

Project Design Document

(DPSort : Distributed, Parallel External MergeSort)

- Milestone 1 -

20180274 Jinho Ko

Progress Report Summary

- Finished deciding project directory structure.
- No progress in coding
 - Still need to study more on network programming
- Designed the architecture of system, especially one when parallelization becomes compulsory.

Project Directory Structure

CS444-project/

bin/ → master, worker shell files (entry point for execution)

conf/ → MASTER-CONF.XML, WORKER-CONF.XML, LOG4J-PROPERTIES

data/ → example data for tests and example execution.

dev/ → files for development

docs/ → documentation

production/ → aggregated files related w. production.

project/

core/

master/

worker/

build.sh

+ source/test files

Package Structure

Package dpsort

package dpsort.master \Rightarrow codes only used in master

↑
Dependency X!

package dpsort.worker \Rightarrow codes only used in worker

package dpsort.common. \Rightarrow commonly used codes.

Execution Behavior

@Master

(.sh)

bin/master [#numWorkersToJoin, args...]

- ↳ parse / validate arguments + attach conf files
(conf/...)
- ↳ find compiled source codes
- ↳ execute JVM.

@Worker

bin/worker (.sh) [MasterIP, -I ~, -O ~]
port

- ↳ //
- ↳ //

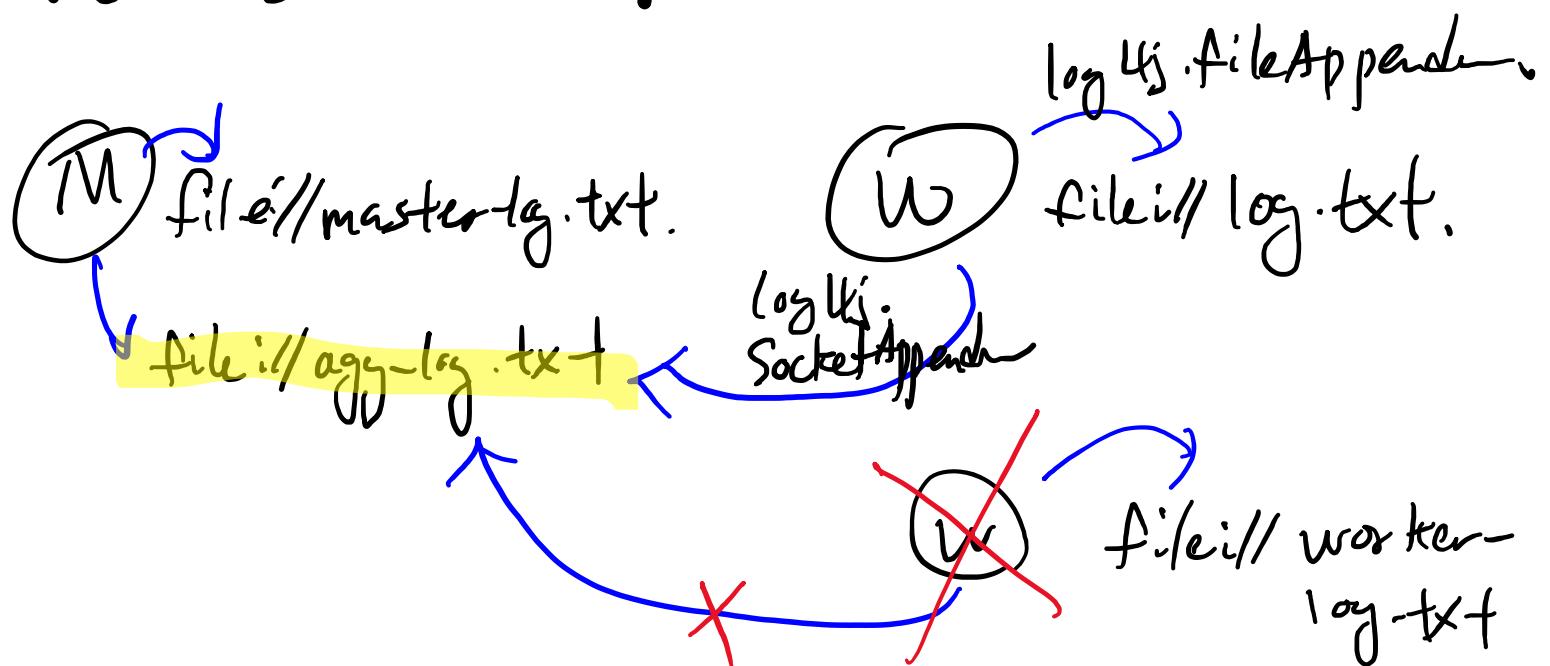
↓
if not exists → err

↓
if not exists →
mkdir.

- ↳ Execute JVM

Log Aggregation

Since logs are gathered from multiple machines, we aggregate them to the master machine.

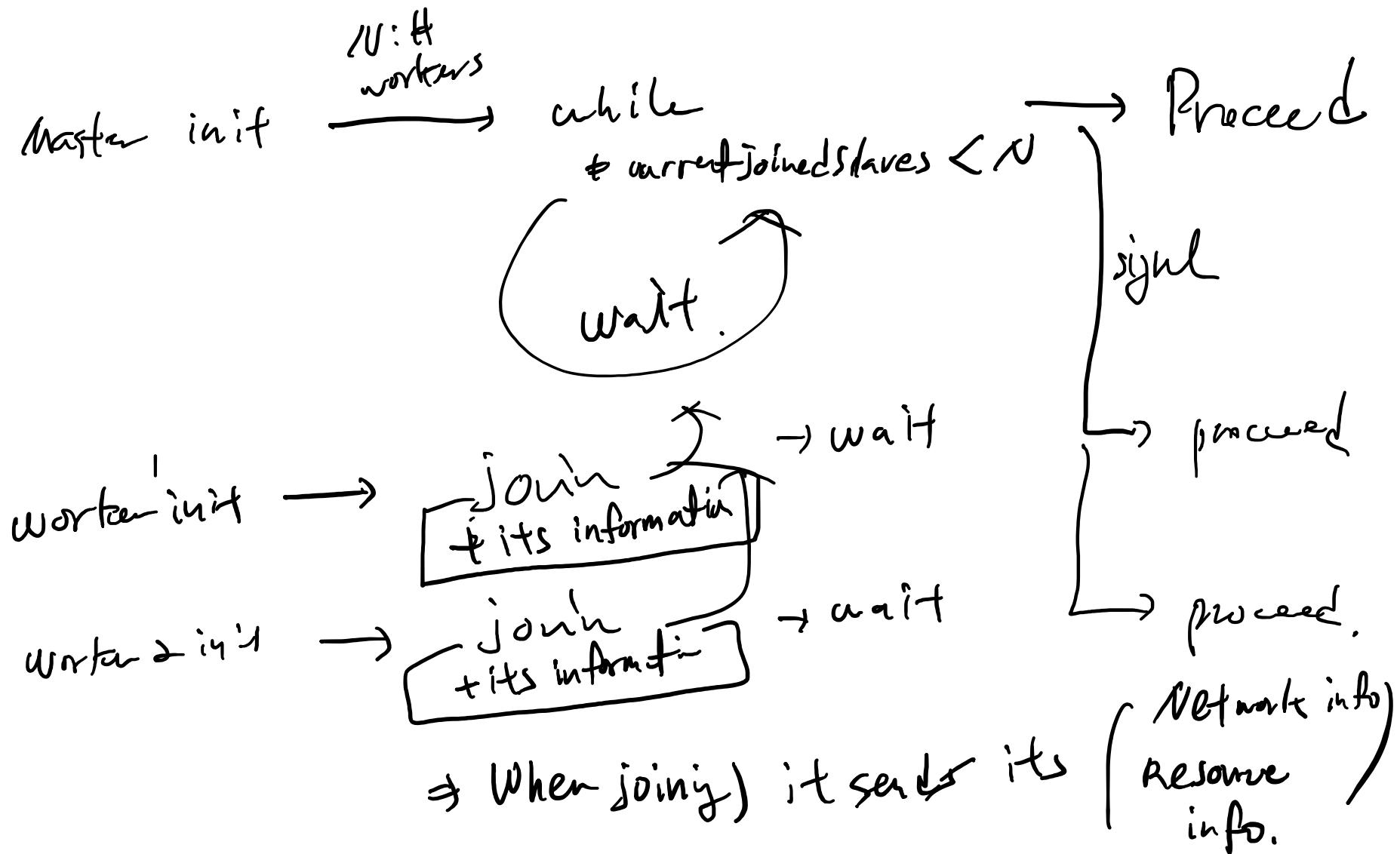


⇒ When network causes error, it can still save log to its local file.

Conf files

- MASTER-CONF.XML
 - common
 - master-only
- WORKER-CONF.XML.
 - common
 - worker-only
 - Parallism behavior, Port range, ...
 - # Partitions
 - Hardware specs (RAM, CORES, ...)
 - ...
- - - - -
 - file
directory / buffer
logging / Data for - mat
Timeout
- | common : includes.

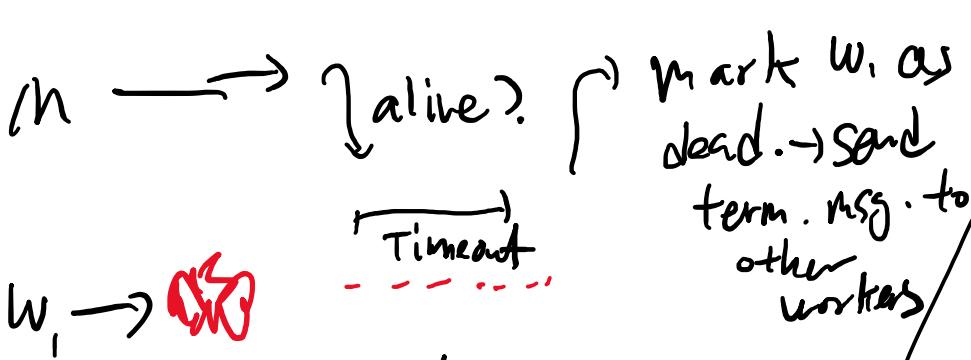
Worker Joining Behavior



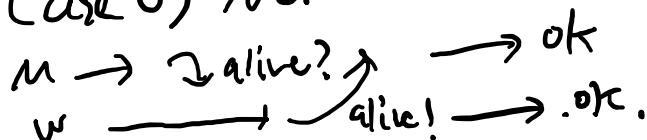
Heartbeat

- Idea: Master never gets to know if worker is running right unless worker sends its heartbeat.
- Why heartbeat? : ~~memory overflow / system lag~~ due to excessive operations.
↳ almost all cause of error.

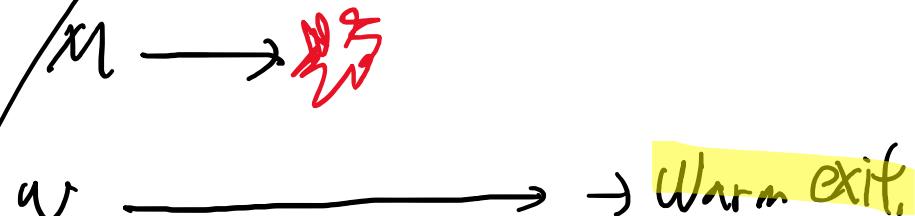
Case 1) worker is missing.



Case 0) Normal.



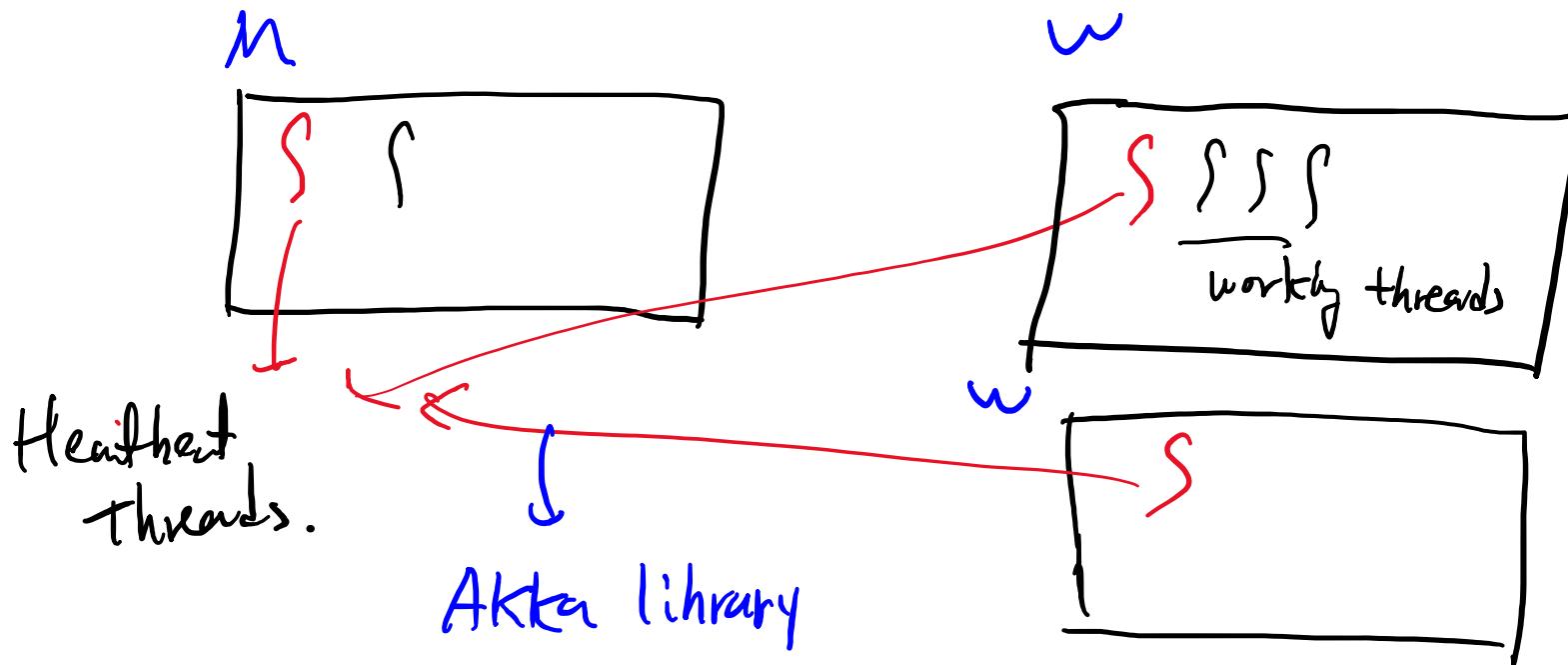
Case 2) master is missing.



No heartbeat check for long time / Network refused

Heartbeat (con't)

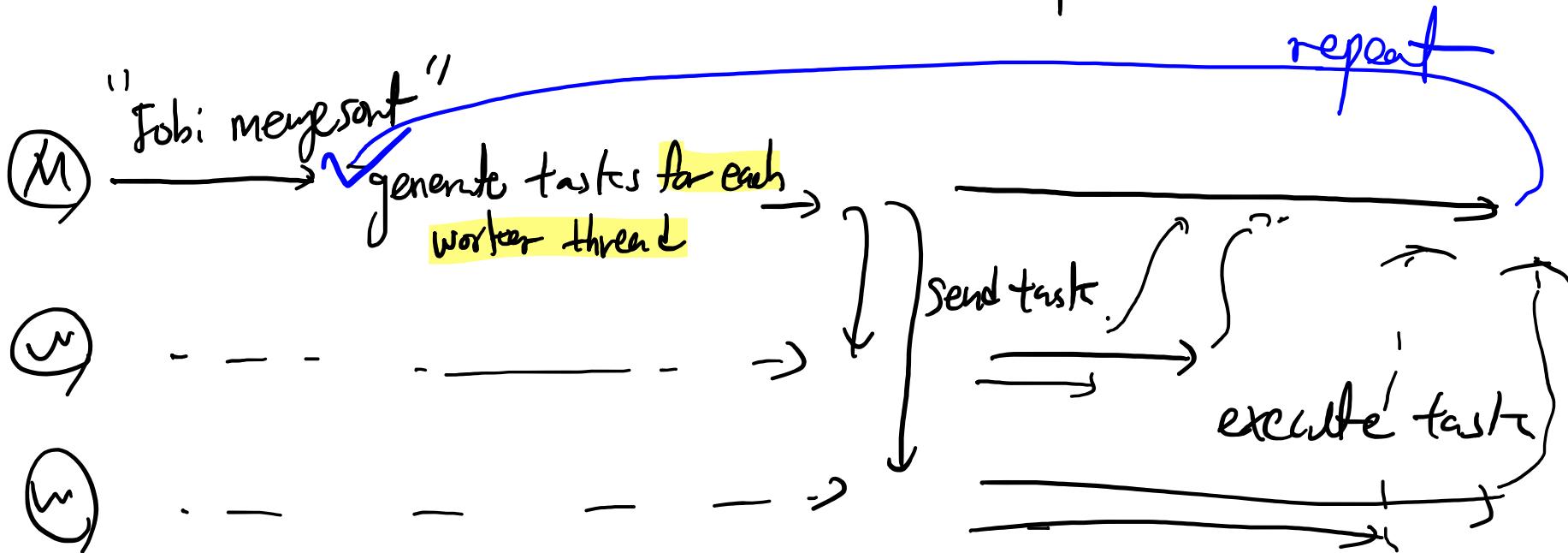
- ✓ Heartbeat "should" run on separate thread
- ✓ The thread also should have enough opportunity to be executed , otherwise, possibility for timeout.



Execution-related Terminologies

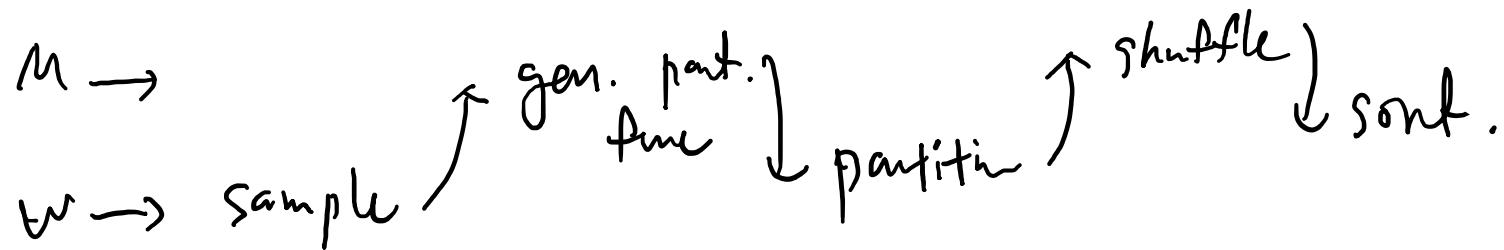
concept referenced
from Apache Spark

- Job: "external merge sort"
- Task: minimum unit of execution allocated to each worker thread (& partition)



Passing Tasks

Idea 1) Concurrency Control.



(+) simplest to implement

(-) hard to revise / upgrade code.
/ change in execution flow.

Passing Tasks (con't)

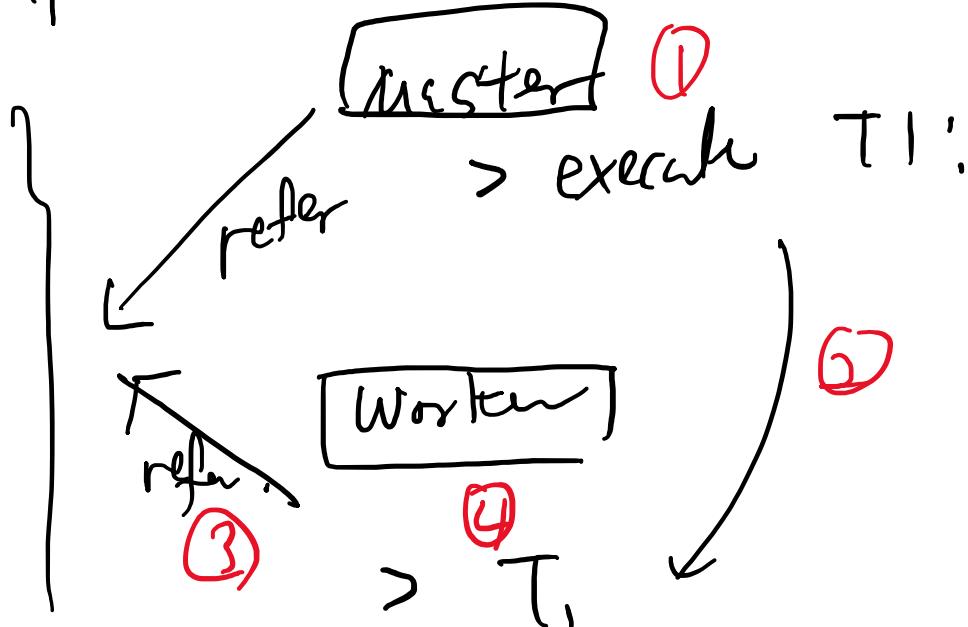
Idea 2) Use task type II.

@ defn. common.

```
g T1 : executable  
T2 : executable
```

```
;
```

```
}
```



(+) much easier to upgrade code →
by modifying

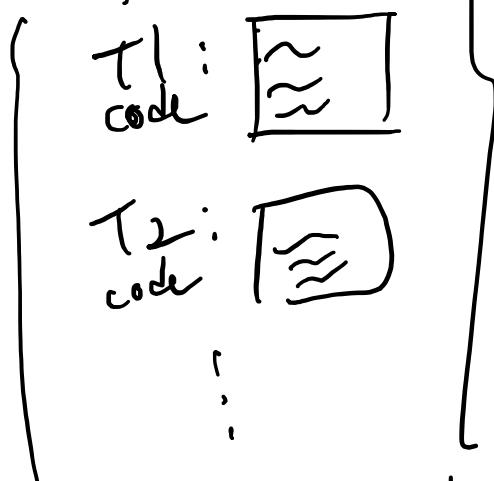
(-) code length ↑

Task types

Passing Tasks (con't)

Idea 3) Worker doesn't know about the

@cdpsnf.core.



v short worker-side code

Worker

> fetch code (+ data)

> execute code.

(+) robust.

v can be easily optimized /
upgraded.

(-) might be
tricky to implement.

Master

> val serCode

= serialize(code)

> execute serCode

Passing Tasks : Decision

- ✓ Considering between Idea 2 & Idea 3.
- ✓ Idea 3 is splitting into two parts.
 - code serialization) function.
 - partition function is broadcasted to all nodes, (or 2nd part)
serialization) function is broadcasted to all nodes.

Serialization

Objects should be serialized in order to be passed via networking.



Data : shuffle / sample data

, Code : generated partition function.

Resource Allocation

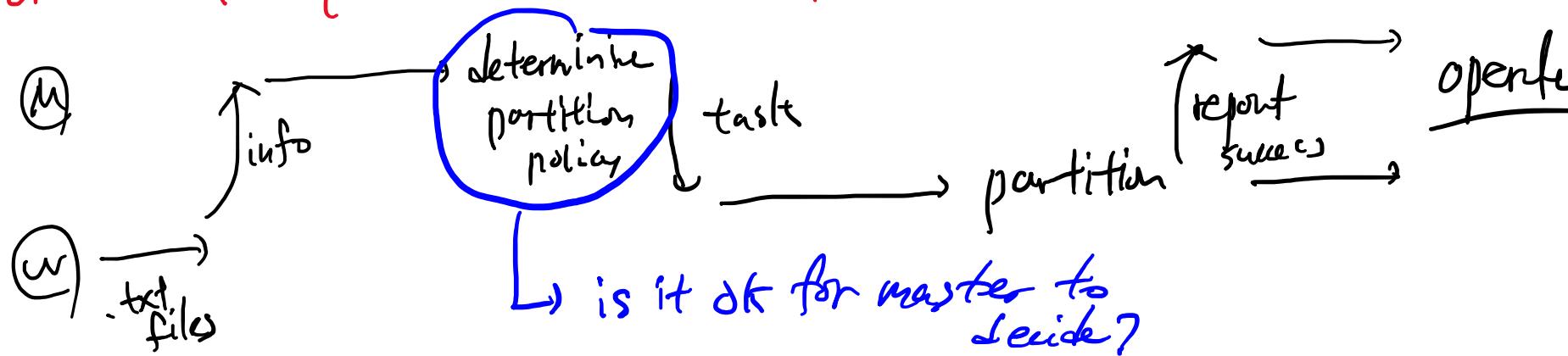
v Multithreading. \Rightarrow 1 partition / thread. golden rule

\Rightarrow Expect equi-memory distribution
to each threads.
(consort is simple!)

Threads \uparrow :
(+) can scale
(-) disk IO \uparrow

Threads \downarrow :
(+) disk IO \downarrow
(-) might not scale.

~~Partition policy~~ \rightarrow not decided yet. (consider file size,
#threads, memory).



Networking

- Use Protobuf (gRPC) for networking
- Need to study more about networking + async programming!

Shuffle

- Might be the most resource-consuming task in whole job
- Might need to consider assigning a thread only for shuffle behavior.
 - Consider case when Partition / shuffle / sort overlaps.

Merge Sort

- When executing merge sort, it is obvious that the whole data won't scale.
- Thus we need mechanism to spill data to disk. (It is differenting from simply writing data directly to file!)
 - A typical DBMS task.
 - Any simple library to use in Scala?