

# IBM Platform LSF培训

[上海天文台]

项目编号： PPSNxxxxxx

[service@paratera.com](mailto:service@paratera.com)

The Right Answer in HPC  
**PARATERA 并行**

# AGENDA

- 计算环境介绍
- LSF基本介绍
- 如何使用LSF
- Trouble Shooting



# 计算环境

- 管理/登录节点: **bright60**      IP: **119.78.226.16**
- 计算节点: **node001-node050**
- **lsf**安装在**/cm/shared/apps/lsf**下;
- 应用软件安装在**/cm/shared/apps**下;

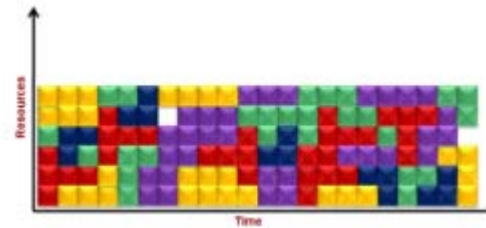
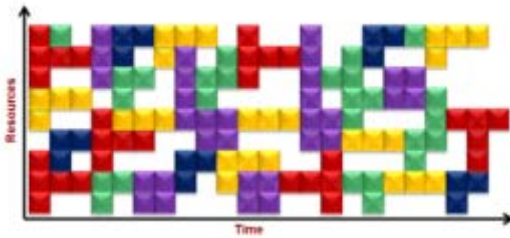
# AGENDA

- 计算环境介绍
- **LSF**基本介绍
- 如何使用LSF
- Trouble Shooting

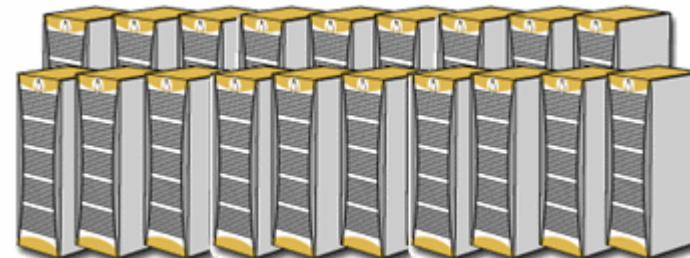


# Platform LSF

By scheduling workloads intelligently according to policy, Platform LSF reduces application run-times while simultaneously optimizing resource use.

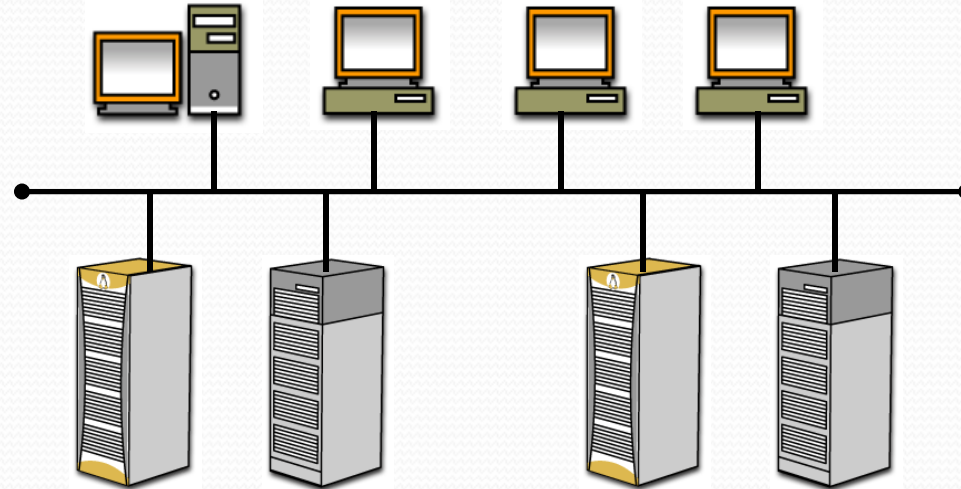
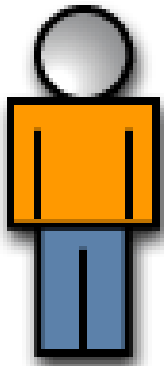


VIRTUALIZED VIEW OF COMPUTE, NETWORK AND STORAGE RESOURCES



# Platform LSF

Which node  
can run my  
job or task?

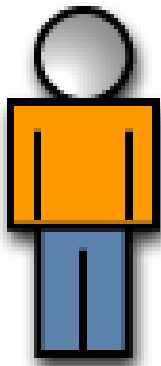


In a Distributed environment (hundreds of hosts)

- Monitoring and control of resources is complex
- Resource usage imbalance
- Users perceive a lack of resources

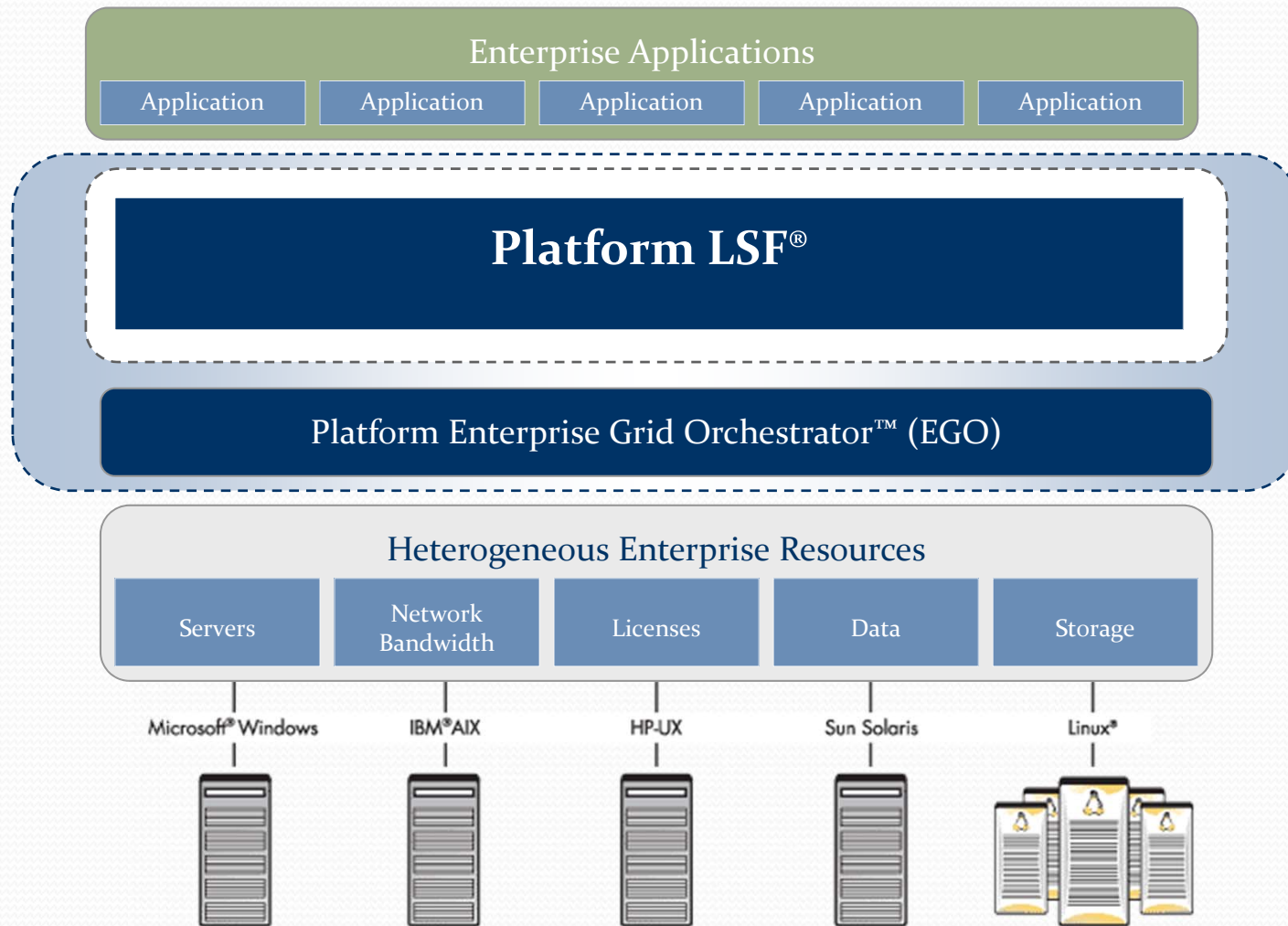
# Platform LSF

Now, Platform LSF will run my job or task on the best node available!



Virtual Pool of computing resources managed by Platform LSF

# Platform LSF





# AGENDA

- 计算环境介绍
- LSF基本介绍
- 如何使用LSF
- Trouble Shooting





# Job Control Command

# Setting up LSF Environment

- Setup the LSF environment before using LSF commands
- Bourne like shells (sh, bash, ksh, etc.)
- `$ . /cm/shared/apps/lsf/conf/profile.lsf`
- Check after setting up environment:
  - `% env | grep -i lsf`
  - `LSF_BINDIR=/cm/shared/apps/lsf/8.3/sparc-sol10-64/bin`
  - `LSF_SERVERDIR=/cm/shared/apps/lsf/8.3/sparc-sol10-64/etc`
  - `LSF_LIBDIR=/cm/shared/apps/lsf/8.3/sparc-sol10-64/lib`
  - `LD_LIBRARY_PATH=/cm/shared/apps/lsf/8.3/sparc-sol10-64/lib`
  - `XLSF_UIDDIR=/cm/shared/apps/lsf/8.3/sparc-sol10-64/lib/uid`
  - `LSF_ENVDIR=/cm/shared/apps/lsf/conf`

# Job Control Command

- `bsub [options] command [cmdargs]`
- `bjobs [-a][-J jobname][-u usergroup|-u all][...] jobID`
- `bhosts [-l/-w][...]`
- `bbot/btop [jobID | "jobID[index_list]"] [position]`
- `bkill [-J jobname] [-m] [-u ] [-q] [-s signalvalue]`
- `bmod [bsub_options] jobID`
- `bpeek [-f] jobID`
- `bstop/bresume jobID`

# Bsub

`bsub` 命令常见用法如下：

**`bsub -n z -q QUEUENAME -J test -o outputfile COMMAND`**

**-n:** 其中`z`代表了提交作业需要的cpu数；

**-q:** 指定作业提交到的队列，如果不采用`-q`选项，系统把作业提交到默认作业队列`normal`，其中一共有两个队列`normal`和`mpi`，**`normal`**队列用来运行串行作业，**`mpi`**用来运行并行作业；

**-J:** 指定作业名称为`test`，缺省为命令名称

**-o:** `outputfile` 代表一个文件，作业提交后标准输出的信息将会保存到这个文件中。

**COMMAND:** 是用户要运行的程序。

**-I:** 运行交互式的作业。

**-a:** 指定作业相关的应用

**-R:** 指定运行作业需要的相关资源

# Bsub

- By script or command
  - `% cd /home/user/project_dir`
  - `% bsub -q mpi -a fluent -n 4 ./my_fluent_launcher.sh`
- By job spooling
  - `% bsub < spoolfile`
- Interactively
  - `% bsub`
  - `bsub> #BSUB -q parallel -n 4`
  - `bsub> #BSUB -a fluent`
  - `bsub> cd /home/user/project_dir`
  - `bsub> ./my_fluent_launcher.sh`
  - `bsub> ^D`
  - `Job <1234> submitted to queue <parallel>`

## Example *spoolfile*

```
#BSUB -q parallel
#BSUB -n 4 -a fluent
cd /home/user/project_dir
./my_fluent_launcher.sh
```

# Bjobs

## bjobs查看作业

普通用户执行bjobs查看自己的作业

```
[efadmin@redhat62 ~]$ bjobs
JOBID  USER  STAT  QUEUE      FROM_HOST  EXEC_HOST  JOB_NAME  SUBMIT_TIME
1956   efadmin RUN   normal    redhat62   redhat62   test      Dec  4 20:55
```

JOBID 为作业号，每个作业有唯一的作业号

USER 作业所属的用户

STAT 作业状态，RUN表示在运行，PEND表示在排队，DONE表示正常完成，EXIT表示异常退出

QUEUE 作业所在队列

EXEC\_HOST 执行作业的节点

JOB\_NAME 作业名

SUBMIT\_TIME 提交作业的时间

bjobs-l 查看作业详细信息

# Bhost

## bhosts查看节点信息

```
[root@bright60 logs]# bhosts
HOST_NAME          STATUS      JL/U    MAX  NJOBS    RUN  SSUSP  USUSP  RSV
bright60.cm.cluste ok          -      12    0      0      0      0      0
node001.cm.cluster ok          -      12    0      0      0      0      0
node002.cm.cluster ok          -      12    0      0      0      0      0
node003.cm.cluster ok          -      12    0      0      0      0      0
node004.cm.cluster ok          -      12    0      0      0      0      0
node005.cm.cluster ok          -      12    0      0      0      0      0
node006.cm.cluster ok          -      12    0      0      0      0      0
node007.cm.cluster ok          -      12    0      0      0      0      0
node008.cm.cluster ok          -      12    0      0      0      0      0
```

STATUS 是节点状态，ok表示正常，可以接受用户提交作业；unavail表示节点lim进程不正常，不能接受用户提交作业；closed，表示节点所有cpu核都已用满或该节点被管理员关闭，此状态下该节点不再接受新作业；

MAX 是节点的cpu核数量

NJOBS 是所有作业在该节点上申请的cpu核数量

RUN 是该节点运行lsf作业的核数



# Other Job Control Commands

- **bbot** – moves a pending job to the bottom of the queue
- **btop** – moves a pending job to the top of the queue
- **kill** – sends a signal to kill, suspend or resume unfinished jobs (use a job ID of “o” to kill all your jobs). New scalability improvements resulting in improved performance and user experience
- **bpeek** – displays the stdout and stderr of an unfinished job
- **bstop** – suspend unfinished jobs
- **bresume** – resumes one or more suspended jobs



# Job Control Example

# Job Submit and Monit

- Set LSF Environment

```
$ . /cm/shared/apps/lsf/conf/profile.lsf
```

- Submit Jobs

```
%bsub -o %J.out "myjob"(myjob is your job command)
```

```
%vi myjob (edit myjob content)
```

```
#!/bin/sh
```

```
df -k
```

- Monit Jobs

```
%bjobs
```

```
%bpeek 1234 (1234 is job ID)
```



# Job Submit and Monit

## ➤ Manage Jobs

%bkill 1234 (1234 is job ID, terminate jobs)

%bstop 1234 (1234 is job ID, pause jobs)

%bresume 1234 (1234 is job ID, resume jobs)

## ➤ Check System

%lshosts (Server Configuration)

%lsload (Server Load)

%bqueues (Queue Status)

%bhosts (Server Job Status)

## ➤ History jobs

%bhist

# Example – HPL

## **hpl.lsf:**

```
#BSUB -J amber_high (作业名称)
#BSUB -o %J.out      (输出结果)
#BSUB -e %J.err      (输出错误信息)
#BSUB -a intelmpi    (指定MPI编译器)
#BSUB -n 16
mpirun ./athena      (athena是应用程序名称)
```

**Run the following command:**

```
bsub < hpl.lsf
```

## 计算节点配置：

- 每刀片2个Intel Xeon5650 2.66GHz 6核处理器
- 每节点配置24GB内存（6条 4GB DDR3RECC 内存）
- 1块500G 3.5' SATA本地硬盘
- 1块DDR IB HBA卡（AOC-IBH-XDS刀片专用IB卡，DDR IB 20G/s接口速度）
- 集成两个千兆网口
- OS Redhat Linux 6.2 64Bit Server Edition
- 计算节点包括2套网络：千兆作业调度管理网络，Infiniband并行互连网络

# 软件使用

Intel® Cluster Studio for  
MPICH,MPICH2,OPENMPI,MVAPICH,MVAPICH2

IDL交互式数据处理开发语言V7.0

Intel C,C++,FORTRAN编译器，使用时需设置环境变量  
> Source /cm/shared/apps/intel-compiler/composerxe-  
2011.5.220/bin/compilervars.sh intel64



# Rules

- Normal user can only do operations on the login node.
- do not run program in login node.
- Use the template job submit script.
- Do not submit jobs to the queues that you are not permitted
- **Do not even try to run jobs pass LSF system!!!**

# AGENDA

- 计算环境介绍
- LSF基本介绍
- 如何使用LSF
- **Trouble Shooting**



# Trouble Shooting

- “Job rejected by LSF”
  - Check the resource require, Check runtime limit
  - Job was submit to a undefined queue or host
  - Job length beyond to the queue length

# Trouble Shooting

- “Job always be pend by LSF”
  - Resource requirement beyond system configuration?  
ex. Memory requirement beyond server memory
  - No user ID on the execution host?
  - Other user may use exclusive operation
  - FCFS police, job was in queue
  - Fairshare police, user have exhaust the resource
  - Use `bjobs -lp` check the reason

# Trouble Shooting

- “Job was kill by LSF”
  - Check queue resource limit
  - Check the execution host can visit the data
  - Check the license
  - Use `bjobs -l` get exit code
  - Exit Code
    - 127 – Can not find command
    - 128 – Can not execute command
    - 130 – Job was terminated by Control-C

# 参考文档

- 官方文档  
lsf\_users\_guide
- 使用手册  
LSF使用手册
- Para文档
  - Paramon用户手册
  - Paratune用户手册