



Сетевые ускорения в комплексе Скала-CP / Postgres Pro: настоящее и будущее





Скала-CP / Postgres Pro



Идея, системная интеграция, сопровождение

Суперкомпьютерное сетевое оборудование с поддержкой RDMA



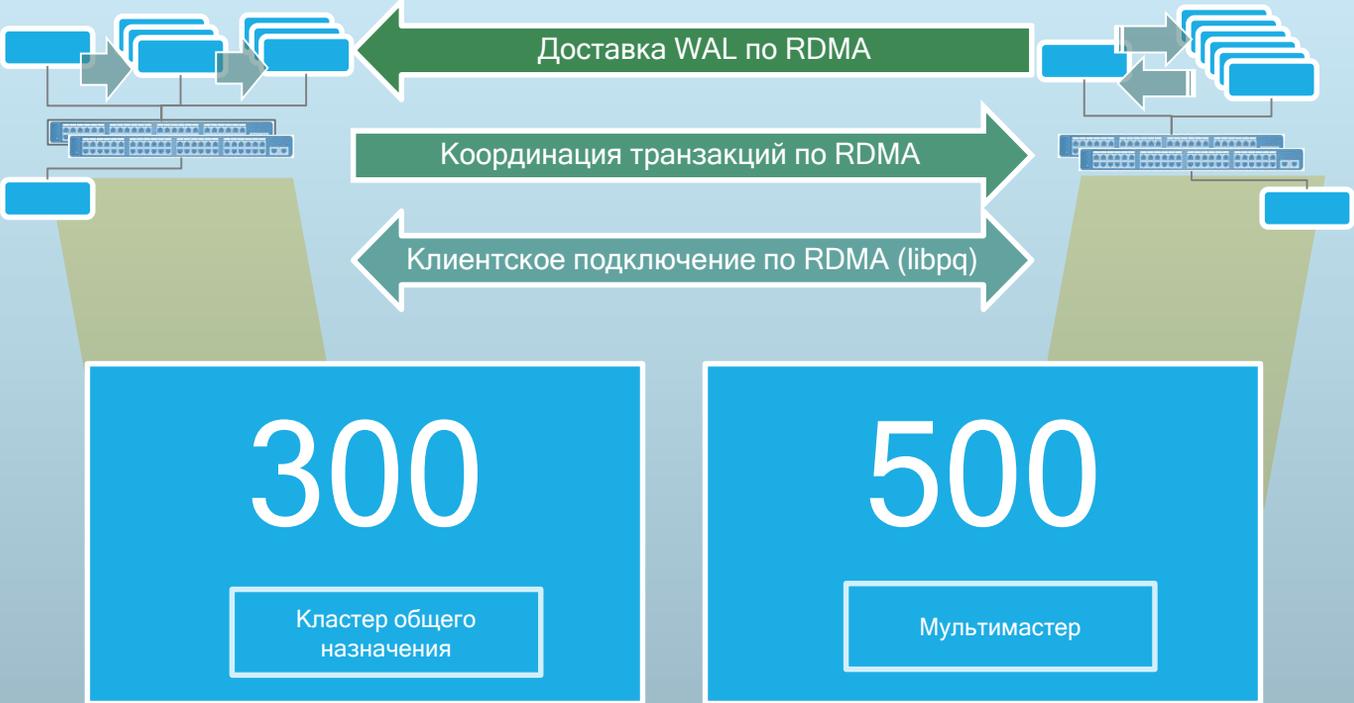
Специальная версия Postgres Pro EE с поддержкой RDMA

Российское серверное оборудование



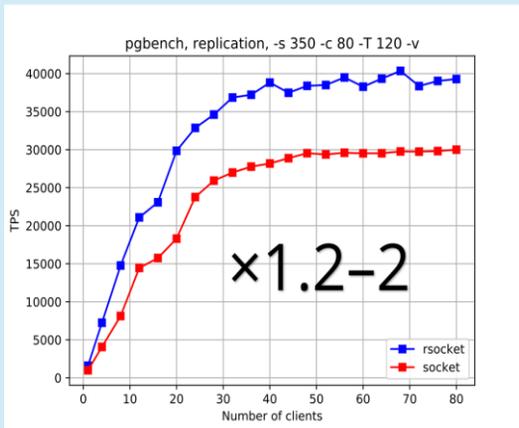


Скала-CP / Postgres Pro: 300 и 500

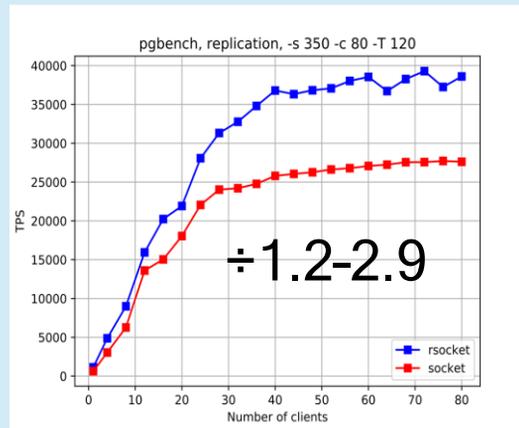


Репликация по RDMA

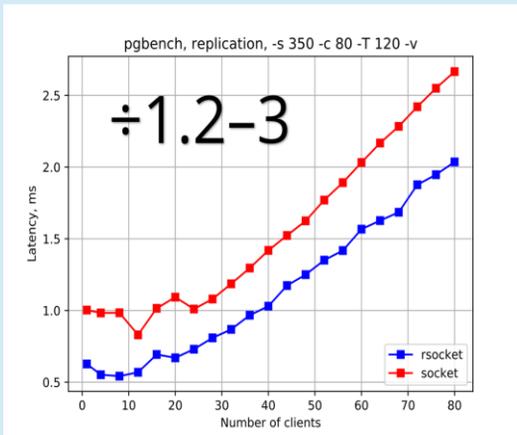
Read-write with synchronous replication, TPS



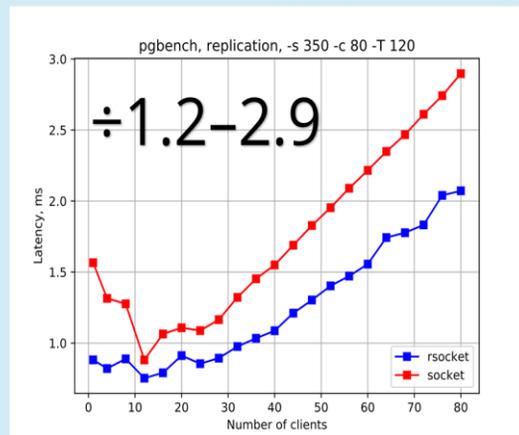
Read-write with remote apply replication, TPS



Read-write with synchronous replication, latency



Read-write with remote apply replication, latency





Аппаратные трюки в мире машин баз данных

OLTP

- ORACLE Exadata Xn-8
- Clustrix
- IBM Pure for Transactions
- DELL DAAD

СУБД общего назначения

- ORACLE Exadata Xn-2
- HUAWEI Fusion Cube / DB
- Скала-CP / Postgres Pro

ROLAP без разделяемых ресурсов

- TERADATA
- IBM Netezza
- PARACCEL
- Pivotal Greenplum DB
- Hewlett Packard Enterprise MICRO FOCUS Vertica

SN-ROLAP + Hadoop

- TERADATA Aster
- Pivotal Greenplum DB + Greenplum HD



- FPGA



- InfiniBand



- RoCE



- PostgreSQL

Почему InfiniBand?



Открытый сетевой стандарт

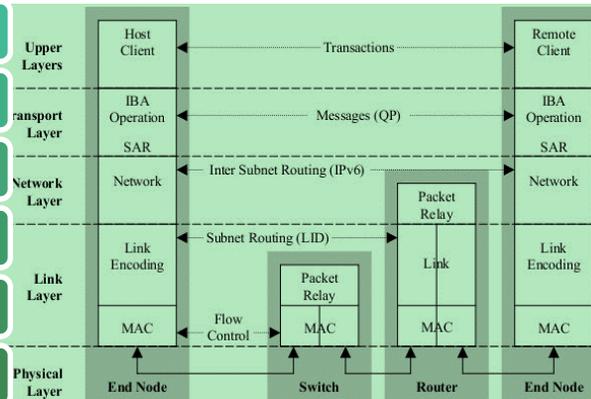
Полностью переделанные все 7 сетевых уровней

Наиболее быстрая сетевая инфраструктура - до 200Гб/с

Самые низкие задержки- 90ns на каждом коммутаторе

Полный HW Offload всех сетевых задач

За счет этого, используя Verbs API, E2E задержки <0.7μs

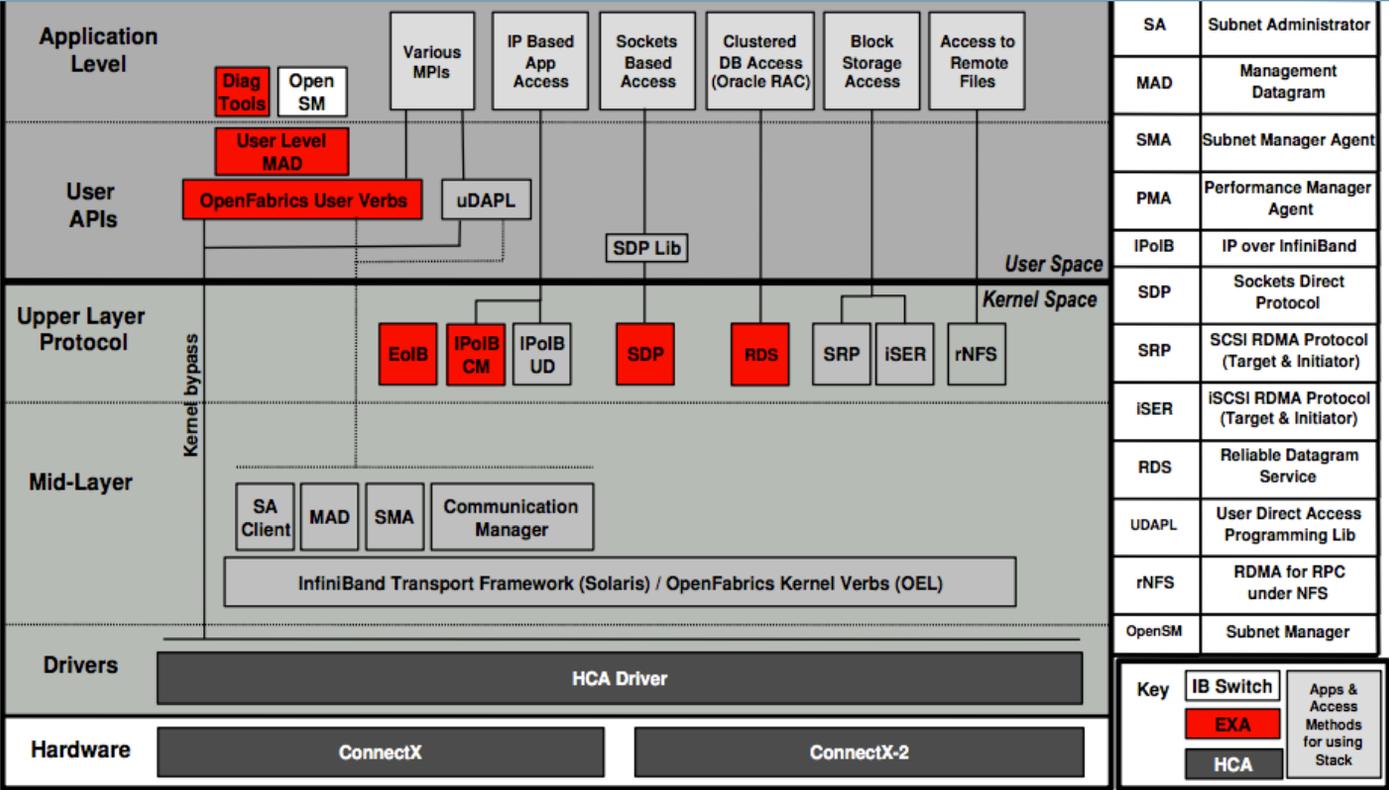


Какой прирост производительности при переходе от 10GbE к 40Gb/s IB?

В 40 раз!



InfiniBand у «больших вендоров»: Oracle iDB и Teradata BYNet



SA	Subnet Administrator
MAD	Management Datagram
SMA	Subnet Manager Agent
PMA	Performance Manager Agent
IPoIB	IP over InfiniBand
SDP	Sockets Direct Protocol
SRP	SCSI RDMA Protocol (Target & Initiator)
ISER	iSCSI RDMA Protocol (Target & Initiator)
RDS	Reliable Datagram Service
UDAPL	User Direct Access Programming Lib
rNFS	RDMA for RPC under NFS
OpenSM	Subnet Manager

ИТ-директор: «сеть должна быть только Ethernet!»

RoCE: свойства Infiniband в Ethernet

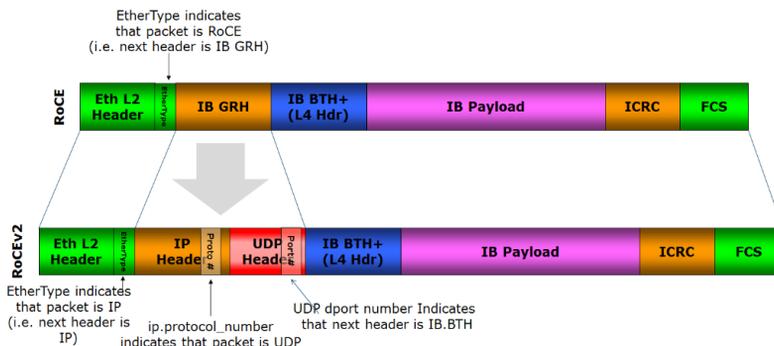


Низкие задержки сети <300µs, 1.2µs E2E

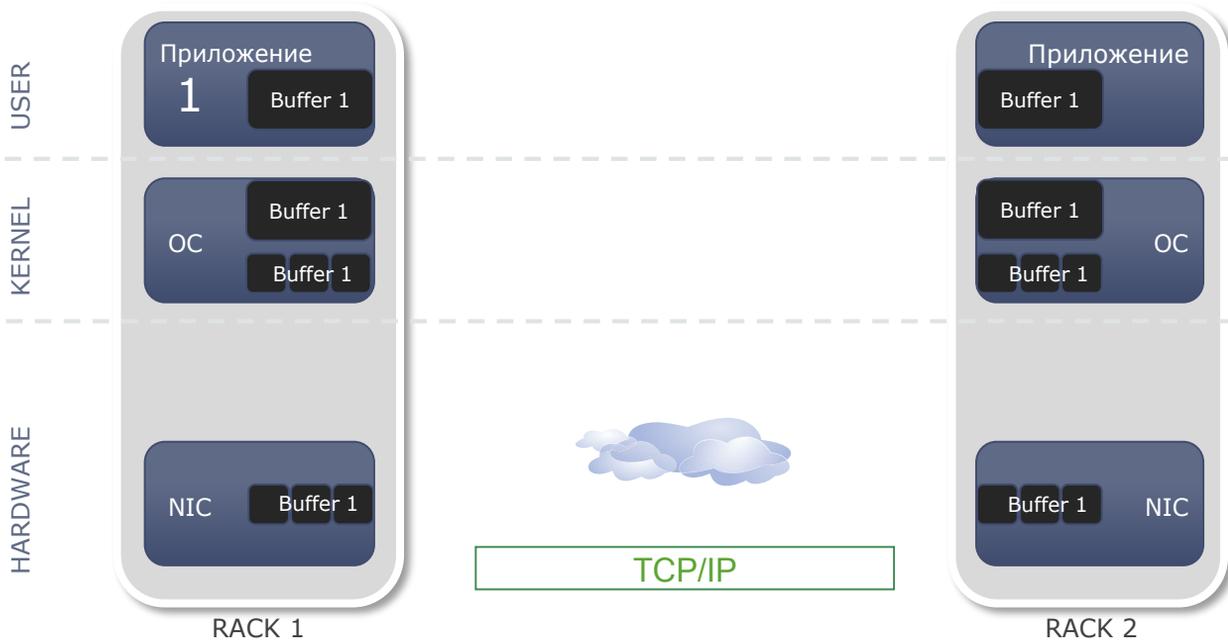
Гарантированная доставка пакетов (Zero Packet Loss)

PFC+ECN для организации lossless-среды

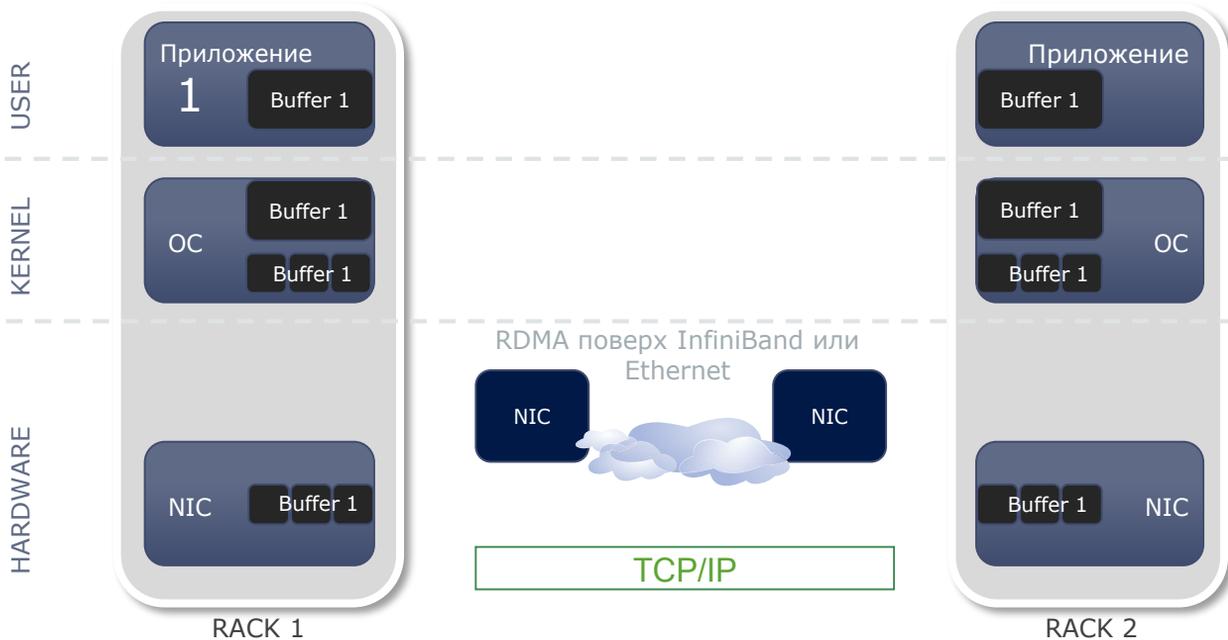
RoCE API на базе Verbs с небольшими отличиями



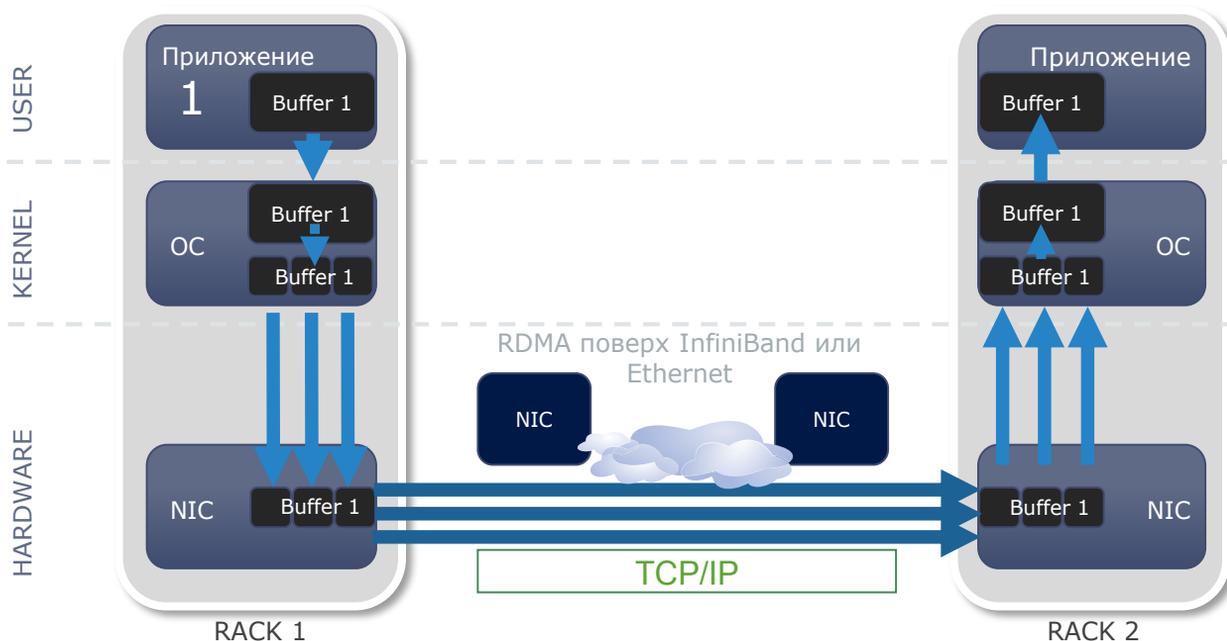
Zero-copy + CPU Offload



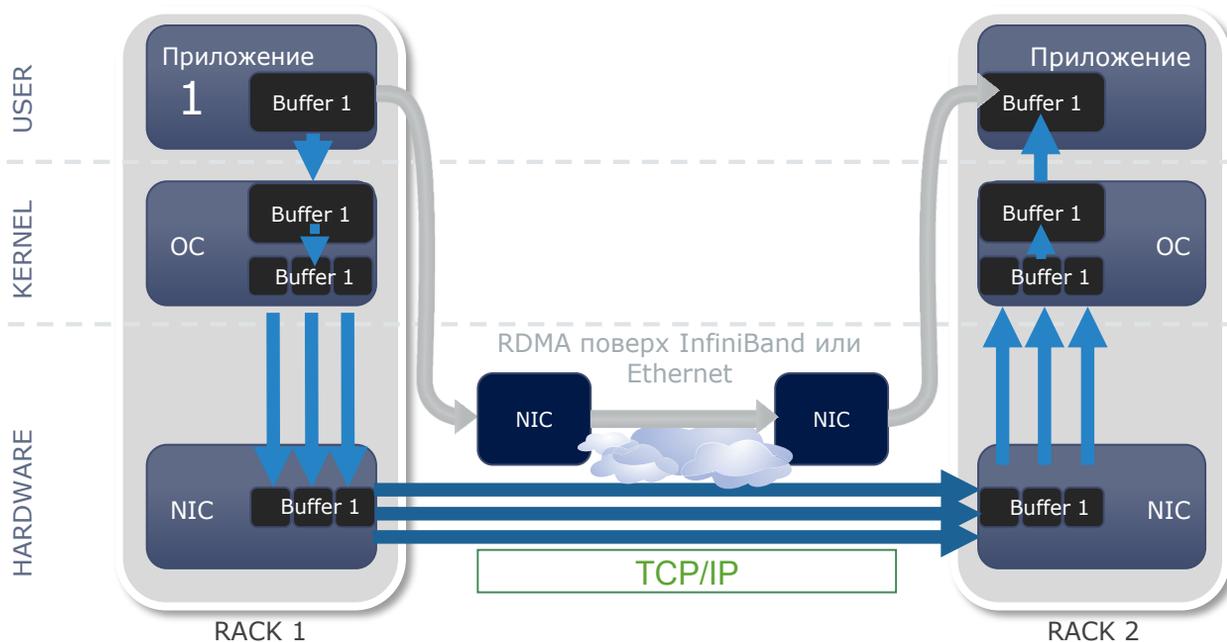
Zero-copy + CPU Offload



Zero-copy + CPU Offload



Zero-copy + CPU Offload





Эффект CPU offload? Много синхронных реплик

Производительность, tps

0 5000 10000 15000 20000 25000 30000

Только мастер

1 реплика

2 реплики

3 реплики

4 реплики

RDMA NoRDMA

Задержка, мс

0 10 20 30 40 50 60 70 80 90 100

Только мастер

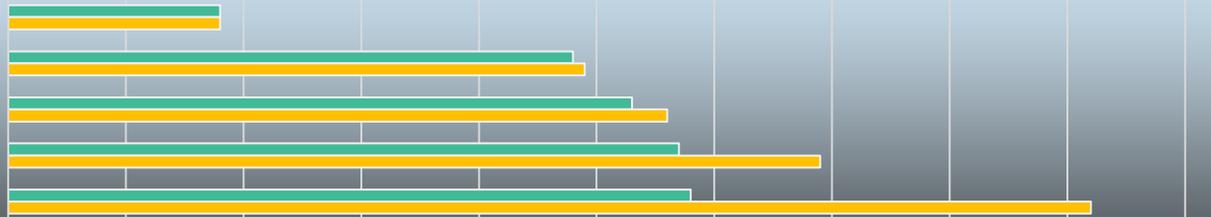
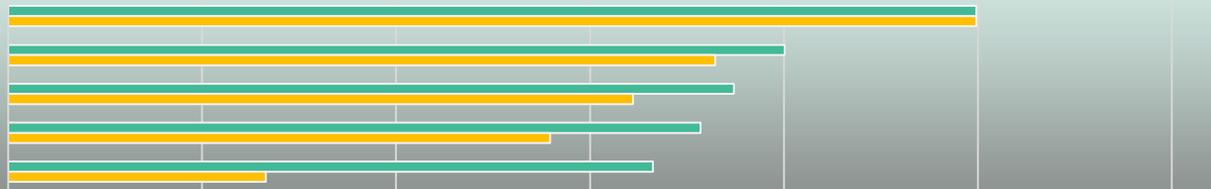
1 реплика

2 реплики

3 реплики

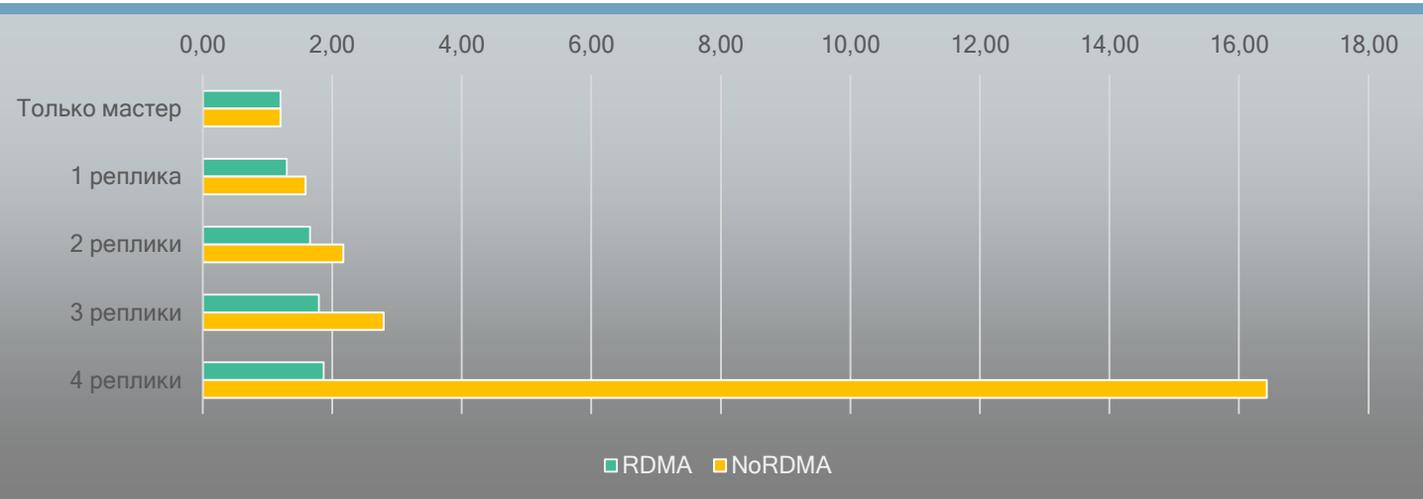
4 реплики

RDMA NoRDMA





Время ЦПУ мастера на транзакцию (мс)



Disclaimer: поверхностный расчёт

(из Load Average на всю систему, с ожиданиями и прочими неаккуратностями)

Но эффект очевидный, значительный

И это далеко не максимум возможностей!

Первая реализация: RSocket

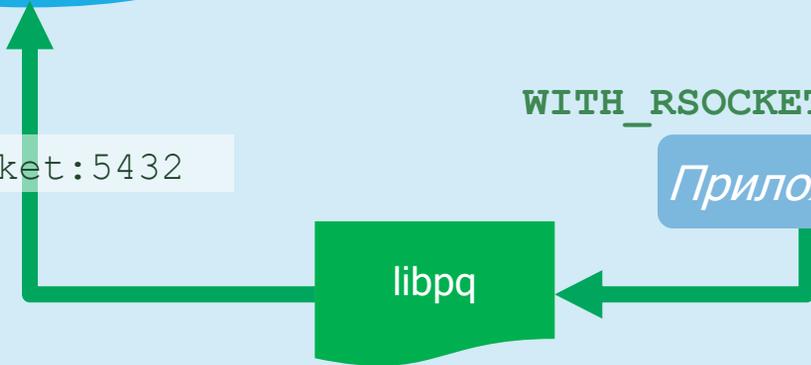


socket:5432

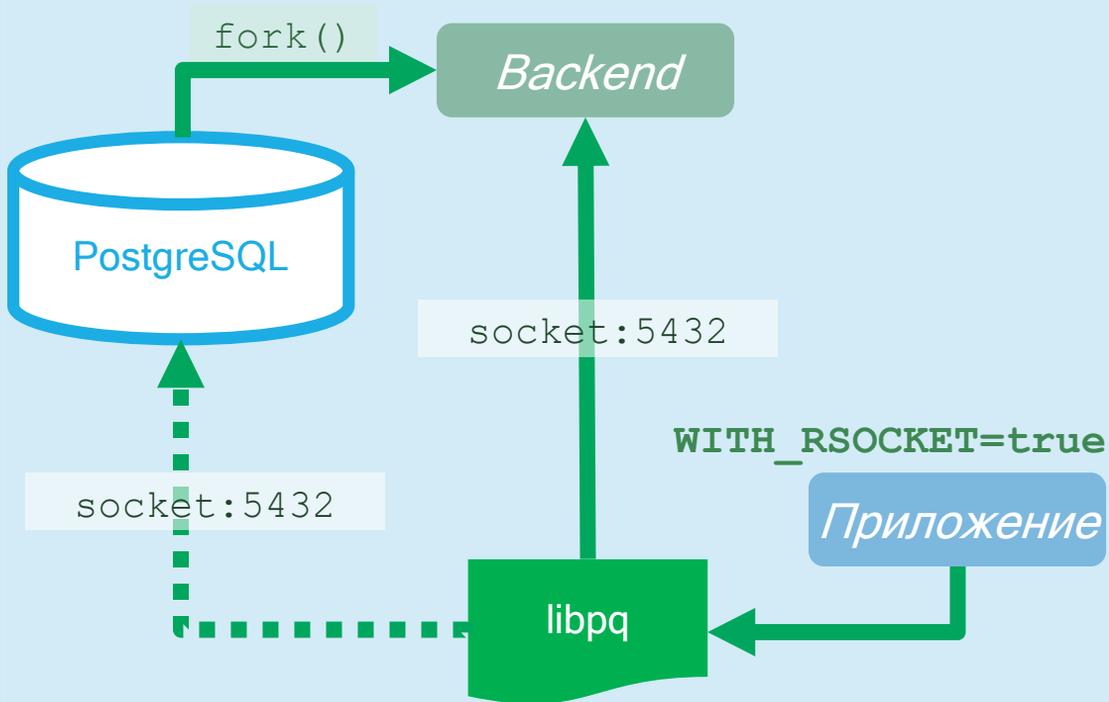
`WITH_RSOCKET=true`

Приложение

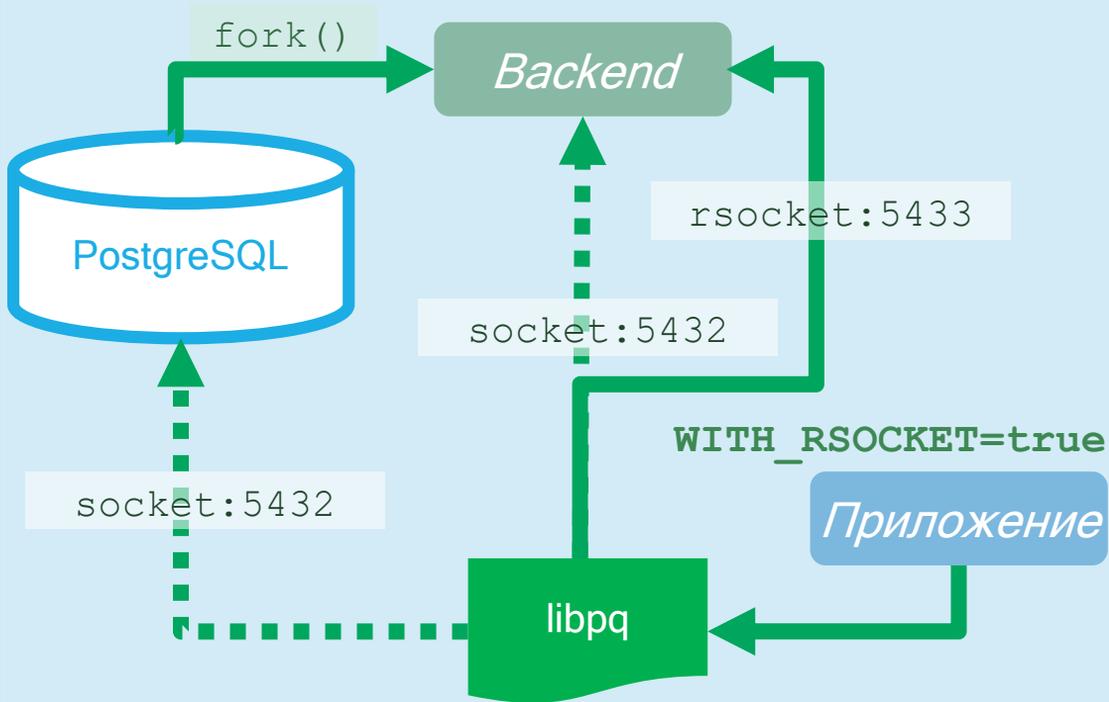
libpq



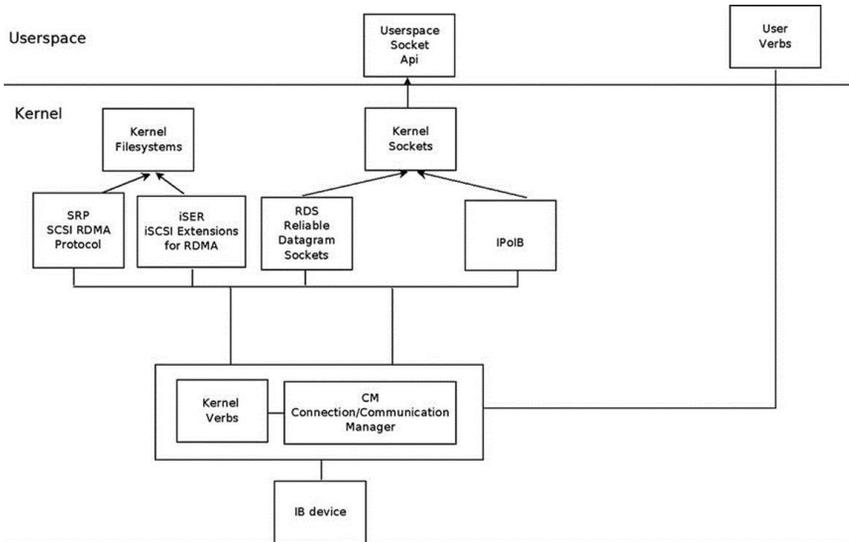
Как установить RDMA-соединение?



...и работаем через RSocket



Verbs: родной язык RDMA



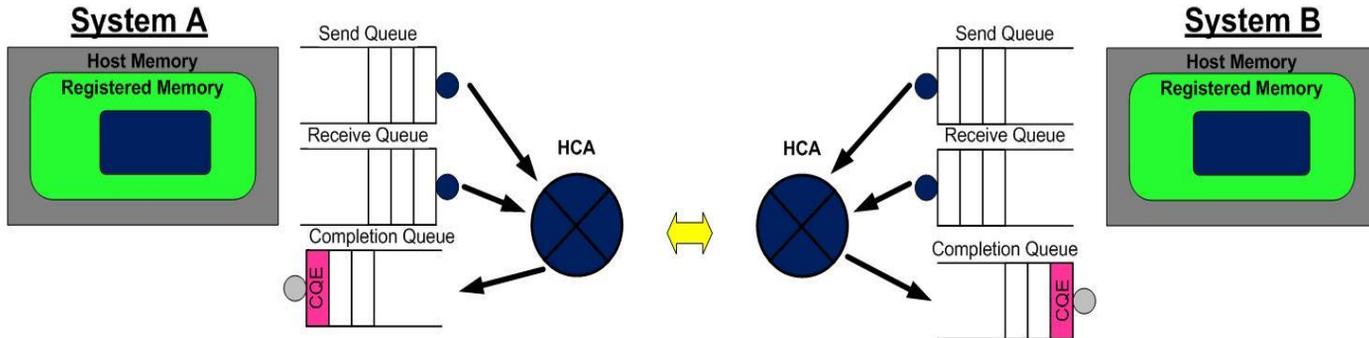
Verbs: родной язык RDMA



- LID - Local Identifier (address which is assigned to any port in subnet)
- GID - Global Identifier (128b address to identify endpoint of multicast group)
- QP - Queue Pair
- CQ - Completion Queue
- PD - Protection Domain
- MR - Memory Region
- MW - Memory Window
- AV - Address Vector
- WR - Work Request
- WQE - Work Queue Entry
- SR - Send Request
- RR - Receive Request
- CQE - Completion Queue Entry

SEND: send only the data (the responder need to post RR)
RDMA Write: send the data + remote address and key
RDMA Read: send a read request to the remote side:
remote address and key + remote size
ATOMIC (Cmp&Swp, Fetch&Add) - Atomic operations
(Read-Modify-Write in one atomic operation) which are
used to implement mutual exclusives objects

Verbs: родной язык RDMA



Возможности приблизиться к «родному языку»



SDP

- Socket Direct Protocol
- Исторически первая реализация «переходника» с TCP-приложений на RDMA
- Устарел, не поддерживается в OFED

VMA

- Mellanox Voltaire Messaging Accelerator
- Та же идея, что RSocket: сокетоподобный язык для миграции приложений на RDMA
- Свободное ПО, активно развивается и поддерживается Mellanox

UCX

- Двухуровневый фреймворк, позволяющий как быстро, так и «тонко» перейти на RDMA
- Свободное ПО, активно развивается и поддерживается Mellanox

Applications

MPICH, Open-MPI, etc.

RPC, Machine Learning, etc.

PGAS/SHMEM, UPC, etc.

SPARK, Hadoop, etc.

UCX

UC-P (Protocols) – High Level API

Transport selection, cross-transport multi-rail, fragmentation, emulation of unsupported operations

Message Passing API Domain:
send/receive, tag matching

I/O API Domain:
Stream

Task Based API Domain:
Active Messages

PGAS API Domain:
Remote memory access

UC-T (Hardware Transports) – Low Level API

Send/Recv, RMA, Atomic, Tag-matching, Active Message

Transport for RoCE/IB Verbs

RC

DCT

UD

Transport for GPU memory access

CUDA

ROCM

Other transports

Shared
Memory

Gemini

UC-S (Services) Common Utilities

Utilities

Data
structures

Memory
management

OFA Verbs Driver

CUDA

ROCM

Hardware

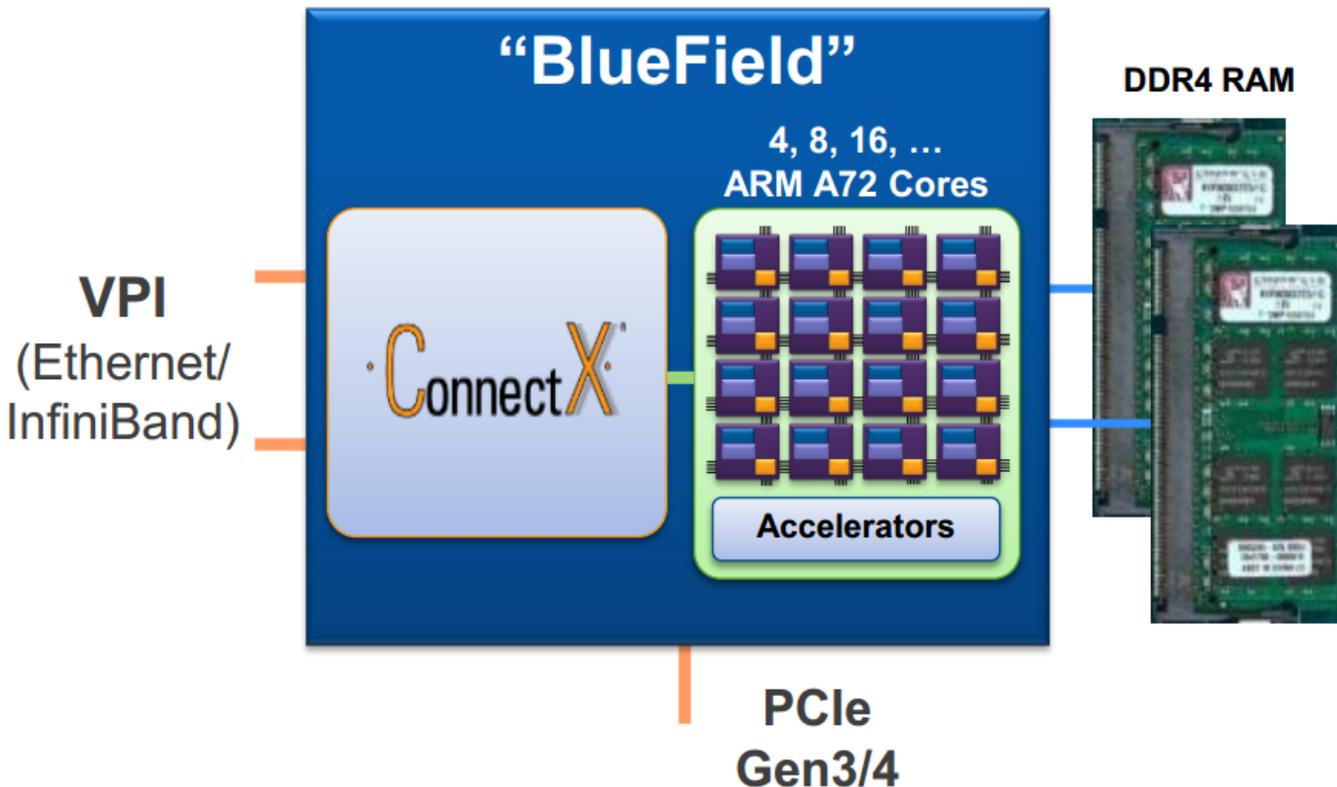
Переход на UCSX

- Перспектива уйти от fork()
- Архитектурные тонкости...

RDMA для pg_shardman

- Возможность мгновенной коммуникации между сегментами

BlueField: компьютер в сетевом адаптере



Система хранения для PostgreSQL

JBOF + NVMeF

16 NVMe SSD
без PCIe-
коммутатора

Система управления на сетевом адаптере

Corosync +
Pacemaker -
внутри хоста, и
на отдельном
чипе

Мониторинг



Спасибо за внимание!

skalar@ibs.ru

