# On the Horton-Strahler Number for Combinatorial Tries

Markus E. Nebel

Johann Wolfgang Goethe-Universität
Fachbereich Informatik
Frankfurt am Main
Germany

## Abstract

In this paper we investigate the average Horton-Strahler number of all possible tree-structures of binary tries. For that purpose we consider a generalization of extended binary trees where leaves are distinguished in order to represent the location of keys within a corresponding trie. Assuming a uniform distribution for those trees we prove that the expected Horton-Strahler number of a tree with $\alpha$ internal nodes and $\beta$ leaves that correspond to a key is asymptotically given by

$$\frac{4^{2\beta - \alpha} \log(\alpha)(2\beta - 1)(\alpha + 1)(\alpha + 2)\binom{2\alpha + 1}{\alpha - 1}}{8\sqrt{\pi}\alpha^{3/2}\log(2)(\beta - 1)\beta\binom{2\beta}{\beta}^2}$$

provided that $\alpha$ and $\beta$ grow in some fixed proportion $\rho$ when $\alpha \to \infty$. A similar result is shown for trees with $\alpha$ internal nodes but with an arbitrary number of keys.

**AMS Subject Classification:** 05A15, 05C05, 68W40.

# 1 Introduction

Let $T$ be a binary tree, i.e. a tree where each node has at most two descendants. Then the *Horton-Strahler number* of $T$ denoted $hs(T)$ is recursively defined by

$$hs(T) := \begin{cases} 0 & : \quad T \text{ is either a leaf or empty} \\ hs(T.l) + 1 & : \quad \text{if } hs(T.l) = hs(T.r) \\ \max(hs(T.l), hs(T.r)) & : \quad \text{otherwise} \end{cases}.$$

Here, $T.l$ (resp. $T.r$) denotes the left (resp. right) subtree of $T$. The Horton-Strahler number was originally introduced to classify river systems (see [Hor45] and [Str52]) but it has also been adopted in computer science, molecular biology, medicine and other disciplines. Ershov [Ers58], for example, has shown that the minimal number of registers needed to evaluate an arithmetic expression $\mathcal{E}$ with binary operators, which is represented as a binary tree $T(\mathcal{E})$ (the syntax-tree), is given by $1 + hs(T(\mathcal{E}))$. If all syntax-trees with $n$ internal nodes ($n$ binary operators) are assumed to be equally likely, then it is known that the average

number of registers that are needed to evaluate an expression optimally is given by

$$\frac{1}{2}\log_2(2\pi^2 n) - \frac{\gamma+2}{2\ln(2)} + F(n) + \mathcal{O}(n^{-1/2+\delta}) \tag{1}$$

for all $\delta > 0$, $n \to \infty$, where $\gamma = 0.577215\ldots$ is Euler's constant and $F$ is a periodic, oscillating function of small amplitude (see e.g. [FRV79] and [Kem79]). Syntax-trees corresponding to expressions built with unary and binary operators were considered in [FlP86]. Furthermore, the minimum stack-size required for a traversal of a binary tree $T$ is also given by $1 + hs(T)$ (e.g. see [FRV79] and [Fra84]). Meir, Moon and Pounder [MMP80] investigated the Horton-Strahler number of channel networks with a fixed number of inputs. The combinatorics of the Horton-Strahler analysis has been used in computer graphics for the creation of faithful synthetic images of trees (see [VEJA89]). The impact of the Horton-Strahler number on molecular biology comes from theoretical considerations about secondary structures of single-stranded nucleic acids (see [Vie90] and the references given there).

All those applications and studies have in common that they deal with ordinary extended binary trees, i.e. trees where each node is either a leaf or has two descendants. All the cited papers which present an average case analysis consider the uniform model, i.e. they assume that all trees of a given size are equally likely.

The Horton-Strahler number of tries has been investigated in a recent work by Devroye and Kruszewski [DeK96]. A trie is a binary tree which is used to store the set of keys $K = \{k_1, \ldots, k_n\}$ in the following manner: Each key $k_i$, considered as a string of 0's and 1's due to its binary representation, defines a path in a binary tree $T$ (0 indicates a left turn, 1 a right turn); the trie defined by $k_1, \ldots, k_n$ is the smallest binary tree for which the paths truncated at the leaves of $T$ are all pairwise different. Thus each leaf of $T$ *stores* exactly one of the keys $k_i$, $1 \leq i \leq n$. Note that $T$ does not need to be an extended binary tree. $T$ might have internal nodes with only one successor. However, the Horton-Strahler number $hs(T)$ remains unchanged when we turn $T$ into an extended binary tree. Devroye and Kruszewski considered random tries constructed from $n$ i.i.d. sequences of Bernoulli random variables with parameter $p$, $0 < p < 1$; they have shown that the Horton-Strahler number $H_n$ of those tries fulfils

$$\frac{H_n}{\log n} \to \frac{1}{\log \frac{1}{\min(p, 1-p)}}$$

in probability as $n \to \infty$. The presented (Bernoulli-) model of a random trie is very realsitic. For example, if we choose $p$ being $\frac{1}{2}$, this model describes exactly the behavior of tries built from random integer data assuming all integers to be equally likely.

In this paper we will adopt a combinatorial point of view by regarding all tree-structures that might be generated by the trie algorithm. For that purpose we

Figure 1: An example for a set of keys $K$, the resulting trie and the corresponding $\mathcal{C}$-trie.

Figure 2: Affecting additive parameters by adding a new node in between two existing ones. All trees shown in the figure have the same Horton-Strahler number but different path lengths.

will consider a class of generalized extended binary trees the so-called $\mathcal{C}$-tries (for combinatorial tries) which were introduced in [Neb99]. Let a $\mathcal{C}$-*trie* be an extended binary tree where each leaf is either colored black or white; each black leaf has to be the brother of an internal node. If we now interpret a white leaf as the location of a key and a black leaf as a NIL-pointer the $\mathcal{C}$-tries resemble all the tree-structures which the trie algorithm might generate. An example for that correspondence can be found in Figure 1. It is obvious that all the correspondences given above between the Horton-Strahler number and the parameters such as the number of registers needed for a syntax-tree evaluation or the stack-size for a traversal remain valid even by coloring the leaves. A $\mathcal{C}$-trie with $\alpha$ internal nodes (with $\alpha$ internal nodes and $\beta$ white leaves) is defined to be of size $\alpha$ (of size $(\alpha, \beta)$) and will be called $\alpha$-trie (($\alpha, \beta$)-trie). Note that $2 \leq \beta \leq \alpha + 1$ must hold. A $\mathcal{C}$-trie of size $(\alpha, \alpha + 1)$ is nothing else but an ordinary extended binary tree. For our investigations we will assume that all $\mathcal{C}$-tries of size $\alpha$ (resp. $(\alpha, \beta)$) are equally likely which is a quite different assumption compared to the before mentioned Bernoulli-model where it is more likely that a trie is balanced than being of a linear structure. This is why our investigation is rather of a combinatorial character.

Note that the idea of coloring leaves is not only useful for introducing a combinatorial equivalent for the tree-structures of binary tries. By introducing the number of white leaves as a second parameter we are more flexible to model natural phenomena. By varying the ratio of $\alpha$ and $\beta$ we can control the average shape of the related trees since a large number of black leaves can only exist if the tree has many linear lists within its inner structure. Thus we hope that the results presented in the rest of this paper will be of interest with respect to geology, molecular biology, synthetic images of trees, channel networks, ... Think for example of a river network modelled by an extended binary tree $T$ and the task to model different lengths of the rivers. By adding a new node with a black leaf as one of its successors in between two arbitrary existing nodes of $T$ (see Figure 2), we do not change the tree's Horton-Strahler number (as

3

Figure 3:
All possible decompositions of a $\mathcal{C}$-trie $T$ with $hs(T) = p$. The number
below a triangle determines the value of $hs$ of the subtrie represented.

it should be since the length of a river should not affect the classification of a
river network). However, we change length-sensitive parametres like e.g. the
external path length.
Thus, even if all our results are presented using the term $\mathcal{C}$-trie instead of gen-
eralized extended binary tree they should be considered as applicable to many
other areas.

## 2   The average Horton-Strahler number

The aim of this section is to derive the average Horton-Strahler number for
uniform random $\mathcal{C}$-tries as defined in Section 1, i.e. the average value of the
function $hs$ applied to a set of $\mathcal{C}$-tries of the same size $(\alpha, \beta)$. We will use gener-
ating functions in order to prove our results. The way these generating functions
are derived is similar to that in [FRV79], the methodology used to determine
asymptotics for the coefficients in question is standard and can be found in
[FlO90]. By $[x_1^{n_1} \cdots x_k^{n_k}] f(x_1, \ldots, x_k)$ we denote the coefficient at $x_1^{n_1} \cdots x_k^{n_k}$ in
an expansion of $f(x_1, \ldots, x_k)$ at $(x_1, \ldots, x_k) = (0, \ldots, 0)$.
We start our investigations by determining the generating function $H_p(x, y)$
which counts those $\mathcal{C}$-tries that have a Horton-Strahler number of exactly $p$.

**Lemma 1** *Let $x$ mark an internal node and let $y$ mark a white leaf. The gen-
erating function $H_p(x, y)$ of $\mathcal{C}$-tries $T$ with $hs(T) = p$ possesses the following
closed form representation:*

$$H_p(x, y) = \frac{\sin(\phi)}{\sin(2^{p-1}\phi)} \cdot \frac{xy^2}{1 - 2x - 2xy}, \tag{2}$$

$$\text{where} \qquad \phi = \arccos\left(\frac{1 - 4x(y + 1) + 2x^2(2 + y(4 + y))}{2x^2y^2}\right). \tag{3}$$

**Proof:** In order to derive a representation for the generating function in ques-
tion we have to distinguish the cases shown in Figure 3. For $p \geq 2$ these cases
translate into the following recurrence for $H_p(x, y)$:

$$H_p(x, y) = (2x + 2xy)H_p(x, y) + xH_{p-1}^2(x, y) + 2xH_p(x, y) \sum_{1 \leq j < p} H_j(x, y). \tag{4}$$

Here, the boundary condition for $p = 1$ is still to be determined yet. It is obvious
that a $\mathcal{C}$-trie $T$ with $hs(T) = 1$ has to have a linear structure, i.e. either the left

4

or the right subtrie of each internal node has to be a leaf. Thus $H_1(x, y)$ must fulfil $H_1(x, y) = xy^2 + (2x + 2xy)H_1(x, y)$ and therefore

$$H_1(x, y) = \frac{xy^2}{1 - 2x - 2xy}.$$

In order to solve this recurrence we divide both sides of (4) by $xH_p(x, y)$, thus

$$\frac{1}{x} = 2 + 2y + \frac{H_{p-1}^2(x, y)}{H_p(x, y)} + 2 \sum_{1 \le j < p} H_j(x, y).$$

Subtracting this from the analogous identity obtained for $p + 1$ eliminates the summation. We find

$$0 = \frac{H_p^2(x, y)}{H_{p+1}(x, y)} + 2H_p(x, y) - \frac{H_{p-1}^2(x, y)}{H_p(x, y)}. \tag{5}$$

Let $V_p(x, y) := \frac{H_{p-1}(x,y)}{H_p(x,y)}$. Dividing (5) by $H_p(x, y)$ our recurrence can be expressed by means of $V_p$:

$$V_{p+1}(x, y) = V_p^2(x, y) - 2, \, p \ge 2,$$

with the initial condition $V_2(x, y) = \frac{H_1(x,y)}{H_2(x,y)}$. We can determine $H_2(x, y)$ by using $H_1(x, y)$ and the recurrence (4) which yields

$$V_2(x, y) = \frac{1 - 4x(y + 1) + 2x^2(2 + y(4 + y))}{x^2 y^2}.$$

This new recurrence can be solved by a trigonometric change of variables. We set $V_2(x, y) = 2\cos(\phi)$ and generally $V_p(x, y) = 2\cos(\phi_p)$. Since $\cos^2(x) = \frac{1}{2}(1 + \cos(2x))$ holds we see that the recurrence translates into

$$2\cos(\phi_{p+1}) = 2\cos(2\phi_p).$$

Therefore, for $p \ge 2$, $\phi_{p+1} = 2\phi_p = 2^{p-1}\phi$ must hold which gives the explicit form

$$\begin{aligned}
V_{p+1}(x, y) &= 2\cos(2^{p-1}\phi), \, p \ge 2, \\
V_2(x, y) &= 2\cos(\phi) = \frac{1 - 4x(y + 1) + 2x^2(2 + y(4 + y))}{x^2 y^2}.
\end{aligned}$$

We can go back to $H_p(x, y)$ by regarding

$$V_p(x, y)V_{p-1}(x, y) \cdots V_2(x) = \frac{H_{p-1}(x, y)}{H_p(x, y)} \frac{H_{p-2}(x, y)}{H_{p-1}(x, y)} \cdots \frac{H_1(x, y)}{H_2(x, y)} = \frac{H_1(x, y)}{H_p(x, y)}.$$

5

By means of the identity $\sin(2x) = 2\sin(x)\cos(x)$ the product on the left-hand side collapses to $\sin(2^{p-1}\phi)$ when multiplied by $\sin(\phi)$. This completes the proof.
□

Next we consider those $\mathcal{C}$-tries that have a Horton-Strahler number of at least $p$.

**Lemma 2** *Let* $S_p(x,y) := \sum_{j \geq p} H_j(x,y)$, $\kappa := \frac{x^2 y^2}{(1-2x-2xy)^2}$, $\varepsilon := \sqrt{1-4\kappa}$ *and* $u := \frac{1-\varepsilon}{1+\varepsilon}$. *We have*

$$S_p(x,y) = y \left[ \frac{\sqrt{u}(1-u)}{u} \cdot \frac{u^{2^{p-1}}}{1 - u^{2^{p-1}}} \right].$$

**Proof:** In order to prove the lemma we use the identity $\sin(x) = \frac{1}{2i}\left(e^{ix} - e^{-ix}\right)$, $i^2 = -1$, which we insert in (2). Together with $t := e^{-i\phi}$ and $r := 2^{p-1}$ we find

$$H_p(x,y) = 2i \sin(\phi) \frac{t^r}{1 - t^{2r}} \frac{xy^2}{1 - 2x - 2xy}.$$

Now consider those parts of the representation that depend on $p$. Summing them up for $j \geq p$ yields

$$\sum_{\substack{j \geq p \\ r=2^{j-1}}} \frac{t^r}{1 - t^{2r}} = \sum_{\substack{j \geq p \\ k \geq 0}} t^{2^{j-1}(1+2k)} = \sum_{\substack{m,k \geq 0 \\ v=t^{2^{p-1}}}} v^{2^m(2k+1)}.$$

Since the mapping $(m,k) \to 2^m(2k+1)$ is a bijection on $\mathbb{N}^2 \to \mathbb{N}$ the last sum equals $\frac{v}{1-v}$. Therefore $S_p(x,y) = 2i \sin(\phi) \frac{xy^2}{1-2x-2xy} \frac{t^r}{1-t^r}$ holds. Returning to trigonometric functions, i.e. setting $t = e^{-i\phi} = \cos(\phi) - i\sin(\phi)$, gives us

$$S_p(x,y) = \left[ -i\sin(\phi) + \frac{\sin(\phi)\cos(2^{p-2}\phi)}{\sin(2^{p-2}\phi)} \right] \frac{xy^2}{1 - 2x - 2xy}.$$

For $p > 1$ it is now possible to express the trigonometric functions by means of Chebyshev polynomials. For $p = 1$ we run into trouble since in that case we would refer to the $\frac{1}{2}$-th polynomial which does not exist. Thus, the next step is to express $S_p(x,y)$ by means of $\frac{1}{2}\phi$ instead of $\phi$. By applying the identity $2\cos(x)\sin(x) = \sin(2x)$ we find that

$$S_p(x,y) = \frac{-i\sin(\phi)xy^2}{1 - 2x - 2xy} + \frac{\cos(\frac{1}{2}\phi)\sin(\frac{1}{2}\phi)\cos(2^{p-1}\frac{1}{2}\phi)}{\sin(2^{p-1}\frac{1}{2}\phi)} \frac{2xy^2}{1 - 2x - 2xy}$$

holds. From equation (3) we derive closed form expressions for $i\sin(\phi)$ and $\cos(\frac{1}{2}\phi)$ which we insert into the last representation of $S_p(x,y)$. Then, applying the following identities for the Chebyshev polynomial of the first-kind $(T_n(x))$

6

(see e.g. [AbSt70] 22.3.15) and the second-kind $(U_n(x))$ (see e.g. [AbSt70] 22.3.16)

$$T_n(\cos(\phi)) = \cos(n\phi),$$

$$U_n(\cos(\phi)) = \frac{\sin((n+1)\phi)}{\sin(\phi)},$$

yield

$$S_p(x,y) = -\frac{\sqrt{(1-2x)(1-2x-4xy)}}{2x} - \frac{yT_{2^{p-1}}(\hat{\kappa})}{U_{2^{p-1}-1}(\hat{\kappa})}.$$

Here $\hat{\kappa} := \cos(\frac{1}{2}\phi) = -\frac{1-2x-2xy}{2xy}$ holds. Now let $T(x,y)$ be the ordinary generating function of all $\mathcal{C}$-tries. In [Neb00] the following representation can be found:

$$T(x,y) = \frac{1-2x-\sqrt{(1-2x)(1-2x-4xy)}}{2x}.$$

This, together with two further identities for Chebyshev polynomials

$$T_n(x) = U_n(x) - xU_{n-1}(x), \text{ (see e.g. [AbSt70, 22.5.6])},$$

$$U_{n+1}(x) = 2xU_n(x) - U_{n-1}(x), \text{ (see e.g. [Kem84, (B77)])},$$

gives us

$$S_p(x,y) = T(x,y) - y + \frac{yU_{2^{p-1}-2}(\hat{\kappa})}{U_{2^{p-1}-1}(\hat{\kappa})},$$

for $U_{-1}(x) := 0$. A closed form representation for $U_n(x)$ is given in [Kem84, (B74)]. By fundamental algebraic manipulations this representation can be transformed into

$$U_n(x) = -\frac{x^n\left[(1-\sqrt{1-x^{-2}})^{n+1} - (1+\sqrt{1-x^{-2}})^{n+1}\right]}{2\sqrt{1-x^{-2}}}.$$

Now, since $1 - \hat{\kappa}^{-2} = 1 - 4\kappa$ holds, we get

$$S_p(x,y) = T(x,y) - y - 2y\sqrt{\kappa}\frac{(1+\sqrt{1-4\kappa})^{2^{p-1}-1} - (1-\sqrt{1-4\kappa})^{2^{p-1}-1}}{(1+\sqrt{1-4\kappa})^{2^{p-1}} - (1-\sqrt{1-4\kappa})^{2^{p-1}}}.$$

We complete the proof by using the substitutions of the lemma and applying some obvious simplifications. □

**Remark:** Besides the Horton-Strahler number there is another monotonic marking of binary trees which is related to the evaluation of arithmetic expressions and the traversal. This is the so called *stack-number* of the tree which was investigated in numerous papers (e.g. [BKR72], [Kem80], [Kem92], [Neb97], [Neb99] and [Neb00]). It corresponds to the stack-size needed to traverse a tree in preorder using the traditional algorithm (see e.g. [Knu97, pp. 319ff]) and the number of cells on a stack needed to evaluate an arithmetic expression by means

of a simple traversal algorithm (see [Kem84] for details). For non-colored extended binary trees we have the following correspondence: The number of trees with $\alpha$ internal nodes, with a stack-number of at most $2^k - 1$, is equal to the number of trees with $\alpha$ internal nodes and a Horton-Strahler number of $k$ (see e.g. [Kem84] Theorem 5.8). If we inspect the generating functions of the previous lemma and of [Neb99] and [Neb00] we see that such a relation does not exist for $\mathcal{C}$-tries neither of size $\alpha$ nor of size $(\alpha, \beta)$.

In order to compute the average Horton-Strahler number we need a representation of the generating function $M(x, y) := \sum_{p \geq 1} p H_p(x, y)$. It is not hard to see that $M(x, y) = \sum_{p \geq 1} S_p(x, y)$ holds. Therefore we have

$$
\begin{aligned}
M(x, y) &= \frac{y \sqrt{u}(1 - u)}{u} \sum_{p \geq 1} \frac{u^{2^{p-1}}}{1 - u^{2^{p-1}}} \\
&= \frac{y \sqrt{u}(1 - u)}{u} \sum_{n \geq 1} (v_2(n) + 1) u^n.
\end{aligned}
\tag{6}
$$

Here $v_2(n)$ denotes the *dyadic valuation* of $n$, i.e. the number of positive divisors of $n$ which are a power of two.

Now, everything is prepared to determine an asymptotic equivalent for the average Horton-Strahler number. We use the *Mellin summation method* as described in [FGD95] to evaluate the sum. For that purpose we set $u = \exp(-t)$ and apply the well-known identity

$$
\exp(-tj) = \frac{1}{2\pi i} \int_{c-i\infty}^{c+i\infty} \Gamma(s) j^{-s} t^{-s} ds, \; i^2 = -1,
$$

for some $c$ in the fundamental strip of the Mellin transform of $\exp(-tj)$ and for $\Gamma(s)$ the complete gamma function. This is how it is possible to express the number-theoretic function $v_2(n)$ by means of the Riemann Zeta function $\zeta(z)$ (see [Apo76] for details). We have (see e.g. [Kem84, p. 155])

$$
\sum_{n \geq 1} v_2(n) n^{-z} = \sum_{n \geq 1} v_2(2n)(2n)^{-z} = \sum_{\substack{j \geq 1 \\ n \geq 1}} (2^j n)^{-z} = \zeta(z)(2^z - 1)^{-1}
$$

and therefore with $u = \exp(-t)$

$$
\sum_{n \geq 1} (v_2(n) + 1) u^n = \frac{1}{2\pi i} \int_{c-i\infty}^{c+i\infty} \zeta(s) 2^s \Gamma(s) t^{-s} (2^s - 1)^{-1} ds.
$$

Now, according to the Mellin summation formula, we have to sum the residues of $\zeta(s) 2^s \Gamma(s) t^{-s} (2^s - 1)^{-1}$ left to the fundamental strip, i.e. the residues with a real part less or equal to one. There are singularities at $s = 1$ and $s =$

8

$-n$, $n \in \mathbb{N} \cup \{0\}$, but we only have to consider those which are larger than $-1$ since the others will only imply terms that can be neglected. There are further singularities at $s = \frac{2\pi i k}{\ln(2)} =: \chi_k$, $k \in \mathbb{Z} \backslash \{0\}$, which would imply some oscillation in the lower order terms. As the known methods for multivariate asymptotics only allow to determine the leading term, the singularities $\chi_k$ will only be considered later in the univariate case. The sum of the residues for $s = 1$ and $s = 0$ is given by

$$\frac{2}{t} + \frac{2\ln(t) + 2\gamma - 2\ln(\pi) - 3\ln(2)}{4\ln(2)}.$$

Here, $\gamma = 0.5772156649...$ denotes Eulers's constant.

In order to approximate the coefficient of $M(x, y)$ at $x^\alpha y^\beta$ we are interested in an expansion of our generating function at its dominant singularity. For that purpose we assume $y$ being a positive constant not equal to $0$ in order to determine the dominant singularity with respect to $x$, i.e. the value of $x$ which has the smallest modulus and which is a singular point of $M(x, y)$. Note that this approach leads to the restriction that our asymptotic will only be valid when $\alpha$ and $\beta$ grow simultaneously in a fixed proportion. By definition of $M(x, y)$ and properties of the Horton-Strahler number we have the following trivial bounds for $[x^\alpha]M(x, y)$:

$$[x^\alpha]T(x, y) \leq [x^\alpha]M(x, y) \leq \log_2(\alpha + 1)[x^\alpha]T(x, y).$$

Under the assumption that $y$ is a constant and since $x = \frac{1}{2+4y}$ is the dominant singularity of $T(x, y)$ the $\mathcal{O}$-transfer method introduced in [FlO90] leads to

$$[x^\alpha]T(x, y) \sim \frac{1 + 2y}{2} \cdot \frac{(2 + 4y)^\alpha}{\sqrt{\pi \alpha^3}}.$$

Therefore the Cauchy-Hadamard formula tells us that both the minorant and the majorant of $M(x, y)$ have a radius of convergence of $\frac{1}{2+4y}$ and we infer that $\frac{1}{2+4y}$ is the radius of convergence of $M(x, y)$ itself. Thus, by the theorem of Pringsheim, we can conclude that $x = \frac{1}{2+4y}$ is a dominant singularity of our generating function. It resides to prove that it is the only dominant singularity. For that purpose we consider the representation (6) of $M(x, y)$. Besides the algebraic singularity at $x = \frac{1}{2+4y}$ implied by our substitution, the factor of the sum is only singular for $u = 0$ i.e. for $x = 0$. This is why it does not extend the set of dominant singularities. The sum $\sum_{n \geq 1}(v_2(n) + 1)u^n$ possesses the minorant $\sum_{n \geq 1} u^n$ and the majorant $\sum_{n \geq 1}(n + 1)u^n$ both with a radius of convergence equal to 1. Therefore the set of solutions of $|u| = 1$ might contribute further dominant singularities. However, $|u| = 1$ has only one solution with an appropriate modulus, namely $x = \frac{1}{2+4y}$. Thus we can conclude that there is only one dominant singularity. Since $t = -\log\left(\frac{1-\varepsilon}{1+\varepsilon}\right)$ and $\varepsilon$ becomes $0$ at our

9

dominant singularity we expand $-\log\left(\frac{1-\varepsilon}{1+\varepsilon}\right)$ about $\varepsilon = 0$ to get $t \sim 2\varepsilon$ and thus $t \sim 2\frac{\sqrt{(1-2x)(1-2x-4xy)}}{1-2x-2xy}$. We conclude that for an expansion at $x = \frac{1}{2+4y}$, $t$ complies with $2\frac{\sqrt{2y(1+2y)}}{y}\sqrt{1 - x(2 + 4y)}$. On the assumption that $y$ is constant and for $u = \exp(-t)$, the factor $\frac{y\sqrt{u}(1-u)}{u}$ possesses the expansion $yt + \mathcal{O}(t^2)$. So, the most significant term of the expansion of $M(x,y)$ around $x = \frac{1}{2+4y}$ is given by

$$\frac{yt\ln(t)}{2\ln(2)} \sim -\frac{\sqrt{2y(1+2y)}}{2\log(2)}\sqrt{1 - x(2+4y)}\log((1 - x(2+4y))^{-1}). \qquad (7)$$

This representation can be used to approximate the coefficients of $M(x,y)$. We find:

**Lemma 3** *Let $\rho := \frac{\alpha}{\beta}$ be fixed. The coefficient of $M(x,y)$ at $x^\alpha y^\beta$ is asymptotically given by*

$$\frac{2^{\alpha(1+\frac{1}{\rho})-2}}{\sqrt{\pi\alpha^3}}\binom{\alpha + \frac{1}{2}}{\frac{\alpha}{\rho} - \frac{1}{2}}\log_2(\alpha),$$

$\alpha \to \infty$.

**Proof:** We use the following well-known expansions

$$\log((1 - x(2+4y))^{-1}) = \sum_{n\geq 1}\frac{x^n}{n}\sum_{k\geq 0}\binom{n}{k}2^{n-k}4^k y^k,$$

$$\sqrt{1 - x(2+4y)} = \sum_{i\geq 0}\binom{\frac{1}{2}}{i}(-1)^i x^i \sum_{k\geq 0}\binom{i}{k}2^{i-k}4^k y^k,$$

$$\sqrt{2y + 4y^2} = \sum_{j\geq 0}\binom{\frac{1}{2}}{j}2^{1-j}y^{1-j},$$

and extract the coefficient at $x^\alpha y^\beta$ in the resulting expansion of the right-hand side of (7). We find

$$[x^\alpha y^\beta]M(x,y) \sim -\frac{2^{\alpha+\beta-1}}{\log(2)}\binom{\alpha + \frac{1}{2}}{\beta - \frac{1}{2}}\underbrace{\sum_{n\geq 1}\frac{(-1)^{\alpha-n}}{n}\binom{\frac{1}{2}}{\alpha - n}}_{=:\sigma(\alpha)}.$$

By induction on $\alpha$ it is possible to prove the following recursion for $\sigma(\alpha)$:

$$\sigma(0) = 0,$$
$$\sigma(\alpha) = \frac{2\alpha - 3}{2\alpha}\sigma(\alpha - 1) + \underbrace{\frac{(-1)^{\alpha+1}\sqrt{\pi}}{2\Gamma(\frac{5}{2} - \alpha)\Gamma(\alpha + 1)}}_{=:\varsigma(\alpha)}.$$

10

This recursion can be solved by using ordinary generating functions. For $A(z) := \sum_{\alpha \geq 0} \sigma(\alpha) z^\alpha$ we get $A(z) = zA(z) - \frac{3}{2} \int_0^z A(t) dt + \sum_{\alpha \geq 1} \varsigma(\alpha) z^\alpha$. Applying the identity $\sum_{\alpha \geq 1} \varsigma(\alpha) z^\alpha = \frac{2}{3}\sqrt{1-z}(z-1) + \frac{2}{3}$, which for instance Zeilberger's "fast algorithm" (see [Zei91]) finds for you, yields a simple differential equation for $A(z)$ which possesses the solution $A(z) = -\sqrt{1-z}\log(1-z)$ and thus

$$\sigma(\alpha) = [z^\alpha]\sqrt{1-z}\log((1-z)^{-1})$$

holds. By applying the $\mathcal{O}$-transfer method we find the approximation $\sigma(\alpha) \sim -\frac{\log(\alpha)}{2\sqrt{\pi \alpha^3}}$ which proves the lemma. $\qquad\square$

In order to get the average value, the coefficient given in Lemma 3 has to be divided by the total number of $\mathcal{C}$-tries of size $(\alpha, \beta)$. We can proceed in the same way as done in the previous proof in order to approximate the coefficient $[x^\alpha y^\beta] T(x,y)$. Thus we factor $T(x,y)$ into

$$\frac{1-2x}{2x} - \frac{1-2x}{2x}\sqrt{1 - 4xy(1-2x)^{-1}} \qquad (8)$$

which corresponds to an expansion of $T(x,y)$ around the singularity $y = \frac{1-2x}{4x}$ (assuming now that x is constant). Since the leftmost term of (8) only possesses coefficients at $y^0$ it can be neglected. Furthermore, $\sqrt{1 - 4xy(1-2x)^{-1}}$ can be expanded using the binomial theorem, which yields

$$\xi(\alpha, \beta) := [x^\alpha y^\beta]\sqrt{1 - 4xy(1-2x)^{-1}} = \binom{\frac{1}{2}}{\beta}\binom{\alpha - 1}{\alpha - \beta}(-1)^\beta 4^\beta 2^{\alpha - \beta}.$$

Now, taking the factor $-\frac{1-2x}{2x}$ into account proves that $[x^\alpha y^\beta] T(x,y) \sim \xi(\alpha, \beta) - \frac{1}{2}\xi(\alpha+1, \beta)$. Thus we conclude that

$$[x^\alpha y^\beta] T(x,y) \sim \frac{2^{\alpha - \beta + 1}}{\beta}\binom{2\beta - 2}{\beta - 1}\binom{\alpha - 1}{\alpha - \beta + 1}.$$

Dividing the coefficient given in Lemma 3 by the previous quantity provides the following theorem:

**Theorem 1** *On the assumption that all $(\alpha, \beta)$-tries are equally likely the average Horton-Strahler number of a $\mathcal{C}$-trie of size $(\alpha, \beta)$ is asymptotically given by*

$$\frac{4^{\alpha(\frac{2}{\rho} - 1)}\log(\alpha)(2\frac{\alpha}{\rho} - 1)(\alpha + 1)(\alpha + 2)\binom{2\alpha + 1}{\alpha - 1}}{8\sqrt{\pi}\alpha^{3/2}\log(2)(\frac{\alpha}{\rho} - 1)\frac{\alpha}{\rho}\left(2\frac{\frac{\alpha}{\rho}}{\frac{\alpha}{\rho}}\right)^2},$$

$\rho := \frac{\alpha}{\beta}$ *fixed,* $\alpha \to \infty$. $\qquad\square$

**Remark:** The asymptotic given for the number of $(\alpha, \beta)$-tries is equal to the exact number of $\mathcal{C}$-tries of this size for $\alpha > 0$. This is due to the fact that we

find a factorization of $T(x, y)$ when expanding it around its dominant singularity. Thus, besides some terms at $y^0$, no terms were neglected when we have developed the leading term and have extracted the coefficients.

Looking at a plot of our average Horton-Strahler number of $\mathcal{C}$-tries (see the last section of this paper) it seems to be hardly dependent on $\beta$. This impression is justified when we use Stirling's formula to approximate the binomial coefficient $\binom{2\beta}{\beta}$ within our result. We find that the average Horton-Strahler number of $\mathcal{C}$-tries is asymptotically given by

$$\left(2 + \frac{1}{\beta - 1}\right) \frac{\log(\alpha)(\alpha + 1)(\alpha + 2)\sqrt{\pi}}{8 \cdot 4^\alpha \alpha^{\frac{3}{2}} \log(2)} \binom{2\alpha + 1}{\alpha - 1}.$$

Thus, only for very sparse $\mathcal{C}$-tries, i.e. $\mathcal{C}$-tries with few white leaves only, we have an influence of $\beta$ on the average Horton-Strahler number. But for every fixed $\rho$ and $\alpha \to \infty$ also $\beta$ tends to infinity and thus $\frac{1}{\beta - 1}$ becomes zero. Thus it becomes possible to express the average Horton-Strahler number of $(\alpha, \beta)$-tries on dependence of $\alpha$ only. We conclude this discussion by noting that

$$\lim_{\alpha \to \infty} \frac{(\alpha + 1)(\alpha + 2)\sqrt{\pi}}{4^{\alpha+1} \alpha^{\frac{3}{2}}} \binom{2\alpha + 1}{\alpha - 1} = \frac{1}{2}$$

holds. Thus we have the following corollary

**Corollary 1** *Under the assumption that the number of internal nodes $\alpha$ and the number of white leaves $\beta$ grow in some fixed proportion the average Horton-Strahler number of $(\alpha, \beta)$-tries is asymptotically given by*

$$\frac{\log(\alpha)}{2 \log(2)}.$$

$\square$

**Remark:** We can also conclude the result of the previous corollary by using the multivariate Darboux-method presented in [Drm94] in order to approximate the coefficient of the leading term in (7) and the number of $\mathcal{C}$-tries of size $(\alpha, \beta)$. In that case, because of side-conditions given by the method, $\rho = \frac{\alpha}{\beta}$ has to be strictly larger that one. However, it is impossible to derive the more accurate results presented in Lemma 3 and Theorem 1 in that way.

We will now use our generating functions to derive an asymptotic equivalent for the average Horton-Strahler number for $\mathcal{C}$-tries of size $\alpha$. As methodology is much more developed for univariate generating functions it is possible to derive results of higher precision. Note, that in the uniform model it is not possible to derive a univariate result with respect to the number of white leaves since all the generating functions would count infinitely many $\mathcal{C}$-tries of any given size $\beta$.

This is due to the fact, that the number of white leaves does not limit the number of internal nodes even if we fix the Horton-Strahler number of the $\mathcal{C}$-tries considered. Thus, we return to (6) and set $y = 1$ in all parts of the generating function. It is obvious that we find the same integral as a representation of $\sum_{n \geq 1}(v_2(n) + 1)u^n$ as in the bivariate case since the substitution $u = \exp(-t)$ yields the same result even if we set $y$ to 1. But now it makes sense to consider terms of lower significance also, because the $\mathcal{O}$-transfer method for univariate generating functions makes it possible to translate them into the right contributions for the asymptotics in question. Therefore we sum the residues of the singularities at $s \in \{1, -n, \chi_k\}$, $n \in \mathbb{N} \cup \{0\}$, $k \in \mathbb{Z} \backslash \{0\}$, and multiply them by the expansion of the factor $\frac{y\sqrt{u}(1-u)}{u}|_{y=1}$ which gives us

$$2 + \left( \frac{2\ln(t) + 2\gamma - 2\ln(\pi) - 3\ln(2)}{4\ln(2)} \right) t + \sum_{k \neq 0} \frac{\Gamma(\chi_k)\zeta(\chi_k)}{\ln(2)} t^{1-\chi_k} + \mathcal{O}(t^2). \quad (9)$$

Again, we are interested in an expansion at the dominant singularity which is $x = \frac{1}{2+4y}|_{y=1} = \frac{1}{6}$. Thus we have to set $t = 2\sqrt{6}\sqrt{1-6x}$ in order to resubstitute $t$ in (9) which yields the desired expansion:

$$-\frac{\sqrt{6}\sqrt{1-6x}\ln((1-6x)^{-1})}{2\ln(2)} + \frac{\sqrt{6}\sqrt{1-6x}(2\gamma - 2\ln(\pi) + \ln(3))}{2\ln(2)}$$

$$+ \sum_{k \neq 0} \frac{\Gamma(\chi_k)\zeta(\chi_k)}{\ln(2)}(2\sqrt{6})^{1-\chi_k}(1-6x)^{\frac{1-\chi_k}{2}} + \mathcal{O}(|1-6x|)$$

$$= -\frac{\sqrt{6}\sqrt{1-6x}\ln((1-6x)^{-1})}{2\ln(2)} + \frac{\sqrt{6}\sqrt{1-6x}(2\gamma - 2\ln(\pi) + \ln(3))}{2\ln(2)}$$

$$+ \frac{2\sqrt{6}}{\ln(2)} \sum_{k \neq 0} \Gamma(\chi_k)\zeta(\chi_k)e^{-\log_2(6)\pi ik}(1-6x)^{\frac{1-\chi_k}{2}} + \mathcal{O}(|1-6x|).$$

Now we can apply the transfer formulæ according to [FlO90]. We have

$$[z^n](1-z)^\alpha \sim \frac{n^{-\alpha-1}}{\Gamma(-\alpha)} + \mathcal{O}(n^{-\alpha-2})$$

and

$$[z^n](1-z)^{\frac{1}{2}}\ln((1-z)^{-1}) \sim -\frac{1}{\sqrt{\pi n^3}} \left( \frac{1}{2}\ln(n) + \frac{\gamma + 2\log(2) - 2}{2} + \mathcal{O}(\frac{\ln(n)}{n}) \right)$$

which provide the following lemma:

**Lemma 4** *The coefficient of $M(x,1)$ at $x^\alpha$ is asymptotically given by*

$$\frac{6^{\frac{1}{2}+\alpha}(\ln(\alpha) - \gamma + 2\ln(2) - 2 + 2\ln(\pi) - \ln(3))}{4 \cdot \alpha^{\frac{3}{2}}\sqrt{\pi}\ln(2)}$$

$$+ \frac{2\sqrt{6}}{\ln(2)} \sum_{k \neq 0} \Gamma(\chi_k)\zeta(\chi_k)e^{-\log_2(6)\pi ik}6^\alpha \alpha^{\frac{\chi_k}{2}-\frac{3}{2}} / \Gamma(\frac{\chi_k}{2} - \frac{1}{2}) + \mathcal{O}(\alpha^{-\frac{5}{2}}).$$

□

Finally, this quantity has to be divided by the asymptotic number of $\alpha$-tries which is known (see [Neb99]) to be given by $6^{\frac{1}{2}+\alpha}\alpha^{-\frac{3}{2}}/(2\sqrt{\pi}) + \mathcal{O}(\alpha^{-\frac{5}{2}})$. After numerous simplifications we find:

**Theorem 2** *On the assumption that all $\mathcal{C}$-tries of the same size are equally likely the average Horton-Strahler number of an $\alpha$-trie is asymptotically given by*

$$\frac{1}{2}\log_2(\frac{4}{3}\pi^2\alpha) - \frac{\gamma+2}{2\ln(2)} + \Delta\left(\log_2\left(\frac{\alpha}{6}\right)\right) + \mathcal{O}(\alpha^{-1}),$$

$\chi_k = \frac{2\pi ik}{\ln(2)}$, $\alpha \to \infty$. *The function $\Delta(x)$ is a periodic function of small modulus ($|\Delta(x)| < 0.041$) and possesses the following representation as a Fourier series:*

$$\Delta(x) := \frac{1}{\ln(2)}\sum_{k\neq 0}(\chi_k - 1)\Gamma(\frac{\chi_k}{2})\zeta(\chi_k)e^{\pi ikx}.$$

□

**Remark:** Note that this is the same result as for non-colored extended binary trees given in (1) when setting $\alpha$ to $\frac{3}{2}\alpha$. The same effect with a different constant can be observed for the average stack-number. In that case we have to set $\alpha$ to $\frac{2}{3}\alpha$ in order to get the same leading term as for non-colored trees. The fact, that there are different constants for different parameters, supports a conjecture stated in [Neb99] which says that it seems to be impossible to conclude the behavior of $\mathcal{C}$-tries with respect to "traversal-parameters" from the well know results for ordinary extended binary trees (e.g. by a simple rearrangement together with an appropriate weighting of the trees). The bound $|\Delta(x)| < 0.041$ can be found by means of numerical studies.

# 3 Visualization and conclusions

Figure 4: The average stack-number (upper graph) and Horton-Strahler number (lower graph) on dependence of the number of internal nodes $\alpha$.

Figure 5: The ratio of the average Horton-Strahler number and the average stack-number.

In this section we will provide some plots of the results presented in [Neb99], in [Neb00] and in this paper. This is how we are going to compare the average stack-size with the average Horton-Strahler number which are related in the following way: If we think of applications of the Horton-Strahler number such

as a tree traversal or the evaluation of an arithmetic expression, the stack-size of the corresponding tree describes the amount of space needed when we apply a usual preorder traversal or a simple traversal strategy for evaluation. Those methods can be optimized with respect to the amount of space needed by easing the restriction that subtrees must be visited in a fixed order. The space requirement of the resulting strategy is described by the Horton-Strahler number. Therefore, we will speak of an *economy of space* when comparing both parameters.

The first plot is presented in Figure 4. It shows the absolute values of the

Figure 6: A plot of the average Horton-Strahler number on dependence of $\rho$ and $\beta$.

Figure 7: The difference of the average stack-number and the average Horton-Strahler number on dependence of $\rho$ and $\beta$.

average stack-number and the average Horton-Strahler number. As we can see the order of growth of both graphs is quite different. Even if the total stack-number for $\mathcal{C}$-tries of size $\alpha$ is small, the relative economy of space implied by the application of the optimized algorithms related to the Horton-Strahler number is remarkable (as we can see in Figure 5). In the bivariate setting a similar behavior can be found. If we take a look at Figure 6 we see that the average Horton-Strahler number for $\mathcal{C}$-tries of size $(\alpha, \beta)$ grows slowly and is of small value even for *sparse* $\mathcal{C}$-tries, i.e. for $\mathcal{C}$-tries with a large internal structure but only a few white leaves. As we would have expected from the univariate case, Figure 7 shows that the economy of space grows the larger the $(\alpha, \beta)$-tries become. We can also observe that the advantage of the optimized algorithms gets larger with $\rho$ growing and not only the relative but also the total economy of space get large when the $\mathcal{C}$-tries become sparse. For example on the assumption of $\rho = 8$ the total economy of space is about 60 for $(\alpha, \beta)$-tries with only 20 white leaves.

# References

[AbSt70]   M. Abramowitz and I. A. Stegun, *Handbook of Mathematical Functions*, Dover (1970)

[Apo76]     T. M. APOSTOL, *Introduction to Analytic Number Theory*, Springer, 1976

[BKR72]     N. G. DE BRUIJN, D. E. KNUTH, S.O. RICE, *The Average Height of Planted Plane Trees*, Graph Theory and Computing (R. C. Read, ed.), Academic Press, 1972

[DeK96]     L. DEVROYE AND P. KRUSZEWSKI, *On the Horton-Strahler Number for Random Tries*, R.A.I.R.O. Theoretical Informatics and Applications **30**, 443-456, 1996

[Drm94]     M. DRMOTA, *Asymptotic Distributions and a Multivariate Darboux Method in Enumeration Problems*, Journal of Combinatorial Theory, Series A **67**, 169-184, 1994

[Ers58]     A. P. ERSHOV, *On Programming of Arithmetic Operations*, Communications of the ACM **1**, 3-6, 1958

[FRV79]     P. FLAJOLET, J.C. RAOULT AND J. VUILLEMIN, *The Number of Registers required for Evaluating Arithmetic Expressions*, Theoretical Computer Science **9**, 99-125, 1979

[FlO90]     P. FLAJOLET AND A. ODLYZKO, *Singularity Analysis of Generating Functions*, SIAM J. Disc. Math. **3**, No. 2, 216-240, 1990

[FlP86]     P. FLAJOLET AND H. PRODINGER, *Register Allocation for Unary-Binary Trees*, SIAM J. Comput. **15**, 629-640, 1986

[FGD95]     P. FLAJOLET, X. GOURDON AND P. DUMAS, *Mellin transforms and asymptotics: Harmonic sums*, Theoretical Computer Science **144**, 3-58, 1995

[Fra84]     J. FRANÇON, *Sur le nombre de registres nécessaires à l'évaluation d'une expression arithmétique*, R.A.I.R.O. Theoretical Informatics and Applications **18**, 355-364, 1984

[Hor45]     R. E. HORTON, *Erosioned development of systems and their drainage basins, hydrophysical approach to quantitative morphology*, Bull. Geol. Soc. of America **56**, 275-370, 1945

[Kem79]     R. KEMP, *The Average Number of Registers Needed to Evaluate a Binary Tree Optimally*, Acta Informatica **11**, 363-372, 1979

[Kem80]     R. KEMP, *A Note on the Stack Size of Regularly Distributed Binary Trees*, BIT **20**, 157-163, 1980

[Kem84]     R. KEMP, *Fundamentals of the Average Case Analysis of Particular Algorithms*, Wiley-Teubner Series in Computer Science, 1984

[Kem92]    R. KEMP, *On the Stack Ramification of Binary Trees*, Random Graphs **2**, 117-138, 1992

[Knu97]    D. E. KNUTH, *The Art of Computer Programming*, Volume 1: Fundamental Algorithms, 3rd ed., Addison-Wesley, 1997

[MMP80]  A. MEIR, J. W. MOON AND J. R. POUNDER, *On the Order of Random Channel Networks*, SIAM J. Alg. Disc. Math. **1**, 25-33, 1980

[Neb97]    M. E. NEBEL, *New Results on the Stack Ramification of Binary Trees*, Journal of Automata, Languages and Combinatorics **2**, 161-175, 1997

[Neb99]    M. E. NEBEL, *The Stack-Size of Tries, A Combinatorial Study*, Theoretical Computer Science, to appear

[Neb00]    M. E. NEBEL, *The Stack-Size of Uniform Random Tries Revisited*, submitted, 2000

[Str52]      A. N. STRAHLER, *Hypsometric (area-altitude) analysis of erosonal topology*, Bull. Geol. Soc. of America **63**, 1117-1142, 1952

[VEJA89] X. G. VIENNOT, G. EYROLLES, N. JANEY AND D. ARQUES, *Combinatorial analysis of ramified patterns and computer imagery of trees*, Proc. SIGGRAPH'89, Computer Graphics **23**, 31-40, 1989

[Vie90]     X. G. VIENNOT, *Trees Everywhere*, Proc. CAAP'90, LNCS **431**, 18-41, 1990

[Zei91]     D. ZEILBERGER, *The method of creative telescoping*, Journal of Symbolic Computation **11**, 195-204, 1991