# System Design of Job scheduler in Golang

# System Design - Intro

Goal setting: design a microservices system work that periodically crawls job information, retry if failed. Then do the comparison between user preference and send out information through email.

Codebase: https://github.com/HarrisonLL/Distributed_Job_Crawler
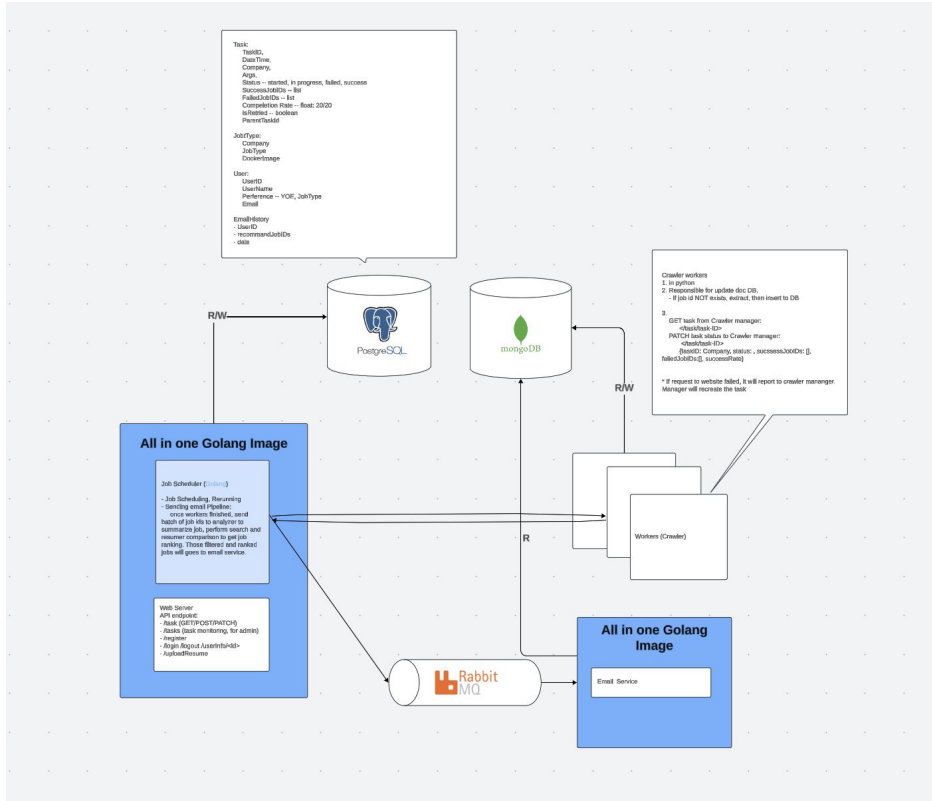
**Design choice:**

**Golang** for job scheduler manager, because it is designed for high concurrent applications, and easy to communicate with kubernetes services.

**RabbitMQ** for emailing queuing, and not used for queuing tasks. Since only the design only targets on spawn and monitor < 10 individual jobs, parallel processes are fine if there are enough CPU cores. However, when numbers of users increase, say hundreds to thousands user, queueing system is needed. This project simulate this situation.

**Python** for crawler. Python has build in easy-to-use crawling and html parser modules like selumni and beautiful soup.

**Docker** for containerization. I separated components into different docker containers, which can be managed by kubernetes or docker-compose.

# System Design



Task:
    TaskID,
    DateTime,
    Company,
    Args,
    Status -- started, in progress, failed, success
    SuccessJobIDs -- list
    FailedJobIDs -- list
    Completion Rate -- float: 20/20
    IsRetried -- boolean
    ParentTaskId

JobType:
    Company
    JobType
    DockerImage

User:
    UserID
    UserName
    Preference -- YOE, JobType
    Email

EmailHistory
 - UserID
 - recommandJobIDs
 - data

Crawler workers
1. in python
2. Responsible for update doc DB,
   - if job id NOT exists, extract, then insert to DB

3.
    GET task from Crawler manager:
        </task/task-ID>
    PATCH task status to Crawler manager:
        </task/task-ID>
        {taskID: Company, status: , successJobIDs: [],
failedJobIDs], successRate}

* If request to website failed, it will report to crawler manager. Manager will recreate the task

R/W

R/W

R

**All in one Golang Image**

Job Scheduler (Golang)

- Job Scheduling, Rerunning
- Sending email Pipeline:
    once workers finished, send batch of job ids to analyrer to summarize job, perform search and resumer comparison to get job ranking. Those filtered and ranked jobs will goes to email service

Web Server
API endpoint
- /task (GET/POST/PATCH)
- /tasks (task monitoring, for admin)
- /register
- /login /logout /user/info/<id>
- /uploadResume

Workers (Crawler)

**All in one Golang Image**

Email Service

# Demo

Step 1. docker-compose up middleware



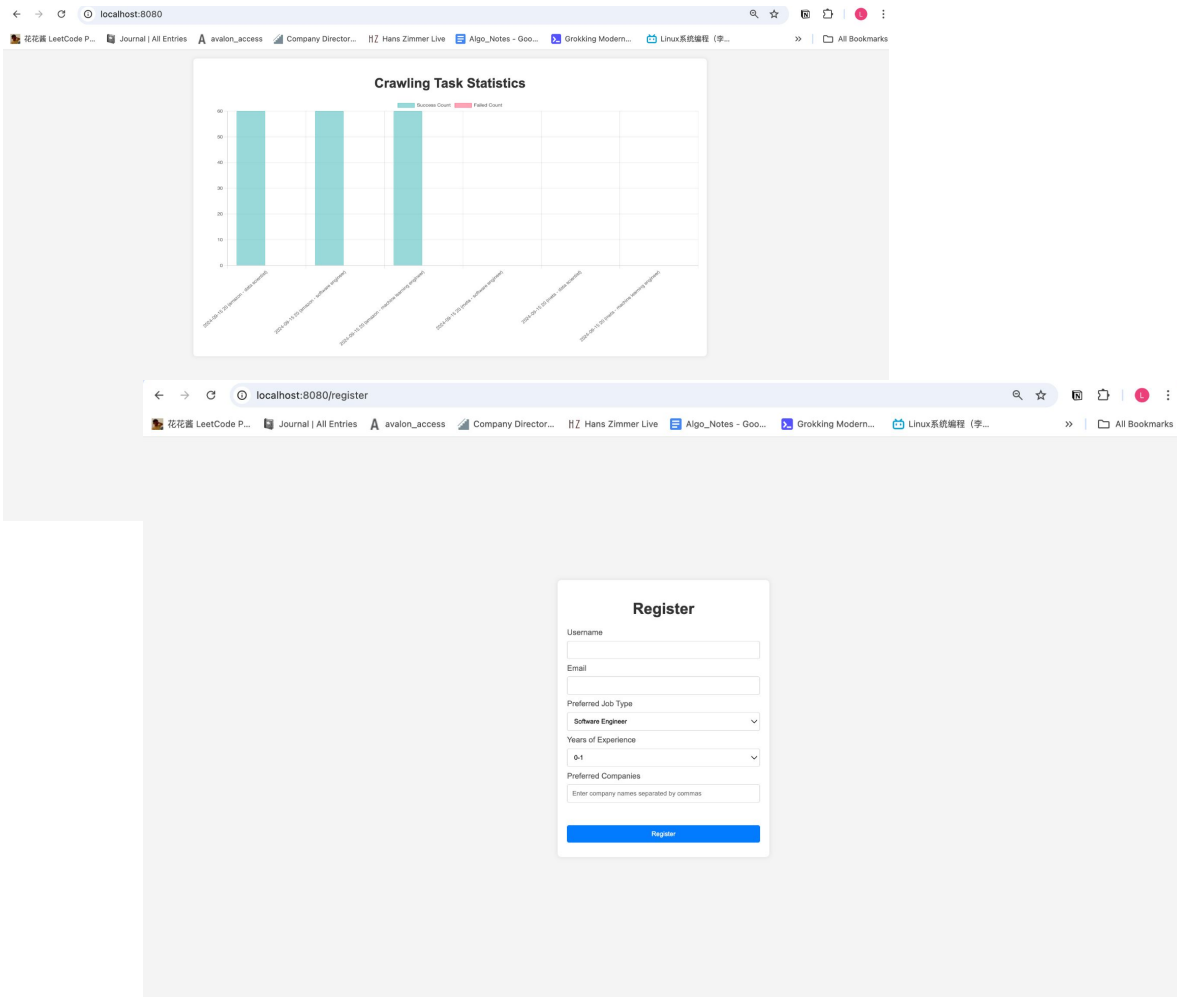Step 2: start two process. Also runnable through docker

# Demo

Steps 3:

Jobs is scheduled every 6 hours, retry happens in hour of the initial job. Once job is found it will send emails to users.

Sample email:

# Demo

Stats and register UI:

# Study notes:

Since I am very new to GoLang, I have learned a lot from doing this project.

1. GoRoutine: go's simple way of doing multithreading
2. GoChannel: go's way to sync threads. We can also run as infinite for loop to listen to something from message queue. (refer: https://www.rabbitmq.com/tutorials/tutorial-one-go)
3. Other than that, I have also experimented some system programming. For example, the code on the right shows I have started many threads to run a python program, for each thread, I also started additional thread to wait till the program finishes and release its resource. (refer: https://pkg.go.dev/os/exec )

```go
go func(jobType models.JobType) {
    pythonCmdDir := os.Getenv("PYTHONFILEPATH")
    pythonCmd := exec.Command("python3", "main.py",
        "--job_type", jobType.JobTypeName,
        "--location", "USA",
        "--company", jobType.CompanyName,
        "--task_id", taskID,
    )
    pythonCmd.Env = append(os.Environ(), envVars...)
    pythonCmd.Dir = pythonCmdDir
    var stderr bytes.Buffer
    pythonCmd.Stderr = &stderr
    if err := pythonCmd.Start(); err != nil {
        log.Printf("Failed to start crawler for company %s: %v", jobType.CompanyName, err, stderr.String()
    } else {
        log.Printf("Started Python crawler for company %s", jobType.CompanyName)
        args := models.JSONMap{
            "job_type": jobType.JobTypeName,
            "location": "USA",
            "company":  jobType.CompanyName,
        }
        err = database.CreateTask(taskID, "", args, false, "")
        if err != nil {
            log.Printf("Failed to create task for company %s: %v", jobType.CompanyName, err)
        }
        // Start a thread to wait till process finishes and release its resource
        go func() {
            if err := pythonCmd.Wait(); err != nil {
                log.Printf("Python crawler for company %s finished with error: %v", jobType.CompanyName, e
                DBerr := database.UpdateTaskStatus(taskID, "", models.Error)
                if err != nil {
                    log.Printf("Failed to update task %s: %v", taskID, DBerr)
                }
            } else {
                log.Printf("Python crawler for company %s finished successfully", jobType.CompanyName)
            }
        }()
    }
}(jobType)
```