

Consortium des Equipements de Calcul Intensif en Fédération Wallonie-Bruxelles

Introduction to Parallel Computing

damien.francois@uclouvain.be
October 2018





Agenda



- 1. General concepts, definitions, blockers
- 2. Hardware for parallel computing
- 3. Programming models
- 4. User tools

1.



General concepts

Why parallel?



Speed up – Solve a problem faster

→ more processing power

(a.k.a. strong scaling)

Scale up – Solve a larger problem

→ more memory and network capacity
(a.k.a. weak scaling)

Scale out – Solve many problems

→ more storage capacity

Parallelization involves:

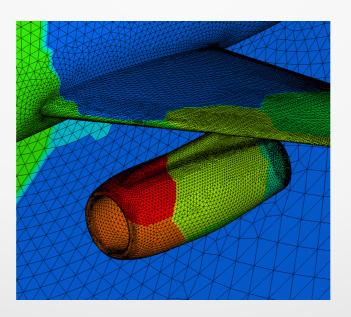


- decomposition of the work
 - distributing instructions to processors
 - distributing data to memories
- collaboration of the workers
 - synchronization of the distributed work
 - communication of data

Decomposition



- Work decomposition: task-level parallelism
- Data decomposition : data-level parallelism
- Domain decomposition: decomposition of work and data is done in a higher model, e.g. in the reality



Collaboration



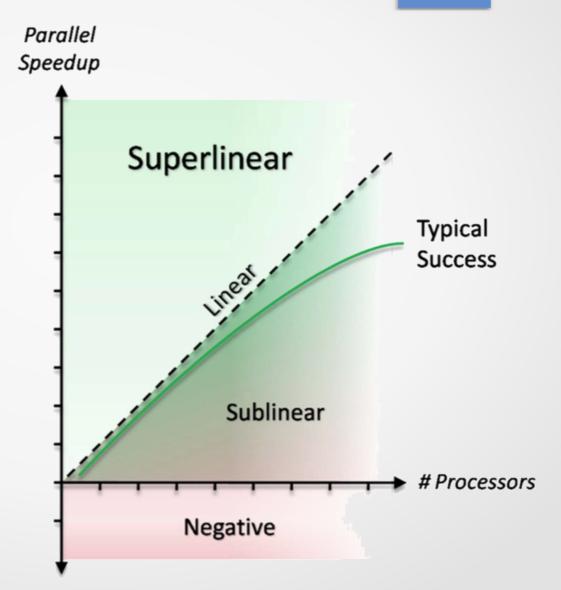
- Synchronous (SIMD) at the processor level
- Fine-grained parallelism if subtasks must communicate many times per second (instruction level); loosely synchronous
- Coarse-grained parallelism if they do not communicate many times per second (function-call level)
- Embarrassingly parallel if they rarely or never have to communicate (asynchronous)

Speedup, Efficiency, Scalability



$$S = \frac{T_S}{T_P}$$

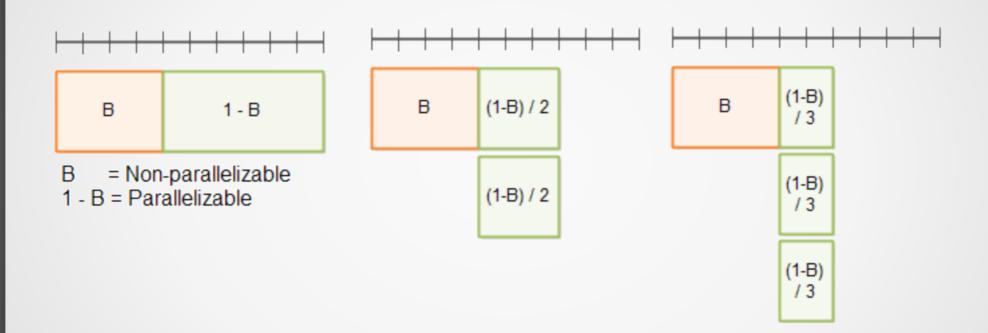
$$E = \frac{S}{p} = \frac{T_S}{pT_p}$$



Blocker 1: Amdahl's Law



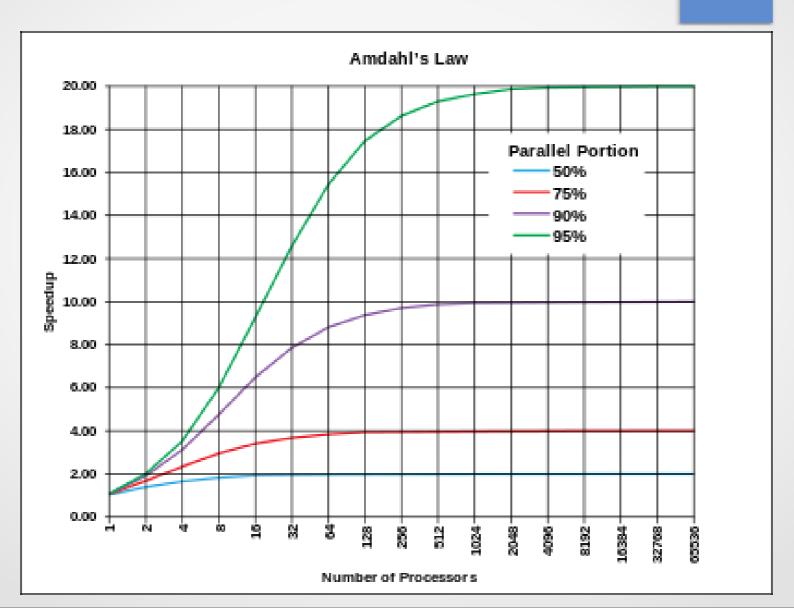
Often, not all the work can be decomposed



In parallel computing, Amdahl's law is mainly used to predict the theoretical maximum speedup for programs using multiple processors.

Blocker 1: Amdahl's Law

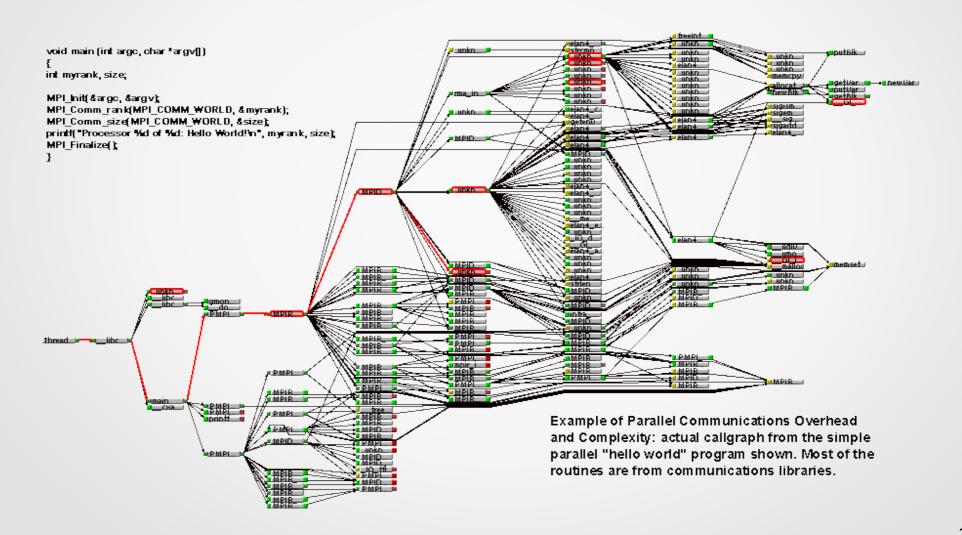




Blocker 2: Parallel overhead



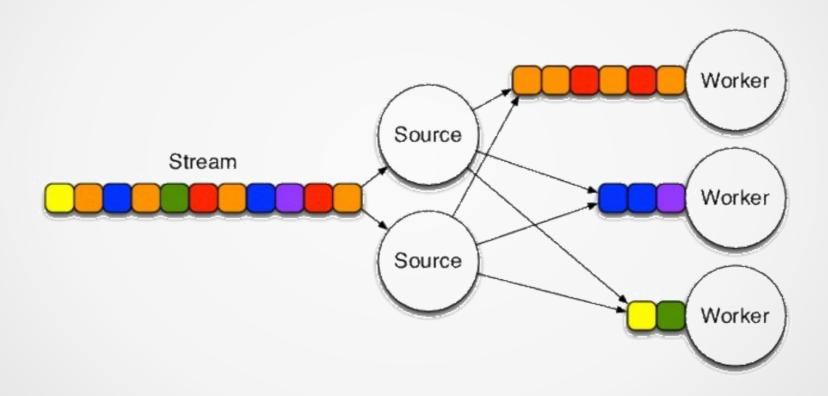
Collaboration means communication and a lot of extra work



Blocker 2: Parallel overhead



Load imbalance



2.



Hardware for parallel computing

At the core level



- Instruction-level parallelism (ILP)
 - Instruction pipelining
 - Superscalar execution
 - Out-of-order execution
 - Speculative execution
- Single Instruction Multiple Data (SIMD)

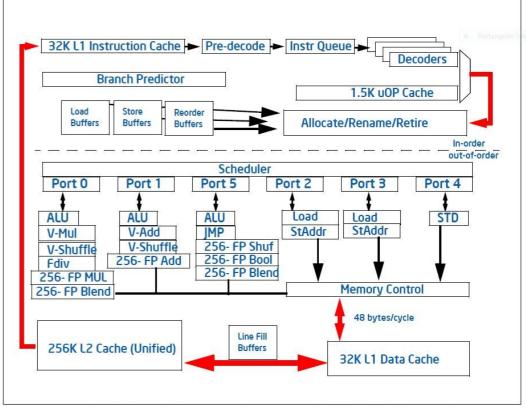
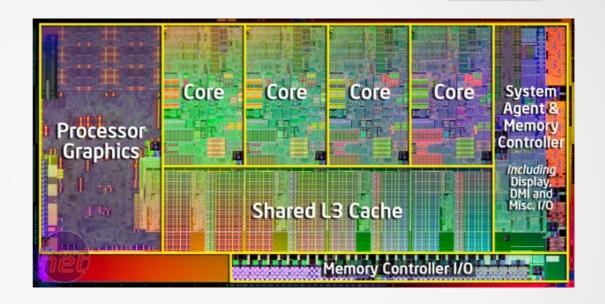


Figure 2-1. Intel microarchitecture code name Sandy Bridge Pipeline Functionality

At the CPU (socket) level



 Multicore parallelism

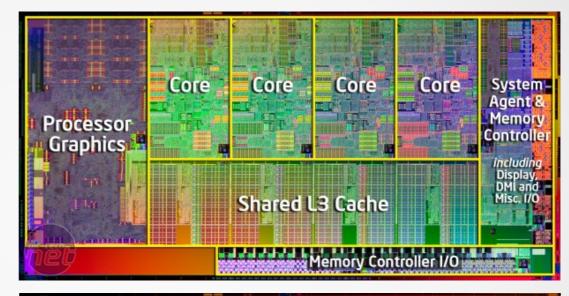


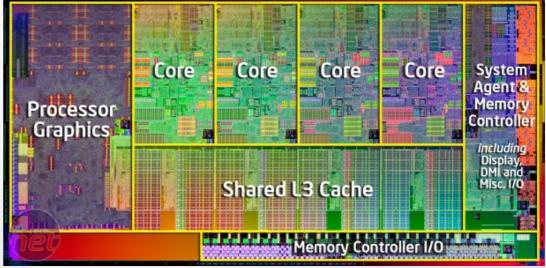
At the computer level



- Multi-socket parallelism
 - SMP
 - NUMA
- Accelerators







At the data center level



Cluster computing

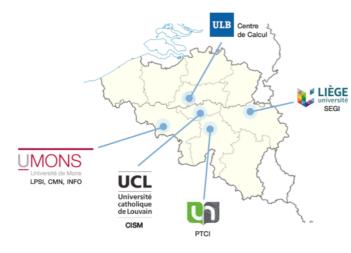


Consortium des Équipements de Calcul Intensif

6 clusters, 10k cores, 1 login, 1 home directory

About

CÉCI is the 'Consortium des Équipements de Calcul Intensif'; a consortium of high-performance computing centers of UCL, ULB, ULg, UMons, and UNamur. Read more.





The common storage is functional!

Have you tried it yet? More info...

Latest News

SATURDAY, 23 SEPTEMBER 2017

A CECI user pictured in the ULiège news!

The ULiège website published a story (in French) about the work of Denis Baurain and his collaborators on the Tier-1 cluster Zenobe that lead to a publication in Nature Ecology & Evolution.

TUESDAY, 01 AUGUST 2017

Ariel Lozano is the new CÉCI logisticien

We are happy to announce the hire of a new CECI logisticien: Ariel Lozano. Welcome Ariel!

At the data center level



Cloud computing











At the world level



Distributed computing



SETI@home is a scientific experiment, based at UC Berkeley, that uses Internet-connected computers in the Search for Extraterrestrial Intelligence (SETI). You can participate by running a free program that downloads and analyzes radio telescope data.

Join SETI@home

User of the Day



The_PC_God

Hello community. My name is Daniel. I am 28 years old and i live in a small village called Kuhardt (Rhineland-Palatinate, Germany)

which is located...

BSRC Student Travel Fundraiser

Berkeley SETI Research Center is holding a fundraiser to raise \$7000 to send our student interns to conferences to present their work.

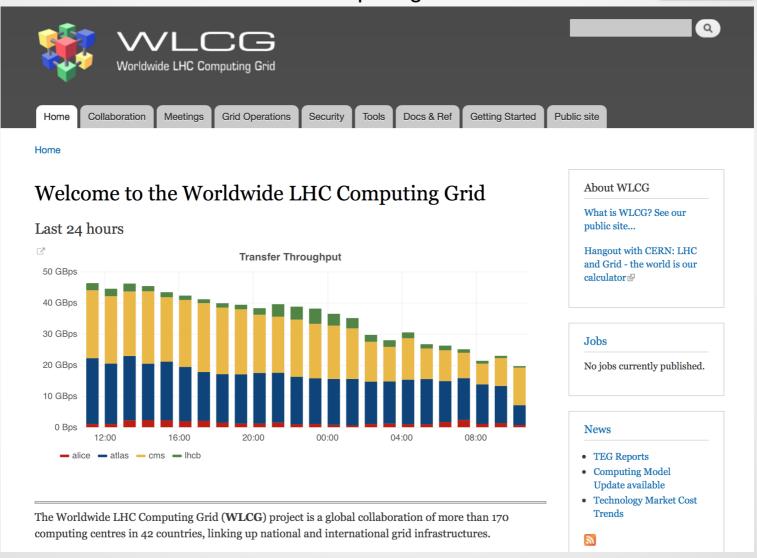
We've been working with some great students at Berkeley SETI, and we're optimistic that some of them will become the scientists and engineers who lead the field in future and maybe even find the signal we're searching for. In the meantime they have been doing amazing work and we'd like to send them to academic conferences to present their results, and for their own professional development. If you would like to help with this effort, we are running a crowdfunding campaign at https://crowdfund.berkeley.edu/SETItravel - every donation counts! We also have some fun perks including the chance to ask questions to members of the Berkeley SETI team, and to attend a party in our lab.

Although this does not directly benefit SETI@home (our annual fundraiser will start in a couple weeks), it's a worthy cause. I'll be contributing! 2 Oct 2017, 18:10:26 UTC • Discuss

At the world level



Grid computing





Programming models

Parallel programming paradigms (C.E.C.I)



How is work organized?

- Task-farming: no communication among workers
 - Master distribute work to workers (master/slave); or
 - Workers pick up tasks from pool (work stealing).
- **SPMD** (Single program multiple data)
 - A single program that contains both the logic for distributing work (master) and the computing part (workers) of which many instances are started and linked together at the same time
- MPMD (Multiple programs multiple data)

Parallel programming paradigms (C.E.C.I)



How is work organized?

- Pipelining (A->B->C, one process per task concurrently)
- Divide and Conquer (processes spawned at need and report their result to the parent)
- Speculative parallelism (processes spawned and result possibly discarded)

Programming models



What programming libraries/syntax constructs, etc. exist?

- Single computer:
 - CPUs: PThreads, OpenMP, TBB, OpenCL
 - Accelerators: CUDA, OpenCL, OpenAcc
- Multi-computer:
 - Distributed memory:
 - Shared storage: MPI (clusters)
 - Distributed storage: MapReduce (clouds)
 - No storage: BOINC (distributed computing)
 - Shared memory: CoArray, UPC

4.



User tools

that GNU/Linux offers

Parallel processes in Bash



```
● ● ●
                                               dfr@hmem00 - bash
dfr@hmem00:~/parcomp $ cat lower.sh
#!/bin/bash
# Usage:
     ./lower.sh [input_file [output_file]]
# Make ACTG chars lower case with extra processing.
# If output_file is not defined, stdout is used
# If input_file and output_file are not defined, stdin and stdout are used.
while read line; do
sleep 1
echo $line | tr ACTG actg >> ${2-/dev/stdout}
done < ${1-/dev/stdin}
dfr@hmem00:~/parcomp $ cat d.txt
dfr@hmem00:~/parcomp $ ./lower.sh d.txt
dfr@hmem00:~/parcomp $
```

Parallel processes in Bash



```
● ● ●
                                               dfr@hmem00 - bash
dfr@hmem00:~/parcomp $ # Foreground: commands end with ';'
dfr@hmem00:~/parcomp $ time { ./lower.sh d1.txt r1.txt ; ./lower.sh d1.txt r2.txt ; };
        0m8.033s
real
        0m0.004s
user
        0m0.019s
SVS
dfr@hmem00:~/parcomp $ # Background, in parallel: commands end with '&' and 'wait' necessary
dfr@hmem00:~/parcomp $ time { ./lower.sh d2.txt r1.txt & ./lower.sh d2.txt r2.txt & wait ; };
[1] 49722
[2] 49723
[1] - Done
                              ./lower.sh d2.txt r1.txt
[2]+ Done
                              ./lower.sh d2.txt r2.txt
real
        0m4.011s
        0m0.004s
user
        0m0.005s
dfr@hmem00:~/parcomp $
```

One program and many files



```
Equivalent to
                                                                                  ./lower.sh d1.txt;
● ● ●
                                           dfr@hmem00 - bash
                                                                                  ./lower.sh d2.txt;
dfr@hmem00:~/parcomp $ 1s d?.txt
d1.txt d2.txt d3.txt d4.txt
                                                                                  ./lower.sh d3.txt:
dfr@hmem00:~/parcomp $ ls d?.txt | xargs -n 1 echo "File: "
                                                                                  ./lower.sh d3.txt;
File: d2.txt
File: d3.txt
File: d4.txt
dfr@hmem00:~/parcomp $ time { ls d?.txt | xargs -n 1 ./lower.sh > /dev/null ; }
       0m16.041s
real
user
       0m0.010s
       0m0.006s
dfr@hmem00:~/parcomp $ time { ls d?.txt | xargs -n 1 -P 4 ./lower.sh > /dev/null ; }
       0m4.014s
real
       0m0.008s
user
SVS
       0m0.016s
dfr@hmem00:~/parcomp $
                                                                                Equivalent to
                                                                              ./lower.sh d1.txt &
                                                                              ./lower.sh d2.txt &
                                                                              ./lower.sh d3.txt &
                                                                              ./lower.sh d3.txt &
                                                                                      wait
```

Several programs and one file



./upper.sh waits for ./lower.sh to finish

```
● ● ●
                                             dfr@hmem00 - bash
dfr@hmem00:~/parcomp $ # Using an intermediay file
dfr@hmem00:~/parcomp $ time { ./lower.sh d.txt tmp.txt ; ./upper.sh tmp.txt res.txt ; }
real
        0m8.033s
        0m0.005s
user
        0m0.017s
SVS
dfr@hmem00:~/parcomp $ # Using pipes (as our programs can handle stdin and stdout)
dfr@hmem00:~/parcomp $ time { ./lower.sh d.txt | ./upper.sh > res.txt ; }
real
        0m5.014s
                                                           ./upper.sh starts as soon as ./lower.sh
       0m0.006s
user
        0m0.009s
                                                                       writes the first output
dfr@hmem00:~/parcomp $ mkfifo tmpfifo
dfr@hmem00:~/parcomp $ ls -l tmpfifo
prw-rw-r-- 1 dfr dfr 0 Oct 7 10:27 tmpfifo
dfr@hmem00:~/parcomp $ time { ./lower.sh d.txt tmpfifo & ./upper.sh tmpfifo res.txt ; }
[1] 65343
[1]+ Done
                             ./lower.sh d.txt tmpfifo
       0m5.013s
real
        0m0.002s
user
        0m0.007s
SVS
dfr@hmem00:~/parcomp $
```

Several programs and one file



```
● ● ●
                                             dfr@hmem00 - bash
dfr@hmem00:~/parcomp $ # Using an intermediay file
dfr@hmem00:~/parcomp $ time { ./lower.sh d.txt tmp.txt ; ./upper.sh tmp.txt res.txt ; }
real
        0m8.033s
        0m0.005s
user
        0m0.017s
SVS
dfr@hmem00:~/parcomp $ # Using pipes (as our programs ca
                                                         If ./upper.sh was not designed to read
dfr@hmem00:~/parcomp $ time { ./lower.sh d.txt | ./upper
                                                        from STDIN, we could use a FIFO file
real
        0m5.014s
user
       0m0.006s
        0m0.009s
dfr@hmem00:~/parcomp $ mkfifo tmpfifo
dfr@hmem00:~/parcomp $ ls -l tmpfifo
prw-rw-r-- 1 dfr dfr 0 Oct 7 10:27 tmpfifo
dfr@hmem00:~/parcomp $ time { ./lower.sh d.txt tmpfifo & ./upper.sh tmpfifo res.txt ; }
[1] 65343
[1]+ Done
                             ./lower.sh d.txt tmpfifo
       0m5.013s
real
        0m0.002s
user
        0m0.007s
SVS
dfr@hmemθθ:~/parcomp $
```

One program and one large file



```
● ● ●
                                             dfr@hmem00 - bash
dfr@hmem00:~/parcomp $ # One process to process the whole file
                                                                  Split the file and start 4 processes
dfr@hmem00:~/parcomp $ time { cat d.txt | ./lower.sh > res.txt
        0m4.014s
real
        0m0.003s
user
        0m0.009s
SVS
dfr@hmem00:~/parcomp $ # Four processes handling one line in round robin fashion
dfr@hmem00:~/parcomp $ time { cat d.txt | split --unbuffered --number r/4 --filter ./lower.sh >res.txt ; }
real
        0m1.011s
       0m0.009s
user
       0m0.021s
dfr@hmem00:~/parcomp $ !! & top -u dfr -bn1 | grep lower
time { cat d.txt | split --unbuffered --number r/4 --filter ./lower.sh >res.txt ; } & top -u dfr -bn1 | gr
ep lower
[1] 12817
12822 dfr
               20
                    0 103m 1252 1052 S 0.0 0.0
                                                    0:00.00 lower.sh
12823 dfr
               20 0 103m 1252 1052 S 0.0 0.0
                                                  0:00.00 lower.sh
12824 dfr
               20
                   0 103m 1252 1052 S
                                                  0:00.00 lower.sh
12825 dfr
               20
                  0 103m 1252 1052 S 0.0 0.0 0:00.00 lower.sh
dfr@hmem00:~/parcomp $
real
       0m1.011s
       0m0.011s
user
SVS
       0m0.019s
                             time { cat d.txt | split --unbuffered --number r/4 --filter ./lower.sh > res
[1]+ Done
.txt; }
dfr@hmem00:~/parcomp $
```

Need recent version of Coreutils/8.22-goolf-1.4.10

Several programs and many files C.E.C.I



```
dfr@hmem00 - bash
# Sample Makefile to process each file with
# lower.sh then upper.sh
all: d1.res d2.res d3.res d4.res
# Build intermediary files
%.tmp: %.txt
    ./lower.sh $< $@
# Build final result
%.res: %.tmp
    ./upper.sh $< $@
"Makefile" 14L, 219C written
                                                                                       14,0-1
```

Several programs and many files C.E.C.D



```
dfr@hmem00 - bash
dfr@hmem00:~/parcomp $ time make
./lower.sh d1.txt d1.tmp
./upper.sh d1.tmp d1.res
./lower.sh d2.txt d2.tmp
./upper.sh d2.tmp d2.res
./lower.sh d3.txt d3.tmp
./upper.sh d3.tmp d3.res
./lower.sh d4.txt d4.tmp
./upper.sh d4.tmp d4.res
rm d1.tmp d2.tmp d4.tmp d3.tmp
        0m32.260s
real
user
        0m0.028s
        0m0.099s
SVS
dfr@hmem00:~/parcomp $ rm *res
dfr@hmem00:~/parcomp $ time make -j 4
./lower.sh d1.txt d1.tmp
./lower.sh d2.txt d2.tmp
./lower.sh d3.txt d3.tmp
./lower.sh d4.txt d4.tmp
./upper.sh d1.tmp d1.res
./upper.sh d2.tmp d2.res
./upper.sh d4.tmp d4.res
./upper.sh d3.tmp d3.res
rm d1.tmp d2.tmp d4.tmp d3.tmp
        0m8.163s
real
user
        0m0.025s
```

Summary



- You have either
 - one very large file to process
 - with one program: split
 - with several programs: pipes, fifo
 - many files to process
 - with one program xargs
 - with many programs make

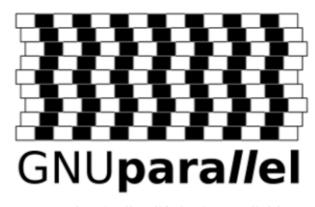


GNU Parallel

GNU **parallel** is a shell tool for executing jobs in parallel using one or more computers. A job can be a single command or a small script that has to be run for each of the lines in the input. The typical input is a list of files, a list of hosts, a list of users, a list of URLs, or a list of tables. A job can also be a command that reads from a pipe. GNU **parallel** can then split the input and pipe it into commands in parallel.

If you use xargs and tee today you will find GNU **parallel** very easy to use as GNU **parallel** is written to have the same options as xargs. If you write loops in shell, you will find GNU **parallel** may be able to replace most of the loops and make them run faster by running several jobs in parallel.

GNU **parallel** makes sure output from the commands is the same output as you would get had you run the commands sequentially. This makes it possible to use output from GNU **parallel** as input for other programs.



For people who live life in the parallel lane.

For each line of input GNU parallel will execute *command* with the line as arguments. If no *command* is given, the line of input is executed. Several lines will be run in parallel. GNU parallel can often be used as a substitute for xargs or cat | bash.

More complicated to use but very powerful



Syntax: parallel command ::: argument list

```
dfr@hmem00 - bash
dfr@hmem00:~/parcomp $ parallel echo ::: 1 2 3 4
dfr@hmem00:~/parcomp $ parallel echo ::: {1..10}
dfr@hmem00:~/parcomp $ time parallel sleep ::: {1..10}
real
        0m11.200s
        0m0.206s
user
        0m0.129s
dfr@hmem00:~/parcomp $ parallel echo ::: d?.txt
d1.txt
d2.txt
d3.txt
dfr@hmem00:~/parcomp $
```



Syntax: {} as argument placeholder. Can be modified

```
dfr@hmem00:~/parcomp $ parallel echo {} ::: d?.txt
d1.txt
d2.txt
d3.txt
d4.txt
dfr@hmem00:~/parcomp $ parallel echo {} {.}.res ::: d?.txt
d1.txt d1.res
d2.txt d2.res
d3.txt d3.res
d4.txt d4.res
dfr@hmem00:~/parcomp $ parallel echo {} ::: ../parcomp/d?.txt
../parcomp/d1.txt
../parcomp/d2.txt
../parcomp/d3.txt
../parcomp/d4.txt
dfr@hmem00:~/parcomp $ parallel echo {/} ::: ../parcomp/d?.txt
d1.txt
d2.txt
d3.txt
d4.txt
dfr@hmem00:~/parcomp $
dfr@hmemθθ:~/parcomp $
dfr@hmem00:~/parcomp $
```



Multiple parameters and --xapply

```
● ● ●
                                               dfr@hmem00 - bash
dfr@hmem00:~/parcomp $ parallel echo ::: 1 2 3 4 ::: A B
1 B
dfr@hmem00:~/parcomp $ parallel --xapply echo ::: 1 2 3 4 ::: A B C D
2 B
dfr@hmem00:~/parcomp $ parallel echo {1} and {2} ::: 1 2 3 4 ::: A B C D
1 and A
1 and B
 and C
  and D
  and A
  and B
  and C
  and D
  and A
 and B
  and C
3 and D
  and A
```



When arguments are in a file: use:::: (4x ':')

```
● ● ●
                                              dfr@hmem00 - bash
dfr@hmem00:~/parcomp $ cat experiments.csv
Number, Letter
2,B
dfr@hmem00:~/parcomp $ parallel --colsep ',' --header '\n' echo {Number} {Letter} ::: experiments.csv
2 B
dfr@hmem00:~/parcomp $
```



Split a file with --pipe

```
dfr@hmem00 - bash
dfr@hmem00:~/parcomp $ cat *.csv
Number, Letter
1.A
2.B
3.B
3.A
5.C
dfr@hmem00:~/parcomp $ tail -n +2 *.csv | parallel -kN1 --recend '\n' --pipe echo -n "JOB{#} : ;cat;"
JOB1 :1.A
JOB2 : 2.B
JOB3 :3,B
JOB4 : 3.A
JOB5 :4,C
JOB6 :5,C
JOB7:5,A
dfr@hmem00:~/parcomp $
```

Other interesting options



- Use remote servers through SSH
- -j n Run n jobs in parallel
- -kKeep same order
- --delay n Ensure there are n seconds between each start
- --timeout n Kill task after n seconds if still running

Author asks to be cited: O. Tange (2011): *GNU Parallel - The Command-Line Power Tool*, The USENIX Magazine, February 2011:42-47.

Exercise



Reproduce the examples with ./lower and ./upper.sh using GNU Parallel

Solutions



One program and many files

```
$ time parallel -k ./lower.sh {} > res.txt ::: d?.txt
```

One program and one large file

```
$ time cat d.txt | parallel -k -N1 --pipe ./lower.sh {} > res.txt
```

Several programs and several files

```
$ time { parallel ./lower.sh {} {.}.tmp ::: d?.txt ; \
> parallel ./upper.sh {} {.}.res ::: d?.tmp ; }
```