



دانشکده برق و کامپیوتر
دانشگاه صنعتی اصفهان

openMosix

ابزاری برای توسعه کلاسترهای مبتنی بر لینوکس

گزارش سمینار درس سیستم عامل پیشرفته

مدرس:

جناب آقای دکتر داورپناه جزی

ارائه دهندگان:

فرجس خاتون حبیبی

زینب زالی

زمستان ۸۵

فهرست مطالب

۳	۱- مقدمه
۳	۲- تعریف کلاستر
۳	۱-۲- تعریف کلاستر کامپیوتری
۳	۲-۲- تقسیم بندی کلاستر های کامپیوتری
۴	۳- معرفی openMosix
۵	۴- نصب openMosix
۵	۴-۱- کامپایل کردن کرنل لینوکس
۶	۴-۲- نصب و تنظیم ابزار ها
۹	۵- تکنولوژی openMosix
۱۰	۵-۱- مکانیزم مهاجرت پروسه (PPM)
۱۱	۵-۲- الگوریتم های اشتراک منابع
۱۱	۵-۲-۱- الگوریتم موازنه بار
۱۲	۵-۲-۲- الگوریتم راهنمای حافظه
۱۲	۵-۳- مهاجرت پروسه (PPM)
۱۳	۵-۴- نکاتی در رابطه با مکانیزم remote/deputy
۱۴	۵-۵- سیستم فایل openMosix
۱۴	۵-۶- API openMosix
۱۵	۶- معرفی openMosixview
۱۵	۶-۱- openMosixview
۱۸	۶-۲- openMosixprocs
۲۱	۶-۳- openMosixcollector
۲۴	۶-۵- openMosixmigmon
۲۴	۷- جمع بندی
۲۵	منابع

۱- مقدمه

تکنولوژی کلاسترینگ امکان می‌دهد که چند سیستم کامپیوتری منابع محاسباتی خود را، برای حل سریعتر مسائل، ترکیب کنند. openMosix یک کلاستر مبتنی بر لینوکس است که با توزیع یکنواخت کارها روی سیستم‌های موجود، کارایی اجرای برنامه‌ها را افزایش می‌دهد. مزیت اصلی openMosix این است که تمام برنامه‌ها به طور شفاف و بدون نیاز به استفاده از کتابخانه‌های خاص، قادر به استفاده از محیط محاسباتی توزیع شده openMosix خواهند بود و نیازی به تغییر در ساختار برنامه‌ها وجود ندارد. در این گزارش نحوه نصب، استفاده و ساختار درونی openMosix تشریح شده است.

۲- تعریف کلاستر

در تعریف کلی یک خوشه یا کلاستر مجموعه‌ای از اجزایی است که به طور مستقل از یکدیگر کار کرده و توسط یک واسطه هماهنگ می‌شوند.

۲-۱- تعریف کلاستر کامپیوتری

کلاستر کامپیوتری گروهی از کامپیوترهای مجزا است که با یکدیگر کار می‌کنند، به گونه‌ای که می‌توان آنها را تقریباً یک کامپیوتر واحد دانست. اجزای یک کلاستر معمولاً به وسیله شبکه‌های محلی پرسرعت به یکدیگر متصلند.

۲-۲- تقسیم بندی کلاسترهای کامپیوتری

کلاسترها به سه دسته کلی تقسیم می‌شوند:

- HAC (High Availability Cluster): در این کلاستر هدف افزایش آماده بودن^۱ سرویس‌های است هائی که کلاستر فراهم می‌کند، به وسیله ندهای تکراری، می‌باشد. هنگام خراب شدن یکی از ندها، ندهای دیگر به سرویس دهی ادامه می‌دهند و کاربر متوجه اشکال به وجود آمده نمی‌شود.
- LBC (Load Balancing Cluster): در این نوع کلاسترها تمام بار کاری به یک یا چند سرور ارجاع شده و آنها بار را میان تعدادی سرور back-end تقسیم می‌کنند. این نوع سرورها معمولاً خصوصیت HAC را نیز دارند. به

^۱ Availability

چنین کلاسترهایی Server Farm نیز گفته می شود. از جمله پیاده سازی های تجاری LBC می توان Sun Grid Engine را نام برد.

- HPC (High Performance Cluster) : این کلاستر به وسیله تقسیم یک وظیفه میان چندین ند کارائی را افزایش می دهد و بیشتر در محاسبات علمی کاربرد دارد. معمولا HPCها را با استفاده از لینوکس و نرم افزارهای آزاد ایجاد می کنند. به این کلاسترها Beowulf نیز می گویند. این کلاسترها برنامه های خاصی را- که برای بهره برداری از توازی سازی HPC طراحی شده اند- اجرا می کنند. این برنامه ها از کتابخانه هایی مانند MPI استفاده می نمایند.

۳- معرفی openMosix

openMosix یک توسعه کرنل لینوکس برای ساخت کلاسترهای HPC است. openMosix امکان می دهد چندین تک پردازنده و یا تعدادی چندپردازنده^۲ که همه از کرنل یکسان استفاده می کنند، با یکدیگر همکاری نمایند. الگوریتم های به اشتراک گذاری openMosix به گونه ای طراحی شده اند که بتواند به تغییرات استفاده از منابع در ندها، در زمان اجرا، پاسخ دهند. این کار به وسیله مهاجرت^۳ پروسه ها از یک ند به ند دیگر، به طور شفاف^۴ و به منظور توزیع بار و جلوگیری از کویدگی^۵ ناشی از مبادله صفحات حافظه صورت می گیرد. هدف بهبود کارایی کل کلاستر و ایجاد محیطی چندکاربره و اشتراک زمانی برای برنامه های موازی و متوالی است.

کلاسترهای openMosix می تواند شامل تعداد محدودی از PCها که با سرعت 10 MBPS به یکدیگر متصلند باشد و یا مجموعه وسیعی از سرورهای SMP که به وسیله Gigabit LAN و یا ATM به هم مرتبط اند. بعد از نصب openMosix بر روی تعدادی ند، ندها شروع به صحبت با یکدیگر می نمایند. یک پروسه بر روی هر ند دلخواه می تواند ایجاد شود و openMosix در صورت نیاز آن را به ندی که مناسب تر باشد منتقل می نماید.

² SMP

³ Migration

⁴ Transparent

⁵ Thrashing

مزیت اصلی openMosix این است که چون تغییرات به کرنل اعمال می شود، برنامه ها به طور شفاف و خودکار قادر به استفاده از محیط محاسباتی توزیع شده openMosix خواهند بود و نیازی به تغییر در ساختار برنامه وجود ندارد.

۴- نصب openMosix

openMosix شامل یک patch کرنل و تعدادی ابزار در سطح کاربری است. patch کرنل برای آن است که کرنل بتواند با دیگر ماشین های با قابلیت openMosix روی شبکه ارتباط برقرار کند. اگر از بسته RPM برای نصب openMosix استفاده می کنیم دیگر نگران patch کرنل نیستیم، زیرا در بسته RPM کرنل patch شده و با گزینه ها و ماژول های معمول کامپایل شده است. ابزارهای سطح کاربر برای استفاده بهینه از قابلیت های کرنل کامپایل شده با openMosix لازم هستند، برای مثال برای شروع و خاتمه دادن به سرویس مهاجرت پروسه ها (migration)، به کارگیری سیستم فایل openMosix و مهاجرت یک پروسه به ند خاص در شبکه.

بنابراین طبق آنچه گفته شد لازم است برای ایجاد یک کلاستر openMosix، دو بسته مختلف را دانلود و نصب کرد: کرنل و ابزارهای سطح کاربر. در این راستا می توان دو بسته باینری و یا سورس patch کرنل و ابزارها را دانلود نمود. Patch کرنل بدین صورت نام گذاری می شود: openMosix-x.y.z-w که در آن x.y.z نشان دهنده نسخه کرنل موردنظر است که patch باید به آن اضافه شود و w نسخه بازبینی شده patch برای همان نسخه کرنل می باشد.

۴-۱- کامپایل کردن کرنل لینوکس

ابتدا patch متناسب با سورس کرنلی که در اختیار دارید را دانلود کنید، مثلا برای سورس linux-2.4.16 نیاز به فایل openMosix-2.4.16.x.tar.gz دارید. اکنون باید patch را به سورس کرنل اضافه و کرنل حاصل را کامپایل نمایید. ابتدا یک کپی از سورس لینوکس خود بگیرید تا در صورت بروز با مشکل در اضافه کردن patch سورس صحیح را از دست ندهید!

```
# cd /usr/src
# mv linux linux.old
```

حال سورس جدید را باز کنید:

```
# tar xzvf /path/to/linux-2.4.16.tar.gz
```

اکنون وارد شاخه سورس شوید و patch ای که دانلود کرده‌اید را با دستور زیر به سورس کرنل اضافه کنید:

```
# cd linux
# cat /path/to/openMosix-2.4.16-3.gz | gzip -d | patch -p1 -
```

پس از این مرحله کرنل آماده کامپایل شدن است. طبق مراحل معمولی که در کامپایل کردن کرنل طی می‌شود می‌توانید این کرنل را کامپایل کنید. تنها تفاوتی که در این جا وجود دارد مربوط به تنظیمات openMosix است که به کرنل اضافه شده است. بنابراین در مرحله‌ای که تنظیمات را با اجرای دستور make menuconfig دنبال می‌کنید، حداقل گزینه‌های زیر را فعال کنید:

```
[*] openMOSIX process migration support
[*] Stricter security on openMOSIX ports
[*] Direct File-System Access
[*] openMOSIX File-System
[*] Poll/Select exceptions on pipes
```

بعد از این که کامپایل کرنل با موفقیت انجام شد، باید ابزارها را نصب کنید.

۴-۲- نصب و تنظیم ابزارها

نصب ابزارها به راحتی با کمک فایل راهنمایی که برای نصب در بسته ابزار وجود دارد انجام می‌شود. برای این که نخواهیم ابزارها را روی یک به یک ندهای کلاستر جداگانه نصب کنیم، یک بار آن‌ها را روی یک ند که کرنل همراه openMosix روی آن نصب شده است نصب می‌کنیم، سپس از فایل‌های حاصل از نصب یک فایل tar می‌سازیم. حال کافی است این فایل tar را روی هر ند باز کنیم. برای ساخت این فایل tar روند زیر را دنبال کنید:

فرض کنید که روی یک ماشین کرنل همراه openMosix را در شاخه /usr/src/linux نصب کرده‌اید. فایل سورس ابزار را باز کنید و وارد شاخه مربوطه شوید:

```
# cd /tmp
# tar xzvf /path/to/openMosixUserland-0.1.3.tgz
# cd openMosixUserland-0.1.3
```

حال فایل تنظیمات را به درستی تنظیم کنید. تنظیمات باید مشابه زیر باشد:

```
OPENMOSIX=/usr/src/linux
PROCDIR=/proc/mosix
MONNAME=mmon
```

```
CC=gcc
INSTALLDIR=/usr
CFLAGS=-I/m/include -I./ -I/usr/include -I$(OPENMOSIX) /include
-O2
INSTALL=/usr/bin/install
```

ابزار را کامپایل کنید:

```
# make clean build man
```

حال ابزار را در یک شاخه موقتی مانند زیر نصب کنید:

```
# mkdir /tmp/openmosix-tools
# make INSTALLDIR=/tmp/openmosix-tools/usr install
# mkdir /tmp/openmosix-tools/etc
```

اکنون یک تصویر از تمام ابزارهای openMosix در /tmp/openMosix-
tools نصب شده است. می توان آن را به راحتی در شاخه / از هر ند باز کرد. قبل از
آن یک سری تغییرات دیگر هم در فایل tarball می دهیم: یک اسکریپت برای
startup/shutdown و یک فایل تنظیمات.

ابتدا فایل /etc/mosix.map را بسازید که فایل تنظیمات کلاستر است:

```
vi /tmp/openmosix-tools/etc/mosix.map
```

حال یک خط برای هر ند در کلاستر مانند زیر اضافه کنید:

```
<node number> <ip address> <span>
```

node number را می توانید از شماره ۱ برای ندها شروع کنید. مقدار IP
Address را هم برای هر ند تعیین کنید و در نهایت مقدار "1" را برای sapan قرار
دهید. مثلا برای یک کلاستر با دو ند فایل map می تواند مانند زیر باشد:

```
1 192.168.1.1 1
2 192.168.1.9 1
```

اگر بخواهید کلاستری با ۱۰۰ ند داشته باشید که آدرس ندها از
192.168.1.1 تا 192.168.1.100 است، با استفاده از گزینه span می توان این
کار را انجام داد:

```
1 192.168.1.1 100
```

سپس فایل map را ذخیره کنید.

اکنون اسکریپتی برای فعال و غیرفعال کردن openMosix می‌نویسیم. ابتدا

در شل دستور زیر وارد کنید و سپس فایل openMosix را بسازید:

```
# install -d /tmp/openmosix-tools/etc/rc.d/init.d  
#vi /tmp/openmosix-oools/etc/rc.d/init.d/openmosix
```

اسکریپت زیر را در آن کپی کنید:

```
#!/bin/bash  
case $1 in  
stop)  
echo "Stopping openMosix"  
echo 0 > /proc/mosix/admin/mospe  
rm -f /var/lock/subsys/mosix  
;;  
start)  
echo "Starting  
openMosix"  
if [ -s /etc/overheads -a -f /proc/mosix/admin/overheads  
]  
then  
grep -v '^#' /etc/overheads >  
/proc/mosix/admin/overheads  
fi  
if [ -s /etc/mfscosts -a -f /proc/mosix/admin/mfscosts ]  
then  
grep -v '^#' /etc/mfscosts >  
/proc/mosix/admin/mfscosts  
fi  
local a1  
local a2  
a1=[ -s /etc/mospe ] && a1="-p `cat /etc/mospe`"  
a2=[ -s /etc/mosgates ] && a2="-g `cat /etc/mosgates`"  
setpe -W $a1 $a2 -f /etc/mosix.map  
touch /var/lock/subsys/mosix  
;;  
esac
```

در نهایت فایل اسکریپت را با دستور زیر قابل اجرا کنید:

```
#chmod +x /tmp/openmosix-tools/etc/rc.d/init.d/openmosix
```

هر ند باید یک شاخه برای mount کردن سیستم فایل داشته باشد، بنابراین شاخه /mnt را اضافه کنید:

```
# mkdir /tmp/openmosix-tools/mfs
```

اکنون شاخه نهایی شده است و آماده ساخت tarball هستید:

```
# cd /tmp/openmosix-tools  
# tar czvf /tmp/openmosix-tools.tar.gz
```

حال یک فایل openmosix-tools.tar.gz در شاخه /tmp دارید که می‌توانید آن را روی هر ند نصب کنید:

```
# tar xzpvf /path/to/openmosix-tools.tar.gz -C /
```

قبل از reboot کردن ندها، در هر ند در شاخه /etc/fstab خط زیر را اضافه کنید تا فایل سیستم openMosix در /mfs mount شود. همچنین openMosix به سرویس‌های قابل اجرا در زمان boot اضافه کنید:

```
#cp / openmosix-tools/etc/rc.d/init.d/openmosix /etc/rc.d/init.d/  
#Chkconfig /etc/rc.d/init.d/openmosix on
```

با restart کردن تمام ندها، کلاستر شما آماده است!
حال اگر به یکی از ندها login کنید و دستور mmon را وارد کنید یک نمودار از بار حاضر روی کلاستر برای شما روی صفحه نمایش داده می‌شود.
برای اجرای پروسه‌ها و مدیریت مهاجرت آن‌ها و دیگر امکانات openMosix باید از دستورات مناسب آن استفاده کنید. البته با نصب openMosixview از درگیر شدن با دستورات انبوه رهایی می‌یابید! در ادامه این ابزار مفید شرح داده می‌شود.

۵- تکنولوژی openMosix

تکنولوژی openMosix شامل دو بخش است:

- مکانیزم مهاجرت پروسه قبضه کننده^۶
- مجموعه ای از الگوریتم های اشتراک منابع

⁶ Preemptive Process Migration

هر دو بخش در سطح کرنل، با استفاده از یک ماژول قابل لود پیاده سازی شده اند، به گونه ای رابط کرنل^۷ تغییر نمی کند و برای برنامه های کاربردی کاملاً شفاف است.

۵-۱- مکانیزم مهاجرت پروسه (PPM)

PPM می تواند هر پروسه را در هر زمان و به هر ند آماده منتقل کند. معمولاً مهاجرت ها بر اساس اطلاعاتی است که توسط یکی از الگوریتم های به اشتراک گذاری منابع فراهم می شود، اما کاربران می توانند به طور دستی پروسه ها را به ندهای دلخواه منتقل کنند.

هر پروسه یک ند خانگی یکتا به نام UHN^۸ دارد که بر روی آن ایجاد می شود. تمام پروسه هایی که روی یک ند ایجاد می شوند، محیط اجرایی UHN را به اشتراک می گذارند. پروسه هایی که به ندهای دیگر منتقل می شوند تا حد امکان از منابع ند محلی (ند دور) استفاده می کنند، اما به وسیله UHN با محیط اصلی تعامل دارند.

تا زمانی که درخواست ها برای منابع (پردازنده، حافظه، ...) کمتر از مقدار آستانه^۹ خاصی باشد، پروسه های کاربر محدود به UHN هستند. بعد از گذشتن از آستانه، بعضی پروسه ها به ندهای دیگر منتقل می شوند. هدف کلی حداکثر کردن کارایی با بهره برداری از کل منابع موجود در کلاستر است.

واحد توزیع کار در openMosix پروسه می باشد و در حال حاضر openMosix قادر به توزیع نخ های^{۱۰} یک پروسه نیست.

هنگامی که یک پروسه ایجاد می شود، openMosix آن را به آماده ترین ند اختصاص می دهد. اگر در طول اجرا منابع جدیدی آماده شدند، الگوریتم های اشتراک منابع ممکن است محل پروسه را تغییر دهند.

در openMosix کنترل مرکزی وجود ندارد، هر ند به عنوان یک سیستم خودمختار عمل می کند و تصمیمات کنترلی را مستقلاً اتخاذ می نماید. ندها می توانند به راحتی کلاستر را ترک کنند یا به آن بپیوندند. به علاوه عدم وجود کنترل مرکزی قابلیت توسعه تعداد ندها را فراهم می کند. قابلیت توسعه به وسیله به کارگیری

⁷ Kernel Interface

⁸ Unique Home-Node

⁹ Threshold

¹⁰ Thread

تصادف^{۱۱} در الگوریتم های کنترلی فراهم می شود. هر ند تصمیماتش را بر اساس اطلاعات جزئی که از سایر ندها دارد می گیرد و نمی کوشد که وضعیت جامع کلاستر را تعیین کند. برای مثال در الگوریتم "توزیع اطلاعات احتمالی"^{۱۲} هر ند به طور منظم اطلاعات منابع آزاد خود را به تعدادی از ندها که به طور تصادفی انتخاب می شوند ارسال می کند.

۵-۲- الگوریتم های اشتراك منابع

مهمترین الگوریتم های اشتراك منابع در openMosix، موازنه باز^{۱۳} و راهنمای حافظه^{۱۴} هستند.

۵-۲-۱- الگوریتم موازنه بار

این الگوریتم پیوسته سعی می کند با مهاجرت پروسه ها از ندهای با بار بیش تر به ندهای با بار کمتر، اختلاف بار میان ندها را کاهش دهد. این الگوریتم غیر متمرکز است و تمام ندها الگوریتم یکسانی را اجرا می کنند. تعداد پروسه ها در هر ند و سرعت ندها دو پارامتر مهم برای الگوریتم توزیع بار هستند. الگوریتم به تغییرات بار ندها و خصوصیات زمان اجرای پروسه ها پاسخ می دهد.

مدل ریاضی این الگوریتم بر اساس تحقیقات انجام شده در حوزه اقتصاد است. تعیین محل بهینه برای یک کار مسئله پیچیده است، زیرا منابع آماده در کلاستر ناهمگون^{۱۵} هستند. هزینه های حافظه، پردازنده، ارتباطات پروسه ها، ... غیر قابل مقایسه اند و حتی بر حسب واحدهای یکسان اندازه گیری نمی شوند. منابع ارتباطی بر حسب پهنای باند، حافظه بر حسب فضا و پردازنده بر حسب سیکل اندازه گیری می شوند. الگوریتمی که توسط openMosix به کار گرفته می شود سعی می کند این تفاوت ها را - بر اساس قوانین اقتصادی و تحلیل رقابتی^{۱۶} - با یکدیگر تطبیق دهد.

ایده اصلی تبدیل هزینه های متفاوت استفاده از منابع ناهمگون به یک هزینه همگون^{۱۷} واحد است. بر اساس نتیجه الگوریتم، پروسه به ماشین با کمترین هزینه نسبت

¹¹ Randomness

¹² Probabilistic Information Dissemination

¹³ Load-balancing

¹⁴ Memory Ushering

¹⁵ Heterogeneous

¹⁶ Competitive Analysis

¹⁷ Homogeneous

داده می شود و الگوریتمی تقریباً بهینه برای تخصیص و به اشتراک گذاری منابع است.

openMosix در زمان boot up، سرعت هر ند را برحسب سرعت یک پردازنده Pentium 1GHz محاسبه و در متغیری به نام speed قرار می دهد. الگوریتم توزیع بار از speed در محاسبات خود استفاده می کند. این متغیر را می توان به طور دستی در /proc/hpc تنظیم نمود.

۲-۲-۵- الگوریتم راهنمای حافظه

این الگوریتم پروسه ها را به گونه ای در حافظه کلاستر^{۱۸} جای می دهد که تا حد امکان از کویدگی ناشی از تعویض صفحات حافظه (swapping) جلوگیری شود.

الگوریتم هنگامی راه اندازی^{۱۹} می شود که یک ند شروع به تعویض مدام صفحات حافظه به خاطر کمبود فضا می کند. در این حالت الگوریتم راهنمای حافظه، الگوریتم توزیع بار را لغو^{۲۰} و سعی می کند پروسه را به ندی که حافظه کافی دارد منتقل نماید، حتی اگر منجر به توزیع بار غیریکنواخت شود.

۳-۵- مهاجرت پروسه (PPM)

openMosix از مهاجرت پروسه به طور قبضه شونده و کاملاً شفاف پشتیبانی می کند. بعد از مهاجرت پروسه به تعامل با محیطش بدون در نظر گرفتن مکان جدید ادامه می دهد.

برای پیاده سازی PPM پروسه به دو بخش تقسیم می شود: زمینه کاربر^{۲۱} که می تواند بر روی ند‌ها جابه جا شود و زمینه سیستم^{۲۲} که وابسته به UHN است و منتقل نمی شود. زمینه کاربر یا remote شامل کد برنامه، داده، نقشه های حافظه و رجیسترهای پروسه است. زمینه سیستم یا deputy شامل شرح منابع وابسته به UHN پروسه و یک استک کرنل برای اجرای کدهای سیستمی از طرف پروسه است. رابط میان remote و deputy کاملاً تعریف شده است.

زمان مهاجرت یک بخش ثابت دارد برای ایجاد ساختارهای لازم اجرای پروسه در remote و یک زمان خطی که متناسب با تعداد صفحات حافظه است که

¹⁸ Cluster-wide Memory

¹⁹ Trigger

²⁰ Override

²¹ User Context

²² System Context

منتقل می شود. به منظور حداقل کردن سربار مهاجرت، فقط جداول حافظه و صفحات تغییر یافته^{۲۳} پروسه منتقل می شوند.

در زمان اجرای پروسه، فراخوانی های سیستمی پروسه بررسی می شوند. اگر فراخوانی وابسته به UHN باشند به deputy ارسال شده، deputy آن را از طرف پروسه اجرا و نتیجه را برمی گرداند. در غیر این صورت بر روی remote اجرا می گردند. این نوع ارتباط میان remote و deputy همگام^{۲۴} است.

نوع دیگر تعامل میان remote و deputy ارسال سیگنال - به عنوان مثال هنگام دریافت داده های شبکه در UHN- است. این ارتباط غیر همگام^{۲۵} می باشد.

۵-۴- نکاتی در رابطه با مکانیزم remote/deputy

همان طور که گفته شد deputy نماینده پروسه منتقل شده در UHN است. چون تمام فضای حافظه در ند remote قرار دارد، deputy نقشه حافظه پروسه مهاجرت یافته را نگه نمی دارد. در اجرای بسیاری از فراخوانی های سیستمی لازم است که داده بین فضای کاربر و کرنل منتقل شود. در openMosix هنگامی که چنین فراخوانی های سیستمی از remote به deputy ارسال می شود، deputy با remote برای انتقال داده ارتباط برقرار می کند. سربار ناشی از کپی کردن داده از remote به deputy می تواند بسیار زیاد باشد. برای بهبود کارایی، یک کش ویژه پیاده سازی شده است که تا حد امکان داده ها را در شروع فراخوانی پیش واکشی^{۲۶} می نماید. برای جلوگیری از حذف یا رونویسی صفحات حافظه (هنگام تعویض صفحات) به علت عدم وجود نقشه حافظه برای پروسه منتقل شده، deputy جدول ویژه ای برای اینگونه صفحات حافظه نگه دارد.

پروسه های مهمان و حافظه آنها در ند remote قابل دستیابی برای پروسه های دیگر آن ند و بالعکس نیستند.

ممکن است نیاز باشد یک پروسه در حالتی که متوقف است بعضی توابع openMosix را اجرا کند. به عنوان مثال مهاجرت پروسه می تواند در زمان توقف پروسه اتفاق بیافتد. پروسه می بایست تابع را اجرا کرده و مجدداً متوقف شود. به این منظور openMosix یک وضعیت منطقی از پروسه ها نگه می دارد که نشان می دهد

²³ Dirty Pages

²⁴ Synchronous

²⁵ Asynchronous

²⁶ Pre-Fetching

بقیه پروسه ها وضعیت پروسه مذکور را چگونه باید ببینند و از دیدن حالت اجرای موقتی پروسه جلوگیری می کند.

۵-۵. سیستم فایل openMosix

ساختار بیشتر سیستم های فایل کلاسترها مبتنی بر یک فایل سرور مرکزی است. openMosix از فایل سیستمی به نام DFSA (Direct File System access) استفاده می کند. DFSA برای کاهش سربار ناشی از اجرای فراخوانی های I/O oriented یک پروسه مهاجرت یافته طراحی شده است.

DFSA بر روی سیستم فایل کلاستر اجرا می شود. سیستم فایل هایی که DFSA می تواند بر روی آنها اجرا شود MFS، GFS و QFS هستند. در حالت عادی هر I/O به UHN برای اجرا فرستاده می شود. اما با استفاده از DFSA اگر فایل مورد نظر بر روی همان ندی باشد که پروسه اجرا می شود، I/O محلی انجام می شود.

علاوه بر DFSA الگوریتم دیگری که عملیات I/O پروسه ها را نیز در نظر می گیرد، به الگوریتم توزیع بار اضافه شده است. در این الگوریتم سعی می شود پروسه به ندی منتقل شود که بیشترین عملیات I/O خود را در آن جا انجام می دهد. در نتیجه برخلاف سیستم فایل های شبکه ای موجود، مانند NSF، که داده را از فایل سرور به ند کلاینت می آورد، کلاستر openMosix سعی می کند پروسه را به ندی منتقل کند که فایل واقعا در آن جا قرار دارد.

۵-۶. API openMosix

API openMosix در ابتدا به وسیله تعدادی فراخوانی سیستمی پیاده سازی شده بود که برای تنظیم، سوال پرسیدن^{۲۷} و راه اندازی openMosix به کار می رفت. به جهت هماهنگی با قرارداد لینوکس، API به گونه ای تغییر یافت که از طریق سیستم فایل مجازی /proc قابل دستیابی باشد. (proc/ یک سیستم فایل مجازی است که از طریق آن می توان اطلاعات کرنل و پروسه های در حال اجرا را مشاهده کرد.)

در openMosix مسیر /proc با افزودن /proc/openMosix توسعه یافته است. از طریق /proc/openMosix فراخوانی های زیر را می توان انجام داد:

- درخواست مهاجرت پروسه

²⁷ Query

- قفل کردن یک پروسه برای جلوگیری از مهاجرت
- یافتن محل اجرای فعلی پروسه
- مدیریت سیستم
- کنترل جمع آوری آمار پروسه های در حال اجرا
- اطلاعات در مورد پروسه های remote
- اطلاعات در مورد محدودیت های مهاجرت

۶- معرفی openMosixview

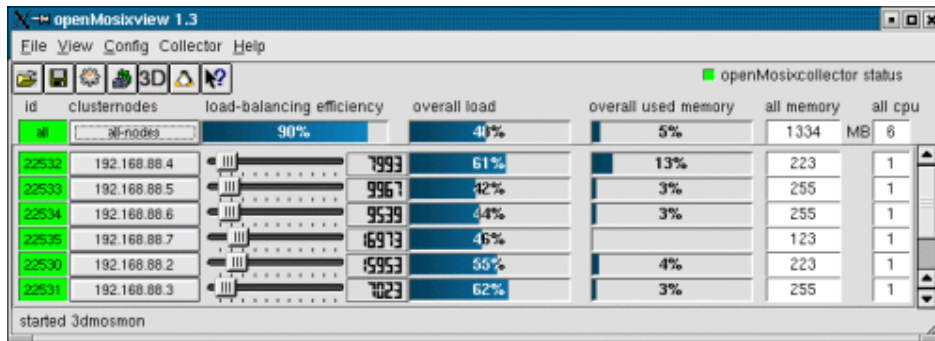
openMosixview یک رابط گرافیکی برای مدیریت کلاسترهای openMosix است. openMosixview نسخه تکمیل شده mosixview می باشد. این بسته شامل ۸ کاربرد مختلف برای نظارت و مدیریت پروسه ها و کلاستر است:

- openMosixprocs : برای مدیریت پروسه ها.
- openMosixcollector : سرویس جمع آوری کننده که اطلاعات ندهای کلاستر را ثبت می کند.
- openMosixanalyzer : برای تحلیل داده های جمع آوری شده توسط ماژول قبل به کار می رود.
- openMosixhistory : تاریخچه پروسه های اجرا شده روی کلاستر.
- openMosixmigmon : جهت مدیریت مهاجرت پروسه ها به کار می رود.
- openMosixpidlog : برای مدیریت یک پروسه خاص به کار می رود.
- 3dmosmon : یک چشم انداز سه بعدی از وضعیت کلاستر.

تمام این کاربردها از صفحه اصلی openMosixview قابل دسترسی است. توسط این بسته مفید تمام اعمالی که از طریق دستورهای مختلف خط فرمان صورت می گیرند، توسط چند کلیک قابل اجرا هستند. در ادامه تعدادی از این ۸ کاربرد توضیح داده می شود.

۶-۱- openMosixview

شکل ۱ تصویری از صفحه اصلی openMosixview را مشاهده می کنید.



شکل ۱- صفحه اصلی openMosixview

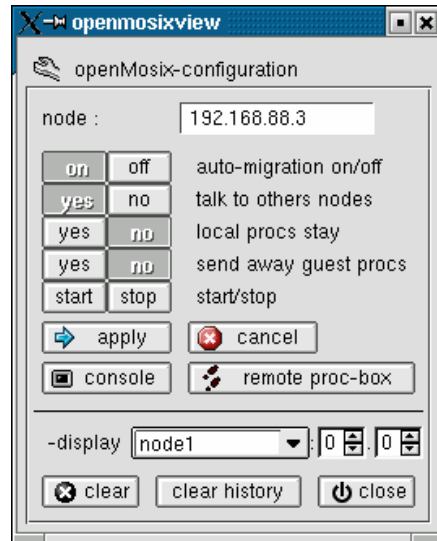
برای هر ند در کلاستر یک سطر حاوی اطلاعات مفیدی راجع به آن ند وجود دارد. در هر سطر اولین خانه چراغی است که openMosix-Id و حالت ند را نشان می دهد، در صورتی که ند در دسترس باشد سبز و در غیر این صورت قرمز است. خانه بعد آدرس IP ند مورد نظر است، اگر بر روی آن کلیک کنید یک دیالوگ تنظیمات ظاهر می شود. slider که بعد از آن مشاهده می کنید جهت تنظیم سرعت openMosix برای هر میزبان به کار می رود و خانه بعدی سرعت را به صورت عددی نشان می دهد. با تغییر این مقدار روی چگونگی موازنه بار کل کلاستر تأثیر می گذارید. پروسه ها روی کلاستر openMosix بیشتر به ندی که سرعتش بیشتر است مهاجرت می کنند. نکته مهم این است که این مقدار می تواند نشان دهنده سرعت فیزیکی نباشد، بلکه سرعتی است که openMosix برای هر ند در نظر می گیرد. روی هر سطر همچنین دو نوار پیشرفتی ملاحظه می شود که اولی درصد بار روی ند و دیگری درصد حافظه جاری مورد استفاده ند را نشان می دهد. خانه بعدی نشان دهنده کل حافظه در دسترس روی هر ند را است. در انتها تعداد پردازنده در هر ند نوشته شده است.

اولین سطر در صفحه اصلی یک دکمه برای تنظیم تمام ندها دارد. با این امکان، می توان ندها را مدیریت و تنظیم کرد. نوار پیشرفت موجود در سطر اول نشان می دهد که کلاستر با چه کیفیتی بار را روی ندها موازنه کرده است. ۱۰۰٪ حالتی است که تمام ندها تقریباً بار یکسان دارند.

در پنجره تنظیمات هر ند (شکل ۲) که از صفحه اصلی قابل دسترسی است، تنظیماتی روی ندها صورت می گیرد که از طریق rsh و ssh این تغییرات اعمال می شوند. دستوراتی که از این صفحه قابل اجرا هستند شامل موارد زیر می باشند:

- automigration on/off
- quiet yes/no
- bring/!stay yes/no

- expsel yes/no
- openMosix start/stop

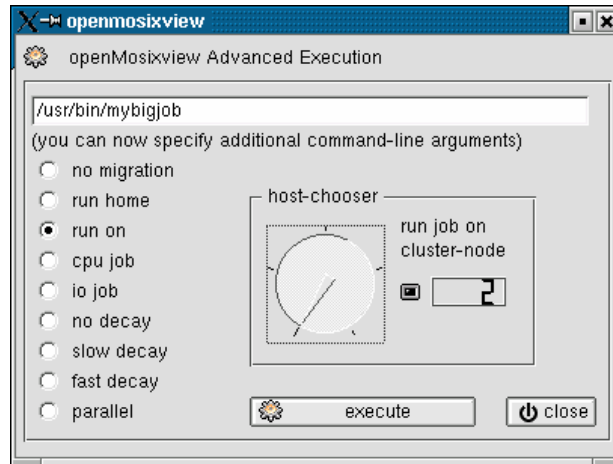


شکل ۲- پنجره تنظیمات مربوط به هر نود

اجرای پروسه‌ها

دیالوگ Advanced Execution برای اجرای برنامه‌ها به کار می‌رود. آن

را در شکل ۳ مشاهده می‌کنید.



شکل ۳- دیالوگ Advanced Execution

برای باز شدن این پنجره بر آیکون open از صفحه اصلی کلیک کنید و پروسه موردنظران را از شاخه مربوطه‌اش انتخاب کنید. حال می‌توانید مشخص کنید که

این پروسه، چگونه و کجا آغاز شود. گزینه های مختلفی برای اجرا وجود دارد. گزینه هایی که در این پنجره وجود دارند عبارتند از:

- no migration: یک کار محلی را آغاز می کند که مهاجرت نخواهد کرد.
 - run home: یک کار محلی را شروع می کند.
 - run on: کار را روی ندی که شما از طریق host-chooser انتخاب می کنید آغاز می کند.
 - cpu job: یک کار با محاسبات بسیار (computation intensive) را روی یک ند آغاز می کند. (host-chooser)
 - io job: یک کار با اعمال ورودی/خروجی بسیار (computation intensive) را روی یک ند آغاز می کند. (host-chooser)
 - parallel: یک کار را به طور موازی روی چند ند یا تمام ندها آغاز می کند.
- البته غیر از این گزینه ها می توانید آرگومان های خط فرمان دیگری را هم در textbox بالای پنجره مشخص کنید.

برای اجرای یک پروسه روی ند خاصی، توسط host-chooser ند مربوطه را مشخص می کنیم. این ابزار dای ندها را نشان می دهد و شما می توانید ند موردنظر را بدین ترتیب انتخاب کنید. برای اجرای یک پروسه روی چند ند به طور موازی، توسط دو علامتگذاری روی ند اول و آخر، پروسه روی تمام ندهای بین ندهای علامتگذاری شده و خود ندهای علامتگذاری شده اجرا خواهد شد.

۶-۲ - openMosixprocs

دیالوگی که لیست تمام پروسه های روی نود را نشان می دهد از صفحه اصلی openMosixview قابل دسترسی است. این دیالوگ برای مدیریت پروسه هایی که از روی نود حاضر اجرا شده اند و در حال اجرا هستند به کار می رود.

pid	n#	lock	nmigs	stat	cmdline	nice	UID
12075	22535	0	2	S	/distkeygen	0	0
12074	22535	0	3	S	/distkeygen	0	0
12072	22535	0	0	S	/distkeygen	0	0
12068	22534	0	3	S	/distkeygen	0	0
12069	22533	0	3	S	/distkeygen	0	0
12067	22533	0	3	S	/distkeygen	0	0
12070	22531	0	2	S	/distkeygen	0	0
11986	22531	0	3	S	/bin/bash	0	0
12073	22530	0	1	S	/distkeygen	0	0
12071	22530	0	1	S	/distkeygen	0	0
12066	22530	0	2	S	/distkeygen	0	0
983	0	1	0	S	/usr/sbin/atd	0	0
947	0	1	0	S	xfs	0	43
852	0	1	0	S	crond	0	0
832	0	1	0	S	aam	0	0

manage procs from remote 67 processes on this system quit

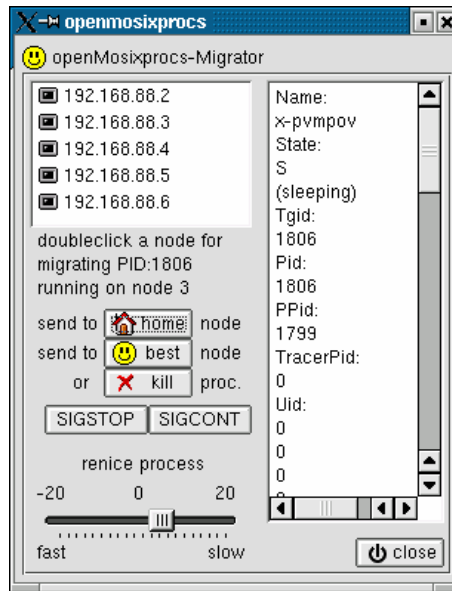
شکل ۴- لیست تمام پروسه‌های در حال اجرا که از نود حاضر اجرا شده‌اند.

لیست پروسه‌ها در این دیالوگ نشان می‌دهد که هر پروسه روی چه ندی در حال اجراست. ستون دوم dای ندی را نشان می‌دهد که پروسه روی آن در حال اجراست. شماره 0 بدین معنی است که پروسه روی ند محلی خود در حال اجراست، بقیه شماره‌ها مربوط به ندهای دور هستند. پروسه‌هایی که مهاجرت کرده‌اند با آیکن سبز علامتگذاری شده‌اند.

با دو کلیک روی یک پروسه، پنجره مدیریت مهاجرت پروسه باز می‌شود. (

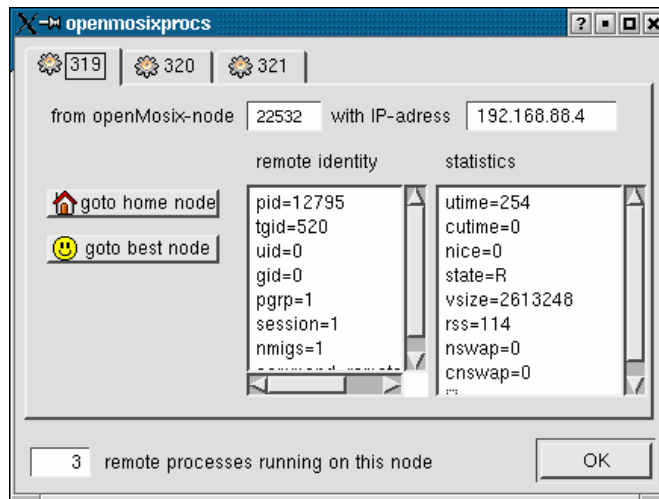
شکل ۵)

این پنجره تمام ندهای روی کلاستر را نشان می‌دهد و برای مدیریت یک پروسه است. با دوبار کلیک روی یک ند، پروسه به آن ند مهاجرت می‌کند. دکمه home پروسه را به ند محلی خود منتقل می‌کند. دکمه best پروسه را به بهترین ند در دسترس منتقل می‌کند.



شکل ۵ - دیالوگ openMosixprocs

با دکمه kill پروسه از بین می‌رود. برای متوقف کردن موقتی پروسه دکمه SIGSTOP و برای اجرای دوباره آن SIGCONT را کلیک کنید. روی دیالوگ openMosixprocs دکمه دیگری با نام manage procs from remote وجود دارد.



شکل ۶- مدیریت پروسه ها از دور

در پنجره مربوطه (شکل ۶) هر یک از برگه‌های بالای پنجره یکی از پروسه‌هایی است که به ند محلی حاضر مهاجرت کرده است.

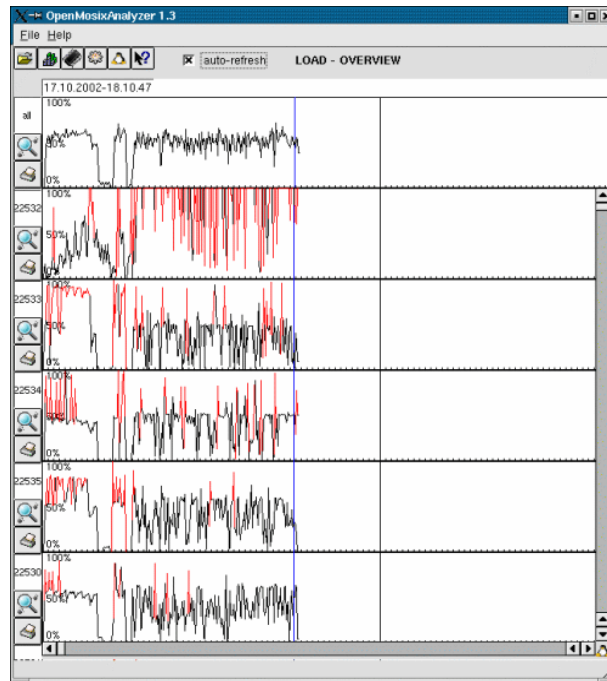
۳-۶ - openMosixcollector

این سرویس می تواند روی هر یک از اعضای کلاستر اجرا شود. openMosixcollector اطلاعات مربوط به بار روی هر ند را در شاخه /tmp/openmosixcollector/* ثبت می کند. فایل های ثبت مربوطه توسط openMosixanalyzer تحلیل و بررسی می شوند. openMosixcollector هر ۱۲ ساعت یک بار اجرا می شود و تاریخچه حاضر از اطلاعات بار را در /tmp/openmosixcollector[date]/* ذخیره می کند. می توان یک نقطه بازرسی (checkpoint) در تاریخچه قرار داد. این نقطه ها به صورت خط های عمودی آبی در رابط گرافیکی openMosixanalyzer علامتگذاری خواهند شد، برای مثال شما می توانید وقتی یک پروسه را اجرا می کنید یک نقطه بررسی قرار دهید و هنگام خاتمه یک نقطه دیگر. دستورات موجود برای این سرویس عبارتند از:

- `openmosixcollector -d` : collector را به عنوان یک سرویس اجرا می کند.
- `openmosixcollector -k` : سرویس collector را خاتمه می دهد.
- `openmosixcollector -n` : یک نقطه بازرسی در تاریخچه قرار می دهد.
- `openmosixcollector -r` : تاریخچه حاضر را ذخیره می کند و یک تاریخچه جدید را شروع می کند.

۴-۶ - OpenMosixanalyzer

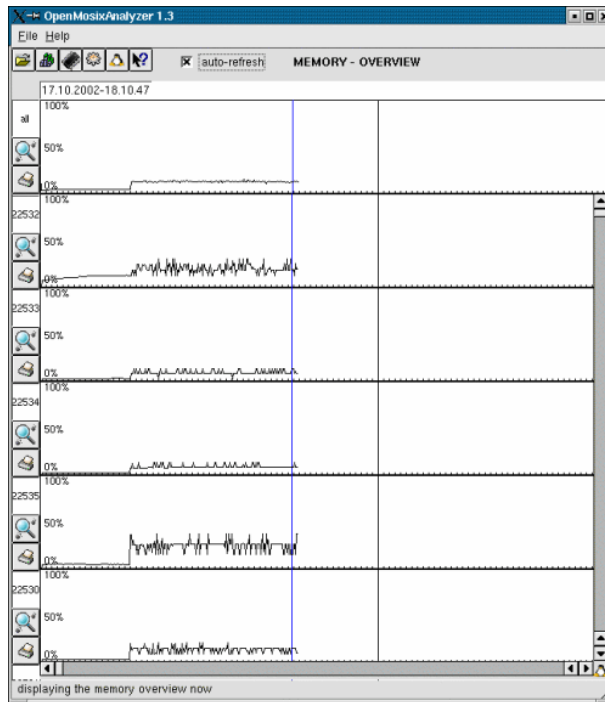
شکل ۷ میزان بار ندها را به طور گرافیکی در openMosixanalyzer نشان می دهد.



شکل ۷- دیاگونگ openMosixanalyzer ، برگه نمودار میزان بار روی نودها

شکلی که نشان داده می شود به طور پیش فرض از فایل ثبت مربوط به بار روی ندها در زمان حاضر است. در صورتی که بخواهیم تاریخچه ای از گذشته را مشاهده کنیم باید فایل ثبت مربوطه آن را از `/tmp/openmosixcollector[date]/` انتخاب کنیم. خط های رسم شده در نمودارها به طور عادی مشکی هستند، ولی در صورتی که بار بیشتر از ۷۵٪ باشد خطها قرمز رسم شده اند. دکمه `find out` مربوط به هر ند مقادیر آماری مفیدی را برای آن ند محاسبه می کند و نمایش می دهد.

همچنین از طریق openMosixanalyzer می توان نمودارهای مربوط به حافظه مورد استفاده پروسه ها در کلاستر را مشاهده کرد.



شکل ۸- دیالوگ openMosixanalyzer ، برگه نمودار میزان حافظه مورد استفاده روی نودها

امکان دیگری که در openMosixanalyzer فراهم آمده مربوط به پروسه‌هاست. در این قسمت اطلاعات بسیاری با جزئیات در رابطه با این که چه پروسه ای، در چه زمان و در چه ندی در حال اجرا بوده است به دست می آید. (شکل ۹)

pid	n#	lock	nmigs	stat	cmdline	nice	UID
12880	22535	0	2	S	/setiathome	-20	0
12874	22534	0	4	S	/setiathome	-20	0
12873	22534	0	7	S	/setiathome	-20	0
12877	22533	0	7	S	/setiathome	-20	0
12876	22531	0	3	S	/setiathome	-20	0
12879	22530	0	5	S	/setiathome	-20	0
12875	22530	0	3	S	/setiathome	-20	0
983	0	1	0	S	/usr/sbin/atd	0	0
947	0	1	0	S	xfs	0	43
852	0	1	0	S	crond	0	0
832	0	1	0	S	gpm	0	0
831	0	1	0	S	-bash	0	0
812	0	1	0	S	sendmail:	0	0
804	0	1	0	S	login	0	0
798	0	1	0	S	ip_dlogind	0	0

شکل ۹- دیالوگ openMosixhistory

زمان را می‌توان با نوار پیشرفتی که در بالای دیالوگ وجود دارد تغییر داد.

۶-۵ - openMosixmigmon

این برنامه یکی از کاربردهای جالب دیگر از openMosixview است. توسط این قسمت می‌توانید تمام ندهای حاضر در کلاستر را با پروسه‌های در حال اجرا روی آن‌ها در یک نمای گرافیکی مشاهده کنید. هر ند به صورت یک پنگوئن که در یک دایره نشسته است دیده می‌شود. پنگوئن اصلی ندی است که openmosixmigmon روی آن اجرا شده است. اگر یک پروسه به ندی مهاجرت کند، پروسه از دایره اصلی خود به دایره مربوط به ند دور منتقل می‌شود، سپس پروسه توسط رنگ سبز علامتگذاری می‌شود و یک خط از ند محلی به ند دور برای نشان دادن مهاجرت رسم می‌شود. اگر موس را روی پروسه نگه‌دارید پروسه نشان داده می‌شود. همچنین می‌توان به راحتی با کشیدن (drag) یک پروسه توسط موس و رها کردن (drop) آن در ندی دیگر به راحتی پروسه را به ند مربوطه مهاجرت داد. اگر پروسه ای را دوبار کلیک کنید به ند محلی خودش برمی‌گردد.

۷- جمع بندی

با توجه به روند رو رشد کاربردهای محاسباتی سنگین، نیاز به تکنولوژی‌های پردازش موازی افزایش می‌یابد. از جمله راه‌حل‌های موجود کلاسترینگ است. با استفاده از سیستم عامل لینوکس و openMosix، می‌توان تعداد دلخواهی از کامپیوترها را تبدیل به کلاستر نمود. امروزه بیش از ۳۰۰۰ کلاستر openMosix در دنیا وجود دارد. متوسط ندهای کلاسترها ۲۵ عدد است و به راحتی تا ۲۰۰۰ عدد گسترش می‌یابد. openMosix دارای جامعه‌ای فعال و حمایتگر از توسعه دهندگان و استفاده‌کنندگان از آن می‌باشد.

- **OpenMosix, Presented by Dr. Moshe Bar et al. Linux Congress 2003.**
- **OpenMosix, OpenSSI and Kerrighed: A Comparative Study. Renaud Lottiaux Pascal Gallard, Geoffroy Valle, Christine Morin. 2005 IEEE International Symposium on Cluster Computing and the Grid.**
- **Migratable sockets in cluster computing, Kamran Malik, Osama Khan, Tahir Mobashir, Mansoor Sarwar. 2004 Elsevier.**
- **openMosix vs. Beowulf: a case study, Moshe Bar, Stefano Gozzini, Alberto Marmodora.**
- **Wikipedia Encyclopedia: en.wikipedia.org**
- **openMosix's Home: www.openMosix.org**
- **openMosixview Home: www.openmosixview.com**