# The Hash-Trie of Knuth & Liang
## — Addendum —

Ştefan Vargyas

**stvar@yahoo.com**

Nov 28, 2015

## Table of Contents

## 1   Introduction

The purpose of the following sections is to provide mathematical proofs for the claims made at the end of section 2.3, *The implementation of* `HashTrie<>` *Class Template*, on page 21 of the technical report *The Hash-Trie of Knuth & Liang: A C++11 Implementation*, **hash-trie-impl.pdf**:

> The replacement of expression "`h + tolerance - mod_x`" with the expression "`add(h, tolerance2) - mod_x`" had to be done, because, as it is easily provable, under cetain conditions, the subexpression "`h + tolerance`" is exceeding the upper bound `trie_size` of the boxed integer `pointer_t`. It is worthy of notice the readily provable fact that both assignments to `last_h` on lines 2523 and 2533 are correct: each of the expression on the right side of these assignments do not exceed upon evaluation the bounds of `pointer_t`.

The claim made by the first phrase is restated and proved by fact 17; the claim of the second phrase – by fact 19. Alongside these two facts, the section containing them establishes by proof a few more facts which concern the inner workings of the core algorithm of class `HashTrie<>`.

## 2   Mathematical Evaluations

**1 Definition** (Defining parameters)**.**

(1)   $trie\_size \in \mathbb{N}$,

(2)   $max\_letter \in \mathbb{N}^*$.

(3)   $tolerance \in \mathbb{N}$.

**2 Definition** (Dependent constants)**.**

(4)   $mod\_x \stackrel{\text{def}}{=} trie\_size - 2 \cdot max\_letter$,

(5)   $alpha \stackrel{\text{def}}{=} \lceil 0.61803 \cdot mod\_x \rceil$,

(6)   $max\_h \stackrel{\text{def}}{=} mod\_x + max\_letter$,

(7)   $max\_x \stackrel{\text{def}}{=} mod\_x - alpha$.

**3 Axioms** (Initiating assertions).

(8)  $alpha > 0$,

(9)  $tolerance > 0$,

(10)  $mod\_x > 0$,

(11)  $max\_h > tolerance$,

(12)  $mod\_x > tolerance$,

(13)  $mod\_x > alpha$.

**4 Proposition.**

(14)  $trie\_size > 0$,

(15)  $max\_h > mod\_x$,

(16)  (12) $\implies$ (11),

(17)  $tolerance < trie\_size$,

(18)  $max\_h = trie\_size - max\_letter$,

(19)  $max\_h < trie\_size$,

(20)  (11) $\iff max\_letter + tolerance < trie\_size$.

*Proof.* For (14): apply (10) and (4).
For (15): apply (6) and (2).
For (16): apply (15).
For (17): $tolerance \stackrel{(12)}{<} mod\_x \stackrel{(4)}{=} trie\_size - 2 \cdot max\_letter \stackrel{(2)}{<} trie\_size$.
For (18): apply (6) and (4).
For (19): apply (18) and (2).
For (20): (11) $\stackrel{(18)}{\iff} trie\_size - max\_letter > tolerance \iff max\_letter + tolerance < trie\_size$.
$\square$

**5 Definition** (The *floor* function).

(21)  $\lfloor x \rfloor \stackrel{\mathrm{def}}{=} \max \mathcal{M}(x) \in \mathbb{Z}$, for $x \in \mathbb{R}$,

where

(22)  $\mathcal{M}(x) \stackrel{\mathrm{def}}{=} \{k \in \mathbb{Z} \mid k \le x\}$.

**6 Proposition** (Basic properties of "$\lfloor \cdot \rfloor$").

(23)  $\lfloor x \rfloor \le x$, *for* $x \in \mathbb{R}$,

(24)  $n = \lfloor n \rfloor$, *for* $n \in \mathbb{Z}$,

(25)  $\lfloor x \rfloor \le \lfloor y \rfloor \iff \mathcal{M}(x) \subseteq \mathcal{M}(y)$, *for* $x, y \in \mathbb{R}$,

(26)  $x \le y \implies \lfloor x \rfloor \le \lfloor y \rfloor$, *for* $x, y \in \mathbb{R}$,

(27)  $n \le x \iff n \le \lfloor x \rfloor$, *for* $n \in \mathbb{Z}$ *and* $x \in \mathbb{R}$,

(28)  $n \le x < n + 1 \iff \lfloor x \rfloor = n$, *for* $n \in \mathbb{Z}$ *and* $x \in \mathbb{R}$,

(29)  $a \bmod b = a - \lfloor a/b \rfloor \cdot b$, *for* $a \in \mathbb{Z}$ *and* $b \in \mathbb{N}^*$.

*Proof.* For (23): $x \in \mathbb{R} \stackrel{(21)}{\implies} \lfloor x \rfloor \in \mathcal{M}(x) \stackrel{(22)}{\implies} \lfloor x \rfloor \le x$.
For (24): by (23) $\lfloor n \rfloor \le n$; $n \le n \in \mathbb{Z} \stackrel{(22)}{\implies} n \in \mathcal{M}(n) \implies n \le \max \mathcal{M}(n) \stackrel{(21)}{=} \lfloor n \rfloor$.
For (25): "$\Rightarrow$": $k \in \mathcal{M}(x) \implies k \le \max \mathcal{M}(x) \stackrel{(21)}{=} \lfloor x \rfloor \stackrel{\mathrm{hyp}}{\le} \lfloor y \rfloor \stackrel{(21)}{\in} \mathcal{M}(y) \stackrel{(22)}{\implies} k \in \mathcal{M}(y)$; "$\Leftarrow$":
$\lfloor x \rfloor \stackrel{(21)}{\in} \mathcal{M}(x) \stackrel{\mathrm{hyp}}{\subseteq} \mathcal{M}(y) \implies \lfloor x \rfloor \in \mathcal{M}(y) \stackrel{(21)}{\implies} \lfloor x \rfloor \le \lfloor y \rfloor$.
For (26): $x \le y \stackrel{(22)}{\implies} \mathcal{M}(x) \subseteq \mathcal{M}(y) \stackrel{(25)}{\implies} \lfloor x \rfloor \le \lfloor y \rfloor$.
For (27): "$\Rightarrow$": $n \le x \stackrel{(26)}{\implies} n \stackrel{(24)}{=} \lfloor n \rfloor \le \lfloor x \rfloor$; "$\Leftarrow$": $n \le \lfloor x \rfloor \stackrel{(21)}{\in} \mathcal{M}(x) \stackrel{(22)}{\implies} n \le x$.

For (28): $n \le x \overset{(27)}{\iff} n \le \lfloor x \rfloor$; $x < n+1 \overset{(27)}{\iff} \lfloor x \rfloor < n+1 \iff \lfloor x \rfloor \le n$.

For (29): firstly remark the uniqueness part of the *Euclidean division theorem*: for $x, y \in \mathbb{Z}$, with $y > 0$: if exists $q', q'', r', r'' \in \mathbb{Z}$ such that $x = q' \cdot y + r'$, $x = q'' \cdot y + r''$, with $0 \le r', r'' < y$, then $q' = q''$ and $r' = r''$. The proof of this goes easily: suppose $r' > r''$. Then $r' - r'' = y \cdot (q'' - q')$; follows that $q'' - q' > 0$; thus $q'' - q' \ge 1$; hence $r' - r'' \ge y$, which contradicts $r' - r'' < y \iff 0 \le r', r'' < y$. Now, this uniqueness property leads to (29) if shown that $0 \le a - b \cdot \lfloor a/b \rfloor < b$. Indeed the relation holds true: $n = \lfloor a/b \rfloor \overset{(28)}{\iff} n \le a/b < n+1 \iff b \cdot n \le a < b \cdot (n+1) \iff 0 \le a - b \cdot n < b$.  □

**7 Proposition.** *For $a$, $b$, $x \in \mathbb{Z}$:*

(30) $0 \le a < b \implies a \bmod b = a,$

(31) $0 < b \le a < 2 \cdot b \implies a \bmod b = a - b,$

(32) $0 \le x < max\_x \implies (x + alpha) \bmod mod\_x = x + alpha,$

(33) $max\_x \le x < mod\_x \implies (x + alpha) \bmod mod\_x = x - max\_x.$

*Proof.* For (30): $0 \le a < b \implies 0 \le a/b < 1 \overset{(28)}{\implies} \lfloor a/b \rfloor = 0 \overset{(29)}{\implies} a \bmod b = a.$

For (31): $0 < b \le a < 2 \cdot b \implies 1 \le a/b < 2 \overset{(28)}{\implies} \lfloor a/b \rfloor = 1 \overset{(29)}{\implies} a \bmod b = a - b.$

For (32): $x < max\_x \overset{(7)}{\iff} x + alpha < mod\_x$; $x \ge 0 \overset{(8)}{\implies} x + alpha > 0$; then (30) implies (32).

For (33): $x \ge max\_x \overset{(7)}{\iff} x + alpha \ge mod\_x$; $alpha \overset{(13)}{<} mod\_x$ and $x \overset{\text{hyp}}{<} mod\_x \implies x + alpha < 2 \cdot mod\_x$; the latter two consequents along with (10) conclude to $(x + alpha) \bmod mod\_x \overset{(31)}{=} x + alpha - mod\_x \overset{(7)}{=} x - max\_x.$  □

**8 Notation.** For $h \in \mathbb{N}^*$ let $h' \overset{\text{def}}{=} h + tolerance \in \mathbb{N}^*$ and $h'' \overset{\text{def}}{=} h' - mod\_x \in \mathbb{Z}$.

**9 Proposition.**

(34) $h \le max\_h - tolerance \iff h' \le max\_h,$

(35) $h \ge max\_letter + 1 \implies h' > max\_letter + 1,$

(36) $h \le max\_h \iff h'' \le max\_letter + tolerance,$

(37) $h > max\_h - tolerance \iff h'' > max\_letter,$

(38) $max\_h - tolerance < h \le max\_h \iff max\_letter < h'' \le max\_letter + tolerance,$

(39) $max\_letter + tolerance \le max\_h,$

(40) $max\_letter < h \le max\_h \iff 0 \le max\_h - h < mod\_x,$

(41) $max\_h + tolerance > trie\_size \iff tolerance > max\_letter.$

*Proof.* For (34): $h \le max\_h - tolerance \iff h' = h + tolerance \le max\_h.$

For (35): $h' = h + tolerance \overset{\text{hyp}}{\ge} max\_letter + 1 + tolerance \overset{(9)}{>} max\_letter + 1.$

For (36): $h \le max\_h \iff h'' = h + tolerance - mod\_x \le max\_h - mod\_x + tolerance \overset{(6)}{=} max\_letter + tolerance.$

For (37): $h > max\_h - tolerance \iff h'' = h + tolerance - mod\_x > max\_h - mod\_x \overset{(6)}{=} max\_letter.$

For (38): apply (36) and (37).

For (39): $max\_letter + tolerance \overset{(12)}{\le} max\_letter + mod\_x \overset{(6)}{=} max\_h.$

For (40): $max\_letter < h \le max\_h \iff 0 \le max\_h - h < max\_h - max\_letter \overset{(6)}{=} mod\_x.$

For (41): $max\_h + tolerance > trie\_size \iff tolerance > trie\_size - max\_h \overset{(18)}{=} max\_letter.$  □

**10 Proposition.**

(42) $max\_letter + 1 \le h \le max\_h - tolerance \implies max\_letter + 1 \le h' \le max\_h,$

(43) $max\_h - tolerance < h \le max\_h \implies max\_letter + 1 \le h'' \le max\_h.$

*Proof.* For (42): by (35) $h' > max\_letter + 1$, therefore $h' \ge max\_letter + 1$. By (34) $h' \le max\_h$.

For (43): by (37) $h'' > max\_letter$, therefore $h'' \ge max\_letter + 1$. By (36) and (39) $h'' \le max\_h$.  □

# 3 Applications to Hash-Trie

**11 Fact.** *The definitions (1), (2) and (3) correspond to lines #2042, #2045 and #2043 respectively of the* **C++** *implementation.*

**12 Fact.** *The definitions (4), (5), (6) and (7) correspond to lines #2267, #2266, #2268 and #2269 respectively of the* **C++** *implementation.*

**13 Fact.** *The axioms (8)–(13) correspond one by one to the* CXX_ASSERT *lines #2271–#2276.*

**14 Fact.** *The assignment on line #2445 is correct.*

*Proof.* Due to the definition of `pointer_t`, the statement on line #2445 is correct if and only if `tolerance` is greater or equal than `0` and less or equal than `trie_size` (these are the limiting bounds of `pointer_t`). From (9) and (17) results that the line #2445 is indeed correct. $\square$

**15 Fact.** *The variable* x *declared, initialized and maintained on lines #2231, #2318 and respectively #2477–#2480 is iterating correctly the elements of the sequence* $(x_n)_{n \in \mathbb{N}}$, *where* $x_n \stackrel{\text{def}}{=} (alpha \cdot n) \bmod mod\_x$ *for* $n \in \mathbb{N}$.

*Proof.* Proceed by induction on $n \in \mathbb{N}$. For $n = 0$ the statement made above holds, since the line #2318 is showing that the initial value of x is `0`. Now suppose that prior to executing line #2477 x has the value of $x_n$ for some $n \in \mathbb{N}$. In view of the relations (32) and (33), upon the execution of lines #2477–#2480, x becomes `(x + alpha)` $mod$ `mod_x`. Taking into account that, by the definition of sequence $(x_n)_{n \in \mathbb{N}}$, $x_{n+1} = (x_n + alpha) \bmod mod\_x$, indeed x is $x_{n+1}$ after the line #2480. $\square$

**16 Fact.** *The assertion stated within the comment on line #2482 is correct.*

*Proof.* Need to prove that upon executing line #2482, $max\_letter < h \leq trie\_size - max\_letter \stackrel{(18)}{=} max\_h$: Fact 15 $\implies 0 \leq x < mod\_x \stackrel{(6)}{=} max\_h - max\_letter \iff -1 < x \leq max\_h - max\_letter - 1 \iff max\_letter < h \stackrel{\#2482}{=} x + max\_letter + 1 \leq max\_h$. $\square$

**17 Fact.** *Under certain conditions, the result of evaluating the expression* add(h, tolerance2) *on line #2523 exceeds the value of* trie_size *for some* h.

*Proof.* By fact 16: $max\_letter < h \leq max\_h$. If let $h \stackrel{\text{def}}{=} max\_h$, then, under the condition that $tolerance > max\_letter$, (41) shows that $h + tolerance > trie\_size$ indeed. $\square$

*18 Remark.* The fact above indicates that the expression `h + tolerance` wouldn't have been a proper choice of coding the line #2523: in the case of $h + tolerance > trie\_size$, the evaluation of the expression `h + tolerance` would have caused the program to halt abruptly (assuming that the configuration parameter `CONFIG_HASH_TRIE_STRICT_TYPES` was `#define`d at compile-time).

**19 Fact.** *The assignments to variable* last_h *on lines #2523 and #2533 are both correct. Upon the execution of either of them,* $max\_letter + 1 \leq last\_h \leq max\_h$.

*Proof.* The statements on lines #2523 and #2533 are correct if and only if each of the expression on the right side of the respective assignments evaluates to an integer not exceeding the bounds of type `pointer_t`. By the fact 16, before executing each of the two lines: $max\_letter < h \leq max\_h$. Now, for the case of line #2523 apply (43) (from line #2506: $h > max\_h - tolerance$) and, respectively, for the case of line #2533 apply (42) (from line #2506: $h \leq max\_h - tolerance$). Both give that $max\_letter + 1 \leq last\_h \leq max\_h$. Consequently, the bounds of `pointer_t` are respected: by (2) and (19), the previous double inequality yields: $0 < last\_h < trie\_size$. $\square$

**20 Fact.** *The inner loops of method* HashTrie<>::find *(not displayed by the listing below) that are based on* h *computed by lines #2545–#2551 are finite.*

*Proof.* By the fact 16, each of these loops start iterating with an h satisfying $max\_letter + 1 \leq h \leq max\_h$. The lines #2545–#2551 show that h is incremented circularly within the boundaries $max\_letter + 1$ and $max\_h$. By the fact 19, $max\_letter + 1 \leq last\_h \leq max\_h$ on each execution of lines #2545–#2551, i.e. `last_h` lies between the same boundaries as h. The implementation code also shows (not seen below, though) that `last_h` is an invariant of each of these loops. Consequently, h has to meet `last_h` upon a finitely many succesive calls of the lambda function `compute_the_next_trial_header_location`. This leads the lambda function to return `false` – thus terminating the iterations. $\square$

# A C++ Implementation Excerpts

```
2006  template<
2007      typename C = char,
2008      template<typename> class T = char_traits_t,
2009      typename S = size_traits_t>
2010  class HashTrie :
2011      private T<C>,
2012      private S
2013  {
2014  public:
2015      typedef S size_traits_t;
2016      typedef T<C> char_traits_t;
...  ...
2034  private:
...  ...
2042      using size_traits_t::trie_size;
2043      using size_traits_t::tolerance;
...  ...
2045      using char_traits_t::max_letter;
...  ...
2230      // x_n = (alpha * n) % mod_x
2231      pointer_t x;
...  ...
2249      pointer_t find(const char_t*);
...  ...
2255      static constexpr size_t make_alpha(size_t trie_size, size_t max_letter)
...  ...
2259      { return std::ceil(0.61803 * (trie_size - 2 * max_letter)); }
...  ...
2266      static constexpr size_t alpha = make_alpha(trie_size, max_letter);
2267      static constexpr size_t mod_x = trie_size - 2 * max_letter;
2268      static constexpr size_t max_h = mod_x + max_letter;
2269      static constexpr size_t max_x = mod_x - alpha;
2270
2271      CXX_ASSERT(alpha > 0);
2272      CXX_ASSERT(tolerance > 0);
2273      CXX_ASSERT(trie_size > 2 * max_letter);
2274      CXX_ASSERT(max_h > tolerance);
2275      CXX_ASSERT(mod_x > tolerance);
2276      CXX_ASSERT(mod_x > alpha);
...  ...
2296  };


2298  template<
2299      typename C,
2300      template<typename> class T,
2301      typename S>
2302  HashTrie<C, T, S>::HashTrie()
2303  {
...  ...
2318      x = 0;
2319  }


2321  template<
2322      typename C,
2323      template<typename> class T,
2324      typename S>
2325  typename
2326      HashTrie<C, T, S>::pointer_t
2327      HashTrie<C, T, S>::find(const char_t* str)
2328  {
...  ...
2445      const pointer_t tolerance2 = tolerance;
2446      // trial header location
2447      pointer_t h;
2448      // the final one to try
2449      pointer_t last_h; // INT: int last_h;
2450
2451      const auto get_set_for_computing_header_locations = [&]() {
2452          // 24. Get set for computing header locations
...  ...
2477          if (x >= max_x)
2478              x -= max_x;
2479          else
2480              x += alpha;
2481
2482          h = x + max_letter + 1; // now max_letter < h <= trie_size - max_letter
...  ...
```

```
2506            if (h > max_h - tolerance) {
…   …
2523                last_h = add(h, tolerance2) - mod_x;
2524            }
2525            else {
…   …
2533                last_h = h + tolerance;
2534            }
2535        };
2536
2537        const auto compute_the_next_trial_header_location = [&]() {
2538            // 25. Compute the next trial header location h, or abort find
…   …
2545            if (h == last_h)
2546                return false;
2547            if (h == max_h)
2548                h = max_letter + 1;
2549            else
2550                h ++;
2551            return true;
2552        };
…   …
2642  }
```