

# Хеш-функция

Хеш-функция — это функция, преобразующая исходные битовые данные произвольной длины (прообраз) в битовую строку фиксированной длины (хеш-образ). Вычисляется за не более чем полиномиальное время.

---

## Что можно хешировать?

- Строки
  - Фотографии
  - Программное обеспечение (например, чексумма для проверки корректности установки)
- 

## Области применения

- Проверка целостности и подлинности сообщений
  - Защита паролей при аутентификации
  - Ускорение поиска данных (хеш-таблицы)
  - Распознавание вредоносных файлов (например, набор данных NSLR)
- 

## Криптографические хеш-функции

Криптографические хеш-функции — это подкласс хеш-функций, обладающий особыми свойствами, необходимыми для использования в криптографии. При разработке современного российского стандарта [ГОСТ Р 34.11-2012 «Стрибог»](#) к ним предъявляются следующие требования:

### 1. Стойкость к вычислению прообраза

Если известно значение хеша  $h$ , должно быть **трудно** найти какое-либо сообщение  $m$ , такое что

$$\text{hash}(m) = h.$$

## 2. Стойкость ко второму прообразу

Если известно сообщение  $m_1$  и его хеш  $\text{hash}(m_1)$ , должно быть **трудно** найти другое сообщение  $m_2 \neq m_1$ , такое что

$$\text{hash}(m_2) = \text{hash}(m_1).$$

## 3. Стойкость к поиску коллизий

Должно быть **трудно** найти **любую** пару различных сообщений  $m \neq m'$ , для которых

$$H(m) = H(m').$$

### Пример коллизии

Пусть длина прообраза — 6 бит, длина хеша — 4 бита. Тогда:

- Всего возможных значений хеша:  $2^4 = 16$
- Всего возможных прообразов:  $2^6 = 64$

Так как  $64 > 16$ , по принципу Дирихле найдутся коллизии: некоторые хеш-значения будут соответствовать сразу нескольким (в нашем примере — четырём) прообразам.

## 4. Стойкость к удлинению прообраза

Если злоумышленник знает только длину сообщения и его хеш, ему должно быть **трудно** подобрать дописку к оригиналу так, чтобы хеш от изменённого сообщения имел **предопределённое** значение. Иными словами, хеш-функция не должна быть «дополняема» в контролируемый способ.

---

# Поиск пароля по хэш-функции с использованием цепочек и функции редукции

## Формальная постановка задачи

Пусть:

- $H$  — хеш-функция, выдающая значения длины  $n$
- $P$  — конечное множество возможных паролей

**Задача:** создать структуру данных, которая для любого значения хеша  $h$  сможет:

- найти такой элемент  $p \in P$ , что  $H(p) = h$ , или

- определить, что такого элемента не существует

**Проблема:** прямой перебор всех  $p \in P$  и вычисление  $H(p)$  слишком трудоёмкий.

## Функция редукции

Вводится функция редукции  $R$ :

$R : H \rightarrow P$ , отображающая в большинстве стандартных библиотек (C++, Rust, Go) используется `ByteBuffer` открытая адресация `ByteBuffer` значения хеша обратно в пространство паролей.

## Цепочка редукций и хешей

Для построения структуры:

- выбирается множество начальных паролей
- для каждого из них строится цепочка чередующихся операций хеширования и редукции:

$$aaaaaa \xrightarrow{H} 281DAF40 \xrightarrow{R} sgfnvd \xrightarrow{H} 920ECF10 \xrightarrow{R} kiebgt$$

**Сохраняются:**

- начальный пароль
- конечное значение цепочки

## Пример поиска пароля

Имеем хеш: `920ECF10`

Находим цепочку, в конце которой получается  $R(920ECF10) = kiebgt$ .

Находим цепочку, заканчивающуюся на `kiebgt`, и восстанавливаем её сначала:

$$aaaaaa \rightarrow H281DAF40 \rightarrow Rsgfnvd \rightarrow H920ECF10$$

Значит, искомый пароль — `sgfnvd`

## Алгоритм поиска пароля

1. Построить цепочки из случайного набора начальных паролей
2. Для каждого заданного хеша `h`:
  - Последовательно вычислять значения  $R(H(R(h)))$ ,  $R(H(R(H(R(h)))))$ , и т.д.
  - Если найденное значение совпадает с концом какой-либо цепочки:

- Берём начало этой цепочки и восстанавливаем её, пока не встретим *hash*
- Предшествующий элемент — искомый пароль

## Проблемы

- **Подбор функции редукции R:**
  - должна равномерно распределять значения в пространстве паролей
  - должна покрывать достаточное подмножество паролей
- **Слияние цепочек:**
  - при коллизиях H и R цепочки могут сливаться, теряя уникальность и ухудшая покрытие

## Взлом с помощью радужных таблиц

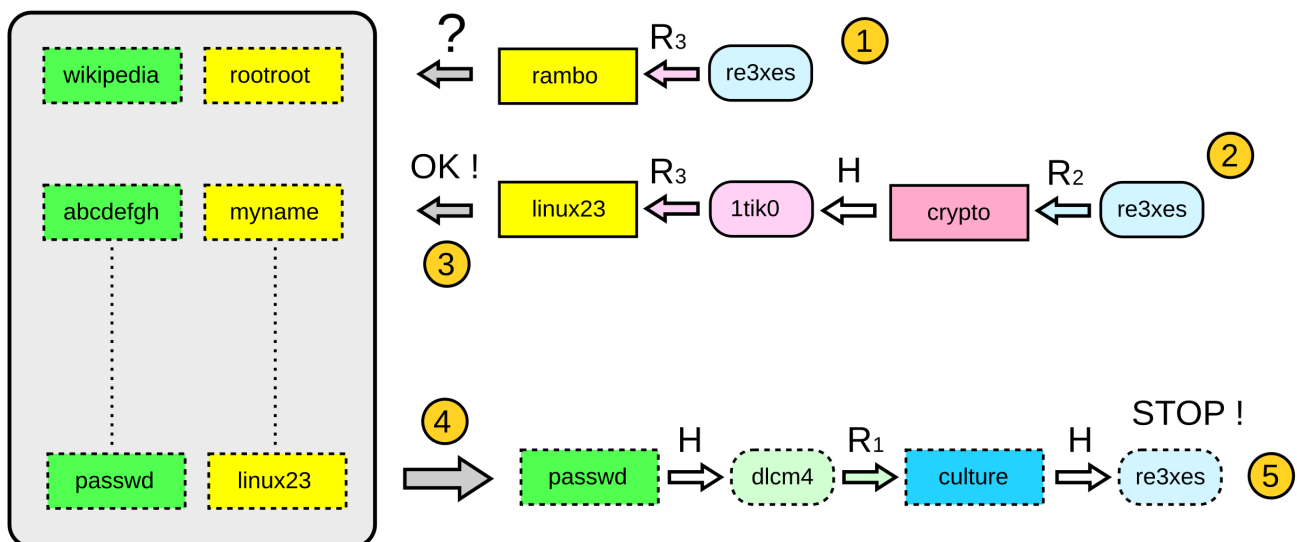
Предположим, что у нас имеется не одна функция редукции  $R : H \rightarrow P$ , а целый набор  $R_k$ . Для каждого пароля мы строим **цепочку** из чередующихся операций хеширования и редукции:

Шаг 1:  $R_k(H)$

Шаг 2:  $R_{k-1}(H(R_k(H(password))))$

Шаг k:  $R_1(H(R_2(\dots(R_k(H(password))\dots)))$

## Поиск исходного пароля в радужной таблице



Если на каком-то шаге результат применения функции редукции совпал с **хвостом** одной из цепочек в радужной таблице, мы переключаемся на *обратный* поиск: продолжаем вычислять пары «хеш → редукция», пока не получим исходный хеш. Формально, мы ищем индекс  $ii$ , такой что

$$H(R_i(\dots H(R_1(H(password)))))) = H_{\text{исходный}}$$

и берём пароль, предшествующий этому хешу (в примере выше — `culture`).

## Криптографическая соль: главное средство защиты

Радужные таблицы становятся практически бесполезными, если вместе с паролем хешируется уникальная для каждого пользователя **соль**:

$$H = H(password + salt)$$

Из-за соли одно и то же слово порождает разные хеши, а предрассчитанную таблицу придётся строить отдельно для каждой соли, что экономически невыгодно.

## Преимущества радужных таблиц

- По времени генерация длиннее, чем простая таблица «хеш → пароль», но **памяти** требуется существенно меньше: с  $O(N)$  до  $O(N^{2/3})$

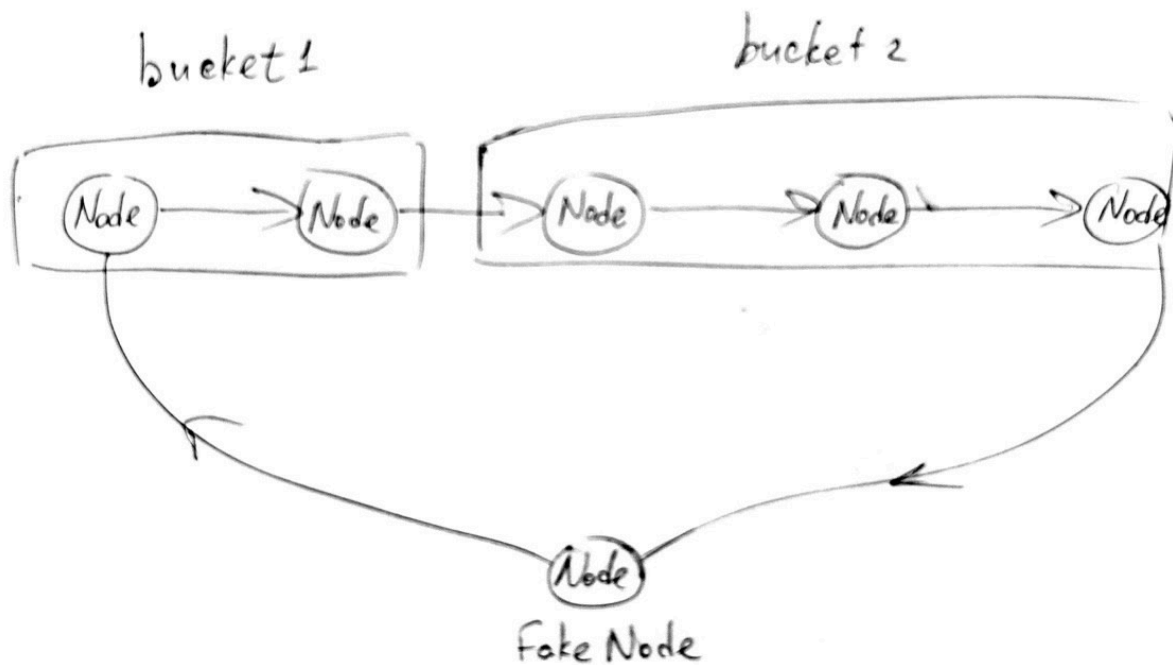
## Недостатки

- Таблица применима **только** к тому алгоритму хеширования, для которого была создана.
- Эффективность быстро падает, если в пароле используется соль или реpper.

---

## Хэш-таблица

### Операции вставки и удаления (метод цепочек)



Базовая реализация использует **односвязный список** ( `forward_list` ) из *бакетов* (цепочек). Алгоритм вставки пары (key, value) выглядит так:

1. Вычислить `hash(key)` .
2. Индекс бакета — `hash % N` .
3. Вставить пару в начало или в конец выбранной цепочки.

```
struct Node {  
    std::pair<const Key, Value> kv;  
    std::size_t hash; // храним хеш, чтобы при коллизии быстро сравнивать  
};
```

При коллизии (несколько ключей с одинаковым индексом) новые элементы просто добавляются в ту же цепочку. Для ускорения вставки полезно хранить указатель на её последний элемент.

## Удаление

По хешу попадаем в нужный бакет, находим элемент, *перешиваем* ссылки и при необходимости обновляем указатель на хвост цепочки.

## Перестройка таблицы

Когда коэффициент заполнения `load_factor` достигает `max_load_factor`, таблица **перехешируется** — размер массива бакетов увеличивается, и для всех ключей пересчитываются индексы.

**Сложность операций:** средняя  $O(1)$ , худшая  $O(n)$  при длинной цепочке.

## Открытая адресация

### Линейное, квадратичное и двойное хеширование

Вместо цепочек все элементы хранятся **внутри массива**. При коллизии выполняется последовательность проб:

$$h_0(x), h_1(x), \dots, h_n(x)$$

- **Линейное пробирование:**  $h_i(x) = (h(x) + i) \bmod N$
- **Квадратичное:**  $h_i(x) = (h(x) + c \cdot i + c \cdot i^2) \bmod N$
- **Двойное хеширование:**  $h_i(x) = (h(x) + i \cdot h_{new}(x)) \bmod N$

### Удаление (tombstone)

Чтобы не разрывать цепочку проб, удалённая ячейка помечается как *занятая-ранее* (tombstone).

- При **поиске** такие ячейки считаются занятыми.
- При **вставке** первая томб-ячейка на маршруте используется под новый элемент.

## Плюсы и минусы двух методов разрешения коллизий

Критерий	Метод цепочек	Открытая адресация
Память	Требуется указатель на следующую запись; массив + список	Всё в одном массиве, без указателей

Критерий	Метод цепочек	Открытая адресация
Производительность при низком <code>load_factor</code>	$O(1)$ вставка/поиск	$O(1)$ вставка/поиск
Производительность при высоком <code>load_factor</code>	Цепочки растут $\rightarrow O(n)$	Увеличивается длина проб $\rightarrow O(n)$
Удаление	Простейшее, за $O(1)$	Требуются tombstones или перестройка
Поддержка итераторов	Легко и стабильно	Сложнее, перестройка может инвалидировать
Подходит для внешней памяти	Да (списки можно хранить на диске)	Обычно нет
Реорганизация (rehash)	Необязательна до разумного <code>load_factor</code>	Критична: поиски резко деградируют

**Выбор метода** зависит от требований к памяти, частоты вставок/удалений и предельного `load_factor`. В большинстве стандартных библиотек (C++, Rust, Go) используется *открытая адресация* либо hybrid-подход с короткими цепочками. Для неизменяемых структур или работы с внешней памятью удобнее классические цепочки.