

Московский физико-технический институт  
(государственный университет)

Физтех-школа прикладной математики и информатики  
Кафедра алгоритмов и технологий программирования

# Система Hjudge

Как автоматизировать проверку  
заданий

при изучении работы с большими  
данными

О.Н. Ивченко

А.А. Драль

М.А. Ройтберг

1. Введение

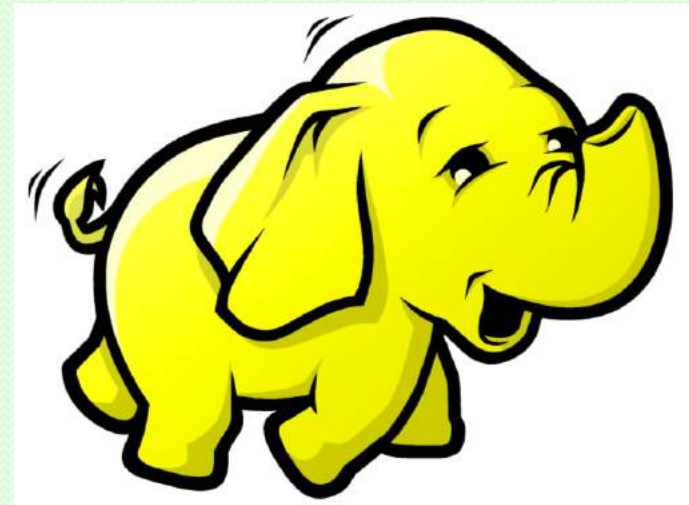
# Большие данные



# Экосистема Hadoop 2008

- Основные компоненты

- Hadoop (2005) – фреймворк, реализующий парадигму MapReduce
  - включает распределенную файловую систему HDFS
- Spark – фреймворк, обеспечивающий более быструю, по сравнению с Hadoop, обработку данных
  - ленивая обработка
  - обработка в реальном времени
- Hive – SQL-обёртка над Hadoop,
- HBase – БД NoSQL-типа, обеспечивает хранение больших разреженных массивов данных



# Экосистема Hadoop

- Где используется

Яндекс

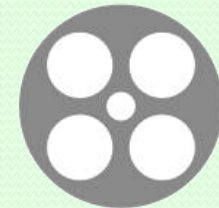
LinkedIn

IBM



Google

ORACLE®



КиноПоиск.Ру  
найди своё кино!

# Описание

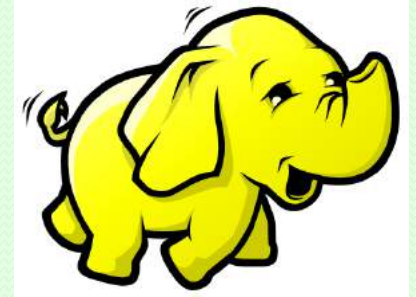
- Хранение и Обработка Больших Объемов Данных
  - <http://atp-fivt.org/kurs-xobod/okurs/>
- Читается с весны 2015 г.
- Направлен на работу с данными, для обработки которых нужно > **1 машины**
- Предмет рассмотрения – экосистема Hadoop
- Практическая часть – программистские задачи по 4 основным модулям:
  - MapReduce (Hadoop Java API, Hadoop Streaming)
  - Hive
  - Spark
  - HBase

## Примеры задач

- **Задача 1:** Посчитать число вхождений слова (term) в каждую статью (term frequency, tf). Слова очистить от знаков пунктуации.

**Данные:** статьи википедии (id <пробел> text)

**Формат результата:** term, article id, tf



- **Задача 2:** Посчитать богатство языка (число уникальных слов) каждого персонажа и число реплик.

**Данные:** текст «Горе от ума» (персонаж <пробел> строка реплики)

**Формат результата:** персонаж, число уникал. слов, число реплик

- **Задача 3:** Вывести TOP10 наиболее посещаемых страниц за указанный интервал времени. На входе - интервал времени, на выходе - список наиболее популярных страниц.

**Данные:** история посещений пользователями страниц.



## Примеры задач

- **Задача 4:** Вывести суммарное распределение количества посещений по часам (для каждого часа в сутках вывести количество посещений, пришедшее в этот час).

**Данные:** история посещений пользователями страниц.



- **Задача 5:** Требуется написать пользовательскую функцию Nive. На вход ей подаётся адрес сети и адрес маски (именно в таком порядке). На выходе получаем множество возможных IP-адресов в этой сети.

**Данные:** история посещений пользователями страниц.



# Данные в задачах

- имеют объём до нескольких десятков Гб
  - ✓ (нельзя быстро обработать на 1 машине)
- читаются из файла в HDFS
- записываются в HDFS или выводятся в КОНСОЛЬ





### 3. Проверка заданий

# Сдача задания студентом: PIAZZA

The screenshot shows the Piazza forum interface. At the top, there are navigation tabs: "Q & A", "Resources", "Statistics", and "Manage Class". The user "VeLKerr" is logged in. Below the navigation, there are filters for "Updated", "Unresolved", and "Following". A search bar is present with the text "Search or add a post...".

The main content area displays a "note" titled "Начало курса" (Course Start) with 62 views. The text of the note includes:

- Добро пожаловать в Большие Данные.
- По всем вопросам про курс давайте общаться здесь, на piazza. Почта - это только для тех, кто не смог сюда зайти (или потерял доступ).
- Общие правила такие. Создавайте новые открытые темы, если у вас появляются вопросы. Знаете ответ на чужой вопрос - помогайте. Пишите индивидуальные сообщения для Instructors при сдаче задания.
- В разделе Resources выложена первая лекция и книжки про Hadoop. The Definitive Guide - наиболее полная и полезная при первом знакомстве с Hadoop.
- Для работы нам понадобится доступ к серверу bds01 по ssh. Под Linux/MacOS для этого есть терминал. Под Windows надо установить любой ssh клиент PuTTY: <http://www.chiark.greenend.org.uk/~sgtatham/putty/download.html>
- Далее нужно выбрать логин по такому правилу:  
LOGIN = s<GROUP><NUMBER>  
GROUP - ваш номер группы: 193 или 194  
NUMBER - любой двузначный номер, от 01 до 50.

On the left side, there is a list of private questions:

- Private **3 задание. Решение.** 11:36PM
- Private **ДЗ-4** 3:30PM
- Private **hw3 done** Fri
- Private **Задание 4** Fri
- Private **Логин s19333** Fri

#### ? private question ☆

### ДЗ-2 Hadoop Streaming

Здравствуйте!

Посылаю ДЗ-2 по Hadoop Streaming

1. Логин **s19444** задача №5
2. Директорию с исходниками на сервере: `/home/gr194/s19444/BIGDATA_HOMEWORK`
3. Как собирать и запускать: В директории с исходниками есть файл `run.sh`
4. Директорию с результатом в HDFS: `/user/s19444/stream_homework`

#### ? private question ☆

### ДЗ 5 HIVE

Мой логин  
s19413

Директория на сервере с исходниками:  
`/home/gr194/s19413/hive`

Как запускать, т.е. какие аргументы у задачи:  
для создания таблицы:

```
./run.sh
```

задание 1:  
`./run2.sh 1`

# Процесс проверки заданий

## Этапы проверки одного задания:

1. Сборка на Hadoop-кластере
2. Запуск программы
3. Мониторинг процесса работы
4. Проверка результата
5. Code review

# Процесс проверки заданий

## Этапы проверки одного задания:

1. Сборка на Hadoop-кластере (~ 5 с)
2. Запуск программы (~ 10 с)
3. Мониторинг процесса работы (~ 20 с)
4. Проверка результата (~ 90 с)
5. Code review (~ 120 с)

# 4. Инструменты автоматизации

## 1. Сборка на Hadoop-кластере



# 4. Инструменты автоматизации

1. Сборка на Hadoop-кластере

2. Запуск программы



*(самописный bash-скрипт)*

# 4. Инструменты автоматизации

1. Сборка на Hadoop-кластере

2. Запуск программы



3. Мониторинг процесса работы

- для Hadoop: MrBench (можно подавать на вход свой Jar):

```
$ hadoop jar /usr/lib/hadoop-mapreduce/hadoop-mapreduce-client-jobclient.jar  
mrbench -help
```

MRBenchmark.0.0.2

Usage: mrbench ... [-jar <local path to job jar file containing Mapper and Reducer implementations, default is current jar file>] ...

- для Hive, Spark и HBase подобных инструментов не найдено



*(самописный bash-скрипт)*

# 4. Инструменты автоматизации

1. Сборка на Hadoop-кластере

2. Запуск программы



3. Мониторинг процесса работы

- для Hadoop: MrBench (можно подавать на вход свой Jar)

4. Проверка результата

- diff (colordiff)



(самописный *bash*-скрипт)

```
~$ colordiff -u optable_loo
--- optable_loop.F      2013-04-30 20:21:57.210318683
+++ optable_loop.2.F    2013-04-30 20:23:29.426318985
@@ -340,6 +340,14 @@
         call timer_start(pardo_tserver_timer)
#endif
         call timer_start(timer_ovrhead)
+
+c Some more code to demonstrate the colored diff
+
+         print *, 'blah blah blah'
+
+c end colored diff exmple
+
+         call exec_thread_server(0)
+         call update_timer(timer_ovrhead)
#ifdef VERY_DETAILED_TIMERS
```

# 4. Инструменты автоматизации

1. Сборка на Hadoop-кластере

2. Запуск программы



3. Мониторинг процесса работы

- для Hadoop: MrBench (можно подавать на вход свой Jar)

4. Проверка результата

- diff (colordiff)



*(самописный bash-скрипт)*

```
~$ colordiff -u optable_loo
--- optable_loop.F      2013-04-30 20:21:57.210318683
+++ optable_loop.2.F    2013-04-30 20:23:29.426318985
@@ -340,6 +340,14 @@
         call timer_start(pardo_tserver_timer)
#endif
         call timer_start(timer_ovrhead)
+
+c Some more code to demonstrate the colored diff
+
+         print *, 'blah blah blah'
+
+c end colored diff exmple
+
+         call exec_thread_server(0)
+         call update_timer(timer_ovrhead)
#ifdef VERY_DETAILED_TIMERS
```

- *в задачах нет однозначного решения*  
=> *diff не всегда информативен*
- *данных очень много*



# 4. Инструменты автоматизации

1. Сборка на Hadoop-кластере

2. Запуск программы



3. Мониторинг процесса работы

- для Hadoop: MrBench (можно подавать на вход свой Jar)

4. Проверка результата

- diff (colordiff)

5. Code review



*(самописный bash-скрипт)*



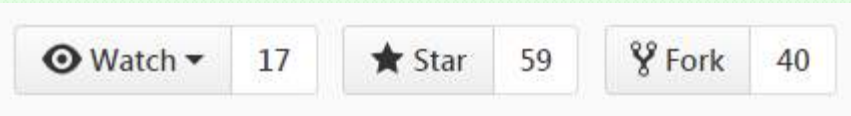
# 4. Инструменты автоматизации

## □ Модульное тестирование

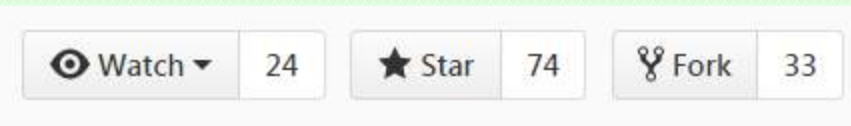
- MapReduce: MRUnit,
- Hive:

Для каждого сервиса существуют свои инструменты модульного тестирования. Они различаются в использовании => их тяжело привести к общему интерфейсу.

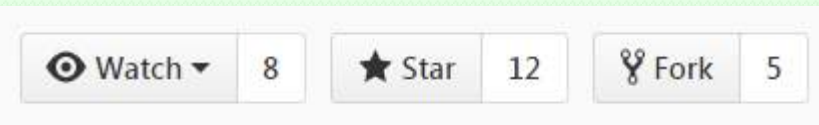
- Hive\_test,



- HiveRunner,



- HiveQLUnit



- Spark

- Spark testing base

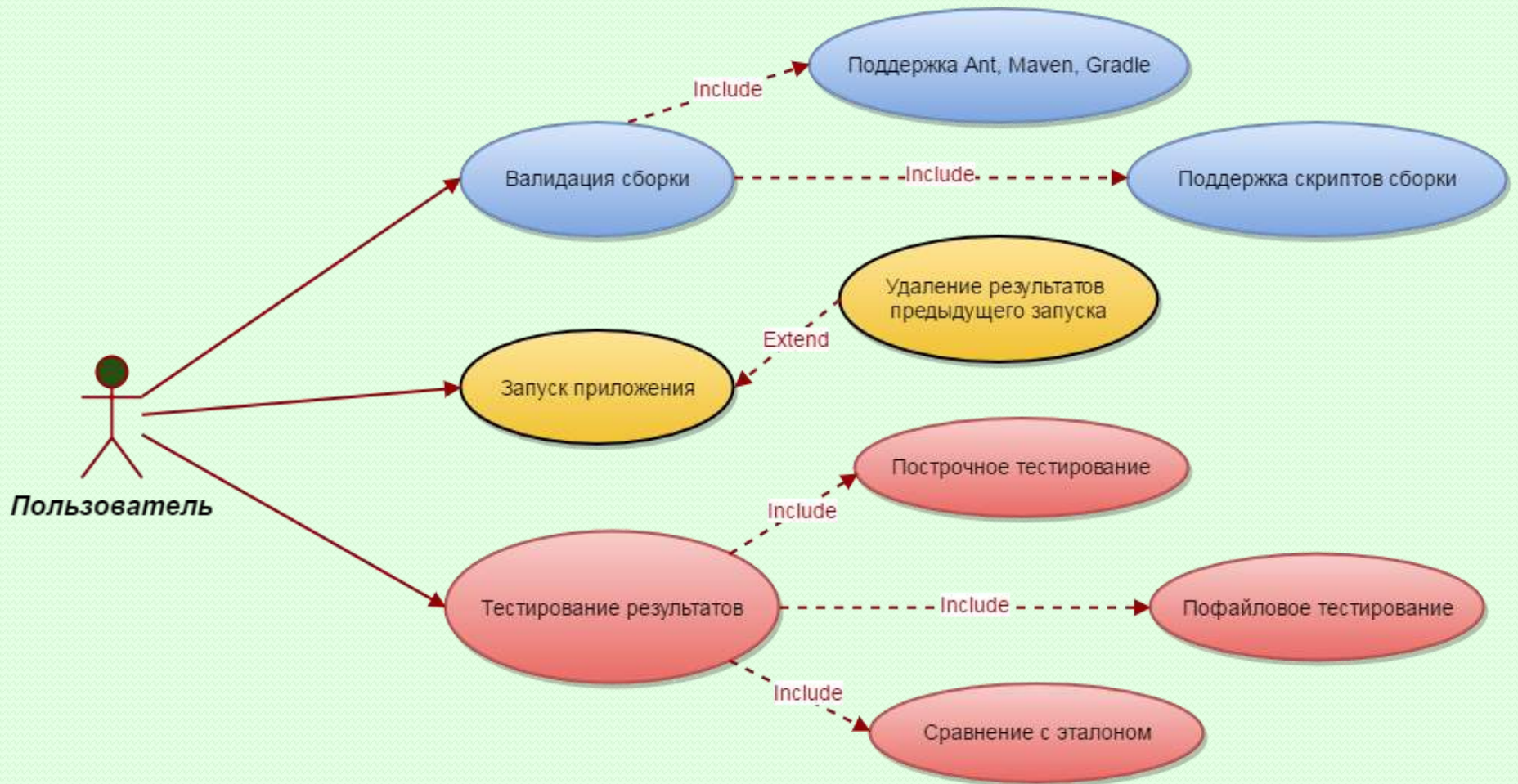
популярность библиотек на GitHub

# 5. NJudge

- Написана на Python 2.7 с использованием библиотеки PyDoop.

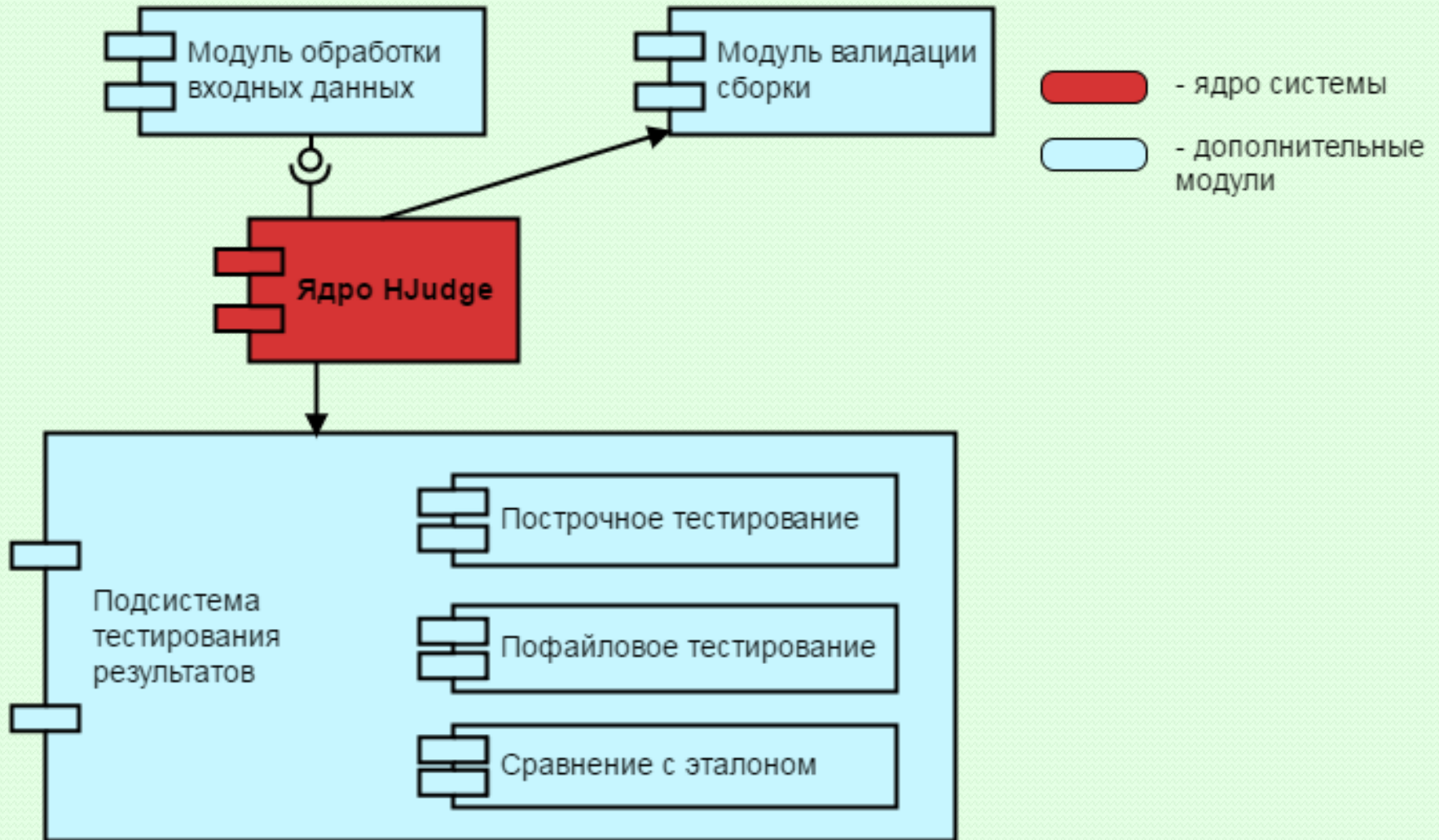


# Схема работы



# Архитектура

- Схема компонентов NJudge



# Архитектура

- Назначение модулей

| Модуль                             | Назначение  |
|------------------------------------|---|
| Модуль обработки<br>входных данных | Получает данные из аргументов CLI и парсит  |
| Модуль валидации<br>сборки         | Контролирует сборку тестируемого приложения   |
| Модуль проверки<br>результатов     | Проверяет выход приложения  |
| Ядро системы                       | <ul style="list-style-type: none"><li>• Запускает тестируемое приложение</li><li>• Передаёт управление от модуля к модулю</li><li>• Ведёт учёт проверенных приложений</li></ul> |

- Процесс Code review на данный момент не автоматизируется

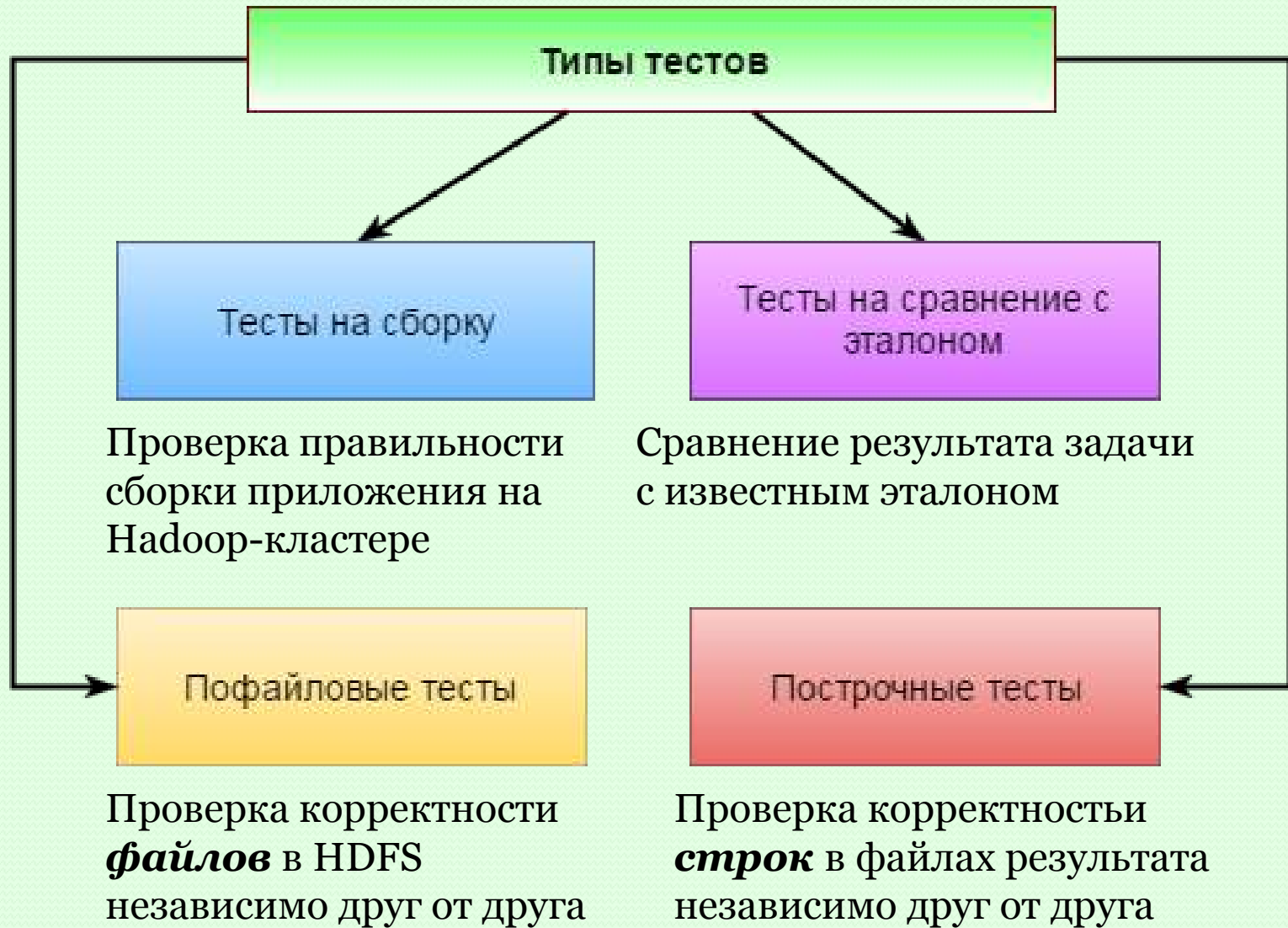
# Архитектура

- Назначение модулей

| Модуль                             | Назначение  |
|------------------------------------|---|
| Модуль обработки<br>входных данных | Получает данные из аргументов CLI и парсит  |
| Модуль валидации<br>сборки         | Контролирует сборку тестируемого приложения<br><i>(можно встраивать свои тесты)</i>   |
| Модуль проверки<br>результатов     | Проверяет выход приложения<br><i>(можно встраивать свои тесты)</i>  |
| Ядро системы                       | <ul style="list-style-type: none"> <li>• Запускает тестируемое приложение</li> <li>• Передаёт управление от модуля к модулю</li> <li>• Ведёт учёт проверенных приложений</li> </ul> |

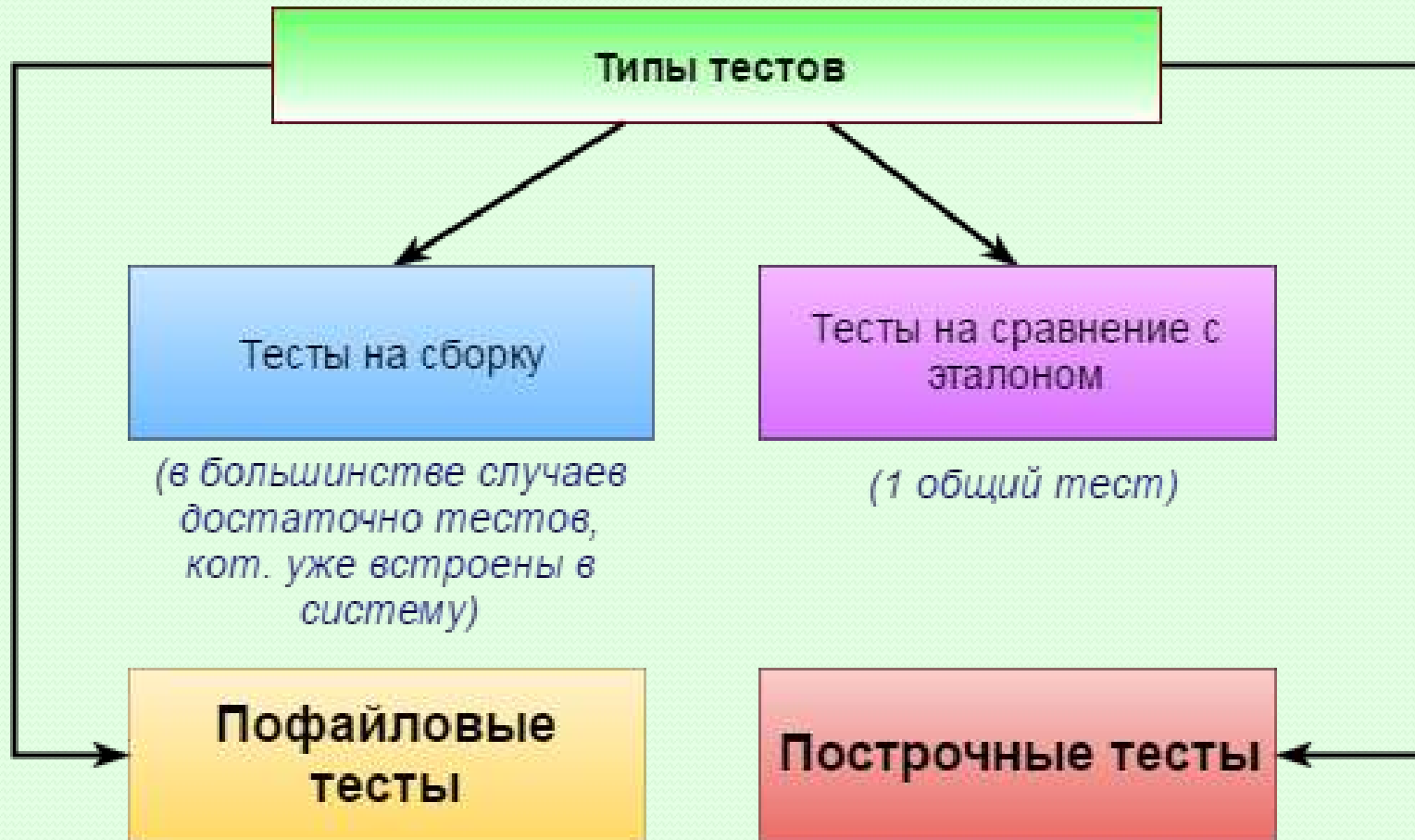
- Процесс Code review на данный момент не автоматизируется
- Пользовательские тесты – это обычные Python-функции, встраивающиеся в специальный класс

# База тестов





# База тестов

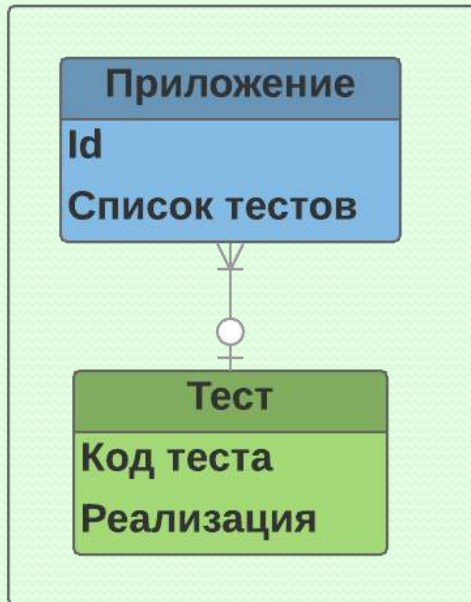


- В основном, нужно встраивать только построчные и пофайловые тесты

# База тестов

- Пофайловые тесты

База пофайловых тестов



*Текстовый файл*

*Python-класс*

|   |     |                  |
|---|-----|------------------|
| 1 | 101 | mrb2,mrb3_1,mrb8 |
| 2 | 102 | mrb2,mrb3_1,mrb8 |
| 3 | 103 | mrb2,mrb3_1,mrb8 |
| 4 | 104 | mrb2,mrb8        |
| 5 | 105 | mrb2,mrb3_1,mrb8 |
| 6 | 106 | mrb2,mrb8        |
| 7 | 107 | mrb2,mrb8        |
| 8 | 108 | mrb7             |
| 9 | 109 | mrb2,mrb3_1,mrb8 |

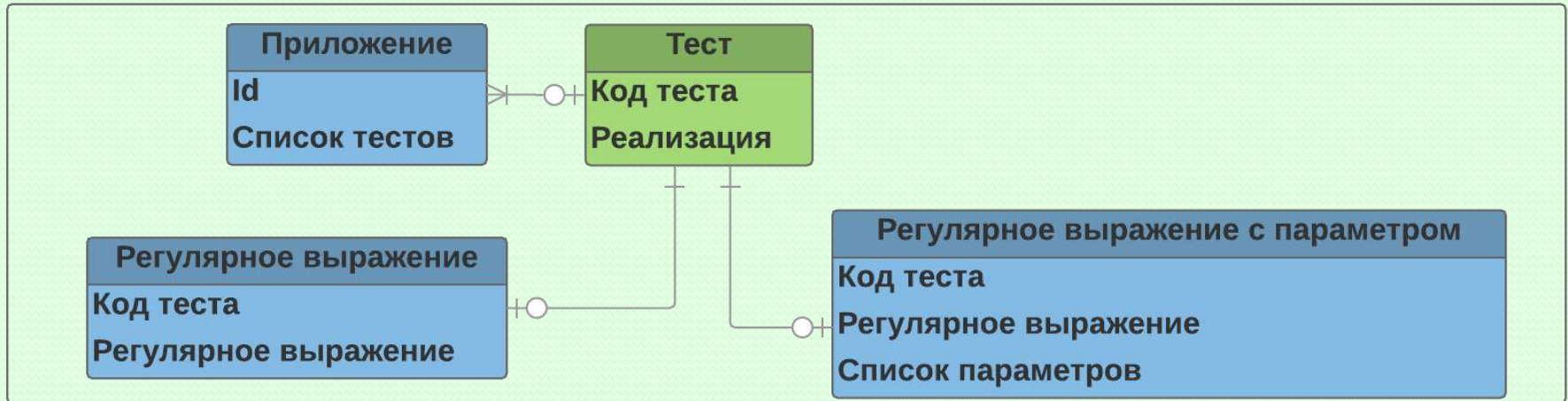
```

99     def mra1_2(self):
100         return self._mra1(5000, 50)
101
102     def mrb12(self):
103         for num, line in enumerate(self.hdfs_file):
104             if num:
105                 return True
106         return False
  
```

# База тестов

- Построчные тесты

База построчных тестов



```

1  mrb1_1    ^\S+\s\d+$
2  mrb1_2    ^1(\.0+)?$|^0\.(?!0+$)\d+$
3  mrb1_3    ^[\S\s]+(?:\s\d+){3}\.?\d*$
4  mrb1_4    ^\S+(?:\s\d+\.?\d*){2}$
5  mrb1_5    ^(?:\S+\s){2}\d+$
6  mrb1_6    ^(?:\S{18}[,\s]){,4}\S{18}$
7  mrb1_7    ^\S+\s\d+\.?\d*(?:e-)?\d{,2}$
8  mrb1_8    ^[\S\s]+(?:\s\d+){3}\.?\d*$

```

*Текстовый файл*  
 *Python-класс*

# Отчётность по тестированию

1. Логи тестирования в реальном времени
2. Отчёт по окончании тестирования приложения

# Отчётность по тестированию

## 1. Логи тестирования в реальном времени

```

instructor@bds02:~/tasks/s19441_task107_v2$ hjudge --run-cmd "bash run.sh" --out-dir /user/s19441/gore_ota
_uma_result/ --build maven --login s19441 --num 107
===== Testing task: 1 (HW: 107): STARTING =====
main.py ##### INFO      [2017-01-24 20:36:00,850] You probably don't have permissions to write into out-dir.
main.py ##### INFO      [2017-01-24 20:36:00,850] Out-dir will be changed to /user/instructor/s19441_task107
main.py ##### WARNING   [2017-01-24 20:36:00,865] Test MRB2 is specified, but not realized.
main.py ##### WARNING   [2017-01-24 20:36:00,865] Test MRB8 is specified, but not realized.
main.py ##### WARNING   [2017-01-24 20:36:00,865] There are no realized full tests for task 107.
main.py ##### WARNING   [2017-01-24 20:36:00,866] There are no specified systems tests for task 107.
main.py ##### INFO      [2017-01-24 20:37:52,486] Testing task: 1 (HW: 107): test MRB14 failed on line "5-я княжна 0 4"!
main.py ##### INFO      [2017-01-24 20:37:52,486] Testing task: 1 (HW: 107): test MRB6 failed on line "5-я княжна 0 4"!
main.py ##### INFO      [2017-01-24 20:37:52,614] Testing task: 1 (HW: 107): test MRB13 passed!
main.py ##### INFO      [2017-01-24 20:37:52,615] Testing task: 1 (HW: 107): test MRB4_1 passed!

```

Нет прав на запись в выходную директорию, система её заменяет

Тесты, кот. есть в списке, но нет в класс с Python-функциями

Список пройденных тестов

# Отчётность по тестированию

## 1. Логи тестирования в реальном времени

```

instructor@bds02:~/tasks/s19441_task107_v2$ hjudge --run-cmd "bash run.sh" --out-dir /user/s19441/gore_ota
_uma_result/ --build maven --login s19441 --num 107
===== Testing task: 1 (HW: 107): STARTING =====
main.py ##### INFO      [2017-01-24 20:36:00,850] You probably don't have permissions to write into out-dir.
main.py ##### INFO      [2017-01-24 20:36:00,850] Out_dir will be changed to /user/instructor/s19441_task107
main.py ##### WARNING   [2017-01-24 20:36:00,865] Test MRB2 is specified, but not realized.
main.py ##### WARNING   [2017-01-24 20:36:00,865] Test MRB8 is specified, but not realized.
main.py ##### WARNING   [2017-01-24 20:36:00,865] There are no realized full tests for task 107.
main.py ##### WARNING   [2017-01-24 20:36:00,866] There are no specified systems tests for task 107.
main.py ##### INFO      [2017-01-24 20:37:52,486] Testing task: 1 (HW: 107): test MRB14 failed on line "5-я княжна 0 4"!
main.py ##### INFO      [2017-01-24 20:37:52,486] Testing task: 1 (HW: 107): test MRB6 failed on line "5-я княжна 0 4"!
main.py ##### INFO      [2017-01-24 20:37:52,614] Testing task: 1 (HW: 107): test MRB13 passed!
main.py ##### INFO      [2017-01-24 20:37:52,615] Testing task: 1 (HW: 107): test MRB4_1 passed!

```

Нет прав на запись в выходную директорию, система её заменяет

Тесты, кот. есть в списке, но нет в класс с Python-функциями

Список пройденных тестов

- Для экономии времени можно настроить систему так, что тестирование будет останавливаться при первом не пройденном тесте.

# Отчётность по тестированию

## 2. Отчёт по окончании тестирования приложения

```
==== Testing task: 1 (HW: 107): SUMMARY ====  
Tests passed: mrb13 mrb4_1  
Tests failed: mrb14 mrb6  
=====
```

- Можно запустить на проверку сразу несколько приложений
  - ❖ Когда тестирование завершится, отчёт будет выдан по каждому

# Установка

## Особенности установки

- Имеет мало дополнительных зависимостей:
  1. Hadoop, Hive, Spark, HBase
  2. Python 2,7
  3. PyDoop 1.2
  4. PyPi Regex



# Установка

## Особенности установки

- Имеет мало дополнительных зависимостей:
    1. Hadoop, Hive, Spark, HBase
    2. Python 2,7
    3. PyDoop 1.2
    4. PyPi Regex //опционально
- } Обычно установлены на Hadoop-кластере

# Установка

## Особенности установки

- Имеет мало дополнительных зависимостей:
  1. Hadoop, Hive, Spark, HBase
  2. Python 2,7
  3. PyDoop 1.2
  4. PyPi Regex //опционально

} Обычно установлены на Hadoop-кластере
- Установка через Pip
  - `pip install git+http://gitlab.vdi.mipt.ru/VeLKerr/hjudge.git`
- Проект пока не выложен в открытый доступ => нужна регистрация на [gitlab.vdi.mipt.ru](http://gitlab.vdi.mipt.ru).

# 6. Апробация NJudge

Система использовалась на курсах:

- Хранение и обработка больших объёмов данных, весна 2016 г.
- Многопроцессорные вычислительные системы, осень 2016 г.

# 6. Апробация HJudge

❖ Кол-во ошибок, обнаруженных на каждом этапе

| Тема ДЗ          | Валидация сборки | Проверка запуска | Проверка результатов |         |        |
|------------------|------------------|------------------|----------------------|---------|--------|
|                  |                  |                  | Построчн.            | Пофайл. | Полное |
| Hadoop Java API  | 3                | 11               | 26                   | 31      | 2      |
| Hadoop Streaming | 1                | 5                | 21                   | 28      | 0      |
| Joins            | 2                | 7                | 25                   | 13      | 1      |
| Spark            | 4                | 32               | 211                  | 114     | 102    |
| <b>Всего</b>     | 10               | 55               | 273                  | 186     | 105    |
| <b>%</b>         | 1,6              | 8,7              | 43,4                 | 29,6    | 16,7   |

# 6. Апробация HJudge

❖ Кол-во ошибок, обнаруженных на каждом этапе

| Тема ДЗ          | Валидация сборки | Проверка запуска | Проверка результатов |         |        |
|------------------|------------------|------------------|----------------------|---------|--------|
|                  |                  |                  | Построчн.            | Пофайл. | Полное |
| Hadoop Java API  | 3                | 11               | 26                   | 31      | 2      |
| Hadoop Streaming | 1                | 5                | 21                   | 28      | 0      |
| Joins            | 2                | 7                | 25                   | 13      | 1      |
| Spark            | 4                | 32               | 211                  | 114     | 102    |
| <b>Всего</b>     | 10               | 55               | 273                  | 186     | 105    |
| <b>%</b>         | 1,6              | 8,7              | 43,4                 | 29,6    | 16,7   |

❖ Больше всего ошибок обнаруживается на проверке результатов

❖ На шаге проверки результатов больше всего ошибок отлавливается построчными и пофайловыми тестами

# 6. Апробация HJudge

## ❖ Замеры времени тестирования

| Тема ДЗ             | Время ручной проверки, с |           |              | Время проверки с HJudge, с |          |             |
|---------------------|--------------------------|-----------|--------------|----------------------------|----------|-------------|
|                     | min                      | mean      | max          | min                        | mean     | max         |
| Hadoop<br>Java API  | 50                       | 100       | 150          | 8                          | 16,5     | 25          |
| Hadoop<br>Streaming | 50                       | 80        | 110          | 3                          | 5        | 7           |
| Joins               | 40                       | 120       | 200          | 5                          | 8,5      | 12          |
| Spark               | 2                        | 24        | 46           | 1                          | 2        | 3           |
| <b>В среднем</b>    | <b>35,5</b>              | <b>81</b> | <b>126,5</b> | <b>4,3</b>                 | <b>8</b> | <b>11,8</b> |

# 6. Апробация HJudge

## ❖ Замеры времени тестирования

| Тема ДЗ             | Время ручной проверки, с |           |              | Время проверки с HJudge, с |          |             |
|---------------------|--------------------------|-----------|--------------|----------------------------|----------|-------------|
|                     | min                      | mean      | max          | min                        | mean     | max         |
| Hadoop<br>Java API  | 50                       | 100       | 150          | 8                          | 16,5     | 25          |
| Hadoop<br>Streaming | 50                       | 80        | 110          | 3                          | 5        | 7           |
| Joins               | 40                       | 120       | 200          | 5                          | 8,5      | 12          |
| Spark               | 2                        | 24        | 46           | 1                          | 2        | 3           |
| <b>В среднем</b>    | <b>35,5</b>              | <b>81</b> | <b>126,5</b> | <b>4,3</b>                 | <b>8</b> | <b>11,8</b> |

## ❖ HJudge даёт выигрыш по времени **в 10 раз**

# 7. Планы по доработке

## 1. Интеграция с облачной платформой Everest

- ❖ позволит пользоваться не только преподавателям, но и студентам, т.к. можно ограничить доступ к тестам

Everest<sup>B</sup> Applications Jobs Resources Groups About Oleg

### Resources

Update + New resource

Filter by state  ANY  ONLINE  OFFLINE  My resources

| Name         | Type   | Total Slots | Free Slots | Max Tasks | Total Tasks | Running Tasks | Owner |
|--------------|--------|-------------|------------|-----------|-------------|---------------|-------|
| MIPT-Cluster | local  | 20          | 20         | 20        | 0           | 0             | Oleg  |
| NetBook-Acer | local  | 3           | 3          | 3         | 0           | 0             | Oleg  |
| test         | docker | 12          | 12         | 12        | 0           | 0             | sol   |



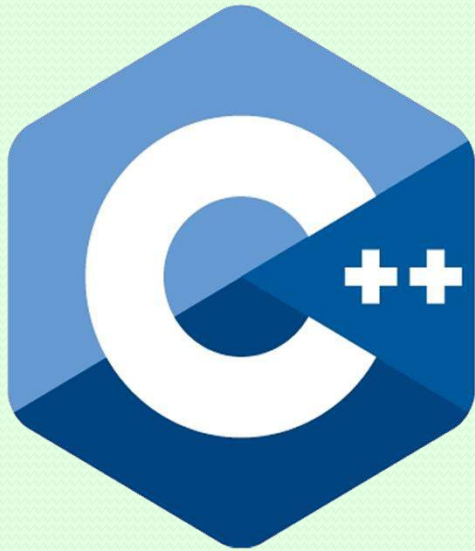
# 7. Планы по доработке

## 2. Интеграция с платформой Piazza

- ❖ студенты уже пользуются платформой
- ❖ при появлении нового сообщения с заданием, система будет получать из него данные автоматически
- ❖ существует неофициальное Python API для взаимодействия с Piazza
  - <http://hfaran.me/reverse-engineering-piazzas-api/>

# 7. Планы по доработке

3. Возможность писать тестирующие функции не на Python



# 8. ВЫВОД

1. Система NJudge реализована и продолжает развиваться
2. Система используется на курсах кафедры АТП и внедрена в учебный процесс



# 8. ВЫВОД

3. Система зарегистрирована в реестре программного обеспечения РФ.



Спасибо за внимание!

Вопросы?