



ЛАБОРАТОРИЯ  
ИНФОРМАЦИОННЫХ  
ТЕХНОЛОГИЙ  
имени М.Г. Мещерякова

# Распределённая параллельная файловая система Lustre для обработки и анализа данных экспериментов физики высоких энергий

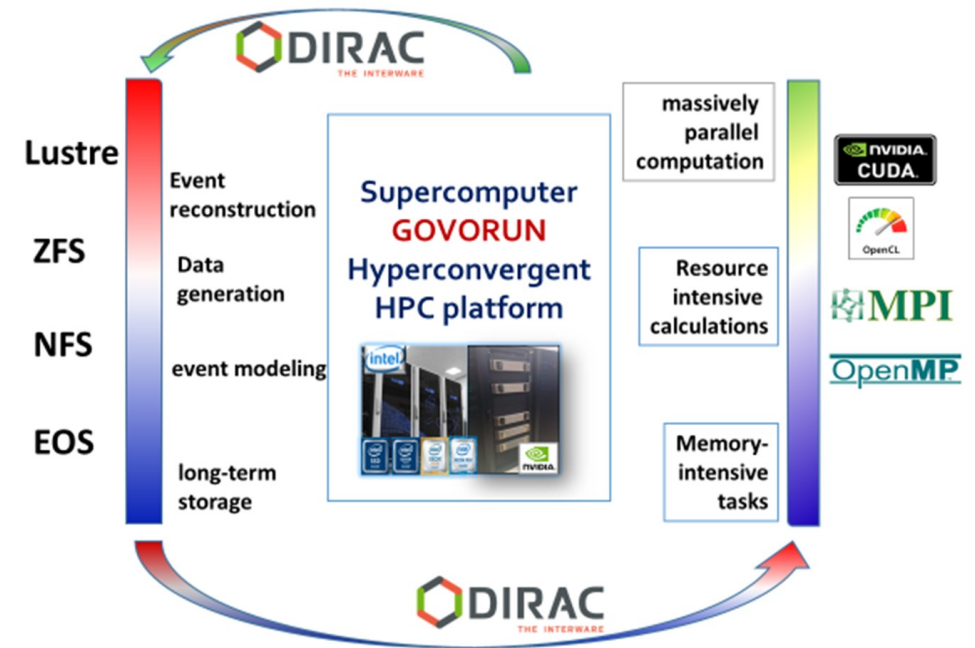
Кокорев А.А., Беляков Д.В., Подгайный Д.В.,

Мошкин А.А., Пелеванюк И.С.

e-mail: [kaa@jinr.ru](mailto:kaa@jinr.ru)



# Nuclotron based Ion Collider facility



**BM@N**

1 PB per year

**MPD**

Planned 20 PB per Year

Now MC generation, reconstruction and physical analysis ~ 6,2PB  
Almost  $10^9$  events

**SPD NICA**

~ 20 PB



# Суперкомпьютер «Говорун»



## CPU компонента

- 21x сервер с Intel Xeon Phi  
Intel Xeon Phi 7290 (72 cores @1.50 GHz), 96 GB RAM
- 76x серверов с Intel Xeon Scalable Gen2 (RSC Tornado TDN511)  
2x Intel Xeon Platinum 8268 (24 Cores @2.90 GHz), 192 GB RAM
- 32x сервера с Intel Xeon Scalable Gen2 (RSC Tornado TDN511S)  
2x Intel Xeon Platinum 8368Q (38 Cores @2.60 GHz), 2 TB RAM

## Пиковая производительность

**1.7 PFLOPS** двойной точности

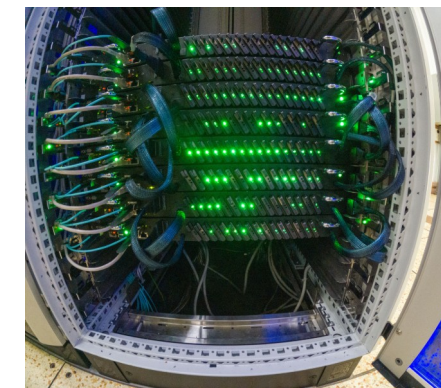
**3.4 PFLOPS** одинарной точности

## GPU компонента

- 5x серверов с NVIDIA V100  
2x Intel Xeon E5-2698 v4 (20 cores @2.20 GHz),  
8x NVIDIA V100 16 GB, 512 GB RAM
- 5x серверов с NVIDIA A100  
2x AMD EPYC 7763 (64 Cores @2.45 GHz),  
8x NVIDIA A100 80 GB, 2 TB RAM



## Система хранения



**8.6 PB**

**26 PFLOPS** половинной точности





# Система хранения и обработки данных гетерогенной платформы HybriLIT



**Warm NFS/ZFS**

**Students Home**  
store[1].hydra.local

**Backups**  
store[2-4].hydra.local

**User Home**  
store[5].hydra.local

**MPD data storage**  
store[6].hydra.local

**BMN data storage**  
store[7].hydra.local

**Warm NFS/ZFS**

**User Scratch (ZFS)**  
s03p002.gvr.local  
/zfs/scratch

**Warm Lustre**

**User Home**  
s03p001.gvr.local  
/lustre/home

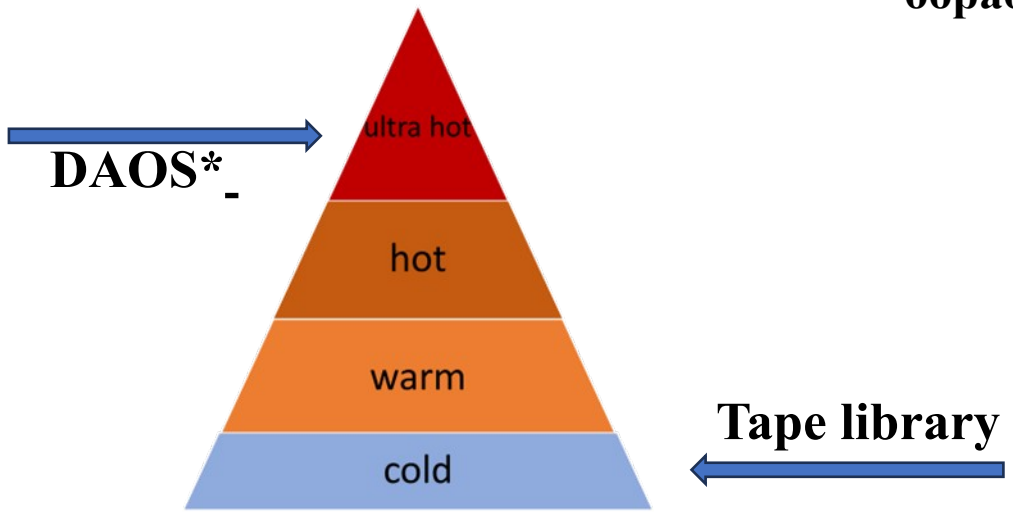
**Hot Lustre**

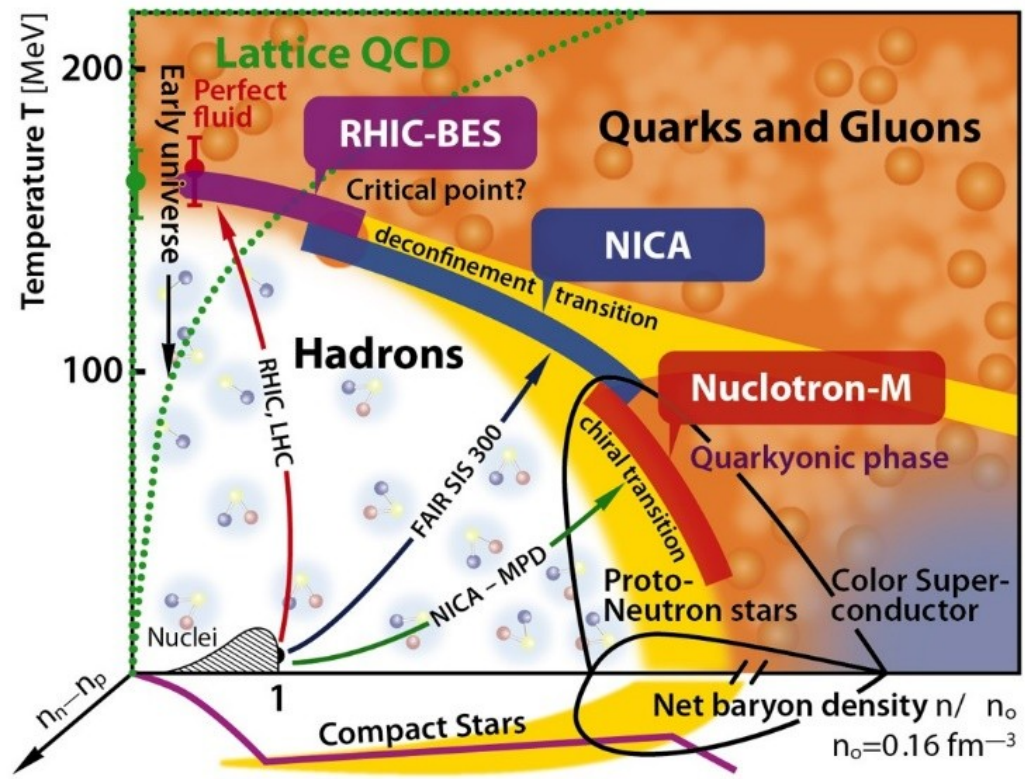
**User Scratch (Lustre) «Lustre 12x12»**  
s01p[001-002].gvr.local  
s02p[001-012].gvr.local  
/lustre/stor1

**MPD data storage «Lustre ruler x4»**  
n04p[001-002].gvr.local  
s03p[003-006].gvr.local  
/lustre/stor2

**Data storage system  
8.6 PB**

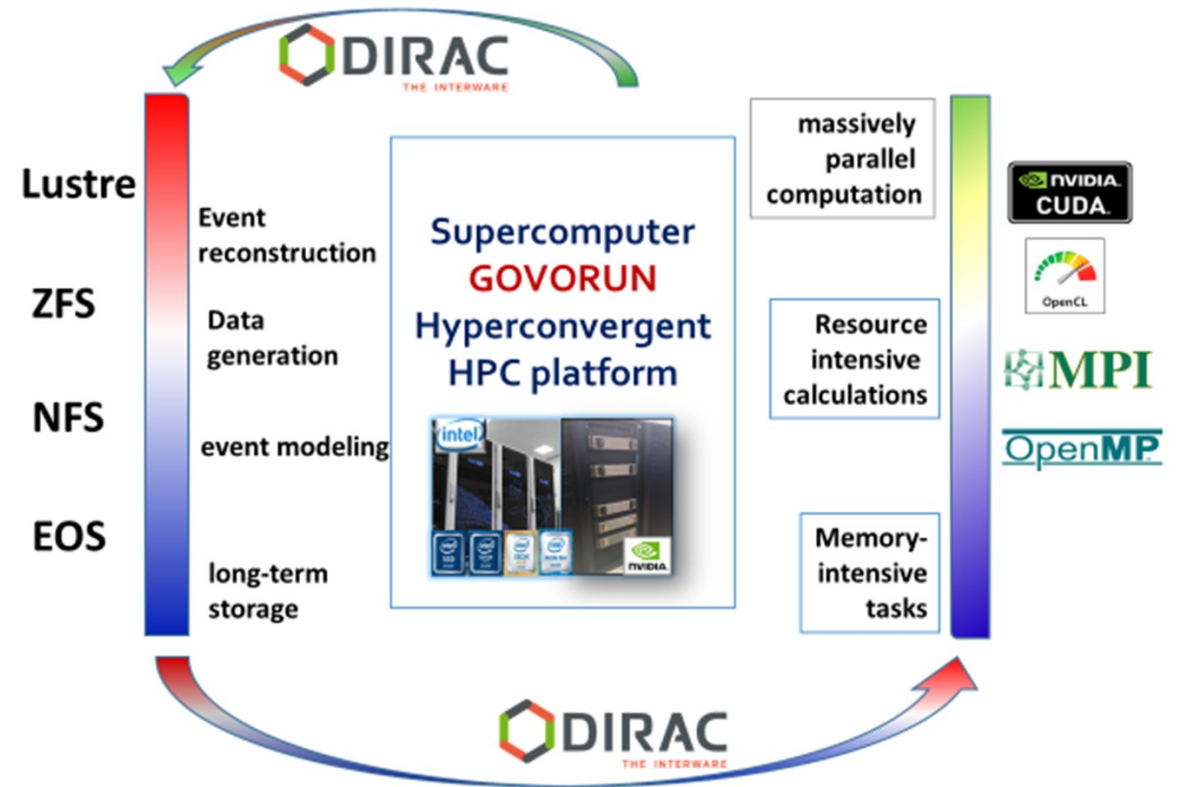
Задачи, решаемые на суперкомпьютере «Говорун», требуют не только больших вычислительных ресурсов, но и привлечения различных технологий для работы с данными, в том числе применения высоко-скоростной внутренней сети с низкой латентностью между компонентами гетерогенной платформы (100 Gbit/sec или выше), а также различных типов систем хранения данных (от "холодных" до "горячих") в соответствии с процедурами обработки входных и выходных данных.





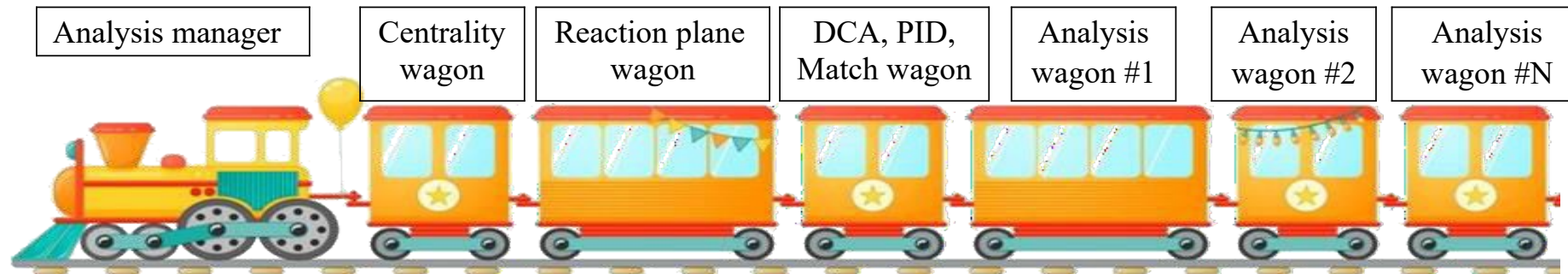
Важнейшие проблемы в изучении свойств адронной материи:

- Природа и свойства сильных взаимодействий между элементарными составляющими Стандартной модели физики элементарных частиц – кварками и глюонами.
- Поиск признаков фазового перехода между адронной материей и КГП; поиск новых фаз барионной материи.
- Исследование основных свойств вакуума сильного взаимодействия и симметрий КХД.



Ресурсы СК «Говорун» используются на всех этапах реализации экспериментов на ускорительном комплексе NICA – от фундаментальных теоретических исследований до обработки экспериментальных данных.

- ❖ Централизованная аналитическая архитектура для доступа и анализа данных → «Поезда анализа», преимущества:
  - ✓ единообразие подходов и результатов совместной работы групп физического анализа, более простое хранение и совместное использование программных кодов и методов.
  - ✓ сокращение количества операций ввода-вывода для дисков и баз данных.
- ❖ Менеджер анализа считывает события в память и вызывает вагоны один за другим для изменения и/или анализа данных:



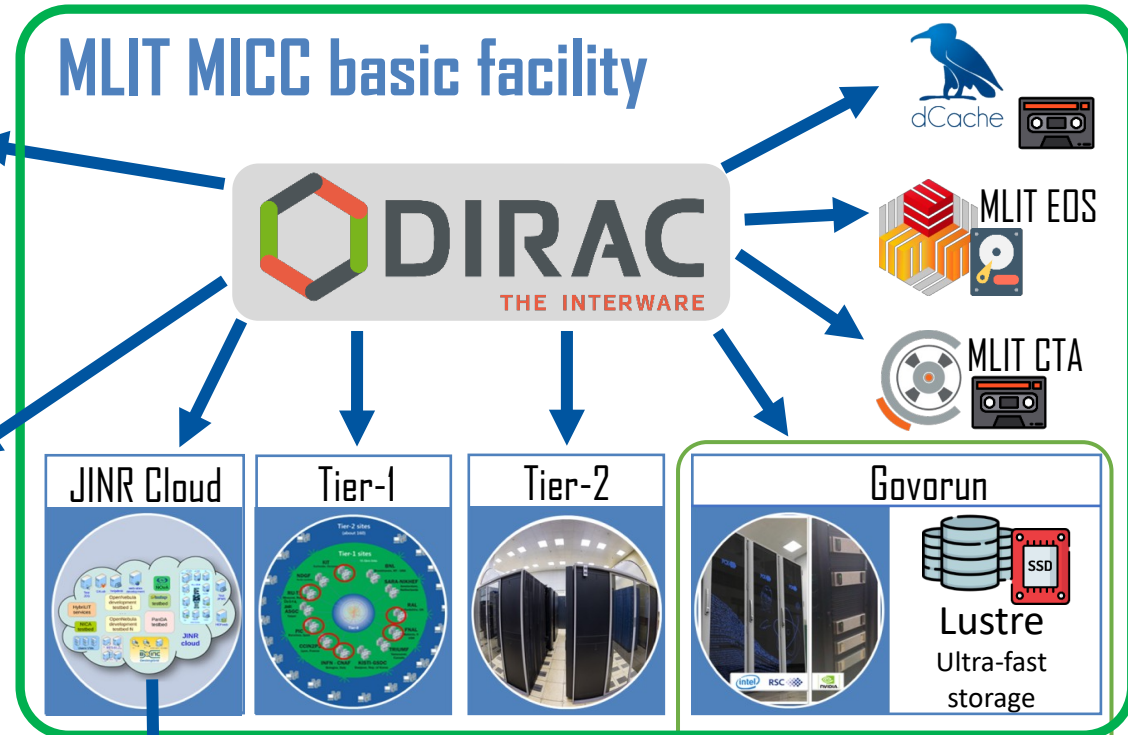
- ❖ Первые запуски анализа по этой схеме произошли в сентябре 2023 г. → с тех пор осуществляются «регулярные рейсы» по запросу групп физического анализа данных коллаборации MPD.
- ❖ Поезду требуется ~ 12 часов (на СК «Говорун») для обработки 50 млн. событий для 10-15 вагонов.
- ❖ По такой схеме, для анализа планируемой физической программы проекта, за 10 запусков были оценены 83 сценария анализа и обработаны 530 млн. моделированных событий коллаборации MPD.



### LHEP resources

DDC cluster	NICA cluster
DAQ computing farm	

### Other/Collaboration resources

### Clouds



- GOVORUN up to 4864 cores in last production
- NICA offline cluster 1000 cores(limit for users)
- Tier1 1500 cores
- Tier2 1000 cores
- Clouds (JINR and JINR Member States)~500 cores
- UNAM (Mexico University) 100 cores
- National Research Computer Network of Russia (now resources from SPBTU and JSCC) 672 cores – New resource, added in 12.2021.

All software packages are centrally stored in /cvmfs and are available on all computing clusters



# Распределенные файловые системы



В ОИЯИ в настоящий момент используются распределенные файловые системы:

- EOS (ЦЕРН), dCACHE (DESY, больше не поддерживается) и CEPH (для облаков).

Общие недостатки – медленные, не позволяют запускать задачи HPC, не позволяют работать с ML/DL

Текущий воркфлоу:

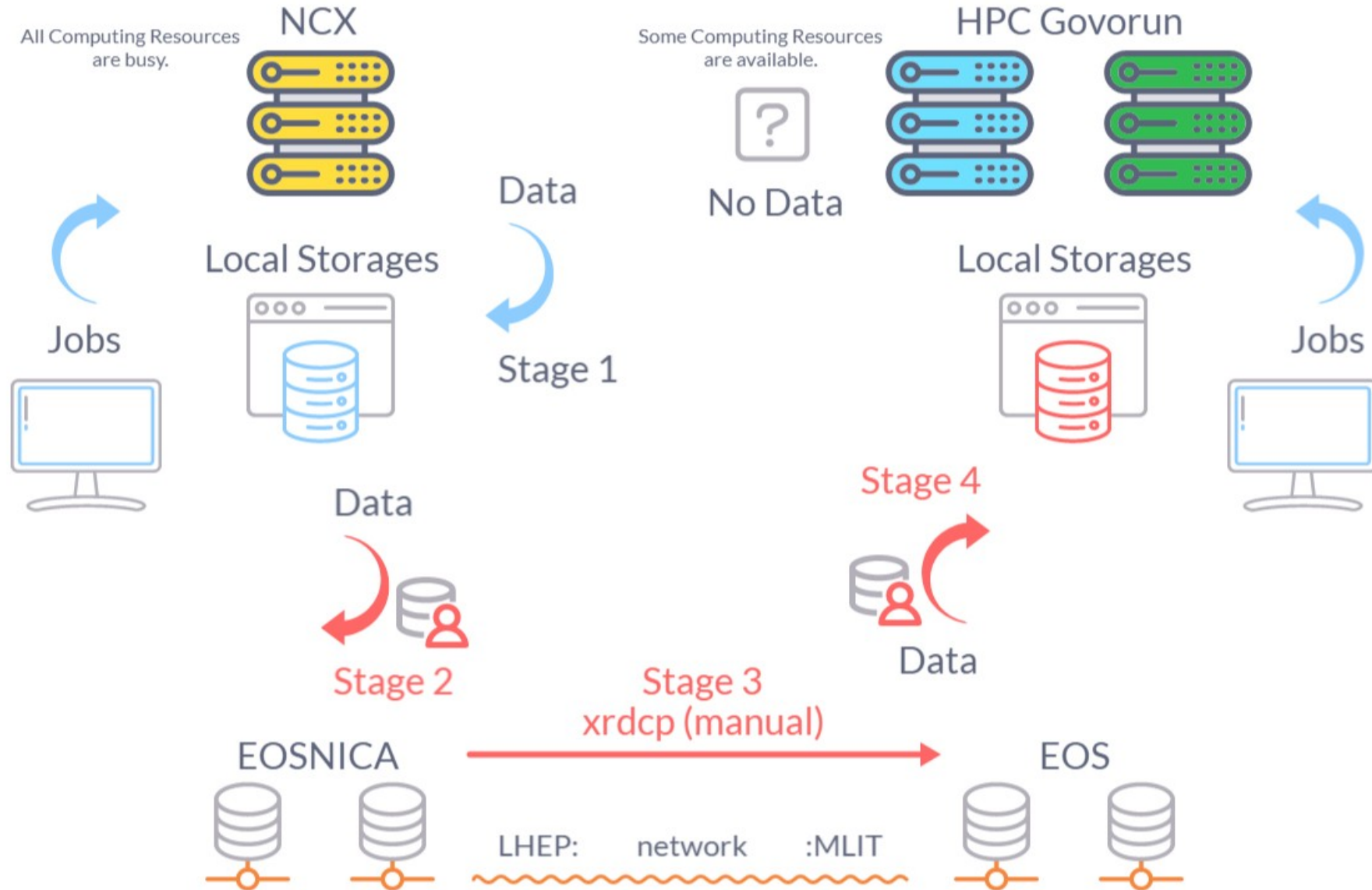
Планируемый воркфлоу:







# Проблема миграции данных между разными вычислительными кластерами

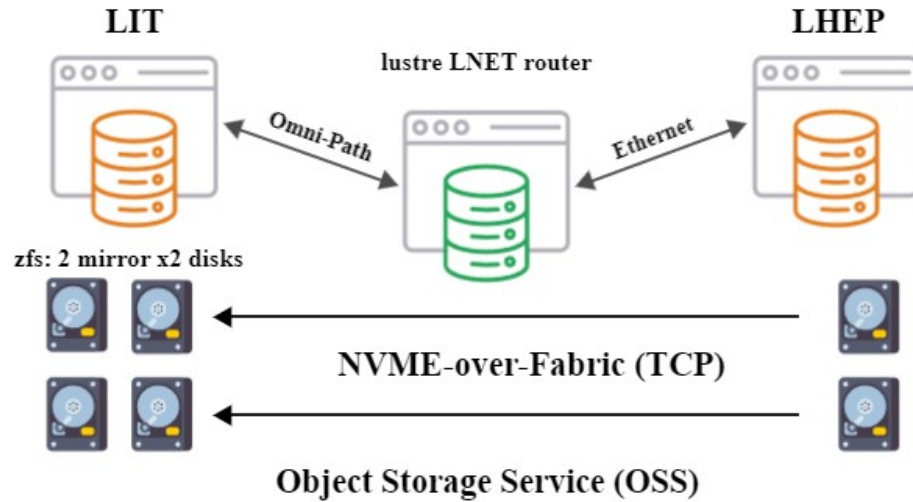




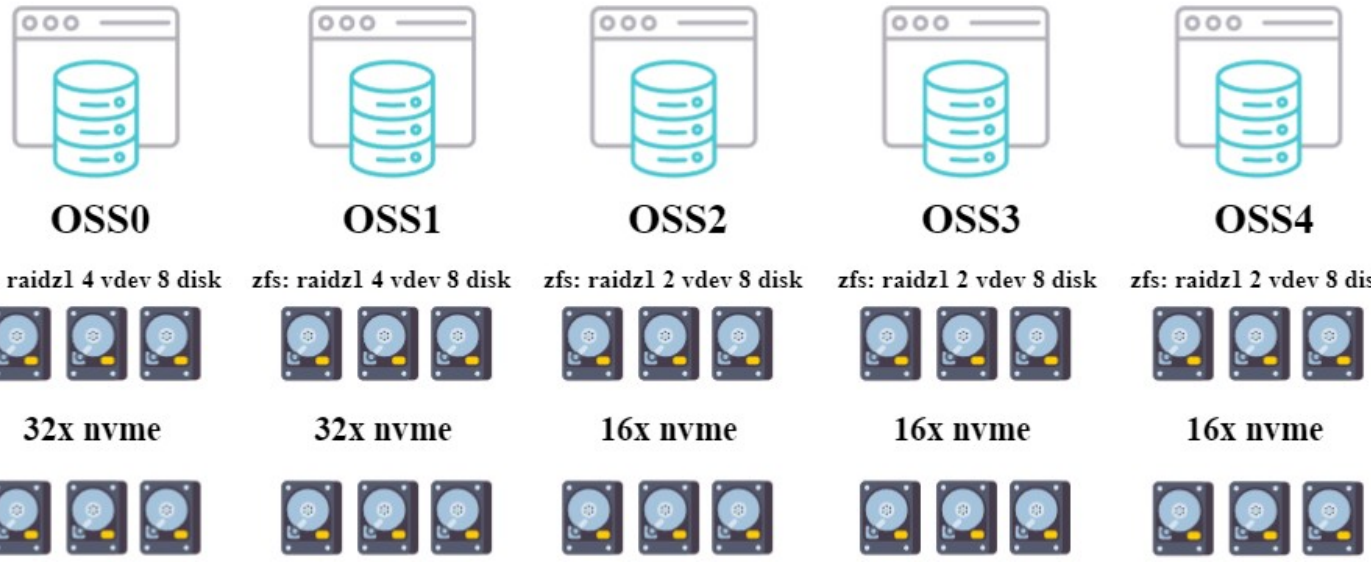
# Схема реализации распределенной параллельной файловой системы Lustre



Management Service (MGS)  
& Metadata Service (MDS)



Распределенная параллельная файловая система Lustre включает в себя два гетерогенных сетевых сегмента. Маршрутизацию между которыми обеспечивает Lustre LNet router.

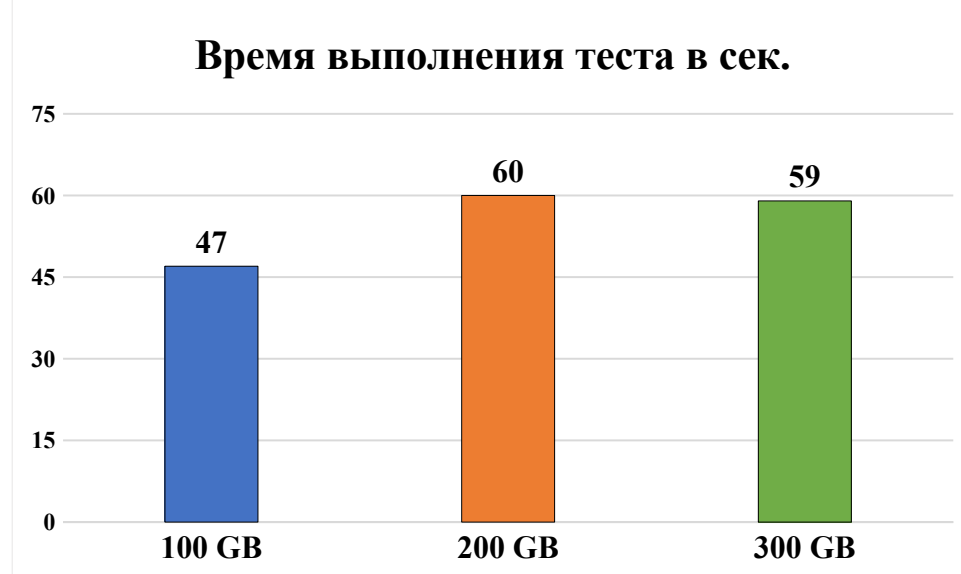
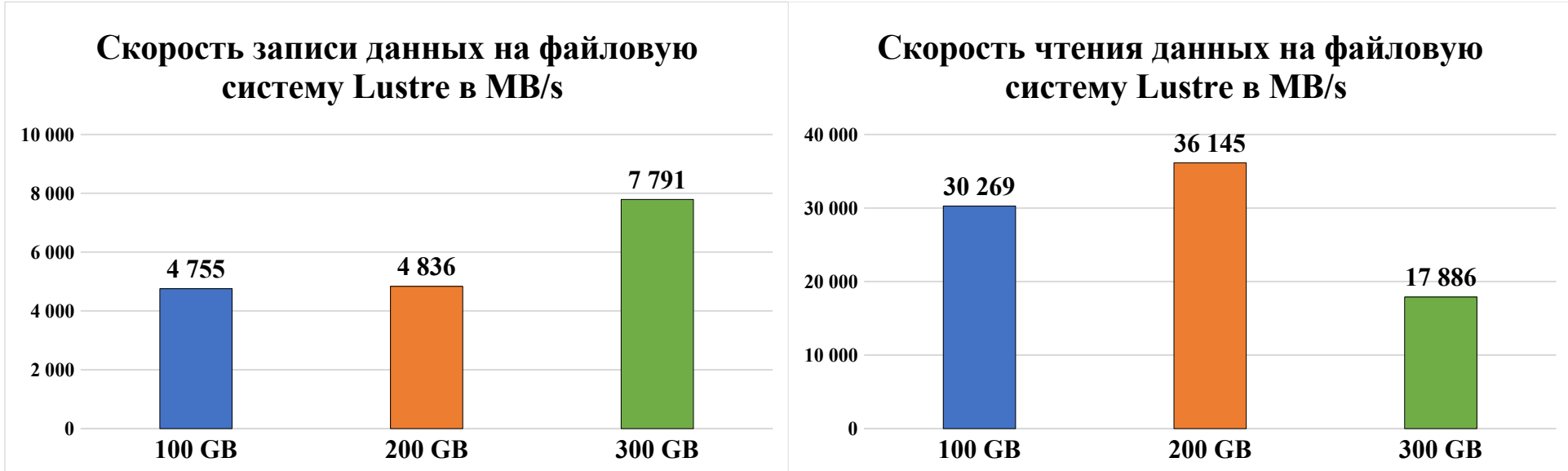


lustre®

2.1 PB



# Синтетический тест распределенной параллельной файловой системы Lustre







# Тестирование Lustre



Задачи обработки данных эксперимента VM@N были выбраны для функционального теста построенной системы. На данный момент эти задачи, с точки зрения работы с данными, являются самыми интенсивными. Описание процесса выполнения этих задач приведено в таблице:

	Step 1 Data load	Step 2 Processing	Step 3 Processing	Step 4 Data Upload	Total
Duration	150 s	600 s	1800 s	10 s	2560 s
CPU load	0	100 %	100 %	0	2400 s
Disk Write (GB)	15	7	1	0	23 GB
Disk Read (GB)	0	15	15	1	31 GB
Network in (GB)	15	0	0	0	15 GB
Network out (GB)	0	0	0	1	1 GB

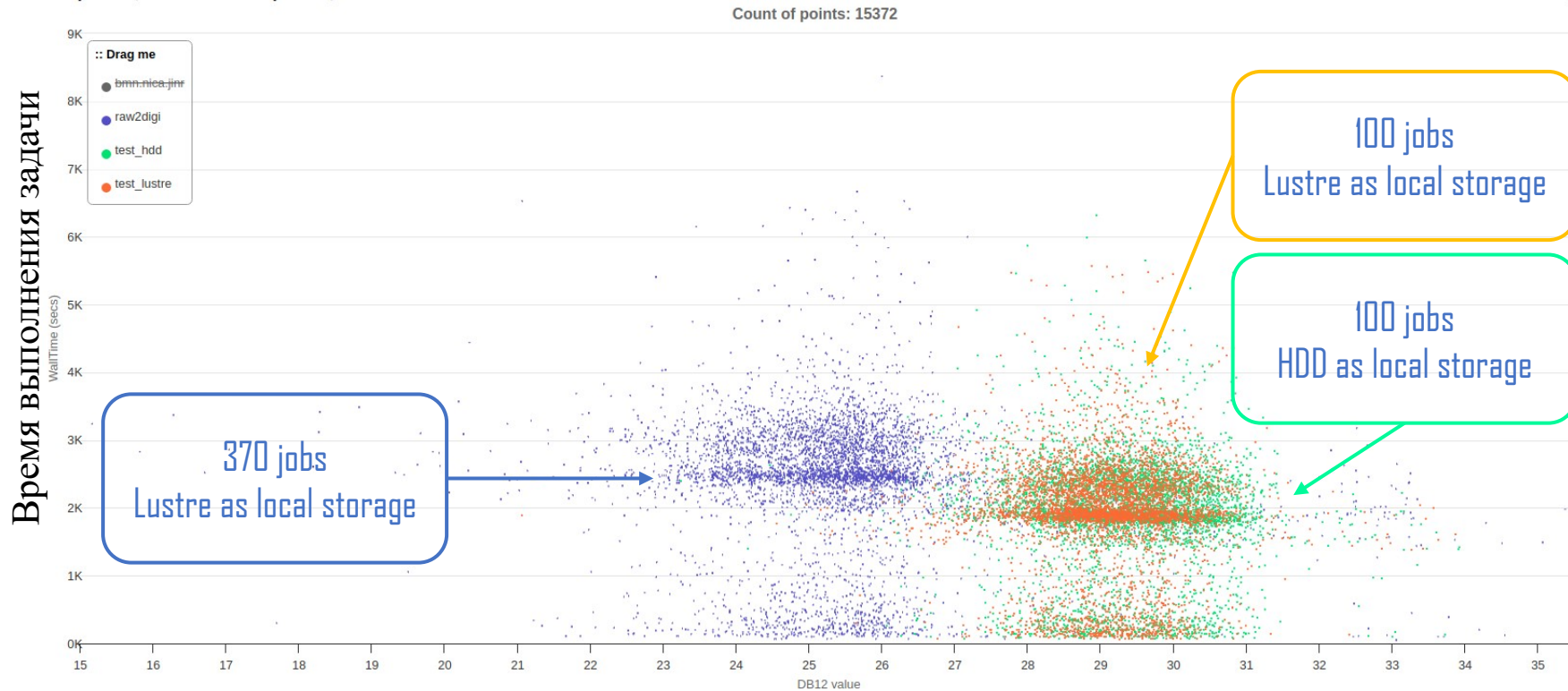
Каждая задача требует около 23 ГБ локального хранилища. Из-за этого невозможно загрузить все доступные на сервере ядра. NCX кластер обладает ограниченным количеством дискового пространства, что не позволяет загружать его ядра задачами обработки данных VM@N более чем на 27%.



# Тестирование Lustre



Как задачи выполняются на NCX кластере.



Значение бенчмарка ядра в условных единицах

1. Задачи работают на HDD. CPU load ~27 %.
2. Задачи работают на Lustre. CPU load 100%. Наблюдается замедление.
3. Задачи работают на Lustre. CPU load ~27%. Производительность такая же как на HDD.



# Тестирование Lustre



## Статистический анализ данных

Название теста (cores_storage)	Средняя длительность задачи	Средняя скорость ядра (относительно 100_hdd)	Скорость кластера (относительно 100_hdd)
100_hdd	1900 s	1.0	1.0
100_lustre	1908 s	0.996	0.996
370_lustre	2390 s	0.798	2.954

При использовании всех ядер бенчмарк каждого ядра падает с 29.5 до 25 (- 15%). Если говорить о времени выполнения задач, то при использовании Lustre мы наблюдаем увеличение времени на 20%. Таким образом мы считаем, что при 370 одновременно работающих задачах, использование Lustre замедляет задачи лишь на 5%. И это при том, что суммарный поток данных во время работы задач примерно 3 ГБ/с на запись и 4 ГБ/с на чтение.

Не смотря на накладные расходы, созданная система позволила ускорить процесс выполнения пакета задач в 3 раза!





- Настройка High Availability для управляющих серверов (MGS/MDS) и серверов хранения данных (OSS).
- Реализация протокола NVMe-over-Fabric (TCP).
- Настройка Lnet router для гетерогенных сетевых сегментов (Ethernet, Omni-Path).
- Настройка LNet сервиса (конфигурирования).
- Выбор оптимальной конфигурации ZFS пулов для серверов хранения данных (OSS).

## **Преимущества данной реализации распределенной файловой системы Lustre:**

- ❖ Возможность проводить расчеты на разных вычислительных ресурсах без переноса данных между кластерами.
- ❖ Высокая производительность выполнения счетных задач с интенсивными Input/Output операциями на систему хранения данных.
- ❖ Возможность пользователя выбирать сервер хранения (OSS), количество реплик данных, количество разбиений файла на части (chunks) для ускорения работы с файлом (работа в режиме RAID).

**Спасибо за внимание!!!**