

npc_hp110的专栏

目录视图 摘要视图 RSS 订阅

个人资料



内胚层

访问: 73次
积分: 10
等级:
排名: 千里之外
原创: 1篇 转载: 0篇
译文: 0篇 评论: 0条

文章搜索

文章存档

2015年08月 (1)

阅读排行

groupcache的设计和实现 (68)

评论排行

groupcache的设计和实现 (0)

推荐文章

- *Networking Named Content 全文翻译
- *边缘检测与图像分割
- *一次mysql慢查询事故分析
- *有关深度学习领域的几点想法
- *Java经典设计模式之七大结构型模式（附实例和详解）
- *网络性能评价方法

【征文】Hadoop十周年特别策划——我与Hadoop不得不说的故事 前端精品课程免费看，写课评赢心动大礼！

groupcache的设计和实现分析

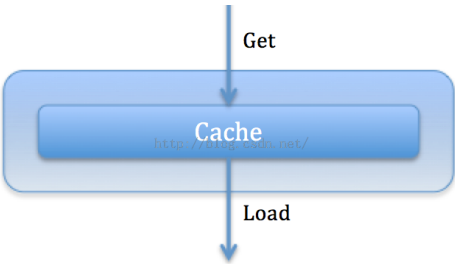
2015-08-27 22:19 70人阅读 评论(0) 收藏 举报

版权声明：本文为博主原创文章，未经博主允许不得转载。

groupcache的设计和实现分析

本文基于groupcache源码，分析分布式缓存系统的设计和实现过程。本文代码大部分是来自groupcache的源码，但根据分析的需要做了少许改动。

1.本地缓存系统

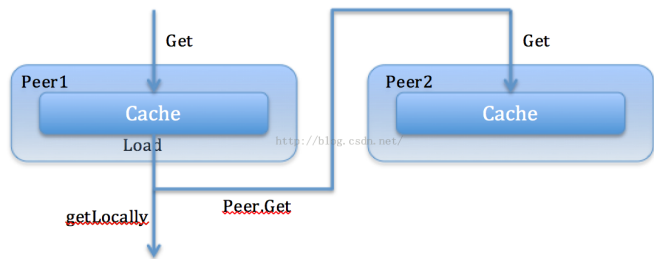


本地缓存系统的基本结构如上图所示。在内存中维护一个cache。查询时，首先查询cache中是否已经缓存查询结果。如果已经缓存，直接返回缓存结果，如果没有缓存，将查询结果Load到cache中，然后返回结果。代码如下：

```
[plain]
01. type Value interface{}
02. type Cache interface {
03.     Get(key string) (Value, bool)
04.     Add(key string, val Value)
05. }
06.
07. type Group struct {
08.     cache Cache
09. }
10.
11. func (g *Group) Get(key string) Value {
12.     //look up cache first
13.     val, cacheHit := g.lookupCache(key)
14.     if cacheHit {
15.         return val
16.     }
17.     //if miss, load value to cache
18.     val = g.lo
19.     return val
20. }
21.
22. func (g *Group) lookupCache(key) (Value, bool) {
23.     val, ok := g.cache.Get(key)
24.     return val, ok
25. }
```

2. 分布式缓存系统

在设计分布式缓存系统的时候，需要让key分布在不同的缓存节点上。当某节点收到查询请求时，如果该key归属于本节点，则在本节点获取查询结果；如果该key归属于其他节点，则本节点向归属节点获取查询结果。如下图所示：



此分布式缓存架构引入了2个新的问题：

- 1. 如何判断查询的key归属哪个节点
- 2. 如何从其他节点获取数据

第1个问题实际上是一个路由问题，即给定key，路由到某个节点。这里忽略具体实现，将问题抽象出来，由2个接口表示（对应于groupcache的ProtoGetter和PeerPicker）：

```
[plain]
01. type Peer interface{
02.     Get(key string) Value
03. }
04.
05. type Router interface{
06.     Route(key string) Peer
07. }
```

接口Peer表示一个远端节点，Get方法从远端节点查询数据；接口Router表示一个路由器，Route方法将给定key路由到其归属的远端节点，如果key归属于本节点，Route方法返回nil。基于此，改写Group结构体和Group.load方法：

```
[plain]
01. type Group struct {
02.     router Router
03.     cache  Cache
04. }
05. func (g *Group) load(key string) Value {
06.     peer := g.router.Route(key)
07.     if peer == nil {
08.         val := g.getLocally(key)
09.         g.cache.Add(key, val) //store result in cache
10.         return val
11.     }
12.     return peer.Get(g.name, key)
13. }
```

至此，一个分布式缓存框架搭建完成了，接下来分析具体细节的实现和功能的优化完善。

2.1. Cache的内存结构

Cache对象是内存中的一个容器，用来存放查询的结果。我们将Cache对象实现为一个List。另外，为了提升访问效率，用一个map结构来索引key及其对应的值：

```
[plain]
01. type ListCache struct {
02.     MaxEntries int
03.     l list.List
04.     m map[string]*list.Element
05. }
```

指定LRU为缓存的淘汰策略：即当Cache满了需要淘汰数据时，优先淘汰最老的数据。我们约定，链表中越靠近表头数据越新，越靠近表尾数据越老。因此在访问cache时，需要把命中的数据移至表头：

[plain]

```
01. type entry struct {
02.     key string
03.     val Value
04. }
05. func (c *ListCache) Get(key string) (Value, bool) {
06.     if e, hit := c.m[key]; hit {
07.         c.l.MoveToFront(e)
08.         return e.Value.(*entry).val, hit
09.     }
10.     return nil, false
11. }
12. func (c *ListCache) Add(key string, val Value) {
13.     if e, ok := c.m[key]; ok {
14.         c.l.MoveToFront(e)
15.     } else {
16.         ee := c.l.PushFront(&entry{key, val})
17.         c.m[key] = ee
18.         if c.l.Len() > c.MaxEntries {
19.             oldest := c.l.Back()
20.             c.l.Remove(oldest)
21.             delete(c.m, oldest.Value.(*entry).key)
22.         }
23.     }
24. }
```

2.2. 路由器的实现——一致性哈希

本节节选自：[一致性哈希算法及其在分布式系统中的应用](#)

路由器的基本功能是将key分布到不同的节点。很多方法可以实现这一点，最常用的方法是计算哈希。例如对于每次访问，可以按如下算法计算其哈希值：

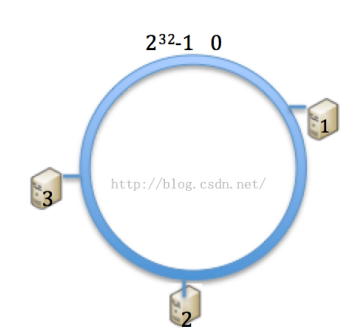
$$h = \text{Hash}(\text{key}) \% N$$

这个算式计算每个key应该被路由到哪个节点，其中N为节点总数，所有节点按照0 – (N-1)编号。

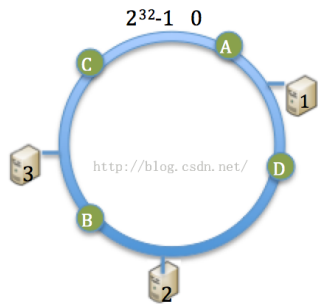
这个算法的问题在于容错性和扩展性不好。假设有一台服务器宕机了，那么为了填补空缺，要将宕机的服务器从编号列表中移除，后面的服务器按顺序前移一位并将其编号值减一，此时每个key就要按 $h = \text{Hash}(\text{key}) \% (N-1)$ 重新计算；同样，如果新增了一台服务器，虽然原有服务器编号不用改变，但是要按 $h = \text{Hash}(\text{key}) \% (N+1)$ 重新计算哈希值。因此系统中一旦有服务器变更，大量的key会被重定位到不同的服务器从而造成大量的缓存不命中。

一致性哈希算法就是解决这个问题的一种哈希方案。

简单来说，一致性哈希将整个哈希值空间组织成一个虚拟的圆环，如假设某哈希函数H的值空间为0 - $2^{32}-1$ （即哈希值是一个32位无符号整形），整个空间按顺时针方向组织。0和 $2^{32}-1$ 在零点中方向重合。将各个节点使用H进行一个哈希，确定其在哈希环上的位置。假设有三台节点使用ip地址哈希后在环空间的位置如下：



接下来使用如下算法将数据访问路由到相应节点：将数据key使用相同的函数H计算出哈希值h，根据h确定此数据在环上的位置，从该位置沿环顺时针“行走”，第一台遇到的节点就是其应该路由到的节点。例如有A、B、C、D四个数据对象，经过哈希计算后，在环空间上的位置如下：



根据一致性哈希算法：数据A，C路由到节点1；数据D路由到节点2；数据B路由到节点3。

下面基于一致性哈希实现路由器。首先给每个节点命名：

```
[plain]
01. type Peer interface {
02.     Get(key string) Value
03.     Name() string
04. }
```

一致性哈希路由HashRouter实现了Router接口（参考groupcache的HTTPPool）：

```
[plain]
01. type Hash func(key string) uint32
02. type HashRouter struct {
03.     self      Peer
04.     replicas  int
05.     hash      Hash
06.     locations []int //Sorted
07.     hashMap   map[int]Peer //from location on the ring to Peer
08. }
09. func (r *HashRouter) Add(peers ...Peer) {
10.     for _, peer := range peers {
11.         for i := 0; i < r.replicas; i++ {
12.             hash := int(r.hash(peer.Name()) + strconv.Itoa(i))
13.             if _, ok := r.hashMap[hash]; !ok {
14.                 r.locations = append(r.locations, hash)
15.             }
16.             r.hashMap[hash] = peer
17.         }
18.     }
19.     sort.Ints(r.locations)
20. }
21. func (r *HashRouter) Route(key string) Peer {
22.     hash := int(r.hash(key))
23.     // Binary search for appropriate replica.
24.     idx := sort.Search(
25.         len(r.locations),
26.         func(i int) bool { return r.locations[i] >= hash },
27.     )
28.     if idx == len(r.locations) {
29.         idx = 0
30.     }
31.     peer := r.hashMap[r.locations[idx]]
32.     if peer == r.self {
33.         return nil
34.     }
35.     return peer
36. }
```

其中HashRouter的locations是一个有序切片，保存所有节点在哈希环上的位置。hashMap保存从哈希环上的位置到节点(Peer)的索引。

2.3. 从远端节点获取缓存数据

我们采用HTTP协议在不同节点间传输缓存数据。数据请求的协议为：GET http://peer/key。HTTPPeer实现了Peer接口：

```
[plain]
```

```

01. <pre name="code" class="plain">type encoder interface {
02.     Marshal(v interface{}) []byte
03.     Unmarshal(data []byte, v interface{})
04. }
05. type HTTPPeer struct {
06.     baseURL string
07.     enc      encoder
08. }
09. func (p *HTTPPeer) Name() string {
10.     return p.baseURL
11. }
12. func (p *HTTPPeer) Get(key string) Value {
13.     u := fmt.Sprintf("%v%v", p.baseURL, url.QueryEscape(key))
14.     resp, _ := http.Get(u)
15.     defer resp.Body.Close()
16.
17.     b := make([]byte, 0)
18.     resp.Body.Read(b)
19.     var val Value
20.     p.enc.Unmarshal(b, &val)
21.     return val
22. }

```

2.4. 处理远端节点的查询请求

每个节点既可以向远端节点请求数据，同时也要处理远端节点的查询请求。HTTPPeer负责处理节点间的数据传输，因此需要处理远端节点的查询请求。



[plain]

```

01. func (p *HTTPPeer) ServeHTTP(w http.ResponseWriter, r *http.Request) {
02.     key := strings.SplitN(r.URL.Path, "/", 1)[0]
03.     group := GetGroup()
04.     val := group.Get(key)
05.
06.     body := p.enc.Marshal(val)
07.     w.Write(body)
08. }

```

2.5. 数据编码

在2.3中我们定义了数据编码器，但未实现它。实际上golang有很多成熟的编码库，如json, gob, protobuf等。groupcache使用的是protobuf。

2.6. 数据Load

当查询数据不在缓存中时，需要从数据源加载数据。加载数据的过程需要用户自定义，因此定义Loader接口，用户自行实现该接口：



[plain]

```

01. type Loader interface {
02.     Load(key string) Value
03. }
04. type Group struct {
05.     loader Loader
06.     router Router
07.     cache  Cache
08. }
09.
10. func (g *Group) getLocally(key string) Value {
11.     return g.loader.Load(key)
12. }

```

2.6. Alltogether and Startup

要在一个节点上启动上述的分布式缓存系统，包含如下步骤：

1. 初始化一个Group对象，指定Loader
2. 初始化一个ListCache对象，并将其赋值给Group对象的cache

3. 初始化一个HashRouter对象，并将其赋值给Group对象的router
4. 为集群的每个节点初始化一个HTTPPeer对象，并将它们添加到HashRouter对象。将本节点对应的HTTPPeer对象赋值给HashRouter对象的self
5. http监听指定端口，并指定处理函数为HTTPPeer对象的ServeHTTP方法

可以做一些适当的封装，让此缓存系统使用更加方便。

2.7. 错误处理

3. 功能优化和完善

至此我们已经实现了一个简陋的分布式缓存系统，接下来我们为这个分布式缓存系统增加一些新的功能（groupcache支持的功能）。

3.1 命名空间

给系统增加命名空间，使不同命名空间的缓存互相独立。为此，给Group对象增加一个名称，每个节点可以创建多个不同名称的Group对象，每个对应一个命名空间。

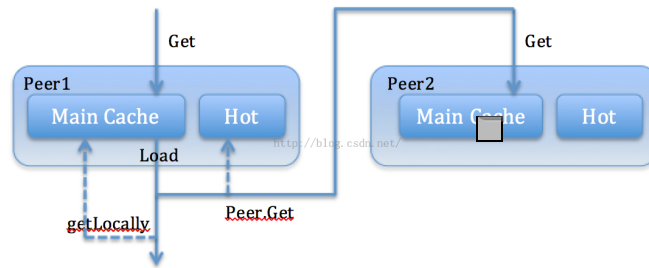
```
[html]
01. var groups = make(map[string]*Group)
02. func GetGroup(name string) *Group {
03.     return groups[name]
04. }
05. type Group struct {
06.     name string
07.     loader Loader
08.     router Router
09.     cache Cache
10. }
```

节点间的数据请求协议变为：GET http://peer:port/groupname/key。

```
[plain]
01. type Peer interface {
02.     Get(group, key string) Value
03.     Name() string
04. }
05. func (p *HTTPPeer) Get(group, key string) Value {
06.     u := fmt.Sprintf("%v/%v", p.baseURL, url.QueryUnescape(group), url.QueryEscape(key))
07.     ...
08. }
09. func (p *HTTPPeer) ServeHTTP(w http.ResponseWriter, r *http.Request) {
10.     parts := strings.SplitN(r.URL.Path, "/", 2)
11.     groupName := parts[0]
12.     key := parts[1]
13.     group := GetGroup(groupName)
14.     val := group.Get(key)
15.
16.     body := p.enc.Marshal(val)
17.     w.Write(body)
18. }
```

3.2 热数据扩散

前面提到，在设计分布式缓存系统的时候，需要让key分布在不同的缓存节点上。如果某些key的访问量特别大，大量请求会路由到该key归属的节点，可能导致该节点无法处理而瘫痪。因此groupcache增加了热数据自动扩散功能。通过在Group中增加一个Cache对象，如果节点的某些key访问特别频繁，而这些key的归属节点不在本节点，此时会将这些key的查询结果缓存到新增的Cache对象中。如下图所示：



代码如下:

```
[plain]
01. type Group struct {
02.     name      string
03.     loader    Loader
04.     router    Router
05.     mailCache Cache
06.     hotCache  Cache
07. }
08. func (g *Group) lookupCache(key string) (Value, bool) {
09.     val, ok := g.mailCache.Get(key)
10.     if ok {
11.         return val, ok
12.     }
13.     val, ok = g.hotCache.Get(key)
14.     return val, ok
15. }
16. func (g *Group) load(key string) Value {
17.     peer := g.router.Route(key)
18.     if peer == nil {
19.         val := g.getLocally(key)
20.         g.mailCache.Add(key, val) //store result in cache
21.         return val
22.     }
23.     val := peer.Get(g.name, key)
24.     if rand.Intn(10) == 0 {
25.         g.hotCache.Add(val)
26.     }
27.     return val
28. }
```

3.3 缓存Load过程合并机制，避免“惊群效应”

如果某节点对某相同的key存在大量并发查询，而该key的值不在缓存中，这些并发查询就会触发大量的Load过程（从数据源或远端节点加载数据）。但这些Load过程都是加载相同的数据，造成了资源的大量浪费。为了避免这种现象，需要在查询触发Load过程前，先判断是否已经有相同的Load过程正在运行。如果存在，本次Load不执行，而是等正在运行的Load过程完成后直接使用其结果。

```
[plain]
01. type loadResult struct {
02.     wg sync.WaitGroup
03.     val interface{}
04. }
05. type LoadGroup struct {
06.     mu sync.Mutex // protects m
07.     m map[string]*loadResult // map key to loading process
08. }
09. type loadFn func() Value
10. func (g *LoadGroup) Do(key string, load loadFn) Value {
11.     g.mu.Lock()
12.     if c, ok := g.m[key]; ok { // if it is loading
13.         g.mu.Unlock()
14.         c.wg.Wait() // wait loaded
15.         return c.val
16.     }
17.     c := new(loadResult)
18.     c.wg.Add(1)
19.     g.m[key] = c // notify that it is loading
20.     g.mu.Unlock()
21.
22.     c.val = load() //start loading
23.     c.wg.Done() //notify loaded
24. }
```

```
25.     g.mu.Lock()
26.     delete(g.m, key)
27.     g.mu.Unlock()
28.
29.     return c.val
30. }
```

上面的代码定义了一个LoadGroup对象，调用它的Do方法可以合并多个Load过程，从而避免“惊群效应”。

```
[plain]
01. type Group struct {
02.     name      string
03.     loader    Loader
04.     router    Router
05.     mailCache Cache
06.     hotCache  Cache
07.     loadGroup LoadGroup
08. }
09. func (g *Group) load(key string) Value {
10.     val := g.loadGroup.Do(key, func() Value {
11.         peer := g.router.Route(key)
12.         if peer == nil {
13.             val := g.getLocally(key)
14.             g.mailCache.Add(key, val) //store result in cache
15.             return val
16.         }
17.         val := peer.Get(g.name, key)
18.         if rand.Intn(10) == 0 {
19.             g.hotCache.Add(key, val)
20.         }
21.         return val
22.     })
23.     return val
24. }
```

3.4. 缓存容量控制

前面提到，Cache对象在内存中是以List的结构组织的。通过ListCache.MaxEntries字段控制缓存的容量。然而该字段只能限制缓存对象的数量，不能限制缓存的实际内存容量。要控制缓存的实际内存容量，需要知道缓存对象的大小。缓存的对象结构如下：

```
[plain]
01. type entry struct {
02.     key string
03.     val Value
04. }
```

key可以通过len(key)计算大小，要计算val的大小，需要 val对象实现一下接口：

```
[html]
01. type ByteView interface{
02.     Len() int
03. }
```

具体实现参考groupcache源码

3.5. 统计功能

参考groupcache源码

猜你在找

- 《C语言/C++学习指南》加密解密篇（安全相关算法）
- 大数据时代的<集装箱式>架构设计与Docker潮流
- Ceph—分布式存储系统的另一个选择
- 疯狂IOS讲义之Objective-C面向对象设计
- C语言系列之 字符串相关算法



短信接口



ios学习路线



30天试用云主



小米笔记本电



秋田犬售价



周杰伦演唱会



红袖女

查看评论

暂无评论

您还没有登录,请[\[登录\]](#)或[\[注册\]](#)

* 以上用户言论只代表其个人观点，不代表CSDN网站的观点或立场

核心技术类目

- 全部主题
- Hadoop AWS 移动游戏 Java Android iOS Swift 智能硬件 Docker OpenStack
- VPN Spark ERP IE10 Eclipse CRM JavaScript 数据库 Ubuntu NFC WAP jQuery
- BI HTML5 Spring Apache .NET API HTML SDK IIS Fedora XML LBS Unity
- Splashtop UML components Windows Mobile Rails QEMU KDE Cassandra CloudStack FTC
- coremail OPhone CouchBase 云计算 iOS6 Rackspace Web App SpringSide Maemo
- Compuware 大数据 aptech Perl Tornado Ruby Hibernate ThinkPHP HBase Pure Solr
- Angular Cloud Foundry Redis Scala Django Bootstrap