

名古屋大学大型計算機センター新システムの運用管理について

長谷川 明生

名古屋大学大型計算機センター

梗概

名古屋大学大型計算機センターの新システムは、汎用機、42台のベクトルプロセッサを持つスーパーコンピュータ、2台のサーバー型ワークステーションおよび画像処理用ワークステーション7台を中心に構成されている。汎用機用オペレーティング・システムおよび複数種類のUNIXオペレーティング・システムが共存した大規模分散処理環境である。

本論文では、このような分散処理環境下での利用者管理や課金管理の方式について議論する。

On the Management Facilities of the Computer System of the Nagoya University Computation Center

Hasegawa Akiumi

Nagoya University Computation Center

Abstract

The system of the Nagoya University Computation Center consists of one mainframe, one super computer with 42 vector processors, two server type workstations, and seven graphic workstations. The operating system MSP and the UXP/M, a kind of UNIX, are operated on the mainframe system simultaneously. The super computer is operating under the UXP/VPP, a kind of UNIX operating system. The Operating systems of the other systems are UNIX of different kinds. These systems are connected with each other by LAN and constitute a large distributed system.

This paper describes the management facilities of our new system.

はじめに

昨年度、名古屋大学大型計算機センターは、大規模な計算機システムの更新を行った。現システムは、42個のベクトル・プロセッサを持つベクトル型並列計算機システム、3CPUの汎用計算機、2台のサーバー型ワークステーション、および7台の画像処理用ワークステーションをLANで相互接続した大規模分散処理システムである。オペレーティング・システムは、汎用機用のMSPおよびUXP/M、スーパーコンピュータ用のUXP/MおよびUXP/VPP、ワークステーション用にSolaris 2およびIRIXと変化にとんでいる。

本論文では、大規模分散処理システムでの利用者管理および課金管理方式について、実際の運用経験に基づいて報告する。

システムの概要

図1に、現システムの構成図を示す。

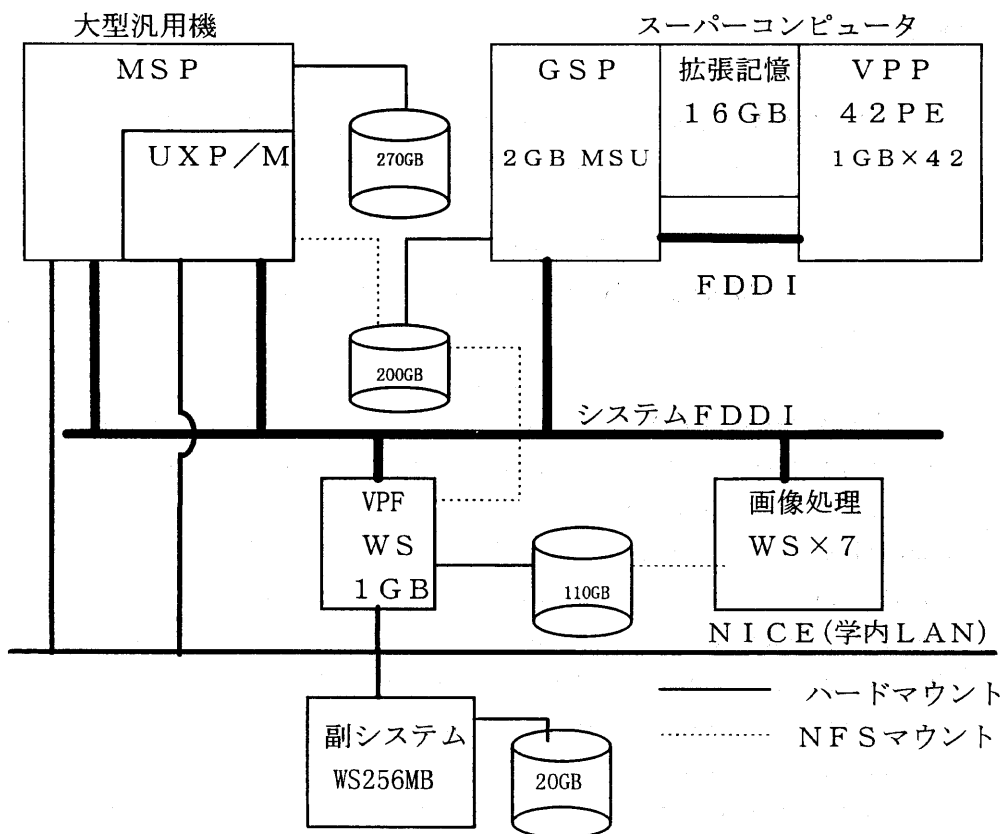


図1 システム構成図

スーパーコンピュータ (VPP 500/42) は、バッチ計算専用のバックエ

ンド・プロセッサとして動作する。スーパーコンピュータへのジョブ投入は、

- 汎用機からM-V P P連携機能による汎用機のバッチ・ジョブ投入
- 汎用機のUNIXからのNQSによるジョブ投入
- V P Fと呼ぶサーバー型ワークステーションからのNQSによるジョブ投入

が可能となっている。

主な機器の構成を表1に示す。

表1 センターの機器構成

機種	OS	ネットワーク	備考
M1800/30 (3CPU, 1GB)	MSP, UXP/M	FDDI, Ether	汎用機
GSP (2CPU+1VU, 2GB, 16GBSSU)	UXP/M	FDDI	VPPサービスプロセッサ
VPP500/42 (42PE, 1GB/PE)	UXP/M VPP	FDDI	ベクトル並列型計算機
S4/1000E (8CPU, 1GB)	Solaris2.4	FDDI, Ether	VPPフロントエンド
S4/1000E (2CPU, 256MB)	Solaris2.4	Ether	UNIX副システム
S4/20H (256MB) ×5	Solaris2.4	FDDI	画像処理システム
S4/10 (256MB)	Solaris2.4	FDDI	画像処理システム
S4/10 (196MB)	Solaris2.4	FDDI	画像処理システム
SGI Reality Engine (256MB) ×2	IRIX5.3	FDDI	画像処理システム

V P Pは単独ではファイル処理機能等を持たないので、ハウスキーピングのためのG S Pシステムと一体で運用する必要がある。これらのシステムはすべてF D D IおよびE t h e r n e tにより結ばれており、大規模な分散処理システムを構成している。

利用者管理システム

これらのシステム間で利用者管理は自動化されている。UNIXシステムでは、N I Sによる利用者管理を行っている。N I Sのマスター・サーバーは、UNIX副システムで、V P Pフロントエンド・システムがスレーブ・サーバーとして動作し、システムF D D I上のUNIXホストにN I Sマップをサービスする。ただし、V P PおよびG S Pシステムはバッチ専用のシステムなのでN I Sマップではなくパスワード・ファイルによる管理としている。

利用者登録の流れを図2に示す。利用者登録は、課金処理とあわせて毎朝実行される。UNIX上では、利用者登録プログラムは/etc/passwdファイルやN I Sマップの更新だけでなく、ホーム・ディレクトリの作成、.loginファイル等の準備、quotaの設定等をおこなう。

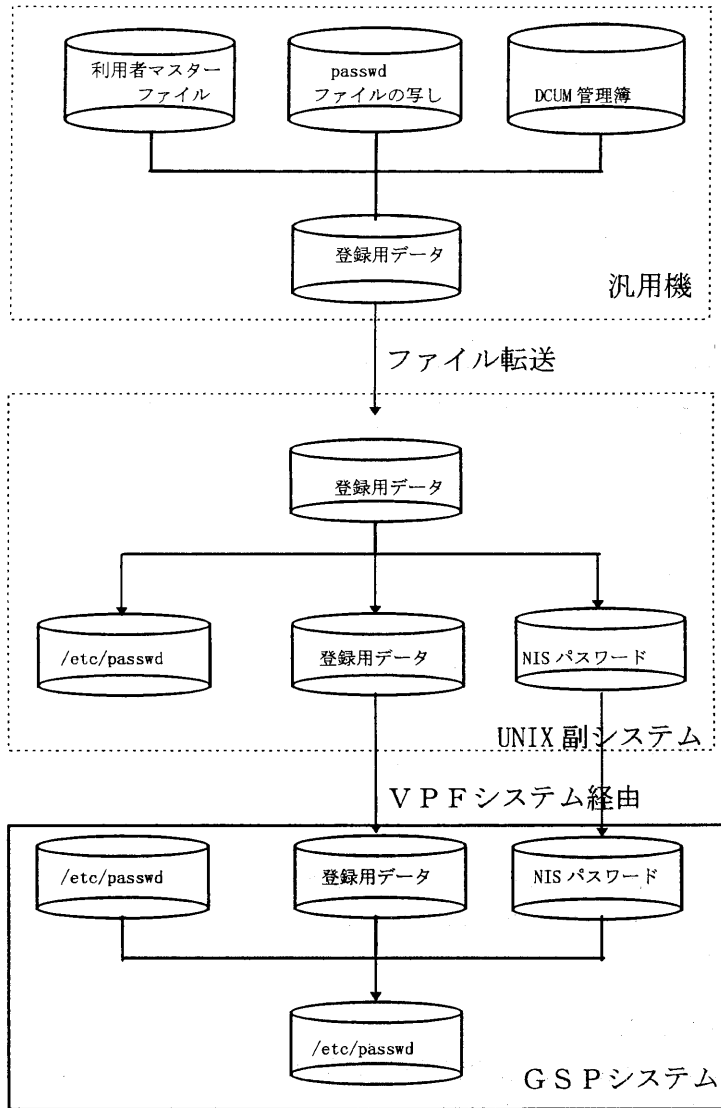


図2 利用者登録の流れ

課金処理

本センターの課金項目は、CPU使用料金、出力負担経費（プリンター出力、XYプロッタ）およびディスク使用料金である。汎用機およびスーパーコンピュータのCPUは、ジョブ単位の課金処理をおこなっている。他のUNIXホストについては、プロセスのCPU時間を1日分について合算処理する。

ディスク使用量は、毎日定時にパトロール・プログラムを走らせて調べる。

UNIX系の印刷出力は、センター内のネットワークプリンタのどれかに出力する。これらのネットワークプリンタは、学内LANとは独立したネットワーク

に接続している。そして、このLANと学内LANとの間にプリンタ・サーバを兼ねたワークステーションを用意した。このような構成により、利用者チェックおよび課金管理が容易になっている。

これらの独立にシステムごとにあがってくる課金データは、すべて汎用機に転送し汎用機上の課金プログラムで課金処理している。

管理システムの設計方針

利用者登録システムおよび課金管理システムの開発にあたって、極力システムの標準機能を手を入れずに利用することを原則とした。また、新システムへの移行およびテスト期間を最小とするために、既存の汎用機上で動作している課金や登録管理システムの機能を前提にUNIXホスト等の利用者登録システムや課金データの形式を設計した。

このような方針で開発したため、システム全体を通してのテスト期間は10日程度しか確保できなかったが、新システムへの移行を短時間で行うことができた。

ジョブクラス設計と運用状況

新システムの中心は、ベクトル並列型スーパーコンピュータである。このようなシステムの大規模計算機センターでの本格運用は、本センターが初めてである。複数プロセッサを並列に動作させるジョブに対する利用者の要求は、システム設計時には、ほとんど予想できなかった。表2に現在のNQSジョブ・クラスとその制限値を示す。

表2 NQSキューと制限値

キュー名	利用目的	CPU時間	メモリー	PE数
c	Fortranの翻訳	30分	100MB	-
x	VPPでの実行	600分	300MB	1
y	VPPでの実行	600分	900MB	1
z	VPPでの実行	600分	900MB/PE	2~16

このキュー設計は、本格的な並列プログラムの開発に時間がかかると予想し、シングルPEジョブ主体を想定して行われた。

実際のジョブ処理状況を表3に示す。

表3 ジョブ処理状況表

システム	年月	計算サービス時間(時:分)	処理件数	CPU時間 時:分	CPU時間/件 (秒)
VPP500	H7.12	599:12	3759	5454:47	5224.05
	H8.1	605:41	4549	7293:53	5772.25
	H8.2	620:12	4677	10012:03	7706.52
	H8.3	445:27	3430	6661:31	6991.68
	計	2270:32	16415	29422:14	6454.64
VP2600	通年	5037:13	25537	4857:25	684.76

また、表 4 に P E の使用状況を示す。

表 4 P E 使用状況表

年月	利用種別	ベクトル化率	平均 P E 数	メモリサイズ
1 月	M-V P P 連携	77.0%	1.29	131.6MB
	N Q S	37.1%	2.63	409.1MB
2 月	M-V P P 連携	72.5%	1.08	108.5MB
	N Q S	29.5%	2.94	376.9MB

ジョブ統計の平均 P E 数の値から、当初の予想と異なり複数 P E を使用するジョブがスーパーコンピュータ利用のかなりの部分を占めていることがわかる。複数 P E を使用した場合、最大 C P U 時間を消費した P E のみ C P U 課金されるという誘導的な課金体系の効果が予想以上に表れたとも言える。また、複数 P E ジョブでは、4 P E 以上が同時に使用されることが多い。導入当初からの並列ジョブの利用者に確認したところ、同一プログラム中で複数データを同時並列処理しているとのことであった。このような形態のジョブでは、プロセスの起動処理とファイル入出力の競合だけが問題となるだけで、ほぼ P E 数に比例して高速化がはかれる。このようなジョブとは、別に複数 P E で並列処理することによって高速化をねらったジョブも投入されている。

スーパーコンピュータ利用者は、現在のところ、シングル P E ジョブ主体の旧システムからの移行者とフロント・エンド・システムや汎用機の U N I X システムから N Q S ジョブを投入する並列ジョブ利用者との 2 分される。現行のジョブ・プロフィールから、N Q S キューの資源配分の再検討をしなければならないが、4 2 プロセッサというのは、並列ジョブの増加を予想すると、けっして十分な数ではないと考えている。

おわりに

本センターの新システムについて運用管理方式を中心に報告した。現在は、新年度の新規利用者登録も済み安定に稼動している。

短期的課題として、スーパー・コンピュータに関連した資源配分の再検討およびジョブ管理パラメータとしての課金体系の見直しがある。

長期的には、汎用大型計算機の役割の検討を含め、将来を見通した計算機環境の検討しなければならない。