

## map/reduce (DIY)

#EmbarassinglyParallelPrograms #Make #SoftwareSelection #ProcessModeling #Redis #Windows

### Embarassingly parallel processing with no setup

We have 100 pairs of big files to be processed by the same program (e.g. a full comparison)

→ *the problem is 'embarassingly parallel'*

The processing time is 6 minutes each.

→ the total processing time seems to be 10 hours.

We want to explore options to make the total time shorter.

we have a Windows domain with 25 idle machines.

Installing distributed-workflow software and agents on our machines can be complicated by long approval chains or skill shortage.

A possible solution is based on a single Windows script that,

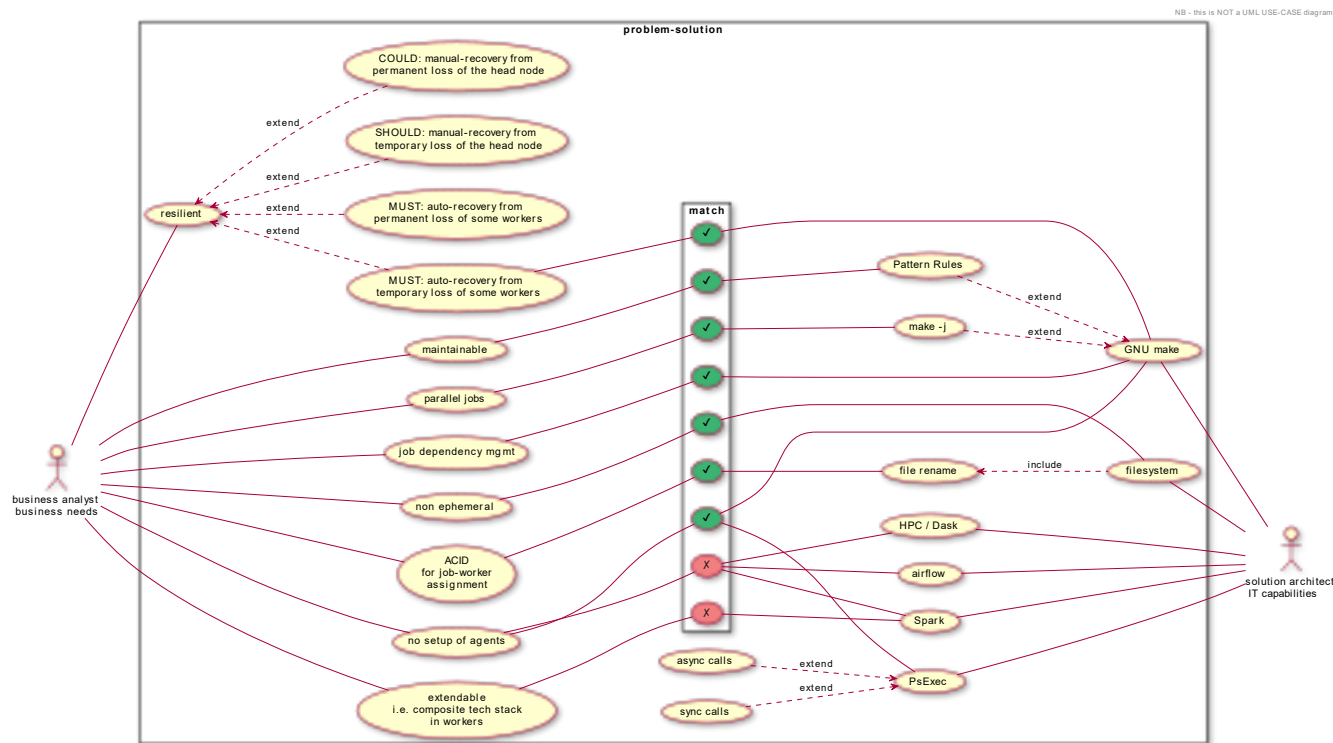
- ▶ with the help of PsExec.exe, submits remote executions on the 25 machines of the domain, and,
- ▶ with the help of the GNU make.exe, performs an automatic parallelization of the map tasks.

An alternative option may be based on Spark. A Spark-based project is detailed inside this website [here](#).

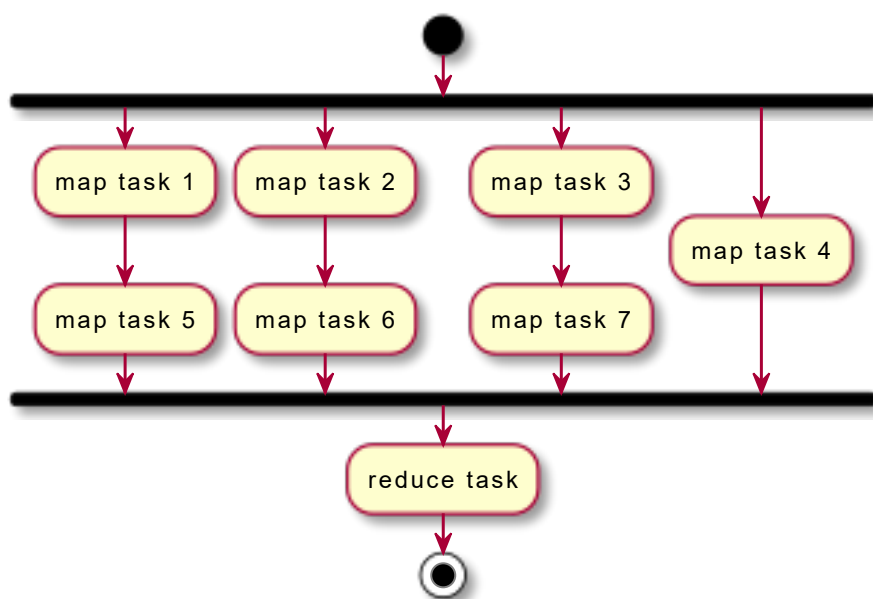
The total execution time is expected to be close to  $6 * 100/25$  minutes = less than half hour (assuming the latency of file-transfers is negligible).

### building the tech stack

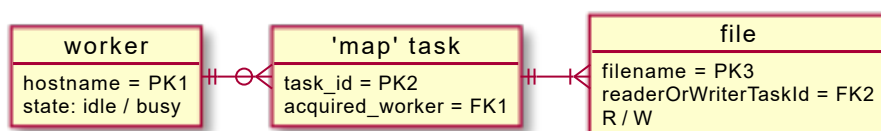
The following diagram recaps our needs on the left, the tech opportunities on the right, what is adopted and what is ruled-out:



a map-reduce work schedule, assuming 7 map-tasks and 4 nodes:



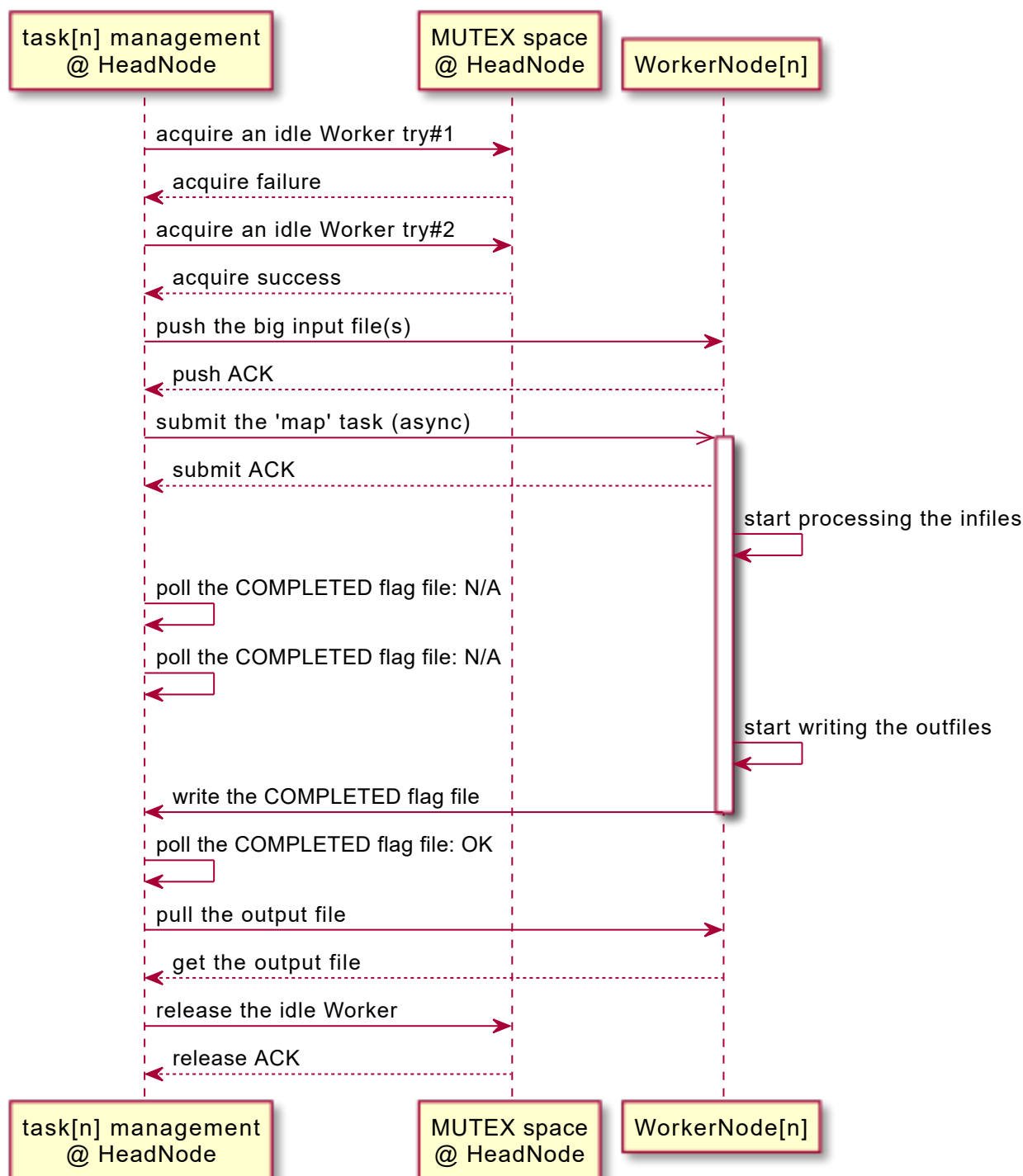
E-R diagram of the entities used by the orchestrator (GNU make):



NB

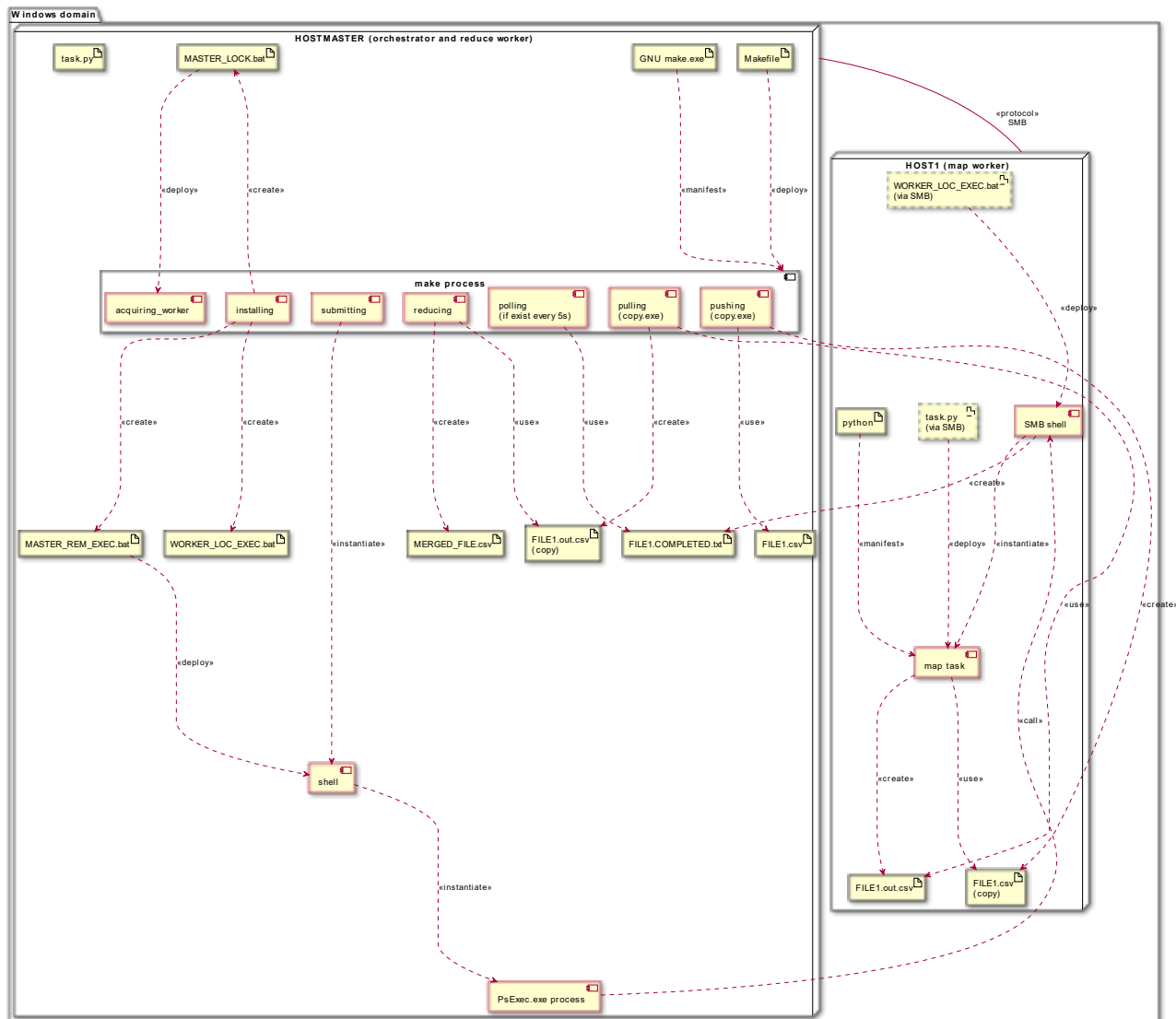
- in such map-reduce scenarios, each file can either be input or output of one and only one task.
- we stick to max one task per worker at any given time.

sequence diagram of the actions performed by each map task:



NB: again, nothing needs to be deployed on the workers machines, except for the input and output data files.

deployment diagram:



The Makefile that defines the steps and the dependencies of the whole process is available at this link: <https://github.com/a-moscattelli/home/blob/main/am-wiki-assets/mapreducewin/Makefile>

Based on the Makefile above, we just run:

```
1 | set PARALLEL_DEGREE=3
2 | make install
3 | make initialize
4 | make -j %PARALLEL_DEGREE% all
```

The activity log, with five 1-minute long map jobs and 3 nodes, is as below:



```
1 | 23:24:56.16 EXEC_ initialize
2 |
3 | 23:25:01.78 BEGIN TASK001.step1.workerAcquisitionLoopFinished.mkcontrol
4 | 23:25:01.84 BEGIN TASK002.step1.workerAcquisitionLoopFinished.mkcontrol
```

```

5 23:25:01.86 BEGIN TASK003.step1.workerAcquisitionLoopFinished.mkcontrol
6 23:25:05.97 END__ TASK001.step1.workerAcquisitionLoopFinished.mkcontrol
7 23:25:06.05 BEGIN TASK004.step1.workerAcquisitionLoopFinished.mkcontrol
8 23:25:06.11 END__ TASK002.step1.workerAcquisitionLoopFinished.mkcontrol
9 23:25:06.13 END__ TASK003.step1.workerAcquisitionLoopFinished.mkcontrol
10 23:25:06.19 BEGIN TASK005.step1.workerAcquisitionLoopFinished.mkcontrol
11 23:25:06.28 BEGIN TASK001.step2.inputFilePushedToWorker.mkcontrol HOST003
12 23:25:06.46 BEGIN TASK002.step2.inputFilePushedToWorker.mkcontrol HOST002
13 23:25:06.64 BEGIN TASK003.step2.inputFilePushedToWorker.mkcontrol HOST001
14 23:25:06.87 BEGIN TASK001.step3.submittedToWorker.mkcontrol HOST003
15 23:25:07.40 BEGIN TASK002.step3.submittedToWorker.mkcontrol HOST002
16 23:25:07.91 BEGIN TASK003.step3.submittedToWorker.mkcontrol HOST001
17
18 23:26:09.53 END__ TASK001.step4.workerCompletionCheckLoopFinished.mkcontrol
19 23:26:09.67 END__ TASK002.step4.workerCompletionCheckLoopFinished.mkcontrol
20 23:26:09.82 END__ TASK003.step4.workerCompletionCheckLoopFinished.mkcontrol
21 23:26:15.17 END__ TASK004.step1.workerAcquisitionLoopFinished.mkcontrol
22 23:26:15.28 BEGIN TASK004.step2.inputFilePushedToWorker.mkcontrol HOST003
23 23:26:15.30 END__ TASK005.step1.workerAcquisitionLoopFinished.mkcontrol
24 23:26:15.44 BEGIN TASK005.step2.inputFilePushedToWorker.mkcontrol HOST002
25 23:26:15.51 BEGIN TASK004.step3.submittedToWorker.mkcontrol HOST003
26 23:26:15.65 BEGIN TASK005.step3.submittedToWorker.mkcontrol HOST002
27
28 23:27:16.85 END__ TASK004.step4.workerCompletionCheckLoopFinished.mkcontrol
29 23:27:16.97 END__ TASK005.step4.workerCompletionCheckLoopFinished.mkcontrol
30 23:27:17.00 BEGIN reduce
31 23:27:17.08 END__ reduce

```

links:

- <https://learn.microsoft.com/en-us/sysinternals/downloads/psexec>   
the agent that makes remote calls
- <https://www.gnu.org/software/make/manual/>   
the agent that enables the parallel executions

the semaphore, required to ensure that only one process at a time can acquire a worker node, is based on the success/failure of a file rename.

Such option does not handle well the possibility that a lock holder will not survive the moment it is supposed to release the lock for other purposes.

a solution can be easily implemented using redis:

```

1 #docker
2 image: "redis:6.0.9"

```

Windows client:

<https://github.com/microsoftarchive/redis/releases> 

example of a working script:

```
1  set REDISCLI=.\Redis-x64-3.0.504\redis-cli.exe
2  set REDISSERVERHOST=DESKTOP-B12345T
3  set REDIS_CALL=%REDISCLI% -h %REDISSERVERHOST% -p 6379
4
5  %REDISCLI% -h %REDISSERVERHOST% -p 6379 PING
6  rem > PONG
7  echo %ERRORLEVEL%
8  rem > 0
9
10 set autoexpire_seconds=600
11
12 set key=HOST1
13 set val=TASK4
14
15 rem simulation of a lock acquisition"
16 %REDIS_CALL% SET %key% %val% EX %autoexpire_seconds% NX
17 rem > OK
18
19 rem simulation of a competing concurrent lock acquisition:
20 %REDIS_CALL% SET %key% %val% EX %autoexpire_seconds% NX
21 rem > (nil)
22
23 %REDIS_CALL% KEYS "*"
24 rem > 1) "HOST1"
25
26 rem simulation of a lock release:
27 %REDIS_CALL% DEL %key%
28 rem > (integer) 1
29
30 %REDIS_CALL% KEYS "*"
31 rem > (empty list or set)
```

*%ERRORLEVEL% after each CLI call above is always 0*

---

back to [Portfolio](#)

---

[backlog](#) 