# PROGRESS PRESENTATION

Team Red

# Team Organization



junseong

dana

jaehwan

# Communication



Separate threads per topics

Threads also act as a part of documentation

# Communication



Weekly meetings



> 📦 You've used all of this workspace's free blocks    Upgrade plan

~~Notion~~ → git for documentation

# Documentation

## High Level

- Draft 1 (10/28)
- Draft 2 (10/31)
- Draft 3 (11/9)
- Draft 4 (11/13)

## Code Level

### Libraries

- Runtime
  - **cats-effect** for asynchronous runtime
- Networking
  - **fs2** for `Stream` abstraction
  - **scalapb** & **fs2-grpc** for protobuf and grpc on top of cats-effect

## Assumptions on Storage

Let combined size of all input files be $M$ and combined size of local storage size of all worker machines be $N$. We assume $N$ to be at least twice as large as $M$. In other words, it must be possible to replicate all input files at least once. No assumptions are made regarding available storage of each worker machines.

## Job Scheduling System

### System Configuration

A cluster consists of multiple slave machines and a single master machine. Slave machines run multiple worker programs and are thus called worker machines. A master machine runs scheduler programs and is called the scheduler machine.
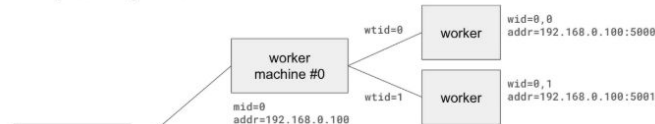
A worker is a program that executes a job upon request from a scheduler. Each worker receives a job request from a scheduler and returns the result using a network. Each worker is assigned a unique IP:PORT pair.

Each worker is assigned a **worker ID** or **wid**. Worker ID is a tuple (mid, wtid) where **mid** (machine ID) identifies a physical machine and **wtid** (worker thread ID) identifies each thread in a machine.
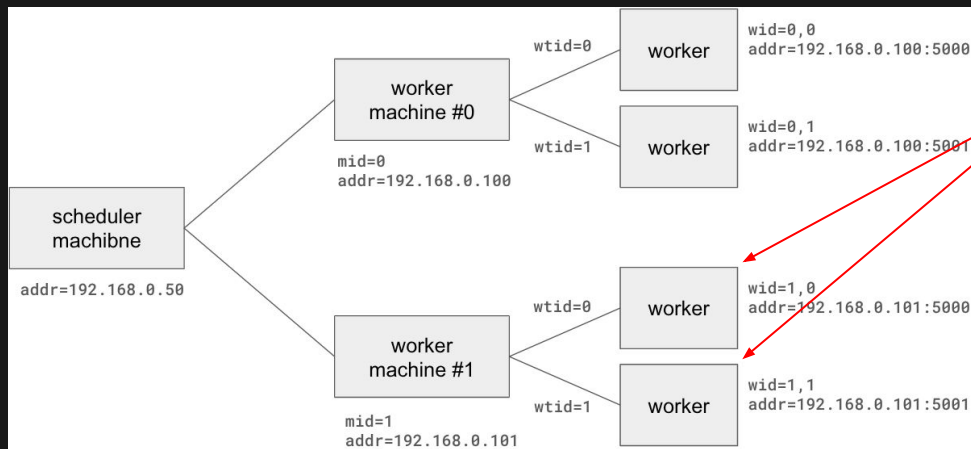
A scheduler is a program that assigns jobs to workers. Its primary role is to distribute jobs so as to minimize communication overhead among workers.

The following figure shows an example system configuration with two worker machines each running two workers as threads.

Example configuration

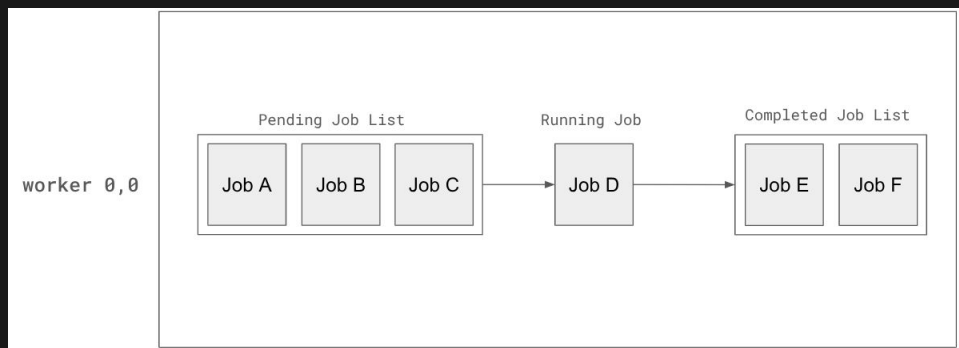worker machine #0
mid=0
addr=192.168.0.100

wtid=0 → worker  wid=0,0  addr=192.168.0.100:5000
wtid=1 → worker  wid=0,1  addr=192.168.0.100:5001

# Design — High Level



naïve

sampling
+
soriting
+
partitioning
+
merging

# *Design — High Level*

machine fault

file replication

error

heartbeat

reschedule

greeting



job created → pending — job sent to worker → running — success message received → completed

rescheduled due to machine fault

# Design — Code Level

```
        ┌──────────────┐
        │    master    │
        │              │
┌──────────┐──────────┘
│  jobs    │
│          │──────┐
└──────────┘      │
        ┌──────────────┐
        │    worker    │
        │              │
        └──────────────┘
```
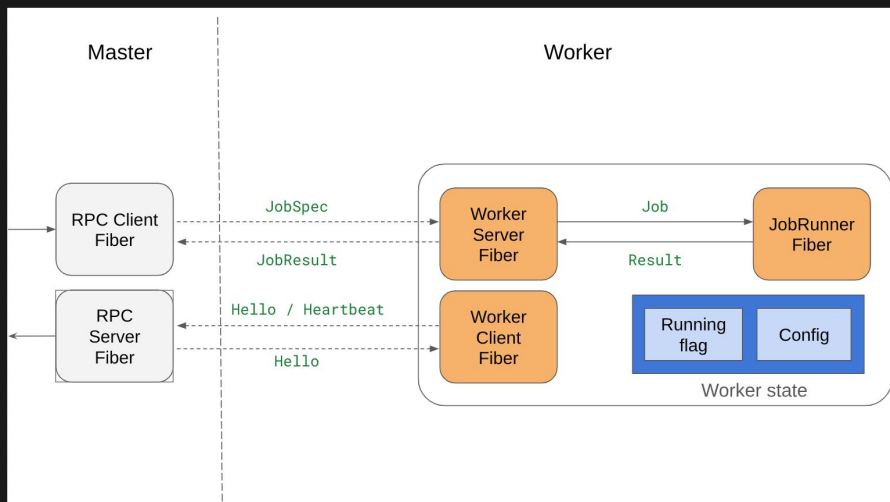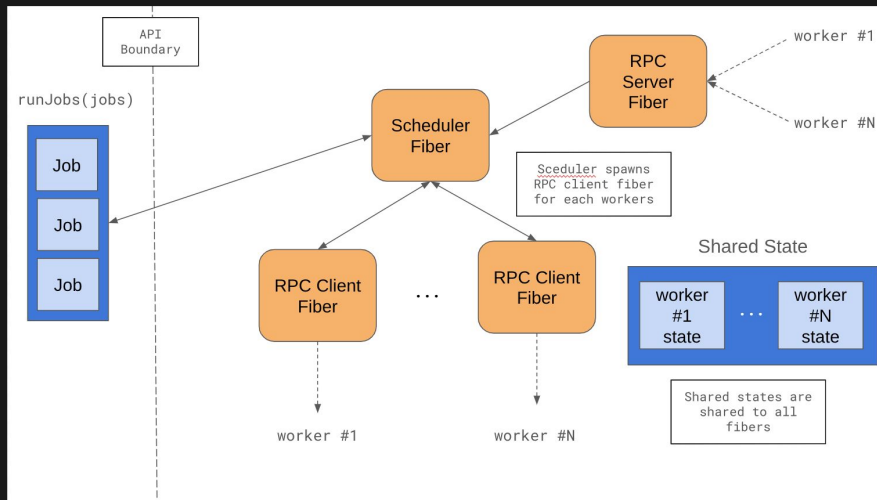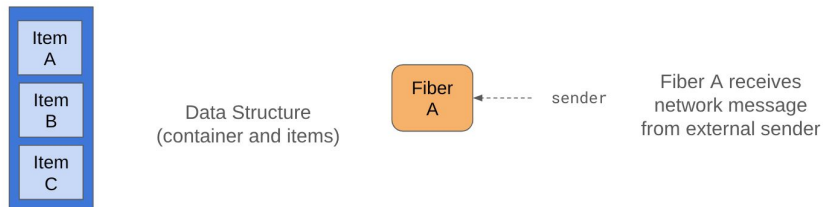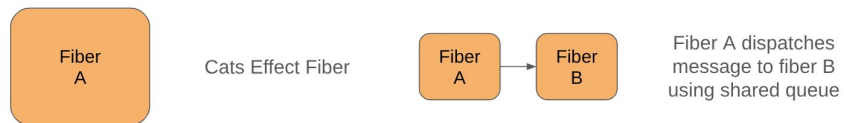
distributed job system
infrastructure

implements distributed
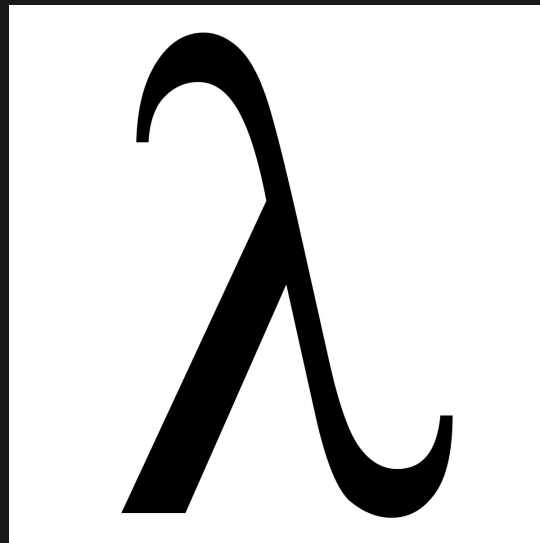sorting on top of **jobs**

# Design - Code Level



Cats Effect Fiber

Fiber A dispatches message to fiber B using shared queue

Data Structure (container and items)

Fiber A receives network message from external sender

API Boundary

runJobs(jobs)

RPC Server Fiber

worker #1

worker #N

Scheduler Fiber

Sceduler spawns RPC client fiber for each workers

RPC Client Fiber

RPC Client Fiber

Shared State

worker #1 state

worker #N state

Shared states are shared to all fibers

worker #1

worker #N

Master

Worker

RPC Client Fiber

JobSpec

Worker Server Fiber

Job

JobRunner Fiber

JobResult

Result

RPC Server Fiber

Hello / Heartbeat

Worker Client Fiber

Hello

Running flag

Config

Worker state

# Design — Code Level

gRPC

File System



Service

λ

dependency injection

functional

# Repository

https://github.com/Quasar-Kim/332project
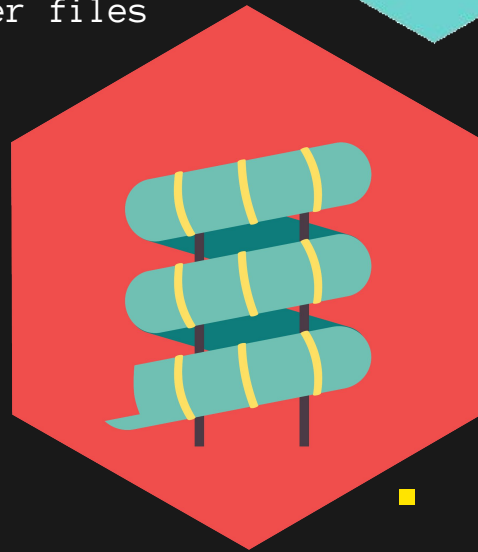
# Libraries & Minor Details

‣ Runtime
  ◦ **cats-effect** for asynchronous runtime
‣ Networking
  ◦ **fs2** for `Stream` abstraction
  ◦ **scalapb** & **fs2-grpc** for protobuf and grpc
    on top of cats-effect
  ◦ (maybe) **fs2-ftp** for transmission of larger files
  ◦ **grpc-netty** for grpc transport
  ◦ **http4s** for web ui
‣ Logging
  ◦ **log4cats** for logging
‣ Testing
  ◦ **scalatest** as testing framework
  ◦ **scalamock** for stubbing and mocking
‣ Programming environment: Scala 2.13.16
‣ **gensort** for input data:
  $ ./gensort -a 500000 ./input.txt # 50MB

# Milestones

‣ Project design finalization (~ ~~week 4.5: 11/13 (Thu)~~ → week 5: 11/16 (Sun))

‣ Being able to sort small data on a local machine (~ week 7: 11/30)

  ○ sorting in an ideal condition. (~ week 6: 11/23)

  ○ presence of faults (short or long period / during each stage)
    (~ week 6.5: 11/27)

  ○ (optional) worst cases (~ week 7: 11/30)

    ▪ lack of disk spaces

    ▪ skewed input key distribution

    ▪ unbalanced disk usage

    ▪ etc.

‣ Running the project on actual (~ week 8: 12/7)

expectations vs. reality