



远古超传

ULTRADATA TRANSMISSION

Rsync增量传输场景

软件定义网络 超传连接未来

- Rsync 增量传输的性能可被近似建模为：

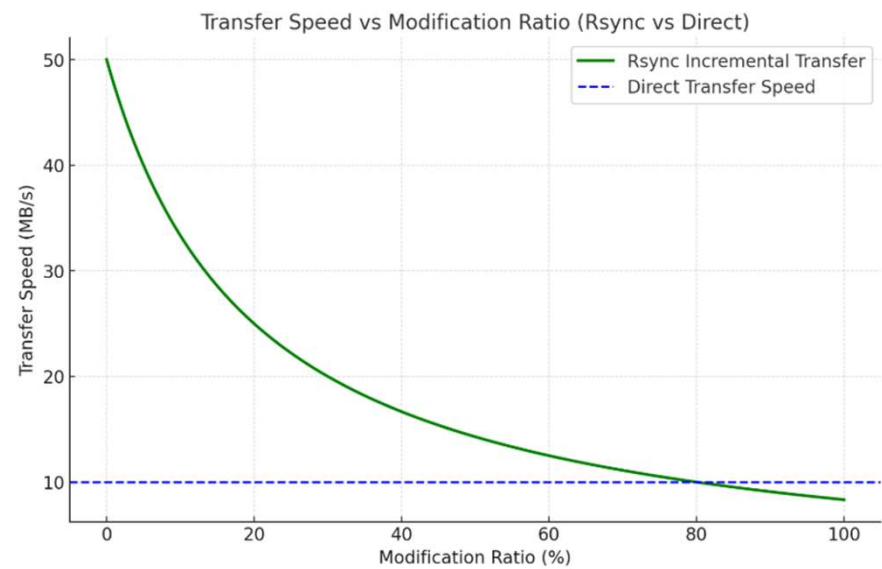
$$S_{\text{rsync}}(m) = S_{\text{base}} \cdot (1 - m) + S_{\text{delta}} \cdot m$$

其中：

- S_{base} 是无修改块的传输速率。
- S_{delta} 是修改块的传输速率。
- m 是修改比例，范围 0-1。
- 蓝线的直接传输模型通常是一个常量 S_{direct} 。

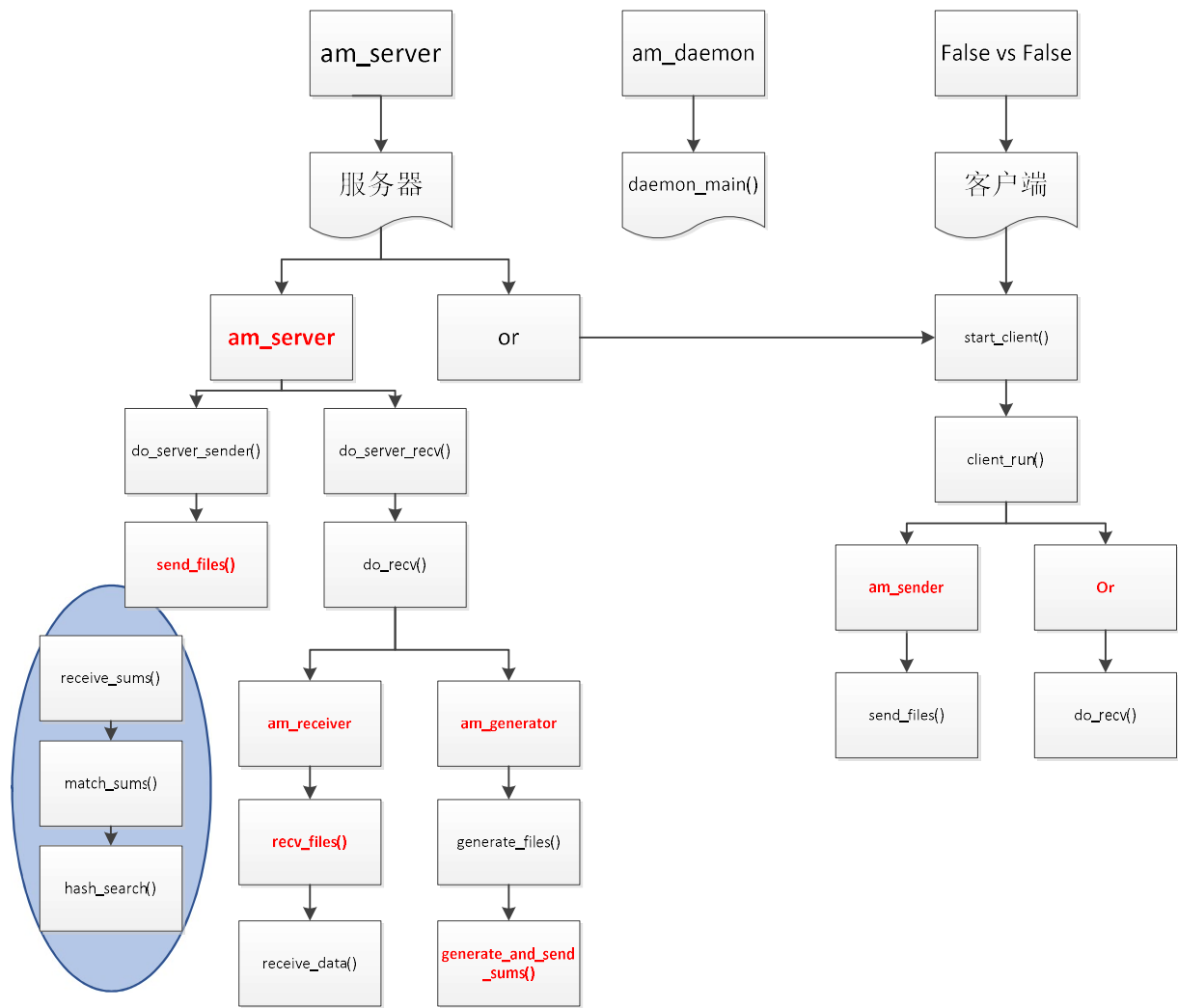
在该模型中：

$$S_{\text{base}} \approx 50 \text{ MB/s}, \quad S_{\text{delta}} \approx 10 \text{ MB/s}.$$



核心问题

- delta (encoding)change algorithm
- rolling checksum algorithm
- **The rsync algorithm**
 - 1. quick check (size, always checksum, mtime)
 - 2. rolling checksum **in send_files()**
 - 3. 3-level search **in send_files()**
- Server, Client and Daemon → sender & receiver & generator
 - send_files()/recv_files()/generate_files()



1. quick_check_ok()

```
int quick_check_ok(enum filetype ftype, const char *fn, struct file_struct *file, STRUCT_STAT *st)
{
    switch (ftype) {
        case FT_REG:
            if (st->st_size != F_LENGTH(file))
                return 0;

            /* If always_checksum is set then we use the checksum instead
             * of the file mtime to determine whether to sync. */
            if (always_checksum > 0) {
                char sum[MAX_DIGEST_LEN];
                file_checksum(fn, st, sum);
                return memcmp(sum, F_SUM(file), flist_csum_len) == 0;
            }

            if (size_only > 0)
                return 1;

            if (ignore_times)
                return 0;

            if (mtime_differs(st, file))
                return 0;
    }
}
```

- **File Size:** If the sizes of the source and destination files differ, the file is marked for synchronization.
- **Always checksum:** Optionally, a full checksum of the file can be calculated for validation.
- **Modification time (mtime):** If the last modification times of the files differ, rsync flags the file for further checks or synchronization.

sum functions

- **receive_sums()**: called in send_files() to receive checksums from receiver
- match_sums(): called in send_files() to match checksums from receiver
- **generate_and_send_sums()**: called in generate_files() in receiver's side
- sum_update()
- get_checksum1()
- get_checksum2()
- sum_end()
- sum_init()
- ...

```
71 static struct sum_struct *receive_sums(int f)
72 {
73     struct sum_struct *s = new(struct sum_struct);
74     int lull_mod = protocol_version >= 31 ? 0 : allowed_lull * 5;
75     OFF_T offset = 0;
76     int32 i;
77
78     read_sum_head(f, s);
79
80     s->sums = NULL;
81
82     if (DEBUG_GTE(DELTA_SUM, 3)) {
83         rprintf(FINFO, "count-%s n-%ld rem-%ld\n",
84             big_num(s->count), (long)s->blength, (long)s->remainder);
85     }
86
87     if (append_mode > 0) {
88         s->flength = (OFF_T)s->count * s->blength;
89         if (s->remainder)
90             s->flength += s->blength - s->remainder;
91         return s;
92     }
93
94     if (s->count == 0)
95         return(s);
96
97     s->sums = new_array(struct sum_buf, s->count);
98     s->sum2_array = new_array(char, (size_t)s->count * xfer_sum_len);
99
100     for (i = 0; i < s->count; i++) {
101         s->sums[i].sum1 = read_int(f);
102         read_buf(f, sum2_at(s, i), s->s2length);
```

```
763 static int generate_and_send_sums(int fd, OFF_T len, int f_out, int f_copy)
764 {
765     int32 i;
766     struct map_struct *mapbuf;
767     struct sum_struct sum;
768     OFF_T offset = 0;
769
770     sum_sizes_sqrtot(&sum, len);
771     if (sum.count < 0)
772         return -1;
773     write_sum_head(f_out, &sum);
774
775     if (append_mode > 0 && f_copy < 0)
776         return 0;
777
778     if (len > 0)
779         mapbuf = map_file(fd, len, MAX_MAP_SIZE, sum.blength);
780     else
781         mapbuf = NULL;
782
783     for (i = 0; i < sum.count; i++) {
784         int32 n1 = (int32)MIN(len, (OFF_T)sum.blength);
785         char *map = map_ptr(mapbuf, offset, n1);
786         char sum2[MAX_DIGEST_LEN];
787         uint32 sum1;
788
789         len -= n1;
790         offset += n1;
791
792         if (f_copy >= 0) {
793             full_write(f_copy, map, n1);
794             if (append_mode > 0)
795                 continue;
796         }
797
798         sum1 = get_checksum1(map, n1);
799         get_checksum2(map, n1, sum2);
800
801         if (DEBUG_GTE(DELTA_SUM, 3)) {
802             rprintf(FINFO,
803                 "chunk[%s] offset-%s len-%ld sum1-%ld\n",
804                 big_num(i), big_num(offset + n1), (long)n1,
805                 (unsigned long)sum1);
806         }
807         write_int(f_out, sum1);
808         write_buf(f_out, sum2, sum.s2length);
809     }
```

2. match_sums()

- Files are divided into fixed-size blocks
- A **rolling checksum** is calculated for each block in the source file. This checksum allows quick recalculations when the comparison window moves by one byte, **using a formula** that updates the checksum incrementally rather than recomputing it from scratch.
- The **destination file's blocks** have both a rolling (weak) checksum and a strong checksum (e.g., MD5) stored. The weak checksum quickly identifies potential matches, and the strong checksum **confirms** these matches.

```
391 if (len > 0 && s->count > 0) {
392     build_hash_table(s);
393
394     if (DEBUG_GTE(DELTA_SUM, 2))
395         rprintf(FINFO, "built hash table\n");
396
397     hash_search(f, s, buf, len);
398
399     if (DEBUG_GTE(DELTA_SUM, 2))
400         rprintf(FINFO, "done hash search\n");
401 } else {
402     OFF_T j;
403     /* by doing this in pieces we avoid too many seeks */
404     for (j = last_match + CHUNK_SIZE; j < len; j += CHUNK_SIZE)
405         matched(f, s, buf, j, -2);
406     matched(f, s, buf, len, -1);
407 }
```

get_checksum1()/get_checksum2()

The weak checksum algorithm we used in our implementation was inspired by Mark Adler's adler-32 checksum. Our checksum is defined by

```
281 /*
282  a simple 32 bit checksum that can be updated from either end
283  (inspired by Mark Adler's Adler-32 checksum)
284  */
285 uint32 get_checksum1(char *buf1, int32 len)
286 {
287     int32 i;
288     uint32 s1, s2;
289     schar *buf = (schar *)buf1;
290
291     s1 = s2 = 0;
292     for (i = 0; i < (len-4); i+=4) {
293         s2 += 4*(s1 + buf[i]) + 3*buf[i+1] + 2*buf[i+2] + buf[i+3] + 10*CHAR_OFFSET;
294         s1 += (buf[i+0] + buf[i+1] + buf[i+2] + buf[i+3] + 4*CHAR_OFFSET);
295     }
296     for (; i < len; i++) {
297         s1 += (buf[i]+CHAR_OFFSET); s2 += s1;
298     }
299     return (s1 & 0xffff) + (s2 << 16);
300 }
```

$$a(k, l) = \left(\sum_{i=k}^l X_i \right) \bmod M$$

$$b(k, l) = \left(\sum_{i=k}^l (l - i + 1) X_i \right) \bmod M$$

$$s(k, l) = a(k, l) + 2^{16} b(k, l)$$

where $s(k, l)$ is the rolling checksum of the bytes $X_k \dots X_l$. For simplicity and speed, we use $M = 2^{16}$.

Hash check not hit

Incremental Rolling Checksum Formula

For a block of size k :

1. Checksum Initialization:

$$s1 = \sum_{i=1}^k \text{block}[i]$$
$$s2 = \sum_{i=1}^k (k - i + 1) \cdot \text{block}[i]$$

2. Rolling Update: When the window shifts one byte:

- Remove the first byte (`map[0]`) and add the new byte (`map[k]` if available):

$$s1_{\text{new}} = s1_{\text{old}} - \text{block}[0] + \text{block}[k]$$
$$s2_{\text{new}} = s2_{\text{old}} - k \cdot \text{block}[0] + s1_{\text{new}}$$

```
null_hash:
    backup = (int32)(offset - last_match);
    /* We sometimes read 1 byte prior to last_match... */
    if (backup < 0)
        backup = 0;

    /* Trim off the first byte from the checksum */
    more = offset + k < len;
    map = (schar *)map_ptr(buf, offset - backup, k + more + backup) + backup;
    s1 -= map[0] + CHAR_OFFSET;
    s2 -= k * (map[0] + CHAR_OFFSET);

    /* Add on the next byte (if there is one) to the checksum */
    if (more) {
        s1 += map[k] + CHAR_OFFSET;
        s2 += s1;
    } else
        --k;

    /* By matching early we avoid re-reading the
       data 3 times in the case where a token
       match comes a long way after last
       match. The 3 reads are caused by the
       running match, the checksum update and the
       literal send. */
    if (backup >= s->blength+CHUNK_SIZE && end_offset > CHUNK_SIZE)
        matched(f, s, buf, offset - s->blength, -2);
    } while (++offset < end);
```

$$a(k+1, l+1) = (a(k, l) - X_k + X_{l+1}) \bmod M$$

$$b(k+1, l+1) = (b(k, l) - (l - k + 1)X_k + a(k+1, l+1)) \bmod M$$

Thus the checksum can be calculated for blocks of length S at all possible offsets within a file in a "rolling" fashion, with very little computation at each point.

3. hash_search()

- Level 1: 16-bit hash table:
 - A hash of the rolling checksum is calculated and used to index a hash table of block checksums.
 - The hash table reduces the number of blocks to consider.
- Level 2: Weak checksum comparison:
 - If the hash table indicates a match, the weak checksum of the current source block is compared to the destination's checksum.
- Level 3: Strong checksum validation:
 - For potential matches, the strong checksum (e.g., MD5) confirms the block match with near certainty.

```

sum = (s1 & 0xffff) | (s2 << 16);
} else {
sum = (s1 & 0xffff) | (s2 << 16);
hash_entry = BIG_SUM2HASH(sum);
if ((i = hash_table[hash_entry]) < 0)
goto null_hash;
}
prev = &hash_table[hash_entry];

hash_hits++;
do {
int32 l;

/* When updating in-place, the chunk's offset must be
 * either >= our offset or identical data at that offset.
 * Remove any bypassed entries that we can never use. */
if (updating_basis_file && s->sums[i].offset < offset
&& !(s->sums[i].flags & SUMFLG_SAME_OFFSET)) {
*prev = s->sums[i].chain;
continue;
}
prev = &s->sums[i].chain;

if (sum != s->sums[i].sum1)
continue;

/* also make sure the two blocks are the same length */
l = (int32)MIN((OFF_T)s->blength, len-offset);
if (l != s->sums[i].len)
continue;

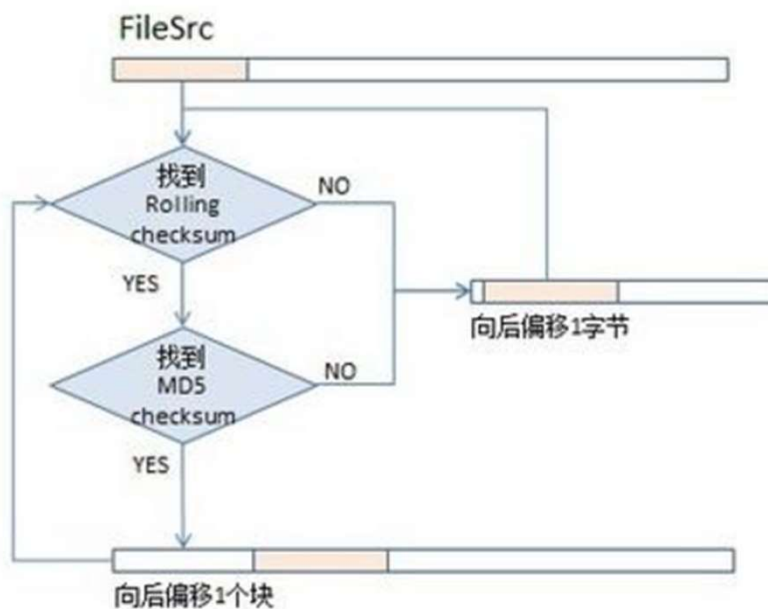
if (DEBUG_GTE(DELTA_SUM, 3)) {
rprintf(FINFO,
"potential match at %s i=%ld sum=%08x\n",
big_num(offset), (long)i, sum);
}

if (!done_csum2) {
map = (schar *)map_ptr(buf, offset, l);
get_checksum2((char *)map, l, sum2);
done_csum2 = 1;
}

if (memcmp(sum2, sum2_at(s, i), s->s2length) != 0) {
false_alarms++;
continue;
}
}

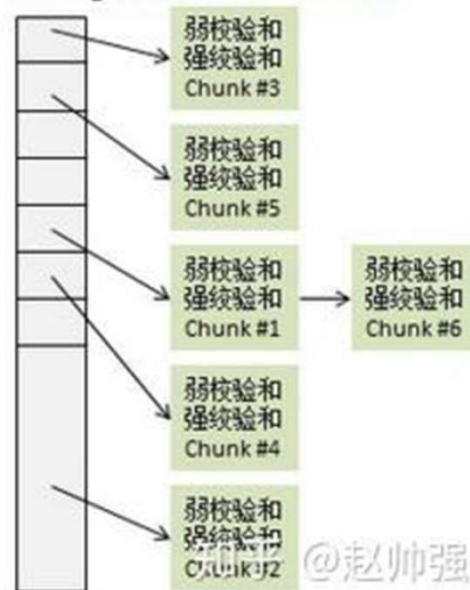
```

FileDst



FileDst

Rolling Checksum Hash Table



@赵帅强

Multiple hits in hash table

```
240 /* When updating in-place, the best possible match is
241 * one with an identical offset, so we prefer that over
242 * the adjacent want_i optimization. */
243 if (updating_basis_file) {
244     /* All the generator's chunks start at blength boundaries. */
245     while (aligned_offset < offset) {
246         aligned_offset += s->blength;
247         aligned_i++;
248     }
249     if ((offset == aligned_offset
250         || (sum == 0 && 1 == s->blength && aligned_offset + 1 <= len))
251         && aligned_i < s->count) {
252         if (i != aligned_i) {
253             if (sum != s->sums[aligned_i].sum1
254                 || 1 != s->sums[aligned_i].len
255                 || memcmp(sum2, sum2_at(s, aligned_i), s->s2length) != 0)
256                 goto check_want_i;
257             i = aligned_i;
258         }
259         if (offset != aligned_offset) {
260             /* We've matched some zeros in a spot that is also zeros
261              * further along in the basis file, if we find zeros ahead
262              * in the sender's file, we'll output enough literal data
263              * to re-align with the basis file, and get back to seeking
264              * instead of writing. */
265             backup = (int32)(aligned_offset - last_match);
266             if (backup < 0)
267                 backup = 0;
268             map = (schar *)map_ptr(buf, aligned_offset - backup, 1 + backup)
269                 + backup;
270             sum = get_checksum1((char *)map, 1);
271             if (sum != s->sums[i].sum1)
272                 goto check_want_i;
273             get_checksum2((char *)map, 1, sum2);
274             if (memcmp(sum2, sum2_at(s, i), s->s2length) != 0)
275                 goto check_want_i;
276             /* OK, we have a re-alignment match. Bump the offset
277              * forward to the new match point. */
278             offset = aligned_offset;
279         }
280         /* This identical chunk is in the same spot in the old and new file. */
281         s->sums[i].flags |= SUMFLG_SAME_OFFSET;
282         want_i = i;
283     }
284 }
```

```
286 check_want_i:
287     /* we've found a match, but now check to see
288     * if want_i can hint at a better match. */
289     if (i != want_i && want_i < s->count
290         && (!updating_basis_file || s->sums[want_i].offset >= offset
291             || s->sums[want_i].flags & SUMFLG_SAME_OFFSET)
292         && sum == s->sums[want_i].sum1
293         && memcmp(sum2, sum2_at(s, want_i), s->s2length) == 0) {
294         /* we've found an adjacent match - the RLL coder
295          * will be happy */
296         i = want_i;
297     }
298     want_i = i + 1;
299
300     matched(f, s, buf, offset, i);
301     offset += s->sums[i].len - 1;
302     k = (int32)MIN((OFF_T)s->blength, len - offset);
303     map = (schar *)map_ptr(buf, offset, k);
304     sum = get_checksum1((char *)map, k);
305     s1 = sum & 0xFFFF;
306     s2 = sum >> 16;
307     matches++;
308     break;
309 } while ((i = s->sums[i].chain) >= 0);
310 }
```

```

1030 /**
1031  * Transmit a verbatim buffer of length @p n followed by a token.
1032  * If token == -1 then we have reached EOF
1033  * If n == 0 then don't send a buffer
1034  */
1035 void send_token(int f, int32 token, struct map_struct *buf, OFF_T offset,
1036                int32 n, int32 tokenlen)
1037 {
1038     switch (do_compression) {
1039     case CPRES_NONE:
1040         simple_send_token(f, token, buf, offset, n);
1041         break;
1042     case CPRES_ZLIB:
1043     case CPRES_ZLIBX:
1044         send_deflated_token(f, token, buf, offset, n, tokenlen);
1045         break;
1046     #ifdef SUPPORT_ZSTD
1047     case CPRES_ZSTD:
1048         send_zstd_token(f, token, buf, offset, n);
1049         break;
1050     #endif
1051     #ifdef SUPPORT_LZ4
1052     case CPRES_LZ4:
1053         send_compressed_token(f, token, buf, offset, n);
1054         break;
1055     #endif
1056     default:
1057         NOISY_DEATH("Unknown do_compression value");
1058     }
1059 }
1060
1061 static void matched(
1062     int32 n = (int32)
1063     int32 j;
1064     if (DEBUG_GTE(DE
1065         rprintf(FIN
1066         "match a
1067         big_num(
1068         (long)s-
1069     )
1070     }
1071     send_token(f, i, buf, last_match, n, i < 0 ? 0 : s->sums[i].len);
1072     data_transfer += n;
1073
1074     if (i >= 0) {
1075         stats.matched_data
1076         n += s->sums[i].ler
1077     }
1078
1079     for (j = 0; j < n; j +=
1080         int32 n1 = MIN(CHUN
1081         sum_update(map_ptr(
1082     )
1083
1084     if (i >= 0)
1085         last_match = offset
1086     else
1087         last_match = offset
1088
1089     if (buf && INFO_GTE(PRC
1090         show_progress(last_match, buf->file_size);
1091 }

```

In recv_data():
while ((i = recv_token(f_in, &data)) != 0){...}

```

280 /* non-compressing recv token */
281 static int32 simple_recv_token(int f, char **data)
282 {
283     static int32 residue;
284     static char *buf;
285     int32 n;
286
287     if (!buf)
288         buf = new_array(char, CHUNK_SIZE);
289
290     if (residue == 0) {
291         int32 i = read_int(f);
292         if (i <= 0)
293             return i;
294         residue = i;
295     }
296
297     *data = buf;
298     n = MIN(CHUNK_SIZE, residue);
299     residue -= n;
300     read_buf(f, buf, n);
301     return n;
302 }

```

```

314 while ((i = recv_token(f_in, &data)) != 0) {
315     if (INFO_GTE(PROGRESS, 1))
316         show_progress(offset, total_size);
317
318     if (allowed_lull)
319         maybe_send_keepalive(time(NULL), MSK_ALLOW_FLUSH | MSK_ACTIVE_RECEIVER);
320
321     if (i > 0) {
322         if (DEBUG_GTE(DELTA_SUM, 3)) {
323             rprintf(FINFO, "data recv %d at %s\n",
324                 i, big_num(offset));
325         }
326
327         stats.literal_data += i;
328         cleanup_got_literal = 1;
329
330         sum_update(data, i);
331
332         if (fd != -1 && write_file(fd, 0, offset, data, i) != i)
333             goto report_write_error;
334         offset += i;
335         continue;
336     }
337
338     i = -(i+1);
339     offset2 = i * (OFF_T)sum.blength;
340     len = sum.blength;
341     if (i == (int)sum.count-1 && sum.remainder != 0)
342         len = sum.remainder;
343
344     stats.matched_data += len;
345
346     if (DEBUG_GTE(DELTA_SUM, 3)) {
347         rprintf(FINFO,
348             "chunk[%d] of size %ld at %s offset-%s\n",
349             1, (long)len, big_num(offset2), big_num(offset),
350             updating_basis_or_equiv && offset == offset2 ? " (seek)" : "");
351     }
352
353     if (mapbuf) {
354         map = map_ptr(mapbuf, offset2, len);
355
356         see_token(map, len);
357         sum_update(map, len);
358     }
359
360     if (updating_basis_or_equiv) {
361         if (offset == offset2 && fd != -1) {
362             if (skip_matched(fd, offset, map, len) < 0)
363                 goto report_write_error;
364             offset += len;
365             continue;
366         }
367     }
368     if (fd != -1 && map && write_file(fd, 0, offset, map, len) != (int)len)
369         goto report_write_error;
370     offset += len;
371 }

```


核心问题

- delta (encoding)change algorithm
- rolling checksum algorithm
- **The rsync algorithm**
 - 1. quick check (size, always checksum, mtime)
 - 2. rolling checksum **in send_files()**
 - 3. 3-level search **in send_files()**
- Server, Client and Daemon → sender & receiver & generator
 - send_files()/recv_files()/generate_files()



远古超传

ULTRADATA TRANSMISSION

谢谢!

THANKS

软件定义网络 超传连接未来