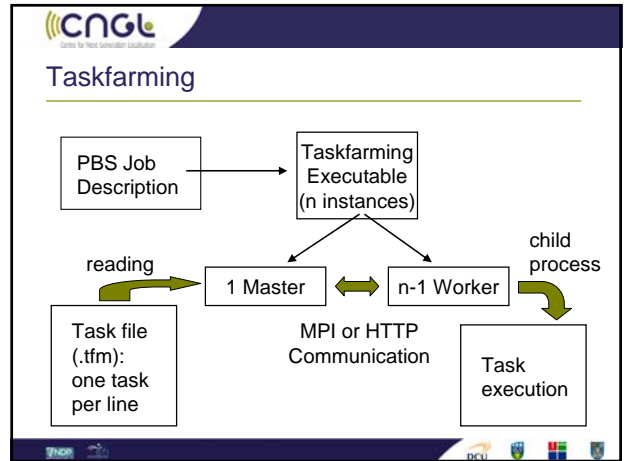


CNGL
Centre for Next Generation Localisation

CPU-Intensive Jobs

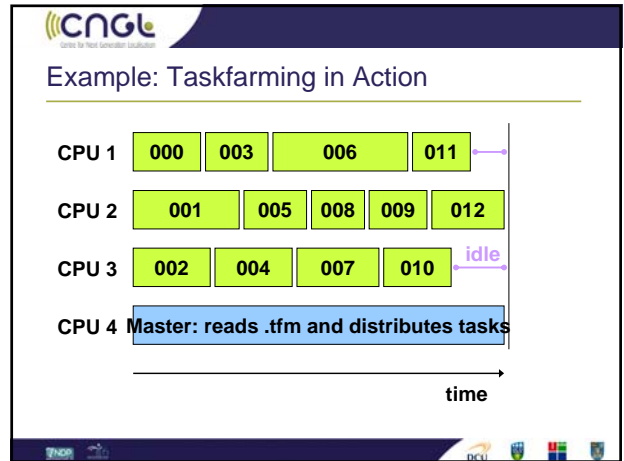
- Parallelisable, for example
 - Sentence by sentence processing
 - Cross-validation runs
 - Parameter search
- Split into parts
 - Run each part on a different CPU core
- Alternatives
 - Submit large number of jobs (ppn=1)
 - Taskfarming



CNGL
Centre for Next Generation Localisation

Taskfarming Options

- Using individual PBS jobs
 - Can only allocate resources in multiples of 1/8 or 1/4
 - Example: 3 GB task -> 4 GB job (ppn=2:cores8:mem16GB)
 - Floods job queue
- MPI-based taskfarming
 - All tasks inside one job
 - Example: 3 GB task -> 5 workers per 16 GB node
 - Master blocks one CPU core
- HTTP/XML-RPC-based taskfarming
 - Master runs on maia login node
 - Workers can run in multiple jobs
 - Example: 3 GB tasks -> one job with 5 workers for 16 GB nodes and one job with 8 workers for 32 GB nodes



CNGL
Centre for Next Generation Localisation

Estimating the PBS Walltime Parameter

- Collect durations from test run
- Usually high variance of execution time
 - Long sentences
 - Parameters
- Don't use #packages x avg. time per package
 - High risk (~50 %) that more time is needed
 - Prefix jobname with, for example, "24h-"
- Random sampling with observed package durations: /home/jwagner/tools/walltime.py

CNGL
Centre for Next Generation Localisation

Questions

?

Contact:
Joachim Wagner
CNGL System Administrator
jwagner@computing.dcu.ie
(01) 700 6915

Installed Software

- OpenMPI
- SRILM
- MaTrEx, Moses, GIZA++
- XLE, Sicstus
- Johnson & Charniak's reranking parser
- In progress:
 - LFG AA, incl. function labeller

PBS Job Management

```

    graph TD
      A[User: Job Description] --> B[Job Submission]
      B --> C[Job Queue]
      C --> D[Job Scheduler]
      D --> E[Job Execution]
  
```

Nodes are allocated job-exclusive for the duration of the job (if ppn = #cores as recommended)

PBS Job Management Commands

- qsub myjob.pbs
 - submits a job
 - PBS description: shell script with #PBS commands (ignored by shell, see next slide)
- qstat, qstat -f jobnumber
- qdel jobnumber
- pbsnodes -a
 - list all nodes with status and properties

PBS Job Description

```

#! /bin/sh

# Common PBS variables
#PBS -l nodes=5:ppn=8:mem16GB
#PBS -N test-mpi
#PBS -M jwagner@computing.dcu.ie
#PBS -m eba
#PBS -l walltime=00:05:00

source ${HOME}/.bashrc

exp_dir=/home/jwagner/
cd ${exp_dir}

# run multiple copies
mpiexec -n 40 /home/jwagner/intro/demo.py
  
```

Annotations:

- Number of nodes #CPU cores/node (points to nodes=5:ppn=8)
- Notification: end, begin and abort (points to -M and -m)
- Maximum runtime (points to walltime=00:05:00)
- Number of processes to start (points to -n 40)

Example: Memory-Intensive Job

```

#! /bin/sh

# Common PBS variables
#PBS -l nodes=1:ppn=8:min32GB
#PBS -N my-test-job
#PBS -M jwagner@computing.dcu.ie
#PBS -m eba
#PBS -l walltime=01:30:00

source ${HOME}/.bashrc

exp_dir=/home/jwagner/
cd ${exp_dir}

# run single process
/home/jwagner/intro/hello.py
  
```

Taskfarming Executable

- If Instance ID == 0
 - Run master code loop:
 - Read .tfm file (arg 1)
 - Send lines to worker
 - Exit if no more task and all worker finished
- Else
 - Run worker loop:
 - Ask master for a task
 - Execute task
 - Exit if master has no more tasks

Example: Taskfarming PBS File

```

#!/bin/bash

# Common PBS variables
#PBS -l nodes=2:ppn=4:min4GB
#PBS -N testxle5
#PBS -M jwagner@computing.dcu.ie
#PBS -m bea
#PBS -l walltime=01:30:00
#PBS -v

source ${HOME}/.bashrc

exp_dir=/home/jwagner/xle_parsing/
cd ${exp_dir}

mpiexec -n 8 /home/jwagner/taskfarm.py \
/home/jwagner/xle_parsing/xle.tfm

```

Example: Taskfarming TFM File

```

/home/jwagner/xle_parsing/run-package.sh 000
/home/jwagner/xle_parsing/run-package.sh 001
/home/jwagner/xle_parsing/run-package.sh 002
/home/jwagner/xle_parsing/run-package.sh 003
/home/jwagner/xle_parsing/run-package.sh 004
/home/jwagner/xle_parsing/run-package.sh 005
/home/jwagner/xle_parsing/run-package.sh 006
/home/jwagner/xle_parsing/run-package.sh 007
/home/jwagner/xle_parsing/run-package.sh 008
/home/jwagner/xle_parsing/run-package.sh 009
/home/jwagner/xle_parsing/run-package.sh 010
/home/jwagner/xle_parsing/run-package.sh 011
/home/jwagner/xle_parsing/run-package.sh 012

```

Example: Taskfarming Helper Script

```

#!/bin/bash

PACKAGE=$1

# never use more than 995 MB
ulimit -v 995000

# keep time
touch /home/jwagner/xle_parsing/Files/$PACKAGE.start

# prepare data
bunzip2 /home/jwagner/xle_parsing/Files/$PACKAGE.bz2

# run parser
/usr/local/xle/bin/xle -noTk -e "create-parser /usr/local/share/xle-grammars/english/english.lfg; parse-testfile /home/jwagner/xle_parsing/Files/$PACKAGE; exit" 2>/dev/null >/dev/null

# delete files not needed anymore
rm -f /home/jwagner/xle_parsing/Files/$PACKAGE
rm -f /home/jwagner/xle_parsing/Files/$PACKAGE.new

# compress results
bzip2 /home/jwagner/xle_parsing/Files/$PACKAGE.state

```

run-package.sh

Example: Non-Terminating Task

CPU 1: 000, 003, 006 (does not terminate)

CPU 2: 001, 005, 008, 009, 011 → idle

CPU 3: 002, 004, 007, 010, 012 → idle

CPU 4: Master: reads .tfm and distributes tasks

Killed at Walltime Limit

Effect of Task Size

- Job will wait for last task to finish (or be killed when walltime limit is reached)
- What if a task crashes?
 - Results are incomplete
 - Next tasks is executed
- What if a task does not terminate?
 - Results are incomplete
 - Fewer CPUs available for remaining tasks
- Overhead of starting tasks