

SBERDATAFUSION

Проксирование данных для Hadoop
и Greenplum



Андрей Ильин

Главный эксперт по технологиям, Сбер



Алексей Тютякин

Главный инженер по разработке, Сбер



Проблемы текущих трактов



Удобно для пользователей, но высокая стоимость сопровождения и hardware, длительные процессы:

Сложный процесс копирования метаданных между кластерами



Ручной процесс копирования прав доступа к объектам Hive между кластерами



Одноуровневый доступ к репликам данных (архивам)



Необходимость доработки приложений на каждом кластере, если требуется отдавать данные в защищенном виде



Текущий альтернативный процесс распространения данных включает в себя промежуточный кластер хранения данных, куда данные периодически копируются с разных источников, что приводит к высокому потреблению hardware и низкой актуальности данных



Большое кол-во копий данных



Сложности в копировании данных между различными стеками, например, из Hadoop в Greenplum





Что хочет клиент?

Не заботится об инфраструктурной составляющей – у меня должны появиться таблицы на моем кластере, про все остальное я не хочу знать

Работать с данными всех источников, различных стеков в онлайн режиме

Высокая скорость передачи данных (сравнимая с скоростью на самом источнике)

0 копий данных

Работать с понятными ему инструментами



Новая архитектура обмена данными между Hadoop кластерами



Для чего нужен продукт?



Защитить реплики и архивы данных от технологических угроз



Расширение технологий защиты данных, включая анонимизацию, шифрование и токенизацию, без доработок на источниках данных



Доступ к метаданным Hive-источников в онлайн



Централизованный аудит запросов к данными



Статистика обращений к данным



Гетерогенный доступ к различным типам источников (GP, Hadoop)



Основные вызовы



Сервисы должны быть доступны 99,95% – иначе кластер не работает



Обновление сервисов без недоступности



Большие задержки (>30%) вносимые проксированием недопустимы



Горизонтальное масштабирование, всплески нагрузки



Совместимость между различными версиями 3.x, Greenplum



Сетевая изоляция источников

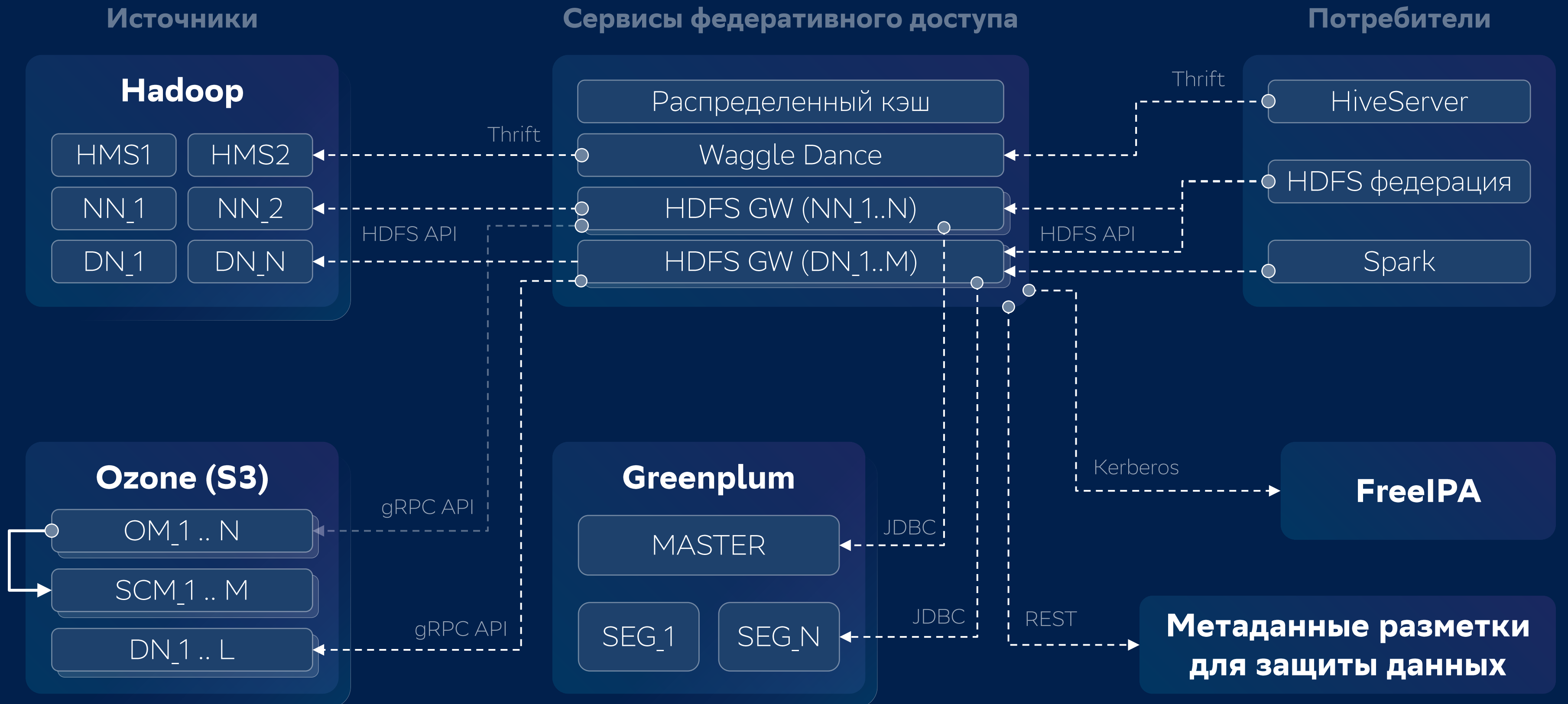


Pushdown на источнике



JOIN-ы работают со скоростью, сравнимой с локальным источником

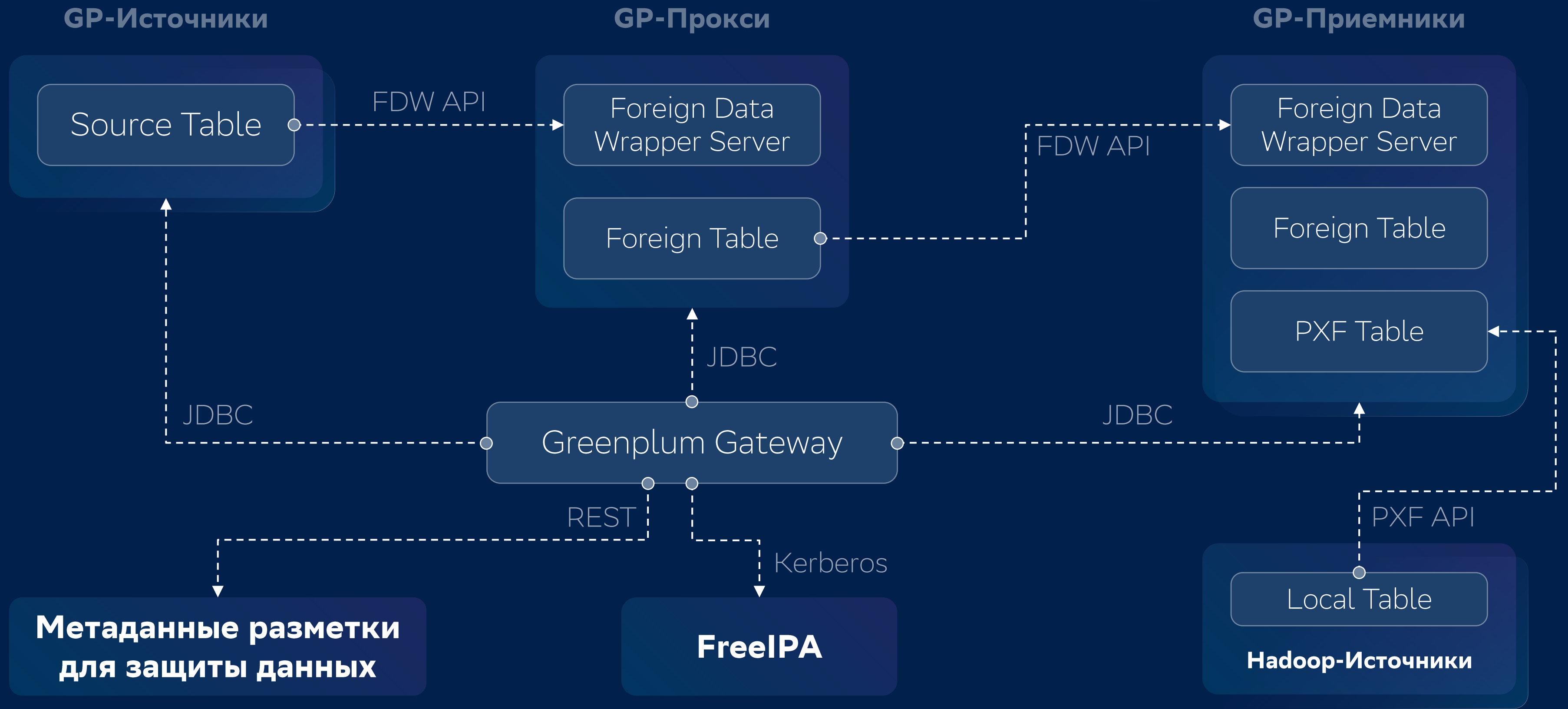
Архитектура HDP-HDP, OZONE-HDP, GP-HDP



Архитектура трактов GP-GP и HDP-GP



С УЧЕТОМ PXF/FDW



Статистика SDF



 КОЛИЧЕСТВО ПРОКСИ-ПОДПИСОК



Статистика SDF



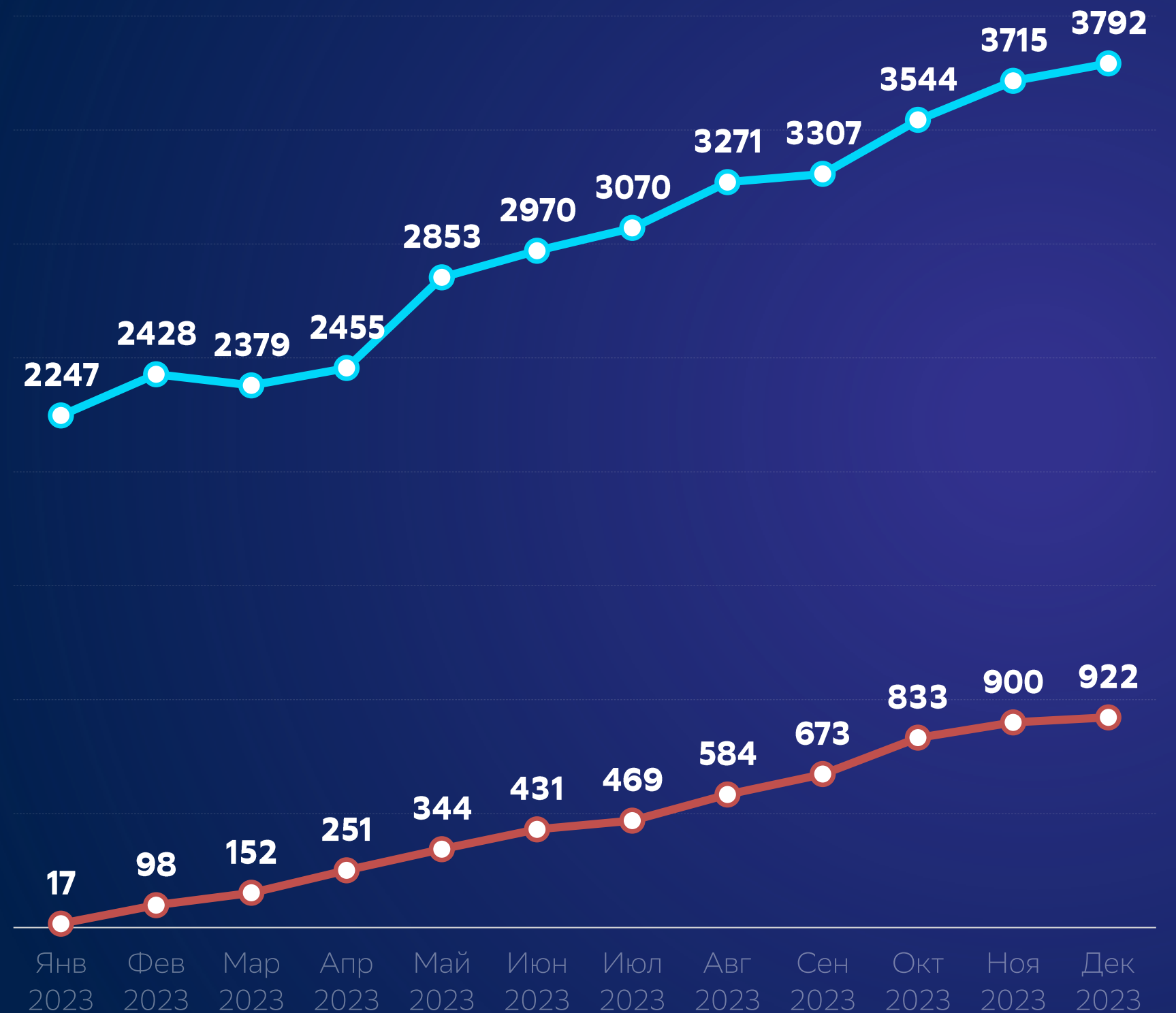
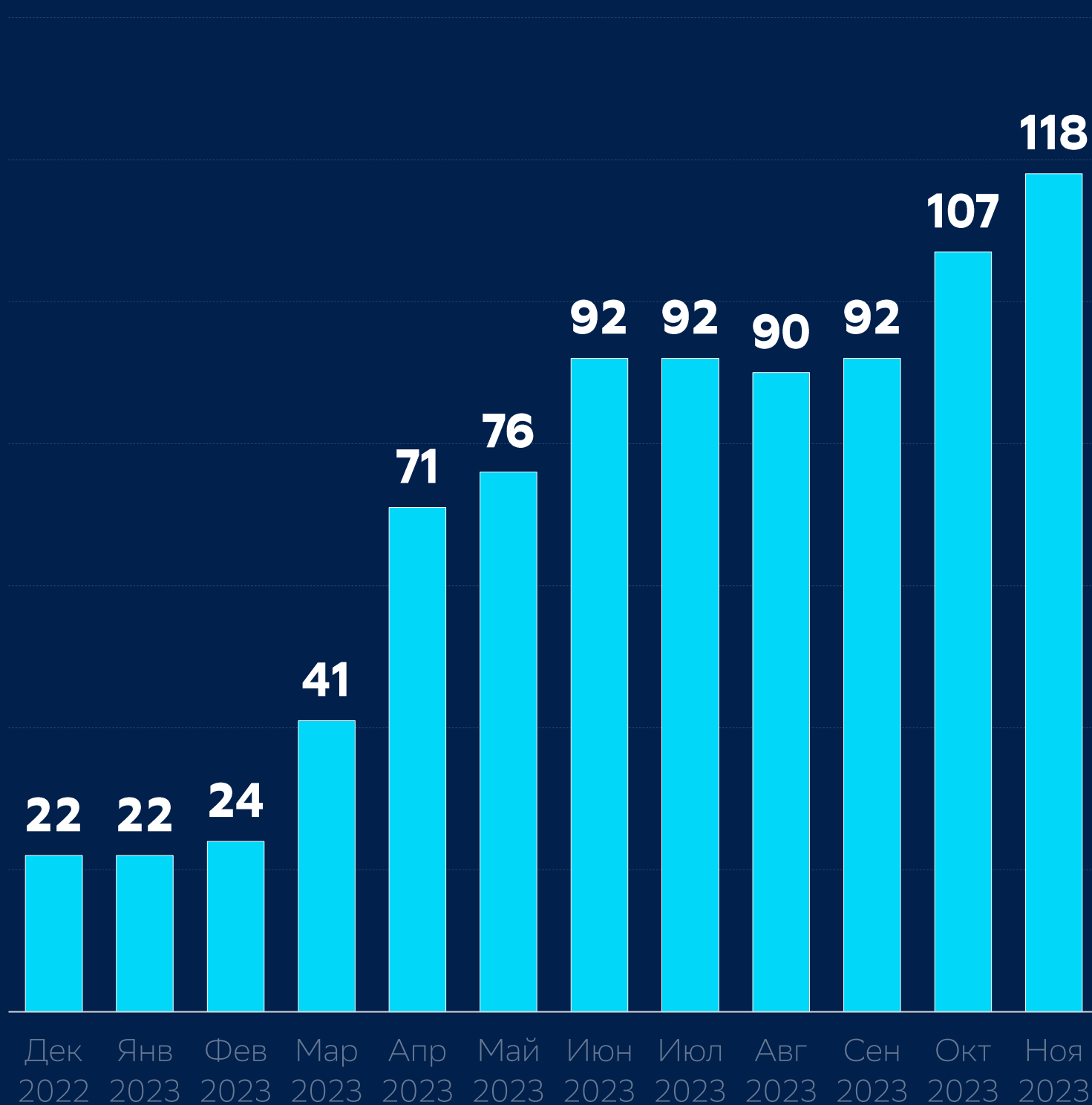
ТРАФИК, ПБ/МЕС

HDP-HDP

MONTHLY ACTIVE USERS

HDP-HDP

GP-GP



Спасибо за внимание!



SberDataFusion

Андрей Ильин

Алексей Тютякин

ilyin.a.vya@sberbank.ru

aatyutyakin@sberbank.ru