



Разработка симулятора вычислительного кластера

Выполнил: Макогон Артём Аркадьевич, БПМИ206

Руководитель: Сухорослов Олег Викторович, к.т.н, доцент НИУ ВШЭ

June 3, 2024



- Вычислительные кластеры активно используются в исследованиях и компаниях.
- От качества алгоритма планирования задач зависит эффективность использования ресурсов.
- При разработке алгоритмов их необходимо тестировать.



- Вычислительные кластеры активно используются в исследованиях и компаниях.
- От качества алгоритма планирования задач зависит эффективность использования ресурсов.
- При разработке алгоритмов их необходимо тестировать.



> Использовать реальный кластер долго и дорого.



- Вычислительные кластеры активно используются в исследованиях и компаниях.
- От качества алгоритма планирования задач зависит эффективность использования ресурсов.
- При разработке алгоритмов их необходимо тестировать.



> Использовать реальный кластер долго и дорого.



✓ Используются **симуляторы**.



Рабочая нагрузка состоит из заданий (job), которые делятся на задачи (tasks).



Рабочая нагрузка состоит из заданий (job), которые делятся на задачи (tasks).

Структуры заданий:

- rigid – фиксированные требования к ресурсам
- moldable – адаптивные требования к ресурсам
- malleable – гибкие требования к ресурсам



Рабочая нагрузка состоит из заданий (job), которые делятся на задачи (tasks).

Структуры заданий:

- rigid – фиксированные требования к ресурсам
- moldable – адаптивные требования к ресурсам
- malleable – гибкие требования к ресурсам

Требования к планированию:

- Раздельное планирование (MapReduce)
- Комплексное планирование (Gang Scheduling, Slurm)



Standard Workload Format (SWF)

- CPU/RAM

Custom workloads

- CPU/RAM/disk/network

Standard Workload Format (SWF)

- CPU/RAM
- Дано время исполнения и требования к ресурсам
- Нагрузка вычисляется тривиально

Custom workloads

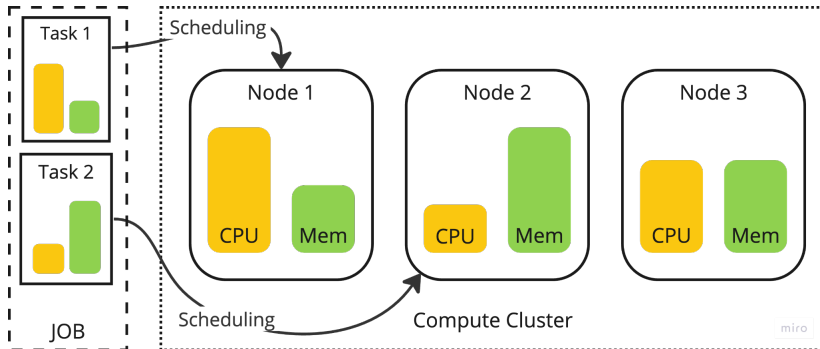
- CPU/RAM/disk/network
- Дана нагрузка и ресурсы
- Время вычисляется с помощью моделей

Standard Workload Format (SWF)

- CPU/RAM
- Дано время исполнения и требования к ресурсам
- Нагрузка вычисляется тривиально
- Используется в популярных трейсах (Google, Alibaba, и т.д.)

Custom workloads

- CPU/RAM/disk/network
- Дана нагрузка и ресурсы
- Время вычисляется с помощью моделей
- Обычно NDA

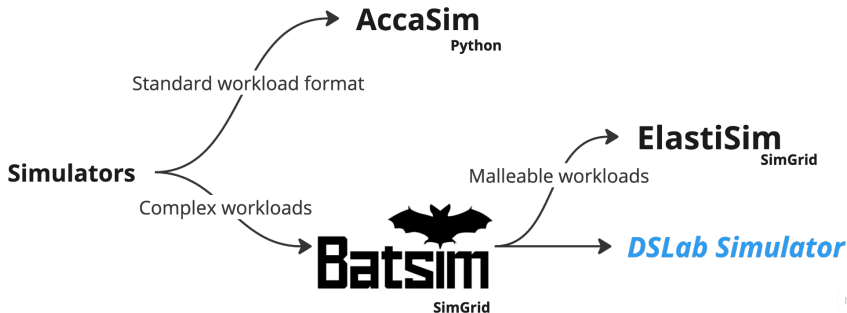


Простая модель архитектуры кластера

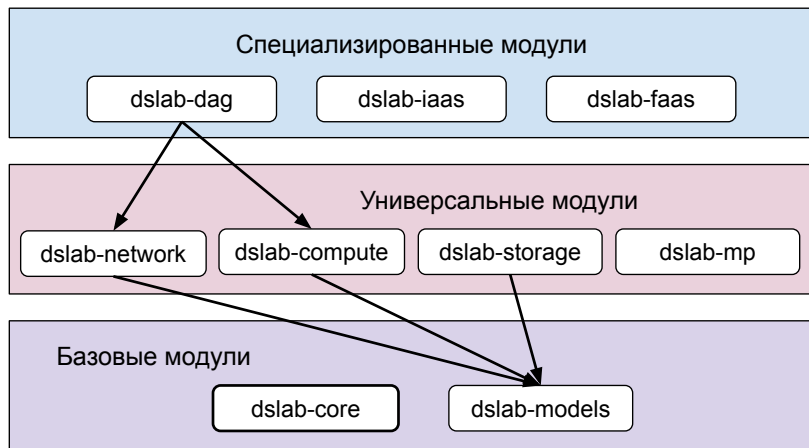


Цель проекта – разработать симулятор вычислительного кластера на базе фреймворка DSLab. Для этого необходимо:

- Изучить литературу по теме и существующие симуляторы, подготовить обзор с анализом их преимуществ и недостатков.
- Реализовать компоненты симулятора.
- Написать документацию и тесты.
- Подготовить и провести эксперименты, демонстрирующие работоспособность симулятора.



Существующие симуляторы



Архитектура фреймворка DSLab

Реализация симулятора

Асинхронное управление событиями. Комбинаторы futures



```
async fn process_task(&self, req: TaskRequest) {  
    let mut task = TaskInfo {req};  
  
    self.download_data(&task).await;  
    self.read_data(&task).await;  
    self.run_task(&task).await;  
    self.write_data(&task).await;  
    self.upload_result(&task).await;  
}
```

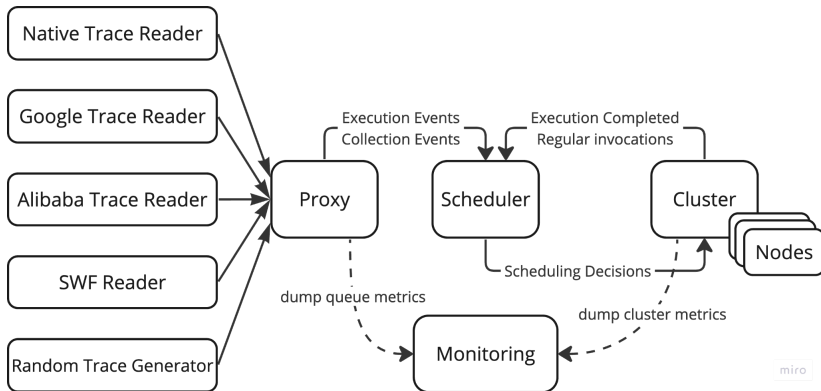
Пример последовательных задач.

```
async fn run(&self, args: JobArgs){  
    futures::join!(  
        self.download_data(args.node_1)  
        self.download_data(args.node_2)  
    )  
}
```

Пример параллельных задач.

Реализация симулятора

Архитектура симулятора



Архитектура симулятора

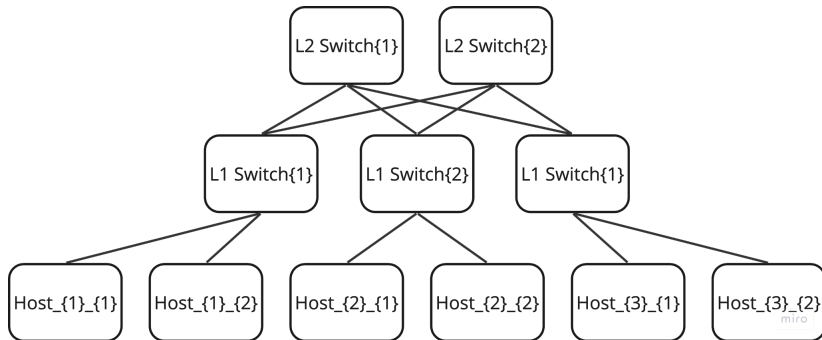


Название модуля	Используемый модуль	Обязательный
compute	dslab_compute::multicore::Compute	да
network	dslab_network::Network	нет
disk	dslab_storage::disk::Disk	нет

Модули каждого сервера в кластере

Реализация симулятора

Модель сети Fat-Tree-Topology



Модель сети Fat-Tree-Topology в кластере

Реализация симулятора

Деление на Collection и Execution



Поле	Тип	Описание
id	u64	Уникальный идентификатор задания
user	Option<String>	Пользователь, который запускает задание.
priority	Option<u64>	Приоритет задания.
...

Описание структуры данных Collection

Поле	Тип	Описание
id	u64	Уникальный идентификатор задачи
collection_id	Option<u64>	Индекс задания
time	f64	Время, когда задача становится доступной для планирования
resources	ResourceRequirements	Структура требуемых ресурсов
profile	Rc<dyn ExecutionProfile>	Профиль нагрузки
...

Описание структуры данных Execution



```
#[async_trait(?Send)]  
pub trait ExecutionProfile {  
    async fn run(self: Rc<Self>, processes: &Vec<HostProcessInstance>)  
}
```

Интерфейс ExecutionProfile

```
async fn sleep(&self, time: f64);  
async fn run_compute_work(&self, compute_work: f64);  
async fn transfer_data_process(&self, size: f64, dst: ProcessId);  
async fn transfer_data_component(&self, size: f64, dst: Id);  
async fn write_data(&self, size: u64);  
async fn read_data(&self, size: u64);  
async fn deallocate(&self);
```

Интерфейс взаимодействия процессов с кластером

Реализация симулятора

Пример реализации ExecutionProfile



```
pub struct MasterWorkersProfile {  
    pub master_compute_work: f64,  
    pub worker_compute_work: f64,  
    pub data_transfer_size: f64,  
}  
  
#[async_trait(?Send)]  
impl ExecutionProfile for MasterWorkersProfile {  
    async fn run(self: Rc<Self>, processes: &Vec<HostProcessInstance>) {  
        let master = &processes[0];  
        let workers = &processes[1..];  
  
        futures::future::join_all(workers.iter().map(|p| async {  
            master.transfer_data_process(self.data_transfer_size, p.id).await;  
            p.run_compute(self.worker_compute_work).await;  
        })))  
        .await;  
  
        master.run_compute(self.master_compute_work).await;  
    }  
}
```

Реализация симулятора

Пример реализации ExecutionProfile



```
pub struct MasterWorkersProfile {
    pub master_compute_work: f64,
    pub worker_compute_work: f64,
    pub data_transfer_size: f64,
}

#[async_trait(?Send)]
impl ExecutionProfile for MasterWorkersProfile {
    async fn run(self: Rc<Self>, processes: &Vec<HostProcessInstance>) {
        let master = &processes[0];
        let workers = &processes[1..];

        futures::future::join_all(workers.iter().map(|p| async {
            master.transfer_data_process(self.data_transfer_size, p.id).await;
            p.run_compute(self.worker_compute_work).await;
        })).await;

        master.run_compute(self.master_compute_work).await;
    }
}
```

Реализация симулятора

Пример реализации ExecutionProfile



```
pub struct MasterWorkersProfile {
    pub master_compute_work: f64,
    pub worker_compute_work: f64,
    pub data_transfer_size: f64,
}

#[async_trait(?Send)]
impl ExecutionProfile for MasterWorkersProfile {
    async fn run(self: Rc<Self>, processes: &Vec<HostProcessInstance>) {
        let master = &processes[0];
        let workers = &processes[1..];

        futures::future::join_all(workers.iter().map(|p| async {
            master.transfer_data_process(self.data_transfer_size, p.id).await;
            p.run_compute(self.worker_compute_work).await;
        })))
        .await;

        master.run_compute(self.master_compute_work).await;
    }
}
```

Реализация симулятора

Пример реализации ExecutionProfile



```
pub struct MasterWorkersProfile {
    pub master_compute_work: f64,
    pub worker_compute_work: f64,
    pub data_transfer_size: f64,
}

#[async_trait(?Send)]
impl ExecutionProfile for MasterWorkersProfile {
    async fn run(self: Rc<Self>, processes: &Vec<HostProcessInstance>) {
        let master = &processes[0];
        let workers = &processes[1..];

        futures::future::join_all(workers.iter().map(|p| async {
            master.transfer_data_process(self.data_transfer_size, p.id).await;
            p.run_compute(self.worker_compute_work).await;
        })))
        .await;

        master.run_compute(self.master_compute_work).await;
    }
}
```




1. Native – нагрузка в формате YAML



1. Native – загрузка в формате YAML
2. Google Trace
3. Alibaba Trace



1. Native – нагрузка в формате YAML
2. Google Trace
3. Alibaba Trace
4. SWF – Standard Workload Format



1. Native – нагрузка в формате YAML
2. Google Trace
3. Alibaba Trace
4. SWF – Standard Workload Format
5. Synthetic Trace Generator



```
pub trait Scheduler {  
    fn on_host_added(&mut self, ctx: &Context, host: HostConfig);  
    fn on_execution_request(&mut self, ctx: &Context, req: Execution);  
    fn on_collection_request(&mut self, ctx: &Context, req: Collection);  
    fn on_execution_finished(&mut self, ctx: &Context, execution_id: u64);  
    fn on_host_notification(&mut self, ctx: &Context, host_id: Id, resources:  
        ResourcesPack);  
}
```

Интерфейс планировщика Scheduler

```
pub trait Scheduler {  
    fn on_host_added(&mut self, ctx: &Context, host: HostConfig);  
    fn on_execution_request(&mut self, ctx: &Context, req: Execution);  
    fn on_collection_request(&mut self, ctx: &Context, req: Collection);  
    fn on_execution_finished(&mut self, ctx: &Context, execution_id: u64);  
    fn on_host_notification(&mut self, ctx: &Context, host_id: Id, resources:  
        ResourcesPack);  
}
```

Интерфейс планировщика Scheduler

- Возможности контекста:

1. Запланировать/отменить задачу (полная отмена недоступна)

```
pub trait Scheduler {  
    fn on_host_added(&mut self, ctx: &Context, host: HostConfig);  
    fn on_execution_request(&mut self, ctx: &Context, req: Execution);  
    fn on_collection_request(&mut self, ctx: &Context, req: Collection);  
    fn on_execution_finished(&mut self, ctx: &Context, execution_id: u64);  
    fn on_host_notification(&mut self, ctx: &Context, host_id: Id, resources:  
        ResourcesPack);  
}
```

Интерфейс планировщика Scheduler

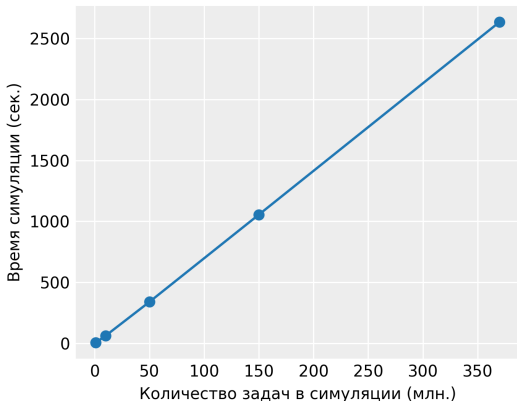
- Возможности контекста:
 1. Запланировать/отменить задачу (полная отмена недоступна)
 2. Доступ к генератору случайных чисел симуляции



```
pub trait Scheduler {  
    fn on_host_added(&mut self, ctx: &Context, host: HostConfig);  
    fn on_execution_request(&mut self, ctx: &Context, req: Execution);  
    fn on_collection_request(&mut self, ctx: &Context, req: Collection);  
    fn on_execution_finished(&mut self, ctx: &Context, execution_id: u64);  
    fn on_host_notification(&mut self, ctx: &Context, host_id: Id, resources:  
        ResourcesPack);  
}
```

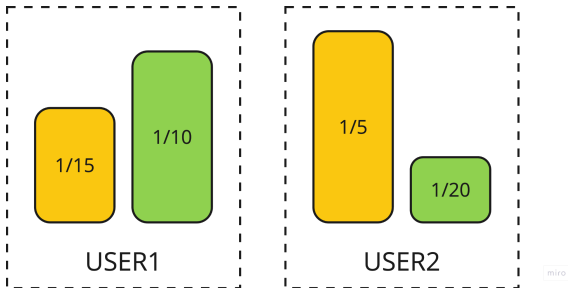
Интерфейс планировщика Scheduler

- Возможности контекста:
 1. Запланировать/отменить задачу (полная отмена недоступна)
 2. Доступ к генератору случайных чисел симуляции
- Можно реализовать планировщик как компонент симуляции и получить больше возможностей (например, асинхронность).

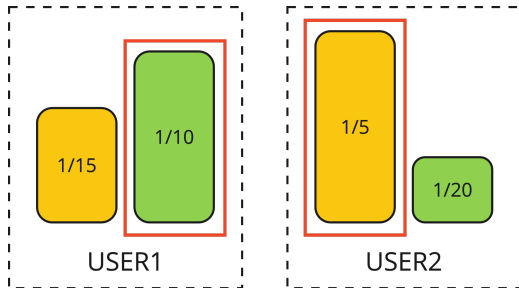


- На трейсе длиной 2.5 дня (350 млн задач, 25GB) время работы 45 минут
- Ускорение в 85 раз относительно реальной работы.
- Необходимо сортировать трейс по времени.

Время работы симуляции в зависимости от количества задач



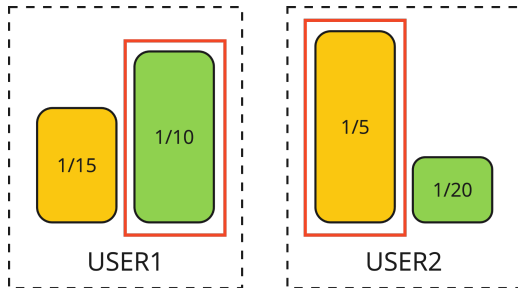
Потребление ресурсов пользователями



Dominant shares:

- USER1 = $1/10$
- USER2 = $1/5$

Потребление ресурсов пользователями



Потребление ресурсов пользователями

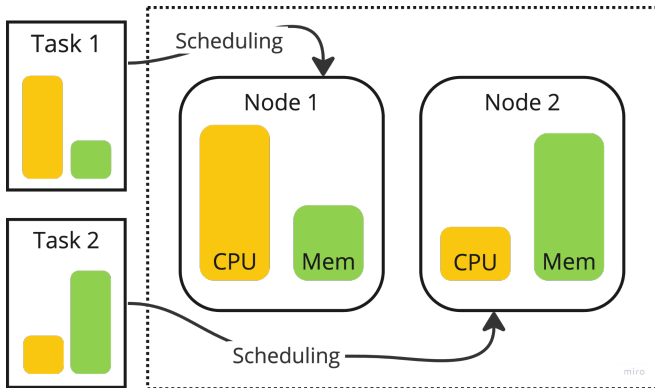
Dominant shares:

- $USER1 = 1/10$
- $USER2 = 1/5$
- USER1 первый получит ресурсы

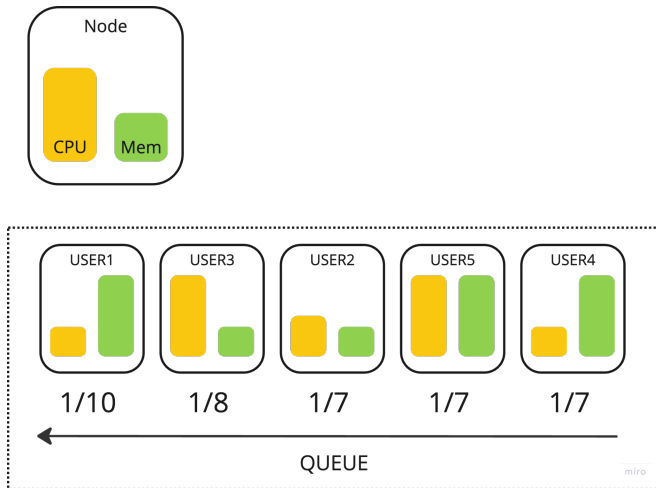


Метрика «упаковки»: $H_{angle} = 1 - \cos(\angle(Task_{res}, Server_{res}))$

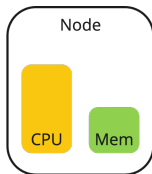
Метрика «упаковки»: $H_{angle} = 1 - \cos(\angle(Task_{res}, Server_{res}))$



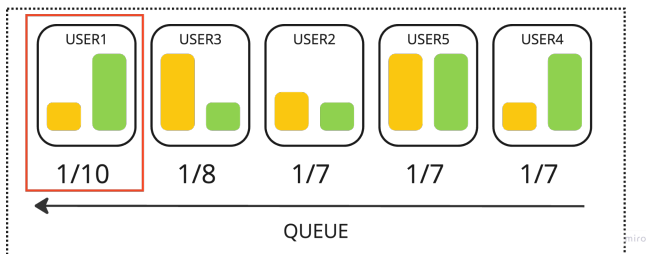
Пример работы алгоритма Tetris



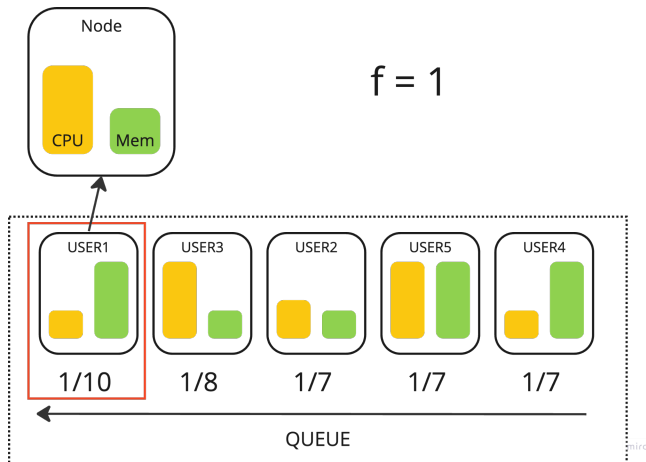
Процесс принятия решения о планировании



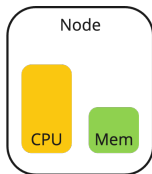
$$f = 1$$



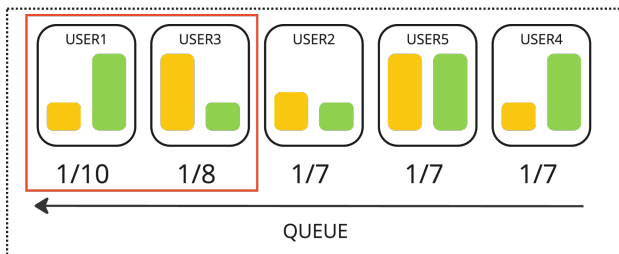
Процесс принятия решения о планировании



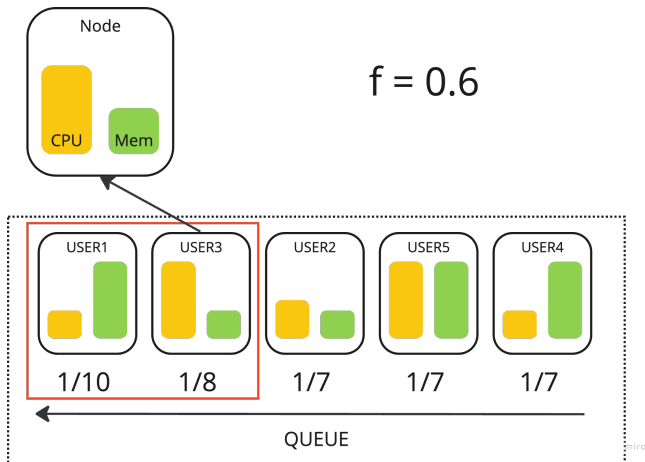
Процесс принятия решения о планировании



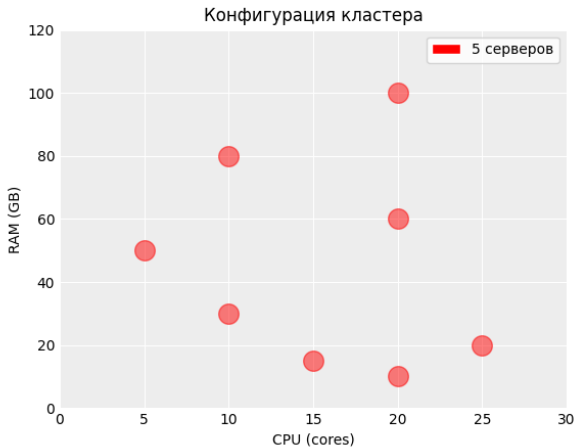
$$f = 0.6$$



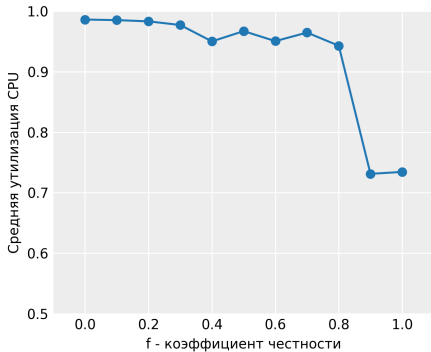
Процесс принятия решения о планировании



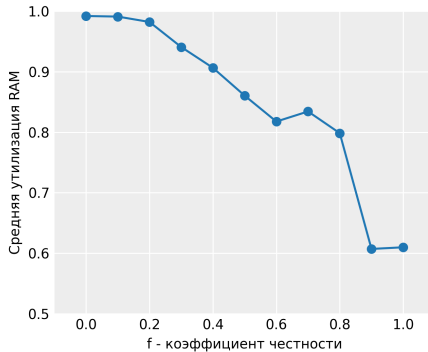
Процесс принятия решения о планировании



Распределение ресурсов на серверах кластера



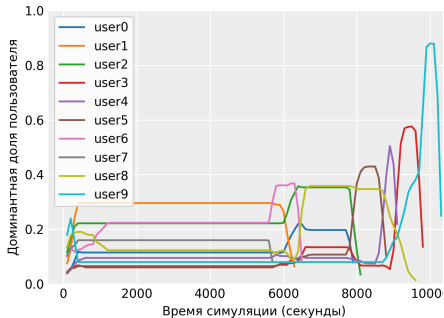
Утилизация CPU при разных f



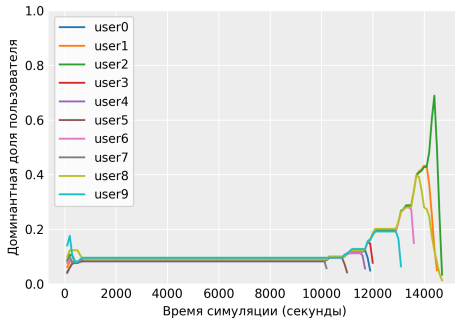
Утилизация RAM при разных f

Эксперименты

Результаты честности



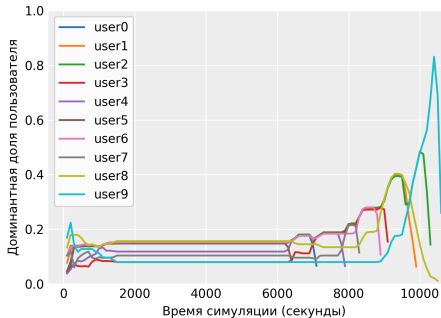
$f = 0$



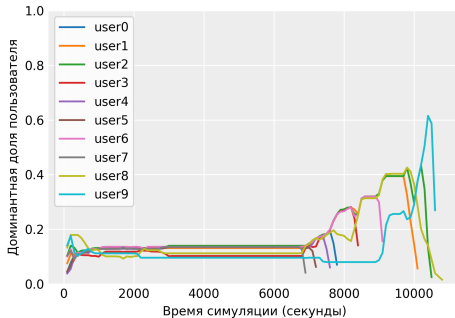
$f = 1$

Эксперименты

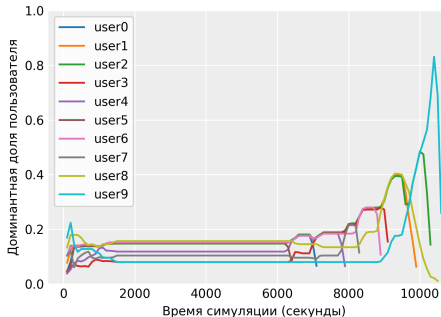
Результаты честности



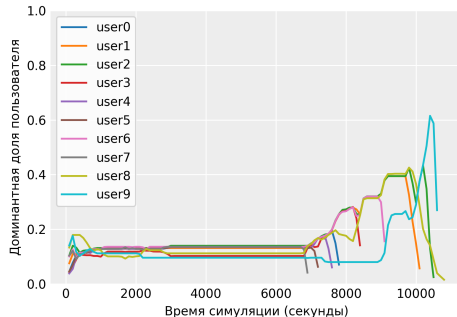
$f = 0.5$



$f = 0.6$



$f = 0.5$



$f = 0.6$

Вывод: $f = 0.6$ – оптимальный параметр честности для такой конфигурации кластера и задач пользователей.



- ✓ Реализовать эффективный симулятор вычислительного кластера.



- ✓ Реализовать эффективный симулятор вычислительного кластера.
- ✓ Поддержать несколько моделей рабочей нагрузки и планирования заданий.



- ✓ Реализовать эффективный симулятор вычислительного кластера.
- ✓ Поддерживать несколько моделей рабочей нагрузки и планирования заданий.
- ✓ Поддерживать чтение трейсов нагрузки компаний (Google, Alibaba).



- ✓ Реализовать эффективный симулятор вычислительного кластера.
- ✓ Поддержать несколько моделей рабочей нагрузки и планирования заданий.
- ✓ Поддержать чтение трейсов нагрузки компаний (Google, Alibaba).
- ✓ Поддержать адаптивный сбор статистики для анализа результатов (нагрузка на кластер, честность, статистика очереди)



- ✓ Реализовать эффективный симулятор вычислительного кластера.
- ✓ Поддержать несколько моделей рабочей нагрузки и планирования заданий.
- ✓ Поддержать чтение трейсов нагрузки компаний (Google, Alibaba).
- ✓ Поддержать адаптивный сбор статистики для анализа результатов (нагрузка на кластер, честность, статистика очереди)
- ✓ Реализовать несколько популярных алгоритмов (FCFS, DRF, Tetris).



- ✓ Реализовать эффективный симулятор вычислительного кластера.
- ✓ Поддержать несколько моделей рабочей нагрузки и планирования заданий.
- ✓ Поддержать чтение трейсов нагрузки компаний (Google, Alibaba).
- ✓ Поддержать адаптивный сбор статистики для анализа результатов (нагрузка на кластер, честность, статистика очереди)
- ✓ Реализовать несколько популярных алгоритмов (FCFS, DRF, Tetris).
- ✓ Доработать реализацию асинхронного ядра `dslab-core` и ускорить его работу в 2 раза.



- ✓ Реализовать эффективный симулятор вычислительного кластера.
- ✓ Поддерживать несколько моделей рабочей нагрузки и планирования заданий.
- ✓ Поддерживать чтение трейсов нагрузки компаний (Google, Alibaba).
- ✓ Поддерживать адаптивный сбор статистики для анализа результатов (нагрузка на кластер, честность, статистика очереди)
- ✓ Реализовать несколько популярных алгоритмов (FCFS, DRF, Tetris).
- ✓ Доработать реализацию асинхронного ядра `dslib-core` и ускорить его работу в 2 раза.

Главное достижение: предоставить возможность описывать рабочую нагрузку в виде асинхронных примитивов языка Rust.



1. Сравнение с симулятором BatSim
2. Входные профили нагрузки через yaml-файл
3. Профили по умолчанию
4. Выходные данные и monitoring
5. Ускорение async-dslab-core в 2 раза
6. Модель мгновенной честной доли с отменой задач
7. Дискретно-событийное моделирование



- Разработан на базе фреймворка SimGrid.
- Поддерживает SWF и пользовательские профили нагрузки через JSON-файл.
- Алгоритмы планирования подключаются к симулятору через IPC (inter-process-communication).

```
"jobs": [  
  {"id": "job1",    ...  "res": 4, "profile": "sequence"},  
],  
  
"profiles": {  
  "homogeneous": {  
    "type": "parallel_homogeneous",  
    "cpu": 10e6,  
    "com": 1e6  
  },  
  "sequence": {  
    "type": "composed",  
    "repeat" : 4,  
    "seq": ["simple", "homogeneous", "simple"]  
  },  
}
```