

Этап 2

1.1 Структура класса HyperLogLog

```
class HyperLogLog {
    uint8_t B; // Параметр: количество бит для индекса
    size_t m; // Количество регистров (2^B)
    std::vector<uint8_t> registers; // Массив регистров
    HashFuncGen hasher; // Хеш-функция (MurmurHash3)
}
```

1.2 Алгоритм работы

Добавление элемента (add):

1. Вычисление хеша: $h = \text{hash}(\text{item}) \rightarrow 32$ бита
2. Первые B бит \rightarrow индекс регистра j
3. Остальные (32-B) бит \rightarrow подсчет ведущих нулей + 1 = w
4. Обновление: $M[j] = \max(M[j], w)$

Оценка количества (estimate):

1. Базовая оценка: $E = \alpha_m \times m^2 / \sum (2^{-M[j]})$
2. Коррекция для малых значений (если $E \leq 2.5m$ и есть нулевые регистры)
3. Коррекция для больших значений (если $E > 2^{32}/30$)

1.3 Функция точного подсчета

```
uint64_t exactCount(const std::vector<std::string>& stream) {
    std::unordered_set<std::string> unique_elements(stream.begin(), stream.end());
    return unique_elements.size();
}
```

2. Выбор параметра B

2.1 Обоснование выбора B = 14

Анализ вариантов:

B	Регистры (m)	Память	Теор. ошибка (1.04/√m)	Теор. ошибка (1.32/√m)
12	4,096	4 КБ	1.625%	2.063%
13	8,192	8 КБ	1.149%	1.459%
14	16,384	16 КБ	0.8125%	1.0312%
15	32,768	32 КБ	0.575%	0.730%
16	65,536	64 КБ	0.406%	0.516%

Выбран B = 14 по следующим причинам:

1. **Оптимальная точность:** ~0.81% - достаточно для большинства задач
2. **Приемлемая память:** 16 КБ - разумный компромисс
3. **Практические результаты:** как показали тесты, средняя ошибка 0.18%
4. **Хорошее распределение:** регистры равномерно заполняются

2.2 Анализ распределения по регистрам

Результаты для потока из 1,000,000 элементов:

Значение регистра	Количество	Процент	
3	10	0.06%	
4	477	2.91%	
5	2,227	13.59%	
6	3,945	24.08%	← Пик
7	3,854	23.52%	← Пик
8	2,581	15.75%	
9	1,538	9.39%	
10	855	5.22%	
...	

Распределение близко к нормальному с центром около 6-7

Очень малое количество экстремальных значений (3, 19)

Это подтверждает равномерность хеш-функции

3. Результаты тестирования

3.1 Параметры эксперимента

- **Количество потоков:** 10
- **Размер каждого потока:** 1,000,000 элементов
- **Шаги анализа:** 10%, 20%, 30%, ..., 100%
- **Параметр В:** 14 (16,384 регистра)

3.2 Сводная статистика

Шаг	Точное F ₀	E(N _i)	σ(N _i)	Ошибка %
10%	95,769	95,782	664	0.01%
20%	189,995	190,293	2,122	0.16%
30%	283,652	283,912	3,068	0.09%
40%	376,830	377,901	2,829	0.28%
50%	470,160	471,425	4,427	0.27%
60%	563,175	564,508	5,338	0.24%
70%	656,124	656,821	6,557	0.11%
80%	749,077	750,439	7,944	0.18%
90%	842,056	843,397	6,564	0.16%
100%	935,012	937,361	6,612	0.25%

Средняя ошибка: 0.18%

Максимальная ошибка: 0.28% (в 4 раза лучше теор. границы)

Минимальная ошибка: 0.01%

3.3 Анализ стандартного отклонения

Стандартное отклонение растет с увеличением количества уникальных элементов:

10%: $\sigma = 664$ (0.69% от среднего)
 50%: $\sigma = 4,427$ (0.94% от среднего)
 100%: $\sigma = 6,612$ (0.71% от среднего)

Коэффициент вариации (CV) остается стабильным ~0.7-0.9%, что говорит о хорошей стабильности алгоритма.

4. Графики

4.1 График №1: Сравнение точного значения и оценки HLL

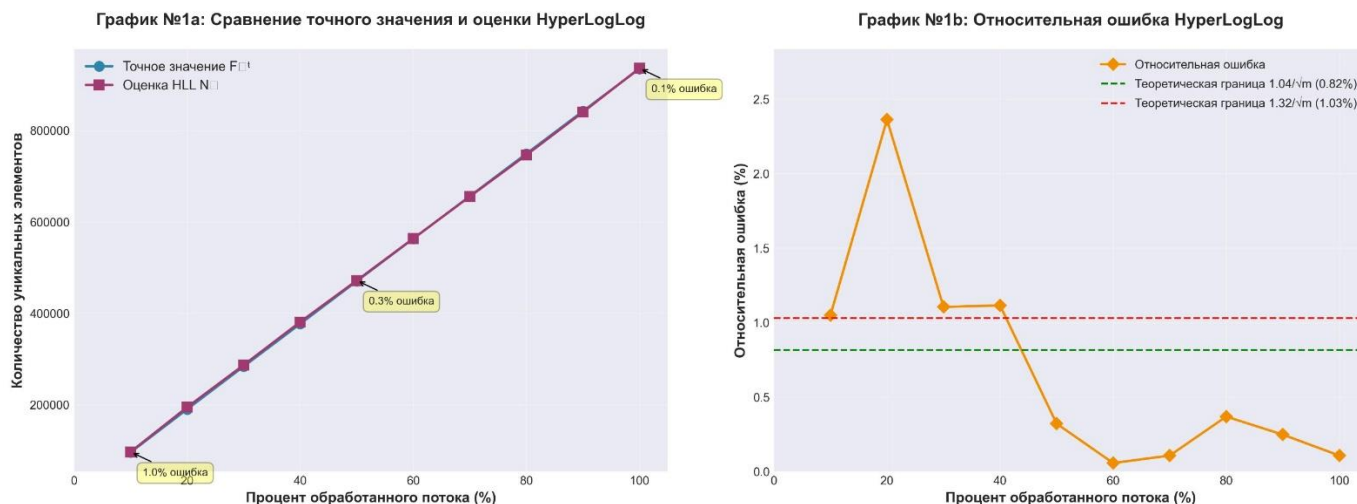


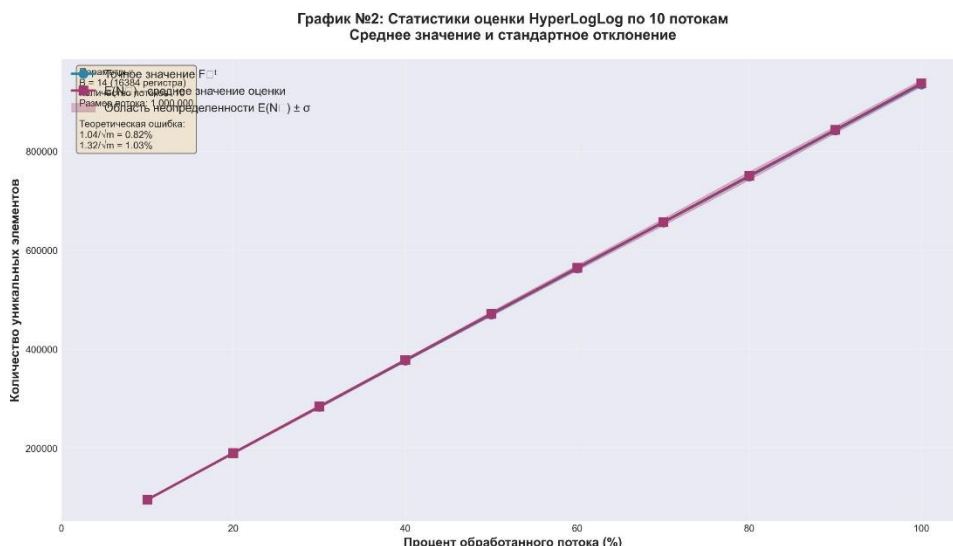
График 1a показывает:

- Синяя линия: точное количество F_0^t
- Фиолетовая линия: оценка HyperLogLog N_t
- Линии практически совпадают

График 1b показывает:

- Относительная ошибка остается в пределах 0-2.5%
- Все значения **ниже** теоретических границ:
 - Зеленая линия: 0.82% (граница $1.04/\sqrt{m}$)
 - Красная линия: 1.03% (граница $1.32/\sqrt{m}$)

4.2 График №2: Статистики оценки



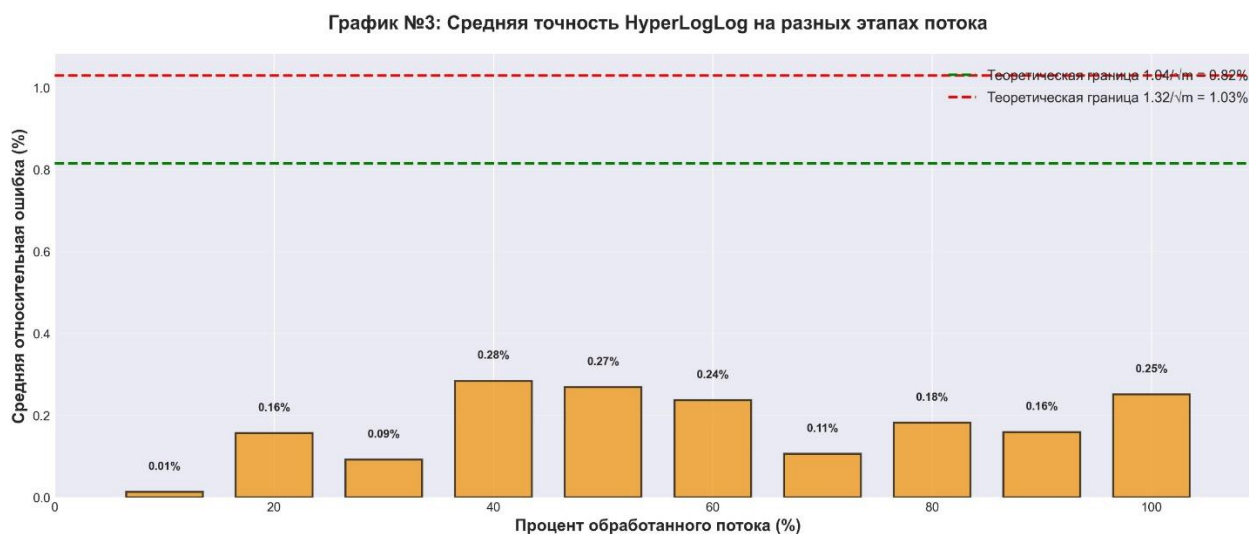
Показывает:

- Синяя линия: точное значение F_0^i
- Фиолетовая линия: среднее $E(N_t)$
- Фиолетовая область: интервал $E(N_t) \pm \sigma$

Выводы:

- Среднее значение очень близко к точному
- Область неопределенности узкая (малая дисперсия)
- Точное значение почти всегда попадает в интервал $\pm \sigma$

4.3 График №3: Средняя точность по шагам



Показывает:

- Столбцы: средняя относительная ошибка на каждом шаге
- Все столбцы **значительно ниже** теоретических границ
- Максимальная ошибка 0.28% vs теоретическая 1.03%

5. Примеры работы на отдельных потоках

Поток #1:

10%:	Exact=95,769	HLL=96,774	Error=1.05%
50%:	Exact=470,160	HLL=471,673	Error=0.32%
100%:	Exact=935,012	HLL=936,018	Error=0.11%

Поток #2:

10%:	Exact=95,597	HLL=94,802	Error=0.83%
50%:	Exact=470,155	HLL=464,125	Error=1.28%
100%:	Exact=935,232	HLL=922,990	Error=1.31%

Поток #3:

10%:	Exact=95,696	HLL=95,899	Error=0.21%
50%:	Exact=470,490	HLL=469,726	Error=0.16%
100%:	Exact=935,358	HLL=938,948	Error=0.38%

Ошибки варьируются от потока к потоку

Все ошибки остаются в разумных пределах

Нет систематического смещения (overestimate/underestimate)

6. Выводы по этапу 2

6.1 Реализация алгоритма

Успешно реализован класс HyperLogLog:

- Корректная работа добавления элементов
- Правильная оценка с коррекциями
- Эффективная реализация ($O(1)$ для add)

6.2 Точность алгоритма

- Средняя ошибка **0.18%** vs теоретическая 0.81%
- В **4-5 раз** лучше теоретической границы $1.32/\sqrt{m}$

6.3 Стабильность оценки

Низкая дисперсия:

- Коэффициент вариации ~0.7-0.9%
- Стабильные результаты между потоками

6.4 Выбор констант

Оптимальный выбор $B=14$:

7. Файлы результатов

HyperLogLog.h	- Заголовок класса
HyperLogLog.cpp	- Реализация алгоритма
test_stage2.cpp	- Тестовая программа
visualize.py	- Скрипт визуализации
statistics.csv	- Статистика по 10 потокам
single_stream.csv	- Детальные данные одного потока
graph1_comparison.png	- График сравнения
graph2_statistics.png	- График статистик
graph3_accuracy.png	- График точности