

Commit 958a2765 authored 4 years ago by  [Ilya Dryomov](#)

crush: straw2 bucket type with an efficient 64-bit crush_ln()

This is an improved straw bucket that correctly avoids any data movement between items A and B when neither A nor B's weights are changed. Said differently, if we adjust the weight of item C (including adding it anew or removing it completely), we will only see inputs move to or from C, never between other items in the bucket.

Notably, there is not intermediate scaling factor that needs to be calculated. The mapping function is a simple function of the item weights.

The below commits were squashed together into this one (mostly to avoid adding and then yanking a ~6000 lines worth of crush_ln_table):


- crush: add a straw2 bucket type
- **crush: add crush_ln to calculate nature log efficently**
- crush: improve straw2 adjustment slightly
- crush: change crush_ln to provide 32 more digits
- crush: fix crush_get_bucket_item_weight and bucket destroy for straw2
- crush/mapper: fix divide-by-0 in straw2
(with div64_s64() for draw = ln / w and INT64_MIN -> S64_MIN - need to create a proper compat.h in ceph.git)

Reflects ceph.git commits 242293c908e923d474910f2b8203fa3b41eb5a53, 32a1ead92efcd351822d22a5fc37d159c65c1338, 6289912418c4a3597a11778bcf29ed5415117ad9, 35fcb04e2945717cf5cfe150b9fa89cb3d2303a1, 6445d9ee7290938de1e4ee9563912a6ab6d8ee5f, b5921d55d16796e12d66ad2c4add7305f9ce2353.


Signed-off-by:  [Ilya Dryomov](#) <idryomov@gmail.com>

 parent [45002267](#)

 master 

 No related merge requests found

Showing 5 changed files ▾ with 316 additions and 2 deletions

▼  include/linux/crush/crush.h			
...	...	@@ -96,13 +96,15 @@ struct crush_rule {	
96	96	* uniform	O(1) poor poor
97	97	* list	O(n) optimal poor
98	98	* tree	O(log n) good good
99		- * straw	O(n) optimal optimal
	99	+ * straw	O(n) better better
	100	+ * straw2	O(n) optimal optimal
100	101	*/	
101	102	enum {	
102	103	CRUSH_BUCKET_UNIFORM = 1,	
103	104	CRUSH_BUCKET_LIST = 2,	
104	105	CRUSH_BUCKET_TREE = 3,	
105		- CRUSH_BUCKET_STRAW = 4	
	106	+ CRUSH_BUCKET_STRAW = 4,	
	107	+ CRUSH_BUCKET_STRAW2 = 5,	
106	108	};	
107	109	extern const char *crush_bucket_alg_name(int alg);	
108	110		
...	...	@@ -149,6 +151,11 @@ struct crush_bucket_straw {	
149	151	__u32 *straws; /* 16-bit fixed point */	
150	152	};	
151	153		
	154	+ struct crush_bucket_straw2 {	
	155	+ struct crush_bucket h;	
	156	+ __u32 *item_weights; /* 16-bit fixed point */	
	157	+ };	
	158	+	
152	159		

153	160	
154	161	/*
...	...	@@ -189,6 +196,7 @@ extern void crush_destroy_bucket_uniform(struct crush_bucket_uniform *b);
189	196	extern void crush_destroy_bucket_list(struct crush_bucket_list *b);
190	197	extern void crush_destroy_bucket_tree(struct crush_bucket_tree *b);
191	198	extern void crush_destroy_bucket_straw(struct crush_bucket_straw *b);
	199	+ extern void crush_destroy_bucket_straw2(struct crush_bucket_straw2 *b);
192	200	extern void crush_destroy_bucket(struct crush_bucket *b);
193	201	extern void crush_destroy_rule(struct crush_rule *r);
194	202	extern void crush_destroy(struct crush_map *map);
...	...	

▼ net/ceph/crush/crush.c

...	...	@@ -17,6 +17,7 @@ const char *crush_bucket_alg_name(int alg)
17	17	case CRUSH_BUCKET_LIST: return "list";
18	18	case CRUSH_BUCKET_TREE: return "tree";
19	19	case CRUSH_BUCKET_STRAW: return "straw";
	20	+ case CRUSH_BUCKET_STRAW2: return "straw2";
20	21	default: return "unknown";
21	22	}
22	23	}
...	...	@@ -40,6 +41,8 @@ int crush_get_bucket_item_weight(const struct crush_bucket *b, int p)
40	41	return ((struct crush_bucket_tree *)b)->node_weights[crush_calc_tree_node(p)];
41	42	case CRUSH_BUCKET_STRAW:
42	43	return ((struct crush_bucket_straw *)b)->item_weights[p];
	44	+ case CRUSH_BUCKET_STRAW2:
	45	+ return ((struct crush_bucket_straw2 *)b)->item_weights[p];
43	46	}
44	47	return 0;
45	48	}
...	...	@@ -77,6 +80,14 @@ void crush_destroy_bucket_straw(struct crush_bucket_straw *b)
77	80	kfree(b);
78	81	}
79	82	
	83	+ void crush_destroy_bucket_straw2(struct crush_bucket_straw2 *b)
	84	+ {
	85	+ kfree(b->item_weights);
	86	+ kfree(b->h.perm);
	87	+ kfree(b->h.items);
	88	+ kfree(b);
	89	+ }
	90	+ }
80	91	void crush_destroy_bucket(struct crush_bucket *b)
81	92	{
82	93	switch (b->alg) {
...	...	@@ -92,6 +103,9 @@ void crush_destroy_bucket(struct crush_bucket *b)
92	103	case CRUSH_BUCKET_STRAW:
93	104	crush_destroy_bucket_straw((struct crush_bucket_straw *)b);
94	105	break;
	106	+ case CRUSH_BUCKET_STRAW2:
	107	+ crush_destroy_bucket_straw2((struct crush_bucket_straw2 *)b);
	108	+ break;
95	109	}
96	110	}
97	111	
...	...	

► net/ceph/crush/crush ln table.h0 → 100644

This diff is collapsed.

▼ net/ceph/crush/mapper.c

...	...	@@ -20,6 +20,7 @@
20	20	
21	21	#include <linux/crush/crush.h>
22	22	#include <linux/crush/hash.h>
	23	+ #include "crush_ln_table.h"
23	24	
24	25	/*
25	26	* Implement the core CRUSH mapping algorithm.
...	...	@@ -237,6 +238,102 @@ static int bucket_straw_choose(struct crush_bucket_straw *bucket,
237	238	return bucket->h.items[high];

```

238     }
239
240
241 + // compute 2^44*log2(input+1)
242 + uint64_t crush_ln(unsigned xin)
243 + {
244 +     unsigned x=xin, x1;
245 +     int iexpon, index1, index2;
246 +     uint64_t RH, LH, LL, xl64, result;
247 +
248 +     x++;
249 +
250 +     // normalize input
251 +     iexpon = 15;
252 +     while(!(x&0x18000)) { x<<=1; iexpon--; }
253 +
254 +     index1 = (x>>8)<<1;
255 +     // RH ~ 2^56/index1
256 +     RH = __RH_LH_tbl[index1 - 256];
257 +     // LH ~ 2^48 * log2(index1/256)
258 +     LH = __RH_LH_tbl[index1 + 1 - 256];
259 +
260 +     // RH*x ~ 2^48 * (2^15 + xf), xf<2^8
261 +     xl64 = (int64_t)x * RH;
262 +     xl64 >>= 48;
263 +     x1 = xl64;
264 +
265 +     result = iexpon;
266 +     result <<= (12 + 32);
267 +
268 +     index2 = x1 & 0xff;
269 +     // LL ~ 2^48*log2(1.0+index2/2^15)
270 +     LL = __LL_tbl[index2];
271 +
272 +     LH = LH + LL;
273 +
274 +     LH >>= (48-12 - 32);
275 +     result += LH;
276 +
277 +     return result;
278 + }
279 +
280 +
281 + /*
282 +  * straw2
283 +  *
284 +  * for reference, see:
285 +  *
286 +  * http://en.wikipedia.org/wiki/Exponential_distribution#Distribution_of_the_minimum_of_exponential_random_variables
287 +  *
288 +  */
289 +
290 + static int bucket_straw2_choose(struct crush_bucket_straw2 *bucket,
291 +                                int x, int r)
292 + {
293 +     unsigned i, high = 0;
294 +     unsigned u;
295 +     unsigned w;
296 +     __s64 ln, draw, high_draw = 0;
297 +
298 +     for (i = 0; i < bucket->h.size; i++) {
299 +         w = bucket->item_weights[i];
300 +         if (w) {
301 +             u = crush_hash32_3(bucket->h.hash, x,
302 +                               bucket->h.items[i], r);
303 +             u &= 0xffff;
304 +
305 +             /*
306 +              * for some reason slightly less than 0x10000 produces
307 +              * a slightly more accurate distribution... probably a
308 +              * rounding effect.
309 +              *
310 +              * the natural log lookup table maps [0,0xffff]
311 +              * (corresponding to real numbers [1/0x10000, 1] to
312 +              * [0, 0xfffffffffffff] (corresponding to real numbers
313 +              * [-11.090355,0])).

```

```
314 +          */
315 +          ln = crush_ln(u) - 0x10000000000011;
316 +
317 +          /*
318 +           * divide by 16.16 fixed-point weight. note
319 +           * that the ln value is negative, so a larger
320 +           * weight means a larger (less negative) value
321 +           * for draw.
322 +           */
323 +          draw = div64_s64(ln, w);
324 +      } else {
325 +          draw = S64_MIN;
326 +      }
327 +
328 +      if (i == 0 || draw > high_draw) {
329 +          high = i;
330 +          high_draw = draw;
331 +      }
332 +  }
333 +  return bucket->h.items[high];
334 + }
335 +
336 +
240 337 static int crush_bucket_choose(struct crush_bucket *in, int x, int r)
241 338 {
242 339     dprintk(" crush_bucket_choose %d x=%d r=%d\n", in->id, x, r);
...   ... @@ -254,12 +351,16 @@ static int crush_bucket_choose(struct crush_bucket *in, int x, int r)
254 351     case CRUSH_BUCKET_STRAW:
255 352         return bucket_straw_choose((struct crush_bucket_straw *)in,
256 353             x, r);
354 +     case CRUSH_BUCKET_STRAW2:
355 +         return bucket_straw2_choose((struct crush_bucket_straw2 *)in,
356 +             x, r);
257 357     default:
258 358         dprintk("unknown bucket %d alg %d\n", in->id, in->alg);
259 359         return in->items[0];
260 360     }
261 361 }
262 362
263 363 +
263 364 /*
264 365  * true if device is marked "out" (failed, fully offloaded)
265 366  * of the cluster
...   ...
```

▼ net/ceph/osdmap.c

```
...   ... @@ -122,6 +122,22 @@ static int crush_decode_straw_bucket(void **p, void *end,
122 122     return -EINVAL;
123 123 }
124 124
125 + static int crush_decode_straw2_bucket(void **p, void *end,
126 +     struct crush_bucket_straw2 *b)
127 + {
128 +     int j;
129 +     dout("crush_decode_straw2_bucket %p to %p\n", *p, end);
130 +     b->item_weights = kcalloc(b->h.size, sizeof(u32), GFP_NOFS);
131 +     if (b->item_weights == NULL)
132 +         return -ENOMEM;
133 +     ceph_decode_need(p, end, b->h.size * sizeof(u32), bad);
134 +     for (j = 0; j < b->h.size; j++)
135 +         b->item_weights[j] = ceph_decode_32(p);
136 +     return 0;
137 + bad:
138 +     return -EINVAL;
139 + }
140 +
125 141 static int skip_name_map(void **p, void *end)
126 142 {
127 143     int len;
...   ... @@ -204,6 +220,9 @@ static struct crush_map *crush_decode(void *pbyval, void *end)
204 220     case CRUSH_BUCKET_STRAW:
205 221         size = sizeof(struct crush_bucket_straw);
206 222         break;
223 +     case CRUSH_BUCKET_STRAW2:
```

	224	+	size = sizeof(struct crush_bucket_straw2);
	225	+	break ;
207	226		default:
208	227		err = -EINVAL;
209	228		goto bad;
...	...		@@ -261,6 +280,12 @@ static struct crush_map *crush_decode(void *pbyval, void *end)
261	280		if (err < 0)
262	281		goto bad;
263	282		break ;
	283	+	case CRUSH_BUCKET_STRAW2:
	284	+	err = crush_decode_straw2_bucket(p, end,
	285	+	(struct crush_bucket_straw2 *)b);
	286	+	if (err < 0)
	287	+	goto bad;
	288	+	break ;
264	289		}
265	290		}
266	291		
...	...		

Please [register](#) or [sign in](#) to comment