

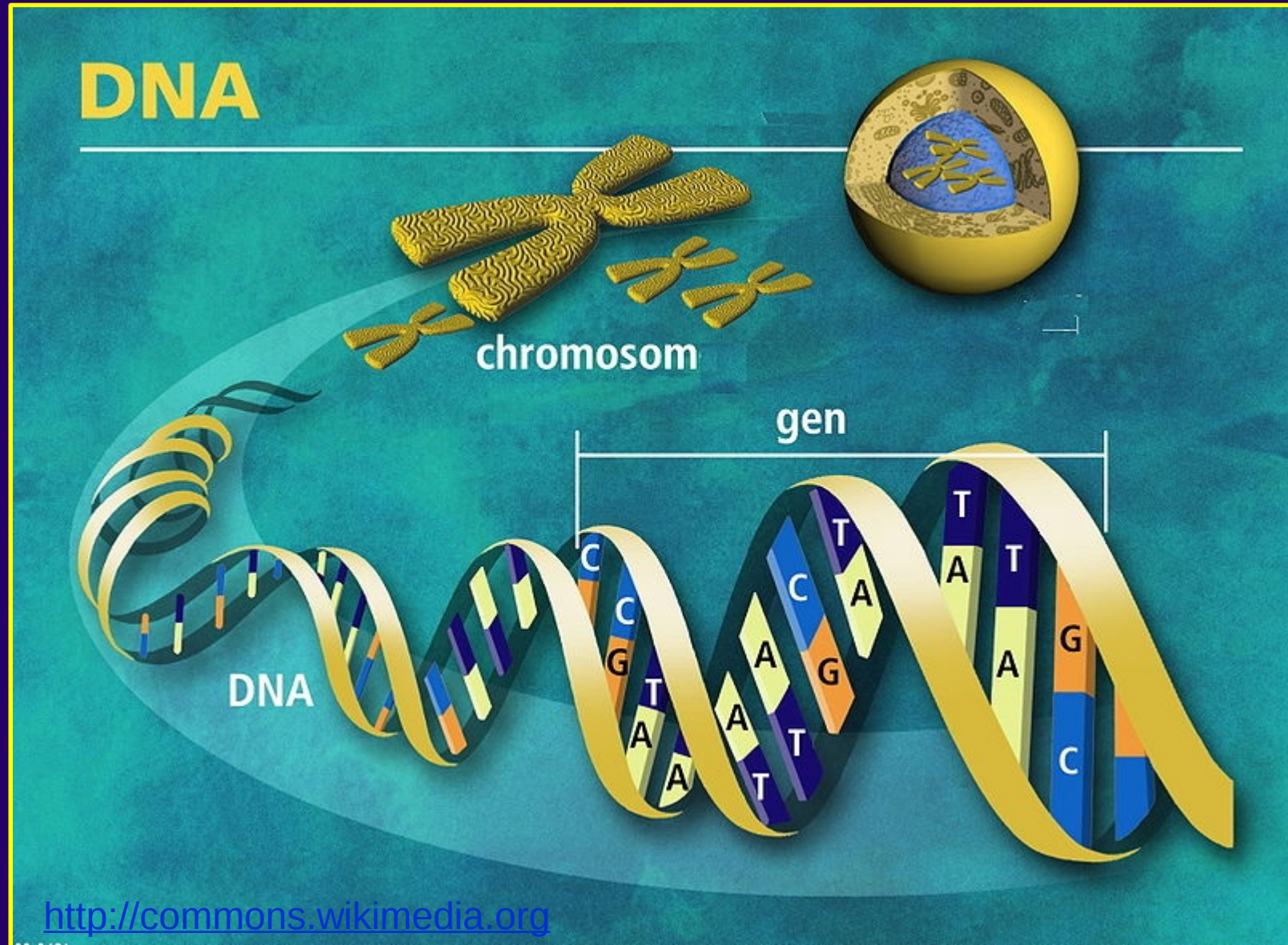
Ubuntu 13.10 Install Fest
МИЭМ НИУ ВШЭ 2.11.2013

Linux для обработки геномов

Сергей Науменко
Ренат Ариффулов
Нина Попова

Лаборатория эволюционной геномики ФББ МГУ
Институт проблем передачи информации РАН
РХТУ им. Д.И. Менделеева

Секвенирование



Секвенирование – прочтение последовательности ДНК или РНК



Проект “геном человека”

- 3,5 миллиарда нуклеотидов
- 13 лет
- 3 миллиарда долларов

Цели проекта:

- Все гены
- Полный геном
- Базы данных
- Новые инструменты
- Передача технологий
- Разработка этических, правовых и социальных аспектов



Высокопроизводительное секвенирование

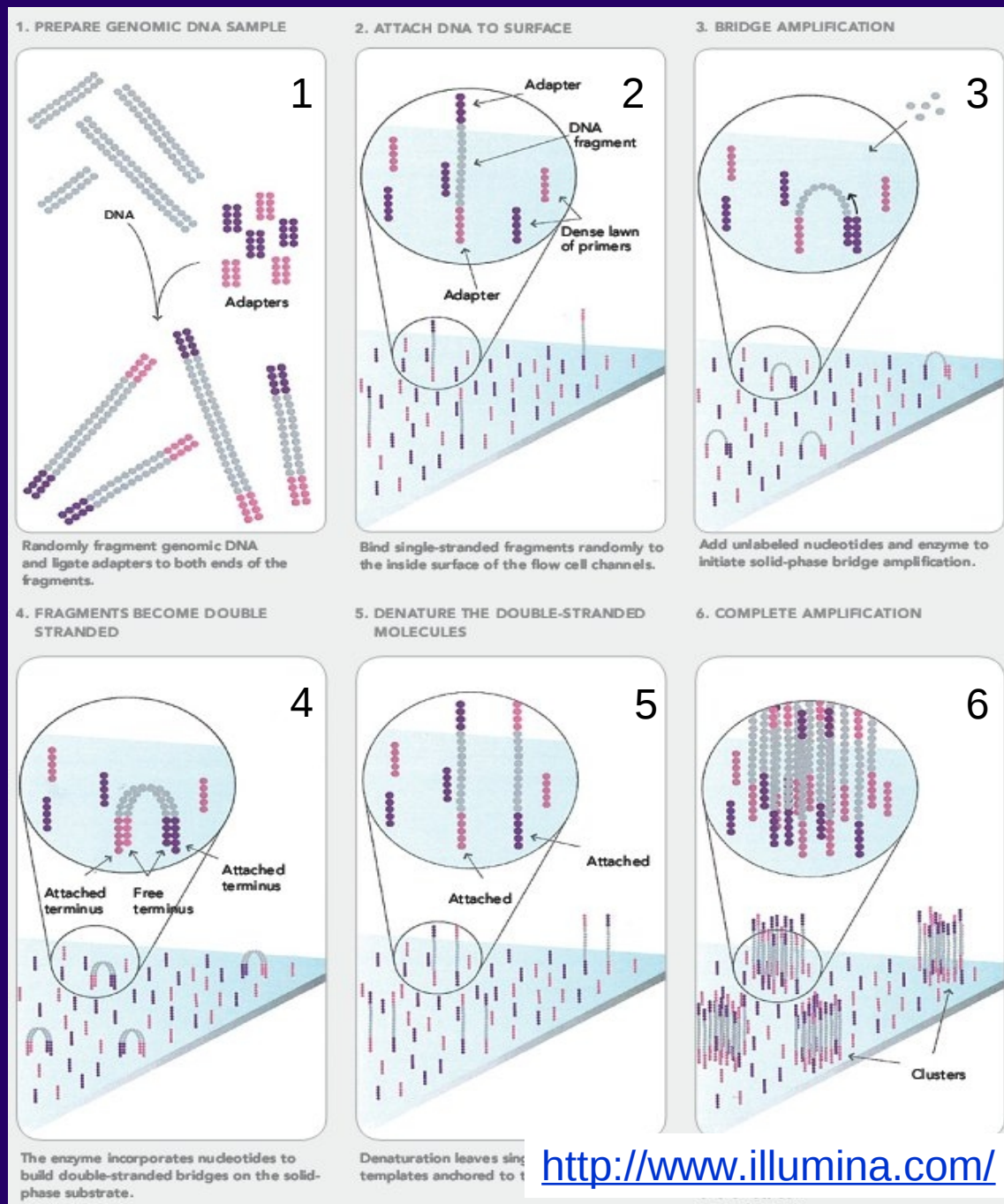
1) Фрагментация ДНК, пришивание адаптеров

2) Связывание фрагментов с подложкой

3-4) Добавление нуклеотидов, образование двойных цепочек

5) Денатурация (разрыв двойных цепочек, отрыв одного конца от подложки)

6) Амплификация - многократное копирование фрагментов, образование кластеров



Высокопроизводительное секвенирование

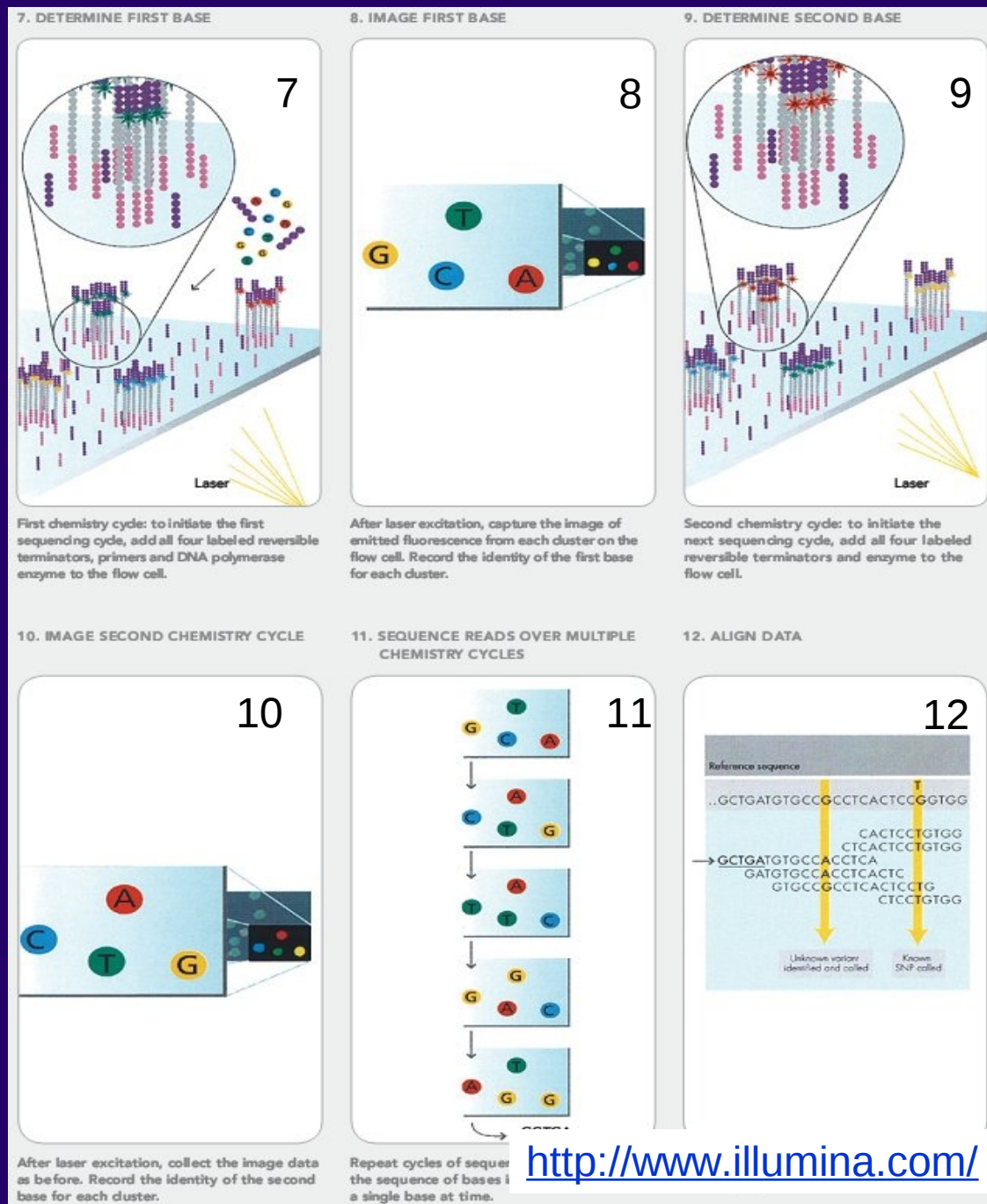
7-8) Добавление меченых нуклеотидов, определение первого нуклеотида при помощи лазера

9-11) Чтение остальных нуклеотидов

12) Анализ данных

Illumina HiSeq 2000:

- Две недели
- 600 миллиардов нуклеотидов
- 17 геномов человека с покрытием 10x
- 25 тысяч долларов США



<http://www.illumina.com/>

Высокопроизводительное секвенирование

в мире

в России

Пекинский институт геномики
(128 секвенаторов)

Лаборатория эволюционной
геномики ФББ МГУ



<http://www.genomics.cn/en/index>



<http://evolgenomics.fbb.msu.ru/>



Институт Сэнгера

<http://www.sanger.ac.uk/>

ЦНИИ эпидемиологии
Роспотребнадзора



<http://www.pcr.ru/>



<http://www.broadinstitute.org/>



Центр геномный исследований
НОЦ СФУ <http://genome.sfu-kras.ru/>



<http://www.jgi.doe.gov/>

ИОГЕН РАН, лаборатория
эволюционной геномики

<http://vigg.ru/>

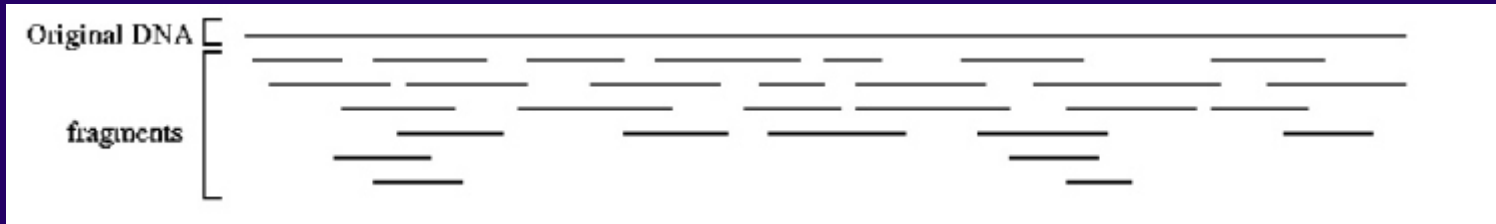


Основные задачи обработки данных

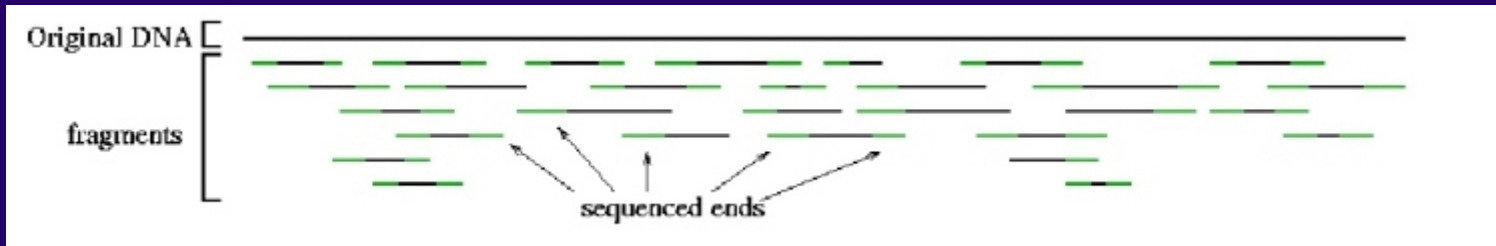
- Приём данных с секвенаторов
- Первичная обработка данных
- Сборка геномов и транскриптомов
- Картирование
- Аннотация
- Выравнивание
- Филогенетика
- Сравнительная геномика
- Популяционная генетика

Задача сборки генома

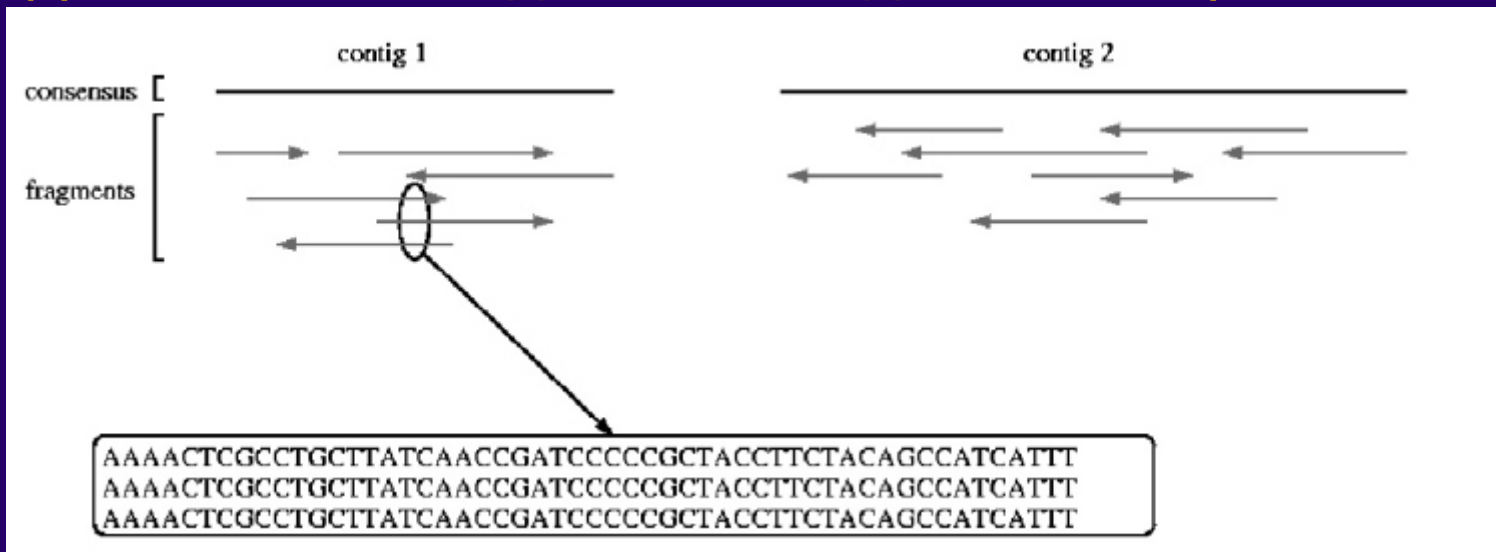
Исходная ДНК разрезается на фрагменты.



Фрагменты прочитываются с двух концов.

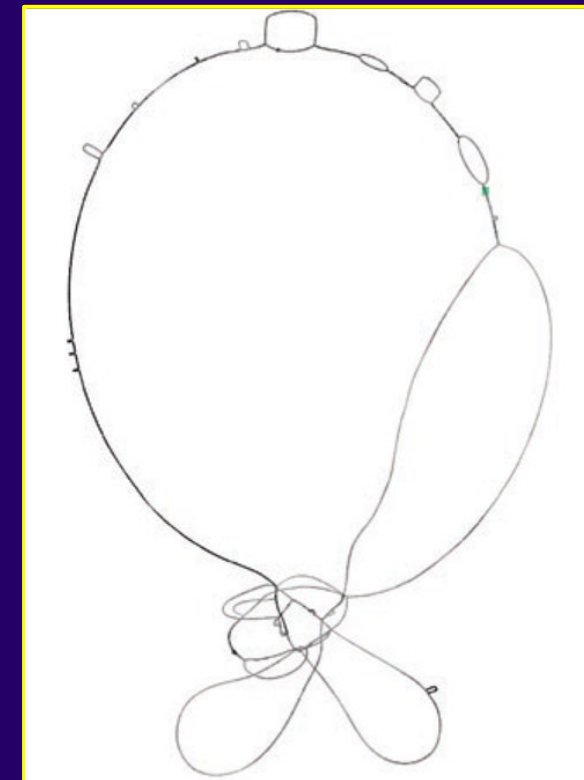
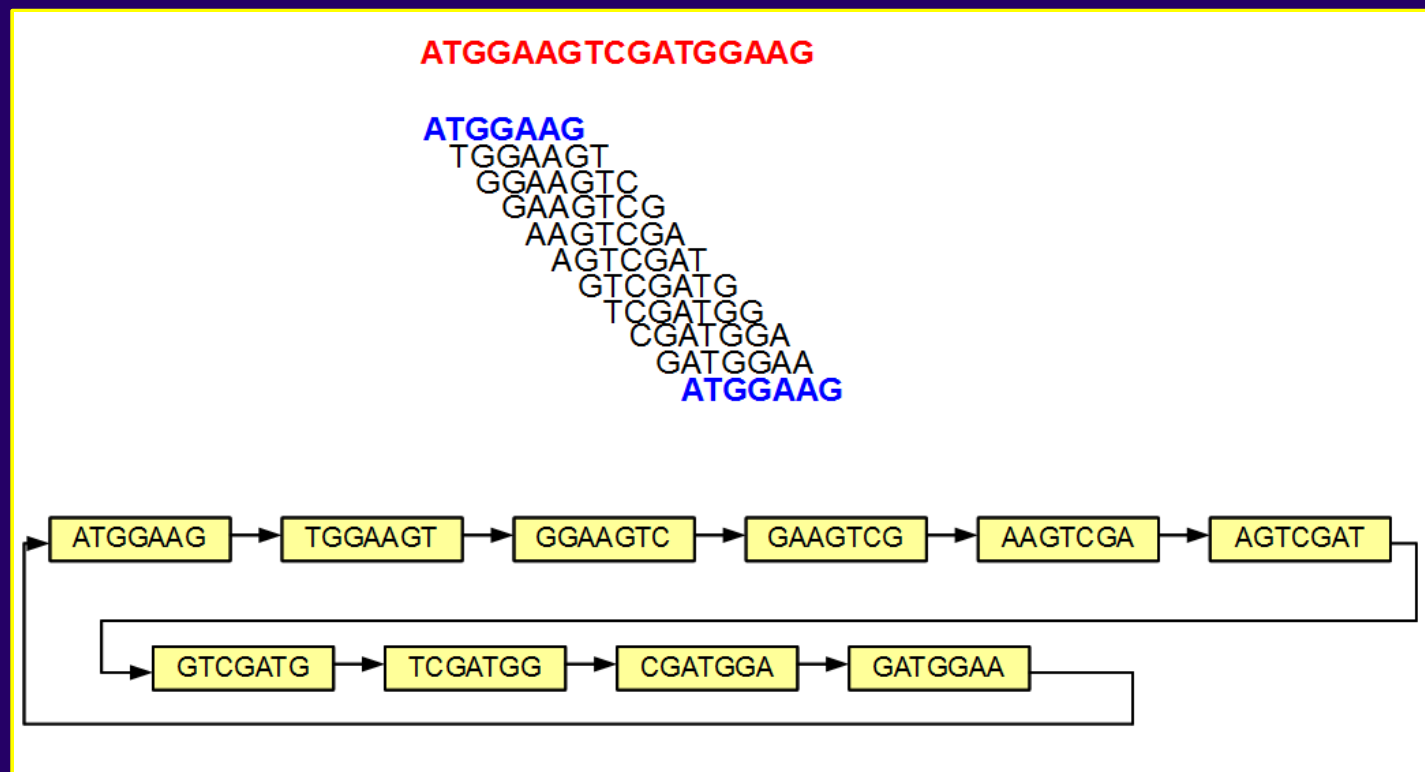


Схожие фрагменты объединяются в длинные строки.



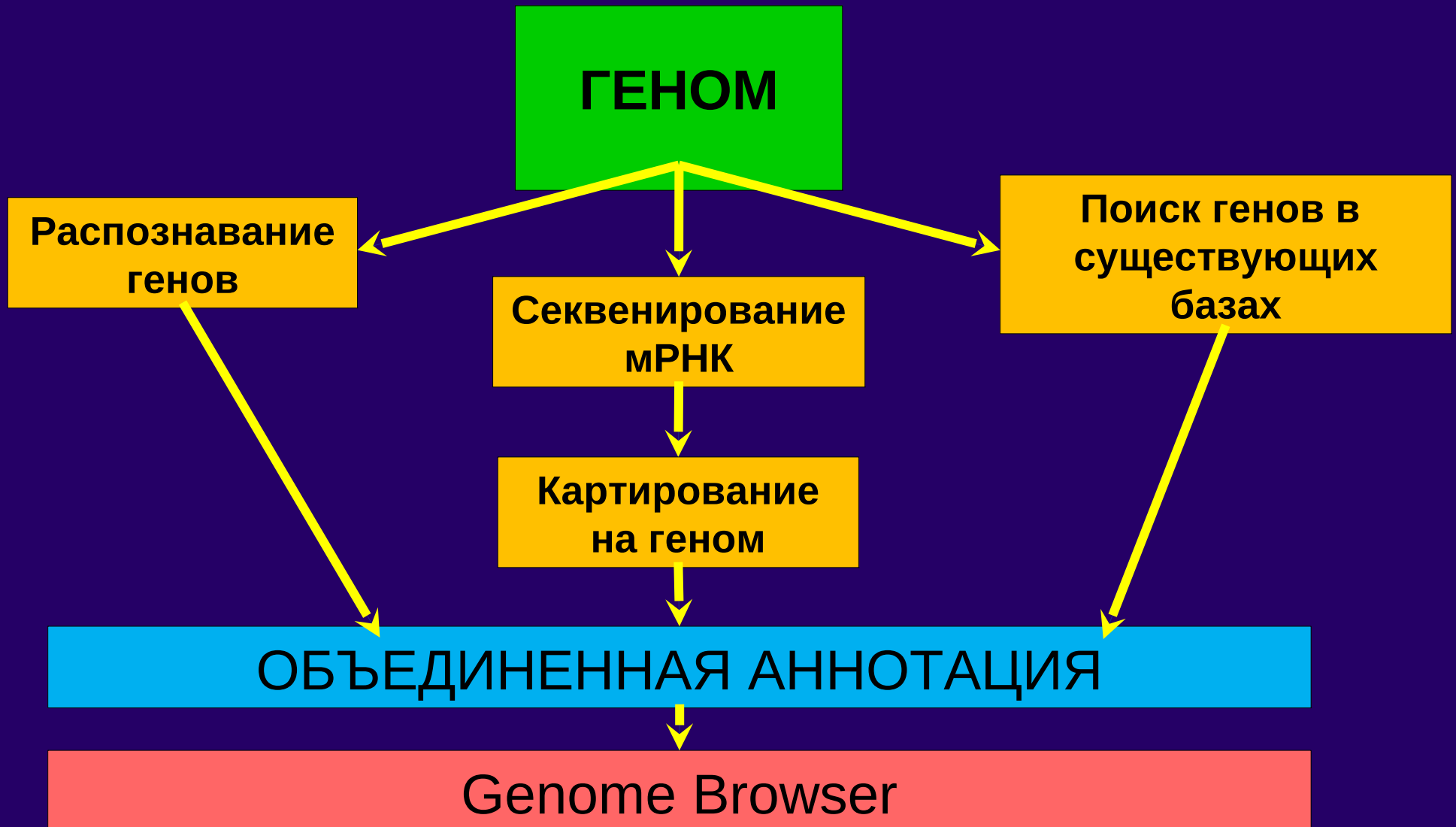
Алгоритм сборки генома

- Фрагменты разбиваются на слова (k-mer)
- Строится строковый граф
- Удаляются дублирующие пути
- На графе находится максимальный путь

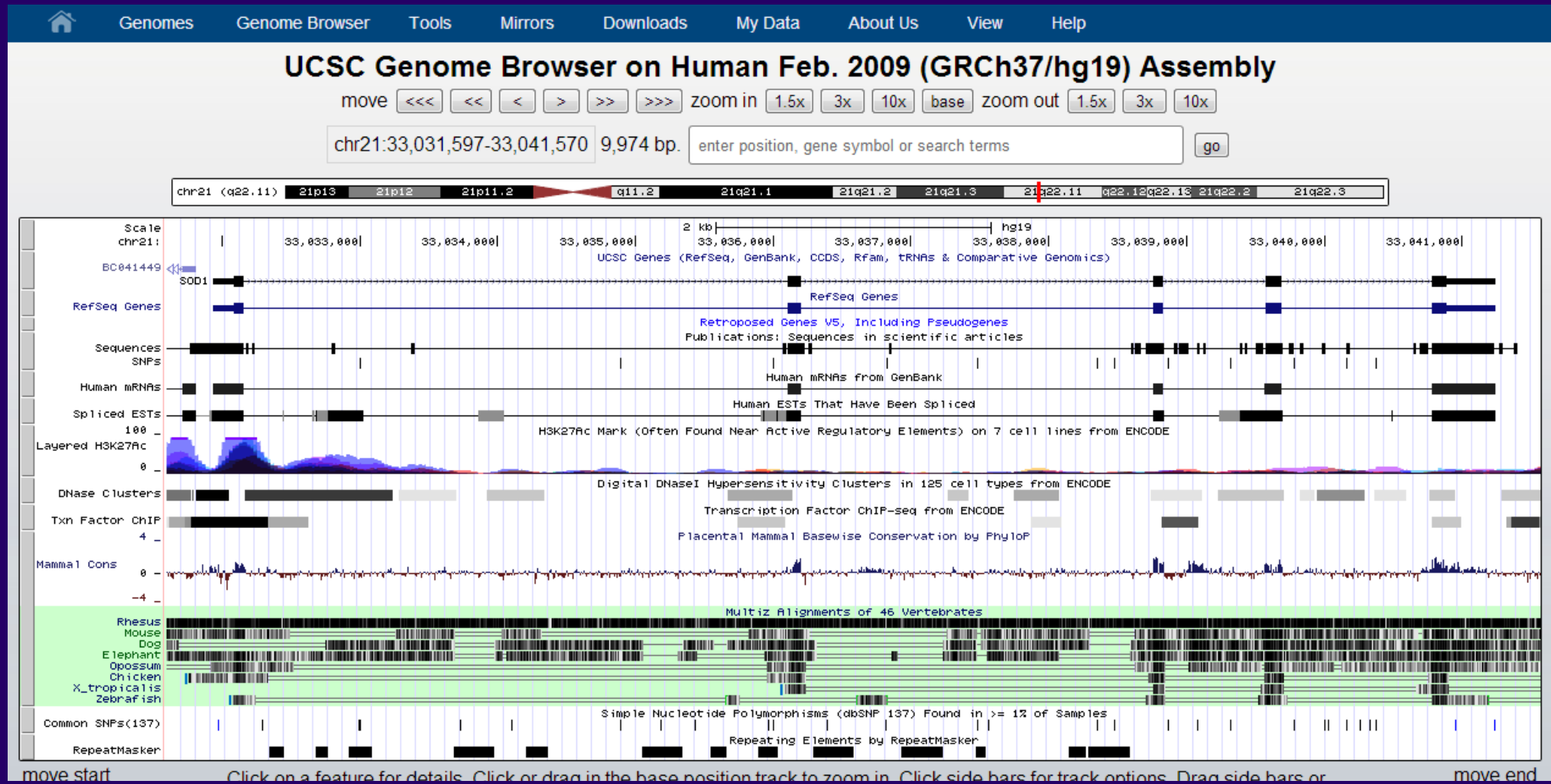


Аннотация генома – «разметка» генома

выделение последовательностей
белок-кодирующих генов и РНК



Собранный и аннотированный геном



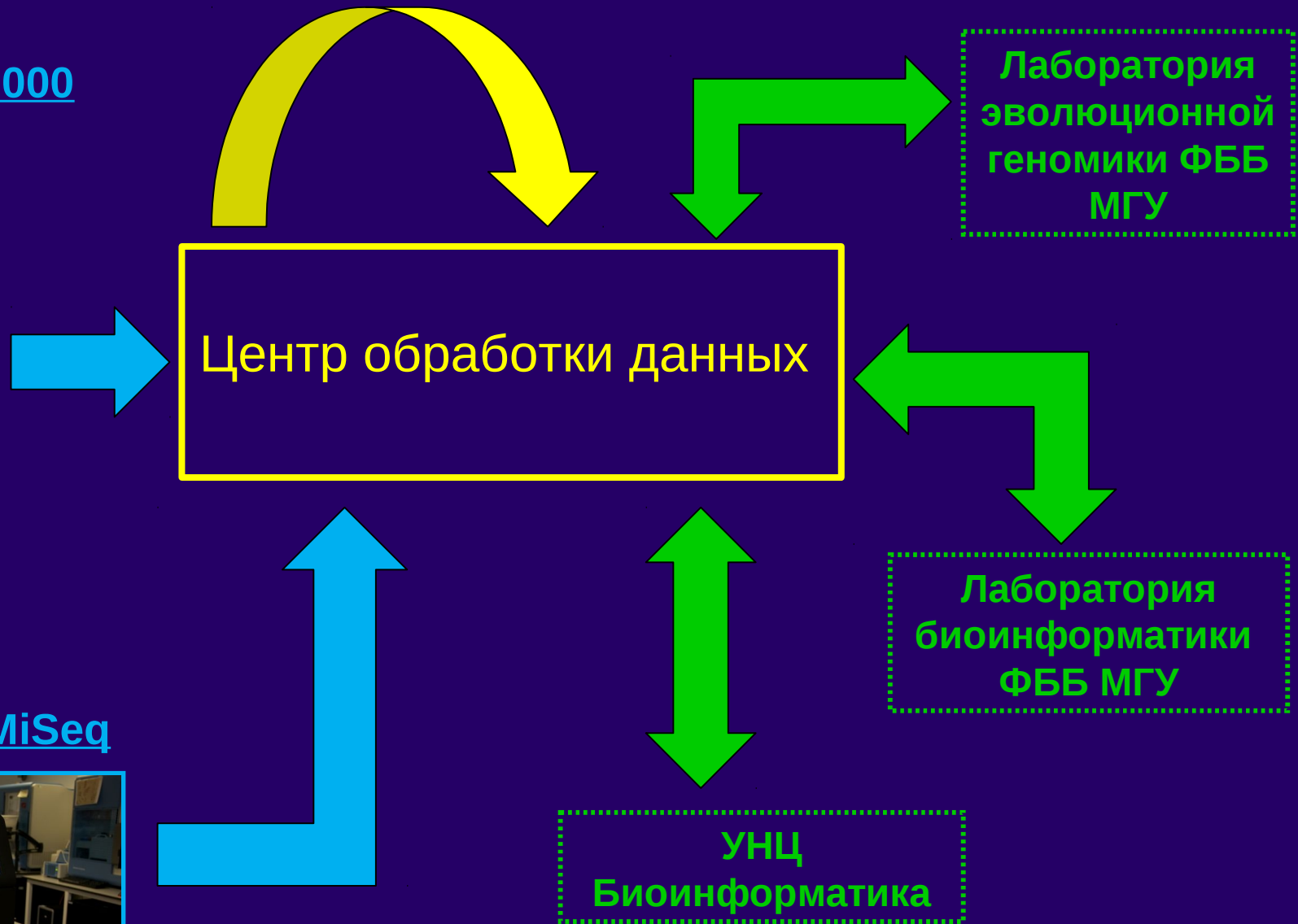
Потоки данных

Пользователи:
(более 50)

Illumina HiSeq2000



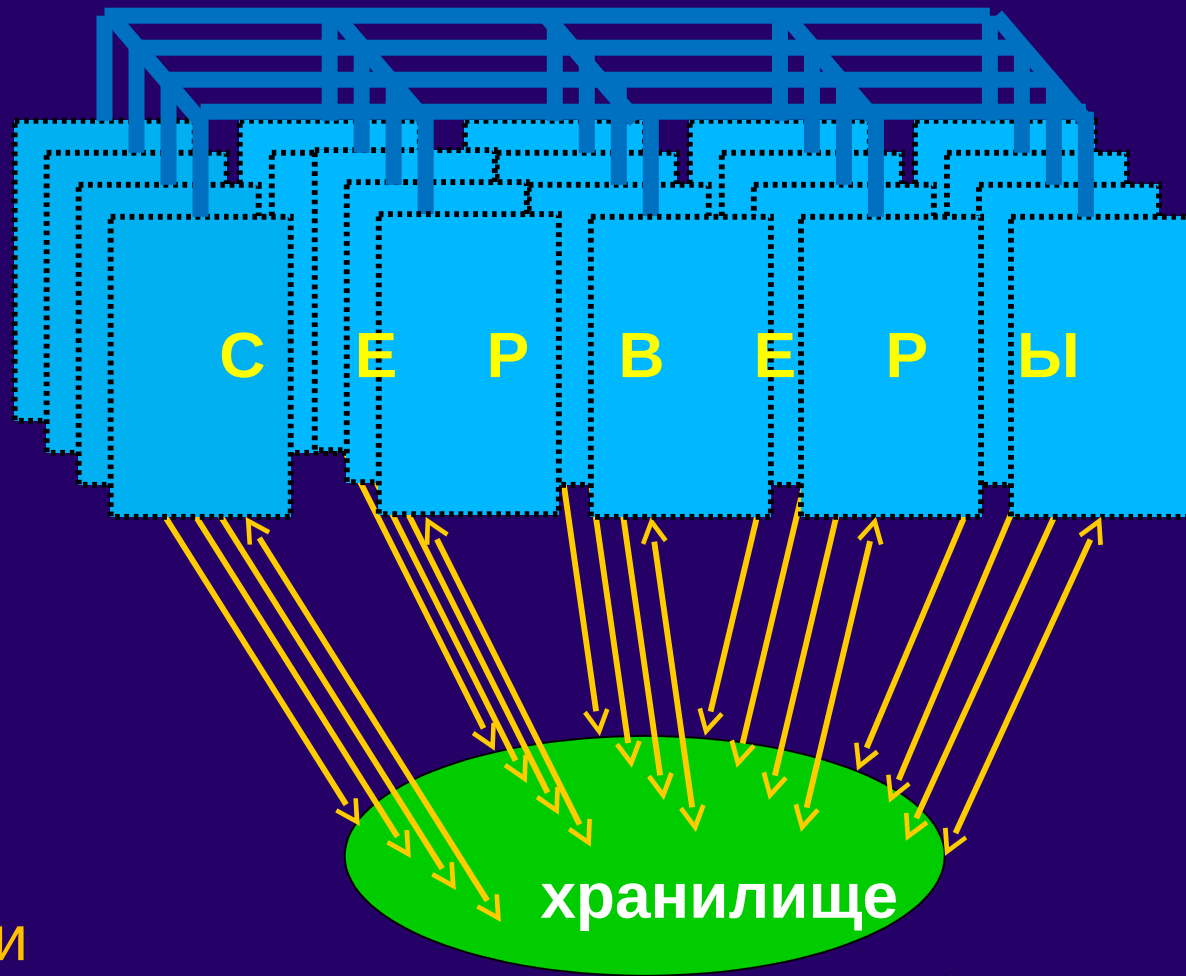
Illumina MiSeq



Вычислительный кластер

Предназначен для решения задач:

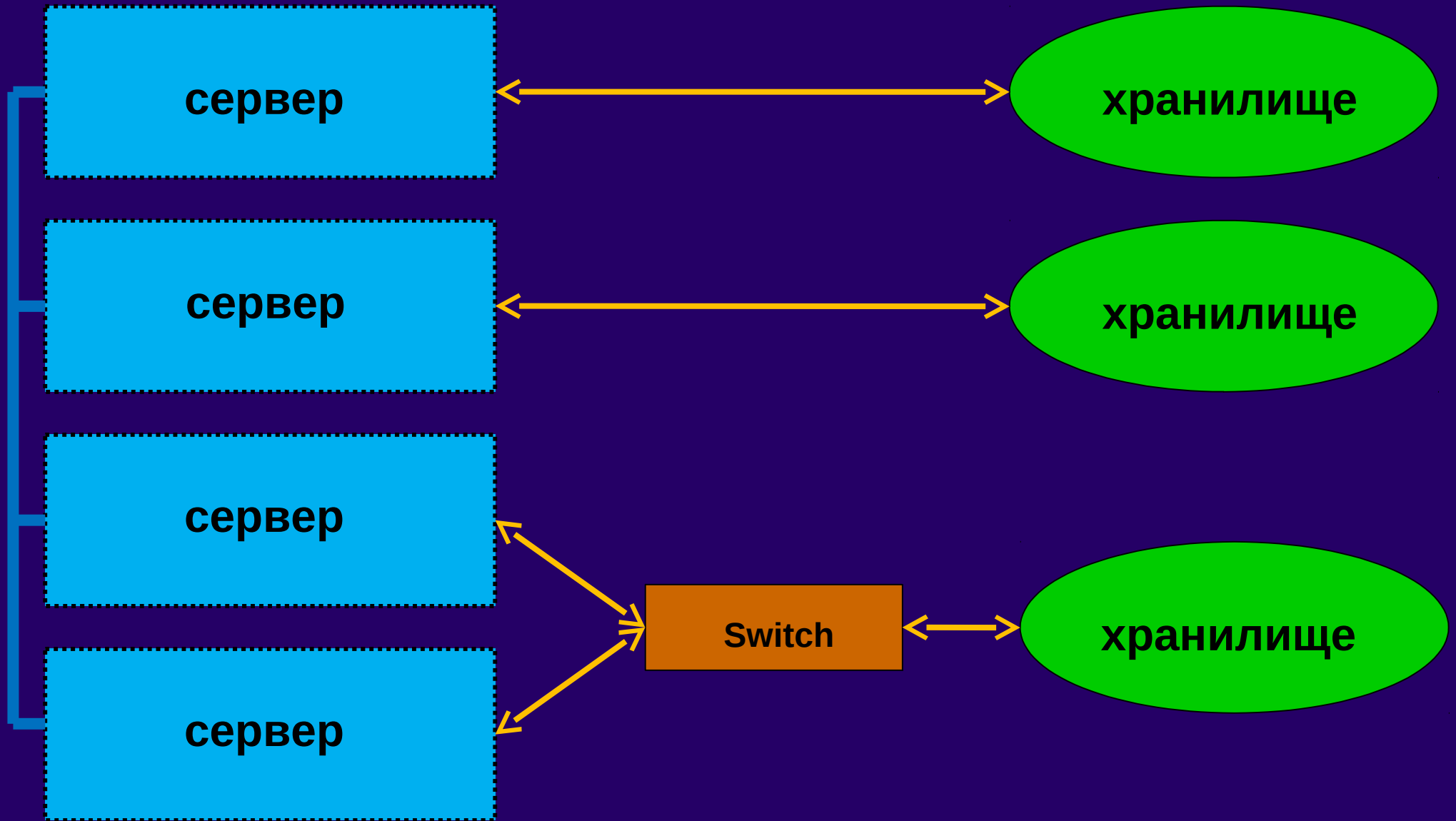
- Вычислительной гидродинамики
- Квантовой химии
- Физики высоких энергий



Особенности:

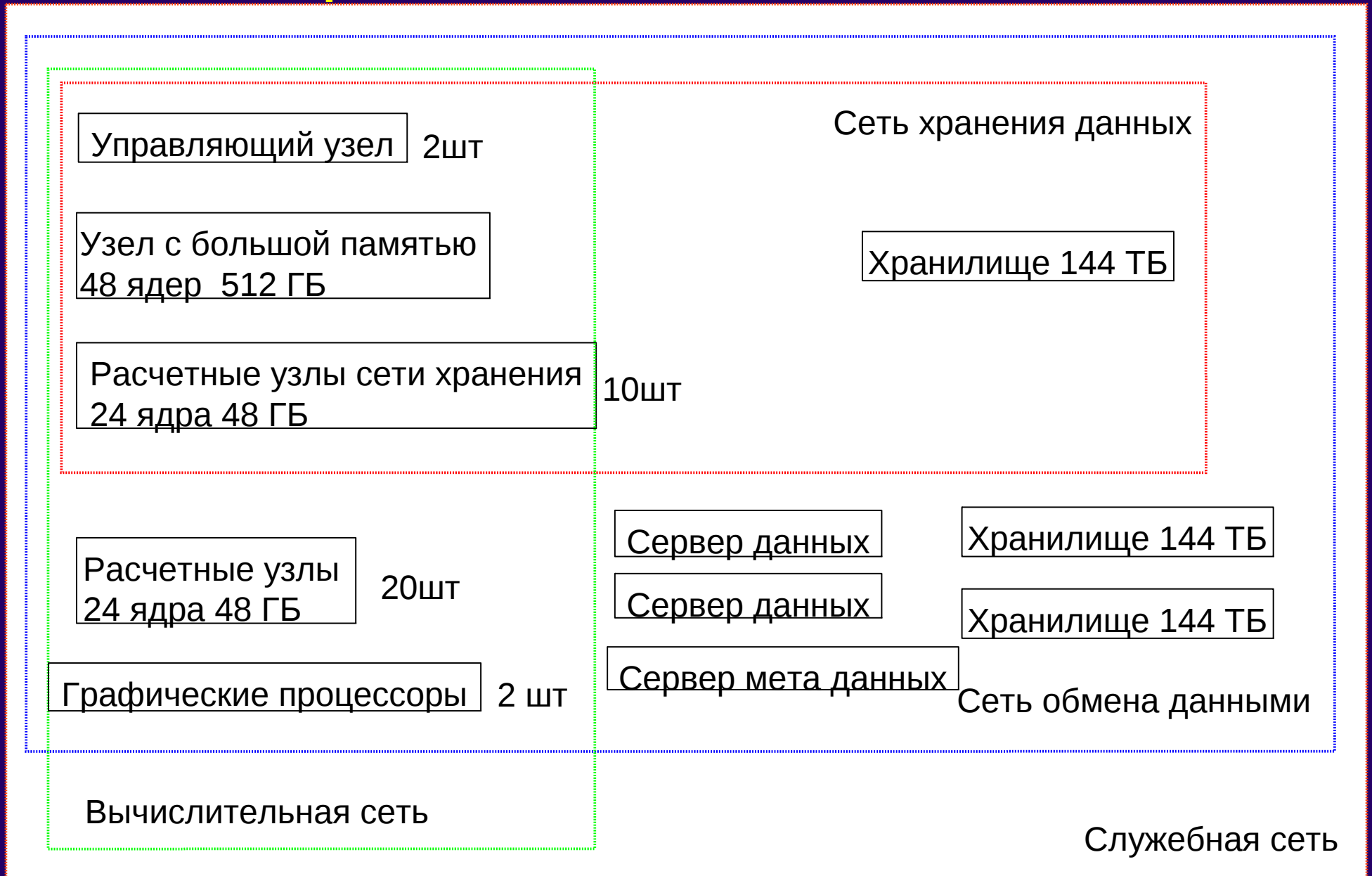
- Много ядер
- Мало оперативной памяти
- Небольшое хранилище

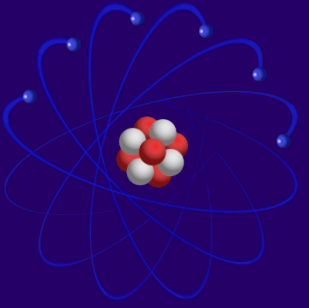
Центр обработки данных



Количество серверов сопоставимо с количеством хранилищ

Схема ЦОД лаборатории эволюционной геномики ФББ МГУ





Программное обеспечение

- Управление ресурсами (очередь задач) – torque
- Мониторинг – nagios
- Управление конфигурациями – puppet
- Файловые системы (XFS, lustre)
- OS – Enterprise Linux (Scientific Linux)
- Биоинформатические пакеты



Nagios[®]



lustre[®]



Биоинформатические пакеты



- Velvet
- Soapdenovo
- Platanus
- GATK
- MUMmer
- Clustal
- R
- Biopython
- BioPerl
- PHYLYP
- SHRiMP
- STAR
- Agalma
- Bambus
- Bamtools
- Blast
- Blat
- RAXML



- MCScanX
- Python
- Pal2nal
- AdapterRemoval
- MCScanX
- HaploMerger
- Paml
- Mrbayes
- geneid
- bowtie
- HaploMerger
- Megan

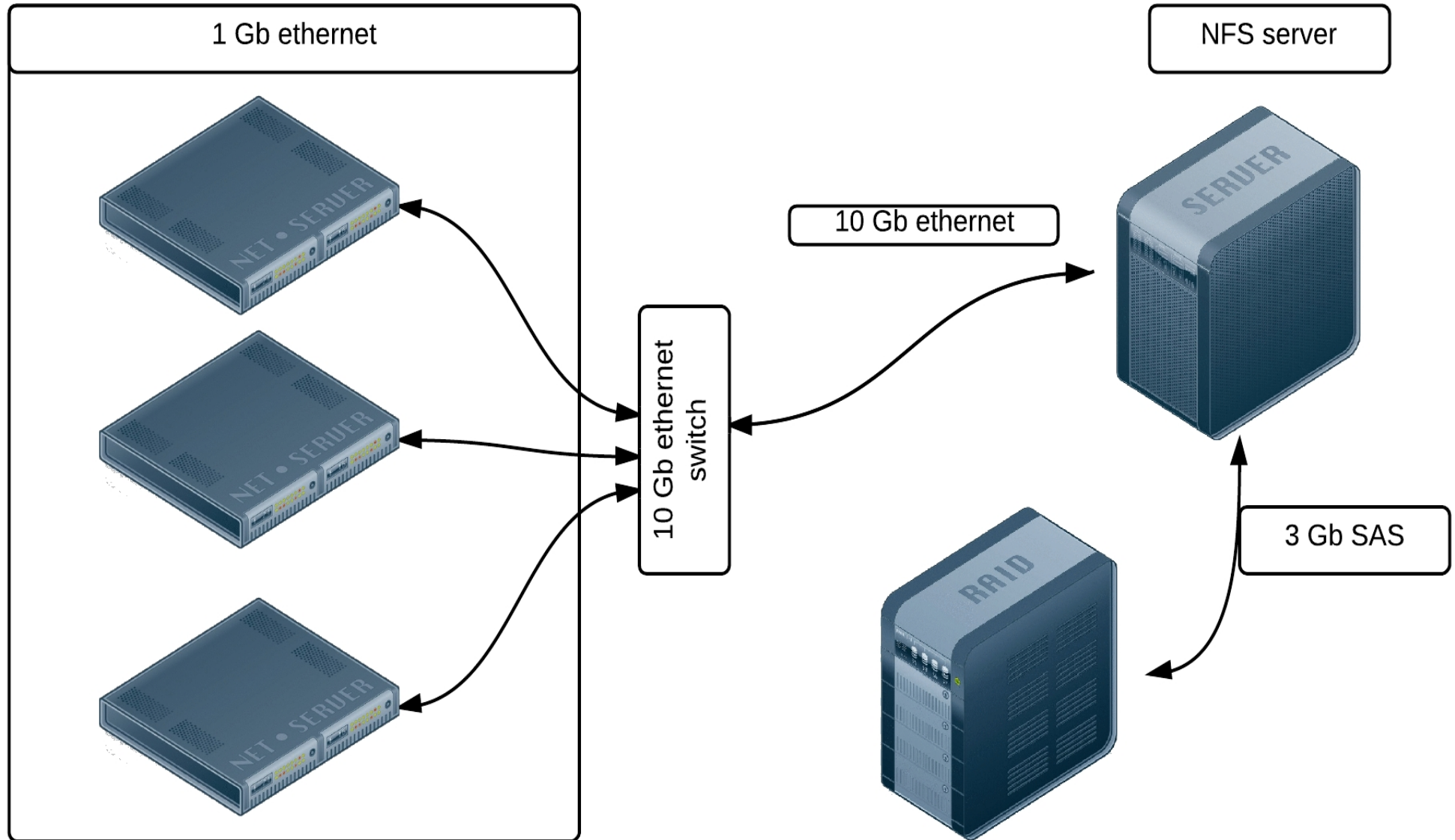
- Cuda
- Orthomcl
- Jellyfish
- Cegma
- Annovar
- Abyss
- AdapterRemoval
- Statistics-Descriptive
- beagle-lib



- libsequence
- BaseSpaceSHRE
- wise
- RepeatMasker
- И др

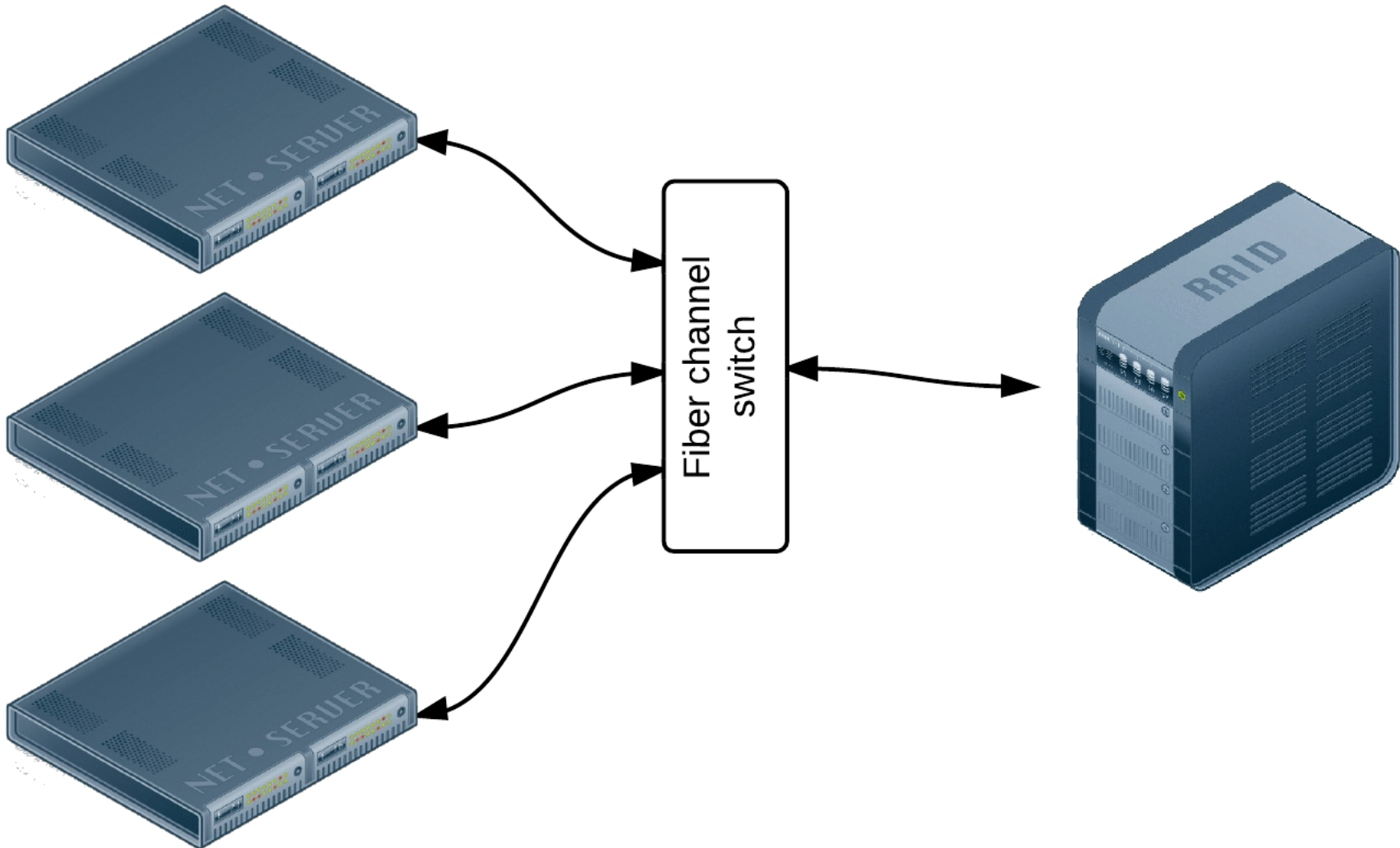


NFS Сервер



Сеть хранения данных

fiber channel network



Результаты тестов чтения/записи для файловых систем NFS, OCFS2, GFS2

dd (write), oflag=direct

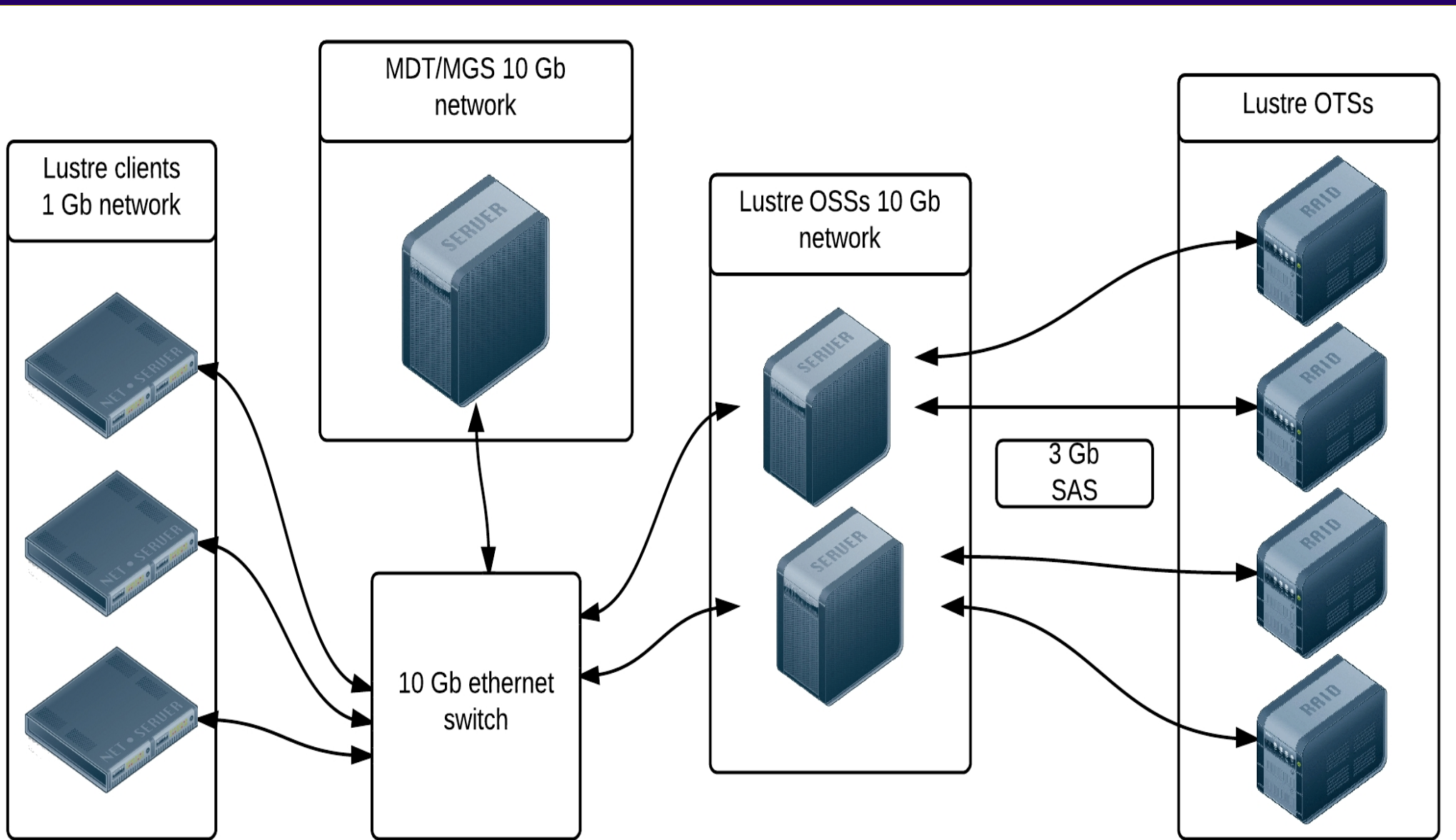
| Nodes | OCFS2 | GFS2 | NFS4 |
|-------|----------|---------|---------|
| 1 | 120 MB/s | 59 MB/s | 57 MB/s |
| 2 | 78 MB/s | 54 MB/s | 53 MB/s |
| 3 | 64 MB/s | 45 MB/s | 51 MB/s |

dd (read), iflag=direct

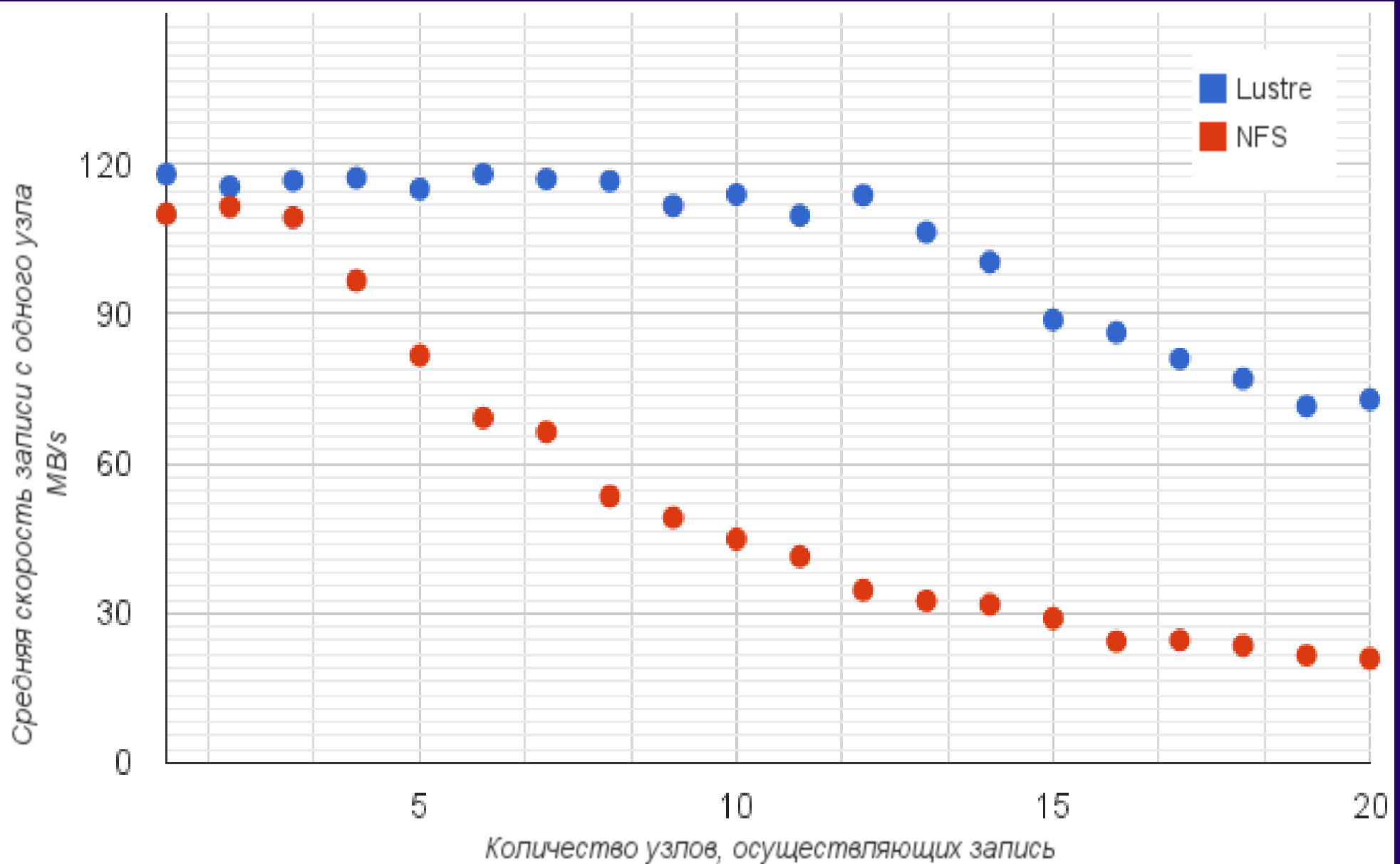
| Nodes | OCFS2 | GFS2 | NFS4 |
|-------|----------|----------|---------|
| 1 | 274 MB/s | 174 MB/s | 29 MB/s |
| 2 | 178 MB/s | 90 MB/s | 29 MB/s |
| 3 | 124 MB/s | 73 MB/s | 29 MB/s |

Распределенная файловая система

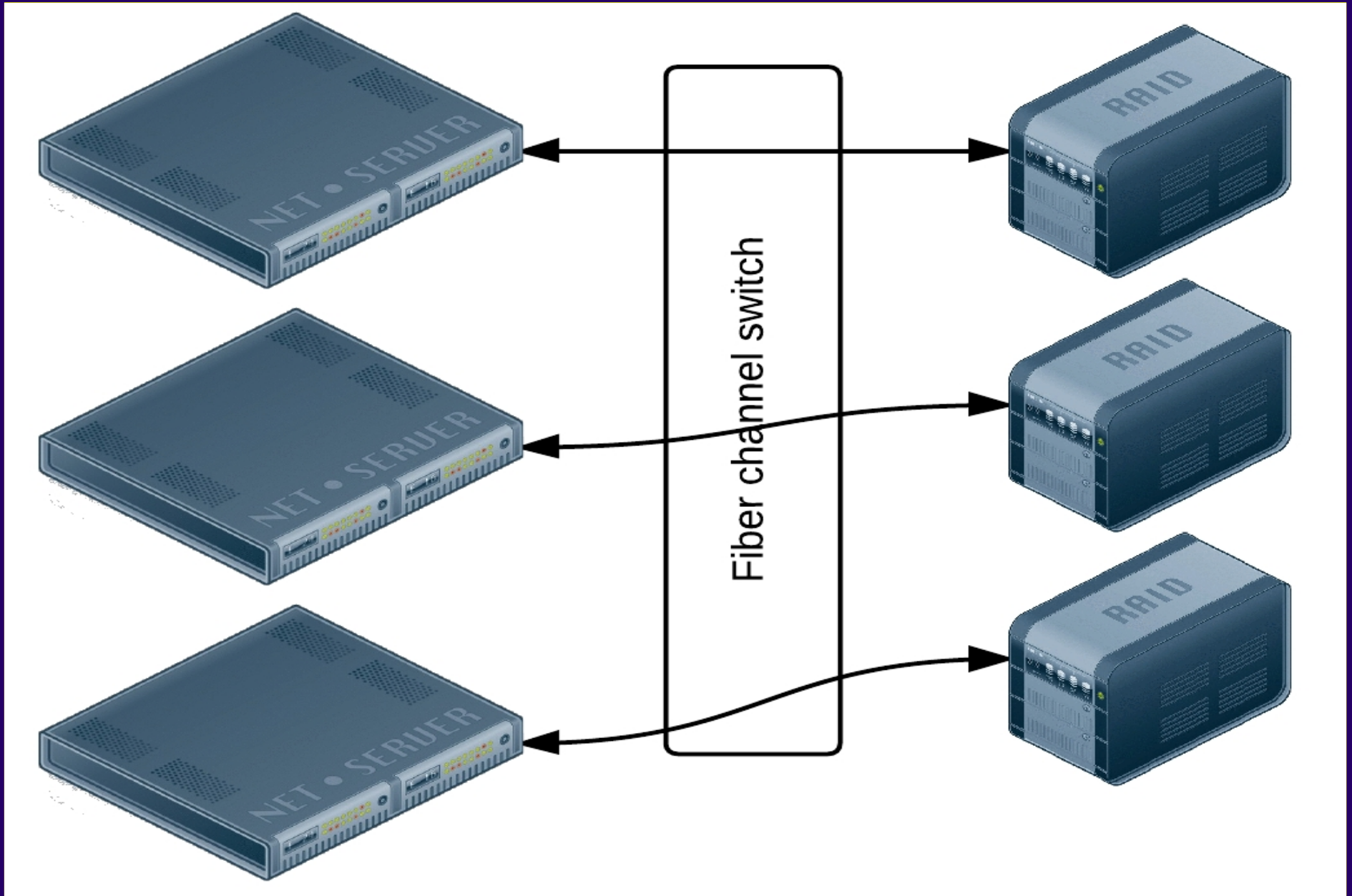
система



Сравнение производительности файловых систем Lustre и NFS



Распределение дисковых ресурсов по проектам



Выводы

- Высокопроизводительное секвенирование – ключевая технология для современной биологии и медицины
- Для обработки геномных данных необходимы соответствующие вычислительные мощности
- Оптимизация потоков данных критична для работы ЦОД
- Сочетание распределенной файловой системы Lustre и инфраструктуры Fiber Channel является оптимальным решением для ЦОД