

# 短网址

短网址 (Short URL) , 是在形式上比较短的网址, 通过映射关系跳转到原有的长网址。

目前已经有许多类似服务, 借助短网址您可以用简短的网址替代原来冗长的网址, 让使用者可以更容易的分享链接。

例如:

短网址: <http://t.cn/R3Krh36>

长网址: <https://blog.mimvp.com/article/25420.html>

生成短网址: <http://dwz.wailian.work/index.php>

看过新浪的短连接服务, 发现后面主要有6个字符串组成, 于是第一个想到的就是原来公司写的一个游戏激活码规则, 也就是下面的算法2,

1) 26个大写字母 26个小写字母, 10个数字, 随机生成6个然后插入数据库对应一个id,

2) 短连接跳转的时候, 根据字符串查询到对应id, 即可实现相应的跳转

62种字符组合成6位字符,  $62^6=568$ 亿个组合数量, 重复的概率是很小的

短链接的好处

- 1、内容需要;
- 2、用户友好;
- 3、便于管理。

为什么要这样做的, 原因我想有这样几点:

1) 微博限制一条字数为140字, 那么如果我们需要发一些连接上去, 但是这个连接非常的长, 以至于将近要占用我们内容的一半篇幅, 这肯定是不能被允许的, 所以短网址应运而生了。

2) 短网址在项目里可以很好的对开放级URL进行管理。有一部分网址可以会涵盖暴力、广告等信息, 这样我们可以通过用户的举报, 完全管理这些链接不出现我们的应用中。因为同样的URL通过加密算法之后, 得到的地址是一样的。

3) 我们可以对一系列的网址进行流量, 点击等统计, 挖掘出大多数用户的关注点, 这样有利于我们对项目的后续工作更好的作出决策。

算法原理

- 算法一

1) 将长网址md5生成32位签名串, 分为4段, 每段8个字节; 52c06085 c4529732 5433e0c7 5b140565

2) 对这4段循环处理, 取8个字节, 将他看成16进制串与0x3fffffff(30位1)与操作, 即前缀超过30位的字符串做忽略处理, 直接舍弃掉了;

3) 这30位分成6段, 每5位的数字作为字母表的索引取得特定字符, 依次进行获得6位字符串;

4) 总的md5串可以获得4个6位串, 取里面的任意一个就可作为这个长url的短url地址;

这种算法, 虽然会生成4个, 但是仍然存在重复几率

PHP 改进型

把固定长度6位, 改进成动态调整的, 如5、6、10、15位等, 使其是30个质数之一

- 算法二

a-zA-Z0-9 这62位取6位组合, 可产生 $62^6=568$ 亿个组合数量, 把数字和字符组合做一定的映射, 就可以产生唯一的字符串, 如第62个组合就是aaaaa9, 第63个组合就是aaaaba, 再利用洗牌算法, 把原字符串打乱后保存, 那么对应位置的组合字符串就会是无序的组合。

把长网址存入数据库, 取返回的id, 找出对应的字符串, 例如返回ID为1, 那么对应上面的字符串组合就是bbb, 同理 ID为2时, 字符串组合为bba, 依次类推, 直至到达62种组合后才会出现重复的可能, 所以如果用上面的62个字符, 任意取6个字符组合成字符串的话, 你的数据存量达到500多亿后才会出现重复的可能。具体参看这里彻底完善新浪微博接口和超短URL算法, 算法四可以算是此算法的一种实现, 此算法一般不会重复, 但是如果是统计的话, 就有很大的问题, 特别是对域名相关的统计, 就抓瞎了。

原理: 指定长度, 做多次循环, 每次从长字符串里随机取出一位字符, 组合成指定长度字符串即可

# 跳转原理

---

当我们生成短链接之后，只需要在数据库MySQL 或 NoSQL表中，存储原始链接与短链接的映射关系即可。

当我们访问短链接时，只需要从映射关系中找到原始链接，即可跳转到原始链接。

例如：

短网址：<http://t.cn/R3Krh36> （存储长链接的映射）

长网址：<https://blog.mimvp.com/article/25420.html>