redis-1.3.7

sds.h/c = simple dynamic string

adlist.h/c = adjacent list

dict.h/c = dictionary / hash table

zipmap =

ae.h/c
├── ae net.h/c
├── ac_epoll.h/c
├── ac_select.h/c
└── ac_kqueue.h/c

redis.h/c = server
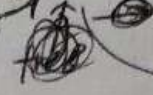
redis-di.c = client

struct sdshdr {
long len;        header
long free;
char buf[];
};

len

← compatible with c string

→ return to client

(free: free space in bytes)
(len: not including header)
and '\0'

typedef char* sds;

len

H /////////// '\0'

free

↓ makeroom
(len + extra) * 2

len

sdsprintf (sds s, const char* fmt, ...) (sds)

va_list ap;
char* buf, *t;
size_t buflen = 16;
while (1) {
    buf = malloc(buflen);
    buf[buflen-2] = '\0';
    va_start(ap, fmt);
    vsnprintf(buf, buflen, fmt, ap);
    va_end(ap);
    if (buf[buflen-2] != '\0') {
        free(buf);
        buflen *= 2;     // doubling
        continue;
    }
    break;
}
t = sdscat(s, buf);
free(buf);   return t;

判断是否溢出,
需要更大空间.

## adlist 双线链表

```
struct listNode {
    struct listNode* prev;
    struct listNode* next;
    void* value;
};

struct ListIter {
    listNode* next;
    int direction;
};

struct list {
    listNode* head;
    listNode* tail;
    unsigned int len;
};
```

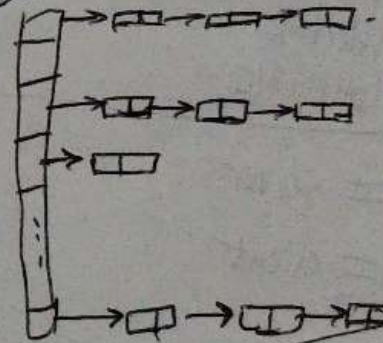## dict.h/c  Hash table

```
struct dictEntry {
    void* key;
    void* val;
    struct dictEntry* next;
};
```

```
struct dictType {
    uint32_t (*hashFuction) (const void* key)
    void* (*keyDup) (void* privdata,
                     const void* key);
    void* (*valDup) (void* privdata
                     const void* obj);
    int (*keyCompare) (void* privdata
                       const void* key,
                       const void* key2);
    void (*keyDestructor) (void* privdata,
                           void* key);
    void (*valDestructor) (void* privdata, void* obj);
};
```

```
struct dict {
    dictEntry ** table;
    dictType* type;
    unsigned long size;    表 slot 个数
    unsigned long sizemask;
    unsigned long used;    结点数
    void* privdata;
};
```

$ratio = \dfrac{used}{size} = 1$
$\Rightarrow$ expand.

```
struct dictIterator {
    dict* ht;
    int index;
    dictEntry* entry, *nextEntry;
};
```

这变量是为了存 entry→next值.
在 iterate 过程中 client 有可能
会 delete Node

生成事 iter
再由 iter 获取下一个 iter.

```
while (1) {
    if (iter→entry == Null) {
        iter→index++;
        if (iter→index >= iter→ht→size)
            break;
        iter→entry = iter→ht→table
                          [iter→index]
    } else {
        iter→entry = iter→nextEntry
    }
    if (iter→entry) {
        iter→nextEntry
            = iter→entry→next;
        return iter→entry;
    }
}
```

## aeEventLoop

```
struct aeEventLoop {
    int maxfd;
    long long timeEventNextId;
    aeFileEvent events[AE_SETSIZE];
    aeFiredEvent fired[AE_SETSIZE];
    aeTimeEvent * timeEventHead; // 链表 首
    int stop;
    void* apidata;
    aeBeforeSleepProc * beforeSleep;
}
```

```
aeMain (aeEventLoop* eventLoop)

eventLoop->stop=0;
while (!eventLoop->stop) {
    if (eventloop->beforeSleep != NULL)
        eventloop->beforeSleep(eventloop);
    aeProcessEvents (eventloop, AE_ALL_EVENTS);
}
```

aeEventLoop*
← aeCreateEventLoop → aeApiCreate(eventloop) ⇒ dispatch 到不同后端
                      ┌──────────┐            ┌──────┐
                      │ aeApiFree │            │ epoll │
                      └──────────┘            └──────┘

```
struct aeApiState {          ← malloc
    int epfd;        = epoll_create(...);
    struct epoll_event events[AE_SETSIZE].
};
```
                                    1024*10 = 10k ( hard coded )

```
aeCreateFileEvent (aeEventLoop* eventloop, int fd, int mask,
                   aeFileProc * proc, void* clientData)
    ↓
aeApiAddEvent ( fd, mask )
                      → Event handler
              → 事件类型
```

```
struct aeFileEvent {  // 事件类型 R/W.
    int mask;
    aeFileProc* rfileProc;
    aeFileProc* wfileProc;
    void* clientData;
};
```

SO_KEEPALIVE → SOL_SOCKET
TCP_NODELAY → IPPROTO_TCP
SO_SNDBUF → SOL_SOCKET

→ 根据 timer event 来设置 timeout, 调用 aeApiPoll { Select / poll / epoll } 都会 timeout 参数。没有 timer event, 则要 blocking wait, 要么 timeout=∞

O(n)[线性] 查找

→ 复制链表 [按照 时间排序。

→ 处理 fired events 不同 backend 将 就绪事件放入 fired events 中. (统一处理), 调用对应 process 处理.

**zip map**

String ⇒ String Map data structure optimized for size ($O(n)$)

在 key-value 比较少时生生 zipmap，较多时，用 hashtable (dict)。

- ZipmapNew ⇒

```
  1   1   (2字节)
┌───┬───┐
│ ∅ │   │
└───┴───┘
  status  ZIPMAP_END (255)
```

- 长度 < 253
```
1字节
┌─────┐  ┌──────────┐
│ len │  │ data key │
└─────┘  └──────────┘
```
```
         1字节↓
┌──────┐ ┌─────┬─────────┐
│▨▨▨▨▨│ │ len │  value  │
└──────┘ └─────┴─────────┘
```

- 长度 ≥ 253
```
1字节  4字节
┌─────┬──────┐  ┌──────────┐
│ 253 │ 长度  │  │   data   │
└─────┴──────┘  └──────────┘
```
```
      1字节  4字节
┌──────┐ ┌─────┬──────┐  ┌────────┐
│▨▨▨▨│ │     │ len  │  │ value  │
└──────┘ └─────┴──────┘  └────────┘
```

- mark = 254
```
┌──────┐ ┌─────────┐ ┌──────────┐
│ 254  │ │1字/5字  │ │ empty    │
│      │ │  len    │ │ ≤ len    │
└──────┘ └─────────┘ └──────────┘
```

(encoded Length + payload)       (encoded Length + single free count + payload)

(1 / 5 字节)

右上角注释：
在已有的 empty slot 上放入一个 key/value, (放新的).
找长度和总长度 ≥ 3 + 全 +1 = 5 ( key 1字, value 1字)
到作为一个 another empty slot 来处理.
重到 put 时 key/value, 空出 free byte
到 value 后面空字节.

右中注释：
(1字/5字) value长度
```
┌─────┐ ┌─────┐ ┌───────┐ ┌────┐
│ len'│ │free │ │ value │ │ ▨▨ │
└─────┘ └─────┘ └───────┘ └────┘
         ↑                  ↘ free space
    trade free space多余 free
         (1字节)
```

"foo" ⇒ "bar", "hello" ⇒ "world"

⟨status⟩ ⟨len⟩ "foo" ⟨len⟩ ⟨free⟩ "bar" ⟨len⟩ "hello" ⟨len⟩ ⟨free⟩ "world"

在改 ⟨free⟩ 是 1字节, 对 key/value 执行 update, value 可接受并, 记录有多少 free bytes
可以对以下次再 update 时长度够的 value.

"\x00|\x03 foo|\x03|\x00bar|\x05hello|\x05|\x00 world|\xff"
                                      ‾‾‾‾‾         ‾‾‾‾
status          free 为空     free 为空      ⇒ end

HMP全 key/value 记录, 或留有 empty slot.
```
          ← empty slot →        len
┌──────┐ ┌──────┐ ┌────────┐    ↓      ┌──────┬──────┐  对 len 从1两后移
│ 254  │ │ len  │ │▨▨▨▨▨▨│         │ 1字 │ 5字 │ ‾‾‾‾‾‾‾‾
└──────┘ └──────┘ └────────┘         └──────┴──────┘
  ↑p      当长度 (len▨▨len)
┌──────┐
│ 254  │ 标志 slot 是 empty 的.
└──────┘
```