# Санкт-Петербургский государственный университет информационных технологий, механики и оптики кафедра вычислительной техники

# Разработка архитектуры и методов организации слабосвязанных архивных систем для автоматизации проектирования

Лукьянов Н.М.

Руководитель:

к.т.н., доц. Тимченко Б.Д.

2011

# Тенденции развития CAD систем

- Растущая производительность рабочих станций
- Переход от суперкомпьютерных архитектур к распределенным высокопроизводительным системам
- Уменьшение цикла разработки изделия
  - По данным Boeing 60 месяцев -> 12 месяцев
- Снижение стоимости разработок
  - По данным Boeing 7 млрд. -> 1 млрд. \$
- Интеграция с PLM и ERP системами
- Увеличение объема хранения периодически используемых данных
- Переход к распределенному хранению данных
- Необходимость снижения совокупной стоимости владения инфраструктурой



# Задачи работы

- Разработать архитектуру распределенной системы архивного хранения и доставки файловых данных, построенную на слабосвязанных гетерогенных узлах с узкими каналами связи
- Предложить подход к разделению программной и аппаратной составляющих системы
- Разработать метод адаптации системы в части управления, как размещением, так и доставкой контента.

# Информационное обеспечение САПР

### Архивируемые данные:

- выполненные проекты
- фотографии объектов, материалов, помещений
- видеоотчеты пусконаладочных работ и т.д.
- 3D модели
- подписанная КД

### Объем данных на 1 организацию:

Средний объем проекта – 500 Гб Проектов в год – 30 Глубина работы с архивом – 5 лет

Общий объем хранения архивной системы (на организацию) – 75 Тб

# Классические решения



Центролизованные NAS, SAN системы:

высокая стартовая стоимость и высокая ТСО



Ленточные (LTO) и DVD-RAM библиотеки:

Низкая скорость доставки короткая жизнь носителей

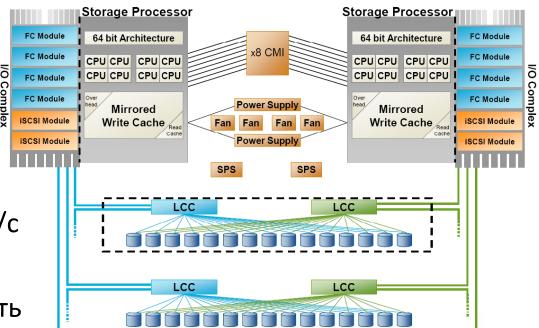
# CXД EMC "Celerra"

### **EMC Celerra CX4-120**

- Средний ценовой диапазон
- TCO 500\$ / Тб / год
- SAN сеть 8 Гбит/с
- NAS NFS, CIFS 10 Гбит/с

### Ограничения

- высокая гранулярность масштабирования
- практически локальная система
- распределенность обеспечивается только дополнительными СХД



### Нет поддержки

- FTP и ACL
- HTTP
- WebDAV

# Слабосвязанные системы

### • слабая связанность

- передача данных по публичным или нестабильным каналам с ограниченной пропускной способностью
- работа за пределами ЦОД

### • гетерогенность

- унаследованные аппаратные компоненты
- различное программное окружение

# • хранение неструктурированных данных

• файлы, а не данные в РБД

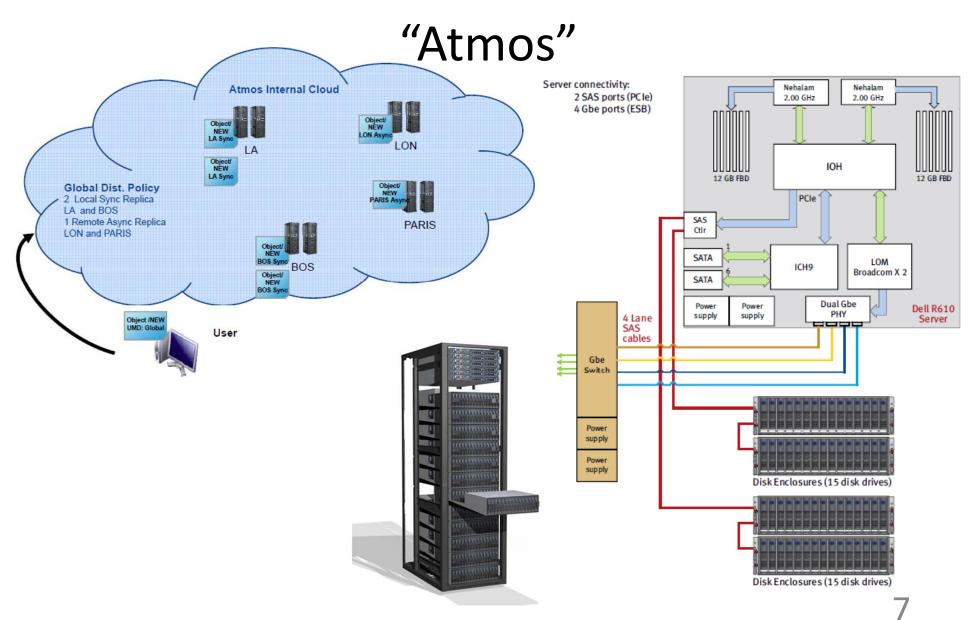
# • эластичное развитие (масштабирование)

- улучшение характеристик без принципиального изменения архитектуры
- плавное увеличение объемов дискового пространства

# • открытость для развития

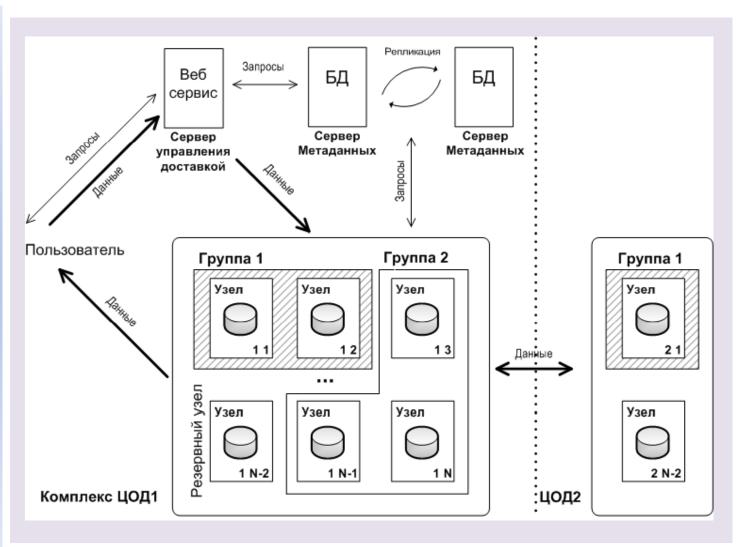
- универсальный (общедоступный) протокол доступа к системе
- возможность интеграции в существующие сервисы

# Слабосвязанная система ЕМС



# Архитектура системы

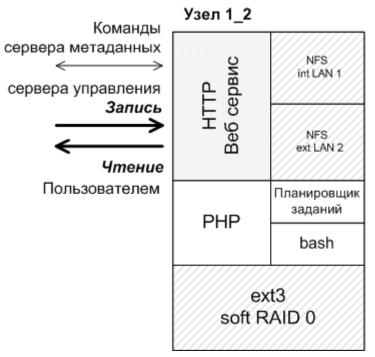
- Открытые протоколы передачи данных:
  - для всей системы **HTTP**
  - между узлами **NFS**
- Ввод/ вывод файлов поверх **HTTP** с помощью программны х компонент
- Открытая БД
  PostgreSQL



# Узел хранения

репликация

с Узлами



### Загрузка файла, запрос 1/3:

PUT /newfolder1/adk45x\_r4324.jpg HTTP/1.1

accept: \*/\*

date: Mon, 02 Aug 2010 18:02:34 GMT+3 content-type: application/octet-stream

host: 172.16.34.51 content-length: 18471

x-access-key: \_N45\$0KL4sdMpDj\_hI8ZRnE= content-md5: tv6qMXw6eTuumbd40yndvJ==

HTTP/1.1 200 OK

Date: Mon, 02 Aug 2010 18:02:34 GMT+3

Загрузка файла, ответ узла:

Server: Apache Content-Length: 0 Connection: close

Content-Type: text/plain;

charset=UTF-8

- именование:
  - node(ЦОД) (номер).domain.name
- нет резервирования данных внутри узла:
  - резервные узлы внутри комплекса
- 2 сетевых интерфейса:
  - внутренний >=1Gb/s
  - внешний 100 Mbit/s
- передача данных:
  - **чтение** пользователем напрямую
  - **запись** сервером управления

# Сервер метаданных

### Основные функции:

- Взаимодействие с серверами управления доставкой
- Хранение ссылок на файлы
- Составление таблицы узлов данных по рейтингу на основе загрузки узла
- Динамическая выборка ссылок

Q	1	Приемник К		Сервер	N.
		Коннектор К	3	метаданны Обработчик	M
<b>A</b>		Обработчик данных		метаданных	
		Локатор соединений	4	Локатор соединений	M
		(6)	8	Модуль статистики	С
		Загрузчик С			
		7			
		Приемник К			
_		Модуль к статистики Локатор			
		соединений			

Сервер управления доставкой

Узел хранения

- Параметры узлаВесовой коэф-тЗагрузка процессора0,15Загрузка канала0,75Заполнение диска0,1
- 10 мегабайт метаданных 1 Гб полезных данных
- 4 запроса в БД метаданных для выборки файла

10

# Управление количеством реплик

Параметр	Балл W <sub>k</sub>	Bec. коэф w <sub>k</sub>
Популярность файла Р	4	0,45
Степень важности І	2	0,22
Тип файла <b>Т</b>	2	0,22
Свободное место носителя <b>F</b>	1	0,11

 $\mathbf{W_k}$  - баллы экспертной оценки тонкой настройки системы

$$w_k = rac{W_k}{\displaystyle \sum_k W_k}$$
 расчета весового коэффициента  $w_k$  для k-ого параметра

 $\mathbf{N}_{\mathsf{add}}$  — наибольшее число реплик, разрешенное в системе

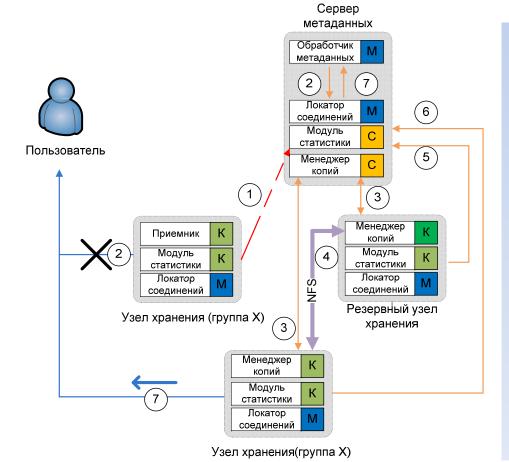
**P, I, F, T**<sub>max</sub> – максимальные значения показателей

Расчет рекомендуемого количества реплик для заданного файла:

$$N_{i} = N_{\min} + N_{add} * (w_{pop} \frac{P_{i}}{P_{\max}} + w_{imp} \frac{I_{i}}{I_{\max}} + w_{free} \frac{F_{i}}{F_{\max}} + w_{type} \frac{T_{i}}{T_{\max}})$$

 $\mathbf{N}_{\mathsf{min}}$  — минимальное число реплик (три)

# Восстановление данных



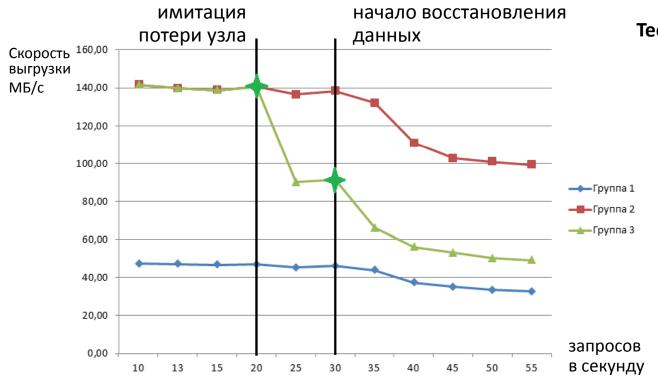
### Особенности алгоритма:

- прямое копирование данных через внутреннее высокоскоростное соединение **int LAN 1** по протоколу NFS
- минимальное взаимодействие с сервером метаданных
  - блокировка поврежденного узла (3 запроса)
  - изменение таблицы предпочтительности узлов (2 запроса)
  - ввод в эксплуатацию узла (З запроса)
- Корректировка только записей об изменении ІР адреса узла и его имени.
- Отсутствие необходимости аппаратного RAID контроллера на узлах

### Группа Х Группа Х Группа Х Исключение Изменение ключей node1 6 ▲ RW node1 2.domain.name A RW node1 2.domain.name A RW node1 2.domain.name ▲ № node1 8.domain.name ▲ № node1 8.domain.name A RW node1 8.domain.name ≜ № node1 6.domain.name node1 7.domain.name на node1\_7 A RW node1 7.domain.name ение данных

# Восстановление данных

Тип	Объем дисков, МБ	Кол-во дисков, шт.	Уровень RAID	Пропускная способность NAS, MБ/с	Время восстановления, ч	Нормированное время, ч
SAS	300	4	5	109	4,7	7,8
SATA	1500	4	5	94	31,5	10,5
Узел	1500	1	нет	81	20,8	20,8
хранения	1500	2	0	97	18,2	18,2



### Тестовые группы:

- Группа 1 − 1 узел
- Группа 2 3 узла
- Группа 3 3 узла / восстановление данных

### Экспериментальный запрос:

синхронный 15 файлов х 100 кб

# Результаты работы

- Разработана архитектура распределенной системы архивного хранения данных, обеспечивающая, в отличии от центролизованных систем:
  - практически неограниченное, плавное масштабирование
  - использование гетерогенных, экономически доступных аппаратных компонент
- Разработан программный метод адаптации системы в части размещения и доставки контента:
  - повышение скорости доставки
  - снижение затрат архивной памяти
- Предложен авторский алгоритм управления количеством реплик, обеспечивающий достижение требуемой доступности хранимых данных
- Разработан программный метод восстановления данных с отказавших узлов системы, исключающий безвозвратную потерю данных
- Реализована и опытно эксплуатируется распределенная система хранения файлового контента в области архивного хранения и представления видеоданных:
  - плавное расширение и соответствующее снижение единовременных капитальных затрат
  - сохранение работоспособности системы при отказах узлов

# СПАСИБО ЗА ВНИМАНИЕ!