM_CMD_VERSION</code>	バージョン番号を返す。
<code>XPMEM_CMD_MAKE</code>	<code>arg->vaddr</code> から始まる長さ <code>arg->size</code> のメモリ領域を共有可能にし、segment id を <code>arg->segid</code> に格納する。メモリ領域のパーミッションは <code>arg->permit_value</code> に設定される。
<code>XPMEM_CMD_REMOVE</code>	<code>arg->segid</code> で指定されたメモリ領域の共有を解除する。
<code>XPMEM_CMD_GET</code>	<code>arg->segid</code> で指定されたメモリ領域に対する <code>arg->permit_value</code> で指定されたパーミッションでのアクセス許可取得を試みる。成功した場合、アクセス許可の id が <code>arg->apid</code> に格納される。
<code>XPMEM_CMD_RELEASE</code>	<code>arg->apid</code> で指定されたメモリ領域に対するアクセス許可を返却する。
<code>XPMEM_CMD_ATTACH</code>	<code>arg->apid</code> で指定された共有メモリ領域のうち <code>arg->offset</code> から始まる長さ <code>arg->size</code> の範囲を <code>arg->vaddr</code> から始まるアドレス範囲にマップする。
<code>XPMEM_CMD_DETACH</code>	<code>arg->vaddr</code> から始まる共有マップを解放する。

403 `XPMEM_CMD_MAKE` コマンドで用いる `xpmem_cmd_make` 構造体は以下のように定義される。

```
404 struct xpmem_cmd_make {
405     __u64 vaddr;
406     size_t size;
407     int permit_type;
408     __u64 permit_value;
409     xpmem_segid_t segid;    /* returned on success */
410 };
```

411 `xpmem_segid_t` は以下のように定義される。

```
412 typedef __s64 xpmem_segid_t;    /* segid returned from xpmem_make() */
```

413 `XPMEM_CMD_REMOVE` コマンドで用いる `xpmem_cmd_remove` 構造体は以下のように定義さ
414 れる。

```
415 struct xpmem_cmd_remove {
416     xpmem_segid_t segid;
417 };
```

418 XPMEM_CMD_GET コマンドで用いる `xpmem_cmd_get` 構造体は以下のように定義される。

```
419 struct xpmem_cmd_get {
420     xpmem_segid_t segid;
421     int flags;
422     int permit_type;
423     __u64 permit_value;
424     xpmem_apid_t apid; /* returned on success */
425 };
```

426 `xpmem_apid_t` は以下のように定義される。

```
427 typedef __s64 xpmem_apid_t; /* apid returned from xpmem_get() */
```

428 XPMEM_CMD_RELEASE コマンドで用いる `xpmem_cmd_release` 構造体は以下のように定義される。

```
430 struct xpmem_cmd_release {
431     xpmem_apid_t apid;
432 };
```

433 XPMEM_CMD_ATTACH コマンドで用いる `xpmem_cmd_attach` 構造体は以下のように定義される。

```
435 struct xpmem_cmd_attach {
436     xpmem_apid_t apid;
437     off_t offset;
438     size_t size;
439     __u64 vaddr;
440     int fd;
441     int flags;
442 };
```

443 XPMEM_CMD_DETACH コマンドで用いる `xpmem_cmd_detach` 構造体を以下のように定義される。

```
445 struct xpmem_cmd_detach {
446     __u64 vaddr;
447 };
```

448 XPMEM_CMD_ATTACH コマンドで用いる `xpmem_addr` 構造体は以下のように定義される。

```
449 struct xpmem_addr {
450     xpmem_apid_t apid; /* apid that represents memory */
451     off_t offset; /* offset into apid's memory */
452 };
```

453 戻り値

454

0	正常終了
-EFAULT	アドレスが不正である
-EINVAL	引数が無効である

455 Chapter 2

456 実装者向けインターフェイス詳細

457 本章の想定読者は以下の通り。

- 458 • McKernel の、アーキテクチャ移植を含む開発を行う開発者

459 2.1 概要

460 McKernel is a lightweight kernel for HPC with the following features.

- 461 • Quickly adapts to the new hardware techniques to provide scalability and full-control
462 of hardware
- 463 • Supports new programming style such as in-situ data analytics and scientific work-flow
- 464 • Provides a complete set of Linux API

465 McKernel is based on a light-weight kernel developed at the University of Tokyo[?]. It
466 works with systems with Intel Xeon processors and systems with Intel Xeon phi processor.
467 Figure ?? shows the architecture of McKernel. Cores and memory of a compute-node are
468 divided into two partitions and Linux runs on one of them and McKernel runs on the other.
469

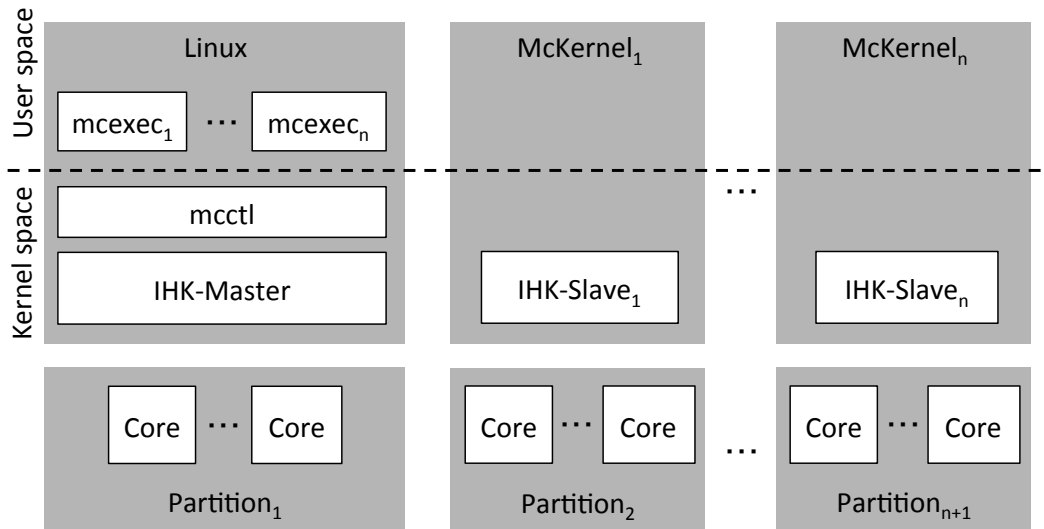


Figure 2.1: The architecture of McKernel

470 Two kernel modules, `mcctl` and `IHK-Master`, and user processes `mcexec` (`mcexec1`,
 471 `mcexec2`, ...) exist in the Linux kernel while McKernel (`McKernel1`, `McKernel2`, ...) and
 472 `IHK-Slave` (`IHK-Slave1`, `IHK-Slave2`, ...) reside in each partition.

473 Linux controls all hardware resources when booting a compute-node. The Interface for
 474 Heterogeneous Kernel, formed by both `IHK-Master` and `IHK-Slave`, implements a com-
 475 munication mechanism between Linux and McKernel, called Inter Kernel Communication
 476 (IKC). In addition of that, the `IHK-Master` has an important role, allocating cores and
 477 memory for McKernel, and booting it. `IHK` is independently designed from McKernel, and
 478 it may be used for other kernels with Linux.

479 The `mcctl` kernel module controls the McKernel. In order to provide Linux API for ap-
 480 plications running on McKernel, OS service requests not provided by McKernel is delegated
 481 to Linux and performed by Linux. The `mcexec` command requests McKernel to launch an
 482 application via `IHK`. After the application's invocation, a `mcexec` process acts as a proxy
 483 or ghost process for the McKernel process in the sense that Linux system calls delegated
 484 from McKernel via `IHK` are issued by this process.

485 In the rest of this section, McKernel features, i.e., McKernel usages, process and mem-
 486 ory management, system calls, and the `procfs/sysfs` file system will be described.

487 McKernel を用いたジョブの実行ステップを図??を用いて説明する。

- 488 1. 運用ソフトが計算ノード上に Linux を起動する (図の (1))
- 489 2. ユーザがジョブキューを指定することで、McKernel と Linux のどちらを使用するか、
 490 また McKernel を使用する場合は様々なチューニングが施されたカーネルイメージのう
 491 ちどれを使用するかを指定する。例えば、ラージページ化が効果のあるアプリ B を実行
 492 しようとしている場合は、その機能を持つイメージを指定するジョブキューにジョブを
 493 投入する。
- 494 3. 運用ソフトウェアがジョブ投入を受けて、資源のパーティショニング、McKernel の起動、
 495 アプリの実行を行う (図の (2))。例では、ラージページ化促進機能を持つ McKernel
 496 が起動され、アプリ B がその上で実行される。

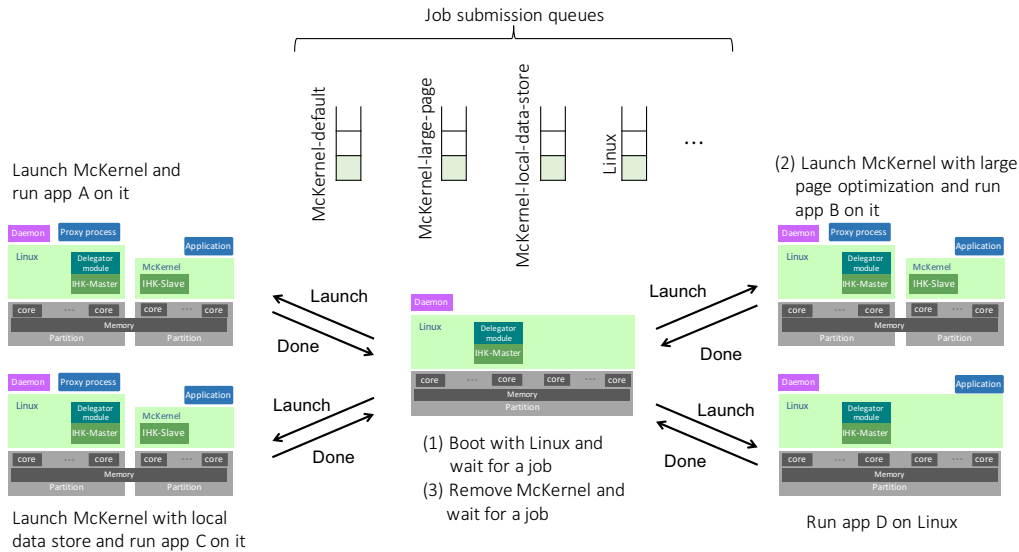


Figure 2.2: McKernel Usages

- 497 4. 運用ソフトウェアが、ジョブ終了時に計算ノード状態を元の状態、すなわち Linux のみ
 498 が動作する状態に戻す (図の (3))

499 2.2 プロセス管理

500 McKernel has a unique process execution model to realize cooperation with Linux. McKernel
 501 processes are primarily spawn by the Linux command line tool `mcexec`¹. For every single
 502 McKernel process there is a corresponding `mcexec` Linux process that exists throughout the
 503 lifetime of the application. `mcexec` serves the following purposes:

- 504 - It provides an execution context for offloaded system calls (explained in Section ??)
 505 so that they can be invoked directly in Linux
- 506 - It enables transparent access to Linux device drivers through the mechanism of unified
 507 address-space (discussed in Section ??) and the ability to map Linux device files
 508 directly to McKernel processes
- 509 - It facilitates Linux to maintain certain application associated kernel state that would
 510 have to be otherwise maintained by McKernel (e.g., open files and the file descriptor
 511 table (see Section ??), process specific device driver state, etc.)

512 Due to its role to providing a gateway to specific Linux features, we call `mcexec` the
 513 *proxy-process*. Figure ?? provides an overview of IHK/McKernel's proxy-process architec-
 514 ture as well as the system call offloading mechanism.

¹An alternative way of creating McKernel processes via the `fork()` system call will be discussed in Section ??.

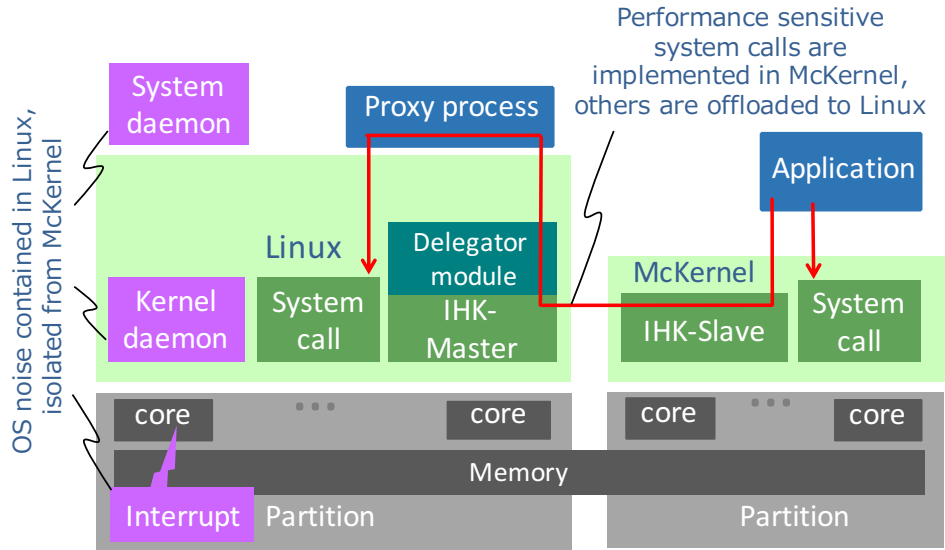


Figure 2.3: Overview of the IHK/McKernel architecture and the system call delegation mechanism.

515 We emphasize that IHK/McKernel runs HPC applications primarily on the LWK but
 516 the full Linux API is available via system call delegation. System call offloading will be
 517 detailed in Section ??.

518 Since the user shell process runs on the Linux side, a signal to an McKernel process
 519 cannot be delivered directly from Linux. Instead, the shell process issues signals to `mcexec`
 520 and `mcexec` forwards the signal to the McKernel process via IKC. For more information on
 521 signaling, see Section ??.

522 2.2.1 Linux からのプロセス起動

523 `mcexec` が Linux からプロセスを起動するステップは以下の通り。

- 524 1. It opens the device `/dev/mcosn` to communicate with McKernel.
- 525 2. It sends the ELF binary description header, the command line and environment
526 variables to the McKernel.
- 527 3. It uploads the application binary to McKernel's memory area.
- 528 4. It creates a Linux thread pool that will serve system call offloading requests. Addi-
529 tionally, one of the workers is designated for waiting for signals from McKernel.
- 530 5. It sends a request for starting the process to McKernel.
- 531 6. The main thread waits for termination of all workers.
- 532 7. When a worker receives the `exit_group()` system call, it terminates all workers in
533 the thread pool.

534 なお、環境変数 `MCEXEC_WL` に McKernel 用実行可能ファイルの（親）ディレクトリを指
 535 定することで、`mcexec` の指定を省略できる。複数ディレクトリを指定する場合は、コロンを

536 デリミタとして指定する。なお、指定ディレクトリ以下に実行可能ファイルが存在しても、以
537 下のケースでは Linux で実行される。

- 538 • McKernel が動作していない場合
- 539 • コマンドが 64 ビット ELF バイナリではない場合
- 540 • コマンド名が `mcexec`, `ihkosctl`, `ihkconfig` である場合

541 この機能は、`mcctrl` が Linux のローダのリストに特別なローダを挿入することで実現される。

542 2.2.2 fork()

543 The `fork()` system call is supported in McKernel and it is an alternative way for spawning
544 new processes. `fork()` is handled as follows:

- 545 1. McKernel allocates a CPU core and memory for the child process.
- 546 2. McKernel creates information on process and virtual memory, and the user execution
547 context.
- 548 3. McKernel copies the parent memory to the child process. Note that the anonymous
549 memory areas such as text, data, bss, are copied without using copy-on-write technique
550 in the current implementation.
- 551 4. McKernel requests `mcexec` to perform a fork system call (i.e., to create a new proxy
552 process for the child) in Linux. `mcexec` executes the following steps:
 - 553 (a) `mcexec` issues the fork system call to create a new Linux process (call it the child
554 proxy).
 - 555 (b) The child proxy closes the device `/dev/mcosn` and reopens it again in order to
556 communicate with McKernel.
 - 557 (c) The child proxy creates the worker thread pool that serve the same role of the
558 parent process's worker threads.
 - 559 (d) The child proxy sends a reply message to McKernel.
- 560 5. When McKernel receives the reply message, it puts the child process into the run-
561 queue.
- 562 6. McKernel returns to its parent process with the child process ID.

563 2.2.3 Files and the File Descriptor Table

564 McKernel does not maintain file system related information (e.g., file caches) and file de-
565 scriptors are managed by the proxy process on Linux. When an McKernel process opens a
566 file, its file descriptor is created in the `mcexec` process and the number is merely returned
567 to the McKernel process.

568 It is worth noting that `mcexec` keeps the IHK device file open for communication with
569 McKernel. Because a file descriptor is an integer value, the IHK device could theoretically
570 be accessed from application code. In order to avoid such scenario, `mcexec` ensures that the
571 IHK device file cannot be accessed by application code.

572 **2.2.4 Signal Handling**

573 Two types of signals are considered: One is signals for the `mcexec` process. An example is
 574 the user sends a signal to the process from the shell. Another one is signals for a McKernel
 575 process, e.g., page fault signal caused by accessing wrong address in the McKernel process.

576 When the `mcexec` process receives a signal, that signal is transferred to the McKernel
 577 process via McKernel. When McKernel receives a signal for the McKernel process from the
 578 `mcexec` process during waiting for completion of a Linux system call, McKernel requests
 579 the `mcexec` process for aborting the system call execution.

580 図??を用いてシグナル中継機能の動作を説明する。ホスト OS の `mcexec` が受け取った
 581 シグナルは、IKC を通じて McKernel に通知され、シグナル登録処理 (`do_kill`) に伝えられ
 582 る。シグナル登録処理では、シグナルを表す `sig_pending` 構造体を作成し、シグナル送付先
 583 の `process` 構造体に登録する。ここで、シグナル送付先がスレッドの場合は `process` 構造体の
 584 `sigpending` に登録するが、スレッドを特定しないシグナルの場合は `process` 構造体の中のス
 585 レッド共通の `sigshared` の `sigpending` に登録する。他の事象により発生したシグナルも同
 586 様にシグナル登録処理 (`do_kill`) によって `process` 構造体にシグナルが登録される。シグナル
 587 を受信するプロセスを実行する CPU では、割り込み処理後やシステムコール処理後などの
 588 ユーザ空間への切り替えのタイミングでプロセスに届いているシグナル (`process` 構造体に登
 589 録されている `sig_pending` 構造体) をチェック (`check_signal`) し、シグナルが届いている場
 590 合には、その処理を行う。シグナルの処理は、`process` 構造体の `sighandler` に従って行う。
 591 `sighandler` のシグナル番号の項目にシグナルハンドラが登録されている場合は、登録されて
 592 いるシグナルハンドラを呼び出す。シグナルを無視する場合は何もしない。それ以外の場合
 593 はプロセスを終了 (シグナルによる終了) する (但し、シグナル番号が `SIGCHLD` と `SIGURG`
 では、シグナルハンドラの登録が無い場合は無視される)。

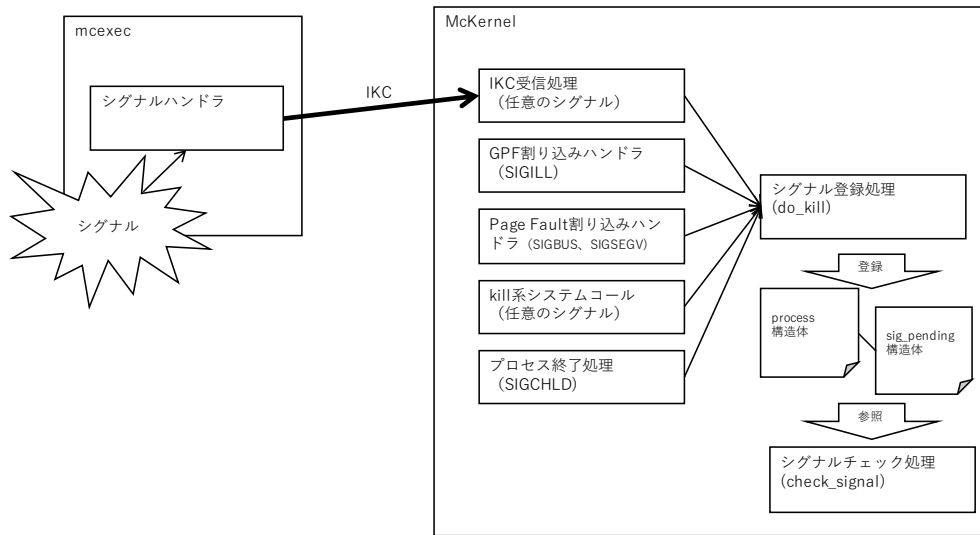


Figure 2.4: シグナル中継処理の動作

594

595 **2.2.5 Process ID**

596 The process ID of a McKernel process is held in the corresponding proxy process and it is
 597 managed via Linux API.

598 2.2.6 Thread ID

599 McKernel スレッドのスレッド ID は、対応する proxy process スレッドで管理される。McKernel
600 スレッド生成時の proxy process スレッドとの対応付けステップは以下の通り。

601 C1 proxy process (`mcexec`) は起動時に生成するスレッド数を決定し、その数だけ生成する。

602 C2 McKernel はスレッド生成時に、そのスレッドと対応付ける proxy process のスレッド
603 を proxy process に問い合わせる。

604 C3 McKernel は新しく生成する McKernel スレッドに当該 proxy process スレッドのスレッ
605 ド ID を割り当てる。また、McKernel はスレッド ID をキャッシュすることでスレッド
606 ID 問い合わせを高速化する。

607 `mcexec` は生成するスレッド数を以下の方法で決定する。

608 S1 `-t <nr_threads>` のオプションが指定された場合はその値を用いる。

609 S2 上記オプションが指定されなかった場合は、環境変数 `OMP_NUM_THREADS` が設定されてい
610 る場合は、環境変数の値を用いる。この環境変数が設定されていない場合は McKernel
611 に割り当てられた CPU 数を用いる。

612 McKernel のスレッド数上限は proxy process がステップ C1 で生成するスレッド数で決
613 まる。このためユーザは上記のステップ S2 で決定される数では足りない場合は `mcexec` の `-t`
614 `<nr_threads>` オプションを用いて十分な数を指定する必要がある。

615 2.2.7 User ID

616 UID 情報取得のオーバーヘッドを削減するため、UID は McKernel と Linux の両方で管理す
617 る。変更の際は McKernel 上の値を変更した後、IKC を用いて Linux 上の値を変更する。

618 2.2.8 Process Groups

619 プロセスグループにシグナルを送付する際のシグナル送付対象プロセス調査のオーバーヘッ
620 ドを削減するため、また、`setpgid` システムコールにおいて、対象プロセスが `execve` を実行
621 したか否かのチェックを行えるようにするため、`pgid` は Linux と McKernel の両方で管理す
622 る。変更の際は McKernel 上の値を変更した後、IKC を用いて Linux 上の値を変更する。

623 2.3 システムコール

624 As already mentioned, one of the proxy process' roles is to facilitate system call offloading
625 by providing an execution context on behalf of the application so that offloaded calls can
626 be directly invoked in Linux.

627 2.3.1 System Call Offloading

628 The main steps of system call offloading (also shown in Figure ??) are as follows. When
629 McKernel determines that a system call needs to be offloaded it marshalls the system
630 call number along with its arguments and sends a message to Linux via a dedicated IKC
631 channel. The corresponding proxy process running on Linux is by default waiting for system
632 call requests through an `ioctl1()` call into IHK's system call delegator kernel module. The

633 delegator kernel module’s IKC interrupt handler wakes up the proxy process, which returns
 634 to userspace and simply invokes the requested system call. Once it obtains the return value,
 635 it instructs the delegator module to send the result back to McKernel, which subsequently
 636 passes the value to user-space.

637 System call offloading internally relies on IHK’s Inter-Kernel Communication (IKC)
 638 facility. For more information on IKC, refer to “IHK Specifications”.

639 2.3.2 Offloading Strategy

640 There are mainly two categories of system calls that need to be implemented by McKernel:

- 641 1. System calls that cannot be offloaded to Linux side, and
- 642 2. Performance critical system calls

643 The first category includes CPU affinity system calls such as `sched_setaffinity()`,
 644 signaling system calls such as `sigaction()`, and memory-related system calls such as
 645 `mmap()` and `fork()`. The second category includes timer-related system calls such as
 646 `gettimeofday()`.

647 System calls, implemented in McKernel or planned to implement, is listed in Table ??.
 648 Other system calls are delated the Linux.

Table 2.1: System calls implemented in McKernel

Category	Implemented	Planned
Proess man- agement	arch_prctl (x86.64 specific), clone, execve, exit, exit_group, fork, futex, get_cpu_id, gete{u,g}id, get{g,p,t,u}id, getppid, getres{g,u}id, {get,set}rlimit, kill, pause, ptrace, rt_sigaction, rt_sigpending, rt_sigprocmask, rt_sigqueueinfo, rt_sigreturn, rt_sigsuspend, set_tid_address, setfs{u,g}id, set{g,u,t}id, setpgid, setre{g,u}id, setres{g,u}id, sigaltstack, tgkill, vfork, wait4, waittid	{get,set}_thread_area, rt_sigtimedwait, signalfd, signalfd4
Memory management	brk, {get,set}_mempolicy, madvise, mincore, mlock, mmap, move_pages, mprotect, mremap, msync, munlock, munmap, process_vm_{readv,writev}, remap_file_pages, shmat, shmctl, shmdt, shmget	{get,set}_robust_list, mbind, migrate_pages, mlockall, modify_ldt, munlockall
Schedule	getcpu, {get,set}itimer, {get,set}timeofday, nanosleep, sched_{get,set}affinity, sched_yield, times	
Performance counter	perf_event_open	

649 2.3.3 gettimeofday()

650 `gettimeofday()` is implemented in user-space by using Virtual Dynamic Shared Object
 651 (vDSO) mechanism (see Section ?? for vDSO).).

652 Table ?? shows the related vDSO pages.

Table 2.2: vDSO pages related to `gettimeofday()`

Name	Description
<code>vdso</code>	System call code and data
<code>vvar</code>	Kernel variables
<code>hpet</code>	Register of the High Precision Event Timer
<code>pvti</code>	Virtual clock updated by virtual machine, such as Xen and KVM

653 2.3.4 `perf_event_open()`

654 `perf_event_open()` is implemented in McKernel by using the technique mentioned in Sec-
655 tion ??.

656 2.4 Memory Management

657 We already described how system call offloading works in the IHK/McKernel architecture.
658 Notice, however, that certain system call arguments may be pointers (e.g., the buffer argu-
659 ment of a `read()` system call) and the actual operation takes place on the contents of the
660 referred memory. Thus, the main problem is how the proxy process on Linux can resolve
661 virtual addresses in arguments so that it can access the memory of the application running
662 on McKernel.

663 In order to overcome this problem McKernel deploys a mechanism called *unified ad-*
664 *dress space*, which essentially ensures that the proxy process can transparently access the
665 same mappings as its corresponding McKernel process. This mechanism is detailed in the
666 following sections.

667 2.4.1 Unified Address Space

668 The unified address space model in IHK/McKernel ensures that offloaded system calls can
669 seamlessly resolve arguments even in case of pointers. This mechanism is depicted in Figure
670 ?? and it is implemented as follows. First, the proxy process is compiled as a position
671 independent binary, which enables us to map the code and data segments specific to the
672 proxy process to an address range which is explicitly excluded from McKernel's user space.
673 The box on the right side of the figure with label "Not used" demonstrates the excluded
674 region. Second, the entire valid virtual address range of McKernel's application user-space
675 is covered by a special mapping in the proxy process for which we use a pseudo file mapping
676 in Linux. This mapping is indicated by the yellow box on the left side of the figure.

677 Note, that the proxy process does not need to fill in any virtual to physical mappings
678 at the time of creating the pseudo mapping and it remains empty unless an address is
679 referenced. Every time an unmapped address is accessed, however, the page fault handler
680 of the pseudo mapping consults the page tables corresponding to the application on the
681 LWK and maps it to the exact same physical page. Such mappings are demonstrated in the
682 figure by the small boxes on the left labeled as *faulted page*. This mechanism ensures that
683 the proxy process, while executing system calls, has access to the same memory content
684 as the application. Needless to say, Linux' page table entries in the pseudo mapping have
685 to be occasionally synchronized with McKernel, for instance, when the application calls
686 `munmap()` or modifies certain mappings.

687 A more detailed sequence of resolving a page fault in Linux for an address in the
688 McKernel process is as follows:

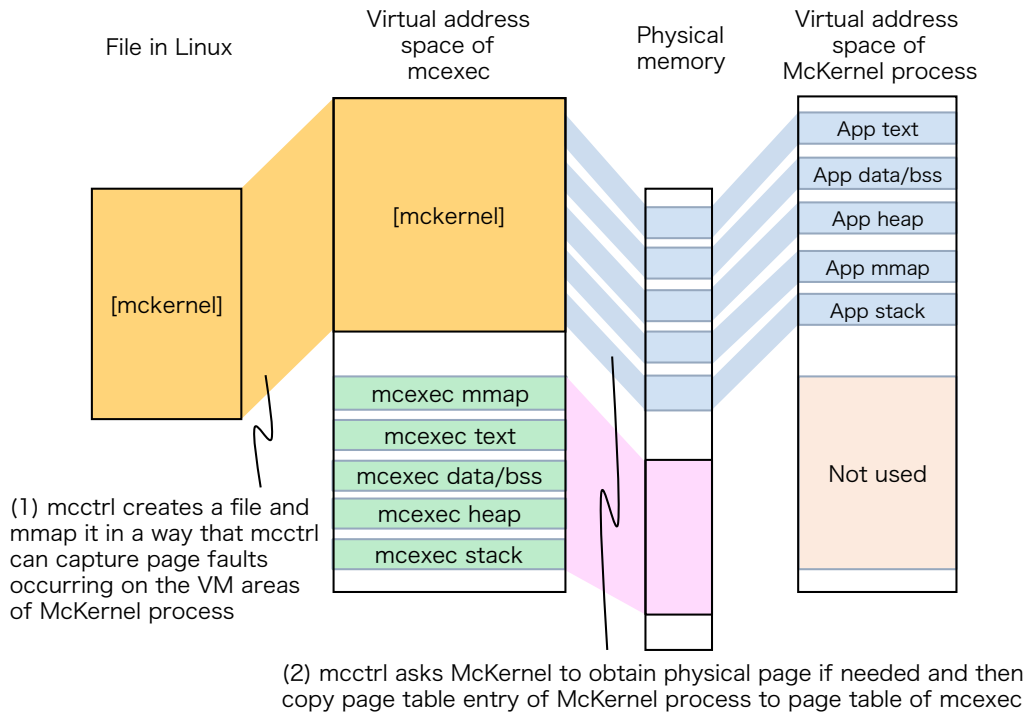


Figure 2.5: Unified Address Space

- 689 1. When `mcexec` accesses a memory area pointed by a pointer variable stored in a system
690 call request a Linux page fault occurs.
- 691 2. The `mcctrl` kernel module captures this page fault. It looks up the page table of the
692 McKernel process to find out the page table entry (PTE) of the physical memory.
- 693 3. In case that PTE is not found, the following sequences of issuing remote page fault
694 are performed as follows.
 - 695 (a) The `mcctrl` module interrupts the system call service. It reports return code
696 `STATUS_PAGE_FAULT` and the faulting address to McKernel.
 - 697 (b) When McKernel receives the return code `STATUS_PAGE_FAULT`, it resolves the
698 page fault.
 - 699 (c) After McKernel finishes page fault processing, it requests resuming the previous
700 system call process by sending an IKC message `SCD_MSG_SYSCALL_ONESIDE` to
701 `mcctrl`.
 - 702 (d) When `mcctrl` receives the request of resuming the previous system call at the
703 IKC message `SCD_MSG_SYSCALL_ONESIDE`, it looks up the page table entry again.
- 704 4. `mcctrl` maps the physical memory pointed by the PTE to the virtual address where
705 the page fault occurred.
- 706 5. `mcctrl` requests resuming the execution of the `mcexec` process.
- 707 6. The `mcexec` process now can access the virtual address requested in the system call.

708 As mentioned above when an McKernel process releases physical pages by issuing sys-
709 tem calls such as `munmap()` or `madvise()` with the option `MADV_REMOVE`, the `mcexec` process
710 clears its page tables to make sure future requests will not resolve an invalid mapping.

711 When the `mcexec` process establishes the pseudo mapping covering the McKernel pro-
712 cess's user space the mapping is read/write enabled except for the text area of the McKernel
713 process. When the McKernel process allocates a read-only memory mapping, e.g., when
714 mapping a shared library, the `mcctrl` kernel module remaps this area with the same access
715 permissions in the Linux side. This remap operation is required because the virtual address
716 space for the McKernel process has been created as one contiguous region whose access per-
717 mission is homogeneous. Most of memory mappings created by the McKernel process are
718 read/write permission, and thus such remap operation happens relatively rarely.

719 **2.4.1.1 McKernel Process Virtual Address Mapping**

720 Theoretically all virtual addresses used in the McKernel process must be mapped to the
721 `mcexec` process's virtual address. There are two issues as follows:

- 722 1. The `mcexec` process has its own text, data and BSS area whose addresses are also
723 used in the McKernel process if those execution binaries have been created in the
724 same way.
- 725 2. If the huge stack area is allocated to `mcexec` via shell environment variable `RLIMIT_STACK`,
726 the virtual address space for the McKernel process cannot be assigned.

727 The solution of those issues on Linux for `x86_64` architectures is described as follows.

728 **2.4.1.1.1 Avoiding Conflict of text, data, and BSS**

729 In the Linux convention for `x86_64` architectures, the text segment starts from virtual ad-
730 dress `0x400000` and the data segment starts from 2 MiB upper address than the text seg-
731 ment. If both an McKernel application and `mcexec` are compiled and linked, those addresses
732 are conflict.

733 As we briefly mentioned above, the `mcexec` binary is created as position independent
734 binary so that each segment's address can be dynamically decided by the runtime. In
735 Linux convention for `x86_64` architectures, by issuing `mmap`, the map address will be the
736 next to the address of the stack area whose address is the highest address in the user address
737 space.

738 **2.4.1.1.2 Huge Stack Size**

739 The virtual address space plan of the McKernel process follows Linux address plan, i.e.,
740 the user space is contiguous and starts from virtual address 0. That is, in order to keep the
741 same address space of the McKernel process in the `mcexec`, the same address space must
742 not be occupied by the `mcexec` process. There is one problem to do so. In Linux for `x86_64`
743 architectures, the start address of a stack area is randomly decided and its size is the lesser
744 of $\frac{5}{6}$ total memory size and size specified by the `RLIMIT_STACK` environment variable. If the
745 huge stack occupies the virtual memory in the `mcexec`, there is no chance to reserve the address
746 space for the McKernel process. In order to eliminate this problem, the `RLIMIT_STACK`
747 environmental variable for `mcexec` and the McKernel process is separated. That is, the

748 `mcexec` checks if `RLIMIT_STACK` is larger than some amount of size (currently 1 GiB), it
749 saves `RLIMIT_STACK` to a temporal environmental variable (`MCKERNEL_RLIMIT_STACK`) and
750 `exec()` itself again with a small stack (10 MiB). The new `mcexec` process restores the
751 original value to `RLIMIT_STACK` so that this environment variable is used for the McKernel
752 process.

753 2.4.2 Physical Pages requiring Linux Management

754 The physical pages of a McKernel process must be under Linux management for I/O related
755 operation (e.g. pin-down). This is because the driver running on the Linux side performs
756 I/O operation and the operation relies on the Linux paging mechanism. For example, when
757 a McKernel process tries to send data in a buffer to a remote host, it calls the Linux driver
758 code and the driver code in turn pins down the physical pages for the buffer using the
759 Linux kernel function. The kernel function in turn assumes that the pages are under Linux
760 management (i.e. managed by `struct page`).

761 Thus, IHK takes the physical pages from physical pages managed by the Linux. That
762 is, IHK reserves physical pages for the co-kernels by using `__get_free_pages()` Linux API.
763 Since the `__get_free_pages()` allocates up to 1024 pages at a time, IHK repeatedly calls
764 this function to get contiguous pages more than 1024 pages.

765 2.4.3 Handling Different Page Sizes

766 There are several implementation options to support different page sizes in Linux:

- 767 1. Linux Transparent Huge Pages (THP)
- 768 2. Hugepage option in System V IPC shared memory
- 769 3. Linux HugeTLBfs
- 770 4. Hugepage option in `mmap()` flags

771 McKernel implments a similar technique to Linux THP, i.e., it automatically maps
772 physical memory with large pages whenever it is possible.

773 2.4.4 `brk()`

774 McKernel の `brk()` システムコールには、ページフォルトオーバーヘッドを削減し、また
775 ラージページ化を促進する機能が追加されている。

776 `brk()` の動作は以下の通り。なお、ヒープ終端アドレスを b 、`brk()` の引数をページ境界
777 で丸め上げたアドレスを r で表す。また、プロセス起動コマンド `mcexec` (第??参照) のオプ
778 ション (`--extend-heap-by=<step>`) で指定されたパラメタを S で表す。

- 779 1. ヒープの縮小が要求された場合、何もせずに戻る。
2. $r - b < S$ の場合、ヒープ終端アドレスを以下の x に設定する。

$$r + S \leq x < r + S + a, x \bmod a = 0, a = \begin{cases} 2^{12} & \text{if } S \leq 4096; \\ 2^{21} & \text{if } S > 4096. \end{cases}$$

- 780 3. $r - b \geq S$ の場合、ヒープの最終アドレスを r に設定する。

781 4. 拡張された部分をプリページングする。

782 なお、この機能は、同一計算ノード上に他ユーザのジョブが存在することはないので、物
783 理ページ利用のフェアネスを考慮する必要がないため、不要になった物理ページを OS に返
784 す必要がない、という HPC アプリの特性を用いている。

785 2.4.5 メモリ割り当てにおける NUMA ノード選択

786 2.4.5.1 ユーザメモリ割り当て

787 ユーザメモリ割り当てにおける NUMA ノード選択については、Linux の機能に対し以下の機
788 能が追加されている。

- 789 1. ヒープ、anonymous mmap 領域だけではなく、text, data, bss, stack の各領域に対しても
790 プロセスのメモリポリシーを用いる。また、これらの領域ごとにプロセスのメモリポ
791 リシーを用いるか否かを指定できる。この指定は、プロセス起動コマンド `mcexec` (第
792 ??参照) のオプション (`--mpol-no-{heap,stack,bss}`) によって行う。
- 793 2. ユーザ指定のメモリポリシーを用いるメモリ要求サイズの閾値 (この値と同じか大きい
794 場合のみユーザ指定のメモリポリシーを用いる) を指定できる。この指定は、プロセス
795 起動コマンド `mcexec` (第??参照) のオプション (`--mpol-threshold=<min>`) によって
796 行う。

797 2.4.5.2 カーネルメモリ割り当て

798 カーネルメモリ割り当てにおける NUMA ノード選択は Linux と同様の方法で行う。すなわ
799 ち、NUMA ノード間の距離行列を用いて、要求元が存在する NUMA ノードから最も距離の
800 短い NUMA ノードからメモリを取得する。

801 2.4.6 Virtual Dynamic Shared Object (vDSO)

802 Mckernel provides the vDSO mechanism, which eliminates the need for switching to kernel-
803 mode when performing some system calls.

804 The steps are in the followings.

- 805 1. The physical addresses of the Linux vDSO pages are compiled by looking into `System.map`
806 when configuring McKernel. They are kept in `mcctrl`.
- 807 2. McKernel adds mappings of the Linux vDSO pages to a process when creating the
808 process. McKernel asks `mcctrl` for their physical addresses.
- 809 3. McKernel passes their virtual addresses to the process via the Auxiliary Vector in the
810 stack.
- 811 4. When a system call is called, first the control is transferred to the `glibc` wrapper
812 function. And then the control is transferred to the function in the vDSO pages
813 without switching processor mode.
- 814 5. The function performs required processing using the data in the vDSO pages.

815 2.4.7 ファイルマップ

816 ファイルマップはファイルと一対一対応する `fileobj` と呼ぶ構造体で管理する。ファイルマッ
817 プに伴うファイル I/O は、`fileobj` と一対一対応する、Linux 側に存在する `pager` と呼ぶ構
818 造体で管理する。ファイルマップの動作を例を用いて説明する。

- 819 1. 第 1 のプロセスが `open()` でファイルディスクリプタを取得する。
- 820 2. 第 1 のプロセスが前記ファイルディスクリプタを引数とした `mmap()` で McKernel にファ
821 イルマップ作成を要求する。
- 822 3. McKernel は `mcctrl` に `pager` を要求する。
- 823 4. `mcctrl` は `pager` のリストをファイルの `inode` で検索する。リストにないため新たな
824 `pager` を作成しリストに挿入し、その `pager` を返す。
- 825 5. McKernel は `fileobj` のリストを `pager` のアドレスで検索する。リストにないため新た
826 に `fileobj` を作成して、取得した `pager` と紐付けた上で、`fileobj` のリストに挿入す
827 る。また、`VM_range` 構造体の `memobj` フィールドにポインタを格納する。
- 828 6. 第 1 のプロセスがページフォールトを起こす。読み込みのページフォールトを起こした
829 とする。
- 830 7. McKernel が `fileobj` の `get_page()` を呼んで、以下のステップで物理ページを取得
831 する。
 - 832 (a) 割り当て済み物理ページを管理するハッシュリストをオフセットで検索する。ハッ
833 シュリストにないためアロケータを呼ぶことで新たな物理ページを取得し、ハッ
834 シュリストに挿入する。
 - 835 (b) `pager` に依頼して、当該物理ページにファイルの対応部分の内容を書き込む。
 - 836 (c) 取得した物理ページのアドレスを返す。
- 837 8. McKernel は取得した物理ページに対応するページテーブルエントリを作成し挿入する。
- 838 9. 第 1 のプロセスが当該物理ページに対する操作を行う。
- 839 10. 第 2 のプロセスが `open()` でファイルディスクリプタを取得する。
- 840 11. 第 2 のプロセスが前記ファイルディスクリプタを引数とする `mmap()` で McKernel にファ
841 イルマップ作成を要求する。
- 842 12. McKernel は `mcctrl` に `pager` を要求する。
- 843 13. `mcctrl` は `pager` のリストを `inode` で検索し、第 1 のプロセスからの依頼によって作成
844 された `pager` を返す。
- 845 14. McKernel は `fileobj` のリストを `pager` のアドレスで検索し、第 1 のプロセスによって
846 作成された `fileobj` を取得し、`VM_range` に記録する。
- 847 15. 第 2 のプロセスがページフォールトを起こす。読み込みのページフォールトを起こした
848 とする。
- 849 16. McKernel が `fileobj` の `get_page()` を呼ぶ。

- 850 17. McKernel は割り当て済み物理ページを管理するハッシュリストをオフセットで検索し、
851 第 1 のプロセスによって取得された物理ページを取得する。
- 852 18. McKernel は取得した物理ページに対応するページテーブルエントリを作成し挿入する。
- 853 19. 第 2 のプロセスが当該物理ページに対する操作を行う。

854 2.4.8 POSIX Shared Memory

855 McKernel の POSIX Shared Memory 機能 (`/dev/shm/*` ファイルのマッピングによる共有メモリ
856 機能) にはプリマップ機能が追加されている。この機能は `mcexec` (第??参照) のオプション
857 (`--mpol-shm-premap`) によって有効にできる。

858 2.4.9 System V 共有メモリ

859 System V 共有メモリ機能によるメモリ領域は、`shmobj` と呼ぶ構造体を用いて、ファイルマッ
860 プと同様に管理する。

861 共有メモリセグメントの属性は、`shmctl` システムコールで要求されたときにそのまま
862 ユーザに渡せるように、カーネル内でも `shmid_ds` 構造体の形で保持する。`shmid_ds` 構造体
863 は、対応する `shmobj` に内包させる。

864 共有メモリのマッピングが `munmap()` で部分解放されたときの共有メモリマッピングの
865 分断や、`fork()` による共有メモリマップ数の増加、プロセス終了による共有メモリマップ数
866 の減少といったマップ数の管理は、共有メモリのアタッチ数 `shm_nattch` で行う。

867 2.4.9.1 実装の制限

868 Linux の `shmget` システムコール仕様のうち、引数 `shmflg` に `SHM_NORESERVE` を指定した、ス
869 ヱップ領域予約なし共有セグメントの作成はサポートしない (指定は無視される)。

870 また、Linux の `shmctl` システムコールの仕様のうち、以下のものをサポートしない。

- 871 1. 引数 `cmd` に `SHM_LOCK` を指定した、共有メモリセグメントのページロック
- 872 2. 引数 `cmd` に `SHM_UNLOCK` を指定した、共有メモリセグメントのページロック解除

873 2.5 procfs/sysfs

874 The `procfs/sysfs` files provided by McKernel are listed in Table ?? and Table ??.

Table 2.3: /proc files provided by McKernel

Full path	Description
/proc/stat	Kernel statistics
/proc/[PID]	Directory containing information of [PID]
/proc/[PID]/auxv	Additional information to ELF loader
/proc/[PID]/cgroup	cgroup it belongs to
/proc/[PID]/cmdline	Command line
/proc/[PID]/cpuset	CPU set
/proc/[PID]/maps	List of memory maps
/proc/[PID]/mem	Memory held by this process
/proc/[PID]/pagemap	Flat page table
/proc/[PID]/smaps	An extension based on <code>maps</code> , showing the memory consumption of each mapping and flags associated with it
/proc/[PID]/stat	Process status
/proc/[PID]/status	Process status in human readable form
/proc/[PID]/task/[THID]	Directory containing information of [THID]
/proc/[PID]/task/[THID]/mem	Memory held by this thread
/proc/[PID]/task/[THID]/stat	Thread status

Table 2.4: /sys files provided by McKernel

Full path	Description
/sys/bus/cpu/devices/cpu*	Symbolic link to /sys/devices/system/cpu/cpu*
/sys/devices/system/cpu/offline	CPUs that are not online because they have been HOT-PLUGGED off or exceed the limit of cpus allowed by the kernel configuration
/sys/devices/system/cpu/online	CPUs that are online and being scheduled
/sys/devices/system/cpu/possible	CPUs that have been allocated resources and can be brought online if they are present
/sys/devices/system/cpu/present	CPUs that have been identified as being present in the system
/sys/devices/system/cpu/cpu*/online	1: Online, 0: Offline
/sys/devices/system/cpu/cpu*/cache/index*/level	Represents the hierarchy in the multi-level cache
/sys/devices/system/cpu/cpu*/cache/index*/type	Type of the cache - data, inst or unified
/sys/devices/system/cpu/cpu*/cache/index*/size	Total size of the cache
/sys/devices/system/cpu/cpu*/cache/index*/coherency_line_size	Size of each cache line usually representing the minimum amount of data that gets transferred from memory
/sys/devices/system/cpu/cpu*/cache/index*/number_of_sets	total number of sets, a set is a collection of cache lines sharing the same index
/sys/devices/system/cpu/cpu*/cache/index*/physical_line_partition	number of physical cache lines sharing the same cachetag
/sys/devices/system/cpu/cpu*/cache/index*/ways_of_associativity	Number of ways in which a particular memory block can be placed in the cache
/sys/devices/system/cpu/cpu*/cache/index*/shared_cpu_map	Set of CPUs sharing this cache in bitmap form
/sys/devices/system/cpu/cpu*/cache/index*/shared_cpu_list	Set of CPUs sharing this cache in human readable form
/sys/devices/system/cpu/cpu*/node*	Symbolic link to /sys/devices/system/node/node*
/sys/devices/system/cpu/cpu*/topology/physical_package_id	Physical package (e.g. socket) ID
/sys/devices/system/cpu/cpu*/topology/core_id	Core ID within a physical package
/sys/devices/system/cpu/cpu*/topology/core_siblings	Logical core set within a physical package in bitmap form. Logical cores include Hyperthreading cores.
/sys/devices/system/cpu/cpu*/topology/core_siblings_list	Logical core set within a physical package in human readable form
/sys/devices/system/cpu/cpu*/topology/thread_siblings	Logical core set within a physical core in bitmap form
/sys/devices/system/cpu/cpu*/topology/thread_siblings_list	Logical core set within a physical core in human readable form
/sys/devices/system/node/online	Numa nodes that are online
/sys/devices/system/node/possible	Nodes that could be possibly become online at some point
/sys/devices/system/node/node*/distance	Distance between the node and all the other nodes in the system
/sys/devices/system/node/node*/cpumap	Logical core set in the node in bitmap form
/sys/devices/system/node/node*/cpu*	Symbolic link to /sys/devices/system/cpu/cpu*
/sys/devices/pci<dom>:/bus/<dom>:<bus>:<slot>.<func>/local_cpus	Nearby CPU mask (logical core set in bitmap form)
/sys/devices/pci<dom>:/bus/<dom>:<bus>:<slot>.<func>/local_cpulist	Nearby CPU mask (logical core set in human readable form)
/sys/devices/system/cpu/num_processors	Number of logical cores (McKernel extension)

875 procfs/sysfs 機能は、以下の 3 機能で実現する。

- 876 ● McKernel がその内容を提供する procfs/sysfs と、Linux のそれとを優先度付きで重
877 ね合わせ、さらに重ね合わせたファイルシステムを /proc や /sys で始まる標準パスで
878 mcexec に見せる機能 (mcoverlayfs)
- 879 ● コールバック関数を mcctrl と McKernel の両方から登録できるようにし、またアクセ
880 ス要求を Linux から McKernel へ転送する機能

881 以下、それぞれの機能を説明する。

882 2.5.1 ファイルシステムの重ね合わせ

883 ファイルシステムの重ね合わせのステップは以下の通り。

- 884 1. McKernel が /proc/mcos0 を作成する。
- 885 2. mcoverlayfs を用いて、/proc/mcos0/と /proc を重ね合わせ /tmp/mcos/mcos0_proc
886 にマウントする。また、mcoverlayfs の機能を用いて、前者と後者に同一ファイルが存
887 在する際には、前者がアクセスされるように設定する。さらに、/tmp/mcos/mcos0_proc
888 を /proc に bind mount する。こうすることで、/proc に存在するファイルであって、
889 McKernel プロセスに Linux プロセスとは異なる内容を見せる必要のないものについ
890 ては /proc/mcos0/ に当該ファイルを準備しないことで元々の /proc のファイルを見せる
891 ことができる。また、異なる内容を見せる必要のあるものについては、/proc/mcos0/ に
892 当該ファイルを準備することでそれを見せることができる。
- 893 3. McKernel が同様に、/sys/devices/virtual/mcos/mcos0/sys を作成し、/sys/devices/
894 virtual/mcos/mcos0/sys と /sys を重ね合わせ /tmp/mcos/mcos0_sys にマウントし、
895 /tmp/mcos/mcos0_sys を /sys に bind mount する。
- 896 4. mcctrl と McKernel が /proc/mcos0 および /sys/devices/virtual/mcos/mcos0/sys
897 のファイル・ディレクトリを作成する。なお、ファイル・ディレクトリの内容は作成時
898 に登録するアクセスコールバック関数によって提供される。
- 899 5. McKernel プロセスが /proc または /sys ファイルにアクセスする。アクセス要求は必要
900 に応じて Linux から McKernel に転送される。

901 2.5.2 アクセス要求の Linux から McKernel への転送

902 アクセス要求の Linux から McKernel への転送の動作を図??を用いて説明する。

- 903 1. アプリは open() などのシステムコールを用いて procfs/sysfs のファイルへのアクセス
904 を試みる。このシステムコールの処理は Linux 側に転送される。(図の (1))
- 905 2. mcexec がシステムコールを代理実行する。mcoverlayfs が優先度に基づいて McKernel
906 が提供するファイルまたは Linux が提供するファイルへのアクセス振り分けを行う。こ
907 の場合は前者に振り分けられたとする。(図の (2))
- 908 3. McKernel が提供するファイルに登録されたコールバック関数が呼び出される。(図の
909 (3))
- 910 4. Linux 側 procfs/sysfs フレームワークがアクセス要求を McKernel 側フレームワーク
911 に転送する。McKernel 側フレームワークはアクセスに応じた処理を行う。例えば、ファ
912 イルの内容を返却する。(図の (4))

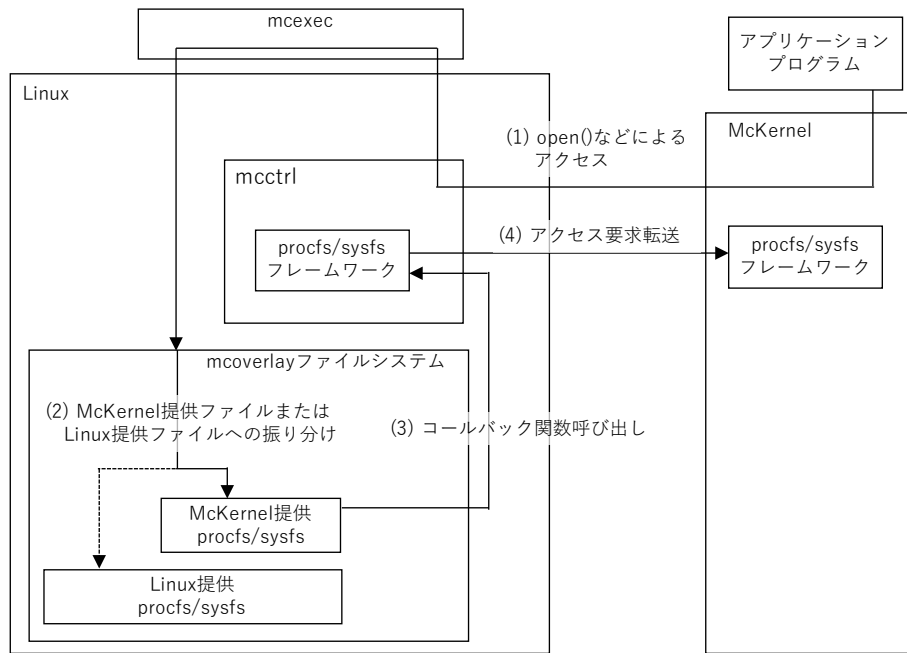


Figure 2.6: procfs/sysfs のアクセス要求転送

2.6 ファイルシステム重ね合わせ

McKernel はカーネルモジュール `mcoverlayfs` によって 2 つのファイルシステムを優先度付きで重ね合わせることができる。本機能は `procfs/sysfs` 機能のために用いられる。

`mcoverlayfs` は `overlayfs` に以下の機能を追加することで実装されている。

1. copyup 処理無効化

`mcoverlayfs` は `lowerdir` と `upperdir` に重ね合わせたいファイルシステムを指定する。McKernel では、`lowerdir` に、McKernel がその内容を提供する `procfs/sysfs` と Linux のそれとを指定して用いる。`overlayfs` では、ライト対象のファイルが `lowerdir` 上のファイルの場合、`copyup` 処理を行い `upperdir` に対象ファイルを作成し、そのファイルをオープンすることで、ライト処理を可能とする。このようにすると、アクセス要求は `procfs/sysfs` に届かない。そのため、`copyup` 処理を無効化し、直接対象ファイルにライト処理する機能を追加する。本機能はオプションに `nocopyupw` を指定することで有効となる。

`nocopyupw` オプションの有無によるライト処理の違いを図 ?? に示す。

- `nocopyupw` オプションなしで、ライト対象のファイルが `lowerdir` 上のファイルの場合、`copyup` 処理を行い `upperdir` に対象ファイルを作成し、そのファイルをオープンすることで、ライト処理を可能とする。
- `nocopyupw` オプションありの場合、ライト対象のファイルが `lowerdir` 上のファイルの場合でも、`copyup` 処理せず、そのファイルをオープンすることで、ライト処理を可能とする。

2. procfs/sysfs サポート

`mcoverlayfs` では `overlayfs` に対して `procfs/sysfs` のディレクトリのマウント機能を追加している。本機能はオプションに `nofscheck` を指定することで有効となる。

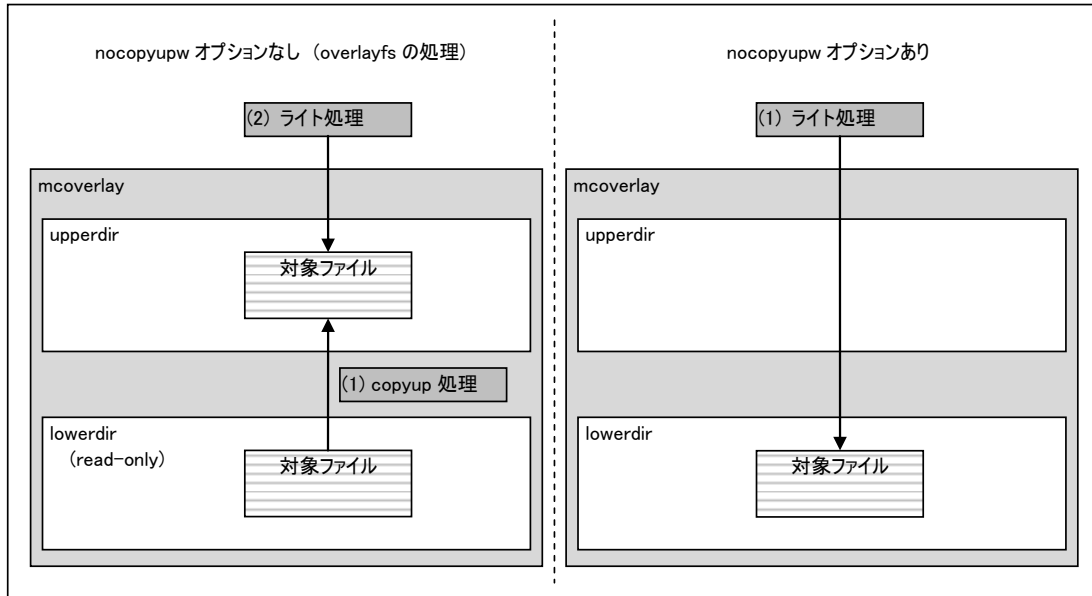


Figure 2.7: nocopyupw オプションの有無によるライト処理の違い

936 mcoverlayfs のマウントオプションを表?? に示す。nocopyupw, nofscheck が overlayfs に対して追加されたオプションである。

Table 2.5: mcoverlayfs のマウントオプション

オプション	説明
lowerdir=<dirs>	lowerdir を指定する。 ';' で区切り複数指定可能 (最大 500)
upperdir=<dir>	upperdir を指定する。
workdir=<dir>	workdir を指定する。 workdir は、upperdir と同じマウント下のディレクトリでなければならない。
default_permissions	デフォルトパーミッションを設定する。
nocopyupw	書き込み時に upperdir にファイルを作成し、それに対して書き込みを行う処理 (copyup 処理) を無効にする。 procfs/sysfs を lowerdir に指定する際は本オプションを指定する必要がある。
nofscheck	procfs/sysfs を lowerdir に指定可能にする。

937

938 2.6.1 詳細

939 overlayfs のデータ構造に対する修正は以下の通り。

- 940 1. ovl_opt_bit
- 941 マウントオプションを追加するために、以下の enum 及び、マクロを追加する。

```
942 enum ovl_opt_bit {
943     __OVL_OPT_DEFAULT    = 0,
```



```

944         __OVL_OPT_NOCOPYUPW      = (1 << 0),
945         __OVL_OPT_NOFSCHECK     = (1 << 1),
946     };
947
948     #define OVL_OPT_NOCOPYUPW(opt) ((opt) & __OVL_OPT_NOCOPYUPW)
949     #define OVL_OPT_NOFSCHECK(opt) ((opt) & __OVL_OPT_NOFSCHECK)

```

2. `ovl_d_fsdata`
`d_fsdata` を格納するために、以下の構造体を追加する。

```

952     struct ovl_d_fsdata {
953         struct list_head list;
954         struct dentry *d;
955         struct ovl_entry *oe;
956     };

```

3. `ovl_config`
マウントオプションを追加するために、`opt` を追加する。

```

959     struct ovl_config {
960         char *lowerdir;
961         char *upperdir;
962         char *workdir;
963         bool default_permissions;
964         unsigned opt;          <-- 追加
965     };

```

4. `ovl_fs`
`d_fsdata` を格納するために、`d_fsdata_list` を追加する。

```

968     struct ovl_fs {
969         struct vfsmount *upper_mnt;
970         unsigned numlower;
971         struct vfsmount **lower_mnt;
972         struct dentry *workdir;
973         long lower_namelen;
974         /* pathnames of lower and upper dirs, for show_options */
975         struct ovl_config config;
976         struct list_head d_fsdata_list; <-- 追加
977     };

```

5. `ovl_tokens`
マウントオプションを追加するために、`OPT_NOCOPYUPW` 及び、`OPT_NOFSCHECK` を追加する。

```

981     enum {
982         OPT_LOWERDIR,
983         OPT_UPPERDIR,
984         OPT_WORKDIR,
985         OPT_DEFAULT_PERMISSIONS,
986         OPT_NOCOPYUPW,          <-- 追加
987         OPT_NOFSCHECK,        <-- 追加
988         OPT_ERR,
989     };
990
991     static const match_table_t ovl_tokens = {
992         {OPT_LOWERDIR,          "lowerdir=%s"},
993         {OPT_UPPERDIR,         "upperdir=%s"},

```

```

994         {OPT_WORKDIR,                "workdir=%s"},
995         {OPT_DEFAULT_PERMISSIONS,    "default_permissions"},
996         {OPT_NOCOPYUPW,              "nocopyupw"},          <-- 追加
997         {OPT_NOFSCHECK,              "nofscheck"},          <-- 追加
998         {OPT_ERR,                    NULL}
999     };

```

1000 6. ovl_fs_type
1001 name の値を”mcoverlay”に変更する。

```

1002     static struct file_system_type ovl_fs_type = {
1003         .owner          = THIS_MODULE,
1004         .name           = "mcoverlay",      <-- 変更
1005         .mount          = ovl_mount,
1006         .kill_sb        = kill_anon_super,
1007     };
1008     MODULE_ALIAS_FS("mcoverlay");         <-- 変更

```

1009 overlayfs に対する関数の修正を表??、表??に示す。

Table 2.6: overlayfs の関数に対する修正 (1)

関数	修正内容
<code>ovl_copy_xattr()</code>	<code>OVL_OPT_NOFSCHECK(opt)</code> が有効の場合、 <code>vfs_getxattr()</code> のエラーを無視する。
<code>ovl_copy_up_locked()</code>	<code>ovl_copy_xattr()</code> 呼び出し時に <code>ovl_get_config_opt()</code> で取得した <code>opt</code> 値を渡す。
<code>ovl_clear_empty()</code>	<code>ovl_copy_xattr()</code> 呼び出し時に <code>ovl_get_config_opt()</code> で取得した <code>opt</code> 値を渡す。
<code>ovl_setattr()</code>	<code>ovl_get_config_opt()</code> で <code>opt</code> 値を取得する。 <code>OVL_OPT_NOCOPYUPW(opt)</code> の場合、処理しない。
<code>ovl_permission()</code>	<code>ovl_reset_ovl_entry()</code> を呼び出してから処理する。
<code>ovl_setxattr()</code>	<code>ovl_get_config_opt()</code> で <code>opt</code> 値を取得する。 <code>OVL_OPT_NOCOPYUPW(opt)</code> の場合、処理しない。
<code>ovl_removexattr()</code>	<code>ovl_get_config_opt()</code> で <code>opt</code> 値を取得する。 <code>OVL_OPT_NOCOPYUPW(opt)</code> の場合、処理しない。
<code>ovl_d_select_inode()</code>	<ol style="list-style-type: none"> <code>ovl_get_config_opt()</code> で <code>opt</code> 値を取得する。 <code>OVL_OPT_NOCOPYUPW(opt)</code> の場合、<code>ovl_open_need_copy_up()</code> を呼び出さない。 <code>OVL_OPT_NOFSCHECK(opt)</code> で対象ファイルが <code>sysfs</code> の場合、<code>ovl_find_d_fsdata()</code> を呼び出して <code>dentry</code> が登録されているか確認する。登録されていない場合には <code>ovl_add_d_fsdata()</code> を呼び出して登録し、<code>dentry->d_fsdata</code> に <code>realpath.dentry->d_fsdata</code> の値を設定する。
<code>ovl_get_config_opt()</code>	<code>opt</code> 値を返す。
<code>ovl_reset_ovl_entry()</code>	<ol style="list-style-type: none"> <code>ovl_get_config_opt()</code> で <code>opt</code> 値を取得する。 <code>OVL_OPT_NOFSCHECK(opt)</code> の場合、<code>ovl_find_d_fsdata()</code> を呼び出して、<code>dentry</code> が登録されている場合には取得した <code>d_fsdata</code> を <code>oe</code> に設定する。
<code>ovl_find_d_fsdata()</code>	<code>dentry->d_sb->s_fs_info</code> の <code>d_fsdata_list</code> に登録されている <code>d_fsdata</code> を検索して、 <code>dentry</code> が登録されていた場合、 <code>dentry</code> の <code>ovl_entry</code> を戻す。
<code>ovl_add_d_fsdata()</code>	<ol style="list-style-type: none"> <code>struct ovl_d_fsdata</code> のメモリ領域を確保して、<code>dentry</code> の登録データを設定する。 <code>dentry->d_sb->s_fs_info</code> の <code>d_fsdata_list</code> に登録する。
<code>ovl_clear_d_fsdata()</code>	<code>d_fsdata_list</code> に登録されている全ての <code>d_fsdata</code> を削除して、 <code>struct ovl_d_fsdata</code> のメモリ領域を解放する。
<code>ovl_path_type()</code>	<code>ovl_reset_ovl_entry()</code> を呼び出してから処理する。
<code>ovl_path_upper()</code>	
<code>ovl_dentry_upper()</code>	
<code>ovl_dentry_lower()</code>	
<code>ovl_dentry_real()</code>	
<code>ovl_dir_cache()</code>	
<code>ovl_set_dir_cache()</code>	
<code>ovl_path_lower()</code>	
<code>ovl_dentry_is_opaque()</code>	
<code>ovl_dentry_set_opaque()</code>	
<code>ovl_dentry_update()</code>	
<code>ovl_dentry_version_inc()</code>	
<code>ovl_dentry_version_get()</code>	
<code>ovl_dentry_release()</code>	
<code>ovl_dentry_revalidate()</code>	
<code>ovl_dentry_weak_revalidate()</code>	

Table 2.7: overlayfs の関数に対する修正 (2)

関数	修正内容
<code>ovl_lookup_real()</code>	<code>OVL_OPT_NOFSCHECK(opt)</code> の場合、 <code>ovl_dentry_weird()</code> を呼び出さない。
<code>ovl_path_next()</code>	<code>ovl_reset_ovl_entry()</code> を呼び出してから処理する。
<code>ovl_lookup()</code>	<ol style="list-style-type: none"> 1. <code>ovl_get_config_opt()</code> で <code>opt</code> 値を取得する。 2. <code>ovl_reset_ovl_entry()</code> を呼び出してから処理する。 3. <code>ovl_lookup_real()</code> を呼び出す際、<code>opt</code> 値を渡す。
<code>ovl_put_super()</code>	<code>ovl_clear_dfsdata()</code> を呼び出してから処理する。
<code>ovl_statfs()</code>	<code>struct kstatfs</code> の <code>f_type</code> に <code>MCOVERLAYFS_SUPER_MAGIC</code> を設定する。
<code>ovl_show_options()</code>	<code>nocopyupw</code> , <code>nofscheck</code> オプションの説明を追加する。
<code>ovl_parse_opt()</code>	<code>nocopyupw</code> , <code>nofscheck</code> オプションの設定を追加する。
<code>ovl_mount_dir_noesc()</code>	<code>OVL_OPT_NOFSCHECK(opt)</code> の場合、 <code>ovl_dentry_weird()</code> を呼び出さない。
<code>ovl_mount_dir()</code>	<code>ovl_mount_dir_noesc()</code> を呼び出す際、 <code>opt</code> を渡す。
<code>ovl_lower_dir()</code>	<code>ovl_mount_dir_noesc()</code> を呼び出す際、 <code>opt</code> を渡す。
<code>ovl_fill_super()</code>	<ol style="list-style-type: none"> 1. <code>struct ovl_fs</code> の <code>dfsdata_list</code> を初期化する。 2. <code>ovl_mount_dir()</code> を呼び出す際、<code>opt</code> を渡す。 3. <code>ovl_lower_dir()</code> を呼び出す際、<code>opt</code> を渡す。 4. <code>OVL_OPT_NOCOPYUPW(opt)</code> の場合、以下の設定を行わない。 <ul style="list-style-type: none"> • <code>mnt->mnt_flags = MNT_READONLY;</code> • <code>sb->s_flags = MS_RDONLY;</code>

1010 2.6.2 実装の制限

1011 McKernel が生成する `/proc/[pid]/` 下のファイルを open して、close せずに open した状態
1012 で exec して、exec したプロセスで同一ファイルを open するとエラー (ENOENT) となる。原因
1013 は、exec() 時には、新たなプロセスの情報を返せるようにするため `/proc/[pid]/` 下のファ
1014 イルを作成し直すが、overlayfs は lower に指定されるディレクトリ下のファイルの inode 番
1015 号が変わった場合、エラーを返すためである。

1016 2.6.3 開発時の留意事項

1017 Linux-4.0 から Linux-4.6 への移行に際する仮想ファイルシステムの以下の仕様変更に従
1018 する必要があった。

- 1019 1. `struct inode_operations` の `dentry_open()` が削除されて、`struct dentry_operations`
1020 の `d_select_inode()` が追加された。
- 1021 2. VFS の `vfs_open()` では、`dentry_open()` が呼ばれずに、`d_select_inode()` が呼び出
1022 されるようになった。

1023 また、以下のバージョンの Linux カーネルでのみ動作する。

- 1024 ● 3.10.0-327 から 3.10.0-693 (RHEL-7.2 から 7.4)
- 1025 ● 4.0.0 から 4.1.0
- 1026 ● 4.6.0 から 4.7.0

1027 2.7 デバイスドライバ

1028 McKernel では、Linux で動作するドライバをそのまま利用可能であるが、システムコール移
1029 譲のオーバーヘッドを削減するために、McKernel 内部で実装することもできる。

1030 以下、それぞれの方法を説明する。

1031 2.7.1 Linux ドライバの利用

1032 Linux ドライバ経由でメモリマップされたデバイスのレジスタを McKernel プロセスからアク
1033 セス可能にすることで、Linux ドライバをそのまま利用できるようにする。

動作を図??を用いて説明する。

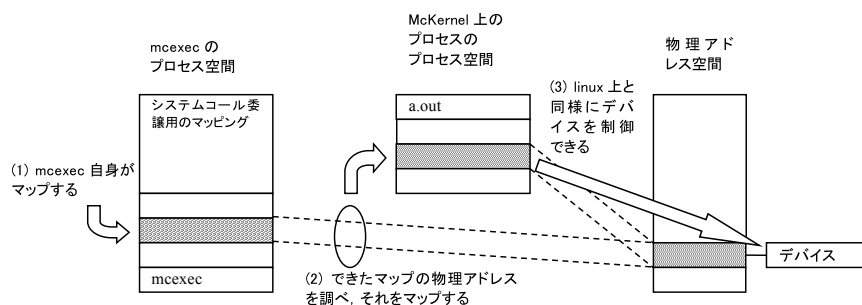


Figure 2.8: Linux ドライバ利用の動作

1034

- 1035 1. システムコール移譲の仕組みを用いてデバイスファイルの `open()`, `ioctl()` を行う。
- 1036 2. レジスタのマッピングについてはシステムコール移譲の仕組みを用いて、`mcexec` 空間への
- 1037 マッピングと、McKernel の仮想メモリ領域構造体への特別なマッピングであることの記録を行
- 1038 う。(図の (1))
- 1039 3. McKernel でのページフォールトの際に Linux に物理ページを問い合わせ、同じ物理ペー
- 1040 ジを参照するマッピングを McKernel 上のプロセス空間に作成する。(図の (2))

1041 2.7.2 McKernel 内部での実装

1042 特定のデバイスファイルに対して、`open()` 時に McKernel 内で処理を行うことをプロセス構

1043 造体に記録しておき、ファイル操作のシステムコールの際にその記録を参照することで、そ

1044 れらのファイルに対する操作を McKernel 内で行う。動作は以下の通り。

- 1045 1. プロセスが `open()` を呼び出した際に対象が McKernel 内で処理を行うデバイスファイ
- 1046 ルであるかをパスにより調べる。そうであった場合は、ダミーの `fd` を取得し、`struct`
- 1047 `process` の `struct mckfd` のリストに `fd` に対応するエントリを挿入する。また、その
- 1048 エントリにファイル操作のコールバック関数を登録する。
- 1049 2. プロセスがファイル操作のシステムコールを呼び出した際に、`struct process` の `struct`
- 1050 `mckfd` のリストに `fd` に対応するエントリが存在するかを調べる。存在する場合は、当
- 1051 該エントリに登録されているコールバック関数を呼び出す。

1052 2.8 XPMEM ドライバ

1053 XPMEM は、あるプロセスがマッピングしたメモリ領域を他のプロセスからマッピングできるよ

1054 うにする。XPMEM はユーザライブラリ部分とドライバ部分に分かれており、ドライバ部分は

1055 McKernel 内部で実装されている。

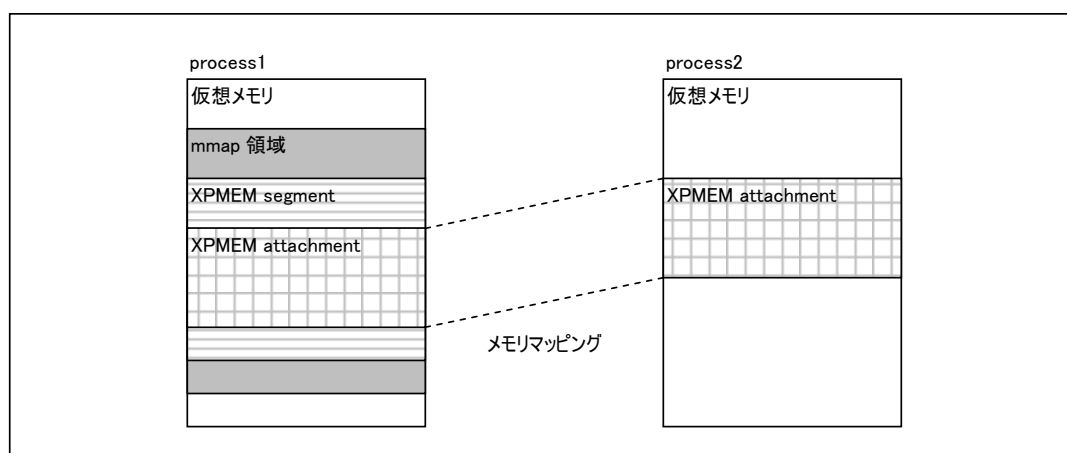


Figure 2.9: XPMEM のメモリマッピング



Figure 2.10: XPMEM の動作フロー

1056 XPMEMのメモリマッピングを図??、動作フローを図??に示す。XPMEMでは、プロ
 1057セス (process1) が mmap したメモリ領域から、xpmem_make() で指定された領域を XPMEM
 1058 segment として管理して他のプロセスからマップできるようにする。マップしたいプロセス
 1059 (process2) は、xpmem_get() でアクセスパーミッションを得て、xpmem_attach() で指定され
 1060 た XPMEM segment のメモリ領域を XPMEM attachment として管理して、マップする。マッ
 1061 プは、プロセス (process2) が XPMEM attachment 領域にアクセスして、ページフォルトが発
 1062 生した際、ページテーブルエントリが示す物理アドレスを、プロセス (process1) の XPMEM
 1063 segment 領域の物理アドレスに置き換えることで実現する。

XPMEMのデータ構造を生成・破棄する関数を図??に示す。

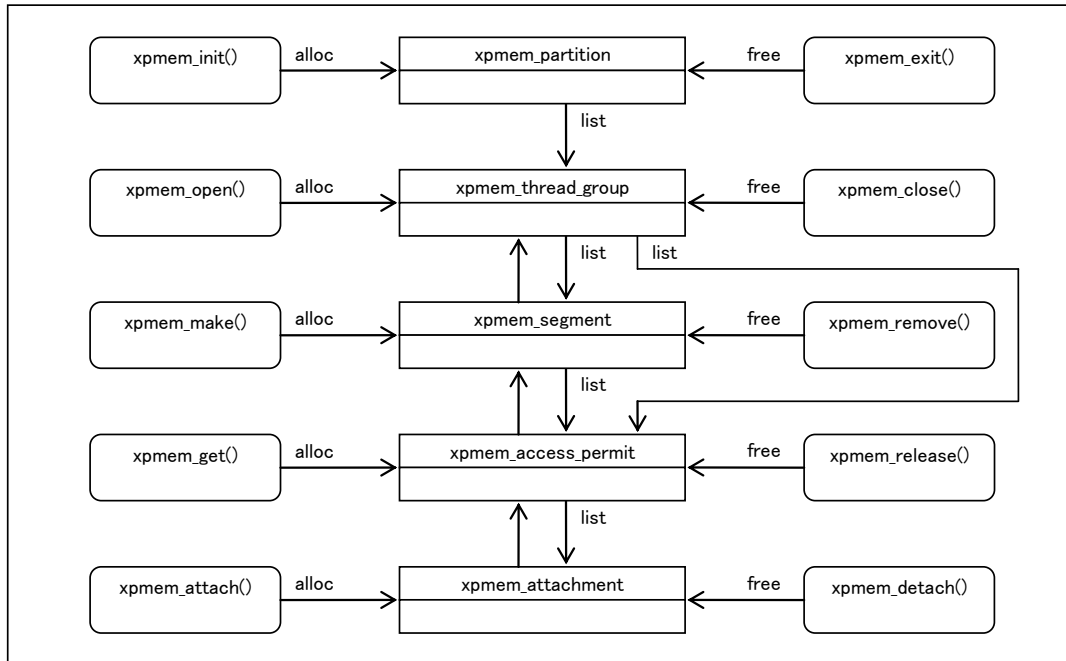


Figure 2.11: XPMEMのデータ構造を生成・破棄する関数

1064

1065

XPMEMでは、以下のデータ構造を管理して機能を実現する。

1066

1. xpmem_partition

1067

2. xpmem_thread_group

1068

3. xpmem_segment

1069

4. xpmem_access_permit

1070

5. xpmem_attachment

1071

XPMEMのデータ構造を図??に示す。

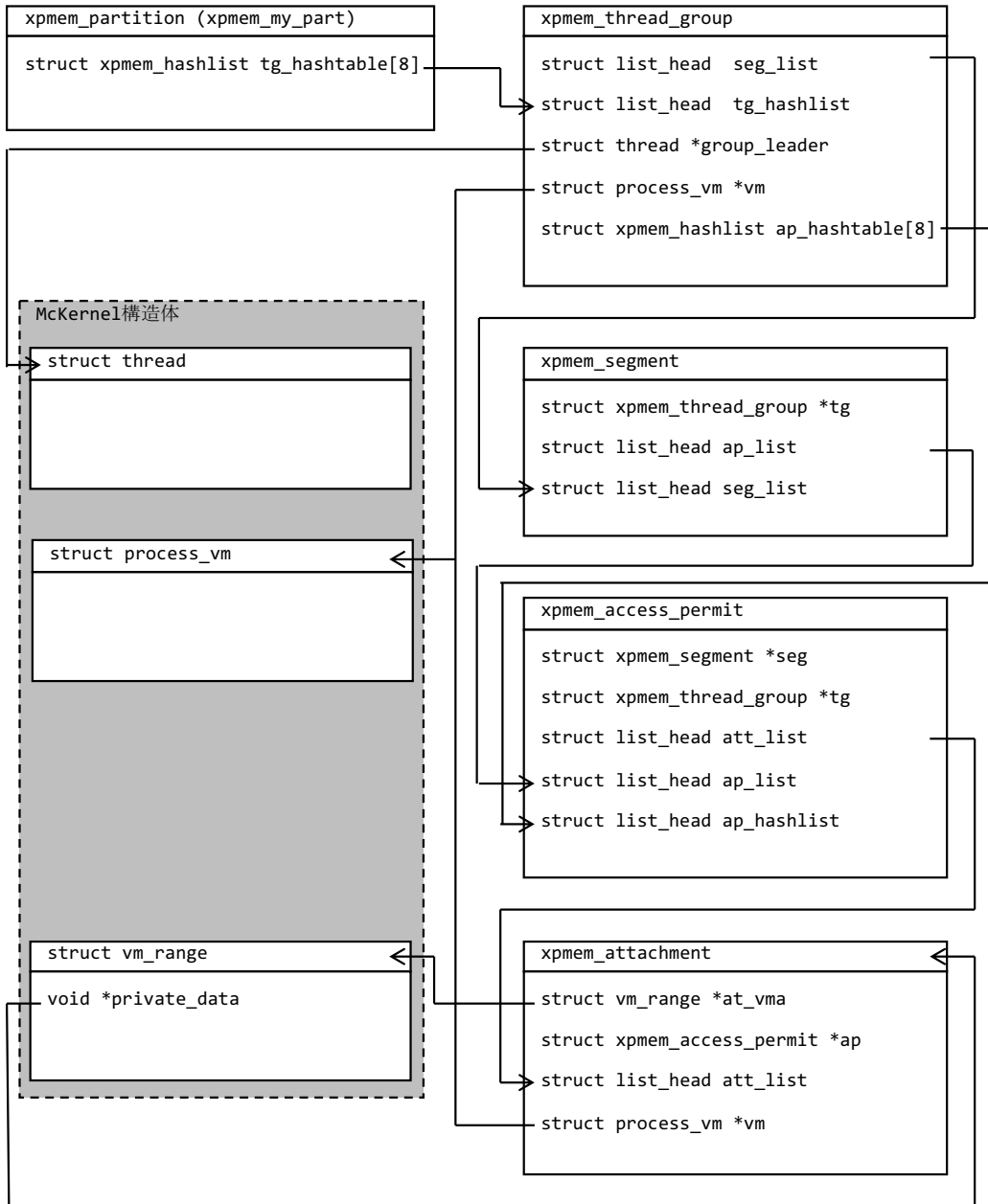


Figure 2.12: XPMEM のデータ構造

1072 以下、各関数のインターフェイスと動作を説明する。

1073 2.8.1 XPMEM デバイスファイルのオープン

1074 書式

```
1075 int xpmem_open(ihk_mc_user_context_t *ctx)
```

1076 説明

1077 xpmem_open は、以下の処理を行う。

- 1078 1. XPMEM を初期化していない場合には、xpmem_init() を呼び出して XPMEM を初期
1079 化する。
- 1080 2. fd を取得する必要があるが、/dev/xpmem デバイスファイルには影響を与えないよう
1081 に、do_syscall() 呼び出しで、/dev/null デバイスファイルをオープンして、その fd を
1082 使用する。fd がマイナス値の場合にはエラー値を戻す。
- 1083 3. _xpmem_open() を呼び出して、自プロセスの xpmem_thread_group が生成されていな
1084 ければ生成する。
- 1085 4. mckfd を生成して、初期設定する。

1086 戻り値

0 以上	ファイルディスクリプタ (正常終了)
-EINVAL	引数が無効である
-ENOMEM	十分な空きメモリ領域が無い

1087

1088 2.8.2 XPMEM デバイスファイルの ioctl 制御

1089 書式

```
1090 static int xpmem_ioctl(struct mckfd *mckfd, ihk_mc_user_context_t *ctx)
```

1091 説明

1092 xpmem_ioctl は、以下の処理を行う。

- 1093 1. cmd を処理する関数を呼び出す。
 - 1094 (a) XPMEM_CMD_VERSION
1095 XPMEM_CURRENT_VERSION を戻す。
 - 1096 (b) XPMEM_CMD_MAKE
1097 xpmem_cmd_make データを取得する。
1098 xpmem_make() を呼び出す。
1099 xpmem_cmd_make データの segid を設定する。

- 1100 (c) XPMEM_CMD_REMOVE
 1101 xpmem_cmd_remove データを取得する。
 1102 xpmem_remove() を呼び出す。
- 1103 (d) XPMEM_CMD_GET
 1104 xpmem_cmd_get データを取得する。
 1105 xpmem_get() を呼び出す。xpmem_cmd_get データの apid を設定する。
- 1106 (e) XPMEM_CMD_RELEASE
 1107 xpmem_cmd_release データを取得する。
 1108 xpmem_release() を呼び出す。
- 1109 (f) XPMEM_CMD_ATTACH
 1110 xpmem_cmd_attach データを取得する。
 1111 xpmem_attach() を呼び出す。xpmem_cmd_attach データの vaddr を設定する。
- 1112 (g) XPMEM_CMD_DETACH
 1113 xpmem_cmd_detach データを取得する。
 1114 xpmem_detach() を呼び出す。

1115 戻り値

0	正常終了
-EFAULT	アドレスが不正である
-EINVAL	引数が無効である

1116

1117 2.8.3 XPMEM デバイスファイルのクローズ

1118 書式

```
1119 static int xpmem_close(struct mckfd *mckfd, ihk_mc_user_context_t *ctx)
```

1120 説明

- 1121 xpmem_close は、以下の処理を行う。
- 1122 1. pid から xpmem_thread_group(tg) を取得する。
 - 1123 2. xpmem_release_aps_of_tg() を呼び出して、xpmem_access_permit、xpmem_attachment
1124 を破棄する。
 - 1125 3. xpmem_remove_segs_of_tg() を呼び出して、xpmem_segment を破棄する。
 - 1126 4. xpmem_destroy_tg() を呼び出して、xpmem_thread_group を破棄する。
 - 1127 5. /dev/xpmem をオープンしているプロセスが存在しない場合には、xpmem_exit() を呼
1128 び出して XPMEM を終了する。
 - 1129 6. xpmem_open() でオープンした/dev/null デバイスファイルについては、sys_close() で
1130 クローズする。

1131 戻り値

1132

0	正常終了
---	------

1133 2.8.4 XPMEM の初期化

1134 書式

```
1135 static int xpmem_init(void)
```

1136 説明

1137 xpmem_init は、以下の処理を行う。

- 1138 1. xpmem_partition を生成して、初期設定する。

1139 戻り値

0	正常終了
-ENOMEM	十分な空きメモリ領域が無い

1140

1141 2.8.5 XPMEM の終了

1142 書式

```
1143 static void xpmem_exit(void)
```

1144 説明

1145 xpmem_exit は、以下の処理を行う。

- 1146 1. xpmem_partition を破棄する。

1147 戻り値

1148 なし。

1149 2.8.6 xpmem_segment の生成

1150 書式

```
1151 static int xpmem_make(unsigned long vaddr, size_t size, int permit_type, void
1152 *permit_value, xpmem_segid_t *segid_p)
```

1153 説明

1154 xpmem_make は、以下の処理を行う。

- 1155 1. 自プロセスの xpmem_thread_group を取得する。
- 1156 2. segid を算出する。
- 1157 3. xpmem_segment を生成して、初期設定する。
- 1158 4. segid_p に segid を設定する。

1159 戻り値

0	正常終了
-EINVAL	引数が無効である
-ENOMEM	十分な空きメモリ領域が無い
-XPMEM_ERRNO_NOPROC	対象プロセスの情報が無い

1160

1161 2.8.7 xpmem_segment の破棄

1162 書式

```
1163 static int xpmem_remove(xpmem_segid_t segid)
```

1164 説明

1165 xpmem_remove は、以下の処理を行う。

- 1166 1. segid から xpmem_thread_group(seg_tg) を取得する。
- 1167 2. seg_tg、segid から xpmem_segment を取得する。
- 1168 3. 取得した xpmem_segment を破棄する。

1169 戻り値

0	正常終了
-EACCES	許可がない
-EINVAL	引数が無効である

1170

1171 2.8.8 xpmem_access_permit の生成

1172 書式

```
1173 static int xpmem_get(xpmem_segid_t segid, int flags, int permit_type, void  
1174 *permit_value, xpmem_apid_t *apid_p)
```

1175 説明

1176 xpmem_get は、以下の処理を行う。

- 1177 1. segid から xpmem_thread_group(seg_tg) を取得する。
- 1178 2. seg_tg、sigid から xpmem_segment を取得する。
- 1179 3. 自プロセスの xpmem_thread_group(ap_tg) を取得する。
- 1180 4. ap_tg から apid を算出する。apid がマイナス値の場合にはエラー値を戻す。
- 1181 5. xpmem_access_permit を生成して、初期設定する。
- 1182 6. apid_p に apid を設定する。

1183 戻り値

0	正常終了
-EACCES	許可がない
-EINVAL	引数が無効である
-ENOMEM	十分な空きメモリ領域が無い
-XPMEM_ERRNO_NOPROC	対象プロセスの情報が無い

1184

1185 2.8.9 xpmem_access_permit の破棄

1186 書式

```
1187 static int xpmem_release(xpmem_apid_t apid)
```

1188 説明

1189 xpmem_release は、以下の処理を行う。

- 1190 1. apid から xpmem_thread_group(ap_tg) を取得する。
- 1191 2. ap_tg、apid から xpmem_access_permit を取得する。
- 1192 3. 取得した xpmem_access_permit を破棄する。

1193 戻り値

0	正常終了
-EACCES	許可がない
-EINVAL	引数が無効である

1194

1195 2.8.10 xpmem_attachment の生成

1196 書式

```
1197 static int xpmem_attach(struct mckfd *mckfd, xpmem_apid_t apid, off_t offset,  
1198 size_t size, unsigned long vaddr, int fd, int att_flags, unsigned long *at_vaddr_p)
```

1199

1200 説明

1201 xpmem_attach は、以下の処理を行う。

- 1202 1. apid から xpmem_thread_group(ap_tg) を取得する。
- 1203 2. ap_tg、apid から xpmem_access_permit(ap) を取得する。
- 1204 3. ap から xpmem_thread_group(seg_tg)、xpmem_segment(seg) を取得する。
- 1205 4. xpmem_attachment を生成して、初期設定する。

1206 5. do_mmap() を呼び出して、メモリ領域 (at_vaddr) を確保する。

1207 6. at_vaddr から vm_range(range) を取得する。

1208 7. range->private_data に xpmem_attachment を設定する。

1209 戻り値

0	正常終了
-EINVAL	引数が無効である
-ENOENT	そのようなファイルやディレクトリは無い
-ENOMEM	十分な空きメモリ領域が無い

1210

1211 2.8.11 xpmem_attachment の破棄

1212 書式

```
1213 static int xpmem_detach(unsigned long at_vaddr)
```

1214 説明

1215 xpmem_detach は、以下の処理を行う。

1216 1. at_vaddr から vm_range(range) を取得する。

1217 2. range->private_data から xpmem_attachment を取得する。

1218 3. xpmem_vm_munmap() を呼び出して、以下の処理を行う。

1219 (a) ihk_mc_clear_range() を呼び出して、メモリ領域を解放する。

1220 (b) range->memobj を解放する。

1221 (c) range を解放する。

1222 4. 取得した xpmem_attachment を破棄する。

1223 戻り値

0	正常終了
-EACCES	許可がない
-EINVAL	引数が無効である

1224

1225 2.8.12 vm_range の fault 処理

1226 書式

```
1227 int xpmem_fault_process_memory_range(struct process_vm *vm, struct vm_range  
1228 *vmr, unsigned long vaddr, uint64_t reason)
```

1229 説明

1230 xpmem_fault_process_memory_range は、以下の処理を行う。

- 1231 1. vmr->private_data から xpmem_attachment(att) を取得する。att が NULL の場合には
- 1232 エラー値 (-EFAULT) を返す。
- 1233 2. att から xpmem_access_permit(ap) を取得する。
- 1234 3. ap から xpmem_thread_group(ap_tg) を取得する。ap->flags または ap_tg->flags が XP-
- 1235 MEMi_FLAG_DESTROYING の場合にはエラー値 (-EFAULT) を返す。
- 1236 4. ap から xpmem_segment(seg) を取得する。
- 1237 5. seg から xpmem_thread_group(seg_tg) を取得する。seg->flags または seg_tg->flags が
- 1238 XPMEM_FLAG_DESTROYING の場合にはエラー値 (-ENOENT) を返す。
- 1239 6. xpmem_remap_pte() を呼び出して、以下の処理を行う。
- 1240 (a) ihk_mc_pt_lookup_pte() を呼び出して、seg の vaddr から pte_t(seg_pte) を取得する。
- 1241 (b) ihk_mc_pt_lookup_pte() を呼び出して、vaddr から pte_t(att_pte) を取得する。
- 1242 (c) ihk_mc_pt_set_pte() を呼び出して、att_pte の物理アドレスを seg_pte の物理アド
- 1243 レスに置き換える。

1244 戻り値

0	正常終了
-EFAULT	アドレスが不正である
-ENOENT	そのようなファイルやディレクトリは無い

1245

1246 2.8.13 vm_range の削除

1247 書式

```
1248 int xpmem_remove_process_memory_range(struct process_vm *vm, struct vm_range
1249 *vmr)
```

1250 説明

1251 xpmem_remove_process_memory_range は、以下の処理を行う。

- 1252 1. vmr->private_data から xpmem_attachment(att) を取得する。
- 1253 2. att が指定されていた場合には、以下の処理を行う。
- 1254 (a) att を解放する。
- 1255 (b) vmr->private_data に NULL を設定する。

1256 戻り値

1257

1258 2.9 ライブラリ切り替え

1259 McKernel は、特定のパスについて、McKernel 上に起動されたプロセスと Linux 上に起動さ
 1260 れたプロセスとに対して異なるファイルを見せる機能を提供する。これは、McKernel での実
 1261 行と Linux での実行とで異なるライブラリファイルをリンクせねばならない例外的なケース
 1262 (例えば、第??節で説明する Utility Thread Offloading のライブラリ) で、ローダ/リンカに
 1263 異なるファイルをリンクさせることを目的とする。

1264 動作は以下の通り。IHK/McKernel のインストールディレクトリを<install>とする。

- 1265 1. unshare コマンドを用いて mcexec の mount name space の設定を変更し、mcexec が
 1266 mcctrl に依頼する bind mount が他プロセスからは見えないようにする。
- 1267 2. mcexec が mcctrl に制御を移す。
- 1268 3. mcctrl がユーザ id を root に変更し、<install>/rootfs/以下のファイルのそれぞれ
 1269 を/に bind mount する。
- 1270 4. mcctrl がユーザ id を元に戻し mcexec に制御を戻す。

1271 2.10 状態監視

1272 McKernel のハングアップ検知は以下のステップで実施される。

- 1273 1. 運用ソフトウェアが IHK の関数を用いて通知のための eventfd を取得する。
- 1274 2. McKernel が CPU ごとの状態と状態遷移回数を記録する。
- 1275 3. Linux 上で動作するスレッド (ihkmond) が上記の状態を監視し、2度同じ状態にあった
 1276 場合、ハングアップと判断し、上記 eventfd を用いて運用ソフトウェアに通知する。

1277 監視スレッドとハングアップ通知のインターフェイスは”IHK Specifications“に記載する。本
 1278 節では第2のステップを説明する。

1279 状態と状態遷移回数の記録には struct ihk_os_monitor 型の変数を用いる。以下の説明
 1280 ではこの型を持つ監視用の変数を monitor と呼ぶ。struct ihk_os_monitor の関連部分は以
 1281 下のように定義される。

```
1282 struct ihk_os_monitor {
1283     ...
1284     int status;                /* OS 状態 */
1285     unsigned long counter; /* OS 状態が変化した回数 */
1286 };
```

1287 状態と状態遷移回数の記録の動作は以下の通り。

- 1288 1. McKernel が以下のようにイベントに応じて状態と状態遷移回数を更新する。
 - 1289 ● カーネルモードからユーザモードへの移行時: monitor.status を IHK_OS_MONITOR_
 1290 USER に設定する。

- 1291 ● ユーザモードからカーネルモードへの移行時: `monitor.status` を `IHK_OS_MONITOR_KERNEL` に設定し、`monitor.counter` をインクリメントする。
- 1292
- 1293 ● システムコール移譲時: 移譲開始直前に `monitor.status` の値を保存し、`IHK_OS_MONITOR_KERNEL_OFFLOAD` に設定する。また、移譲完了後に `monitor.status` の値を保存しておいた値に戻し、`monitor.counter` をインクリメントする。
- 1294
- 1295
- 1296 ● `rt_sigtimedwait()`、`do_sigsuspend()`、`futex()`、`nanosleep()` 呼び出し時: 関数に入った直後に `monitor.status` を `IHK_OS_MONITOR_KERNEL_HEAVY` に設定する。なお、この状態に長時間滞在してもハングアップとは判定しない。
- 1297
- 1298
- 1299 ● `idle()` 呼び出し時: 関数に入った直後に `monitor.status` を `IHK_OS_MONITOR_IDLE` に設定する。なお、この状態に長時間滞在してもハングアップとは判定しない。その後、`cpu_safe_halt()` から復帰したタイミングで `monitor.status` を `IHK_OS_MONITOR_KERNEL` に設定し、`monitor.counter` をインクリメントする。
- 1300
- 1301
- 1302

1303 2.11 Non-Maskable Interrupt

1304 Non-Maskable Interrupt (NMI) は対象 CPU に以下の動作をさせるために用いられる。

- 1305 ● カーネルダンプの準備
- 1306 ● 一時停止状態への遷移
- 1307 ● 一時停止状態からの復帰

1308 なお、カーネルダンプの準備については第??節に、一時停止状態への遷移及びそこから
1309 の復帰については第??節に記載する。

1310 NMI の利用ステップは以下の通り。

- 1311 1. McKernel がブート時に `ihk_set_nmi_mode_addr()` で NMI の動作を指定する McKernel
1312 の変数 `nmi_mode` の物理アドレスを IHK-master に伝える。
- 1313 2. IHK-master driver が `nmi_mode` の値を上記の動作のいずれかを示す値に設定し、`smp_`
1314 `ihk_os_send_nmi()` を呼び、各 CPU に NMI を送る。
- 1315 3. 各 CPU が以下を実行する。
 - 1316 (a) NMI を受けて、NMI ハンドラ `nmi()` に制御を移す。
 - 1317 (b) `nmi()` で `nmi_mode` の値に応じた処理を行う。

1318 以下、関連関数の動作を説明する。

1319 2.11.1 NMI 動作設定

1320 書式

```
1321 int ihk_set_nmi_mode_addr(unsigned long addr)
```

1322 説明

1323 `addr` で指定される物理アドレスを NMI の動作を規定する McKernel の変数 `nmi_mode` の
1324 物理アドレスとして IHK に登録する。こうすることで、IHK から McKernel の NMI ハンド
1325 ラの動作を切り替えることができるようになる。`nmi_mode` の値と NMI ハンドラの動作の対
1326 応は以下の通り。

値	動作
0	NMI ハンドラで各 CPU のカーネルダンプの準備を行う。
1	NMI ハンドラで各 CPU の状態を一時停止状態へ遷移させる。
2	NMI ハンドラで各 CPU の状態を一時停止状態から復帰させる。

1327 戻り値

0	正常終了
---	------

1328 2.11.2 NMI 送信

1329 書式

```
1330 static int smp_ikh_os_send_nmi(ikh_os_t ikh_os, void *priv, int mode)
```

1331 説明

1332 nmi_mode を mode に設定した上で各 CPU に NMI を発行する。

1333 戻り値

0	正常終了
-EINVAL	エラー

1334

1335 2.11.3 NMI ハンドラ

1336 書式

```
1337 void nmi()
```

1338 説明

1339 nmi_mode に指定された値に従った動作を行う。nmi_mode の値と動作の対応は第??に示す。

1340 2.12 全 CPU 一時停止

1341 McKernel は全 CPU を FROZEN と呼ぶ一時停止状態に遷移させる機能および FROZEN から
 1342 復帰させる機能を提供する。この機能と全 CPU を低電力状態に遷移させる機能とを組み合
 1343 わせることで、ジョブ単位での低電力状態への遷移とそこからの復帰を実現する。

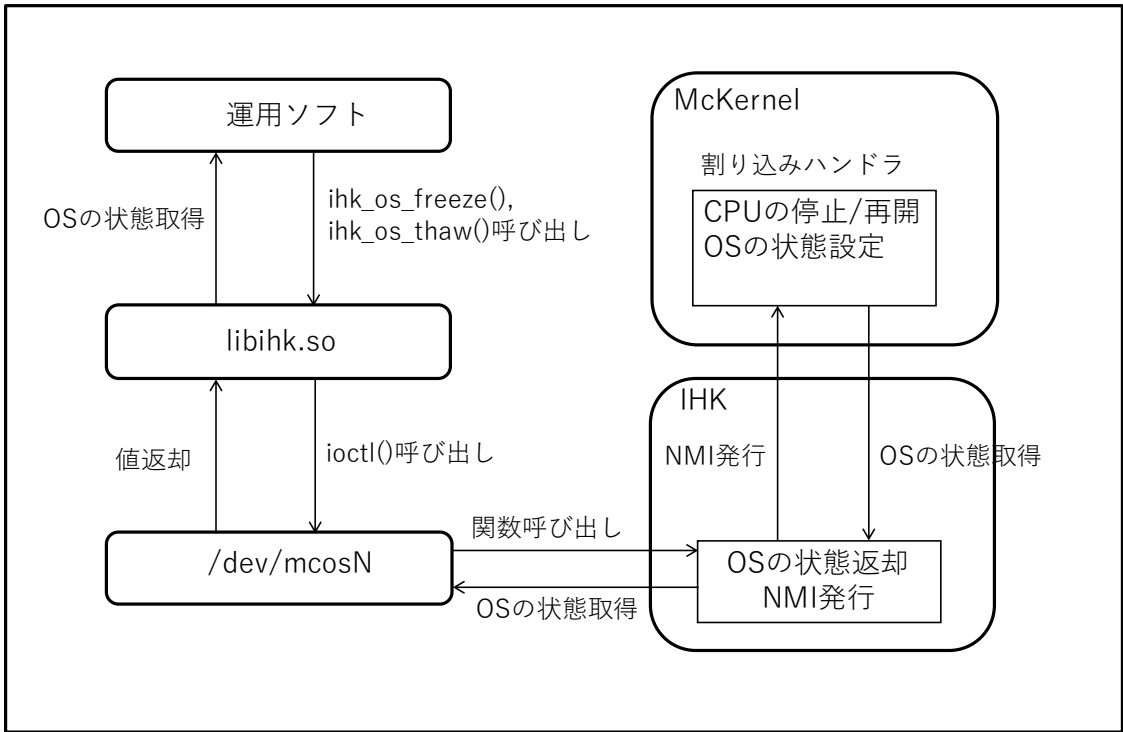


Figure 2.13: 構成要素関連図

1344

全 CPU 一時停止機能の構成を図?? に示す。

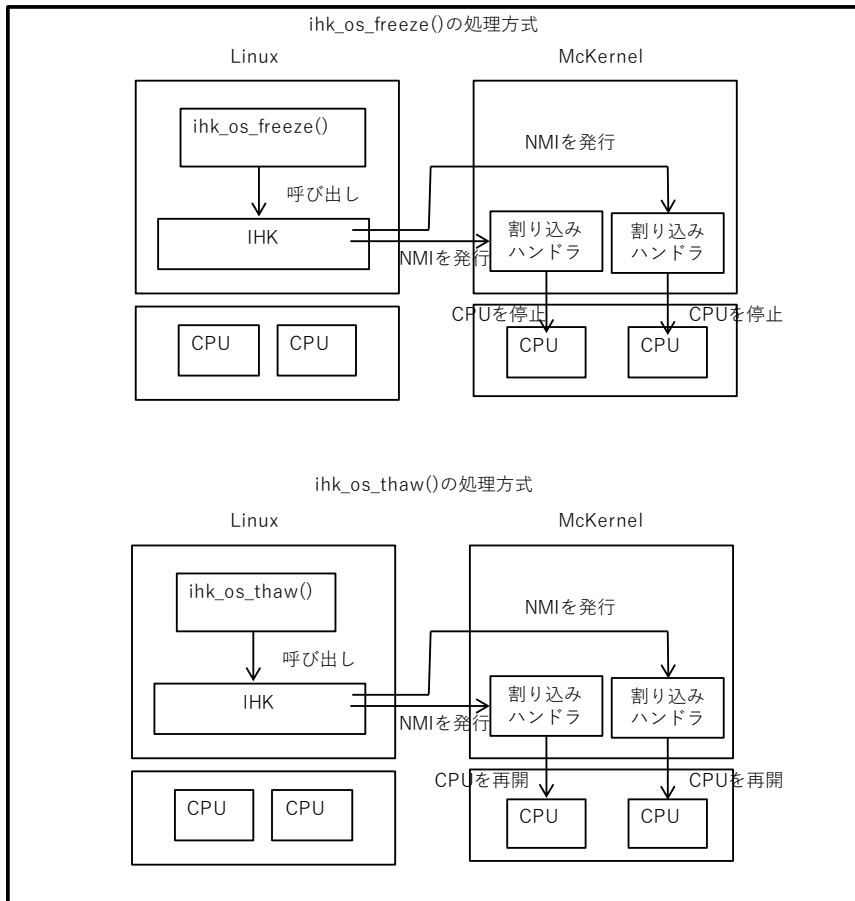


Figure 2.14: 全 CPU 一時停止および一時停止からの復帰のフロー

- 1345 全 CPU 一時停止の動作を図??を用いて説明する。
- 1346 1. バッチジョブスケジューラが `ihk_os_freeze()` 経由で `IHK_OS_FREEZE` コマンドを指定
- 1347 して `ioctl()` を呼ぶ。
- 1348 2. IHK-master core が `_ihk_os_freeze()` 経由で IHK-master driver の `smp_ihk_os_freeze()`
- 1349 を呼ぶ。
- 1350 3. IHK-master driver が `nmi_mode` に一時停止状態への遷移を示す値を設定し、`smp_ihk_`
- 1351 `os_send_nmi()` を呼び、各 CPU に NMI を送る。
- 1352 4. 各 CPU が以下を実行する。
- 1353 (a) NMI を受けて、NMI ハンドラ `nmi()` に制御を移す。
- 1354 (b) `nmi()` で `nmi_mode` に設定された指示に従った処理を行う。この場合は一時停止状
- 1355 態への遷移であるため、`freeze_thaw()` を呼ぶ。
- 1356 (c) `freeze_thaw()` で `nmi_mode` に設定された指示に従った処理を行う。この場合は
- 1357 一時停止状態への遷移であるため、`mod_nmi_ctx()` を用いて `iret` 命令後のジャンプ
- 1358 先を `_freeze()` にする。

1359 (d) `__freeze()` は `freeze()` を呼び出す。`freeze()` は以下を実行する。

1360 i. CPU の状態をバックアップ用変数に保持する。

1361 ii. CPU の状態を一時停止状態に設定する。

1362 iii. CPU を停止させる。x86_64 アーキでは `hlt` 命令を実行する。

1363 一時停止からの復帰の動作を図??を用いて説明する。

1364 1. バッチジョブスケジューラが `ihk_os_thaw()` 経由で `IHK_OS_THAW` コマンドを指定して
1365 `ioctl()` を呼ぶ。

1366 2. IHK-master core が `__ihk_os_thaw()` 経由で IHK-master driver の `smp_ihk_os_thaw()`
1367 を呼ぶ。

1368 3. IHK-master driver が `nmi_mode` に一時停止状態からの復帰を示す値を設定し、`smp_ihk_`
1369 `os_send_nmi()` を呼び、各 CPU に NMI を送る。

1370 4. 各 CPU が以下を実行する。

1371 (a) NMI を受けて、NMI ハンドラ `nmi()` に制御を移す。

1372 (b) `nmi()` で `nmi_mode` に設定された指示に従った処理を行う。この場合は一時停止状
1373 態からの復帰であるため、`freeze_thaw()` を呼ぶ。

1374 (c) `freeze_thaw()` で `nmi_mode` に設定された指示に従った処理を行う。この場合は
1375 一時停止状態からの復帰であるため、CPU の状態をバックアップ用変数を用いて復
1376 元する。

1377 以下、関連関数のインターフェイスと動作を説明する。

1378 2.12.1 一時停止指示 (IHK-master core)

1379 書式

```
1380 static int __ihk_os_freeze(struct ihk_host_linux_os_data *data)
```

1381 説明

1382 アーキ依存の一時停止指示関数を呼ぶ。smp-x86 では `smp_ihk_os_freeze()` を呼び出す。

1383 戻り値

0	正常終了
---	------

1384

1385 2.12.2 一時停止からの復帰指示 (IHK-master core)

1386 書式

```
1387 static int __ihk_os_thaw(struct ihk_host_linux_os_data *data)
```

1388 説明

1389 アーキ依存の一時停止からの復帰指示関数を呼ぶ。smp-x86 では smp_ihk_os_thaw() を呼
1390 び出す。

1391 戻り値

0	正常終了
---	------

1392

1393 2.12.3 一時停止指示 (IHK-master driver)

1394 書式

```
1395 static int smp_ihk_os_freeze(ihk_os_t ihk_os, void *priv)
```

1396 説明

1397 smp_ihk_os_send_nmi() を呼び出して各 CPU に NMI を送り、CPU の状態を一時停止状
1398 態へ遷移させ、また CPU を NMI を受けるまで停止させる。

1399 戻り値

0	正常終了
---	------

1400

1401 2.12.4 一時停止からの復帰指示 (IHK-master driver)

1402 書式

```
1403 static int smp_ihk_os_thaw(ihk_os_t ihk_os, void *priv)
```

1404 説明

1405 smp_ihk_os_send_nmi() を呼び出して各 CPU に NMI をを送り、NMI 待ちで停止している
1406 CPU の処理を再開させ、また CPU の状態を元の状態に戻す。

1407 戻り値

0	正常終了
---	------

1408

1409 2.12.5 一時停止および一時停止からの復帰指示

1410 書式

```
1411 long freeze_thaw(void *nmi_ctx)
```

1412 説明

1413

- 1414 1. 変数 `nmi_mode` が一時停止状態への遷移を意味する場合、`mod_nmi_ctx()` を呼び出すこと
1415 で `__freeze()` を呼び出し、CPU の状態を一時停止状態に遷移させ、また CPU を NMI
1416 を受けるまで停止させる。
- 1417 2. 変数 `nmi_mode` が一時停止状態からの復帰を意味する場合、CPU の状態を一時停止前の
1418 状態に戻す。

1419 戻り値

0	一時停止を行った
1	一時停止からの復帰を行った

1420

1421 2.12.6 NMI ハンドラからの復帰時の指定関数へのジャンプ設定

1422 書式

```
1423 void mod_nmi_ctx(void *nmi_ctx, void (*func)())
```

1424 説明

1425 NMI ハンドラからの復帰時 (x86_64 アーキテクチャでは `iret` 命令実行時) に割り込み発
1426 生命令に戻らず、`func` で指定した、NMI 受け付けが必要な関数にジャンプするようにスタッ
1427 クの内容を変更する。このような処理が必要なのは、NMI ハンドラ内では NMI を受け付け
1428 ないためである。`func` に `__freeze()` を指定することで、CPU を NMI 待ちの状態に停止さ
1429 せることができる。

1430 2.12.7 一時停止指示 (ラッパー)

1431 書式

```
1432 void __freeze()
```

1433 説明

1434 `freeze()` を呼び出して CPU を一時停止させ、その後割り込みハンドラから復帰する。
1435 x86_64 アーキでは割り込みハンドラからの復帰には `iret` 命令を用いる。

1436 2.12.8 一時停止指示

1437 書式

```
1438 void freeze()
```


1439 説明

1440 ステップは以下の通り。

- 1441 1. CPU 状態を保存する。
- 1442 2. CPU 状態を `IHK_OS_MONITOR_KERNEL_FROZEN` に遷移させる。
- 1443 3. `cpu_halt()` を呼び CPU を停止させる。なお、CPU は NMI を受けると処理を再開する。
- 1444 4. CPU が処理を再開した後、CPU 状態を保存しておいた値に戻す。

1445 2.13 カーネルダンプ

1446 カーネルダンプの採取と解析のステップは以下の通り。

- 1447 1. 以下のいずれかの方法でダンプファイルを作成する。
 - 1448 (a) IHK の関数 `ihk_os_makedumpfile()` または IHK のコマンド `ihkosctl` を用いて、
 - 1449 `McKernel` 形式のダンプファイルを作成する（以降、`McKernel` 主導ダンプと呼ぶ）。
 - 1450 (b) Linux の `panic` を契機に `makedumpfile` 形式のダンプファイルを作成する。また、
 - 1451 コマンド `vmcore2mckdump` を用いて `McKernel` 形式に変換する（以降、Linux 主導
 - 1452 ダンプと呼ぶ）。
- 1453 2. `eclair` と呼ぶコマンドを用いてダンプファイルを解析する。

1454 以下、詳細を説明する。

1455 2.13.1 全体の処理の流れ

1456 `McKernel` 主導ダンプの場合のダンプ採取機能とダンプ形式変換機能の処理の流れを図??を
1457 用いて説明する。

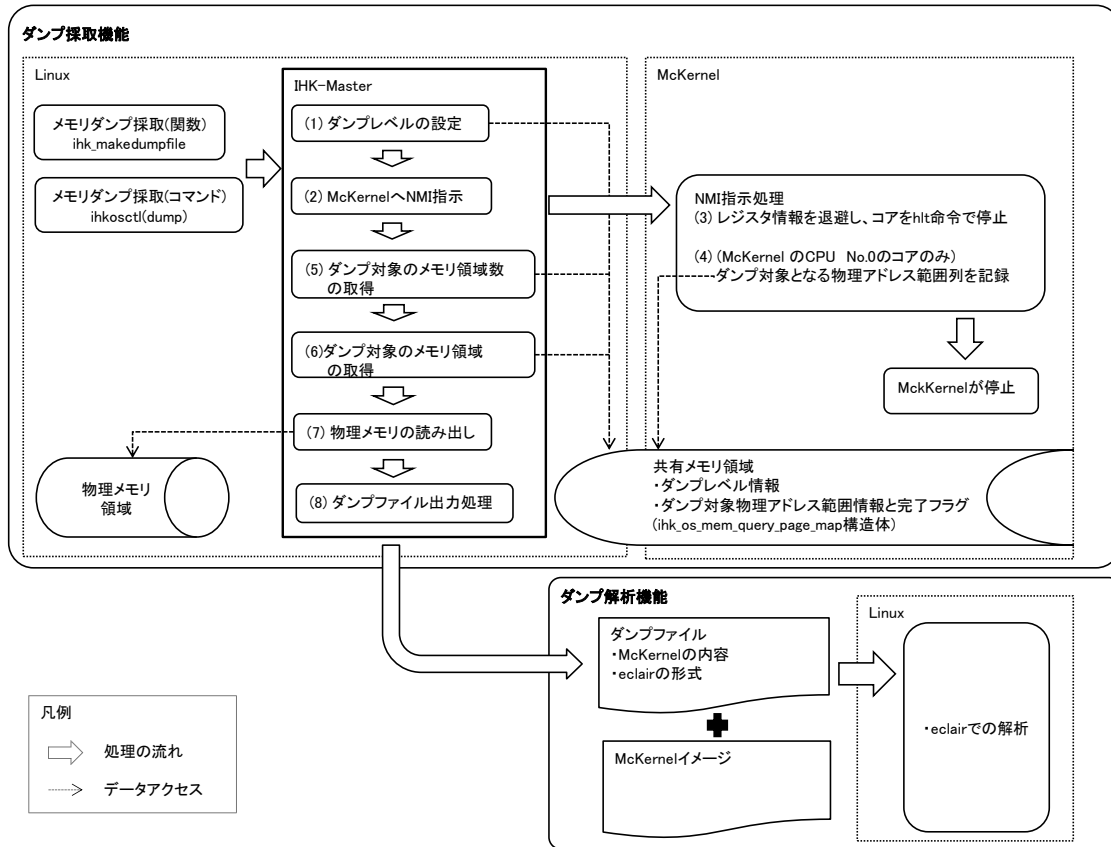


Figure 2.15: McKernel 主導ダンプの場合のダンプ採取機能とダンプ形式変換機能の処理の流れ

- 1458 1. IHK が OS ブート時に、各物理ページがダンプ対象であることを示す情報（以下、ダンプ
 1459 対象ページリストと呼ぶ）を Linux と McKernel とで共有しているメモリ領域（以降、
 1460 共有メモリと呼ぶ）に確保する。また、IHK が McKernel に割り当てた物理アドレス範囲
 1461 をダンプ対象とするように初期化する。
- 1462 2. 管理者が `ihk_os_makedumpfile()` でダンプを指示する。
- 1463 3. IHK が共有メモリにダンプレベルを記録する。また、共有メモリ上の、ダンプ対象ペー
 1464 ジリストの設定完了を表すフラグ（以降、完了フラグと呼ぶ）を 0 に設定する。（図
 1465 の (1)）
- 1466 4. IHK が McKernel の各コアへ NMI を送る。（図の (2)）
- 1467 5. McKernel の第 0 CPU 以外の CPU はレジスタ情報を退避した後 `hlt` 命令で停止する
 1468 （図の (3)）。McKernel の第 0 CPU は以下を実行する。（図の (3)、(4)）
- 1469 (a) レジスタ情報を退避する。
- 1470 (b) 共有メモリを参照してダンプレベルを取得する。
- 1471 (c) ダンプからユーザ領域を除外する指定がされている場合は、ユーザメモリ領域情
 1472 報を取得し、ダンプ対象ページリストの対応ビットを 0 にする。

- 1473 (d) ダンプから未使用領域を除外する指定がされている場合は、未使用メモリ領域情
1474 報を取得し、ダンプ対象ページリストの対応ビットを0にする。
- 1475 (e) 完了フラグに1をセットする。
- 1476 (f) hlt 命令で停止する。

1477 6. IHK が完了フラグが1になるまで待ち、`ioctl()` でダンプ対象のメモリ領域数を取得
1478 し、領域情報を格納するメモリ領域を確保し、さらに `ioctl()` でダンプ対象の領域情
1479 報を前記メモリ領域に記録する。(図の (5)、(6))

1480 7. IHK がダンプ対象のメモリ領域を `ioctl()` で読み出し、ファイルに書き込む。(図の
1481 (7)、(8))

1482 8. 管理者は `eclair` を用いてダンプファイルの解析を行う。

1483 ダンプ対象は `ihk_dump_page_set` で表現する。定義は以下の通り。

```
1484 struct ihk_dump_page_set {
1485     unsigned int completion_flag; /* 書き込み完了フラグ */
1486     unsigned int count;          /* ダンプ対象のページ情報数 */
1487     unsigned long page_size;     /* ダンプ対象のページ情報の全体サイズ */
1488     unsigned long phy_page;     /* ダンプ対象のページ情報の物理アドレス
1489                                 (struct ihk_dump_page の配列) */
1490 }
1491
1492 struct ihk_dump_page {
1493     unsigned long start;        /* マップ情報の開始物理アドレス */
1494     unsigned long map_count;    /* マップ情報の領域数 (map[] の配列数) */
1495     unsigned long map[];       /* マップ情報 (ビットマップ形式) */
1496 };
```

1497 Linux 主導ダンプの場合のダンプ採取機能とダンプ形式変換機能の処理の流れを図??を
1498 用いて説明する。

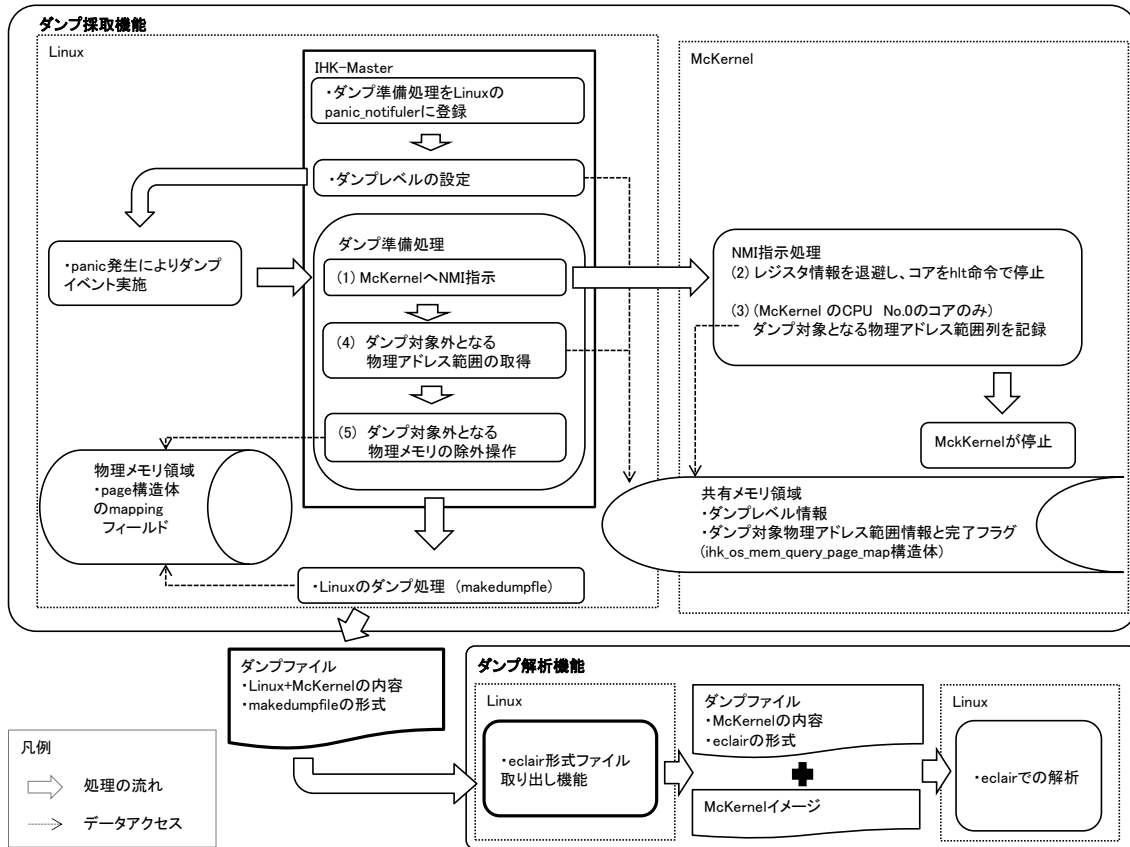


Figure 2.16: Linux 主導ダンプの場合のダンプ採取機能とダンプ形式変換機能の処理の流れ

- 1499 1. IHK が OS ブート時に、ダンプ対象ページリストを共有メモリに確保する。また、IHK
1500 が McKernel に割り当てた物理アドレス範囲をダンプ対象とするように初期化する。
- 1501 2. IHK が McKernel 起動時にダンプレベル設定オプション (-d) でダンプレベルを指定す
1502 る。IHK は共有メモリにこのダンプレベルを記録する。また、完了フラグを 0 に設定
1503 する。
- 1504 3. IHK がダンプ準備処理関数を Linux の panic_notifier に登録する。
- 1505 4. Linux で panic が発生し、登録されているダンプ準備処理関数が呼び出される。
- 1506 5. IHK が McKernel の各コアへ NMI を送る。(図の (1))
- 1507 6. McKernel の第 0 CPU 以外の CPU は、レジスタ情報を退避し、halt 命令で停止する (図
1508 の (3))。McKernel の第 0 CPU は以下を実行する。(図の (2)、(3))
- 1509 (a) レジスタ情報を退避する。
- 1510 (b) 共有メモリを参照してダンプレベルを取得する。
- 1511 (c) ダンプからユーザ領域を除外する指定がされている場合は、ユーザメモリ領域情
1512 報を取得し、ダンプ対象ページリストにダンプからの除外を記録する。
- 1513 (d) ダンプから未使用領域を除外する指定がされている場合は、未使用メモリ領域情
1514 報を取得し、ダンプ対象ページリストのダンプからの除外を記録する。

- 1515 (e) 完了フラグに1をセットする。
 1516 (f) hlt 命令で停止する。
- 1517 7. IHK が完了フラグが1になるまで待ち、`ioctl()` でダンプ対象外の物理アドレス範囲
 1518 に該当する Linux の page 構造体の `mapping` フィールドを操作し `anonymous` に設定す
 1519 る。(図の (4)、(5))
- 1520 8. Linux が `makedumpfile` コマンドを実行する。
- 1521 9. Linux が Linux と McKernel の両方の情報を含むダンプファイルを作成する。
- 1522 10. 管理者が `ldump2mcdump` コマンドで、`makedumpfile` 形式のダンプファイルを `eclair` 形
 1523 式に変換する。
- 1524 11. 管理者は `eclair` を用いてダンプファイルの解析を行う。

1525 2.13.2 ユーザメモリ領域情報取得

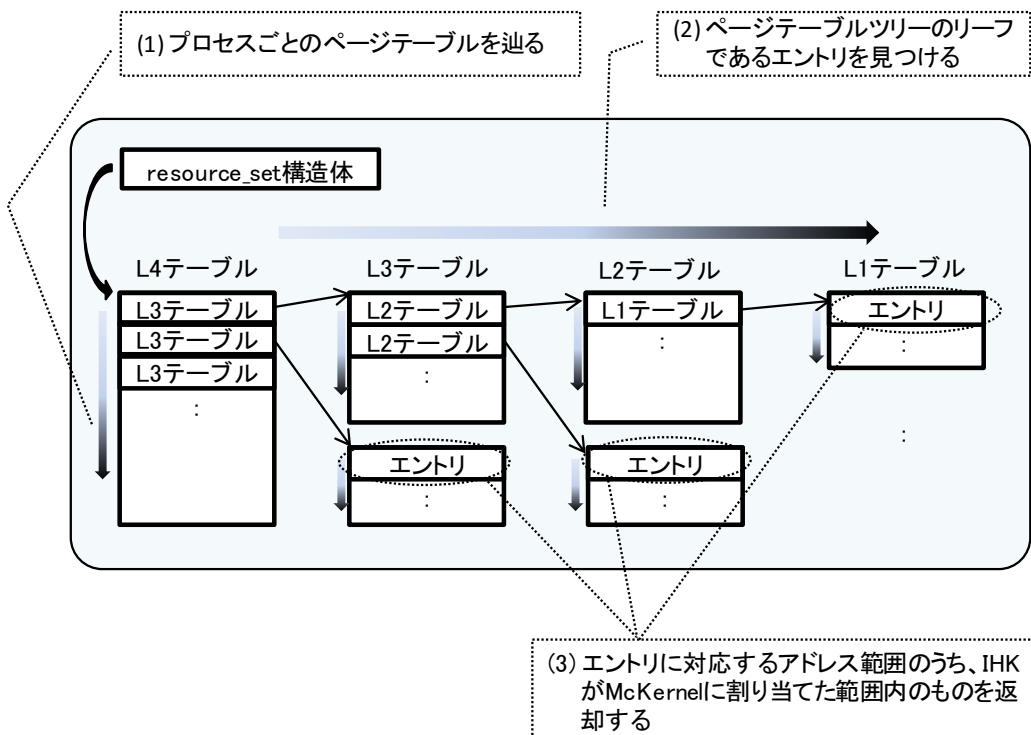


Figure 2.17: ユーザメモリ領域情報取得処理の流れ

- 1526 ユーザメモリ領域情報取得処理の流れを図??を用いて説明する。
- 1527 1. `resource_set` を参照して全プロセスを辿り、プロセスごとのページテーブルについて
 1528 以下を行う。(図の (1))
- 1529 (a) ページテーブルツリーのリーフであるエントリを見つける。(図の (2)) なお、4
 1530 段目のエントリは 4 KB ページのエントリ、3 段目かつ `PageSize` フラグが1のエ
 1531 ントリは 2 MB、2 段めかつ `PageSize` フラグが1のエントリは 1 GB ページのエ
 1532 ントリである。

1533 (b) エントリに対応するアドレス範囲のうち、IHK が McKernel に割り当てた範囲に
 1534 収まるものをユーザメモリ領域として返却する。収まらないものはエントリが破
 1535 壊されているとみなし破棄する。(図の (3))

1536 2.13.3 未使用メモリ領域情報取得

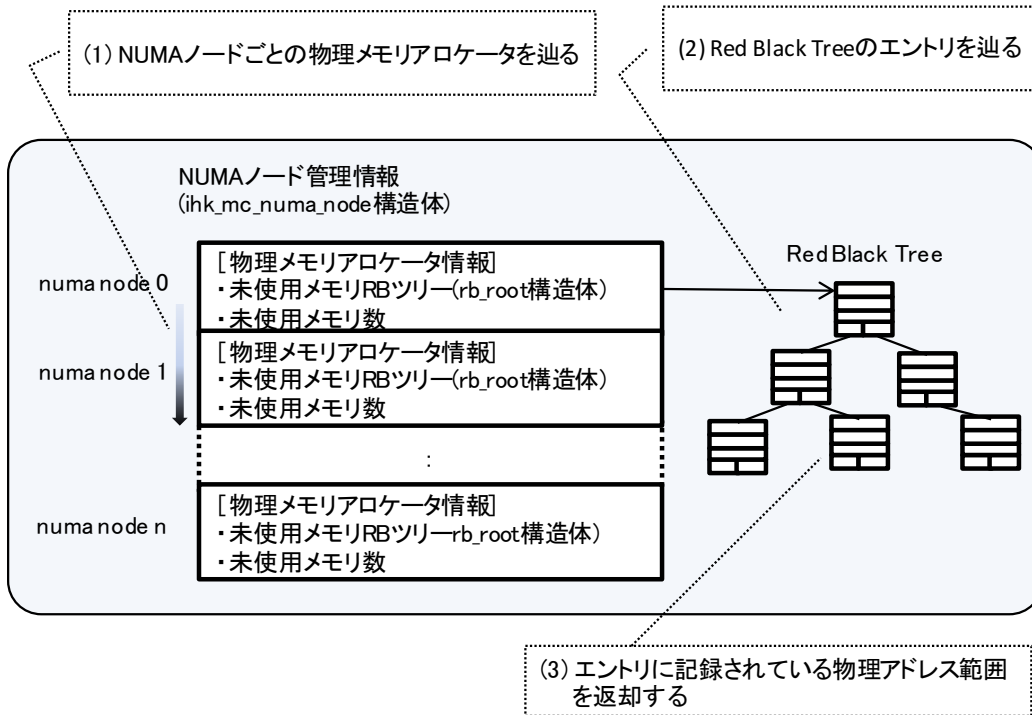


Figure 2.18: 未使用メモリ領域情報取得処理の流れ

1537 未使用メモリ領域情報取得の処理の流れを図??を用いて説明する。

1538 1. NUMA ノード管理情報 (ihk_mc_numa_node 構造体) を参照して NUMA ノードごとの物
 1539 理メモリアロケータについて以下を行う。(図の (1))

1540 (a) 未使用メモリを管理する Red Black tree (rb_root 構造体) のエントリを辿る。(図
 1541 の (2))

1542 (b) エントリ (free_chunk 構造体) に記録されている物理アドレス範囲を返却する。(図
 1543 の (3))

1544 2.13.4 ダンプ処理用 ioctl() コマンド

1545 書式

```
1546 int ioctl(int fd, IHK_OS_DUMP, struct ihk_dump_args *args)
```

1547 説明

1548 fd で指定された OS インスタンスに対して、args->cmd に指定されたダンプ関連処理を
 1549 行う。

1550 dumpargs_t は以下のように定義される。

```

1551 struct ihk_dump_args {
1552     int cmd; /* コマンド */
1553     unsigned int level; /* ダンプレベル */
1554     long start; /* 開始物理アドレス */
1555     long size; /* サイズ */
1556     void *buf; /* メモリ内容 */
1557     int num_mem_chunks; /* メモリ領域数 */
1558     struct ihk_dump_mem_chunk *mem_chunks; /* メモリ領域情報 */
1559 };

```

1560 struct ihk_dump_mem_chunk は以下のように定義される。

```

1561 struct ihk_dump_mem_chunk {
1562     unsigned long addr;
1563     unsigned long size;
1564 };

```

1565 args->cmd ごとの処理は以下の通り。

args->cmd	動作				
DUMP_QUERY_NUM_MEM_AREAS	ダンプ対象メモリ領域数を返す。				
DUMP_QUERY_MEM_AREAS	ダンプ対象メモリ領域の情報を args->mem_chunks に格納する。呼び出し元が args->mem_chunks の領域を用意する。				
DUMP_READ	args->start, args->size で指定された物理メモリ領域の内容を args->buf で指定されたバッファにコピーする。				
DUMP_SET_LEVEL	<p>ダンプ対象とするメモリ領域の種類を args->level に設定する。設定可能な値は以下の通り。</p> <table border="1"> <tr> <td>0</td> <td>IHK が McKernel に割り当てたメモリ領域を出力する。</td> </tr> <tr> <td>24</td> <td>カーネルが使用しているメモリ領域を出力する。</td> </tr> </table> <p>なお、args->level が設定可能でない値であった場合は-EINVAL を返却する。</p>	0	IHK が McKernel に割り当てたメモリ領域を出力する。	24	カーネルが使用しているメモリ領域を出力する。
0	IHK が McKernel に割り当てたメモリ領域を出力する。				
24	カーネルが使用しているメモリ領域を出力する。				
DUMP_NMI	全 CPU に NMI を発行し、ダンプの準備を指示する。				
DUMP_SET_ANONYMOUS	(IHK が McKernel に割り当てたメモリ領域) から (args->mem_chunks, args->num_mem_chunks で指定したメモリ領域) を除いた領域に対し、Linux の struct page の mapping フィールドの最下位ビットをセットし anonymous テーブルに見せかける。こうすることで、Linux の makedumpfile が該当領域をダンプ対象から除外できるようになる。				
DUMP_QUERY	IHK によって割り当てられた物理メモリ領域の情報を args->start, args->size に格納する。本機能は、IHK が McKernel に割り当てたメモリ領域の全てをダンプする際に使用する。				

1566 戻り値

0	正常終了
-EFAULT	アドレスが不正である
-EINVAL	引数が無効である

1567 2.13.4.1 ダンプファイルの形式

1568 ダンプファイルは ELF 形式を採用している。ダンプファイルで使用しているセクションは以下
1569 の通り。

セクション名	説明
Date	ダンプ採取日時 例： Thu Mar 3 21:42:35 2016
hostname	ダンプ採取ホスト名 例： kncc08
User	ダンプ採取 実ユーザ名 例： nakamura
physmem	物理メモリダンプ

1570 なお、レジスタの値はダンプファイルには格納しない。その代わりに、スレッドを表現す
1571 る構造体に格納されている退避コンテキストから値を取得する。スレッドを表現する構造体
1572 の位置は、まず各コアの run queue の位置をシンボル情報から取得し、そこに挿入されてい
 るエントリを見つけることで取得する。objdump での出力例を図??に示す。

eclair形式のファイルフォーマット

```

mcdump_20160303_214235:   file format elf64-little

Sections:
Idx Name          Size      VMA           LMA           File off  Algn
  0 date           00000018  0000000000000000  0000000000000000  00000040  2**0
   CONTENTS
  1 hostname       00000006  0000000000000000  0000000000000000  00000058  2**0
   CONTENTS
  2 user           00000008  0000000000000000  0000000000000000  0000005e  2**0
   CONTENTS
  3 physmem       20000000  0000000771800000  0000000771800000  00000066  2**0
   CONTENTS、 ALLOC、 LOAD、 DATA

Contents of section date:
0000 54687520 4d617220 20332032 313a3432  Thu Mar  3 21:42
0010 3a333520 32303136                :35 2016

Contents of section hostname:
0000 6b6e6363 3038                kncc08

Contents of section user:
0000 6e616b61 6d757261                username

Contents of section physmem:
77180000 00000000 00f0ffff 00000000 00400000  .....@..
77180010 21000000 0c000000 48004800 00000000  !.....H.H....
77180020 ffffffff ffffffff ffffffff ffffffff  .....
(以下略)

```

Figure 2.19: ダンプファイルの objdump での出力例

1573

1574 2.13.5 ダンプ解析コマンドと gdb コマンドとの連携方法

1575 ダンプ解析コマンドと gdb コマンドとの間の remote serial protocol は、IPv4 を使った TCP
1576 通信でやり取りする。unix ドメインソケットなども利用可能とは思いますが、異常終了時にごみ
1577 ファイルが残ることを回避するために TCP/IP 通信を選択した。ユーザが直接 gdb を起動し
1578 て毎回リモートデバッグの設定を行うことは、難しくはないが面倒である。そこで、ユーザ
1579 には、gdb ではなくダンプ解析コマンドを起動してもらおう。ダンプ解析コマンドが gdb コマ

1580 ンドの起動とリモートデバッグの設定を行う。ダンプ解析コマンドによるリモートデバッグ
1581 の設定から、実際にユーザからの解析コマンドを受け取る gdb に、スムーズに端末を受け渡
1582 すため、以下の手順で動作する。

1583 1. ユーザから起動されたダンプ解析コマンドは、コマンドラインオプションを解析して
1584 gdb エージェントとしての初期化をする。

1585 2. ダンプ解析コマンドは、remote serial protocol 通信用の TCP ソケットを作成する。

1586 3. ダンプ解析コマンドは、gdb を fork() と exec() で起動する。この時、以下のようなコマ
1587 ンドライン引数としてリモートデバッグの設定に必要なコマンドを与える。
1588

```
1589 -q -ex set prompt (eclair) -ex target remote :< TCP ポート番号> <カーネルイメージファイ  
1590 ル名>
```

1591 4. ダンプ解析コマンドは、TCP ソケットに gdb が接続してくるのを待つ。

1592 5. ダンプ解析コマンドは、TCP ソケット接続後、端末からの入力をせずに gdb エージェ
1593 ントとしての動作に専念する。

1594 上記の手順によって、ダンプ解析コマンドと gdb とが同じ端末を共有した状態になる。共有
1595 していても、標準入力の読み出しをダンプ解析コマンドが一切実行しなければ、gdb が標準
1596 入力を占有しているのと同じ動作をさせることができる。

1597 2.13.6 ダンプ形式変換 (crash プラグイン)

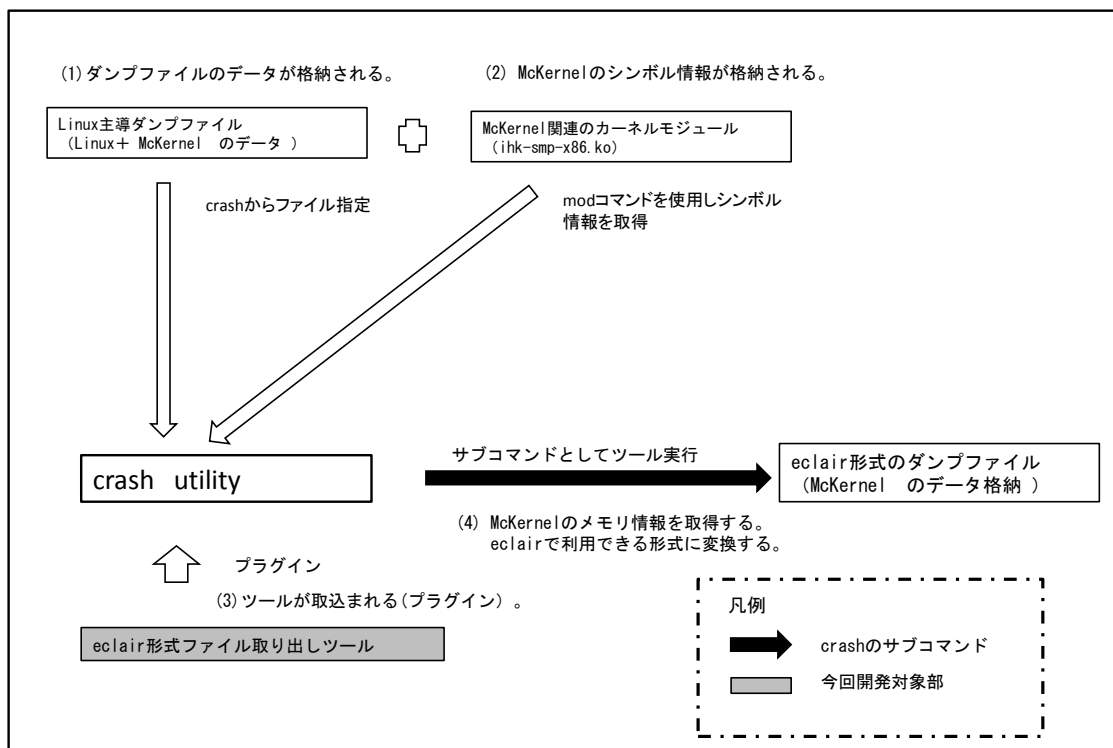


Figure 2.20: ダンプ形式変換処理の流れ

1598 ダンプ形式変換処理の流れを図??を用いて説明する。

- 1599 (1) crash 内領域にダンプファイルのデータが格納される。
 1600 コマンド : crash <vmlinux のパス> <ダンプファイル (vmcore) のパス>
- 1601 (2) crash 内領域に IHK および McKernel のカーネルモジュールのシンボル情報が格納される。
 1602
 1603 コマンド : mod -s ihk-smp-x86 <ihk-smp-x86.o のパス>
 1604 シンボル情報 : dump_page_set_addr (ダンプ対象ページリストのアドレス情報)
- 1605 (3) crash 内にプラグインが取込まれる。
 1606 コマンド : extend <crash utility extension のパス (例 : dump2mcdump.so)>
- 1607 (4) ダンプ形式変換ツール (ldump2mcdump) を実行し、(2) のシンボル情報を用いて、McK-
 1608 ernel に割り当てられた物理アドレス領域の情報を取得する。
 1609 取得した情報を用いて、(1) から McKernel 関連情報を取り出し、eclair 形式に整形して
 1610 出力する。

1611 2.13.7 利用時の留意事項

1612 Linux 主導ダンプは RHEL-7.4 以降のバージョンの Linux カーネルでのみ動作する。また、
 1613 Linux のブートパラメタに crash_kexec_post_notifiers を指定して、Linux に、ファイル
 1614 を生成する前に panic_notifier に登録された関数を呼ばせる必要がある。

1615 2.14 プロセスダンプ

1616 McKernel では Linux と同様の方法でプロセスのダンプファイルを生成できる。

1617 2.14.1 実装の制限

1618 出力される情報についての制限は以下の通り。

- 1619 1. 複数スレッドが存在しても、親プロセスの情報しか出力しない。
- 1620 2. NOTE セグメントに格納されるプロセス状態 (struct elf_prstatus64 型) のうち、以
1621 下のフィールドに対応する情報は格納しない。

```
1622     struct elf_siginfo pr_info;    /* シグナル情報 */
1623     short int pr_cursig;          /* 現在のシグナル */
1624     a8_uint64_t pr_sigpend;      /* ペンディングされているシグナル */
1625     a8_uint64_t pr_sighold;      /* hold されているシグナル */
1626     pid_t pr_pid;
1627     pid_t pr_ppid;
1628     pid_t pr_pgrp;
1629     pid_t pr_sid;
1630     struct prstatus64_timeval pr_utime; /* ユーザ時間 */
1631     struct prstatus64_timeval pr_stime; /* システム時間 */
1632     struct prstatus64_timeval pr_cutime; /* 累積ユーザ時間 */
1633     struct prstatus64_timeval pr_cstime; /* 累積システム時間 */
```

1634 なお、struct elf_siginfo は以下のように定義される。

```
1635     struct elf_siginfo {
1636         int si_signo; /* signal number */
1637         int si_code; /* extra code */
1638         int si_errno; /* errno */
1639     };
```

- 1640 3. NOTE セグメントに格納されるプロセス情報 (struct elf_prpsinfo64 型) のうち、以
1641 下のフィールドに対応する情報は格納しない。

```
1642     char pr_sname;          /* プロセス状態 (文字列) */
1643     char pr_zomb;          /* Zombie か否か */
1644     char pr_nice;          /* Nice 値 */
1645     a8_uint64_t pr_flag;   /* フラグ */
1646     unsigned int pr_uid;
1647     unsigned int pr_gid;
1648     int pr_ppid, pr_pgrp, pr_sid;
1649     char pr_fname[16];     /* 実行可能ファイル名 */
1650     char pr_psargs[ELF_PRARGSZ]; /* 引数リスト先頭部分 */
```

1651 2.15 Utility Thread Offloading

1652 McKernel は、スレッドを Linux の CPU にマイグレートする機能を提供する。この機能によ
1653 り、通信のプログレススレッドなどのヘルパースレッド (utility thread と呼ぶ) を、計算用
1654 CPU 資源を利用することなく実行することができる。

1655 Linux CPU へのマイグレートは、以下の処理の組み合わせによって実現する。

- 1656 1. スレッドマイグレート処理
1657 McKernel で実行しているスレッドを Linux CPU にマイグレートする処理
- 1658 2. システムコール処理
1659 Linux にマイグレートしたスレッドの発行するシステムコールと McKernel スレッドの
1660 発行するシステムコールとの一貫性を担保する処理
- 1661 3. シグナル受信処理
1662 Linux にマイグレートしたスレッドへシグナルを中継する処理
- 1663 4. スレッド終了処理
1664 Linux にマイグレートしたスレッドを正しく終了させる処理

1665 以下、それぞれの処理の概要を説明する。

1666 2.15.1 スレッドマイグレート処理

1667 スレッドマイグレートの処理のうち、McKernel 側の処理は以下の通り。

- 1668 1. McKernel で実行中のユーザスレッドが自スレッドを Linux にマイグレートすることを
1669 システムコールを用いて指示する。
- 1670 2. McKernel はシステムコールを受けて以下の処理を行う。
 - 1671 (a) マイグレート対象のユーザスレッドのコンテキストを取得する。
 - 1672 (b) McKernel から Linux へのシステムコール委譲を用いて、`mcexec` に対してスレッ
1673 ドのマイグレート指示を行う。このとき、取得したコンテキストを引き渡す。
- 1674 3. McKernel のスレッドは Linux へマイグレートされたスレッドが終了するまでシステム
1675 コール完了を待ってスリープする。
- 1676 4. McKernel はシステムコールが完了すると当該スレッドを起床する。
- 1677 5. McKernel のスレッドはシステムコールの戻り値を引数として `_exit` を呼び出し、スレッ
1678 ドを終了する。

1679 スレッドマイグレートの処理のうち、`mcexec` 側の処理は以下の通り。

- 1680 1. McKernel からスレッドマイグレート指示を受ける。
- 1681 2. システムコールワーカースレッドを新規に生成する。これは、自スレッドでマイグレー
1682 トしたコンテキストを処理するため、システムコールワーカースレッドが不足するため
1683 である。生成したスレッドはシステムコール委譲待ちとなる。なお、マイグレートは一
1684 回のみ可能であるため、システムコールワーカースレッドが必要以上に生成されること
1685 はない。

- 1686 3. 当該スレッドの孫プロセスを生成する。当該孫プロセスは当該スレッドに `ptrace` シス
1687 テムコールを用いて接続し、当該スレッドが発行するシステムコールを捕捉する（次節
1688 で説明する）。
- 1689 4. 自スレッドにおいて、コンテキストをマイグレート対象スレッドのコンテキストに切り
1690 替える。このとき、切り替え前のコンテキストを保存しておく。切り替え前コンテキス
1691 トを保存するのは、シグナル受信時やスレッド終了時に一時的に `mcexec` のコンテキス
1692 トに復帰する必要があるためである。
- 1693 5. コンテキスト切り替え後、スレッドはマイグレート対象のスレッドとして、マイグレー
1694 ト指示のシステムコールからの戻りアドレスから処理を再開する。

1695 2.15.2 システムコール処理

1696 Linux にマイグレートしたスレッドが発行するシステムコールは捕捉し、必要に応じて `McKernel`
1697 に処理を依頼する。これは、システムコールの中には、`futex()` や `mmap()` など、`McKernel`
1698 の状態を操作するものがあるためである。

1699 システムコールは `ptrace` を用いて捕捉する。具体的には、マイグレート時にマイグレー
1700 トしたスレッド (`tracee`) を監視する `tracer` プロセスを Linux 上で生成し、`tracee` にシステム
1701 コールを発行を報告させる。また、`tracer` が `tracee` のレジスタを操作することで必要に応じ
1702 て Linux 上でのシステムコール発行をスキップさせる。

システムコールごとの処理を??に示す。

Table 2.8: マイグレートされたスレッドが発行するシステムコールの処理

システムコール	処理
<code>mmap</code> , <code>mprotect</code> , <code>munmap</code> , <code>brk</code> , <code>futex</code>	IKC を用いて <code>McKernel</code> に処理を依頼する
<code>getpid</code> , <code>gettid</code>	<code>mcexec</code> に記録しておいた <code>id</code> を返す
<code>open</code> , <code>read</code> , <code>write</code> など <code>McKernel</code> からは システムコール移譲を行うもの	Linux 上でシステムコールを発行する
<code>exit_group</code> , <code>_exit</code>	<code>mcexec</code> で処理する
それ以外	エラー (<code>ENOSYS</code>) とする

1703

1704 2.15.3 シグナル受信処理

1705 2.15.3.1 シグナル送信処理

1706 既存の `McKernel` に、Linux にシステムコールマイグレートしているプロセスに対して `McK-`
1707 `ernel` からシグナルを送信し、処理を中断する処理が存在する。この処理を応用して、Linux
1708 にマイグレートしたスレッドへのシグナル配送を実現する。

1709 具体的には以下の処理を行う。

- 1710 1. シグナル送信処理において、シグナル配送先スレッドが Linux にマイグレートされてい
1711 る場合、システムコールオフロード中スレッドへのシグナル送信と同様に IKC を通じ
1712 て `mcexec` にシグナル送信を依頼する。
- 1713 2. `mcexec` は `McKernel` のスレッド ID (リモートスレッド ID) から Linux 上のスレッド
1714 ID (ローカルスレッド ID) への変換を行い、ローカルスレッド ID に対してシグナルを
1715 送信する。

1716 2.15.3.2 mcexec のシグナルハンドラ処理

1717 mcexec は Linux から McKernel のスレッドに送られたシグナルを McKernel へ中継するため
1718 に、特別なシグナルハンドラを登録している。このため、McKernel のスレッドから Linux に
1719 マイグレートされたスレッドに送られるシグナルに対応するシグナルハンドラの呼び出しは、
1720 直接行うことができず、この特別なシグナルハンドラから行う必要がある。

1721 mcexec のシグナルハンドラの入り口と出口では、TLS を切り替える必要がある。TLS
1722 を切り替えない場合、Linux にマイグレートされたスレッドの `errno` やその他の TLS 領域を
1723 破壊する可能性があるためである。TLS の切り替えは、libc の `arch_prctl` や `syscall` を使用で
1724 きない (これらは `errno` を更新する)。TLS の切り替えはシステム依存の手段でシステムコール
1725 を呼び出す必要がある (例えば、x86_64 では `syscall` 命令の発行)。

1726 TLS を切り替えるため、Linux にマイグレートしているスレッドに対して `mcctrl` におい
1727 て mcexec の TLS と McKernel スレッドの TLS を保持しておく。TLS 切り替え要求に対して
1728 通常の mcexec のスレッドは何もしないが、Linux にマイグレートされたスレッドでは、シグ
1729 ナルハンドラの入りの入り口では mcexec の TLS に切り替え、出口では McKernel スレッドの TLS
1730 に切り替える。

1731 2.15.4 スレッド終了処理

1732 マイグレートされたスレッドの終了処理のステップは以下の通り。

- 1733 1. スレッド終了の捕捉
1734 マイグレートしたスレッドが `_exit()` を呼び出した場合は `tracer` が `ptrace` で補足し、
1735 シグナルによってスレッドが終了する場合は `mcctrl` が Linux のスレッド終了フック
1736 (`trace_sched_process_exit`) を用いて捕捉する。
- 1737 2. McKernel 側でのスレッド終了
1738 mcexec が IKC を用いて McKernel にスレッドの終了ステータスを通知し、McKernel
1739 側のスレッドを終了させる。
- 1740 3. Linux 側でのスレッド終了
1741 以下のステップで Linux 側のスレッドを終了させる。
 - 1742 (a) mcexec のスレッドのコンテキストを、マイグレートした McKernel スレッドのもの
1743 のから、`mcctrl` に記録しておいたマイグレート処理前のものへ切り替える。コン
1744 テキストの切り替えによって、mcexec の `ioctl()` の直後から実行が再開される。
 - 1745 (b) `tracer` プロセスが終了する。
 - 1746 (c) mcexec のスレッドが `_exit()` を呼び出すことで終了する。

1747 2.15.5 実装詳細

1748 2.15.5.1 Linux CPU へのスレッド生成の構成

1749 Linux CPU へのスレッド生成の構成を図??に示す。

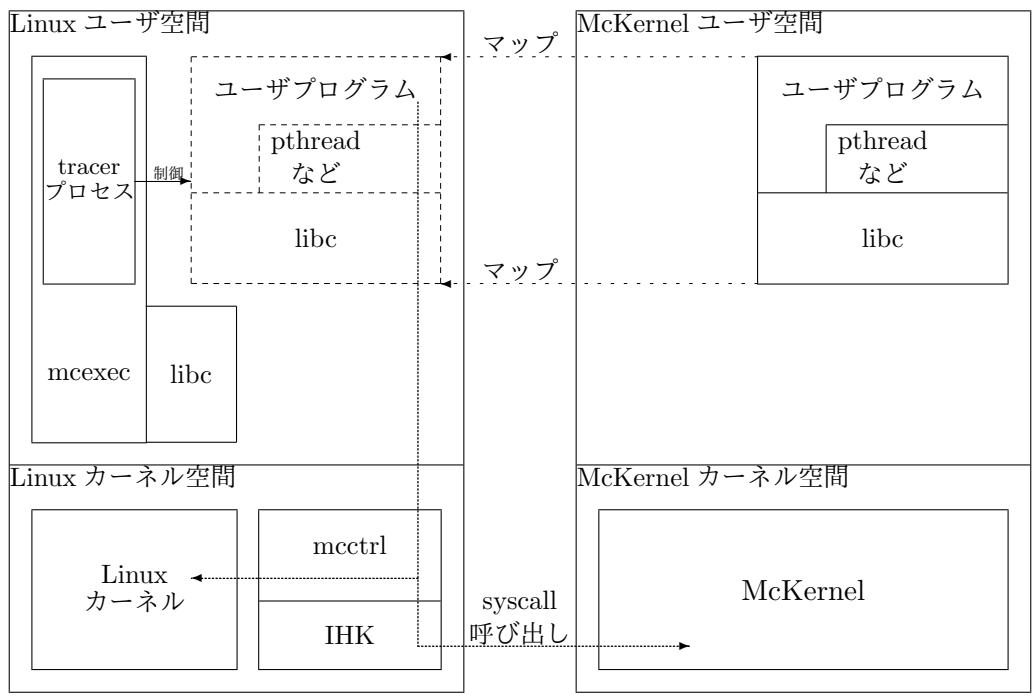


Figure 2.21: Linux CPU へのスレッド生成の構成

1750 ユーザプログラムのスレッドを Linux CPU に生成 (実際はマイグレート) すると、ユー
1751 ザプログラムが占めるメモリが Linux ユーザ空間にマップされ、Linux ユーザ空間の `mcexec`
1752 から参照可能となる。このとき、Linux ユーザ空間内には、`mcexec` の `libc` とユーザプログラ
1753 ムの `libc` が異なる実体として配置されている。

1754 McKernel のスレッドを Linux に生成するとき、`mcexec` の子プロセスとして `tracer` プロ
1755 セスを生成する。`tracer` プロセスは Linux に生成したスレッドのシステムコールを監視し、
1756 Linux CPU のユーザプログラムが発行したシステムコールの一部 (`mmap` など) を `mcctrl` 経
1757 由で McKernel 上で処理するように制御する。

1758 2.15.5.2 `util_indicate_clone` システムコール

1759 `util_indicate_clone` システムコールは自スレッドの `thread` 構造体に `util_indicate_clone`
1760 システムコールの引数 `mod` と `arg` (カーネル空間にコピー済の `arg`) を設定する。

1761 `thread` 構造体の関連フィールドは以下のように定義される。

```
1762 struct thread {  
1763     ... 略 ...  
1764     int mod_clone;                // 生成対象 OS  
1765     void *mod_clone_arg;          // CPU 位置の指示、スレッドの振る舞いの記述  
1766 };
```

1767 2.15.5.3 `util_migrate_inter_kernel` システムコール

1768 `util_migrate_inter_kernel` システムコールは以下の処理を行う。

- 1769 1. `arg` が非 NULL の場合、`arg` の内容をユーザ空間からカーネル空間にコピーする。コピー
1770 に失敗した場合、`EFAULT` を返却する。
- 1771 2. コンテキスト退避用ページを確保する。
- 1772 3. コンテキスト退避用ページにユーザコンテキストを退避する。
- 1773 4. コンテキスト退避用ページの物理アドレスを引数として、`sched_setaffinity` をオフロー
1774 ドする²。
- 1775 5. コンテキスト退避用ページを解放する。
- 1776 6. `sched_setaffinity` のオフロードの戻り値が正 (成功) の場合、以下を行う。
 - 1777 (a) 戻り値の `0x100000000` ビットが立っている場合、プロセスの終了を表すため、`terminate()`
1778 でプロセスを終了する。
 - 1779 (b) それ以外の場合、`sched_setaffinity()` の処理をもう一度 McKernel に依頼する
1780 ことで、コンテキスト保存領域を `unmap` し、`do_exit()` でスレッドを終了する。
- 1781 7. `sched_setaffinity` のオフロードの戻り値が負 (エラー) の場合、戻り値を返却する。

1782 2.15.5.4 `get_system` システムコール

1783 `get_system` システムコールは McKernel 上で実行すると 0 を返却する。Linux 上で実行する
1784 と、当該システムコールは存在しないため、`ENOSYS` でエラーリターンする。

²他のシステムコール番号と被らないように、オフロード対象ではない `sched_setaffinity` のシステムコール番号を使用している。`sched_setaffinity` をオフロードすると、`mcexec` にてスレッドオフロード処理を行う

1785 **2.15.5.5** clone システムコール

1786 clone システムコールにて、子スレッドの thread 構造体を runq に接続する (runq_add_thread
1787 呼び出し) 前に以下の処理を行う。

- 1788 1. 親スレッドの mod_clone に SPAWN_TO_REMOTE が設定されている場合、子スレッドの
1789 thread の mod_clone に SPAWNING_TO_REMOTE を設定する。これにより、子スレッドを
1790 schedule が処理するときに util_migrate_inter_kernel の処理が行われる。

1791 **2.15.5.6** schedule の処理

1792 schedule に対して、以下の変更を行う。

- 1793 1. next を探す処理において、mod_clone に SPAWNING_TO_REMOTE が設定されているスレッ
1794 ドが runq に有る場合、そのスレッドを優先して next に設定する。

1795 **2.15.5.7** enter_user_mode の処理

1796 enter_user_mode に対して、以下の変更を行う。

- 1797 1. check_signal を呼び出した後で、auto_utilthr_migrate を呼び出す。この処理は current
1798 スレッドに SPAWNING_TO_REMOTE が設定されている場合、スレッドを開始する前に Linux
1799 にマイグレートする。

1800 **2.15.5.8** auto_utilthr_migrate の処理

1801 auto_utilthr_migrate は以下の処理を行う。

- 1802 1. current スレッドの mod_clone に SPAWNING_TO_REMOTE が設定されている場合、
1803 mod_clone に SPAWN_TO_LOCAL を設定し、util_migrate_inter_kernel を呼び出す。こ
1804 れによって、新しいスレッドの開始前に util_migrate_inter_kernel システムコールの処
1805 理が実行される。

1806 **2.15.5.9** do_syscall の処理

1807 do_syscall に対して、以下の変更を行う。

- 1808 1. オフロードしたシステムコールの状態が、STATUS_SYSCALL の場合 (Linux から McK-
1809 ernel へのシステムコール委譲)、以下の処理を行う。

- 1810 (a) システムコール番号が rt_sigreturn の場合、シグナルハンドラの内容を返却する。
1811 (b) システムコール番号が rt_sigreturn 以外の場合、以下を行う。
1812 i. syscall_table を検索し、システムコール番号が登録されているか調べ、登録さ
1813 れていない場合は ENOSYS を返却する。
1814 ii. システムコールコンテキストを作成し、syscall_table に登録されているシステ
1815 ムコール処理を呼び出す。結果を返却する。
1816 (c) システムコールの結果は、send_syscall 呼び出しによって、IKC を通じて Linux に
1817 通知する。この処理は remote page fault と同様である。

1818 **2.15.5.10 mcexec の処理**

1819 mcexec は以下の処理を行う。

- 1820 1. sched_set_affinity に対するオフロード処理として、以下を行う。
- 1821 (a) create_worker_create を呼び、新しいワーカースレッドを作成する。このスレッド
 - 1822 は、不足するシステムコールオフロードスレッドを補完するものである。
 - 1823 (b) mcexec スレッドと McKernel スレッドのコンテキスト退避領域を作成する。
 - 1824 (c) mcctrl に McKernel スレッドのコンテキストを McKernel スレッドのコンテキス
 - 1825 ト退避領域へコピーさせる (MCEXEC_UP_UTIL_THREAD1)。
 - 1826 (d) tracer プロセスを生成する。tracer プロセスの詳細は 4 に示す。
 - 1827 (e) uti_attr が指定されている場合、mcctrl に uti_attr 処理を依頼する
 - 1828 (MCEXEC_UP_UTI_ATTR)。
 - 1829 (f) 以下に示す switch_ctx を行う。
 - 1830 i. mcexec スレッドのコンテキスト退避領域に現在のコンテキストを退避する。
 - 1831 ii. mcctrl に Linux にマイグレートされたスレッドの情報を登録する
 - 1832 (MCEXEC_UP_UTIL_THREAD2)。
 - 1833 iii. コンテキストをマイグレートされたスレッドのコンテキストに切り替える。
 - 1834 (g) マイグレートしたスレッドが完了した後に元のコンテキストに戻る。
 - 1835 (h) mcexec スレッドのコンテキスト退避領域を解放する。
 - 1836 (i) スレッドを終了する。
- 1837 2. シグナルを受信した際、シグナルハンドラにて以下の処理を行う。
- 1838 (a) mcctrl に MCEXEC_UP_SIG_THREAD を要求し、mcexec の TLS に切り替える。
 - 1839 (b) mcexec のスレッドの場合、McKernel に受信したシグナルを通知する。
 - 1840 (c) Linux にマイグレートしたスレッドの場合は以下の処理を行う。
 - 1841 i. mcctrl に rt_sigaction で MCEXEC_UP_SYSCALL_THREAD 要求を行い、シ
 - 1842 グナルハンドラの設定を取得する。
 - 1843 ii. シグナルハンドラが SIG_IGN の場合、何もしない。
 - 1844 iii. シグナルハンドラが SIG_DFL の場合、シグナルが SIGCHLD、SIGURG、SIG-
 - 1845 CONT 以外の場合はシグナルハンドラを解除し、自プロセスにシグナルを送
 - 1846 付する。これによって、当該シグナルを受信して自プロセスが終了する。
 - 1847 iv. シグナルハンドラがアドレスの場合、一時的に TLS を元に戻してアドレスの
 - 1848 関数 (シグナルハンドラ) を呼び出す。
 - 1849 (d) mcctrl に MCEXEC_UP_SIG_THREAD を要求し、元の TLS に切り替える。
- 1850 3. tracer プロセスは以下の処理を行う。
- 1851 (a) tracee にて、tracer と待ち合わせに使用するパイプを作成する。
 - 1852 (b) tracer プロセスを tracee の孫プロセスとして fork する。孫プロセスとするのは
 - 1853 tracee が tracer を wait しないまま終了する場合に対応するためである。
 - 1854 (c) tracee は以下の処理を行う。
 - 1855 i. tracer の予期せぬ終了を検知できるように、パイプの出力側を閉じる。

- 1856 ii. 子プロセスの終了を待つ。子プロセスはすぐに終了する (孫プロセスが tracer
- 1857 になる)。
- 1858 iii. パイプの入カイベントの発生を (最大)1 秒待つ (select)。
- 1859 iv. イベントが発生せずに 1 秒経過した場合、タイムアウトでエラーリターン。
- 1860 v. select がエラーの場合は、そのエラーコードでエラーリターン。
- 1861 vi. パイプから 1 バイト読み込み、パイプを閉じる。
- 1862 vii. パイプから 1 バイト読み込めなかった場合 (EOF)、EAGAIN でエラーリターン。
- 1863 viii. 正常にリターン。(以下、tracee スレッドはオフロード処理を継続する。)
- 1864 (d) パイプの入力側を閉じる。
- 1865 (e) 子プロセスを fork し、親プロセスは終了 (exit) する。子プロセスが tracer となる。
- 1866 (f) /dev/mcos 以外のファイルディスクリプタを全て閉じる。
- 1867 (g) 標準入出力を /dev/null に割り当てる。
- 1868 (h) tracee スレッドに PTRACE_ATTACH する。
- 1869 (i) tracee の停止を待つ (wait)。
- 1870 (j) PTRACE_SYSCALL 後の停止理由がシステムコールかシグナルかの区別を付ける
- 1871 ために、PTRACE_SETOPTIONS で PTRACE_O_TRACESYSGOOD を指定
- 1872 する。
- 1873 (k) パイプに 1 バイト書き出し、パイプを閉じる。
- 1874 (l) 以下、無限ループ。
 - 1875 i. PTRACE_SYSCALL により tracee を再開する。
 - 1876 ii. tracee の停止を待つ (wait)。
 - 1877 iii. tracee が終了した場合、終了コードを McKernel に通知し、tracer を終了する。
 - 1878 iv. 停止以外の場合、continue³。
 - 1879 v. システムコールで停止した場合、以下を行う。
 - 1880 A. PTRACE_GETREGS を行い、tracee のレジスタを得る。
 - 1881 B. システムコール番号が ioctl で引数に MCEXEC_UP_SYSCALL_THREAD
 - 1882 が指定されている場合、戻り値を逆オフロード結果に書き換える。
 - 1883 C. システムコール番号が逆オフロード対象で、戻り値 (x86 の場合、rax) が-
 - 1884 ENOSYS の場合 (システムコール呼び出し時)、システムコール番号を ioctl
 - 1885 に変更し、システムコール逆オフロードの引数を設定する。
 - 1886 D. PTRACE_SETREGS を行い、tracee のレジスタを更新する。
 - 1887 vi. システムコール以外 (つまりシグナル) で停止した場合、次回 PTRACE_SYSCALL
 - 1888 に指定するシグナルとして、停止シグナルを設定する。

1889 2.15.5.11 mcctrl の処理

1890 MCEXEC_UP_UTIL_THREAD1 コマンドに対しては、McKernel から渡された物理アドレスで示さ
 1891 れる McKernel スレッドのコンテキストを、mcexec のコンテキスト退避領域にコピーする。
 1892 MCEXEC_UP_UTIL_THREAD2 コマンドに対しては、host_thread 構造体を作成する。
 1893 host_thread 構造体は Linux にマイグレートされたスレッドの情報を保持し、以下のよ
 1894 うに定義される。

³停止以外の可能性としては、SIGCONT による処理再開が考えられる。

```

1895 struct host_thread {
1896     struct host_thread *next;           // 同一 PID 内のリスト
1897     struct mcos_handler_info *handler;  // LWK 情報へのハンドラ
1898     int pid;                            // プロセス ID
1899     int tid;                            // スレッド ID
1900     unsigned long usp;                  // mcexec コンテキストの SP
1901     unsigned long lfs;                  // mcexec の TLS ベース
1902     unsigned long rfs;                  // ユーザプログラムの TLS ベース
1903 };

```

1904 `mcos_handler_info` は LWK の情報を保持し、以下のように定義される。

```

1905 struct mcos_handler_info {
1906     int pid;
1907     int cpu;
1908     struct mcctrl_usrdata *ud;
1909     struct file *file;
1910 };

```

1911 2.15.5.11.1 MCEXEC_UP_UTI_ATTR の処理

1912 MCEXEC_UP_UTI_ATTR の要求に対して、以下の処理を行う。

- 1913 1. 初回の場合、以下を行う。
 - 1914 (a) Linux カーネル内の `sched_setaffinity` と `sched_setscheduler_nocheck` のアドレスを
1915 解決する。
 - 1916 (b) ラウンドロビン管理用配列を確保し、0 で初期化する。
- 1917 2. `uti_attr` のフラグに、背反する組み合わせが指定されている場合、エラーリターンする
1918 (`-ENOMEM`)。
- 1919 3. `mcctrl_usrdata` の `cpu_topology_list` を辿って、McKernel にて clone を発行したスレッ
1920 ド(親スレッド)の McKernel の CPU ID を持つ CPU を検索する。存在しない場合はエ
1921 ラーリターンする (`-EINVAL`)。
- 1922 4. 作業用の `cpumask` を確保し、`cpu_active_mask` で初期化する。確保できない場合、エ
1923 ラーリターンする (`-ENOMEM`)。
- 1924 5. 以下の処理によって、`uti_attr` のフラグ指定に従い、割り当てる CPU の候補を求め、作
1925 業用 `cpumask` に設定する。
 - 1926 (a) フラグに `UTI_FLAG_NUMA_SET` が設定されている場合、以下の処理を行う。
 - 1927 i. `mcctrl_usrdata` の `node_topology_list` を辿って、`numa_set` に設定されている
1928 NUMA ID と一致する NUMA ノードを検索し、見付かった NUMA ノードに
1929 属す CPU の和集合を求める。
 - 1930 ii. 作業用 `cpumask` と NUMA ノードに属す CPU 集合の積集合を求め、作業用
1931 `cpumask` に設定する。
 - 1932 (b) フラグに `UTI_FLAG_SAME_NUMA_DOMAIN` か `UTI_FLAG_DIFFERENT_NUMA_DOMAIN` が
1933 設定されている場合、以下の処理を行う。

- 1934 i. 全ての NUMA ドメインについて、親スレッドが属しているかどうかを調べ、
1935 UTI_FLAG_SAME_NUMA_DOMAIN が指定されている場合は当該ドメインに属す CPU
1936 集合の、また、UTI_FLAG_DIFFERENT_NUMA_DOMAIN が指定されている場合には
1937 親スレッドが属していないドメインの CPU 集合の和集合を求める。
- 1938 ii. 作業用 cpumask と NUMA ノードに属す CPU 集合の積集合を求め、作業用
1939 cpumask に設定する。
- 1940 (c) フラグに UTI_FLAG_SAME_L3 か UTI_FLAG_DIFFERENT_L3 が設定されており、且つ、
1941 親スレッドの CPU が L3 キャッシュと持つ場合、以下を行う。
- 1942 i. UTI_FLAG_SAME_L3 が指定されている場合、キャッシュを共有する CPU の集合
1943 を求める。
- 1944 ii. UTI_FLAG_DIFFERENT_L3 が指定されている場合、キャッシュを共有する CPU
1945 の集合の補集合を求める。
- 1946 iii. 求めた CPU 集合と作業用 cpumask の積集合を求め、作業用 cpumask に設定
1947 する。
- 1948 (d) フラグに UTI_FLAG_SAME_L2 か UTI_FLAG_DIFFERENT_L2 が設定されており、且つ、
1949 親スレッドの CPU が L2 キャッシュと持つ場合、以下を行う。
- 1950 i. UTI_FLAG_SAME_L2 が指定されている場合、キャッシュを共有する CPU の集合
1951 を求める。
- 1952 ii. UTI_FLAG_DIFFERENT_L2 が指定されている場合、キャッシュを共有する CPU
1953 の集合の補集合を求める。
- 1954 iii. 求めた CPU 集合と作業用 cpumask の積集合を求め、作業用 cpumask に設定
1955 する。
- 1956 (e) フラグに UTI_FLAG_SAME_L1 か UTI_FLAG_DIFFERENT_L1 が設定されており、且つ、
1957 親スレッドの CPU が L1 キャッシュと持つ場合、以下を行う。
- 1958 i. UTI_FLAG_SAME_L1 が指定されている場合、キャッシュを共有する CPU の集合
1959 を求める。
- 1960 ii. UTI_FLAG_DIFFERENT_L1 が指定されている場合、キャッシュを共有する CPU
1961 の集合の補集合を求める。
- 1962 iii. 求めた CPU 集合と作業用 cpumask の積集合を求め、作業用 cpumask に設定
1963 する。
- 1964 6. 以下の処理によって、CPU アフィニティの設定、及び、スケジューラの設定を行う。
- 1965 (a) 作業用 cpumask が空集合の場合は何もしない。
- 1966 (b) フラグに UTI_FLAG_EXCLUSIVE_CPU が設定されている場合、以下の処理を行う。
- 1967 i. 作業用 cpumask から CPU を 1 つ選んで CPU アフィニティに設定する。(CPU
1968 の選択方法は後述)
- 1969 ii. スケジューラに SCHED_FIFO を設定する。
- 1970 (c) フラグに UTI_FLAG_CPU_INTENSIVE が設定されている場合、以下の処理を行う。
- 1971 i. 作業用 cpumask から CPU を 1 つ選んで CPU アフィニティに設定する。(CPU
1972 の選択方法は後述)
- 1973 (d) フラグに UTI_FLAG_HIGH_PRIORITY が設定されている場合、以下の処理を行う。
- 1974 i. 作業用 cpumask から CPU を 1 つ選んで CPU アフィニティに設定する。(CPU
1975 の選択方法は後述)

- 1976 ii. スケジューラに SCHED_FIFO を設定する。
- 1977 (e) フラグに UTI_FLAG_NON_COOPERATIVE が設定されている場合、以下の処理を行う。
- 1978 i. 作業用 cpumask から CPU を 1 つ選んで CPU アフィニティに設定する。(CPU
- 1979 の選択方法は後述)
- 1980 (f) 以上に該当しない場合、作業用 cpumask の内容を CPU アフィニティに設定する。

1981 7. 作業用 cpumask を開放する。

1982 作業用 cpumask から CPU を 1 つ選択する手順を以下に示す。

- 1983 1. ラウンドロビン管理用配列は CPU ID ごとのマイグレートスレッド数を記録している。
- 1984 この配列を用いて、作業用 cpumask に含まれかつマイグレートスレッド数が最小とな
- 1985 る CPU ID を求める。
- 1986 2. compare and swap によってラウンドロビン管理用配列を更新 (1 加算) する。失敗した
- 1987 場合は、1 から再度行う。
- 1988 3. 更新した CPU ID 以外の作業用 cpumask のビットをクリアする。

1989 **2.15.5.11.2 MCEXEC_UP_SIG_THREAD の処理**

1990 MCEXEC_UP_SIG_THREAD の要求に対して、以下の処理を行う。

- 1991 1. 要求元スレッドが Linux にマイグレートされたスレッドでない場合、EINVAL でエラー
- 1992 リターンする。
- 1993 2. 引数に従い、スレッドの FS ベースアドレスを切り替える。

1994 **2.15.5.11.3 MCEXEC_UP_SYSCALL_THREAD の処理**

1995 MCEXEC_UP_SYSCALL_THREAD の要求に対して、以下の処理を行う。

- 1996 1. システムコール番号とシステムコール引数を syscall_request 構造体に設定する。
- 1997 2. util_migrate_inter_kernel の要求を検索する。存在しない場合は ENOENT でエラーリ
- 1998 ターンする。
- 1999 3. wait_queue_head_list_node を作成し、システムコールの結果待ちに備える。
- 2000 4. util_migrate_inter_kernel のレスポンスに syscall_request の物理アドレスを設定する。ま
- 2001 た、レスポンスの状態をシステムコール要求に設定する。
- 2002 5. _notify_syscall_requester を呼び出して、McKernel 側にレスポンスの変更を通知する。
- 2003 6. wait_event_interruptible によって、システムコール完了を待つ。
- 2004 7. 結果を返却する。

2005 2.15.6 実装の制限

2006 制限は以下の通り。

- 2007 ● Linux CPU にマイグレートしたスレッドに対して ptrace システムコールによるトレー
2008 スは行えない。
- 2009 ● Linux CPU にマイグレートしたスレッドが発行可能なシステムコールの種類は上記で
2010 説明したもののみである。
- 2011 ● Linux CPU にマイグレートしたスレッドを、再度 McKernel に移動することはできない。

2012 2.16 高速プロセス起動

2013 McKernel は、複数種の MPI プログラムを起動しさらにそれを繰り返すジョブにおいて MPI
2014 プログラム起動時間を短縮する機能を提供する。利用例としては、アンサンブルシミュレー
2015 ションとデータ同化を繰り返す気象アプリケーションが挙げられる。

2016 起動時間の短縮は、それぞれの MPI プログラムを常駐させて、起動を停止状態からの復
2017 帰で置き換えることで実現する。高速プロセス起動は以下の機能から構成される。

- 2018 ● プロセス実行停止および停止からの再開機能
2019 MPI プログラムが、本来プロセスとして終了するタイミングで終了せずに再開指示待
2020 ち状態で停止できるようにする。また、再開指示を受けて停止状態から復帰できるよ
2021 うにする。ライブラリ関数として実装する。
- 2022 ● MPI プログラムの繰り返し起動指示機能
2023 MPI プログラムの実行回数を把握し、初回は `mpiexec` を用いて起動し、2 回目以降は
2024 停止しているプロセスを再開する。`ql_mpiexec.start` と呼ぶユーザコマンドとして実
2025 装する。
- 2026 ● MPI プログラムの繰り返し起動からの終了指示機能
2027 再開指示待ち状態で停止している MPI プロセスを終了させる。`ql_mpiexec.finalize`
2028 と呼ぶユーザコマンドとして実装する。

2029 以下、これらの機能の詳細を説明する。

2030 2.16.1 詳細

2031 関連プロセスを表??示す。

Table 2.9: プロセス一覧

プロセス	説明
ql_mpiexec_start	ユーザがジョブスクリプトに記述して用いる、各回の計算開始を指示するコマンドである。指示は ql_talker と ql_server を経由して MPI プログラムに送られる。なお、一回の計算が完了すると、本コマンドは終了する。
ql_mpiexec_finalize	ユーザがジョブスクリプトに記述して用いる、MPI プログラムの実行終了を指示するコマンドである。指示は ql_talker と ql_server を経由して MPI プログラムに送られる。なお、MPI プログラムの終了と共に本コマンドは終了する。
mpiexec 監視	mpiexec プロセスの起動、死活監視、標準入出力およびエラー出力のリダイレクトを行う。本プロセスは ql_mpiexec_start コマンドより起動され常駐する。また、mpiexec プロセス終了と共に終了する。
mpiexec	MPI プロセスを生成する。本プロセスは mpiexec 監視プロセスの子プロセスとして起動される。すべてのランク終了と共に終了する。
mcexec	ホスト Linux 上で McKernel のユーザプログラムプロセスを生成・管理する。
ql_server	高速プロセス起動対象の MPI プログラムを記録し、ql_mpiexec_{start,finalize} コマンドの指示を MPI プロセスに送る。ql_server は ql_mpiexec_start コマンドから ssh で起動され常駐する。また、高速プロセス起動対象の MPI プログラムが全て終了した時点で終了する。
ql_talker	ql_mpiexec_{start,finalize} と ql_server との間の通信を仲介する。ql_mpiexec_{start,finalize} コマンドから ssh で ql_server が実行されている計算ノードに起動される。本プロセスは指示完了と共に終了する。

プロセス構成を図??に示す。

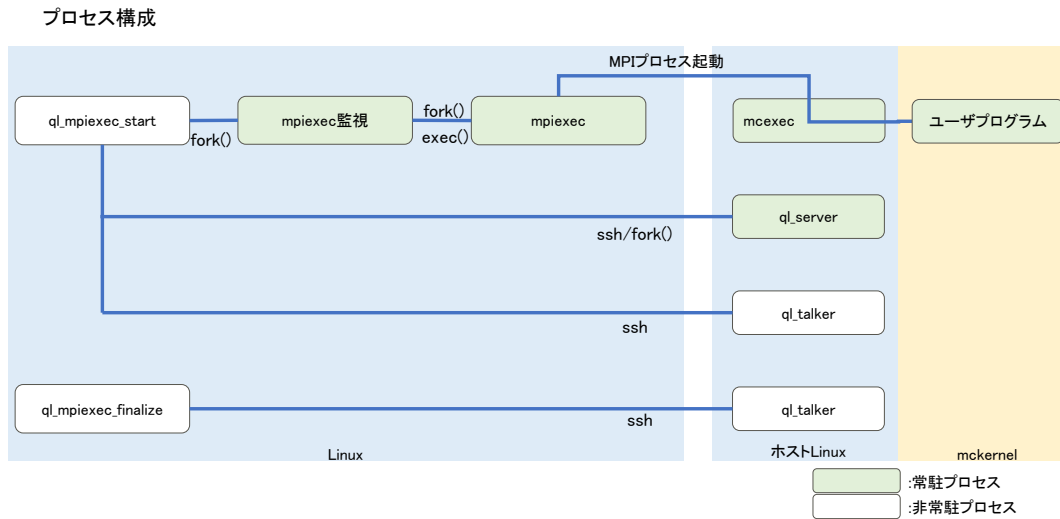


Figure 2.22: プロセス構成

2032
2033

プロセス間通信で用いるコマンドは以下の通り。

コマンド	マクロ	説明
E コマンド	QL_EXEC_END	mcexec が ql_server 経由で ql_mpiexec.start へ各回の計算完了を通知する際に用いる
F コマンド	QL_RET_FINAL	mpiexec 監視プロセスが ql_server へ MPI プログラムの終了を通知する際に用いる
R コマンド	QL_RET_RESUME	ql_server が mcexec へ待ち状態からの起床を指示する際に用いる
N コマンド	QL_COM_CONN	ql_mpiexec.start が ql_server に MPI プログラムの登録を依頼する際に用いる
A コマンド	QL_AB_END	ql_server が他プロセスからのコマンドを処理する際に、コマンド転送先プロセスを見つけれなかった場合に返答として用いる

2034 プロセス間通信の通信電文フォーマットは以下の通り。

2035 <コマンド> <データサイズ> <データ>

各フィールドのサイズ、意味は以下の通り。

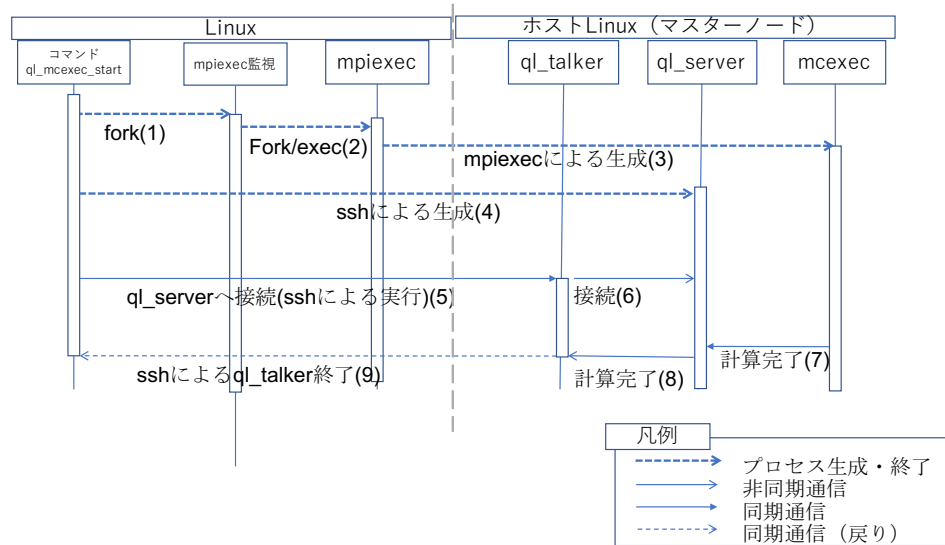
フィールド	サイズ	説明
<コマンド>	1 byte	コマンド
<データサイズ>	4 byte	byte 単位のデータサイズを表す 16 進数文字列
<データ>	可変	データを表す文字列

2036

2037

関連コマンドと関連プロセスの動作フローを図??を用いて説明する。

MPIプログラム初回実行



MPIプログラム再開から終了まで

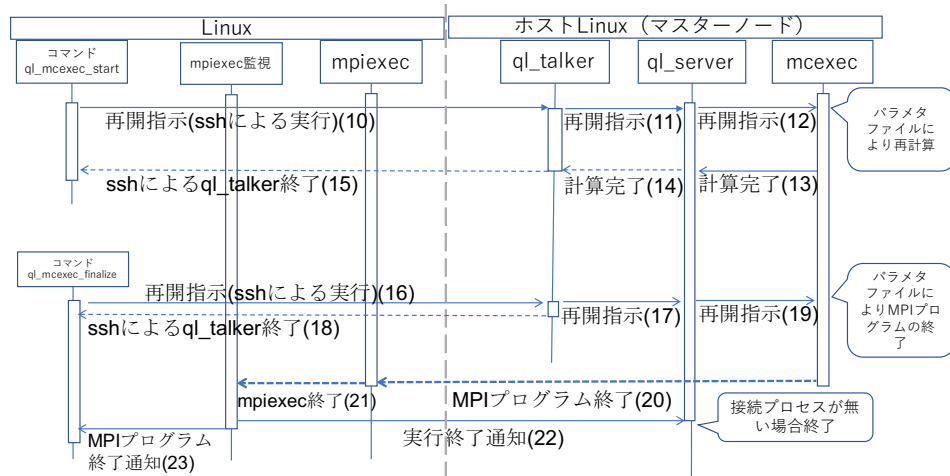


Figure 2.23: 関連コマンドと関連プロセスの動作フロー

- 2038 01 MPIプログラムの初回実行時は、`ql_mpiexec_start` から、`mpiexec` 監視プロセスを
- 2039 `fork()` で常駐プロセスとして起動する。(図の (1))
- 2040 02 `mpiexec` 監視プロセスは、`mpiexec` を `fork()/exec()` で起動する。また、`mpiexec` の
- 2041 標準入出力およびエラー出力を `ql_mpiexec_start` へリダイレクトする。(図の (2))
- 2042 03 `mpiexec` が `mcexec` を用いて MPI プロセスを McKernel 上に起動する。(図の (3))
- 2043 04 `ql_mpiexec_start` が `ssh` で Host Linux 上に `ql_server` を常駐プロセスとして起動す
- 2044 る。(図の (4))
- 2045 05 `ql_mpiexec_start` が `ssh` で Host Linux 上に `ql_talker` を起動する。(図の (5))
- 2046 06 `ql_talker` は `ql_server` に N コマンドを送信し、`ql_server` から計算完了の返信を待
- 2047 つ。`ql_server` は当該 MPI プログラムの存在を管理表に記録する。(図の (6))

- 2048 07 `mcexec` は MPI プログラムの一回の計算完了後 `ql_server` へ計算完了を意味する E コ
2049 マンドを送信し、返信を待つ。(図の (7))
- 2050 08 `ql_server` が `ql_talker` へ計算完了を意味する E コマンドを送信する。(図の (8))
- 2051 09 `ql_talker` は E コマンドを受け取り、正常終了する。`ql_mpiexec_start` は `ql_talker`
2052 の終了を受けてリダイレクトしている標準入出力およびエラー出力をクローズし終了す
2053 る。(図の (9))
- 2054 10 `ql_mpiexec_start` は MPI プログラムの次の計算の開始時、`mpiexec` 監視プロセスに依
2055 頼して `mpiexec` の標準入出力およびエラー出力を自身にリダイレクトする。また、`ssh`
2056 でホスト Linux 上に `ql_talker` を起動する。(図の (10))
- 2057 11 `ql_talker` は `ql_server` へ再開指示を意味する R コマンドを送信し、`ql_server` から
2058 の返信を待つ。(図の (11))
- 2059 12 `ql_server` は `mcexec` へ R コマンドを送信する。MPI プログラムはパラメタファイル
2060 を読み、再開指示であることを確認し、引数と環境変数をパラメタファイルに指定され
2061 たものに置き換え、次の回の計算を行う。(図の (12))
- 2062 13 `mcexec` は一回の計算完了後 `ql_server` へ計算完了 (E コマンド) を送信し、返信を待
2063 つ。(図の (13))
- 2064 14 `ql_server` が `ql_talker` へ計算完了 (E コマンド) を送信する。(図の (14))
- 2065 15 `ql_talker` は E コマンドを受け取り、正常終了する。`ql_mpiexec_start` は終了を受けて
2066 リダイレクトしている標準入出力およびエラー出力をクローズし終了する。(図の (15))
- 2067 16 `ql_mpiexec_finalize` は `mpiexec` 監視プロセスに依頼して `mpiexec` の標準入出力およ
2068 びエラー出力を自身へリダイレクトする。また `ssh` でホスト Linux 上に `ql_talker` を
2069 起動する。(図の (16))
- 2070 17 `ql_talker` は `ql_server` へ再開指示を意味する R コマンドを送信し、終了する。(図の
2071 (17) (18))
- 2072 18 `ql_server` は `mcexec` へ R コマンドを送信する。MPI プロセスは R コマンドを受けて、
2073 パラメタファイルを読み終了指示であることを確認し終了処理を行う。(図の (19))
- 2074 19 `mcexec` は MPI プロセスの終了と共に終了する。(図の (20))
- 2075 20 `mpiexec` は全ランクの終了を待って終了する。`mpiexec` 監視プロセスは `mpiexec` プロ
2076 セス終了を検知し、戻り値を取得する。(図の (21))
- 2077 21 `mpiexec` 監視プロセスは `ql_talker` 経由で `ql_server` へ実行終了を意味する F コマン
2078 ドを送信する。`ql_server` は当該 MPI プログラムを管理表から削除する。また、管理
2079 表が空になった場合は終了する。(図の (22))
- 2080 22 `mpiexec` 監視プロセスは `ql_mpiexec_finalize` へ MPI プログラムの終了を通知し、戻
2081 り値を渡し、終了する。(図の (23))
- 2082 23 `ql_mpiexec_finalize` は MPI プログラムの終了通知を受けて、リダイレクトしている
2083 標準入出力およびエラー出力をクローズし、`mpiexec` の戻り値を自身の戻り値として終
2084 了する。

2085 2.16.2 MPI プロセス起動指示コマンド

2086 書式

2087 `ql_mpiexec_start -machinefile <hostfile_path> [<mpiopts>...] <exe> [<args>...]`

2088 説明

2089 処理ステップは以下の通り。

- 2090 1 ホストファイルの内容、`mpiexec` への引数、実行可能ファイル名から md5 ハッシュに
2091 より MPI プログラム ID を作成する。ID を環境変数 `QL_NAME` に記録する。
- 2092 2 `ql_server` との通信のためのソケットファイルのパスを環境変数 `QL_SOCKET_FILE` に記
2093 録する。また、`ssh` でホストファイルの先頭のホスト（以降、マスターノードと呼ぶ）
2094 上に `ql_server` を起動する。起動失敗した場合は終了コード (-1) で終了する。
- 2095 3 `mpiexec` 監視プロセスとの通信のためのソケットファイルが存在しない場合は、ソケッ
2096 トファイルを作成後、`mpiexec` 監視プロセスを `fork` する。`mpiexec` 監視プロセスは、
2097 `mpiexec` を `fork/exec` で生成する。また、`mpiexec` の標準入出力およびエラー出力を
2098 無名パイプ（以降、リダイレクト用パイプと呼ぶ）の片方の端に接続する。
- 2099 4 再開指示のためパラメタファイルを作成する。
- 2100 5 `mpiexec` 監視プロセスに、自身の標準入出力、エラー出力のファイルディスクリプタ番
2101 号を渡す。`mpiexec` 監視プロセスは当該ファイルディスクリプタをリダイレクト用パイ
2102 プの空いている方の端に接続する。
- 2103 6 第 1 回の計算開始時は `ssh` で `ql_talker` をマスターノード上に起動する。`ql_talker`
2104 は、`ql_server` へ接続を意味する N コマンドを送信し、各回の計算完了を意味する E
2105 コマンドを受信するまで待機する。
- 2106 7 第 2 回目以降の計算開始時は、`ssh` で `ql_talker` をマスターノード上に起動する。`ql_talker`
2107 は、`ql_server` へ再開指示を意味する R コマンドを送信し、各回の計算完了を意味す
2108 る E コマンド受信まで待機する。
- 2109 8 `ql_talker` コマンドがその終了をもって `ql_mpiexec_start` へ計算完了を通知する。
2110 `ql_mpiexec_start` は `mpiexec` 監視プロセスと通信を行って `mpiexec` が終了していない
2111 ことを確認する。`mpiexec` が終了している場合は、`mpiexec` の終了コードを取得する。
- 2112 9 各回の計算完了の場合はパラメタファイルを削除し 0 を返し終了する。`mpiexec` が終了
2113 していた場合は、その終了コードを自身の終了コードに設定して終了する。

`ql_mpiexec_start` が使用する環境変数は以下の通り。

名前	説明	作成・参照
<code>QL_NAME</code>	MPI プログラム ID	作成
<code>QL_SOCKET_FILE</code>	<code>mcexec</code> と <code>ql_server</code> との接続に用いるソケットファイル名	作成

2114

2115 `ql_mpiexec_start` が使用するファイルは以下の通り。

ファイル名	説明	作成・参照
$\${QL_SOCKET_PATH}/ql_sock/<MPI \text{プログラム ID}>.s$	ql_talker と ql_server との間の通信に用いるソケットファイル	作成／参照
$\${QL_PARAM_PATH}/<MPI \text{プログラム ID}>.param$	ql_mpiexec.start から mcexec への指示と、次回の計算に使用する引数と環境変数を記載するコマンド・パラメタファイル	作成

2116 パラメタファイルは、ql_mpiexec_{start,finalize} から mcexec への指示を記載する。
 2117 内容は、起床後の動作および次の回の計算に必要な引数などのデータである。
 2118 フォーマットは以下の通り。

2119 <ヘッダ部>
 2120 <データ部>
 2121 [<データ部>...]

2122 <ヘッダ部>のフォーマットは以下の通り。

2123 0 COM=<mcexec への指示> <引数の数> <環境変数定義の数>

それぞれのフィールドの意味及び取りうる値は以下の通り。

フィールド	説明
<mcexec への指示>	R:次の回の計算開始、F:MPI プロセスの終了
<引数の数>	データ部に存在する引数の数
<環境変数定義の数>	データ部に存在する環境変数定義の数

2124 <データ部>のフォーマットは以下の通り。
 2125

2126 <種別> <データ長> <データ値>

それぞれのフィールドの意味及び取りうる値は以下の通り。

フィールド	説明
<種別>	1:引数、2:環境変数定義
<データ長>	データ長
<データ値>	文字列

2127

2.16.3 MPI プロセス終了指示コマンド

書式

2130 ql_mpiexec_finalize -machinefile <hostfile> [<mpiopts>...] <exe>

説明

2132 処理ステップは以下の通り。

- 2133 1 ホストファイルの内容、mpiexec への引数、実行可能ファイル名から md5 ハッシュにより MPI プログラム ID を作成し、環境変数 QL_NAME に記録する。
- 2134
- 2135 2 mpiexec 監視プロセスと通信を行うソケットファイルの存在を確認し、存在しない場合は ql_mpiexec.start が実行されていないと判断し 1 を返し終了する。
- 2136

- 2137 3 終了指示のためのパラメタファイルを作成する。
- 2138 4 `ql_mpiexec_finalize` は自身の標準入出力とエラー出力のファイルディスクリプタ番
2139 号を `mpiexec` 監視プロセスに渡す。`mpiexec` 監視プロセスは当該ファイルディスクリ
2140 プタをリダイレクト用パイプの空いている方の端に接続する。
- 2141 5 `ssh` でマスターノード上に `ql_talker` を起動する。`ql_talker` は、`ql_server` へ再開指
2142 示を意味する R コマンドを送信する。`ql_mpiexec_start` の場合と異なり `ql_talker` は
2143 `ql_server` からの返答を待つことなく終了する。
- 2144 6 `mpiexec` 監視プロセスは `mpiexec` の終了時にその終了コードを自身の終了コードに設
2145 定し終了する。
- 2146 7 `mpiexec` 監視プロセスの終了を受けてパラメタファイルを削除する。また `mpiexec` 監視
2147 プロセスから渡された `mpiexec` の終了コードを自身の終了コードに設定して終了する。

`ql_mpiexec_finalize` で作成／参照する環境変数は以下の通り。

名前	説明	作成・参照
QL_NAME	MPI プログラム ID	作成

2148

`ql_mpiexec_finalize` で作成／参照するファイルは以下の通り。

ファイル名	説明	作成・参照
<code>#{QL_SOCKET_PATH}/ql_sock/<MPI プログラム ID>.s</code>	<code>ql_talker</code> と <code>ql_server</code> との間の通信に用いるソケットファイル	参照
<code>#{QL_PARAM_PATH}/<MPI プログラム ID>.param</code>	<code>ql_mpiexec_finalize</code> から <code>mcexec</code> への指示を記載するコマンドファイル。	作成

2149

2150 2.16.4 MPI 実行環境初期化関数 (C 言語)

2151 書式

2152 `int MPI_Init(int *argc, char ***argv)`

2153 説明

2154 `argc`, `argv` を用いて高速プロセス起動の初期化を行う。本関数は PMPI インタフェース
2155 により `MPI_Init()` を置き換える。
2156 処理のステップは以下の通り。

- 2157 1 `PMPI_init()` 関数を呼び出し、MPI 環境を初期化する。また、引数情報を取得する。
- 2158 2 `PMPI_init()` が正常終了した場合、`ql_init()` 関数を呼び出し、高速プロセス起動を初
2159 期化する。
- 2160 3 `PMPI_init()` の戻り値自身の戻り値に設定して戻る。

2161 戻り値

2162

戻り値	説明
MPI_SUCCESS	正常終了
MPI_ERR_OTHER	MPI_init() が複数回実行された

2163 2.16.5 MPI 実行環境初期化関数 (fortran)

2164 書式

2165 subroutine MPI_INIT(INT ierr)

2166 説明

2167 Fortran 環境において、高速プロセス起動のための初期化を行う。本関数は、PMPI インタ
2168 フェースにより MPI_INIT を置き換えることで実装される。処理のステップは以下の通り。

- 2169 1 pmpi_init_() が存在していない場合、ierr に MPI_ERR_OTHER をセットして戻る。
- 2170 2 pmpi_init_() を呼び出し、MPI 環境を初期化する。
- 2171 3 戻り値 ierr が MPI_SUCCESS の場合、ql_init_() 関数を呼び出し、高速プロセス起動を
2172 初期化する。

2173 なお、Fortran コンパイラは GNU Fortran Compiler もしくは Intel Fortran Compiler をサ
2174 ポートする。Intel Fortran Compiler を使用する場合は、コンパイルオプションに -shared-intel
2175 を指定する必要がある。

2176 戻り値

戻り値	説明
MPI_SUCCESS	正常終了
MPI_ERR_OTHER	MPI_init() が複数回実行された

2177

2178 2.16.6 計算の再開・終了関数 (C 言語)

2179 書式

2180 ql_client(int *argc, char ***argv)

2181 説明

2182 処理のステップは以下の通り。

- 2183 1 当該プロセスが ql_mpiexec_start により起動されていない場合は、QL_EXIT を返す。
- 2184 2 スレッドの停止を行う。また PMI_Barrier() で計算完了同期を行う。
- 2185 3 システムコールによりカーネルモードに移行し、mcexec に ql_mpiexec_{start,finalize}
2186 による指示待ちを依頼する。
- 2187 4 指示待ちから復帰し、パラメタファイルを参照して指示を確認する。指示が次の回の計
2188 算開始の場合、パラメタファイルを用いて計算のための引数と環境変数を設定する。

- 2189 5 スレッドの再開を行う。
- 2190 6 指示が次の回の計算開始の場合 QL_CONTINUE、MPI プロセスの終了の場合 QL_EXIT を
2191 返す。

2192 2.16.7 計算の再開・終了関数 (Fortran)

2193 書式

2194 subroutine QL_CLIENT(ierr)

2195 説明

2196 ql_client() を呼び、その戻り値を ierr に格納して戻る。

2197 2.16.8 初期化関数

2198 書式

2199 int ql_init(int argc, char **argv)

2200 説明

2201 MPI_Init() から呼びされ、高速プロセス起動の初期化を行う。

2202 処理ステップは以下の通り。

2203 1 環境変数 QL_NAME から MPI プログラム ID を取得する。取得できなかった場合、ql_mpiexec_start
2204 から起動されていないと判断し、QL_NORMAL を返す。

2205 3 MPI プログラム ID から、パラメタファイルのパスを作成する。

2206 4 QL_SUCCESS を返す。

2207 戻り値

戻り値	説明
QL_SUCCESS	高速プロセス起動の初期化成功
QL_NORMAL	当該プロセスが ql_mpiexec_start から起動されていない

2208

2209 2.16.9 計算ノードの管理サーバ

2210 書式

2211 ql_server

2212 説明

2213 ql_server は、ql_mpiexec_start により RANK#0 が存在する計算ノード上に起動され、以
2214 下の処理を行う。

- 2215 1 既に ql_server が起動されている場合は、-1 を返して終了する。
- 2216 2 mcexec、ql_talker との通信に用いるユニックスドメインソケットをオープンする。
- 2217 3 select() で当該ソケットを監視する。
- 2218 4 電文を読み込み、コマンドとデータを取得する。
- 2219 5 ql_talker から N コマンドを受け取った際は、対応する MPI プログラムを管理表に登
2220 録する。また、MPI プロセス ID をインデックスとし ql_server に接続しているプロセ
2221 スを返すマップ（接続マップと呼ぶ）に ql_talker を登録する。
2222 接続マップは struct client_fd で実装される。struct client_fd は以下のように定
2223 義される。

```
2224 struct client_fd {  
2225     int fd;           // 接続元プロセスのファイルディスクリプタ  
2226     int client;       // 接続元プロセスの種別  
2227     char *name;      // MPI プログラム ID  
2228     int status;      // 現在実行中の通信コマンド  
2229 };
```

- 2230 6 mcexec から E コマンドを受けとった際は、ql_mpiexec_{start,finalize} の指示が
2231 あるまで待たせる。また、mcexec を接続マップに登録する。さらに、接続マップを用
2232 いて対応する ql_talker を見つけ、それに対して E コマンドを送信する。

- 2233 7 ql_talker から R コマンドを受けとった際は、ql_talker を接続マップに登録する。ま
2234 た、接続マップを用いて対応する mcexec プロセスを見つけ、それに対して R コマンド
2235 を送信することで mcexec を起床する。

- 2236 8 mpiexec 監視プロセスから F コマンドを受け取った際は、対応する MPI プログラムを
2237 管理表から削除する。管理表が空になった場合は ql_server 自身も終了する。

2238 ql_talker や mcexec が ql_server と通信するために使用するソケットファイルは、環
2239 境変数 QL_SOCKET_PATH が定義されている場合は\${QL_SOCKET_PATH}/ql_sock 下に、定義
2240 されていない場合は/run/user/ユーザ ID/ql_sock 下に作成される。当該ディレクトリは
2241 ql_mpiexec_start コマンドが実行されるノードとランク#0 が実行されるノードからアクセ
2242 スできる必要がある。

2243 2.16.10 指示中継コマンド

2244 書式

```
2245 ql_talker <send_command> <receive_command> <MPI_Program_ID>
```

2246 引数

2247

引数	説明
<send_command>	ql_server へ送信するコマンド (1文字) を指定する。
<receive_command>	ql_server からの受信を期待するコマンド (1文字) を指定する。受信を待たずに終了する場合は、"-n" を指定する。
<MPI_Program_ID>	MPI プログラム ID を指定する。

2248 説明

2249 ql_mpiexec_{start,finalize} から ql_server が動作するノード上に起動され、ql_server
2250 に <send_command> で指定されたコマンドを送り、<receive_command> で指定された応答を待
2251 つ。ql_server とはユニックスドメインソケットを用いて通信する。

2252 処理ステップは以下の通り。

- 2253 1 argc の数をチェックし、4 未満の場合は終了コード-1 で終了する。
- 2254 2 環境変数を参照して ql_server との接続に用いるユニックスドメインソケットを見つ
2255 け、ql_server に接続する。
- 2256 3 <send_command> と <MPI_Program_ID> より電文を作成し、ql_server へ電文を送信する。
2257 失敗した場合は終了コード-1 で終了する。
- 2258 4 <receive_command> に "-n" が指定されていた場合、終了コード 0 で終了する。
- 2259 5 <receive_command> を受信した場合、終了コード 0 で終了する。<receive_command> 以
2260 外の文字列を受信した場合終了コード-2 で終了する。

2261 戻り値

戻り値	説明
0	正常終了
-1	ソケット通信エラー
-2	<receive_command> 以外の文字列を受信

2262

2263 2.16.11 swapout システムコール

2264 書式

2265 int swapout(char *filename, void *workarea, size_t size, int flag)

2266 引数

引数	説明
filename	スワップファイル名へのポインタ
workarea	作業領域へのポインタ
size	作業領域のサイズ
flag	swapout の動作制御用フラグ

2267

2268 説明

2269

2270 A. スワップアウト処理

2271

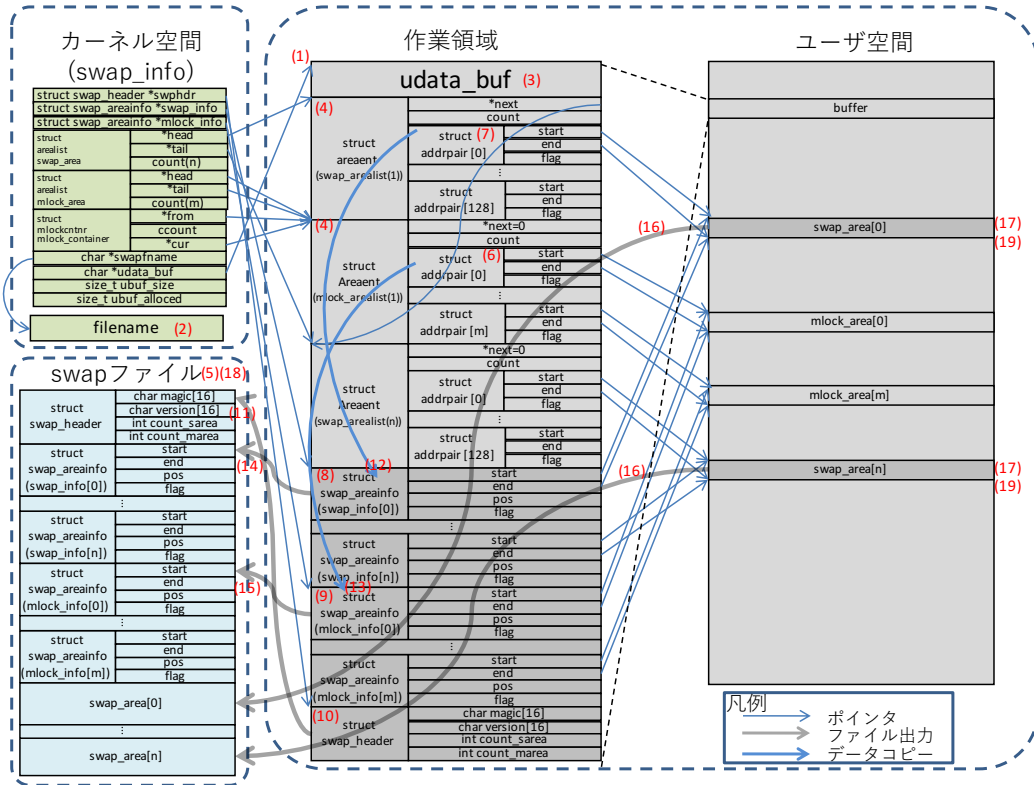


Figure 2.24: スワップアウトの処理フロー

2272

2273 スワップアウトの処理フローを図??を用いて説明する。

2274

2275 1. McKernel のユーザから渡された作業領域を `mlock()` によりロックする。swapout 情報
2276 を管理する `swap_info` 構造体の `udata_buf` メンバに作業領域の先頭アドレスを記録する。(図の (1))

2277

2278 2. 引数で指定されたファイル名を `copy_from_user` でカーネル空間にコピーする。`swap_info`
2279 構造体の `swapfname` メンバにファイル名のアドレスを記録する。(図の (2))

2279

2280 3. 作業領域に汎用バッファ `udata_buf` を割り当てる。(図の (3))

2280

2281 4. 作業領域にスワップエリア管理用リスト `swap_arealist` と `mlock` エリア管理用リスト
2282 `mlock_arealist` の領域を割り当てる。(図の (4))

2282

2283 5. `swap` ファイルを `open()` でオープンする。(図の (5))

2283

2284 6. `lookup_process_memory_range` および `next_process_memory_range` を用いて、ユーザ
2284 プロセスのメモリ領域を検索し、それぞれについて以下を行う。

- 2285 (a) `mlock()` されている領域の開始アドレス、終了アドレス、`flag` を作業領域の `mlock_`
2286 `arealist` に記録する。 (図の (6))
- 2287 (b) `mlock()` されていない領域の開始アドレス、終了アドレス、`flag` を作業領域の
2288 `swap_arealist` に記録する。 (図の (7))
- 2289 7. 作業領域の `swap_arealist` のエントリ数と同数のエントリを持つ `swap_info` 配列を作
2290 業領域に割り当てる。カーネル領域の `swap_info` 構造体の `swap_info` メンバに作業領
2291 域の `swap_info` 配列の先頭アドレスを記録する。 (図の (8))
- 2292 8. 作業領域の `mlock_arealist` のエントリ数と同数のエントリを持つ `mlock_info` 配列を
2293 作業領域に割り当てる。カーネル領域の `swap_info` 構造体の `mlock_info` メンバに作業
2294 領域の `mlock_info` 配列の先頭アドレスを記録する。 (図の (9))
- 2295 9. 作業領域に `swap_header` を割り当てる。カーネル領域の `swap_info` 構造体の `swphdr` メ
2296 ンバに先頭アドレスを記録する。 (図の (10))
- 2297 10. 作業領域の `swap_header` の `magic` メンバに "McKernel swap"、`version` メンバに "0.9.0"、
2298 `count_sarea` メンバに `swap_arealist` のエントリ数、`count_marea` メンバに `mlock_`
2299 `arealist` のエントリ数を記録する。
- 2300 11. 作業領域の `swap_header` を `write()` を用いてスワップファイルへ書き出す。 (図の (11))
- 2301 12. 作業領域の `swap_arealist` のリスト形式データを作業領域の `swap_info` 配列へコピー
2302 する。 (図の (12))
- 2303 13. 作業領域の `mlock_arealist` のリスト形式データを作業領域の `mlock_info` 配列へコ
2304 ピーする。 (図の (13))
- 2305 14. 作業領域の `swap_info` 配列を `write()` を用いてスワップファイルへ書き出す。 (図の
2306 (14))
- 2307 15. 作業領域の `mlock_info` 配列を `write()` を用いてスワップファイルへ書き出す。 (図の
2308 (15))
- 2309 16. 作業領域の `swap_info` の情報を用いて、ユーザプロセスのメモリ領域のうち、スワップ
2310 アウト対象となっているものを `write()` を用いてスワップファイルへ出力する。 (図の
2311 (16))
- 2312 17. スワップアウト対象となっているメモリ領域のうち、McKernel 側でマップされている
2313 ものを `ihk_mc_pt_free_range()` でアンマップする。 (図の (17))
- 2314 18. スワップファイルを `close()` を用いてクローズする。 (図の (18))
- 2315 19. スワップアウト対象となっているメモリ領域のうち、Linux 側でマップされているもの
2316 を `mcexec` に依頼することでアンマップする。 (図の (19))

2317 B. スワップイン処理

2318

2319

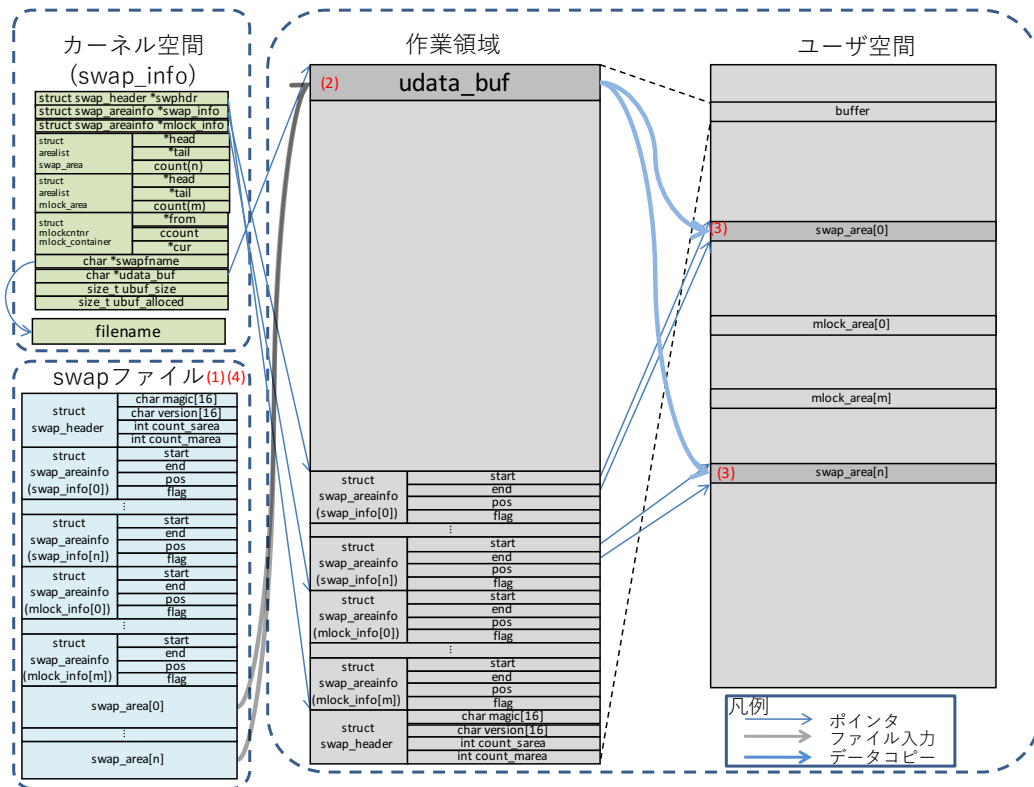


Figure 2.25: スワップインの処理フロー

2320 スワップインの処理フローを図??を用いて説明する。

2321 1. スワップファイルを `open()` を用いてオープンする。 (図の (1))

2322 2. スワップイン対象アドレス範囲を記録している `swap_info` 配列の各エントリに対して
2323 以下を行う。なお、ユーザ空間の作業領域はスワップアウトを経ても残っているため、
2324 `swap_info` 配列をファイルから取得する必要はない。

2325 (a) `read()` を用いてスワップファイルから作業領域の `udata_buf` へスワップイン対象
2326 のメモリ内容をコピーする。 (図の (2))

2327 (b) `copy_to_user` を用いて、作業領域の `udata_buf` からユーザプロセスのメモリ領域
2328 へ、スワップイン対象のメモリ内容をコピーする。 (図の (3))

2329 3. スワップファイルを `close()` を用いてクローズする。 (図の (4))

2330 2.17 Portability

2331 IHK/McKernel has been designed not only for post K computer but also for other manycore
2332 architectures, including Intel Xeon phi. In order to make the source code portable as much
2333 as possible. The following is coding convention of IHK/McKernel.

2334 The directories for architecture dependent and independent source codes are created and
2335 codes are separately stored into those two directories. That is, source codes, including
2336 header files, for some specific architecture are located in its architecture dependent directory.

2337 The source codes, accessing some hardware registers, are hardware specific, and thus
2338 those are machine dependent. Low-level interrupt handlers, some memory management
2339 codes, context switch codes, and signaling codes are the examples. Those source codes are
2340 located in an architecture dependent directory.

2341 Any program code and header files must not include any machine dependent codes
2342 including conditional compile macros, such as `#ifdef ARCH`. As much as possible, we define
2343 machine independent interfaces so that those interfaces are implemented for each architec-
2344 ture.

2345 2.18 Formal Verification

2346 Some of the behaviors of McKernel is verified in a formal way by embedding behaviors in
2347 code and running a verification engine. We employ an extended version of the ANSI/ISO
2348 C Specification Language, whose extensions[?] were developed at the project “Dependable
2349 Operating Systems for Embedded Systems Aiming at Practical Applications” in the research
2350 area named Core Research for Evolutional Science and Technology (CREST), sponsored by
2351 Japan Science and Technology Agency (JST).

2352 2.18.1 Specification Language

2353 The following are expressions defined in the formal specification language. The behavior
2354 of each function is formally specified by using those expressions that are written as C
2355 comments.

2356 `\result` specifies return value.

2357 `\interrupt_disabled` the CPU is interruptable if 0, the CPU is not interruptable if 1 or
2358 more.

2359 `\process_env` the execution is under the user context if 1 or more, the execution is under
2360 the kernel context if 0.

2361 `\atomicity` the execution is not allowed to block if 1 or more, the execution may be
2362 suspended if 0.

2363 `\dont_call_schedule` the context switch is not allowed if 1 or more, the context switch is
2364 allowed if 0.

2365 `is_locked(<pointer variable>)` returns true if a memory block pointed by the pointer
2366 variable is the lock status, otherwise returns false.

2367 `requires <condition expression>` The condition expression must be satisfied at the be-
2368 ginning of the function execution.

2369 `ensures <condition expression>` The condition expression must be satisfied at the end
2370 of the function execution.

2371 `invariant <condition expression>` The condition expression must be satisfied during the
2372 function execution.

2373 Here is a sample code.

```

2374  /*@
2375   @ behavior valid_vector:
2376   @   assumes 32 <= vector <= 255;
2377   @   requires \valid(h);
2378   @   assigns handlers[vector-32];
2379   @   ensures \result == 0;
2380   @ behavior invalid_vector:
2381   @   assumes (vector < 32) || (255 < vector);
2382   @   assigns \nothing;
2383   @   ensures \result == -EINVAL;
2384   @*/
2385  int ihk_mc_register_interrupt_handler(int vector,
2386                                       struct ihk_mc_interrupt_handler *h)
2387  {
2388      if (vector < 32 || vector > 255) {
2389          return -EINVAL;
2390      }
2391      list_add_tail(&h->list, &handlers[vector - 32]);
2392      return 0;
2393  }

```

2394 2.19 Limitations

2395 Certain system calls are only partially implemented in McKernel or not conforming Linux
2396 API. These are either due to design restrictions of the proxy approach or because their
2397 support is intentionally omitted. Table ?? shows the limitations.

Table 2.10: Limitations of McKernel

Function	Description
<code>arch_prctl</code>	It returns the <code>EOPNOTSUPP</code> error when <code>ARCH_SET_GS</code> is passed.
<code>brk</code>	It extends the heap more than <code>requestd</code> when <code>-h (--extend-heap-by=)<step></code> option of <code>mcexec</code> is used with the value larger than 4 KiB.
<code>clone</code>	It supports only the following flags. All other flags cause <code>clone()</code> to return error or are simply ignored. <ul style="list-style-type: none"> • <code>CLONE_CHILD_CLEAR_TID</code> • <code>CLONE_CHILD_SET_TID</code> • <code>CLONE_PARENT_SET_TID</code> • <code>CLONE_SETTLS</code> • <code>CLONE_SIGHAND</code> • <code>CLONE_VM</code>
<code>getrusage</code>	The time spent is measured in a different way than Linux for <code>RUSAGE_THREAD</code> . That is, time spent in user-mode and kernel-mode are updated when CPU mode changes (i.e. when switching from user-mode to kernel-mode and vice versa).
<code>mbind</code>	Per-memory-range policy can be set but it is not used when allocating physical pages.
<code>set_mempolicy</code>	<ul style="list-style-type: none"> • <code>MPOL_F_RELATIVE_NODES</code> and <code>MPOL_INTERLEAVE</code> flags are not supported. • <code>MPOL_BIND</code> works in the same way as <code>MPOL_PREFERRED</code>. That is, <code>MPOL_BIND</code> doesn't return an error when there is no space left in the NUMA nodes specified, but continues to search space in the other nodes.
<code>migrate_pages</code>	It returns the <code>ENOSYS</code> error.
<code>msync</code>	Only the modified pages mapped by the calling process are written back.
<code>setpriority</code> , <code>getpriority</code>	They could set/get the priority of a random <code>mcexec</code> thread. This is because there's no fixed correspondence between a McKernel thread which issues the system call and a <code>mcexec</code> thread which handles the offload request.
<code>set_rlimit</code>	It sets the limit values but they are not enforced.
<code>set_robust_list</code>	It returns the <code>ENOSYS</code> error.
<code>signalfd</code>	It returns the <code>EOPNOTSUPP</code> error.
<code>signalfd4</code>	It returns a <code>fd</code> , but signal is not notified through the <code>fd</code> .
<code>setfsuid</code> , <code>setfsgid</code>	It cannot change the id of the calling thread. Instead, it changes that of the <code>mcexec</code> worker thread which takes the system-call offload request.
<code>mmap</code> (hugeTLBfs)	The physical pages corresponding to a map are released when no McKernel process exist. The next map gets fresh physical pages.
Sticky bit on executable file	It has no effect.
Anonymous shared mapping	Mixing page sizes is not allowed. <code>mmap</code> creates <code>vm_range</code> with one page size. And <code>munmap</code> or <code>mremap</code> that needs the reduced page size changes the sizes of all the pages of the <code>vm_range</code> .
<code>ihk_os_getperfevent</code>	It could time-out when invoked from Fujitsu TCS (job-scheduler).
<code>madvise</code> , <code>mbind</code>	The behaviors of <code>madvise</code> and <code>mbind</code> are changed to do nothing and report success as a workaround for Fugaku.
<code>mmap</code>	It allows unlimited overcommit. Note that it corresponds to setting <code>sysctl vm.overcommit_memory</code> to 1.
<code>mlockall</code>	It is not supported and returns <code>-EPERM</code> .
<code>munlockall</code>	It is not supported and returns zero.

2398 Chapter 3

2399 運用ガイド

2400 本章の想定読者は以下の通り。

- 2401 • McKernel を用いたシステムを運用するシステム管理者

2402 SMP プロセッサ向け、x86_64 アーキ向けの関連ファイルの場所は以下の通り。なお、IHK/McKernel のインストールディレクトリを<install>とする。

インストール先	説明
<install>/kmod/ihk.ko	IHK-master core
<install>/kmod/ihk-smp-x86.ko	IHK-master driver
<install>/kmod/mcctrl.ko	Delegator module
<install>/kmod/mcoverlayfs.ko	/sys, /proc のためのファイルシステム重ね合わせカーネルモジュール
<install>/smp-x86/kernel/mckernel.img	カーネルイメージ

2403 運用向けコマンド・デーモンのファイルの場所は以下の通り。なお、IHK/McKernel の
2404 インストールディレクトリを<install>とする。

インストール先	説明
<install>/sbin/mcreboot.sh	ブートスクリプト
<install>/sbin/mcstop+release.sh	シャットダウンスクリプト
<install>/bin/mcexec	プロセス起動コマンド
<install>/bin/eclair	ダンプ解析ツール
<install>/bin/vmcore2mckdump	ダンプ形式変換ツール

2405 以下、関連コマンドおよび関連関数のインターフェイスを説明する。
2406

2407 3.1 インターフェイス

2408 3.1.1 カーネル引数

2409 McKernel のカーネル引数を表??に示す。

Table 3.1: McKernel のカーネル引数

引数	説明				
<code>hidos</code>	IKC を有効にする。				
<code>dump_level= <dump_level></code>	Linux の panic ハンドラ経由でダンプを行った場合の、ダンプ対象とするメモリ領域の種類を<dump_level>に設定する。設定可能な値は以下の通り。 <table border="1" data-bbox="516 380 1357 443"> <tr> <td>0</td> <td>IHK が McKernel に割り当てたメモリ領域を出力する。</td> </tr> <tr> <td>24</td> <td>カーネルが使用しているメモリ領域を出力する。</td> </tr> </table> 指定がなかった場合は 24 が用いられる。	0	IHK が McKernel に割り当てたメモリ領域を出力する。	24	カーネルが使用しているメモリ領域を出力する。
0	IHK が McKernel に割り当てたメモリ領域を出力する。				
24	カーネルが使用しているメモリ領域を出力する。				
<code>allow_oversubscribe</code>	McKernel に割り当てられた CPU 数より大きい数のスレッドまたはプロセスの生成を許可する。この引数が指定されない場合に、CPU 数より大きい数のスレッドまたはプロセスを <code>clone()</code> , <code>fork()</code> , <code>vfork()</code> など生成しようとすると、当該システムコールが <code>EINVAL</code> エラーを返す。				

2410 3.1.2 ブートスクリプト

2411 書式

```
2412 mcreboot.sh [-c <cpulist>] [-r <ikcmap>] [-m <memlist>] [-f <facility>] [-o
2413 <chownopt>] [-i <mon_interval>] [-k <redirct_kmsg>] [-q <irq>] [-t] [-d <dump_level>]
2414 [-0]
```

2415 オプション

2416

オプション	説明				
-c <cpulist>	McKernel に割り当てる CPU のリストを指定する。フォーマットは以下の通り。<CPU logical id>[,<CPU logical id>...] または<CPU logical id>-<CPU logical id>[,<CPU logical id>-<CPU logical id>...] または両者の混合。				
-r <ikcmap>	McKernel の CPU が IKC メッセージを送る Linux CPU を指定する。フォーマットは以下の通り。<CPU list>:<CPU logical id>[+<CPU list>:<CPU logical id>...] <CPU list>のフォーマットは-c オプションにおけるものと同じである。各<CPU list>:<CPU logical id>は<CPU list>で示される McKernel の CPU が<CPU logical id>で示される Linux の CPU に IKC メッセージを送信することを意味する。				
-m <memlist>	McKernel に割り当てるメモリ領域を指定する。フォーマットは以下の通り。<サイズ> @ <NUMA-node 番号> [, <サイズ> @ <NUMA-node 番号> ...]。				
-f <facility>	ihkmond が使用する syslog プロトコルの facility を指定する。デフォルトは LOG_LOCAL6。				
-o <chownopt>	IHK のデバイスファイル (/dev/mcd*, /dev/mcos*) のオーナーとグループの値を<user>[:<group>] の形式で指定する。デフォルトは mcreboot.sh を実行したユーザ。				
-i <mon_interval>	ihkmond がハングアップ検知のために OS 状態を確認する時間間隔を秒単位で指定する。-1 が指定された場合はハングアップ検知を行わない。指定がない場合はハングアップ検知を行わない。				
-k <redirect_kmsg>	カーネルメッセージの/dev/log へのリダイレクト有無を指定する。0 が指定された場合はリダイレクトを行わず、0 以外が指定された場合はリダイレクトを行う。指定がない場合はリダイレクトを行わない。				
-q <irq>	IHK が使用する IRQ 番号を指定する。指定がない場合は 64-255 の範囲で空いているものを使用する。				
-t	(x86_64 アーキテクチャ固有) Turbo Boost をオンにする。デフォルトはオフ。				
-d <dump_level>	Linux の panic ハンドラ経由でダンプを行った場合の、ダンプ対象とするメモリ領域の種類を<dump_level>に設定する。設定可能な値は以下の通り。 <table border="1" data-bbox="516 1031 1357 1094"> <tbody> <tr> <td>0</td> <td>IHK が McKernel に割り当てたメモリ領域を出力する。</td> </tr> <tr> <td>24</td> <td>カーネルが使用しているメモリ領域を出力する。</td> </tr> </tbody> </table> 指定がなかった場合は 24 が用いられる。	0	IHK が McKernel に割り当てたメモリ領域を出力する。	24	カーネルが使用しているメモリ領域を出力する。
0	IHK が McKernel に割り当てたメモリ領域を出力する。				
24	カーネルが使用しているメモリ領域を出力する。				
-0	McKernel に割り当てられた CPU 数より大きい数のスレッドまたはプロセスの生成を許可する。指定がない場合は許可しない。すなわち、CPU 数より大きい数のスレッドまたはプロセスを生成しようとするとエラーとなる。				

2417 説明

2418

2419 McKernel 関連カーネルモジュールを insmod し、<cpulist>で指定された CPU と<memlist>で
2420 指定されたメモリ領域からなるパーティションを作成し、IKC map を<ikcmap>に設定し、前
2421 記パーティションに McKernel をブートする。

2422 戻り値

0	正常終了
0 以外	エラー

2423 3.1.3 シャットダウンスクリプト

2424 書式

2425 mcstop+release.sh

2426 オプション

2427 なし

2428 説明

2429 McKernel をシャットダウンし、McKernel 用パーティションを削除し、関連カーネルモ
2430 ジュールを rmmmod する。

2431 戻り値

0	正常終了
0 以外	エラー

2432 3.1.4 プロセス起動コマンド

2433 インターフェイスは第??節に記載する。

2434 3.1.5 統計情報取得

2435 バッチジョブスケジューラは、IHK の関数 `ihk_os_getrusage()` を呼ぶことでジョブの統計
2436 情報を取得できる（インターフェイスは”IHK Specifications” 参照）。

2437 `ihk_os_getrusage()` は `void *rusage` という引数で結果を返す。McKernel では `rusage`
2438 の実際の型は `struct mckernel_rusage` 型で、以下のように定義される。

```
2439 struct mckernel_rusage {
2440     unsigned long memory_stat_rss[IHK_MAX_NUM_PGSIIZES];
2441     /* ユーザのページサイズごとの anonymous ページ使用量現在値 (バイト単位) */
2442     unsigned long memory_stat_mapped_file[IHK_MAX_NUM_PGSIIZES];
2443     /* ユーザのページサイズごとの file-backed ページ使用量現在値 (バイト単位) */
2444     unsigned long memory_max_usage;
2445     /* ユーザのメモリ使用量最大値 (バイト単位) */
2446     unsigned long memory_kmem_usage;
2447     /* カーネルのメモリ使用量現在値 (バイト単位) */
2448     unsigned long memory_kmem_max_usage;
2449     /* カーネルのメモリ使用量最大値 (バイト単位) */
2450     unsigned long memory_numa_stat[IHK_MAX_NUM_NUMA_NODES];
2451     /* NUMA ごとのユーザのメモリ使用量現在値 (バイト単位) */
2452     unsigned long cpuacct_stat_system;
2453     /* システム時間 (USER_HZ 単位) */
2454     unsigned long cpuacct_stat_user;}
2455     /* ユーザ時間 (USER_HZ 単位) */
2456     unsigned long cpuacct_usage;}
2457     /* ユーザの CPU 時間 (ナノ秒単位) */
2458     unsigned long cpuacct_usage_percpu[IHK_MAX_NUM_CPUS];
2459     /* コアごとのユーザの CPU 時間 (ナノ秒単位) */
2460     int num_threads;
2461     /* スレッド数現在値 */
2462     int max_num_threads;
2463     /* スレッド数最大値 */
2464 };
```

2465 `memory_stat_rss` および `memory_stat_mapped_file` のインデックスはサイズによるペー
2466 ジ種であり、x86_64 アーキでは以下のように定義される。

```
2467 #define IHK_OS_PGFSIZE_4KB 0
2468 #define IHK_OS_PGFSIZE_2MB 1
2469 #define IHK_OS_PGFSIZE_1GB 2
```

2470 3.1.6 ダンプ解析コマンド

2471 インターフェイスは第??節に記載する。

2472 3.1.7 ダンプ形式変換コマンド

2473 インターフェイスは第??節に記載する。

2474 3.2 ブート手順

2475 `mcreboot.sh` を用いてブート手順を説明する。
スクリプトは以下の通り。

```
1  #!/bin/bash
2
3  # IHK SMP-x86 example boot script.
4  # author: Balazs Gerofi <bgerofi@riken.jp>
5  #       Copyright (C) 2014 RIKEN AICS
6  #
7  # This is an example script for loading IHK, configuring a partition and
8  # booting McKernel on it. Unless specific CPUs and memory are requested,
9  # the script reserves half of the CPU cores and 512MB of RAM from
10 # NUMA node 0 when IHK is loaded for the first time.
11 # Otherwise, it destroys the current McKernel instance and reboots it using
12 # the same set of resources as it used previously.
13 # Note that the script does not output anything unless an error occurs.
14
15 prefix="/home/takagi/project/os/install"
16 BINDIR="${prefix}/bin"
17 SBINDIR="${prefix}/sbin"
18 ETCDIR=/home/takagi/project/os/install/etc
19 KMODDIR="${prefix}/kmod"
20 KERNDIR="${prefix}/smp-x86/kernel"
21 ENABLEMCOVERLAYFS="yes"
22
23 mem="512M@0"
24 cpus=""
25 ikc_map=""
26
27 if [ "${BASH_VERSINFO[0]}" -lt 4 ]; then
28     echo "You need at least bash-4.0 to run this script." >&2
29     exit 1
30 fi
31
32 redirect_kmsg=0
33 mon_interval="-1"
34 DUMPLEVEL=24
35 facility="LOG_LOCAL6"
36 chown_option='logname 2> /dev/null '
37
38 if [ "${systemctl status irqbalance_mck.service 2> /dev/null |grep -E 'Active: active' |}" \
39     != "" -o "${systemctl status irqbalance.service 2> /dev/null |grep -E 'Active: active' |}" \
40     != "" ]; then
41     irqbalance_used="yes"
42 else
43     irqbalance_used="no"
44 fi
45
```

```

46 turbo=""
47 ihk_irq=""
48
49 while getopts :tk:c:m:o:f:r:q:i:d: OPT
50 do
51     case ${OPT} in
52         f)     facility=${OPTARG}
53             ;;
54         o)     chown_option=${OPTARG}
55             ;;
56         k)     redirect_kmsg=${OPTARG}
57             ;;
58         c)     cpus=${OPTARG}
59             ;;
60         m)     mem=${OPTARG}
61             ;;
62         r)     ikc_map=${OPTARG}
63             ;;
64         q)     ihk_irq=${OPTARG}
65             ;;
66         t)     turbo="turbo"
67             ;;
68         d)     DUMP_LEVEL=${OPTARG}
69             ;;
70         i)     mon_interval=${OPTARG}
71             ;;
72         *)     echo "invalid option -${OPTARG}" >&2
73             exit 1
74     esac
75 done
76
77 # Start ihkmond
78 pid='pidof ihkmond'
79 if [ "${pid}" != "" ]; then
80     sudo kill -9 ${pid} > /dev/null 2> /dev/null
81 fi
82 if [ "${redirect_kmsg}" != "0" -o "${mon_interval}" != "-1" ]; then
83     ${SBINDIR}/ihkmond -f ${facility} -k ${redirect_kmsg} -i ${mon_interval}
84 fi
85 #
86 # Revert any state that has been initialized before the error occurred.
87 #
88 error_exit() {
89     local status=$1
90
91     case $status in
92         mcos_sys_mounted)
93             if [ "$enable_mcoverlay" = "yes" ]; then
94                 umount /tmp/mcos/mcos0_sys
95             fi
96             ;&
97         mcos_proc_mounted)
98             if [ "$enable_mcoverlay" = "yes" ]; then
99                 umount /tmp/mcos/mcos0_proc
100            fi
101            ;&
102        mcoverlayfs_loaded)
103            if [ "$enable_mcoverlay" = "yes" ]; then
104                rmmod mcoverlay 2>/dev/null
105            fi
106            ;&
107        linux_proc_bind_mounted)
108            if [ "$enable_mcoverlay" = "yes" ]; then
109                umount /tmp/mcos/linux_proc
110            fi
111            ;&
112        tmp_mcos_mounted)
113            if [ "$enable_mcoverlay" = "yes" ]; then

```

```

114             umount /tmp/mcos
115         fi
116     ;&
117 tmp_mcos_created)
118     if [ "$enable_mcoverlay" = "yes" ]; then
119         rm -rf /tmp/mcos
120     fi
121     ;&
122 os_created)
123     # Destroy all LWK instances
124     if ls /dev/mcos* 1>/dev/null 2>&1; then
125         for i in /dev/mcos*; do
126             ind='echo $i|cut -c10-';
127             if ! ${SBINDIR}/ihkconfig 0 destroy $ind; then
128                 echo "warning: failed to destroy LWK instance $ind" >&2
129             fi
130         done
131     fi
132     ;&
133 mcctrl_loaded)
134     rmmod mcctrl 2>/dev/null || echo "warning: failed to remove mcctrl" >&2
135     ;&
136 cpus_reserved)
137     cpus='${SBINDIR}/ihkconfig 0 query cpu'
138     if [ "${cpus}" != "" ]; then
139         if ! ${SBINDIR}/ihkconfig 0 release cpu $cpus > /dev/null; then
140             echo "warning: failed to release CPUs" >&2
141         fi
142     fi
143     ;&
144 mem_reserved)
145     mem='${SBINDIR}/ihkconfig 0 query mem'
146     if [ "${mem}" != "" ]; then
147         if ! ${SBINDIR}/ihkconfig 0 release mem $mem > /dev/null; then
148             echo "warning: failed to release memory" >&2
149         fi
150     fi
151     ;&
152 ihk_smp_loaded)
153     rmmod ihk_smp_x86 2>/dev/null || echo "warning: failed to remove ihk_smp_x86" >&2
154     ;&
155 ihk_loaded)
156     rmmod ihk 2>/dev/null || echo "warning: failed to remove ihk" >&2
157     ;&
158 irqbalance_stopped)
159     if [ "$(systemctl status irqbalance_mck.service 2> /dev/null |'\
160 'grep -E 'Active: active' " != "" ); then
161         if ! systemctl stop irqbalance_mck.service 2>/dev/null; then
162             echo "warning: failed to stop irqbalance_mck" >&2
163         fi
164         if ! systemctl disable irqbalance_mck.service >/dev/null 2>/dev/null; then
165             echo "warning: failed to disable irqbalance_mck" >&2
166         fi
167         if ! etcdir=/home/takagi/project/os/install/etc perl -e \
168 '$etcdir=${ENV}'{etcdir}'; @files = grep { -f } glob "$etcdir/proc/irq/*/smp-affinity";'\
169 ' foreach $file (@files) { $dest = substr($file, length($etcdir));'\
170 ' if(0) {print "cp $file $dest\n";} system("cp $file $dest 2>/dev/null"); }'; then
171             echo "warning: failed to restore /proc/irq/*/smp-affinity" >&2
172         fi
173         if ! systemctl start irqbalance.service; then
174             echo "warning: failed to start irqbalance" >&2;
175         fi
176     fi
177     ;&
178 initial)
179     # Nothing more to revert
180     ;;
181 esac

```

```

182
183         exit 1
184     }
185
186     ihk_ikc_irq_core=0
187
188     release='uname -r'
189     major='echo ${release} | sed -e 's/^\([0-9]*\)*/\1/'
190     minor='echo ${release} | sed -e 's/^[0-9]*\.\([0-9]*\)*/\1/'
191     patch='echo ${release} | sed -e 's/^[0-9]*\.[0-9]*\.\([0-9]*\)*/\1/'
192     linux_version_code='expr \( ${major} \* 65536 \) + \( ${minor} \* 256 \) + ${patch}'
193     rhel_release='echo ${release} | sed -e 's/^[0-9]*\.[0-9]*\.[0-9]*-\([0-9]*\)*/\1/'
194     if [ "${release}" = "${rhel_release}" ]; then
195         rhel_release=""
196     fi
197
198     enable_mcoverlay="no"
199
200     if [ "${ENABLEMCOVERLAYFS}" = "yes" ]; then
201         if [ "${rhel_release}" = "" ]; then
202             if [ ${linux_version_code} -ge 262144 -a ${linux_version_code} -lt 262400 ]; then
203                 enable_mcoverlay="yes"
204             fi
205             if [ ${linux_version_code} -ge 263680 -a ${linux_version_code} -lt 263936 ]; then
206                 enable_mcoverlay="yes"
207             fi
208         else
209             if [ ${linux_version_code} -eq 199168 -a ${rhel_release} -ge 327 -a ${rhel_release} -le 327 ]; then
210                 enable_mcoverlay="yes"
211             fi
212             if [ ${linux_version_code} -ge 262144 -a ${linux_version_code} -lt 262400 ]; then
213                 enable_mcoverlay="yes"
214             fi
215         fi
216     fi
217
218     # Figure out CPUs if not requested by user
219     if [ "$cpus" = "" ]; then
220         # Get the number of CPUs on NUMA node 0
221         nr_cpus='lscpu --parse | awk -F" ," '{if ($4 == 0) print $4}' | wc -l'
222
223         # Use the second half of the cores
224         let nr_cpus="$nr_cpus / 2"
225         cpus='lscpu --parse | awk -F" ," '{if ($4 == 0) print $1}' | tail -n $nr_cpus | \
226 ' xargs echo -n | sed 's/ /,/g'
227         if [ "$cpus" = "" ]; then
228             echo "error: no available CPUs on NUMA node 0?" >&2
229             exit 1
230         fi
231     fi
232
233     # Remove mcoverlay if loaded
234     if [ "$enable_mcoverlay" = "yes" ]; then
235         if grep mcoverlay /proc/modules &>/dev/null; then
236             if [ "$(cat /proc/mounts | grep /tmp/mcos/mcos0_sys)" != "" ]; \
237 then umount -l /tmp/mcos/mcos0_sys; fi
238             if [ "$(cat /proc/mounts | grep /tmp/mcos/mcos0_proc)" != "" ]; \
239 then umount -l /tmp/mcos/mcos0_proc; fi
240             if [ "$(cat /proc/mounts | grep /tmp/mcos/linux_proc)" != "" ]; \
241 then umount -l /tmp/mcos/linux_proc; fi
242             if [ "$(cat /proc/mounts | grep /tmp/mcos)" != "" ]; then umount -l /tmp/mcos; fi
243             if [ -e /tmp/mcos ]; then rm -rf /tmp/mcos; fi
244             if ! rmmmod mcoverlay 2>/dev/null; then
245                 echo "error: removing mcoverlay" >&2
246                 exit 1
247             fi
248         fi
249     fi

```



```

250
251 # Stop irqbalance
252 if [ "${irqbalance_used}" == "yes" ]; then
253     systemctl stop irqbalance_mck.service 2>/dev/null
254     if ! systemctl stop irqbalance.service 2>/dev/null ; then
255         echo "error: stopping irqbalance" >&2
256         exit 1
257     fi;
258
259     if ! etcdir=/home/takagi/project/os/install/etc perl -e \
260 'use File::Copy qw(copy); $etcdir=$ENV{'etcdir'}; '\
261 '@files = grep { -f } glob "/proc/irq/*/smp_affinity"; foreach $file (@files) { '\
262 '$rel = substr($file, 1); $dir=substr($rel, 0, length($rel)-length("/smp_affinity")); '\
263 'if(0) { print "cp $file $etcdir/$rel\n"; } if(system("mkdir -p $etcdir/$dir")){ exit 1;} '\
264 'if(!copy($file, "$etcdir/$rel")){ exit 1;} }'; then
265     echo "error: saving /proc/irq/*/smp_affinity" >&2
266     error_exit "mcos_sys_mounted"
267 fi;
268
269 # Prevent /proc/irq/*/smp_affinity from getting zero after offlining
270 # McKernel CPUs by using the following algorithm.
271 # if (smp_affinity & mck_cores) {
272 #     smp_affinity = (mck_cores ^ -1);
273 # }
274 ncpus=`lscpu | grep -E '^CPU\(s\):' | awk '{print $2}'`
275 smp_affinity_mask=`echo $cpus | ncpus=$ncpus perl -e \
276 'while(<>){@tokens = split /,/;foreach $token (@tokens) '\
277 '@nums = split /-/, $token; for($num = $nums[0]; $num <= $nums[$#nums]; $num++) {'\
278 '$ndx=int($num/32); $mask[$ndx] |= (1<<($num % 32))}}'\
279 '$nint32s = int(($ENV{'ncpus'}+31)/32); for($j = $nint32s - 1; $j >= 0; $j--) {'\
280 ' if($j != $nint32s - 1){print ",";}'\
281 '$nblks = ($j != $nint32s - 1) ? 8 : ($ENV{'ncpus'} % 32 != 0) ? '\
282 'int(($ENV{'ncpus'} + 3) % 32) / 4) : 8;'\
283 ' for($i = $nblks - 1;$i >= 0;$i--){ printf("%01x",($mask[$j] >> ($i*4)) & 0xf);}'\
284 #     echo cpus=$cpus ncpus=$ncpus smp_affinity_mask=$smp_affinity_mask
285
286     if ! ncpus=$ncpus smp_affinity_mask=$smp_affinity_mask perl -e \
287 '@dirs = grep { -d } glob "/proc/irq/*"; foreach $dir (@dirs) {'\
288 '$hit = 0; $affinity_str = `cat $dir/smp_affinity`; chomp $affinity_str;'\
289 '@int32strs = split /,/, $affinity_str; @int32strs_mask=split /,/, $ENV{'smp_affinity_mask'};'\
290 ' for($i=0;$i <= $#int32strs_mask; $i++) {'\
291 '$int32strs_inv[$i] = sprintf("%08x",hex($int32strs_mask[$i])^0xffffffff);'\
292 ' if($i == 0) { $len = int(((($ENV{'ncpus'}%32)+3)/4); if($len != 0) {'\
293 '$int32strs_inv[$i] = substr($int32strs_inv[$i], -$len, $len); } } }'\
294 '$inv = join(",", @int32strs_inv); $nint32s = int(($ENV{'ncpus'}+31)/32);'\
295 ' for($j = $nint32s - 1; $j >= 0; $j--) {'\
296 ' if(hex($int32strs[$nint32s - 1 - $j]) & hex($int32strs_mask[$nint32s - 1 - $j])) {'\
297 '$hit = 1; } } if($hit == 1) {'\
298 '$cmd = "echo $inv > $dir/smp_affinity 2>/dev/null"; system $cmd;}'\
299     echo "error: modifying /proc/irq/*/smp_affinity" >&2
300     error_exit "mcos_sys_mounted"
301 fi
302
303 fi
304
305 # Load IHK if not loaded
306 if ! grep -E 'ihk\s' /proc/modules &>/dev/null; then
307     if ! taskset -c 0 insmod ${KMODDIR}/ihk.ko 2>/dev/null; then
308         echo "error: loading ihk" >&2
309         error_exit "irqbalance_stopped"
310     fi
311 fi
312
313 # Increase swappiness so that we have better chance to allocate memory for IHK
314 echo 100 > /proc/sys/vm/swappiness
315
316 # Drop Linux caches to free memory
317 sync && echo 3 > /proc/sys/vm/drop_caches

```

```

318
319 # Merge free memory areas into large , physically contiguous ones
320 echo 1 > /proc/sys/vm/compact_memory 2>/dev/null
321
322 sync
323
324 # Load IHK-SMP if not loaded and reserve CPUs and memory
325 if ! grep ihk_smp-x86 /proc/modules &>/dev/null; then
326     if [ "$ihk_irq" = "" ]; then
327         for i in `seq 64 255`; do
328             if [ ! -d /proc/irq/$i ] && \
329 [ "`cat /proc/interrupts | grep ":" | awk '{print $1}' | grep -o '[0-9]*' | grep -e '^$i$'`\
330 = "" ]; then
331                 ihk_irq=$i
332                 break
333             fi
334         done
335     if [ "$ihk_irq" = "" ]; then
336         echo "error: no IRQ available" >&2
337         error_exit "ihk_loaded"
338     fi
339 fi
340 if ! taskset -c 0 insmod ${KMODDIR}/ihk-smp-x86.ko ihk_start_irq=$ihk_irq\
341 ihk_ikc_irq_core=$ihk_ikc_irq_core 2>/dev/null; then
342     echo "error: loading ihk-smp-x86" >&2
343     error_exit "ihk_loaded"
344 fi
345
346 # Offline-reonline RAM (special case for OFP SNC-4 mode)
347 if [ "`hostname | grep "c[0-9][0-9][0-9][0-9].ofp" != "" ] && [ "`cat /sys/devices/system/nod
348     for i in 0 1 2 3; do
349         find /sys/devices/system/node/node$i/memory*/ -name "online" |\
350 while read f; do
351             echo 0 > $f 2>&1 > /dev/null;
352         done
353     find /sys/devices/system/node/node$i/memory*/ -name "online" |\
354 while read f; do
355         echo 1 > $f 2>&1 > /dev/null;
356     done
357     done
358     for i in 4 5 6 7; do
359         find /sys/devices/system/node/node$i/memory*/ -name "online" |\
360 while read f; do
361             echo 0 > $f 2>&1 > /dev/null;
362         done
363     find /sys/devices/system/node/node$i/memory*/ -name "online" |\
364 while read f; do
365         echo 1 > $f 2>&1 > /dev/null;
366     done
367     done
368 fi
369
370 if ! ${SBINDIR}/ihkconfig 0 reserve mem ${mem}; then
371     echo "error: reserving memory" >&2
372     error_exit "ihk_smp_loaded"
373 fi
374 if ! ${SBINDIR}/ihkconfig 0 reserve cpu ${cpus}; then
375     echo "error: reserving CPUs" >&2;
376     error_exit "mem_reserved"
377 fi
378 fi
379
380 # Load mcctrl if not loaded
381 if ! grep mcctrl /proc/modules &>/dev/null; then
382     if ! taskset -c 0 insmod ${KMODDIR}/mcctrl.ko 2>/dev/null; then
383         echo "error: inserting mcctrl.ko" >&2
384         error_exit "cpus_reserved"
385     fi

```

```

386 fi
387
388 # Destroy all LWK instances
389 if ls /dev/mcos* 1>/dev/null 2>&1; then
390     for i in /dev/mcos*; do
391         ind='echo $i|cut -c10-';
392         # Retry when conflicting with ihkmond
393         nretry=0
394         until ${SBINDIR}/ihkconfig 0 destroy $ind || [ $nretry -lt 4 ]; do
395             sleep 0.25
396             nretry=$(( $nretry + 1 ))
397         done
398         if [ $nretry -eq 4 ]; then
399             echo "error: destroying LWK instance $ind failed" >&2
400             error_exit "mcctrl_loaded"
401         fi
402     done
403 fi
404
405 # Create OS instance
406 if ! ${SBINDIR}/ihkconfig 0 create; then
407     echo "error: creating OS instance" >&2
408     error_exit "mcctrl_loaded"
409 fi
410
411 # Assign CPUs
412 if ! ${SBINDIR}/ihkosctl 0 assign cpu ${cpus}; then
413     echo "error: assign CPUs" >&2
414     error_exit "os_created"
415 fi
416
417 if [ "$ikc_map" != "" ]; then
418     # Specify IKC map
419     if ! ${SBINDIR}/ihkosctl 0 set ikc_map ${ikc_map}; then
420         echo "error: assign CPUs" >&2
421         error_exit "os_created"
422     fi
423 fi
424
425 # Assign memory
426 if ! ${SBINDIR}/ihkosctl 0 assign mem ${mem}; then
427     echo "error: assign memory" >&2
428     error_exit "os_created"
429 fi
430
431 # Load kernel image
432 if ! ${SBINDIR}/ihkosctl 0 load ${KERNDIR}/mckernel.img; then
433     echo "error: loading kernel image: ${KERNDIR}/mckernel.img" >&2
434     error_exit "os_created"
435 fi
436
437 # Set kernel arguments
438 if ! ${SBINDIR}/ihkosctl 0 kargs "hidos $turbo dump_level=${DUMP_LEVEL}"; then
439     echo "error: setting kernel arguments" >&2
440     error_exit "os_created"
441 fi
442
443 # Boot OS instance
444 if ! ${SBINDIR}/ihkosctl 0 boot; then
445     echo "error: booting" >&2
446     error_exit "os_created"
447 fi
448
449 # Set device file ownership
450 if ! chown ${chown_option} /dev/mcd* /dev/mcos*; then
451     echo "warning: failed to chown device files" >&2
452 fi
453

```

```

454 # Overlay /proc, /sys with McKernel specific contents
455 if [ "$enable_mcoverlay" = "yes" ]; then
456     if [ ! -e /tmp/mcos ]; then
457         mkdir -p /tmp/mcos;
458     fi
459     if ! mount -t tmpfs tmpfs /tmp/mcos; then
460         echo "error: mount /tmp/mcos" >&2
461         error_exit "tmp_mcos_created"
462     fi
463     if [ ! -e /tmp/mcos/linux_proc ]; then
464         mkdir -p /tmp/mcos/linux_proc;
465     fi
466     if ! mount --bind /proc /tmp/mcos/linux_proc; then
467         echo "error: mount /tmp/mcos/linux_proc" >&2
468         error_exit "tmp_mcos_mounted"
469     fi
470     if ! taskset -c 0 insmod ${KMODDIR}/mcoverlay.ko 2>/dev/null; then
471         echo "error: inserting mcoverlay.ko" >&2
472         error_exit "linux_proc_bind_mounted"
473     fi
474     while [ ! -e /proc/mcos0 ]
475     do
476         sleep 0.1
477     done
478     if [ ! -e /tmp/mcos/mcos0_proc ]; then
479         mkdir -p /tmp/mcos/mcos0_proc;
480     fi
481     if [ ! -e /tmp/mcos/mcos0_proc_upper ]; then
482         mkdir -p /tmp/mcos/mcos0_proc_upper;
483     fi
484     if [ ! -e /tmp/mcos/mcos0_proc_work ]; then
485         mkdir -p /tmp/mcos/mcos0_proc_work;
486     fi
487     if ! mount -t mcoverlay mcoverlay -o\
488     lowerdir=/proc/mcos0:/proc,upperdir=/tmp/mcos/mcos0_proc_upper,\
489     workdir=/tmp/mcos/mcos0_proc_work,nocopyupw,nofscheck /tmp/mcos/mcos0_proc; then
490         echo "error: mounting /tmp/mcos/mcos0_proc" >&2
491         error_exit "mcoverlayfs_loaded"
492     fi
493     # TODO: How de we revert this in case of failure??
494     mount --make-rprivate /proc
495
496     while [ ! -e /sys/devices/virtual/mcos/mcos0/sys/setup_complete ]
497     do
498         sleep 0.1
499     done
500     if [ ! -e /tmp/mcos/mcos0_sys ]; then
501         mkdir -p /tmp/mcos/mcos0_sys;
502     fi
503     if [ ! -e /tmp/mcos/mcos0_sys_upper ]; then
504         mkdir -p /tmp/mcos/mcos0_sys_upper;
505     fi
506     if [ ! -e /tmp/mcos/mcos0_sys_work ]; then
507         mkdir -p /tmp/mcos/mcos0_sys_work;
508     fi
509     if ! mount -t mcoverlay mcoverlay -o\
510     lowerdir=/sys/devices/virtual/mcos/mcos0/sys:/sys,upperdir=/tmp/mcos/mcos0_sys_upper,\
511     workdir=/tmp/mcos/mcos0_sys_work,nocopyupw,nofscheck /tmp/mcos/mcos0_sys; then
512         echo "error: mount /tmp/mcos/mcos0_sys" >&2
513         error_exit "mcos_proc_mounted"
514     fi
515     # TODO: How de we revert this in case of failure??
516     mount --make-rprivate /sys
517
518     touch /tmp/mcos/mcos0_proc/mckernel
519
520     rm -rf /tmp/mcos/mcos0_sys/setup_complete
521

```

```

522     # Hide NUMA related files which are outside the LWK partition
523     for cpuid in \
524 'find /sys/devices/system/cpu/* -maxdepth 0 -name "cpu[0123456789]*" -printf "%f "'; do
525         if [ ! -e "/sys/devices/virtual/mcos/mcos0/sys/devices/system/cpu/$cpuid" ]; then
526             rm -rf /tmp/mcos/mcos0_sys/devices/system/cpu/$cpuid
527             rm -rf /tmp/mcos/mcos0_sys/bus/cpu/devices/$cpuid
528             rm -rf /tmp/mcos/mcos0_sys/bus/cpu/drivers/processor/$cpuid
529         else
530             for nodeid in \
531 'find /sys/devices/system/cpu/$cpuid/* -maxdepth 0 -name "node[0123456789]*" -printf "%f "'; do
532                 if [ ! -e \
533 "/sys/devices/virtual/mcos/mcos0/sys/devices/system/cpu/$cpuid/$nodeid" ]; then
534                     rm -f \
535 /tmp/mcos/mcos0_sys/devices/system/cpu/$cpuid/$nodeid
536                 fi
537             done
538         fi
539     done
540     for nodeid in \
541 'find /sys/devices/system/node/* -maxdepth 0 -name "node[0123456789]*" -printf "%f "'; do
542         if [ ! -e "/sys/devices/virtual/mcos/mcos0/sys/devices/system/node/$nodeid" ]; \
543 then
544             rm -rf /tmp/mcos/mcos0_sys/devices/system/node/$nodeid/*
545             rm -rf /tmp/mcos/mcos0_sys/bus/node/devices/$nodeid
546         else
547             # Delete non-existent symlinks
548             for cpuid in \
549 'find /sys/devices/system/node/$nodeid/* -maxdepth 0 -name "cpu[0123456789]*" -printf "%f "'; do
550                 if [ ! -e \
551 "/sys/devices/virtual/mcos/mcos0/sys/devices/system/node/$nodeid/$cpuid" ]; then
552                     rm -f \
553 /tmp/mcos/mcos0_sys/devices/system/node/$nodeid/$cpuid
554                 fi
555             done
556             rm -f /tmp/mcos/mcos0_sys/devices/system/node/$nodeid/memory*
557         fi
558     done
559     rm -f /tmp/mcos/mcos0_sys/devices/system/node/has_*
560     for cpuid in \
561 'find /sys/bus/cpu/devices/* -maxdepth 0 -name "cpu[0123456789]*" -printf "%f "'; do
562         if [ ! -e "/sys/devices/virtual/mcos/mcos0/sys/bus/cpu/devices/$cpuid" ]; then
563             rm -rf /tmp/mcos/mcos0_sys/bus/cpu/devices/$cpuid
564         fi
565     done
566 fi
567
568 # Start irqbalance with CPUs and IRQ for McKernel banned
569 if [ "${irqbalance_used}" = "yes" ]; then
570     banirq='cat /proc/interrupts | \
571 perl -e 'while(<>) { if (/^\s*(\d+).*IHK\--SMP\s*$/) {print $1;}}'
572     sed "s/%mask%/$smp_affinity_mask/g" $ETCDIR/irqbalance.mck.in | \
573 sed "s/%banirq%/$banirq/g" > /tmp/irqbalance_mck
574     systemctl disable irqbalance_mck.service >/dev/null 2>/dev/null
575     if ! systemctl link $ETCDIR/irqbalance.mck.service >/dev/null 2>/dev/null; then
576         echo "error: linking irqbalance_mck" >&2
577         error_exit "mcos-sys-mounted"
578     fi
579 fi
580
581 if ! systemctl start irqbalance_mck.service 2>/dev/null ; then
582     echo "error: starting irqbalance_mck" >&2
583     error_exit "mcos-sys-mounted"
584 fi
585 # echo cpus=$cpus ncpus=$ncpus banirq=$banirq
586 fi

```

2476 手順は以下の通り。

- 2477 1. `ihkmond` を起動する。`ihkmond` は任意のタイミングで起動してよい。これは、`ihkmond`
2478 は OS インスタンスの作成を検知して動作を開始するためである。(83 行目)
- 2479 2. Linux のカーネルバージョンが、`mcoverlayfs` が動作するものであるかを確認する。
2480 (200–216 行目)
- 2481 3. `irqbalance` を停止する。(251–257 行目)
- 2482 4. `/proc/irq/*/affinity` の設定を保存した上で McKernel CPU を担当から外す。担当
2483 CPU が無くなる場合は、全ての Linux CPU を指定する。(269–303 行目)
- 2484 5. `ihk.ko` を `insmod` する。(307 行目)
- 2485 6. Linux によるメモリフラグメンテーションを緩和するために以下を実施する。(313–320
2486 行目)
- 2487 (a) アクティブでないプロセスを積極的にスワップアウトするように設定する
- 2488 (b) クリーンなページキャッシュを無効化し、また `dentries` や `inode` の slab オブジェ
2489 クトのうち可能なものを破棄する
- 2490 (c) 連続する空き領域を結合してより大きな空き領域にまとめる
- 2491 7. `ihk-smp-x86.ko` を `insmod` する。(340 行目) `ihk-smp-x86.ko` は関数を `ihk.ko` に登
2492 録する。このため、`ihk-smp-x86.ko` は `ihk.ko` を `insmod` した後に `insmod` する必要が
2493 ある。
- 2494 8. メモリを予約する。(370 行目)
- 2495 9. CPU を予約する。(374 行目)
- 2496 10. McKernel のカーネルモジュール `mcctrl.ko` を `insmod` する。(382 行目) `mcctrl.ko`
2497 は McKernel ブート時に呼ばれる関数を `ihk.ko` に登録する。このため、`mcctrl.ko` の
2498 `insmod` は `ihk.ko` の `insmod` の後に、またブートの前に行う必要がある。
- 2499 11. OS インスタンスを作成する。(406 行目)
- 2500 12. OS インスタンスに CPU を割り当てる。(412 行目)
- 2501 13. McKernel CPU の IKC メッセージ送信先の Linux CPU を設定する。(419 行目)
- 2502 14. OS インスタンスにメモリを割り当てる。(426 行目)
- 2503 15. カーネルイメージをロードする。(432 行目)
- 2504 16. カーネル引数をカーネルに渡す。(438 行目)
- 2505 17. カーネルをブートする。(444 行目)
- 2506 18. `/proc`, `/sys` ファイルの準備をする。また、その中で `mcoverlayfs.ko` を `insmod` す
2507 る。`mcoverlayfs.ko` は他モジュールとの依存関係を持たない。(454 行目から 567 行
2508 目) なお、関数インターフェイスでの対応関数は `ihk_os_create_pseudofs()` である。
- 2509 19. `irqbalance` を、Linux CPU のみを対象とする設定で開始する。(569–587 行目)

2510 3.3 シャットダウン手順

2511 mcstop+release.sh を用いてシャットダウン手順を説明する。
スクリプトは以下の通り。

```
1 #!/bin/bash
2
3 # IHK SMP-x86 example McKernel unload script.
4 # author: Balazs Gerofi <bgerofi@riken.jp>
5 # Copyright (C) 2015 RIKEN AICS
6 #
7 # This is an example script for destroying McKernel and releasing IHK resources
8 # Note that the script does no output anything unless an error occurs.
9
10 prefix="/home/takagi/project/os/install"
11 BINDIR="/home/takagi/project/os/install/bin"
12 SBINDIR="/home/takagi/project/os/install/sbin"
13 ETCDIR="/home/takagi/project/os/install/etc"
14 KMODDIR="/home/takagi/project/os/install/kmod"
15 KERNDIR="/home/takagi/project/os/install/smp-x86/kernel"
16
17 mem=""
18 cpus=""
19 irqbalance_used=""
20
21 # No SMP module? Exit.
22 if ! grep ihk_smp_x86 /proc/modules &&>/dev/null; then exit 0; fi
23
24 if [ "'systemctl status irqbalance_mck.service 2> /dev/null |grep -E 'Active: active'" \
25 != "" ]; then
26     irqbalance_used="yes"
27     if ! systemctl stop irqbalance_mck.service 2>/dev/null; then
28         echo "warning: failed to stop irqbalance_mck" >&2
29     fi
30     if ! systemctl disable irqbalance_mck.service >/dev/null 2>/dev/null; then
31         echo "warning: failed to disable irqbalance_mck" >&2
32     fi
33 fi
34
35 # Destroy all LWK instances
36 if ls /dev/mcos* 1>/dev/null 2>&1; then
37     for i in /dev/mcos*; do
38         ind='echo $i|cut -c10-';
39         # Retry when conflicting with ihkmond
40         nretry=0
41         until ${SBINDIR}/ihkconfig 0 destroy $ind || [ $nretry -lt 4 ]; do
42             sleep 0.25
43             nretry=$(( $nretry + 1 ))
44         done
45         if [ $nretry -eq 4 ]; then
46             echo "error: destroying LWK instance $ind failed" >&2
47             exit 1
48         fi
49     done
50 fi
51
52 # Query IHK-SMP resources and release them
53 if ! ${SBINDIR}/ihkconfig 0 query cpu > /dev/null; then
54     echo "error: querying cpus" >&2
55     exit 1
56 fi
57
58 cpus='${SBINDIR}/ihkconfig 0 query cpu '
59 if [ "${cpus}" != "" ]; then
60     if ! ${SBINDIR}/ihkconfig 0 release cpu $cpus > /dev/null; then
61         echo "error: releasing CPUs" >&2
62         exit 1
```

```

63         fi
64     fi
65
66     if ! ${SBINDIR}/ihkconfig 0 query mem > /dev/null; then
67         echo "error: querying memory" >&2
68         exit 1
69     fi
70
71     mem='${SBINDIR}/ihkconfig 0 query mem'
72     if [ "${mem}" != "" ]; then
73         if ! ${SBINDIR}/ihkconfig 0 release mem $mem > /dev/null; then
74             echo "error: releasing memory" >&2
75             exit 1
76         fi
77     fi
78
79     # Remove delegator if loaded
80     if grep mcctrl /proc/modules &>/dev/null; then
81         if ! rmmmod mcctrl 2>/dev/null; then
82             echo "error: removing mcctrl" >&2
83             exit 1
84         fi
85     fi
86
87     # Remove mcoverlay if loaded
88     if grep mcoverlay /proc/modules &>/dev/null; then
89         if [ "$(cat /proc/mounts | grep /tmp/mcos/mcos0_sys)" != "" ]; \
90     then umount -l /tmp/mcos/mcos0_sys; fi
91         if [ "$(cat /proc/mounts | grep /tmp/mcos/mcos0_proc)" != "" ]; \
92     then umount -l /tmp/mcos/mcos0_proc; fi
93         if [ "$(cat /proc/mounts | grep /tmp/mcos/linux_proc)" != "" ]; \
94     then umount -l /tmp/mcos/linux_proc; fi
95         if [ "$(cat /proc/mounts | grep /tmp/mcos)" != "" ]; then umount -l /tmp/mcos; fi
96         if [ -e /tmp/mcos ]; then rm -rf /tmp/mcos; fi
97         if ! rmmmod mcoverlay 2>/dev/null; then
98             echo "warning: failed to remove mcoverlay" >&2
99         fi
100     fi
101
102     # Remove SMP module
103     if grep ihk_smp_x86 /proc/modules &>/dev/null; then
104         if ! rmmmod ihk_smp_x86 2>/dev/null; then
105             echo "error: removing ihk_smp_x86" >&2
106             exit 1
107         fi
108     fi
109
110     # Remove core module
111     if grep -E 'ihk\s' /proc/modules &>/dev/null; then
112         if ! rmmmod ihk 2>/dev/null; then
113             echo "error: removing ihk" >&2
114             exit 1
115         fi
116     fi
117
118     # Stop ihkmond
119     pid='pidof ihkmond'
120     if [ "${pid}" != "" ]; then
121         sudo kill -9 ${pid} > /dev/null 2> /dev/null
122     fi
123
124     # Start irqbalance with the original settings
125     if [ "${irqbalance_used}" != "" ]; then
126         if ! etcdir=/home/takagi/project/os/install/etc perl -e \
127     '$etcdir=$ENV{'etcdir'}; @files = grep { -f } glob "$etcdir/proc/irq/*/smp-affinity";'\
128     ' foreach $file (@files) { $dest = substr($file, length($etcdir));'\
129     ' if(0) {print "cp $file $dest\n";} system("cp $file $dest 2>/dev/null"); }'; then
130             echo "warning: failed to restore /proc/irq/*/smp-affinity" >&2

```



```
131         fi
132         if ! systemctl start irqbalance.service; then
133             echo "warning: failed to start irqbalance" >&2;
134         fi
135     fi
136
137 # Set back default swappiness
138 echo 60 > /proc/sys/vm/swappiness
```

2512 手順は以下の通り。

- 2513 1. ブート時に Linux CPU のみを対象とする設定で開始された irqbalance を停止する。
2514 (24-33 行目)
- 2515 2. 全ての OS インスタンスを破壊する。OS インスタンスに割り当てられていた資源は IHK
2516 が LWK のために予約した状態に移行する。(35-50 行目)
- 2517 3. IHK が LWK のために予約していた資源を開放する。(52-77 行目)
- 2518 4. mcctrl.ko を rmmod する。(81 行目)
- 2519 5. /proc, /sys ファイルの準備をする。また、その中で mcoverlayfs.ko を rmmod する。
2520 (87-100 行目) なお、関数インターフェイスでの対応関数は ihk_os_destroy_pseudofs()
2521 である。
- 2522 6. ihk-smp-x86.ko を rmmod する。(104 行目)
- 2523 7. ihk.ko を rmmod する。(112 行目)
- 2524 8. ihkmond を停止する。(121 行目)
- 2525 9. /proc/irq/*/affinity の設定をブート時に保存しておいたものに戻し、ブート前の設
2526 定で irqbalance を開始する。(124-135 行目)
- 2527 10. Linux カーネルのスワップアウト積極度の設定をデフォルトの値に戻す。(138 行目)

Bibliography

- 2529 [1] H. Fujita, M. Matsuda, T. Maeda, S. Miura, and Y. Ishikawa. P-Bus: Programming
2530 Interface Layer for Safe OS Kernel Extensions. In *Pacific Rim International Symposium*
2531 *on Dependable Computing (PRDC)*, pages 235–236, 2010.
- 2532 [2] T. Shimosawa, B. Gerofi, M. Takagi, G. Nakamura, T. Shirasawa, Y. Saeki, M. Shimizu,
2533 A. Hori, and Y. Ishikawa. Interface for Heterogeneous Kernels: A Framework to Enable
2534 Hybrid OS Designs targeting High Performance Computing on Manycore Architectures.
2535 *In Proc. of IEEE International Conference on High Performance Computing (HiPC)*,
2536 2014.