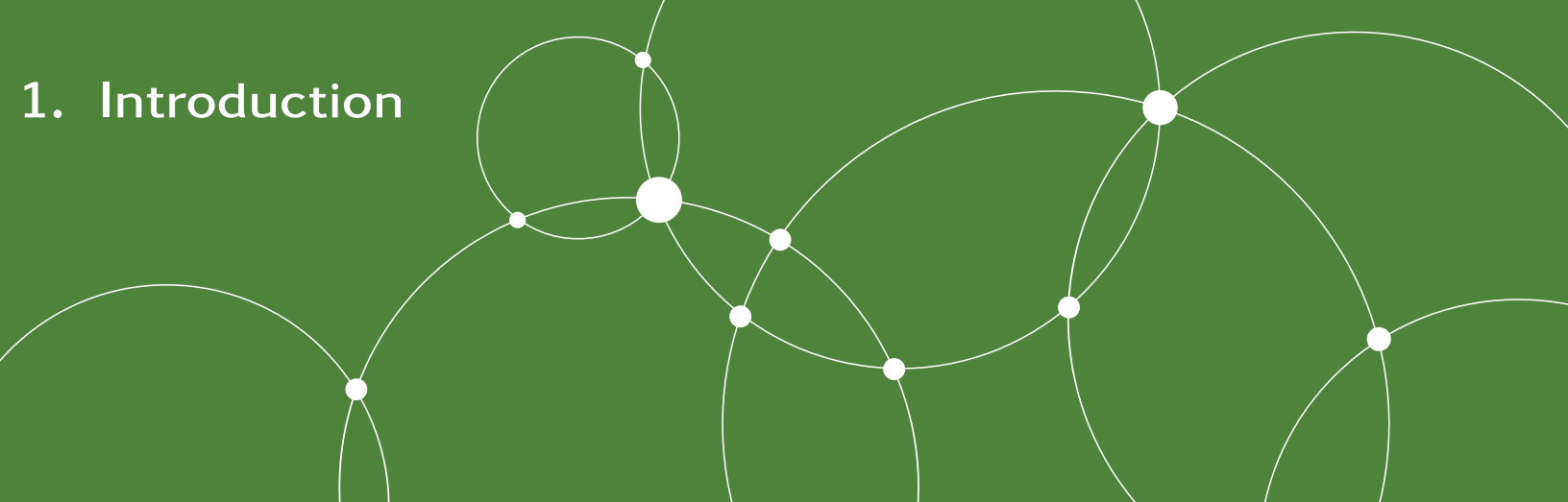




Contents

1	Introduction	3
1.1	Motivation	3
1.2	Related Research	3
2	The Collatz Tree	7
2.1	The Connection between Groups and Graphs	7
2.2	Defining the Tree	8
2.3	Relationship of successive nodes in $H_{C,3}$	10
2.4	A note on the functions left-child and right-sibling	13
2.5	Relationship of sibling nodes in $H_{C,3}$	14
2.6	A vertex's left-child, n -fold right-sibling in $H_{C,3}$	14
2.7	A vertex's left-child, n -fold right-sibling in $H_{C,5}$	17
2.8	A vertex's left-child, n -fold right-sibling in $H_{C,7}$	18
2.9	Generalizing the relationship of sibling nodes for $H_{C,k}$	19
2.10	Generalizing a vertex's left-child, n -fold right-sibling for $H_{C,k}$	20
3	Binary Collatz Tree	21
3.1	Some essentials on binary trees	21
3.2	Transforming the Collatz tree into a binary tree	21
4	Cycles in the Collatz Tree	25
4.1	A remark about cycles	25
4.2	Which variants of H_C have non-trivial cycles?	26
5	Conclusion and Outlook	29
5.1	Summary	29
5.2	Further Research	29

A	Appendix	31
A.1	A brief note on the tree $H_{C,1}$	31
A.2	The sum of reciprocated vertices depending only on v_1	31
A.3	The product of reciprocated vertices incremented by one	32
	Literature	33
	About our approach	39
	Acknowledgements	41
	Communities	41
	Contributors in the same field	41
	About Us	43



1. Introduction

It is well known that the inverted Collatz sequence can be represented as a graph or a tree. Similarly, it is acknowledged that in order to prove the Collatz conjecture, one must demonstrate that this tree covers all odd natural numbers. A structured reachability analysis is hitherto not available. This paper investigates the problem from a graph theory perspective. We define a tree that consists of nodes labeled with Collatz sequence numbers. This tree will be transformed into a sub-tree that only contains odd labeled nodes. Furthermore, we derive and prove several formulas that can be used to traverse the graph. The analysis covers the Collatz problem both in its original form $3x + 1$ as well as in the generalized variant $kx + 1$. Finally, we transform the Collatz graph into a binary tree, following the approach of Kleinnijenhuis, which could form the basis for a comprehensive proof of the conjecture.

1.1 Motivation

The Collatz conjecture is a number theoretical problem, which has puzzled countless researchers using myriad approaches. Presently, there are scarcely any methodologies to describe and treat the problem from the perspective of the Algebraic Theory of Graphs. Such an approach is promising with respect to facilitating the comprehension of the Collatz sequence's "mechanics".

The current gap in research forms the motivation behind the present contribution. The authors are convinced that exploring the Collatz conjecture in an algebraic manner, relying on the findings and fundamentals of Graph Theory, will contribute to a simplification of the problem as a whole.

Key results of this manuscript have been achieved using Data Science techniques. Our main tool was a Python-API, which implements the later introduced theorems and is optimized for processing arbitrarily big integers [1], see chapter [About our approach](#).

1.2 Related Research

The following literature study is largely based on one given by a similar earlier essay [2] which deals with the Collatz conjecture from the vantage of automata theory.

The Collatz conjecture is one of the unsolved "Million Buck Problems" [3]. When Lothar Collatz began his professorship in Hamburg in 1952, he mentioned this problem to his colleague Helmut Hasse. From 1976 to 1980, Collatz wrote several letters but missed referencing that he first proposed the problem in 1937. He introduced a function $g : \mathbb{N} \rightarrow \mathbb{N}$ as

follows:

$$g(x) = \begin{cases} 3x+1 & 2 \nmid x \\ x/2 & \text{otherwise} \end{cases} \quad (1.1)$$

This function is surjective, but it is not injective (for example $g(3) = g(20)$) and thus is not reversible. The Collatz conjecture states that for each start number $x_1 > 0$ the sequence $x_1, x_2 = g(x_1), x_3 = g(x_2), \dots$ will at some point enter the so called trivial cycle 1, 4, 2. One example is the sequence 17, 52, 26, 13, 40, 20, 10, 5, 16, 8, 4, 2, 1 starting at $x_1 = 17$. The assumption has not yet been proven. If the conjecture were wrong, then for a starting number x_1 the sequence either would diverge indefinitely or enter a cycle different from the trivial one (a so called non-trivial cycle).

In order to specify compressed Collatz sequences containing only the odd members, Bruckman [4] for instance used the more convenient function that opts out all even integers:

$$f(x) = (3x+1) \cdot 2^{-\alpha(x)}, \text{ where } 2^{\alpha(x)} \parallel (3x+1) \quad (1.2)$$

Note that $\alpha(x)$ is the largest possible exponent for which $2^{\alpha(x)}$ exactly divides $3x+1$. Especially for prime powers, one often says p^α divides the integer x exactly, denoted as $p^\alpha \parallel x$, if p^α is the greatest power of the prime p that divides x .

In his book “The Ultimate Challenge: The $3x+1$ Problem” [5], along with his annotated bibliographies [6], [7] and other manuscripts like an earlier paper from 1985 [8], Lagarias researched and put together different approaches from various authors intended to describe and solve the Collatz conjecture.

For the integers up to 2,367,363,789,863,971,985,761 the conjecture holds valid. For instance, see the computation history given by Kahermanes [9] that provides a timeline of the results which have already been achieved.

Inverting the Collatz sequence and constructing a Collatz tree is an approach that has been carried out by many researchers. It is well known that inverse sequences [10] arise from all functions $h \in H$, which can be composed of the two mappings $q, r : \mathbb{N} \rightarrow \mathbb{N}$ with $q : m \mapsto 2m$ and $r : m \mapsto (m-1)/3$:

$$H = \{h : \mathbb{N} \rightarrow \mathbb{N} \mid h = r^{(j)} \circ q^{(i)} \circ \dots, i, j, h(1) \in \mathbb{N}\}$$

An argumentation that the Collatz Conjecture cannot be formally proved can be found in the work of Craig Alan Feinstein [11], who presents the position that any proof of the Collatz conjecture must have an infinite number of lines and thus no formal proof is possible. However, this statement will not be acknowledged in depth within this study.

Treating Collatz sequences in a binary system can be performed as well. For example, Ethan Akin [12] handles the Collatz sequence with natural numbers written in base 2 (using the Ring \mathbb{Z}_2 of two-adic integers), because divisions by 2 are easier to deal with in this method. He uses a shift map σ on \mathbb{Z}_2 and a map τ :

$$\sigma(x) = \begin{cases} (x-1)/2 & 2 \nmid x \\ x/2 & \text{otherwise} \end{cases} \quad \tau(x) = \begin{cases} (3x+1)/2 & 2 \nmid x \\ x/2 & \text{otherwise} \end{cases}$$

The shift map’s fundamental property is $\sigma(x)_i = x_{i+1}$, noting that $\sigma(x)_i$ is the i -th digit of $\sigma(x)$. This property can easily be comprehended by an example $x = 5 = 1010000\dots = x_0x_1x_2\dots$, containing $\sigma(x) = 2 = 0100000\dots$.

Akin then defines a transformation $Q : \mathbb{Z}_2 \rightarrow \mathbb{Z}_2$ by $Q(x)_i = \tau^i(x)_0$ for non-negative integers i which means $Q(x)_i$ is zero if $\tau^i(x)$ is even and then it is one in any other instance. This transformation is a bijective map that defines a conjugacy between τ and σ : $Q \circ \tau = \sigma \circ Q$ and it is equivalent to the map denoted Q_∞ by Lagarias [8] and it is the inverse of the map Φ introduced by Bernstein [13]. Q can be described as follows: Let x be a 2-adic integer. The transformation result $Q(x)$ is a 2-adic integer y , so that $y_n = \tau^{(n)}(x)_0$. This means, the first bit y_0 is the parity of $x = \tau^{(0)}(x)$, which is one, if x is odd and otherwise zero. The next bit y_1 is the parity of $\tau^{(1)}(x)$, and the bit after next y_2 is parity of $\tau \circ \tau(x)$ and so on. The conjugacy $Q \circ \tau = \sigma \circ Q$ can be demonstrated by transforming the expression as follows: $(\sigma \circ Q(x))_i = Q(x)_{i+1} = \tau^{(i+1)}(x)_0 = \tau^{(i)}(\tau(x))_0 = Q(\tau(x))_i$

A simulation of the Collatz function by Turing machines has been presented by Michel [14]. He introduces Turing machines that simulate the iteration of the Collatz function, where he considers them having 3 states and 4 symbols. Michel examines both turing machines, those that never halt and those that halt on the final loop.

A function-theoretic approach to this problem has been provided by Berg and Meinardus [15], [16] as well as Gerhard Opfer [17], who consistently relies on the Berg's and Meinardus' idea. Opfer tries to prove the Collatz conjecture by determining the kernel intersection of two linear operators U, V that act on complex-valued functions. First he determined the kernel of V , and then he attempted to prove that its image by U is empty. Benne de Weger [18] contradicted Opfer's attempted proof.

At the number of divisions by two Paul S. Bruckman [4] has taken a deeper look, who has attempted to provide an elementary proof by contradiction. He repeatedly applies the Collatz function using a starting value n_0 and defines:

$$\{e_k\} : n_1 = (3n_0 + 1) \cdot 2^{-e_1}, n_2 = (3n_1 + 1) \cdot 2^{-e_2} = (3^2 n_0 + 3 + 2^{e_1}) \cdot 2^{-(e_1+e_2)}, \dots$$

Denoting the sum of exponents as $E_k = e_1 + e_2 + \dots + e_k$ Bruckman obtains the following equation:

$$2^{E_k} n_k = 3^k n_0 + \sum_{j=0}^{k-1} 3^{k-1-j} 2^{E_j}$$

Reachability Considerations based on a Collatz tree exist as well. It is well known that the inverted Collatz sequence can be represented as a graph; to be more specific, they can be depicted as a tree [19], [20]. It is acknowledged that in order to prove the Collatz conjecture, one needs to demonstrate that this tree covers all odd natural numbers.

The Stopping Time theory has been introduced by Terras [21], it has been taken up and continued, inter alia, by Silva [22] and Idowu [23]. Terras introduces another notation of the Collatz function $T(n) = (3^{X(n)}n + X(n))/2$, where $X(n) = 1$ when n is odd and $X(n) = 0$ when n is even, and defined the stopping time of n , denoted by $\chi(n)$, as the least positive k for which $T^{(k)}(n) < n$, if it exists, or otherwise it reaches infinity. Let L_i be a set of natural numbers, it is observable that the stopping time exhibits the regularity $\chi(n) = i$ for all n fulfilling $n \equiv l \pmod{2^i}$, $l \in L_i$, $L_1 = \{4\}$, $L_2 = \{5\}$, $L_4 = \{3\}$, $L_5 = \{11, 23\}$, $L_7 = \{7, 15, 59\}$ and so on. As i increases, the sets L_i , including their elements, become significantly larger. Sets L_i are empty when $i \equiv l \pmod{19}$ for $l = 3, 6, 9, 11, 14, 17, 19$. Additionally, the largest element of a non-empty set L_i is always less than 2^i .

Dynamical systems provide a wide basis for examining the Collatz sequence as well [24]. A dynamical system [25, p. 464] is a triple (M, G, Φ) for a set M , a group $(G, +)$ and a map $\Phi : M \times G \rightarrow M$ for which $\Phi(\cdot, 0) = id_M(\cdot)$ firstly applies and secondly $\Phi(\Phi(m, s), t) = \Phi(m, s+t)$ for all $m \in M, s, t \in G$. The set M is called phase space. Terence Tao [26] considers orbits of the dynamical system generated by the Collatz map (an orbit, also called trajectory, is a subset of the phase space). For an integer function $f : \mathbb{Z} \rightarrow \mathbb{Z}$, we denote by $f^i = f \circ f^{i-1}$ the i -fold iterate of f with the convention $f^0 = id_{\mathbb{Z}}$. If $n \in \mathbb{Z}$, the orbit (trajectory) of n under f is the sequence $T_f(n) = (n, f(n), f \circ f(n), f \circ f \circ f(n), \dots)$, see [24, p. 10]. Tao proved that almost all of these orbits attain almost bounded values. To achieve this, he advanced the results of Allouche [27] and Korec [28]. Their main idea was to prove that the set of positive integers with finite stopping time has a density one, in this case the term density refers to the concept of *natural density* (also known as *asymptotic density*). It measures how large a subset of the set of natural numbers is. The natural density of a set $M \subseteq \mathbb{N}$ is defined as:

$$\lim_{n \rightarrow \infty} \frac{\#\{m \in M : m < n\}}{n}$$

In this context, the authors used the Collatz map as the map Φ . They proved that the set $\{x \in \mathbb{N} : (\exists t \in \mathbb{N})(\Phi(x, t) < x)\}$ has a natural density one.

Many other approaches exist as well. From an algebraic perspective Trümper [29] analyzes the Collatz problem in the light of an Infinite Free Semigroup. Kohl [30] generalized the problem by introducing residue class-wise affine mappings, in short rcwa mappings. A polynomial analogue of the Collatz Conjecture has been provided by Hicks et al. [31] [32] and there are also stochastic, statistical and Markov chain-based and permutation-based approaches to proving this elusive theory.



2. The Collatz Tree

2.1 The Connection between Groups and Graphs

Let (a_k) be a numerical sequence with $a_k = g^{(k)}(m)$, then a reversion produces an infinite number of sequences of reversely-written Collatz members [10].

Let S be a set containing two elements q and r , which are bijective functions over \mathbb{Q} :

$$\begin{aligned} q(x) &= 2x \\ r(x) &= \frac{1}{3}(x-1) \end{aligned} \tag{2.1}$$

Let a binary operation be the right-to-left composition of functions $q \circ r$, where $q \circ r(x) = q(r(x))$. Composing functions is an associative operation. All compositions of the bijections q and r and their inverses q^{-1} and r^{-1} are again bijective. The set, whose elements are all these compositions, is closed under that operation. It forms a free group F of rank 2 with respect to the free generating set S , where the group's binary operation \circ is the function composition and the group's identity element is the identity function $id_{\mathbb{Q}} = e$. We call e an *empty string*. F consists of all expressions (strings) that can be concatenated from the generators q and r . The corresponding Cayley graph $Cay(F, S) = G$ is a regular tree whose vertices have four neighbors [33, p. 66]. A tree is called *regular* or *homogeneous* when every vertex has the same degree, in this case, $d(v) = 4$ for every vertex v in G . The Cayley graph's set of vertices is $V(G) = F$, and its set of edges is $E(G) = \{ \{f, f \circ s\} \mid f \in F, s \in (S \cup S^{-1}) \setminus \{e\} \}$ [33, p. 57]. More precisely, the vertices are *labeled* by the elements (strings) of F .

In conformance with graph-theoretical precepts [34], [35], [36] we specify a subgraph H of G as a triple $(V(H), E(H), \psi_H)$ consisting of a set $V(H)$ of vertices, a set $E(H)$ of edges, and an incidence function ψ_H . The latter is, in our case, the restriction $\psi_G|_{E(H)}$ of the Cayley graph's incidence function to the set of edges that only join vertices, which are labeled by a string over alphabet $\{r, q\}$ without the inverses: $E(H) = \{ \{f, f \circ s\} \mid f \in F, s \in S \setminus \{e\} \}$.

This subgraph corresponds to the monoid S^* , which is freely generated by S follows related thoughts [29] that examine the Collatz problem in terms of a free semigroup on the set S^{-1} of inverse generators. Note that this semigroup is not to be confused with an *inverse semigroup* "in which every element has a unique inverse" [37, p. 26], [33, p. 22].

Let $Y^X = \{f \mid f \text{ is a map } X \rightarrow Y\}$ be the set of functions, which in category theory is referred to as the *exponential object* for any sets X, Y . The evaluation function $ev : Y^X \times X \rightarrow Y$ sends the pair (f, x) to $f(x)$. For a detailed description of this concept, see [38, p. 127], [39, p. 155], [40, p. 54] and [41, p. 188]. We define the evaluation function $ev_{S^*} : S^* \times \{1\} \rightarrow \mathbb{Q}$ that evaluates an element of S^* , id est a composition of q and r , for the given input value 1. Furthermore we define the corestriction $ev_{S^*}^0$ of ev_{S^*} to \mathbb{N} . Since a corestriction of a function restricts the function's codomain [42, p. 3], the function $ev_{S^*}^0$ operates on a subset $T \subset S^*$ that

Definition 2.1 The graph H_U possess the following key properties:

- **H_U is a directed graph (digraph):** Fundamentally, when we consider the more general case, an undirected graph as a triple (V, E, ψ) , the incidence function maps an edge to an arbitrary vertex pair $\psi : E \rightarrow \{X \subseteq V : |X| = 2\}$. In a digraph, the set $V \times V$ represents ordered vertex pairs. Accordingly the incidence function is more specifically defined, namely as a mapping of the edges to that set $\psi : E \rightarrow \{(v, w) \in V \times V : v \neq w\}$, see [45, p. 15].
- **H_U is a rooted tree:** According to Rosen [44, p. 747], a rooted tree is "a tree in which one vertex has been designated as the root and every edge is directed away from the root." Peculiarly, this definition considers the directionality as an inherent part of rooted trees. Unlike Mehlhorn and Sanders [46, p. 52], for example, who distinguish between an undirected and directed rooted tree.

Note: As long as we do not stipulate that vertices may collapse, it is absolutely guaranteed that the graph is a tree.

- **H_U is an out-tree:** There is exactly one path from the root to every other node [46, p. 52], which means that edge directions go from parents to children [47, p. 108]. This property is implied in Rosen's definition for a rooted tree as well by saying "every edge is directed away from the root." An out-tree is sometimes designated as *out-arborescence* [47, p. 108].
- **H_U is a labeled tree:** For defining a labeled graph, Ehrig et al. [48, p. 23] use a label alphabet consisting of a vertex label set and an edge label set. Since we only label the vertices, in our case the specification of a vertex label set L_V together with the vertex label function $l_V : V \rightarrow L_V$ is sufficient. Originally, we said vertex labels are strings over the alphabet $S = \{q, r\}$, through which the free monoid S^* is generated. We illustrate labeling H_U by defining $l_{V(H_U)}(v) = ev_{S^*}^0(l_{V(G)}(\iota(v)), 1)$, whereby $\iota : V(H_U) \hookrightarrow V(G)$ is the inclusion map [49, p. 142] from the set of vertices of H_U to the set of vertices from the previously defined Cayley graph G .

We define a tree $H_{C,3}$ by taking the tree H_U as a basis and for every vertex $v \in V(H_U)$ satisfying $2 \mid l_{V(H_U)}(v)$, we contract the incoming edge. We attach the label of the parent of v to the new vertex, which results by replacing (merging) the two overlapping vertices that the contracted edge used to connect. Visually, we obtain $H_{C,3}$ by contracting all edges in H_U that have an even-labeled target vertex, which (due to contraction) gets "merged into its parent." Edge contraction is occasionally referred to as *collapsing an edge*. For more details and examples on edge contraction, one can see Voloshin [50, p. 27] and Loehr [51].

The tree $H_{C,3}$ is well known as the so-called *Syracuse tree*, see for example [52], [53], and [54]. It is a *minor* of H_U , since it can be obtained from H_U "by a sequence of any vertex deletions, edge deletions and edge contractions" [50, p. 32]. The sequence of contracting the edges between adjacent (in our case even-labeled) vertices is called *path contraction*.

A small section of the tree $H_{C,3}$, the Syracuse tree, is shown in figure 2.2. Other definitions of the same tree exist, see for example Conrow [55], Bauer [56, p. 379], Batang [57] or Jan Kleinnijenhuis and Alissa M. Kleinnijenhuis [52], [58].



Figure 2.2: Section of the Syracuse tree $H_{C,3}$ (displaying the trivial cycle is waived)

2.3 Relationship of successive nodes in $H_{C,3}$

Let v_1 and v_{n+1} be two vertices of $H_{C,3}$, where v_1 is reachable from v_{n+1} with $\text{depth}(v_1) - \text{depth}(v_{n+1}) = n$. Hence, a path (v_{n+1}, \dots, v_1) exists between these two vertices. Theorem 2.1 specifies the following relationship between v_1 and v_{n+1} , empirically identified by Koch [1].

Theorem 2.1 $l_{V(H_{C,3})}(v_{n+1}) = 3^n l_{V(H_{C,3})}(v_1) \prod_{i=1}^n \left(1 + \frac{1}{3l_{V(H_{C,3})}(v_i)}\right) 2^{-\alpha_i}$. In order to simplify readability, we waive writing down the vertex label function and put it shortly: $v_{n+1} = 3^n v_1 \prod_{i=1}^n \left(1 + \frac{1}{3v_i}\right) 2^{-\alpha_i}$. The value $\alpha_i \in \mathbb{N}$ is the number of edges which have been contracted between v_i and v_{i+1} in H_U .

In order to demonstrate the construction produced by theorem 2.1 in an illustrative fashion, example 2.1 runs through a concrete path in $H_{C,3}$.

Example 2.1 For example, the two vertices $v_1 = 45$ and $v_{1+3} = v_4 = 5$ are connected via the path $(5, 13, 17, 45)$, see figure 2.2. Furthermore, one can retrace in figure 2.3 the uncontracted path between these two nodes within H_U . When applied to this example, theorem 2.1 produces the following:

$$5 = v_{1+3} = 3^3 \cdot 45 \cdot \left(1 + \frac{1}{3 \cdot 45}\right) \cdot 2^{-3} \cdot \left(1 + \frac{1}{3 \cdot 17}\right) \cdot 2^{-2} \cdot \left(1 + \frac{1}{3 \cdot 13}\right) \cdot 2^{-3}$$

Proof. This relationship of successive nodes can simply be proven inductively. For the base case, we set $n = 1$ and retrieve

$$v_{1+1} = 3v_1 \left(1 + \frac{1}{3v_1}\right) 2^{-\alpha_1} = (3v_1 + 1) 2^{-\alpha_1} = v_2$$

The path from v_2 to v_1 can conformly be expressed by a string $rq \cdots q$ of S^* , because of $v_1 = r \circ q^{\alpha_1}(v_2)$. We set $n = n + 1$ for the step case, which leads to

$$\begin{aligned}
 v_{n+2} &= 3^{n+1} v_1 \prod_{i=1}^{n+1} \left(1 + \frac{1}{3v_i}\right) 2^{-\alpha_i} \\
 &= 3^{n+1} v_1 \left(1 + \frac{1}{3v_{n+1}}\right) 2^{-\alpha_{n+1}} \prod_{i=1}^n \left(1 + \frac{1}{3v_i}\right) 2^{-\alpha_i} \\
 &= 3 \left(1 + \frac{1}{3v_{n+1}}\right) 2^{-\alpha_{n+1}} 3^n v_1 \prod_{i=1}^n \left(1 + \frac{1}{3v_i}\right) 2^{-\alpha_i} \\
 &= 3 \left(1 + \frac{1}{3v_{n+1}}\right) 2^{-\alpha_{n+1}} v_{n+1} \\
 &= (3v_{n+1} + 1) 2^{-\alpha_{n+1}}
 \end{aligned}$$

In this case the path from v_{n+2} to v_{n+1} is conformly expressable by a string $rq \cdots q$ of S^* too, since $v_{n+1} = r \circ q^{\alpha_{n+1}}(v_{n+2})$. \square

Even though the tree may theoretically contain two or more identically labeled vertices, it is essential to emphasize that we only consider such paths (v_{n+1}, \dots, v_1) whose vertices are all labeled differently. Later in section 4.1, we even require that identically labeled nodes are one and the same. In order to correctly determine successive nodes using theorem 2.1, we must consider the halting conditions. These are specified in definition 2.2.

Definition 2.2 When determining successive nodes starting at v_1 according to theorem 2.1, we halt if one of the following two conditions is fulfilled:

1. $v_{n+1} = 1$
2. $v_{n+1} \in \{v_1, v_2, \dots, v_n\}$

If the first condition applies, the Collatz conjecture is true for a specific sequence. When the second condition is fulfilled, the sequence has led to a cycle. For every starting node, except the root node (labeled with 1), the Collatz conjecture is consequently falsified. Let us consider the example $v_1 = 13$, where the algorithm halts after two iterations, because the first condition is met:

$$v_{n+1} = 3^2 \cdot \left(1 + \frac{1}{3 \cdot 13}\right) \left(1 + \frac{1}{3 \cdot 5}\right) \cdot 2^{-7} = 1$$

If we examine the case $v_1 = 1$, we realize that the algorithm finishes after the first iteration, since both halting conditions are true. The sequence stops because the final node labeled with 1 is reached. Furthermore, the sequence has led to a cycle:

$$v_{n+1} = 3 \cdot \left(1 + \frac{1}{3}\right) 2^{-2} = 1$$

The trivial cycle is the only sequence where both conditions are fulfilled.

Theorem 2.1 can be used for specifying the condition of a cycle as follows:

$$\begin{aligned}
 v_1 &= 3^n v_1 \prod_{i=1}^n \left(1 + \frac{1}{3v_i}\right) 2^{-\alpha_i} \\
 2^{\alpha_1 + \dots + \alpha_n} &= \prod_{i=1}^n \left(3 + \frac{1}{v_i}\right)
 \end{aligned} \tag{2.2}$$

A similar condition has been formulated by Hercher [59] and Eric Roosendaal [60]. Taking a first look at equation 2.2, we are able to recognize the trivial cycle for $n = 1$. One might easily come to the false conclusion that the term only results in a natural number for this trivial cycle, since we are multiplying fractions. The following counterexample, starting at $v_1 = 31$, disproves this assumption:

$$20480 = \left(3 + \frac{1}{31}\right) \left(3 + \frac{1}{47}\right) \left(3 + \frac{1}{71}\right) \left(3 + \frac{1}{107}\right) \left(3 + \frac{1}{161}\right) \left(3 + \frac{1}{121}\right) \left(3 + \frac{1}{91}\right) \left(3 + \frac{1}{137}\right) \left(3 + \frac{1}{103}\right)$$

According to OESIS [61], the integer $v_1 = 31$ is called *self-contained*. The term self-contained is based on the fact that the node $v_{n+1} = v_{10} = 155$ is divisible by the starting node $v_1 = 31$. Moreover, v_{10} results from applying one and the same function (in this case the Collatz function) using v_1 as input, see also Guy [62, p. 332]. For such a case equation 2.2 leads to a natural number, but not necessarily to a cycle. A cycle only occurs if the term results in a power of two. One example is the trivial cycle. We find another case when we choose the factor 5 instead of 3:

$$128 = 2^7 = \left(5 + \frac{1}{13}\right) \left(5 + \frac{1}{33}\right) \left(5 + \frac{1}{83}\right)$$

The above example shows that non-trivial cycles can be found if we generalize the Collatz conjecture by replacing the factor 3 with the variable k . We study this generalized form and the occurrence of cycles in section 4.1. A detailed elaboration of the divisibility and a deeper understanding of the tree $H_{C,3}$ needs to be performed in order to get towards any proof of the Collatz conjecture.

Generally, for any variant $kx + 1$ it applies that if $v_1 \mid v_{n+1}$, the product $\prod_{i=1}^n (k + 1/v_i)$ is natural.

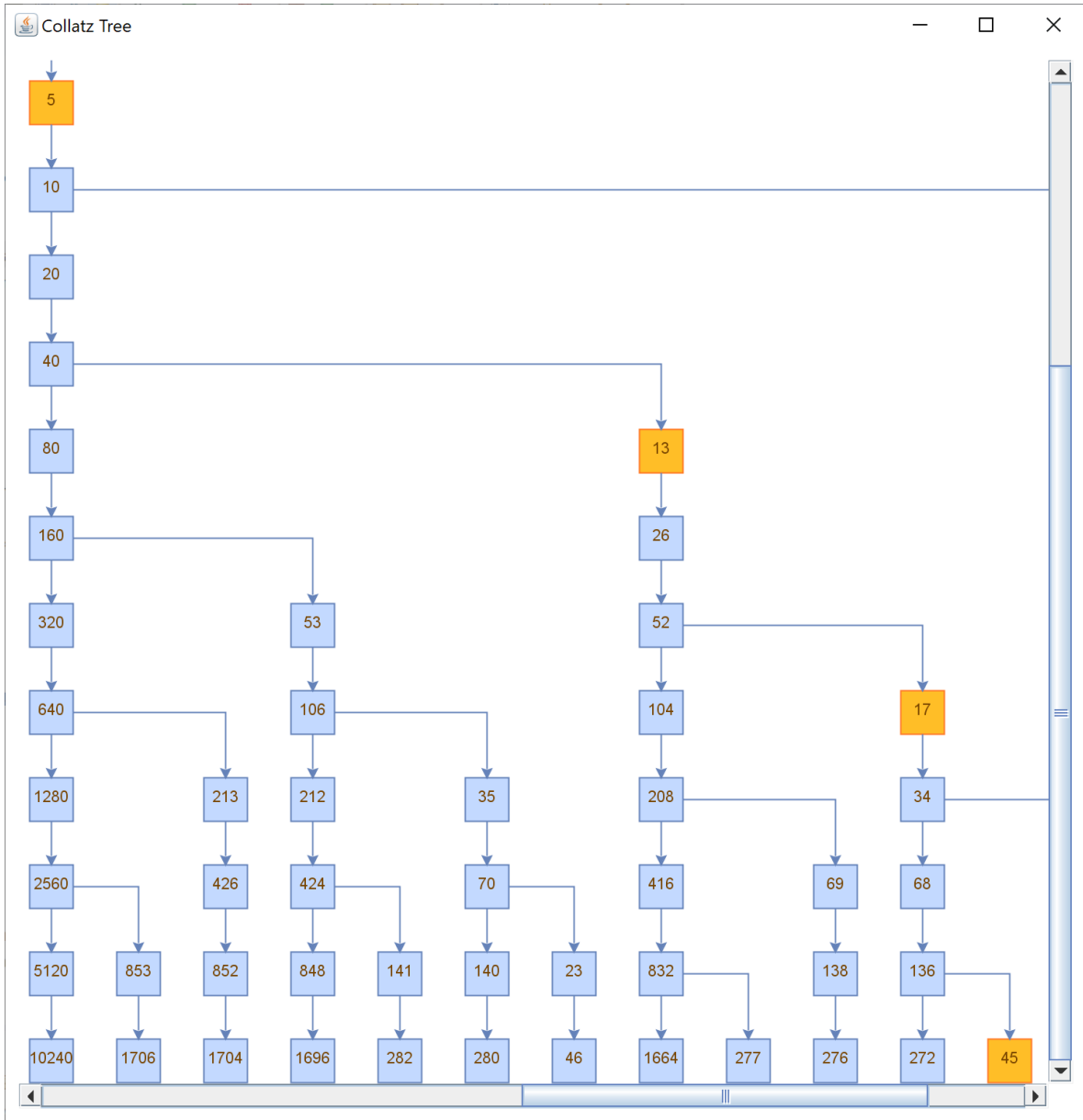


Figure 2.3: Section of H_U containing the path from 5 to 45

2.4 A note on the functions left-child and right-sibling

Referring to the "left-child, right-sibling representation" of rooted trees [63, p. 246], the function $\text{left-child} : V \rightarrow V$ returns the leftmost child of a vertex v . Nesting this function n times leads to the definition of a vertex's n -fold left-child, which is given by $\text{left-child}^n(v)$. As shown in figure 2.2, for example $\text{left-child}^3(13) = 7$.

The function $\text{right-sibling} : V \rightarrow V$ points to the sibling of a vertex v immediately to its right [63, p. 246]. If this function is nested n times, we get a vertex's n -fold right-sibling de-

finied by $\text{right-sibling}^n(v)$. One example is $\text{right-sibling}^2(113) = 1813$ which has been demonstrated in figure 2.2 too.

2.5 Relationship of sibling nodes in $H_{C,3}$

In a rooted tree, vertices which have the same parent are called "siblings" [38, p. 702], [44, p. 747]. Sibling vertices accordingly have the same depth and thus the same level.

Let w be a vertex, from which a path exists to the vertex v_1 . Let v_2 be the immediate right-sibling of v_1 , then $l_{V(H_{C,3})}(v_2) = 4 \cdot l_{V(H_{C,3})}(v_1) + 1$. This fact has been expressed differently by Kak [20] as follows: "If an odd number a leads to another odd number (after several applications of the Collatz transformation) b , then $4a + 1$ also leads to b ."

Applied to our approach, consider w as the parent of v_1 and v_2 . Suppose, in H_U , a path consisting of $n + 1$ edges goes from w to v_1 . Then we can straightforwardly show that n edges in H_U have been contracted between both nodes w and v_1 and $n + 2$ edges between w and v_2 (for simplicity we again omit writing the label function):

$$\begin{aligned} v_1 &= \frac{w \cdot 2^n - 1}{3} \\ v_2 &= \frac{w \cdot 2^{n+2} - 1}{3} = 4 \cdot v_1 + 1 \end{aligned} \quad (2.3)$$

For example, $n = 3$ edges in H_U have been contracted between $w = 5$ and $v_1 = 13$ and $n + 2 = 5$ edges between w and $v_2 = 53$, whereby in $H_{C,3}$, the vertex v_2 is the right-sibling of v_1 and these two sibling vertices are immediate children of w .

Batang [57] demonstrated that using the geometric series $s_n = 1 + 4 + 4^2 + \dots + 4^{n-1} = 4^n - 1/3$ we are able to calculate the sibling nodes directly (see [64, p. 191-192] for more details on geometric series). Let us consider the sibling nodes 5, 21, 85, 341. The first sibling of 5 is calculated by $s_1 + 4^1 \cdot 5 = 21$, the second sibling is $s_2 + 4^2 \cdot 5 = 85$, and the third is $s_3 + 4^3 \cdot 5 = 341$.

The same principle applies to the siblings 13, 53, 213, 853. The first sibling of 13 is calculated by $s_1 + 4^1 \cdot 13 = 53$, the next one is $s_2 + 4^2 \cdot 13 = 213$, and the third is $s_3 + 4^3 \cdot 13 = 853$.

2.6 A vertex's left-child, n -fold right-sibling in $H_{C,3}$

Let w be a vertex in $H_{C,3}$ and v_0 the left-child of w . Using Koch's data science approach [1], we empirically we have found out, that the n -fold right-sibling of v_0 can be calculated as follows:

$$v_n = \text{right-sibling}^n(v_0) = \frac{1}{3} \left(w \cdot 2^{2n+\pi_3(w \bmod 3)} - 1 \right) \quad (2.4)$$

Hereby the function π_3 , which appears in the exponent, is the self-inverse permutation (involution):

$$\pi_3 = \begin{pmatrix} 1 & 2 \\ 2 & 1 \end{pmatrix} \quad (2.5)$$

We consider permutations of the set $\{1, 2\}$ and not of $\{0, 1, 2\}$, due to the fact that $w \bmod 3$ cannot be zero. A node w in $H_{C,3}$, which is labeled by an integer divisible by 3 is a leaf; and therefore such node has no left-child, more specifically it has no children at all. When setting $n = 0$, we trivially retrieve the vertex's w left-child:

$$v_0 = \text{left-child}(w) = \frac{1}{3} \left(w \cdot 2^{\pi_3(w \bmod 3)} - 1 \right)$$

Example 2.2 Let us refer to figure 2.2 again and pick out $w = 5$. Then the vertex's w left-child is $v_0 = 3$ and the threefold right-sibling $v_3 = 213$:

$$\begin{aligned} v_0 &= \frac{1}{3} \left(5 \cdot 2^{\pi_3(5 \bmod 3)} - 1 \right) = 3 \\ v_3 &= \frac{1}{3} \left(5 \cdot 2^{2 \cdot 3 + \pi_3(5 \bmod 3)} - 1 \right) = 213 \end{aligned}$$

We will now explain the reasons that are underlying this behavior. Basically two integers a and b are congruent modulo n if their difference $a - b$ is divisible by n or, to put it differently, if a and b have the same remainder when divided by n [65, p. 15], [66, p. 44], [67, p. 19]:

$$n \mid (a - b) \rightarrow a \equiv b \pmod{n} \quad (2.6)$$

In modular arithmetic we are allowed to interpret integers as names, or to be more specific as *representatives*, for their equivalence class and therefore reduce (or expand) a congruence as follows:

$$(a + n) \equiv b \pmod{n} \rightarrow a \equiv b \pmod{n} \quad (2.7)$$

This means, that in modular arithmetic both operations – addition and multiplication – are independent from the choice of representatives in the residue classes [65, p. 16].

The residue class (also termed congruence class) of the integers for a modulus n is the set $[a]_n = \{a + k \cdot n \mid k \in \mathbb{Z}\}$ and sometimes denoted by \bar{a}_n or by $a + n\mathbb{Z}$, see [65, p. 15], [68, p. 120], [67, p. 25]. Let us put all possible remainders that arise from the division modulo n together into a new set – the set of all residue classes $[a]_n$. This set is known as the ring of integers modulo n and denoted by $\mathbb{Z}/n\mathbb{Z} = \{[a]_n \mid a \in \mathbb{Z}\}$ and trivially $\mathbb{Z}/0\mathbb{Z} = \mathbb{Z}$ and for all $n \neq 0$ we have $\mathbb{Z}/n\mathbb{Z} = \{[0], [1], \dots, [n-1]\}$, see [65, p. 15], [67, p. 25], [64, p. 81].

Now there is one more tool that we will make use of later, and that is the *Congruence Power Rule (CPR)*. It states that we are allowed to raise both sides of an congruence to the n -th power [67, p. 19], [69, p. 117]:

$$a \equiv b \pmod{n} \rightarrow a^n \equiv b^n \pmod{n} \quad (2.8)$$

Let G be a group and $a \in G$. If there exist an integer $d > 0$ with $a^d = e$, then the smallest such d is called the *order* of a , written $d = \text{ord}(a)$ [65, p. 35], [66, p. 50], [70, p. 240]. If no such d exists, we formally write $\text{ord}(a) = \infty$. The number of elements of G is called the *order* of G , written $\text{ord}(G)$ [65, p. 26], [66, p. 50].

Let G be a group and $a \in G$ an element of G . We consider the set of all elements $a^n \in G$ with $n \in \mathbb{Z}$. Since $a^n a^m = a^{n+m} = a^m a^n$ and $(a^n)^{-1} = a^{-n}$, this set forms an abelian subgroup H_a of G . This subgroup H_a is also written $\langle a \rangle$ and called the subgroup of G *generated* by a [66, p. 50]. A group G is called *cyclic*, if an $a \in G$ exists so that G consists only of powers of a (with exponents in \mathbb{Z}), thus if $G = \langle a \rangle$ [65, p. 34], [66, p. 50], [70, p. 240]. In this case $\text{ord}(a) = \text{ord}(G)$, id est the order of an element $a \in G$ is equal to the order of the cyclic subgroup $\langle a \rangle$ [66, p. 50].

Let us consider the cyclic group $\langle a \rangle = \{e, a, a^2, \dots, a^{n-1}\}$ and an element $b \in \langle a \rangle$. Now let us face the question "How do we find the unique integer $0 \leq j \leq n-1$, such that $a^j = b$?" This

integer j is denoted by $j = \log_a b$ and it is called the *discrete logarithm* of b [71, p. 255-256]. To make it more clear that we are talking about the discrete logarithm we write $j = \text{dlog}_a b$ or more specifically $j = \text{dlog}_{a,k} b$ if we solve the congruence $a^j \equiv b \pmod{k}$ which means we solve the equation $a^j \bmod k = b$, see [72].

The multiplicative group of integers modulo n , denoted as $(\mathbb{Z}/n\mathbb{Z})^\times$ or briefly as \mathbb{Z}_n^* contains those elements from $\mathbb{Z}/n\mathbb{Z}$ whose representatives are coprime to n [64, p. 87]:

$$\mathbb{Z}_n^* = \{a \in \mathbb{Z}/n\mathbb{Z} \mid \gcd(a, n) = 1\} \quad (2.9)$$

This group \mathbb{Z}_n^* is a finite abelian group, which contains only the elements from the ring $\mathbb{Z}/n\mathbb{Z}$ that are invertible with respect to multiplication – the units of $\mathbb{Z}/n\mathbb{Z}$. Within the ring $\mathbb{Z}/n\mathbb{Z}$ an element a is invertible exactly if there exists an element b such that $a * b \equiv 1 \pmod{n}$. An element inverse to a exists exactly if $\gcd(a, n) = 1$.

The *multiplicative order* $\text{ord}(a)$ of an element $a \in \mathbb{Z}_n^*$ is the smallest natural exponent d which satisfies $a^d = 1$. In other words, for a positive integer n we say that an integer a has multiplicative order d modulo n if $a^d \equiv 1 \pmod{n}$ where again d is the smallest possible exponent [73, p. 76], [74, p. 32]. To indicate that the order of a refers to the modulus n , it is also often written $d = \text{ord}_n(a)$. Recall that $\gcd(a, n) = 1$, since $a \in \mathbb{Z}_n^*$.

But what about the order of an multiplicative group of integers modulo n ? It is exactly given by *Euler's totient function*, see [65, p. 27]:

$$\text{ord}(\mathbb{Z}_n^*) = \phi(n) \quad (2.10)$$

Euler's totient function $\phi(n)$ counts the positive integers up to a given integer n that are coprime to n [66, p. 49]. The fact that equation 2.10 is correct follows directly from the definition of $\phi(n)$ – we include into \mathbb{Z}_n^* exactly those integers (representatives) from $\mathbb{Z}/n\mathbb{Z}$ that are coprime to n .



The elements of the ring of integers modulo n do not form a group with respect to multiplication, because the element 0 can not be inverted. But also $\mathbb{Z}/n\mathbb{Z} \setminus \{0\}$ does not form a group for a composite n , since there are always products of elements $a \neq 0, b \neq 0$ with $a * b = 0$, which means that the "closure" property is not given [75].

An important theorem related to Euler's totient function, which we will use at a later stage, is Euler's theorem. Euler's theorem states that for an integer $a \geq 2$ coprime to the modulus n the following congruence holds [67, p. 37], [66, p. 56], [64, p. 104]:

$$a^{\phi(n)} \equiv 1 \pmod{n} \quad (2.11)$$

This means that given $\langle a \rangle = \mathbb{Z}_n^*$, for any generator a coprime to the modulus n the congruence $a^{\phi(n)} \equiv 1 \pmod{n}$ becomes true, where again $\phi(n)$ is the order of \mathbb{Z}_n^* and thus the order of a (see 2.10). If $\phi(n) \equiv 2 \pmod{4}$ then the group \mathbb{Z}_n^* is cyclic. Consequently the multiplicative group of integers modulo n is cyclic for $n \in \{1, 2, 4, p^j, 2p^j\}$, where p being an odd prime and $j \geq 1$ [75].

At this point it is appropriate that we explain the mapping (permutation) given by 2.5 in more detail. A helpful tool that we can use as a point of departure is the multiplicative group of integers modulo 3. The element 2 is a generator $\langle 2 \rangle = \{1, 2\} = \mathbb{Z}_3^*$, since $2 \equiv 2 \pmod{3}$ and $2^2 \equiv 4 \equiv 1 \pmod{3}$. The order of 2 is 2, since $2^2 \equiv 1 \pmod{3}$. Now we use the CPR given by 2.8 and obtain $(2^2)^{n+1} \equiv 1^{n+1} \pmod{3}$ and the following generic congruence:

$$2^j 2^{2n+2-j} \equiv 1 \pmod{3} \quad (2.12)$$

Setting $j = 0, 1$ leads to the following behavior, which explains the formula 2.4 and the mapping 2.5:

j	congruence 2.12	node w	setting w as per 2.7	divisibility as per 2.6
0	$1 \cdot 2^{2n+2} \equiv 1$	$w \in [1]_3$	$w \cdot 2^{2n+\pi_3(w \bmod 3)} \equiv 1$	$3 (w \cdot 2^{2n+\pi_3(w \bmod 3)} - 1)$
1	$2 \cdot 2^{2n+1} \equiv 1$	$w \in [2]_3$		

If addition is the group operation, as it is the case for example with the additive group of integers modulo 3, denoted as $(\mathbb{Z}/3\mathbb{Z}, +)$ or as $(\mathbb{Z}_3, +)$, then for an element a the term a^n means add (and not multiply) a to itself n times. In this specific case the group contains three elements $\mathbb{Z}_3 = \{0, 1, 2\}$ and the identity element is $e = 0$. The element 2 is a generator $\langle 2 \rangle = \{0, 1, 2\} = \mathbb{Z}_3$, since $2 \equiv 2 \pmod{3}$ and $2 + 2 \equiv 4 \equiv 1 \pmod{3}$ and $2 + 2 + 2 \equiv 6 \equiv 0 \pmod{3}$. The order of 2 is 3, because $2^3 \equiv 0 \pmod{3}$.



2.7 A vertex's left-child, n -fold right-sibling in $H_{C,5}$

In the following we take a look at the tree $H_{C,5}$ – the $5x+1$ variant of H_C , elaborated by Koch [1] using Data Science too. We must note that it is not a tree and moreover that not all of its vertices are reachable from the root. We define the permutation π_5 as follows:

$$\pi_5 = \begin{pmatrix} 1 & 2 & 3 & 4 \\ 4 & 3 & 1 & 2 \end{pmatrix}$$

Next, by letting w be a vertex in $H_{C,5}$ and v_0 the left-child of w we obtain the n -fold right-sibling of v_0 by the function that is slightly different to the one defined by 2.4:

$$v_n = \text{right-sibling}^n(v_0) = \frac{1}{5} (w \cdot 2^{4n+\pi_5(w \bmod 5)} - 1) \quad (2.13)$$

Analogous to 2.5 only permutations on the set without zero $\{1, 2, 3, 4\}$ need to be considered, since $w \bmod 5$ cannot be zero. Otherwise, if $w \equiv 0 \pmod{5}$ which means that w were labeled by an integer divisible by 5, then the node w has no successor in $H_{C,5}$. By setting $n = 0$, the function (above given by 2.13) returns the left child of w :

$$v_0 = \text{left-child}(w) = \frac{1}{5} (w \cdot 2^{\pi_5(w \bmod 5)} - 1)$$

Equation 2.13 has been identified empirically as well and can be explained using the cyclic group $\langle 2 \rangle = \{1, 2, 3, 4\} = \mathbb{Z}_5^*$. First of all, it is obvious that 2 generates this group, since $2 \equiv 2 \pmod{5}$ and $2^2 \equiv 4 \pmod{5}$ and $2^3 \equiv 8 \equiv 3 \pmod{5}$ and $2^4 \equiv 16 \equiv 1 \pmod{5}$. The order is 4 and according to 2.11 and 2.10 we have $2^{\text{ord}(\mathbb{Z}_5^*)} \equiv 2^{\phi(5)} \equiv 2^4 \equiv 1 \pmod{5}$. Again we use the CPR given by 2.8 and obtain $(2^4)^{n+1} \equiv 1^{n+1} \pmod{5}$ and the following generic congruence:

$$2^j 2^{4n+4-j} \equiv 1 \pmod{5} \quad (2.14)$$

Setting $j = 0, 1, 2, 3$ leads to the following behavior, which explains the formula 2.13 and the mapping 2.7:

j	congruence 2.14	node w	setting w as per 2.7	divisibility as per 2.6
0	$1 \cdot 2^{4n+4} \equiv 1$	$w \in [1]_5$	$w \cdot 2^{4n+\pi_5(w \bmod 5)} \equiv 1$	$5 (w \cdot 2^{4n+\pi_5(w \bmod 5)} - 1)$
1	$2 \cdot 2^{4n+3} \equiv 1$	$w \in [2]_5$		
2	$4 \cdot 2^{4n+2} \equiv 1$	$w \in [4]_5$		
3	$8 \cdot 2^{4n+1} \equiv 1$	$w \in [3]_5$		

Figure 2.4 illustrates a small section of $H_{C,5}$ starting at its root. The particularly interesting thing about the graph $H_{C,5}$ is that it contains three cycles, the trivial cycle starting from the root 1, 3 and two non-trivial cycles 43, 17, 27 and 83, 33, 13. To be precise, three cycles are known (as it will become apparent later in section 4.2), and on the basis of present knowledge it cannot be ruled out with any certainty that other cycles exist.

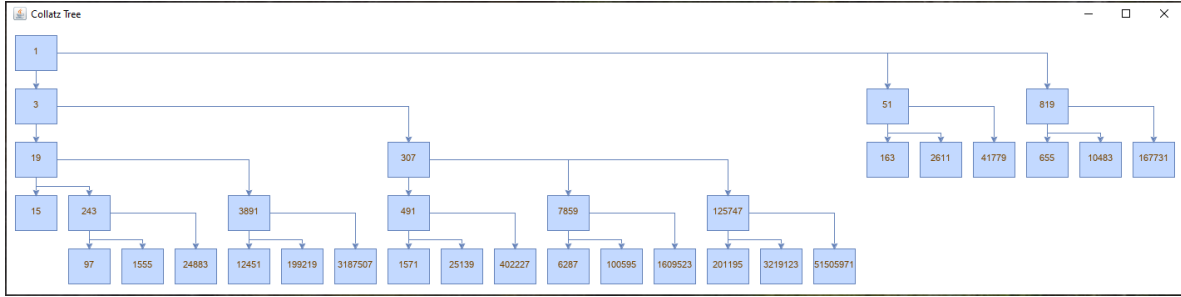


Figure 2.4: Section of the graph $H_{C,5}$ starting at its root (without branches that reflect a subsequence containing the trivial cycle)

2.8 A vertex's left-child, n -fold right-sibling in $H_{C,7}$

Now we are able to develop the formula on our own that calculates for a given node w the left-child and right-sibling in $H_{C,7}$. We refer to the cyclic group $\mathbb{Z}_7^* = \{1, 2, 3, 4, 5, 6\}$. Note that in this case 2 is not a generator of this group. But nevertheless \mathbb{Z}_7^* is cyclic and $\text{ord}(2) = 3$ which gives $2^3 \equiv 1 \pmod{7}$. Again we use the CPR given by 2.8 and obtain $(2^3)^{n+1} \equiv 1^{n+1} \pmod{7}$ and the following generic congruence:

$$2^j 2^{3n+3-j} \equiv 1 \pmod{7} \quad (2.15)$$

Setting $j = 0, 1, 2$ leads to the following behavior, which produces the formula 2.16 and the mapping 2.17:

j	congruence 2.15	node w	setting w as per 2.7	divisibility as per 2.6
0	$1 \cdot 2^{3n+3} \equiv 1$	$w \in [1]_7$	$w \cdot 2^{3n+\pi_7(w \bmod 7)} \equiv 1$	$7 (w \cdot 2^{3n+\pi_7(w \bmod 7)} - 1)$
1	$2 \cdot 2^{3n+2} \equiv 1$	$w \in [2]_7$		
2	$4 \cdot 2^{3n+1} \equiv 1$	$w \in [4]_7$		

$$v_n = \text{right-sibling}^n(v_0) = \frac{1}{7} \left(w \cdot 2^{3n+\pi_7(w \bmod 7)} - 1 \right) \quad (2.16)$$

The mapping 2.17 is not a permutation as in the case of π_3 and π_5 , it is defined as follows:

$$\pi_7(n) = \begin{cases} 3 & n = 1 \\ 2 & n = 2 \\ 1 & n = 4 \end{cases} \quad (2.17)$$

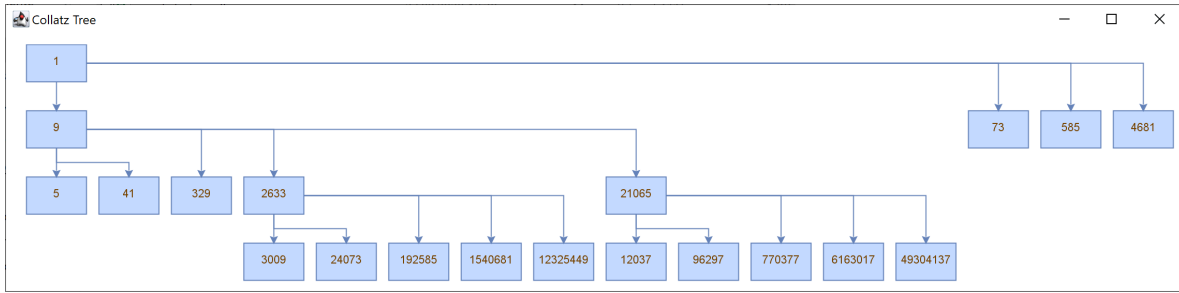


Figure 2.5: Section of the graph $H_{C,7}$ starting at its root (without branches that reflect a subsequence containing the trivial cycle)

2.9 Generalizing the relationship of sibling nodes for $H_{C,k}$

In section 2.5 we have taken a closer look at the relationship of sibling nodes in $H_{C,3}$. But what is the formula for calculating the n -fold right-sibling of a given node v_0 generalized to the $kx + 1$ variant of H_C ? We remember that the multiplicative order $d = \text{ord}_k(2)$ is the smallest natural exponent d such that $2^d \equiv 1 \pmod{k}$. In the case $k = 3$ we can calculate the next sibling v_1 of a given node v_0 as follows: $v_1 = v_0 \cdot 4 + 1$, see 2.3. Within $H_{C,7}$, the next sibling of v_0 is given by $v_1 = v_0 \cdot 8 + 1$. In the case of $k = 5$, we calculate the next sibling $v_1 = v_0 \cdot 16 + 3$ and for $k = 9$ we obtain the next sibling v_1 of a given node v_0 by $v_1 = v_0 \cdot 64 + 7$. The general formula for calculating a given node's v_0 immediate right sibling is:

$$v_1 = v_0 \cdot 2^{\text{ord}_k(2)} + \frac{1}{k} \left(2^{\text{ord}_k(2)} - 1 \right) \quad (2.18)$$

In order to calculate the n -fold right sibling of a given node v_0 , we need to repeatedly nest n times the (linear) function 2.18. For the sake of simplicity, let us substitute $a = 2^{\text{ord}_k(2)}$ and $b = \frac{1}{k}(2^{\text{ord}_k(2)} - 1)$. Then the n -fold right sibling of v_0 is obtained by the following structure:

$$\begin{aligned}
v_n = \text{right-sibling}^n(v_0) &= (((v_0 \cdot a + b) \cdot a + b) \cdot a + b) \cdots \\
&= v_0 \cdot a^n + b(a^{n-1} + \dots + a^2 + a + 1) = v_0 \cdot a^n + b \frac{a^n - 1}{a - 1}
\end{aligned}$$

Note that the term $b(a^{n-1} + \dots + a^2 + a + 1)$ can be simplified using the n -th partial sum of a geometric series ([64, p. 192]). The resubstitution of both coefficients a and b leads us to the final generalized formula that calculates the n -fold right sibling of a node v_0 in $H_{C,k}$:

$$v_n = \text{right-sibling}^n(v_0) = v_0 \cdot 2^{n \cdot \text{ord}_k(2)} + \frac{1}{k} \left(2^{\text{ord}_k(2)} - 1 \right) \cdot \frac{2^{n \cdot \text{ord}_k(2)} - 1}{2^{\text{ord}_k(2)} - 1} \quad (2.19)$$

This can be verified by inserting $n = 0$ and $n = 1$ into formula 2.20 that calculates the a vertex's left-child, n -fold right-sibling of $H_{C,k}$:

$$\begin{aligned}
v_0 &= \frac{1}{k} \left(w \cdot 2^{\text{ord}_k(2) - \text{dlog}_{2,k} w} - 1 \right) & kv_0 + 1 &= w \cdot 2^{\text{ord}_k(2) - \text{dlog}_{2,k} w} \\
v_1 &= \frac{1}{k} \left(w \cdot 2^{2 \cdot \text{ord}_k(2) - \text{dlog}_{2,k} w} - 1 \right) & kv_1 + 1 &= w \cdot 2^{2 \cdot \text{ord}_k(2) - \text{dlog}_{2,k} w}
\end{aligned}$$

This brings us to the following quotient leading to the basic relationship between two sibling nodes that is given by equation 2.18 in the form of $v_1 = v_0 \cdot a + b$:

$$\frac{kv_1 + 1}{kv_0 + 1} = \frac{2^{2 \cdot \text{ord}_k(2) - \text{dlog}_{2,k} w}}{2^{\text{ord}_k(2) - \text{dlog}_{2,k} w}} = 2^{\text{ord}_k(2)}$$

Here we point out that the equation 2.19 of course only works for such k , where the order of two is not infinity $\text{ord}_k(2) \neq \infty$. This means that, for instance, it does not work for $k = 1$, id est for the $1x + 1$ variant of H_C . However, this case is trivial and for the sake of completeness we have added a picture including some few words about $H_{C,1}$ in appendix A.1.

2.10 Generalizing a vertex's left-child, n -fold right-sibling for $H_{C,k}$

In the following, we generalize the formulas, that has been developed in sections 2.6 - 2.8 to calculate the left-child, n -fold right-sibling for a given node w that is the direct parent node of v_0 :

$$v_n = \text{right-sibling}^n(v_0) = \frac{1}{k} \left(w \cdot 2^{\text{ord}_k(2) \cdot n + \text{ord}_k(2) - \text{dlog}_{2,k} w} - 1 \right) \quad (2.20)$$

The discrete logarithm $\text{dlog}_{2,k} w = j$ finds the smallest exponent j such that $2^j \equiv w \pmod{k}$ respectively it solves the equation $2^j \bmod k = w$.

3. Binary Collatz Tree

3.1 Some essentials on binary trees

A binary tree is a rooted tree, where each node has at most two immediate successors. Those nodes, from which no edge goes out downward, are called leaves, the others are called internal nodes. In a full binary tree, all internal nodes have exactly two children [76, p. 102]. Full binary trees have an odd number $2n + 1$ of nodes. Of these $n + 1$ are leaves and n are inner nodes [77, p. 134]. Each node in a binary tree has a left subtree and a right subtree, which is why a binary tree is inherently recursive, since the left and right subtrees of the root are themselves binary trees [78, p. 246-247]. As it often pops up in combinatorial problems, the famous n -th Catalan number, named after the Belgian mathematician Eugène Catalan, comes in connection with binary trees into play. For $n \geq 1$ it specifies the number of binary trees on n vertices [78, p. 247]:

$$B_n = \sum_{i=0}^{n-1} B_i B_{n-1-i} = \sum_{i=1}^n B_{i-1} B_{n-i} = \frac{1}{n+1} \binom{2n}{n}$$

There is an interesting property that trees exhibit regarding abstract algebra. Let's have a look at the algebraic structure of magmas. Consider an element x of a magma $(M, *)$ which is an iterated product of other elements in M . Such an element can be described by a planar (no edges cross each other) rooted binary tree whose n leaves are labelled by these other elements $x_1, \dots, x_n \in M$ [79, p. 96].

Binary trees make well-suited data structures for storing information. With about 2^m data points (nodes), a search of a binary tree takes only about m steps, compared to about 2^{m-1} steps which are required to search a simple list [69, p. 84].

3.2 Transforming the Collatz tree into a binary tree

Jan Kleinnijenhuis and Alissa M. Kleinnijenhuis [52] introduced a binary tree $T_{\geq 0}$ by transforming the original Collatz tree H_U into the Syracuse tree $H_{C,3}$, which in turn is transformed into the binary tree $T_{\geq 0}$ as described next. The edges are changed according to the following procedure: whenever a parent node w has edges to its child nodes v_0, v_1, \dots, v_n , on the tree $H_{C,3}$, we draw an edge from w to v_0 , and edges from v_i to v_{i+1} for each $i = 1, \dots, n-1$, in the binary new tree. Note that the nodes v_1, v_2, \dots, v_n are sorted in increasing order of label $v_0 < v_1 < \dots < v_n$, which is already given by 2.19. Figure 3.1 and 3.2 display that tree – once in our standard layout and once reversed (from bottom to top).

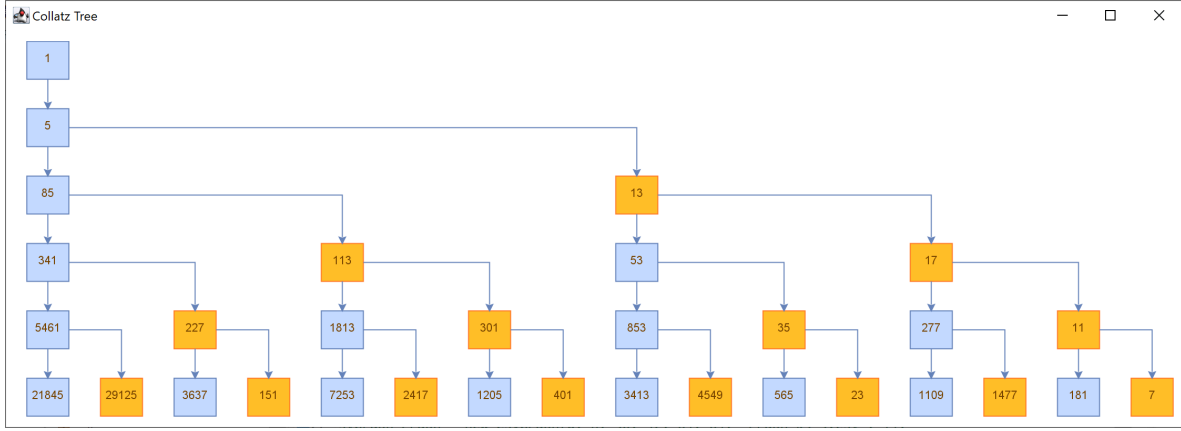


Figure 3.1: The Collatz Tree transformed to the binary tree $T_{\geq 0}$

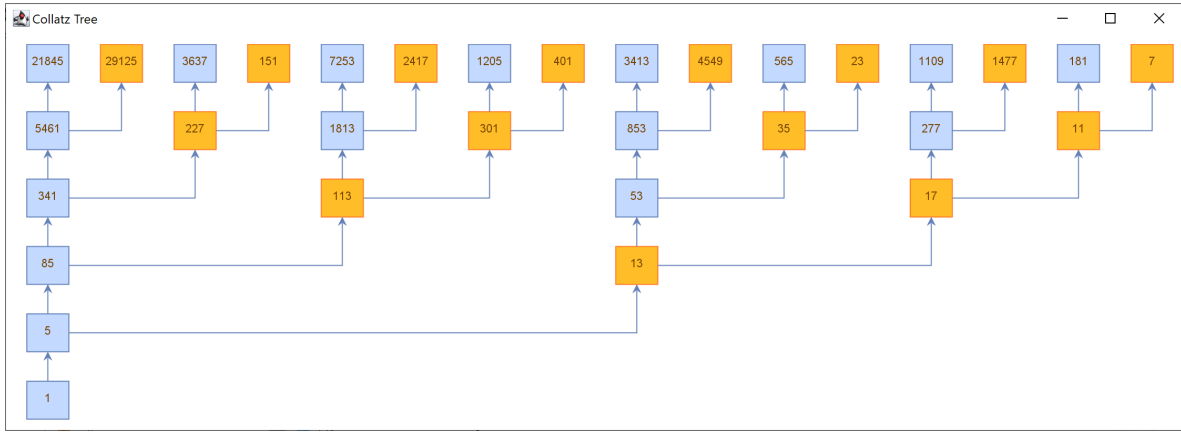


Figure 3.2: The binary tree $T_{\geq 0}$ with *bottom-to-top* layout orientation



To clarify the terminology, it should be mentioned that Jan and Alissa M. Kleinnijenhuis in their manuscripts [52], [58] denote the original Collatz tree T_C while we call it H_U . They denote the Syracuse Tree T_T which in our nomenclature is referred to as $H_{C,3}$.

Nodes that are highlighted orange in figures 3.1, 3.2 are called *prunable* and they are exactly those nodes resulting as output of the *Rightward* function. For navigating within this binary tree, Jan Kleinnijenhuis and Alissa M. Kleinnijenhuis [52] defined an *Upward* function $U(n)$ and a *Rightward* function $R(n)$ as follows:

$$U(n) = \begin{cases} 4n+1 & n \equiv 1 \pmod{6} \\ 16n+5 & n \equiv 5 \pmod{6} \end{cases} \quad R(n) = \begin{cases} (2^2 n - 1)/3 & n \equiv 1 \pmod{18} \\ (2^3 n - 1)/3 & n \equiv 5 \pmod{18} \\ (2^4 n - 1)/3 & n \equiv 7 \pmod{18} \\ (2^1 n - 1)/3 & n \equiv 11 \pmod{18} \\ (2^2 n - 1)/3 & n \equiv 13 \pmod{18} \\ (2^1 n - 1)/3 & n \equiv 17 \pmod{18} \end{cases} \quad (3.1)$$

The domain and codomain of both functions consist of the two residue classes $[1]_6, [5]_6$, which form the multiplicative (cyclic) group $\mathbb{Z}_6^* = \{1, 5\} = \langle 5 \rangle$. Consequently, the domain and codomain exclude all integers divisible by 2 and 3, which is due to the fact that this binary tree (just like our tree $H_{C,3}$) does not contain even numbers and additionally all leaves – namely those nodes labeled with an integer divisible by three – were deleted. The function $U(n)$ is very similar to the function 2.3 and to the more general function 2.19 (when setting $n = 1, k = 3$) which both calculate the right-sibling of a given vertex. This is clear, since siblings (parallel) in $H_{C,3}$ are successors (serial) in the binary tree $T_{\geq 0}$. In the end, for a node v_0 having a leaf as right-sibling in $H_{C,3}$, the function $U(v_0)$ is defined as $v_1 = 4v_0 + 1$ executed twice $v_1 = 4(4v_0 + 1) + 1 = 16v_0 + 5$, because we must skip this leaf. Recall that all leaves in $H_{C,3}$ are excluded from the binary tree without exception. For any $n \in [5]_6$ it applies that $U(n) \equiv 16n + 5 \equiv 1 \pmod{6}$ since $6 \mid 16n + 5 - 1$ resulting in $6 \mid 16(5 + k \cdot 6) + 5 - 1$, see 2.6, and analogously for any $n \in [1]_6$ it applies that $U(n) \equiv 4n + 1 \equiv 5 \pmod{6}$ since $6 \mid 4n + 1 - 5$ resulting in $6 \mid 4(1 + k \cdot 6) + 1 - 5$. Therefore executing the Upward function twice in a row leads unconditionally to $U^2(n) = 16(4n + 1) + 5 = 4(16n + 5) + 1 = 64n + 21$.

While we read (and displayed) trees from top to down, it is sometimes usual to describe trees in a bottom-to-top fashion as for example Kleinnijenhuis [52], [58] do. This means.



Jan and Alissa M. Kleinnijenhuis [52] defined the set $N(T_C) = N(H_U)$ that contains the labels of all nodes, to which a path from the root in H_U exists, in other words, this set contains all integers n for which the orbit of n under the (uncompressed) Collatz function 1.1 converges to 1. Furthermore they introduced $S_{\geq 0}$ as the node set containing integers that are neither divisible by 2 nor by 3. The set S_{-1} comprises on the contrary all numbers, which are divisible by 2 or 3. In order to comprehend the structure of these sets S , let us take a look at the following list showing which tree includes which node set, see also the ancillary files of [52], [58]:

Original Collatz tree	$N(T_C) = N(H_U)$	=	\mathbb{N}^+ if the Collatz conjecture holds
Syracuse tree	$N(T_T) = N(H_{C,3})$	=	$N(T_C) \setminus 2\mathbb{N}$
Binary tree $T_{\geq 0}$	$N(T_{\geq 0}) = S_{\geq 0}$	=	$N(T_C) \setminus S_{-1} = S_0 \cup S_1 \cup S_2 \dots$
Binary tree $T_{\geq 1}$	$N(T_{\geq 1}) = S_{\geq 1}$	=	$N(T_C) \setminus \bigcup_{i=-1}^0 S_i = S_1 \cup S_2 \cup S_3 \dots$
Binary tree $T_{\geq j}$	$N(T_{\geq j}) = S_{\geq j}$	=	$N(T_C) \setminus \bigcup_{i=-1}^{j-1} S_i = \bigcup_{i=j}^{\infty} S_i$

Let us describe these sets using multiplicative groups. The set $S_{\geq 0} = \mathbb{Z}_6^*$ can be understood as the multiplicative group modulo 6 and the set $S_{-1} = \mathbb{Z}/6\mathbb{Z} \setminus \mathbb{Z}_6^* = \{0, 2, 3, 4\}$ as the set of all non-invertible elements (non-units) of $\mathbb{Z}/6\mathbb{Z}$.

The set S_0 consists of all nodes resulting as output of $R(n)$ within the binary tree $T_{\geq 0}$. These are the orange highlighted nodes displayed by figures 3.1, 3.2. In other words, S_0 is the codomain of the function $R(n)$ operating on nodes within $T_{\geq 0}$. The binary tree $T_{\geq 0}$ can be transformed to a (pruned) binary tree $T_{\geq 1}$. For this, the prunable nodes will be deleted and their neighbors reconnected using the algorithm given by the listing 3.1. The upward neighbor of a pruned node will then be identified as pruning candidate for a later transformation of the resulting tree $T_{\geq 1}$ to a more pruned tree $T_{\geq 2}$.

```

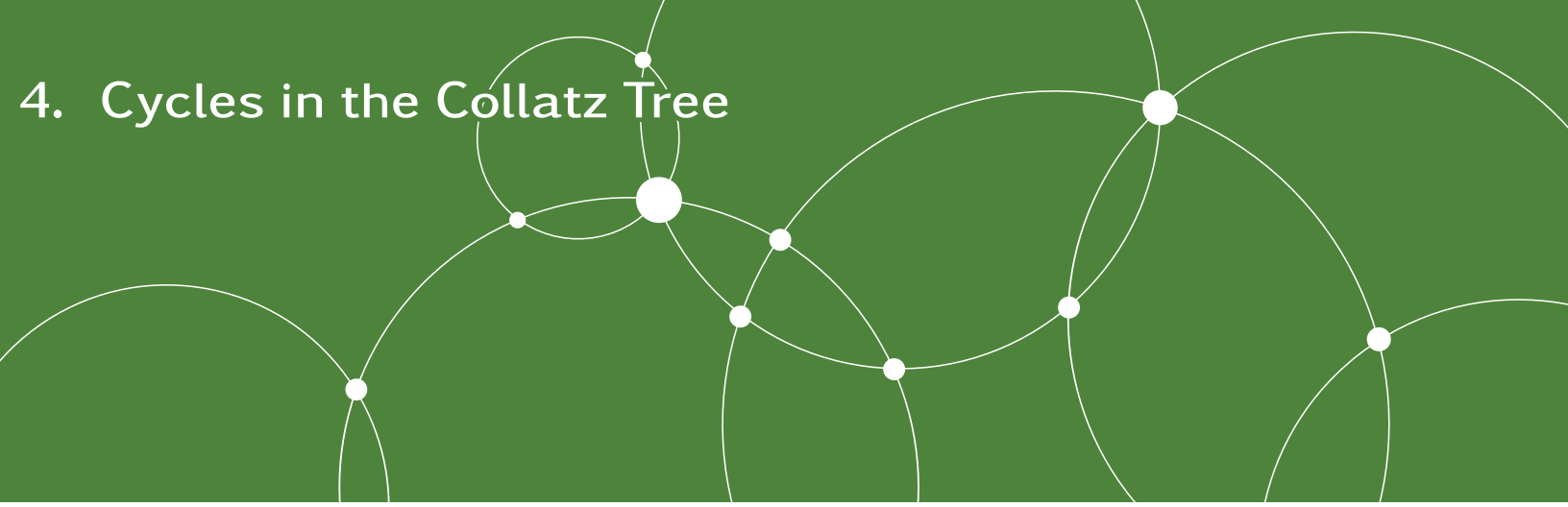
1  def shrinkRoot(self):
2      self.root = self.root.successors[0]
3      self.root.predecessor = None
4
5  def prune(self):
6      self.shrinkRoot()
7      new_prunables = []
8      for node in self.prunable_nodes:
9          if node.predecessor is not None:
10             if len(node.successors) > 1 and len(node.predecessor.successors) > 1:
11                 node.predecessor.successors[1].successors.append(node.successors[1])
12                 node.predecessor.successors[1].successors.reverse()
13                 node.successors[1].predecessor = node.predecessor.successors[1]
14                 node.successors[1].prunable = True
15                 new_prunables.append(node.successors[1])
16                 node.predecessor.successors.remove(node)
17                 node.predecessor = None
18                 self.labels.remove(node.label)
19 self.prunable_nodes = new_prunables
20 return self

```

Listing 3.1: Python function for pruning a binary tree $T_{\geq j}$ [86]

The set S_1 contains all nodes that are (as per description above) identified as pruning candidate for the next transformation of $T_{\geq 1}$ to $T_{\geq 2}$. After having transformed $T_{\geq 1}$ to $T_{\geq 2}$, the more pruned binary tree $T_{\geq 2}$ contains nodes that are identified as pruning candidates for another upcoming transformation of $T_{\geq 2}$ to $T_{\geq 3}$ – these nodes are elements of the set S_2 . This pruning algorithm is repeatedly applied in the same pattern. And in this way we obtain the sets S_1, S_2, S_3, \dots and so forth. Generally, we can write these sets in the form $S_j = \{n \in N(T_{j-1}) \mid U^{-j} \in S_0\}$

The (cyclic) multiplicative group $(\mathbb{Z}/18\mathbb{Z})^\times = \mathbb{Z}_{18}^* = \{1, 5, 7, 11, 13, 17\} = \langle 5 \rangle$ has an order $\text{ord}(\mathbb{Z}_{18}^*) = 6$ and based on Euler's theorem we can derive the following congruences from $5^j 5^{6n+6-j} \equiv 1 \pmod{18}$:



4. Cycles in the Collatz Tree

4.1 A remark about cycles

In graph theory, a path of length $n \geq 1$ that starts and ends at the same vertex is called a circuit. A circuit, in which no vertex is repeated with the sole exception that the initial vertex is the terminal vertex, is called a cycle. A cycle of length n is referred to as an n -cycle. For these definitions, we rely on [44, p. 599], [80, p. 35] and [81, p. 445]. Furthermore, we call a cycle originating from the root a trivial cycle.

In order for the cycles to become graphically visible, we now require that in a graph H two vertices v_1 and v_2 are one and the same if the label of both nodes are identical: $l_{V(H)}(v_1) = l_{V(H)}(v_2) \rightarrow v_1 = v_2$. As a consequence, there is no guarantee that the graph precisely refers to the algebraic structure of a free monoid anymore. A free monoid requires that each of its elements can be written in one and only one way.



When different nodes collapse on one, the graph is no longer necessarily a tree. Let us point to the monoid S^* , which we introduced in section 2.1. Take for example four of its elements, the empty string e , the strings qqr , $qqrqqr$, and $qqrqqrqqr$. These elements lie as well within the subset $U \subset T \subset S^*$, and they are represented by nodes of the tree H_U that all have the same label $1 = ev_{S^*}(qqr, 1) = ev_{S^*}(qqrqqr, 1) = ev_{S^*}(qqrqqrqqr, 1)$. These nodes are one and the same, the root of H_U . Visually, then in H_U a directed edge goes from the vertex labeled with 4 back to the root node. Analogously, in $H_{C,3}$ a loop connects the root to itself, since due to the path contraction even labeled nodes do not exist in $H_{C,3}$. The aforementioned example reflects the trivial cycle of the Collatz sequence.

Figure 4.1 depicts a section of $H_{C,5}$, which includes the 3-cycle 43, 17, 27. Because of the two non-trivial cycles 43, 17, 27 and 83, 33, 13, in $H_{C,5}$ there does not exist a path between the root and the vertex 43 and between the root and the vertex 83. Hence, $H_{C,5}$ is said to be a disconnected graph. Generally, a graph is called a disconnected graph if it is impossible to walk (along its edges) from any vertex to any other [80, pp. 46-47].

The following considerations focus on non-trivial cycles, and therefore on cycles that do not originate from the root, but cause the graph to be a disconnected graph. Utilizing the example of the graph $H_{C,5}$ we are able to deduct from the cycle 43, 17, 27 the simple and self-evident equality $left-child^3(43) = 43$:

$$left-child(43) = \frac{1}{5} * (43 * 2^1 - 1) = 17$$

$$left-child(17) = \frac{1}{5} * (17 * 2^3 - 1) = 27$$

$$left-child(27) = \frac{1}{5} * (27 * 2^3 - 1) = 43$$

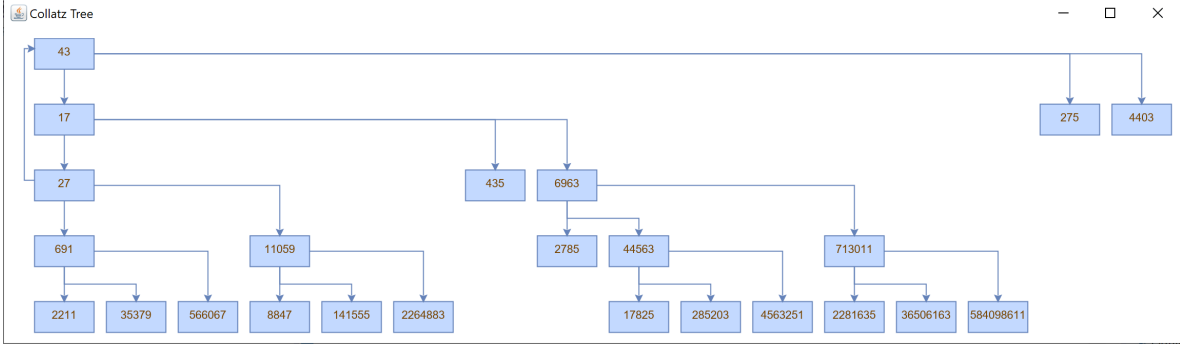


Figure 4.1: Section of $H_{C,5}$ including the 3-cycle 43, 17, 27

Obviously, the authors note, it would be interesting to find out what circumstances enable a graph to have non-trivial cycles, whether it be the $5x + 1$ variant, the $7x + 1$ variant of H_C or any variant $H_{C,k}$ with $k \geq 1$.

4.2 Which variants of H_C have non-trivial cycles?

Let us refer to $H_{C,k}$. By having introduced and proven theorem 2.1 we already started an assertion about the reachability of successive nodes in $H_{C,3}$. This reachability relationship can be generalized for any graph $H_{C,k}$ as follows:

$$v_{n+1} = k^n v_1 \prod_{i=1}^n \left(1 + \frac{1}{kv_i}\right) 2^{-\alpha_i} \quad (4.1)$$

This generalization leads to the condition for an existence of an n -cycle in any $kx + 1$ variant of H_C , which looks analogous to the condition given by equation 2.2 that specifies $H_{C,3}$ has a cycle:

$$2^\alpha = \prod_{i=1}^n \left(k + \frac{1}{v_i}\right) \quad (4.2)$$

The natural number α is the sum of edges that have been contracted between the vertices v_i forming the cycle, in other words α is the number of divisions by 2 within the sequence. The natural number n is the cycle length and k obviously specifies the variant of H_C . Since between each vertex at least one edge has been contracted (at least one division by 2 took place), we know that our exponent alpha is greater than or equal to the sequence length:

$$\alpha \geq n \quad (4.3)$$

Using incremental search, Koch et al. [82] calculated cycles through trial and error. The authors list all empirically discovered cycles having a length up to 100 which appear in $H_{C,k}$ for $k \in [1, 1000]$. Within each of these variants, the cycles have been searched at potential starting nodes v_1 with a label between 1 and 1000. Based on their results they stated the following theorem 4.1 that renders more precisely the prerequisite for cycles that may occur in variants of H_C .

Theorem 4.1 An n -cycle can only exist in a graph $H_{C,k}$, if the following equation holds:

$$2^{\bar{\alpha}} = 2^{\lfloor n \log_2 k \rfloor + 1} = \prod_{i=1}^n \left(k + \frac{1}{v_i} \right)$$

The statement behind theorem 4.1 consists in the claim that, in order for an n -cycle to occur, the exponent α has to be $\bar{\alpha} = \lfloor n \log_2 k \rfloor + 1$. This statement is true if the following general condition for the validity of the cycle-alpha's upper limit always holds (see [82]):

$$n \log_2 k - \lfloor n \log_2 k \rfloor < 2 - \log_2 \left(\prod_{i=1}^n \left(1 + \frac{1}{kv_i} \right) \right) \quad (4.4)$$

A product $\prod(1 + a_n)$ with positive terms a_n is convergent if the series $\sum a_n$ converges, see Knopp [83, p. 220]. A similar statement provides Murphy [84], who write the factors in the form $c_n = 1 + a_n$ and explains that if $\prod c_n$ is convergent then $c_n \rightarrow 1$ and therefore if $\prod(1 + a_n)$ is convergent then $a_n \rightarrow 0$. Thus, to verify whether the product in condition 4.4 is converging towards a limiting value, it is sufficient to examine the following sum:

$$\sum_{i=1}^n \frac{1}{kv_i}$$

The sum of reciprocal vertices depending only from v_1 is given in appendix A.2.

5. Conclusion and Outlook

5.1 Summary

We defined an algebraic graph structure that expresses the Collatz sequences in the form of a tree. Next, the vertex reachability properties were unveiled by examining the relationship between successive nodes in H_C . Moreover, we dealt with graphs that represent other variants of Collatz sequences, for instance $5x+1$ or $181x+1$. The interesting part of both variants just mentioned is that for these sequences the existence of cycles is known. They serve as the basis for further investigations of the problem.

5.2 Further Research

In subsequent studies, the properties of vertices in H_C might be elaborated upon more closely by taking into account a vertex's label as well as its properties.

A.1 A brief note on the tree $H_{C,1}$

A special case of Collatz trees is the graph $H_{C,1}$ – the $x + 1$ variant of H_C . In this case, any sequence of successive nodes along the path from v_n down to v_1 is strictly monotonically increasing. If we run reverse to the edge direction (towards the root), then of course the node sequence is strictly monotonically decreasing. Figure A.1 shows a portion of the graph $H_{C,1}$ starting in its root.

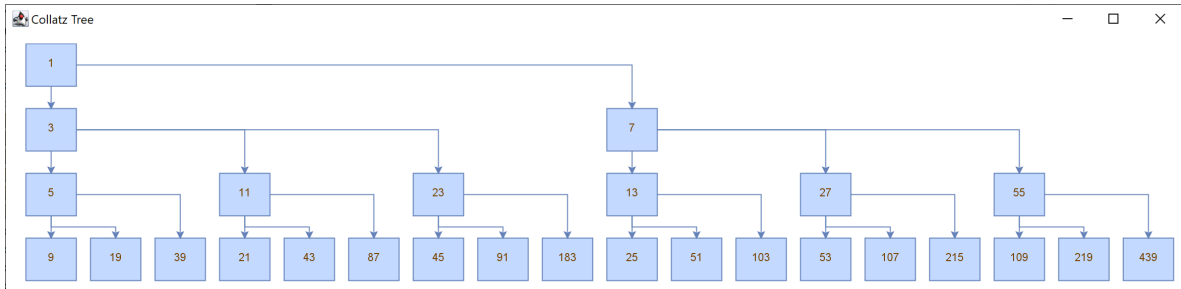


Figure A.1: Section of the graph $H_{C,1}$ starting at its root (without branches that reflect a subsequence containing the trivial cycle)

A.2 The sum of reciprocated vertices depending only on v_1

One condition deduced from theorem 2.1 is the product condition 4.4, which specifies the validity of the cycle-alpha's upper limit. This condition requires the sum $\frac{1}{kv_1} + \frac{1}{kv_2} + \frac{1}{kv_3} + \dots$ to be limited. In order to formulate this sum independently from the successive vertices v_2, v_3, \dots , we substitute these as follows:

$$\begin{aligned}
v_1 &= v_1 \\
v_2 &= \frac{kv_1 + 1}{2^{\alpha_1}} \\
v_3 &= \frac{k^2v_1 + k + 2^{\alpha_1}}{2^{\alpha_1 + \alpha_2}} \\
v_4 &= \frac{k^3v_1 + k^2 + k \cdot 2^{\alpha_1} + 2^{\alpha_1 + \alpha_2}}{2^{\alpha_1 + \alpha_2 + \alpha_3}} \\
&\vdots \\
v_{n+1} &= \frac{k^n v_1 + \sum_{j=1}^n k^{j-1} 2^{\alpha_1 + \dots + \alpha_n - \sum_{l>n-j} \alpha_l}}{2^{\alpha_1 + \dots + \alpha_n}}
\end{aligned} \tag{A.1}$$

$$v_{n+1} = \frac{k^n v_1 + \sum_{j=1}^n k^{j-1} 2^{\alpha_1 + \dots + \alpha_n - \sum_{l>n-j} \alpha_l}}{2^{\alpha_1 + \dots + \alpha_n}} \tag{A.2}$$

The sum of the reciprocated vertices can be expressed as a term that depends from v_1 and from the number of contracted edges, id est the number of divisions by two, between two successive vertices $\alpha_1, \alpha_2, \alpha_3, \dots$:

$$\sum_{i=1}^{n+1} \frac{1}{kv_i} = \frac{1}{k} \left(\frac{1}{v_1} + \sum_{i=1}^n \frac{1}{v_{i+1}} \right) = \frac{1}{k} \left(\frac{1}{v_1} + \sum_{i=1}^n \frac{2^{\alpha_1 + \dots + \alpha_i}}{k^i v_1 + \sum_{j=1}^i k^{j-1} 2^{\alpha_1 + \dots + \alpha_n - \sum_{l>i-j} \alpha_l}} \right)$$

A.3 The product of reciprocated vertices incremented by one

In a similar way to deduce the sum of reciprocated vertices depending only on v_1 as performed in A.2, we evolve the product formula depending only on v_1 :

$$\prod_{i=1}^{n+1} \left(1 + \frac{1}{kv_i} \right) = 1 + \frac{2^{\alpha_1 + \dots + \alpha_n} + k \cdot 2^{\alpha_1 + \dots + \alpha_{n-1}} + \dots + k^{n-1} \cdot 2^{\alpha_1} + k^n}{k^{n+1} v_1} \tag{A.3}$$

$$= 1 + \frac{2^{\alpha_1 + \dots + \alpha_n} + k \cdot \sum_{j=1}^i k^{j-1} 2^{\alpha_1 + \dots + \alpha_n - \sum_{l>i-j} \alpha_l}}{k^{n+1} v_1} \tag{A.4}$$

$$\begin{aligned}
&= 1 + \frac{2^{\alpha_1 + \dots + \alpha_n} + k \cdot (v_{n+1} \cdot 2^{\alpha_1 + \dots + \alpha_n} - k^n v_1)}{k^{n+1} v_1} \\
&= \frac{2^{\alpha_1 + \dots + \alpha_n} (1 + kv_{n+1})}{k^{n+1} v_1}
\end{aligned} \tag{A.5}$$

We inserted the sum used in equation A.2 into the above-given equation A.3 and then obtained equation A.4. Let us divide this product by the last factor and consider the product in the condition for cycle-alpha's upper limit, which iterates to n instead of $n+1$:

$$\prod_{i=1}^n \left(1 + \frac{1}{kv_i} \right) = \frac{\prod_{i=1}^{n+1} \left(1 + \frac{1}{kv_i} \right)}{\frac{kv_{n+1} + 1}{kv_{n+1}}} = \frac{2^{\alpha_1 + \dots + \alpha_n} (1 + kv_{n+1}) kv_{n+1}}{k^{n+1} v_1 (kv_{n+1} + 1)} = \frac{2^{\alpha_1 + \dots + \alpha_n} v_{n+1}}{k^n v_1} \tag{A.6}$$

The above-shown equation A.6 becomes simplified, when we replaced the numerator by equation A.5. The question which sequence maximizes its last member v_{n+1} ties into the question: Which sequence maximizes the product? The product formula A.6 does not depend from all vertices v_1, v_2, \dots, v_n , it depends only from $2^\alpha = 2^{\alpha_1 + \dots + \alpha_n}$, from the first vertex v_1 and the final one v_{n+1} .

- [1] Christian Koch. *Collatz Python Library*. <https://github.com/c4ristian/collatz>. 2020.
- [2] E. Sultanow, D. Volkov, and S. Cox. “Introducing a Finite State Machine for Processing Collatz Sequences”. In: *International Journal of Pure Mathematical Sciences* 19 (2017), pp. 10–19.
- [3] S. W. Williams. “Million Buck Problems”. In: *National Association of Mathematicians Newsletter* 31.2 (2000), pp. 1–3.
- [4] P. S. Bruckman. “RETRACTED ARTICLE: A proof of the Collatz conjecture”. In: *International Journal of Mathematical Education in Science and Technology* 39.3 (2008), pp. 403–407. doi: [10.1080/00207390701691574](https://doi.org/10.1080/00207390701691574).
- [5] J. C. Lagarias. *The Ultimate Challenge: The $3x+1$ Problem*. Providence, RI: American Mathematical Society, 2010. ISBN: 978-0821849408.
- [6] J. C. Lagarias. “The $3x + 1$ Problem: An Annotated Bibliography (1963-1999)”. In: *ArXiv Mathematics e-prints* (2011). eprint: [math/0309224v13](https://arxiv.org/abs/math/0309224v13).
- [7] J. C. Lagarias. “The $3x + 1$ Problem: An Annotated Bibliography, II (2000-2009)”. In: *ArXiv Mathematics e-prints* (2012). eprint: [math/0608208v6](https://arxiv.org/abs/math/0608208v6).
- [8] J. C. Lagarias. “The $3x + 1$ Problem and Its Generalizations”. In: *The American Mathematical Monthly* 92.1 (1985), pp. 3–23.
- [9] S. Kahermanes. *Collatz Conjecture*. Tech. rep. Math 301 Term Paper. San Francisco State University, 2011.
- [10] M. Klisse. “Das Collatz-Problem: Lösungs- und Erklärungsansätze für die 1937 von Lothar Collatz entdeckte $(3n+1)$ -Vermutung”. 2010.
- [11] C. A. Feinstein. “The Collatz $3n+1$ Conjecture is Unprovable”. In: *Global Journal of Science Frontier Research Mathematics and Decision Sciences* 12.8 (2012), pp. 13–15.
- [12] E. Akin. “Why is the $3x + 1$ Problem Hard?” In: *Chapel Hill Ergodic Theory Workshops*. Ed. by I. Assani. Vol. 356. Contemporary Mathematics. Providence, RI: American Mathematical Society, 2004, pp. 1–20. doi: <http://dx.doi.org/10.1090/conm/364>.
- [13] D. J. Bernstein and J. C. Lagarias. “The $3x + 1$ Conjugacy Map”. In: *Canadian Journal of Mathematics* 48 (1996), pp. 1154–1169.
- [14] P. Michel. “Simulation of the Collatz $3x + 1$ function by Turing machines”. In: *ArXiv Mathematics e-prints* (2014). eprint: [1409.7322v1](https://arxiv.org/abs/1409.7322v1).

- [15] L. Berg and G. Meinardus. “Functional Equations Connected With The Collatz Problem”. In: *Results in Mathematics* 25.1 (1994), pp. 1–12. doi: [10.1007/BF03323136](https://doi.org/10.1007/BF03323136).
- [16] L. Berg and G. Meinardus. “The $3n+1$ Collatz Problem and Functional Equations”. In: *Rostocker Mathematisches Kolloquium*. Vol. 48. Rostock, Germany: University of Rostock, 1995, pp. 11–18.
- [17] G. Opfer. “An analytic approach to the Collatz $3n + 1$ Problem”. In: *Hamburger Beiträge zur Angewandten Mathematik* 2011-09 (2011).
- [18] B. de Weger. *Comments on Opfer’s alleged proof of the $3n + 1$ Conjecture*. Tech. rep. Eindhoven University of Technology, 2011.
- [19] S. Andrei and C. Masalagiu. “About the Collatz conjecture”. In: *Acta Informatica* 35.2 (1998), pp. 167–179. doi: [10.1007/s002360050117](https://doi.org/10.1007/s002360050117).
- [20] S. Kak. *Digit Characteristics in the Collatz $3n+1$ Iterations*. Tech. rep. Oklahoma State University, 2014. URL: <https://subhask.okstate.edu/sites/default/files/collatz4.pdf>.
- [21] R. Terras. “A stopping time problem on the positive integers”. In: *Acta Arithmetica* 30.3 (1976), pp. 241–252.
- [22] T. Oliveira e Silva. “Maximum Excursion and Stopping Time Record-Holders for the $3x + 1$ Problem: Computational Results”. In: *Mathematics of Computation* 68.225 (1999), pp. 371–384.
- [23] M. A. Idowu. “A Novel Theoretical Framework Formulated for Information Discovery from Number System and Collatz Conjecture Data”. In: *Procedia Computer Science* 61 (2015), pp. 105–111.
- [24] G. J. Wirsching. *The Dynamical System Generated by the $3n+1$ Function*. Springer, 1998. doi: [10.1007/BFb0095985](https://doi.org/10.1007/BFb0095985).
- [25] G. Walz, ed. *Lexikon der Mathematik*. 2nd ed. Vol. 1. Springer, 2017. doi: [10.1007/978-3-662-53498-4](https://doi.org/10.1007/978-3-662-53498-4).
- [26] T. Tao. “Almost all orbits of the collatz map attain almost bounded values”. In: *ArXiv Mathematics e-prints* (2020). eprint: [1909.03562v3](https://arxiv.org/abs/1909.03562v3).
- [27] J.-P. Allouche. “Sur la conjecture de “Syracuse-Kakutani-Collatz””. In: *Séminaire de Théorie des Nombres de Bordeaux* (1978), pp. 1–15. issn: 09895558. URL: <https://www.jstor.org/stable/44166344>.
- [28] I. Korec. “A density estimate for the $3x+1$ problem”. In: *Mathematica Slovaca* 44.1 (1994), pp. 85–89. URL: <http://dml.cz/dmlcz/133225>.
- [29] M. Trümper. “The Collatz Problem in the Light of an Infinite Free Semigroup”. In: *Chinese Journal of Mathematics* (2014), pp. 105–111. doi: <http://dx.doi.org/10.1155/2014/756917>.
- [30] S. Kohl. “On conjugates of Collatz-type mappings”. In: *International Journal of Number Theory* 4.1 (2008), pp. 117–120. doi: <http://dx.doi.org/10.1142/S1793042108001237>.
- [31] K. Hicks et al. “A Polynomial Analogue of the $3n + 1$ Problem”. In: *The American Mathematical Monthly* 115.7 (2008), pp. 615–622.
- [32] B. Snapp and M. Tracy. “The Collatz Problem and Analogues”. In: *Journal of Integer Sequences* 11.4 (2008).

- [33] C. Löh. *Geometric Group Theory: An Introduction*. Springer, 2017. DOI: <https://doi.org/10.1007/978-3-319-72254-2>.
- [34] J. A. Bondy and U. S. R. Murty. *Graph Theory with Applications*. Elsevier Science, 1976. ISBN: 0-444-19451-7.
- [35] C. P. Bonnington and C. H.C. Little. *The Foundations of Topological Graph Theory*. Springer, 1995. DOI: [10.1007/978-1-4612-2540-9](https://doi.org/10.1007/978-1-4612-2540-9).
- [36] E. A. Bender and S. G. Williamson. *Mathematics for Algorithm and System Analysis*. Dover, 2005. ISBN: 0-486-44250-0.
- [37] J. Almeida. “Profinite semigroups and applications”. In: *Structural Theory of Automata, Semigroups, and Universal Algebra*. Ed. by V. B. Kudryavtsev and I. G. Rosenberg. Vol. 207. NATO Science Series II: Mathematics, Physics and Chemistry. Dordrecht, Netherlands: Springer, 2005, pp. 1–45.
- [38] R. Johnsonbaugh. *Discrete Mathematics*. 8th ed. Pearson, 2017. ISBN: 0-321-96468-3.
- [39] S. Mac Lane and G. Birkhoff. *Algebra*. 3rd ed. AMS Chelsea Publishing, 1999. ISBN: 0821816462.
- [40] V. Novák, I. Perfilieva, and J. Močkoř. *Mathematical Principles of Fuzzy Logic*. Springer, 1999. DOI: [10.1007/978-1-4615-5217-8](https://doi.org/10.1007/978-1-4615-5217-8).
- [41] R. Angot-Pellissier. “The Relation Between Logic, Set Theory and Topos Theory as It Is Used by Alain Badiou”. In: *The Road to Universal Logic: Festschrift for the 50th Birthday of Jean-Yves Beziau*. Ed. by A. Koslow and A. Buchsbaum. Vol. 2. Birkhäuser, 2015, pp. 181–200. DOI: [10.1007/978-3-319-15368-1](https://doi.org/10.1007/978-3-319-15368-1).
- [42] A. Ya. Helemskii. *Lectures and Exercises on Functional Analysis*. American Mathematical Society, 2006. ISBN: 0-8218-4098-3.
- [43] R. Sedgewick and K. Wayne. *Algorithms*. 4th ed. Upper Saddle River, NJ: Addison-Wesley, 2011. ISBN: 978-0-321-57351-3.
- [44] K. H. Rosen. *Discrete Mathematics and Its Applications*. 7th ed. McGraw-Hill, 2011. ISBN: 978-0-07-338309-5.
- [45] B. Korte and J. Vygen. *Combinatorial Optimization: Theory and Algorithms*. 6th ed. Springer, 2018. DOI: <https://doi.org/10.1007/978-3-662-56039-6>.
- [46] K. Mehlhorn and P. Sanders. *Algorithms and Data Structures: The Basic Toolbox*. Springer, 2008. DOI: [10.1007/978-3-540-77978-0](https://doi.org/10.1007/978-3-540-77978-0).
- [47] D.-Z. Du, K.-I Ko, and Z. Hu. *Design and Analysis of Approximation Algorithms*. Springer, 2012. DOI: [10.1007/978-1-4614-1701-9](https://doi.org/10.1007/978-1-4614-1701-9).
- [48] H. Ehrig et al. *Fundamentals of Algebraic Graph Transformation*. Springer, 2006. DOI: [10.1007/3-540-31188-2](https://doi.org/10.1007/3-540-31188-2).
- [49] L. N. Childs. *A Concrete Introduction to Higher Algebra*. 3rd ed. Springer, 2006. DOI: [10.1007/978-0-387-74725-5](https://doi.org/10.1007/978-0-387-74725-5).
- [50] V. I. Voloshin. *Introduction to Graph and Hypergraph Theory*. Nova Science Publishers, 2011. ISBN: 978-1-61470-112-5.
- [51] N. A. Loehr. *Combinatorics*. 2nd ed. CRC Press, 2017. ISBN: 978-1-4987-8025-4.
- [52] J. Kleinnijenhuis and A. M. Kleinnijenhuis. “Pruning the binary tree, proving the Collatz conjecture”. In: *ArXiv Mathematics e-prints* (2020). eprint: [arXiv:2008.13643v1](https://arxiv.org/abs/2008.13643v1).

- [53] I. Aberkane. “On the Syracuse conjecture over the binary tree”. In: *Hyper Articles en Ligne* (2017). eprint: hal-01574521.
- [54] I. Aberkane. *At least almost all orbits of the Collatz map attain bounded values, and other significant corollaries on the Syracuse problem*. <http://idrissaberkane.org/index.php/2020/01/28/derniere-publication-en-anglais/>. Jan. 2020.
- [55] K. Conrow. *The Structure of the Collatz Graph; A Recursive Production of the Predecessor Tree; Proof of the Collatz $3x+1$ Conjecture*. 2010. URL: <http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.423.3396>.
- [56] F. L. Bauer. *Historische Notizen zur Informatik*. Springer, 2009. DOI: [10.1007/978-3-540-85790-7](https://doi.org/10.1007/978-3-540-85790-7).
- [57] Z. B. Batang. “Integer patterns in Collatz sequences”. In: *ArXiv Mathematics e-prints* (2019). eprint: [arXiv:1907.07088v2](https://arxiv.org/abs/1907.07088v2).
- [58] J. Kleinnijenhuis and A. M. Kleinnijenhuis. “The Collatz tree is a Hilbert hotel: a proof of the $3n + 1$ conjecture”. In: *ArXiv Mathematics e-prints* (2020). eprint: [arXiv:2008.13643v2](https://arxiv.org/abs/2008.13643v2).
- [59] C. Hercher. “Über die Länge nicht-trivialer Collatz-Zyklen”. In: *Die Wurzel* (Hefte 6 und 7 2018).
- [60] E. Roosendaal. *On the $3x + 1$ problem*. 2020. URL: <http://www.ericr.nl/wondrous/>.
- [61] The OEIS Foundation. *Self-contained numbers: odd numbers n whose Collatz sequence contains a higher multiple of n* . 2020. URL: <https://oeis.org/A005184>.
- [62] R. K. Guy. *Unsolved Problems in Number Theory*. 3rd ed. Springer, 2004. ISBN: 978-1-4419-1928-1. DOI: [10.1007/978-0-387-26677-0](https://doi.org/10.1007/978-0-387-26677-0).
- [63] T. H. Cormen et al. *Introduction to Algorithms*. 3rd ed. The MIT Press, 2009. ISBN: 978-0-262-03384-8.
- [64] G. Teschl and S. Teschl. *Mathematik für Informatiker*. 4th ed. Vol. 1. Springer Vieweg, 2013. DOI: [10.1007/978-3-642-37972-7](https://doi.org/10.1007/978-3-642-37972-7).
- [65] J. Wolfart. *Einführung in die Zahlentheorie und Algebra*. 2nd ed. Wiesbaden, Germany: Vieweg+Teubner, 2011. ISBN: 978-3-8348-1461-6.
- [66] O. Forster. *Algorithmische Zahlentheorie*. 2nd ed. Wiesbaden, Germany: Springer Spektrum, 2015. ISBN: 978-3-658-06539-3. DOI: [10.1007/978-3-658-06540-9](https://doi.org/10.1007/978-3-658-06540-9).
- [67] S. Müller-Stach and J. Piontowski. *Elementare und algebraische Zahlentheorie*. 2nd ed. Wiesbaden, Germany: Vieweg+Teubner, 2011. ISBN: 978-3-8348-1256-8.
- [68] M. Schubert. *Mathematik für Informatiker*. Wiesbaden, Germany: Vieweg+Teubner, 2009. ISBN: 978-3-8351-0157-9.
- [69] A. T. Benjamin. *Discrete Mathematics*. Chantilly, VA: The Great Courses, 2009.
- [70] F. Modler and M. Kreh. *Tutorium Analysis 2 und Lineare Algebra 2*. 2nd ed. Heidelberg, Germany: Spektrum Akademischer Verlag, 2012. ISBN: 978-3-8274-2895-0.
- [71] D. R. Stinson and M. B. Paterson. *Cryptography: Theory and Practice*. 4th ed. Boca Raton, FL: CRC Press, 2019. ISBN: 978-1-1381-9701-5.
- [72] R. Jain. *Number Theory*. Saint Louis, MO: Washington University in Saint Louis, 2011. URL: https://www.cse.wustl.edu/~jain/cse571-11/ftp/l_08mnt.pdf.

- [73] B. Hutz. *An Experimental Introduction to Number Theory*. Vol. 31. Pure and Applied Undergraduate Texts. Providence, RI: American Mathematical Society, 2018. ISBN: 978-1-4704-3097-9.
- [74] V. Shoup. *A Computational Introduction to Number Theory and Algebra*. 2nd ed. Cambridge, UK: Cambridge University Press, 2008. ISBN: 978-0-521-51644-0. URL: <https://shoup.net/ntb/>.
- [75] K. Schwalen. *Prime Restklassengruppen: Aufbau und Eigenschaften*. Version 1.12. 2014. URL: <http://www.primath.homepage.t-online.de/Homepagedateien/PR.pdf>.
- [76] N. J. Higham. *The Princeton Companion to Applied Mathematics*. Princeton, NJ: Princeton University Press, 2015. ISBN: 978-0-691-15039-0.
- [77] G. Kersting and A. Wakolbinger. *Elementare Stochastik*. Basel, Switzerland: Birkhäuser, 2008. ISBN: 978-3-7643-8430-2.
- [78] D. R. Mazur. *Combinatorics: A Guided Tour*. MAA Textbooks. Washington, DC: The Mathematical Association of America, 2010. ISBN: 978-0-88385-762-5.
- [79] A. Kalka. “Non-associative public-key cryptography”. In: *Algebra and Computer Science*. Ed. by D. Kahrobaei, B. Cavallo, and D. Garber. Vol. 677. Contemporary Mathematics. American Mathematical Society, 2016. Chap. 5, pp. 85–112. ISBN: 978-1-4704-2303-2. DOI: <http://dx.doi.org/10.1090/conm/677/13623>.
- [80] A. Benjamin, G. Chartrand, and P. Zhang. *The Fascinating World of Graph Theory*. Princeton University Press, 2015. ISBN: 978-0-691-16381-9.
- [81] G. Chartrand and P. Zhang. *Discrete Mathematics*. Waveland Press, Inc., 2011. ISBN: 978-1-57766-730-8.
- [82] C. Koch, E. Sultanow, and S. Cox. *Divisions by Two in Collatz Sequences: A Data Science Approach*. Tech. rep. Version 2. Technische Hochschule Nürnberg, 2020. DOI: <https://doi.org/10.34646/thn/ohmdok-620>.
- [83] K. Knopp. *Theorie und Anwendung der Unendlichen Reihen*. 2nd ed. Springer, 1924. ISBN: 978-3-662-41730-0. DOI: [10.1007/978-3-662-41871-0](https://doi.org/10.1007/978-3-662-41871-0).
- [84] T. Murphy. 2006 Course 4281: Prime Numbers. 2006. URL: <https://www.maths.tcd.ie/pub/Maths/Courseware/428/>.
- [85] Eldar Sultanow. *A Java Tool for Visualizing Collatz Trees*. https://github.com/Sultanow/collatz_java. 2020.
- [86] Eldar Sultanow. *Sources for the exploration of Collatz Sequences: TeX, Mathematica and Python*. <https://github.com/Sultanow/collatz>. 2020.

About our approach

The results published in this paper have been achieved with an interdisciplinary approach. Not suprising, we applied classic mathematical theory and reasoning. Since we are convinced that the Collatz problem cannot be solved with classical maths alone, we furthermore used techniques and tools of modern data science. We combined the two fields in different ways. Firstly, we analyzed Collatz sequences and related features empirically, to derive new formulas and theorems. On the other hand, we used data science to challenge our proofs. As suggested by Karl Popper, we tried to falsify them with counterexamples. In the course of our work, we have learned that the combination of the two fields leads to a very efficient working mode.

Key findings have been explored empirically using techniques of data science. Our main tool was a Python-API, which implements the theorems of this article and is optimized for processing arbitrarily big integers within milliseconds [1]:

🔗 <https://github.com/c4ristian/collatz>

After the generated data has been exported into a comma-separated values (CSV) file, a Java tool reads that file and carries out the visualization of the corresponding Collatz trees [85]:

🔗 https://github.com/Sultanow/collatz_java

For quick experiments some notebooks may provide an efficient playground [86], but it should be noted that these are not designed for large amounts of data and professional use like Christians API does:

🔗 <https://github.com/Sultanow/collatz>

Acknowledgements



Communities

Mathematics is not the easiest discipline. Mistakes happen quickly and symbols, formulas or definitions can sometimes appear contradictory in the literature. The ambiguity/overriding of the notation for powers of an ideal and for the repeated direct product of a ring is only one example. All the more we are grateful for the numerous help from mathematics communities like the [Stack Exchange Network](#), the [MatheBoard Community](#), and [Matroids Matheplanet](#). As an example let us mention [Bill Dubuque](#), for whose instant help in matters of abstract algebra we are very grateful as for the help by many other kind people from the math communities.

Contributors in the same field

We would like to thank Jan Kleinnijenhuis and Alissa M. Kleinnijenhuis for the interesting conversations, which we are still maintaining and for proofreading our chapter on binary trees and for the interesting input on pruning these trees, which we incorporated into our research.

About Us



Eldar Sultanow is an Architect at Capgemini. In 2015, he completed his doctoral studies at the Chair of Business Information Systems and Electronic Government at the University of Potsdam. He likes mathematics, computer science and scuba diving.

✉ eldar.sultanow@capgemini.com



Christian Koch is a Data Architect at TeamBank AG and lecturer at the Institute of Technology (Technische Hochschule Georg Simon Ohm) in Nuremberg. He began his career as an IT-Consultant and has since developed analytical systems for various European banks and public institutions. Christian likes programming, data science and swimming.

✉ christian.koch@th-nuernberg.de



Sean Cox is an Analyst at RatPac-Dune Entertainment. He is an expert in mathematics of finance and econometrics. Previously he worked for The Blackstone Group as mathematician and analyst. In his spare time he pursues outdoor activities.

✉ sean.cox@ratpacent.com