

J-PARC MR 制御での仮想マシンの応用

VIRTUAL MACHINES IN J-PARC MR CONTROL

上窪田紀彦^{#, A)}, 吉田奨^{B)}, 本橋重信^{B)}, 飯塚上夫^{B)}, 根本弘幸^{C)},
高橋大輔^{B)}, 山本昇^{A)}, 山田秀衛^{A)}, 佐藤健一^{A)}

Norihiko Kamikubota^{#, A)}, Susumu Yoshida^{B)}, Shigenobu Motohashi^{B)}, Takao Iitsuka^{B)}, Hiroyuki Nemoto^{C)},
Daisuke Takahashi^{B)}, Noboru Yamamoto^{A)}, Shuei Yamada^{A)}, Kenichi C.Sato^{A)}

^{A)} J-PARC Center, KEK and JAEA

^{B)} Kanto Information Service (KIS)

^{C)} ACMOS Co. Ltd.

Abstract

In J-PARC MR, a scheme of virtual ioc (an EPICS I/O-controller which runs on a virtual machine) was introduced successfully in 2011. Following on this success, we intended to introduce more virtual machines even for server services. Scientific Linux 6 and KVM are selected as an OS and virtual environment for parent servers. Details of our virtual environment setups and experiences for the past two years are reported

1. はじめに

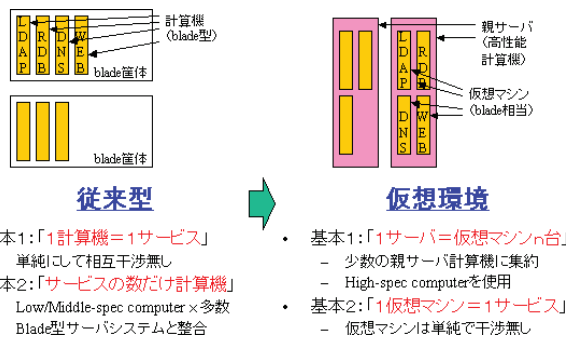
最近、計算機の仮想化技術が著しく発展・普及し、さまざまな機会に耳にするようになった。仮想環境を利用するには、かつては専門的な知識やそれなりの予算が前提であった。今や Windows には Virtual PC が、Linux では KVM が、仮想環境として標準で用意されている。仮想環境が身近になったのは、急速な 64 ビット PC の普及で、幅広いユーザが旧 OS (特に 32bit 時代) の backward-compatibility を必要としたためと思われる。

仮想化技術を導入すると、サーバ計算機の構成には構造的な変化が起こる。従来は 1 台の計算機には 1 種のサービスのみを受け持たせ、個々の計算機の管理を単純化しようと考えた。その結果多数の Low/Middle-spec の計算機が必要になるが、blade 型計算機システムの導入で計算機機体の扱いは効率的に行えた (図 1 左)。一方仮想環境では、1 台の仮想マシンが 1 種のサーバを受け持つ (管理の単純さは従来と同じ) が、少数の High-spec な親計算機が複数の仮想マシンを持つ構成となる (図 1 右)。仮想環境では、負荷状態や保守都合に応じ仮想マシンを親サーバ間で柔軟に移動出来るため、少数の親サーバ機のハードウェアに集中投資して総コストを下げる事が出来る。

J-PARC MR は、EPICS ベースの制御システムで運用している^[1]。100 個以上の IOC (In Out Controller) が使用されているが、IOC の半数はこの分野で標準的な VME-bus 計算機である^[2]。2010 年頃から VMWARE や XEN などを使用して仮想環境を加速器制御に応用する検討を行っていた。2011 年には仮想環境に KVM を導入し、仮想マシン上で IOC を動かす”Virtual IOC (vioc)”を開発した^[3]。ソフトウェアやネットワークデバイスのみを扱う IOC

(Real I/O 無し) の一部は、VME-bus 計算機から vioc に移行した。現在(2013 年)、約 30 台の vioc を運用している。

vioc の成功に引き続き、2012 年に J-PARC MR 制御に必要なサーバ機能の仮想マシンへの移動が試みられた^[4]。本稿では、J-PARC MR の仮想環境を解説し、2 年間の運用経験を報告する。



- 基本1:「1 計算機 = 1 サービス」
 - 単純にして相互干渉無し
- 基本2:「サービスの数だけ計算機」
 - Low/Middle-spec computer × 多数
 - Blade型サーバシステムと整合
- 基本1:「1 サーバ = 仮想マシン n 台」
 - 少数の親サーバ計算機に集約
 - High-spec computer を使用
- 基本2:「1 仮想マシン = 1 サービス」
 - 仮想マシンは単純で干渉無し

Figure 1: Server structures with/without virtualization.

2. J-PARC MR の仮想環境運用

2.1 仮想環境の概要

J-PARC MR では、主たる OS として Linux 系の Scientific Linux (SL) を採用している。運転が始まった 2008 年は SL4 であったが、2012 年夏季にサーバ系計算機を SL4 から SL6 に update した。当時、KVM で vioc の運用経験を積んでいたこと、SL6 に KVM が装備されたこと、などから、KVM を標準的に利用する仮想環境とした。

現在 (2013 年) サーバ系と vioc 系の 2 系統の仮想環境を運用している。それぞれの系は、親サーバ計算機 3 台で構成する (ただし 1 台は両系で重複し総数 5 台)。通常時は 3 台で負荷が平均化するよう

[#] norihiko.kamikubota@kek.jp

仮想マシンを分散しているが、1台が故障した際は残りの2台で運転を継続する。図2(a)にサーバ系、図2(b)にvioc系の構成を示す。

親サーバ計算機は、20GBのメモリを積んだ比較的High-specな計算機(IBM Blade HS22, Xeon E5504 (4core 2cpu) 2GHz)を使用している。これに対し仮想マシンに割り当てる1台あたりの典型的なメモリは、サーバ系で2-4GB、vioc系で512MBである。

図2(a)で示すように、サーバ系が実計算機から仮想環境に移行したことで、全体として8台の計算機を3台に集約したことになる。Archive engineは今後負荷が増えると思われるが、分割や親サーバ追加で対応する予定である。保守serverは特殊な仮想マシンで、SL6移行以前のSL4.4やサーバ系仮想マシンで使用するSL5.4の開発環境を保持する(この事情は2.3節も参照されたい)。昔のアプリの動作検証などに重宝している。

サーバ系

親サーバ機(親)	仮想マシン(子)
SL6.0, mem=20GB	SL5.4, mem=2-4GB (typical spec for each v-machine)
jkjblade3a	<ul style="list-style-type: none"> Archive engine CA-gateway
Jkjblade3b	<ul style="list-style-type: none"> 管理server (dhcp, ftp, ldap slave) App server (cron, zlog, cacti) RDB server (postgres, mysql) 保守server (SL5.4環境)
jkjblade3c	<ul style="list-style-type: none"> 管理server (ldap master) 保守server (旧OS:SL4.4環境)

Figure 2(a): Assignment for virtual machines (servers).

Vioc系

親サーバ機(親)	仮想マシン(子)
SL6.0, mem=20GB	SL6.3, mem=512MB (typical spec for each vioc)
jkjblade3a	<ul style="list-style-type: none"> cont-group 3台
Jkjblade3e	<ul style="list-style-type: none"> mag-group 4台 mon-group 3台 rf-group 1台 inj-group 2台 accom-group 1台
jkjblade3f	<ul style="list-style-type: none"> cont-group 6台 mon-group 1台 sx-group 2台 個人試験・開発用 8台

Figure 2(b): Assignment for virtual machines (viocs).

2.2 仮想環境の運用・管理

典型的な親サーバ機(jkjblade3b、jkjblade3e)の運用中の負荷状態(10分間のCPUとメモリ)を図3に示す。負荷監視は、標準的な管理ツールを使っている。親サーバ機のうち1台(jkjblade3b)では、メモリ使用率がやや高くなっていることがわかる。

我々は全部で5台の親サーバ機を持っているが、ある仮想マシン(子)の親を変更したい場合、やはり標準的な管理ツールで作業する。図4は、ある仮想マシンを選択し、親サーバ機を変更しようとする

作業の画面である。



Figure 3: Typical CPU and memory load status.

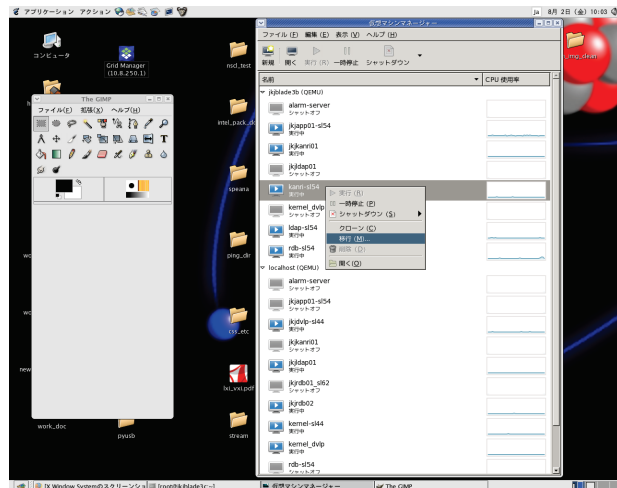


Figure 4: Management tool for virtual machines.

2.3 仮想環境の障害

サーバ系の仮想環境運用を開始した2012年は、仮想マシン(子)のOSに親サーバ機と同じSL6.0を使用していた。ところが、保守や故障でネットワークの一時停止が起こった後、仮想マシンのNFSが正常に復帰せず、read-onlyとなる現象が見られた。この現象は、仮想でないSL6計算機では起こらない。また、不調となった仮想マシンを復活させるには再起動が必要であり、サーバ系サービスにとっては大きな不具合であった。

このNFS不調問題は、仮想マシンのNFSパラメータ調整などを試みたが効果が無く、原因ははっきりしていない。SL6のNFSがv3からv4に改修された余波のように思えるが、情報収集を続ける。現時点では、仮想マシンのOSをあえてSL5.4にdowngradeし、問題を回避している。

2013年2月9日、vioc系の親サーバ機の1台(jkblade3f)が突然停止する事件があった。その親サーバ機で稼動していた仮想マシンは全滅したため、別の親サーバ機でそれらの仮想マシンを手作業で再起動した。この作業には約3時間かかった。このように、親サーバ機が死んだ場合の対応は自動ではない。当面は手作業での対応手順の効率化検討を行うが、さらなる信頼度を期待するならば、自動で親を移動する仕組みなどが望まれる。

なお、この時の親サーバ停止の原因を追及すると、kernel configのsoft_lockup設定に行き当たった。仮想マシンが親から実CPUを割り当ててもらえない場合、SL6.0の初期設定では親がいきなり落ちる仕様になっている。親サーバ機全部でsoft_lockupの見直し・再設定を行い、その後半年間この問題は発生していない。

4. まとめ

J-PARC MRでは、計算機制御に仮想環境KVMを導入した。2011年からIOC系仮想マシンを、2012年からサーバ系仮想マシンの運用を始めた。

いくつかの障害を経験したが、おおむね運用は順調である。仮想の親サーバ機が停止した場合の復帰手順の確立が望まれる。

参考文献

- [1] N.Kamikubota, et al., "J-PARC CONTROL TOWARD FUTURE RELIABLE OPERATION", ICALEPCS2011, Grenoble, France, Oct. 2011, p.378-381
- [2] 根本弘幸、他、"J-PARC MR 加速器制御システムにおけるIOC統合管理"、加速器学会(大阪)、Aug. 2012、p.745-747
- [3] N.Kamikubota, et al., "VIRTUAL IO CONTROLLERS AT J-PARC MR USING XEN", ICALEPCS2011, Grenoble, France, Oct. 2011, p.1165-1167
- [4] 上窪田紀彦、他、"J-PARC MR 制御計算機の進展"、加速器学会(大阪)、Aug. 2012、p.741-744