

The distribution of cycles in breakpoint graphs of signed permutations

Simona Grusea

Institut de Mathématiques de Toulouse, INSA de Toulouse, Université de Toulouse.

Anthony Labarre

Department of Computer Science, K. U. Leuven, Celestijnenlaan 200A - bus 2402, 3001 Heverlee, Belgium.

Abstract

Breakpoint graphs are ubiquitous structures in the field of genome rearrangements. Their cycle decomposition has proved useful in computing and bounding many measures of (dis)similarity between genomes, and studying the distribution of those cycles is therefore critical to gaining insight on the distributions of the genomic distances that rely on it. We extend here the work initiated by Doignon and Labarre [1], who enumerated unsigned permutations whose breakpoint graph contains k cycles, to *signed* permutations, and prove explicit formulas for computing the expected value and the variance of the corresponding distributions, both in the unsigned case and in the signed case. We also show how our results can be used to derive simpler proofs of other previously known results. Finally, we compare the distribution of the number of cycles in breakpoint graphs of unsigned and signed permutations to the distributions of several well-studied genomic distances, emphasising the cases where approximations obtained in this way stand out.

Keywords: Genome rearrangements, Hultman numbers, Permutations

1. Introduction

The field of comparative genomics is concerned with quantifying similarity or divergence between organisms. Several measures have been proposed to that

end, including pattern matching based approaches or edit distances relying on a given set of biologically relevant operations. A standard example of such a method, and a *de facto* standard in phylogenetics, is the approach based on *sequence alignment*, which is motivated by the observation that genomes evolve by point mutations and aims at explaining evolution by replacements, insertions or deletions of single nucleotides (see e.g. Li and Homer [2] for a recent account of sequence alignment techniques and their uses).

However, genomes also evolve by large-scale mutations that act on whole segments of the genome, as opposed to point mutations. Examples of such mutations include *reversals*, which reverse the order of elements along a segment, *transpositions*, which move segments to another location, and *translocations*, which exchange segments that belong to different chromosomes. Many models have been proposed for studying those *genome rearrangements*, which vary according to the kinds of mutations one wants to take into account, how these should be weighted, or which objects are best suited for representing genomes (see e.g. Fertin et al. [3] for an extensive survey). Nonetheless, a striking similarity between all these models is how heavily they rely on variants of a graph first introduced by Bafna and Pevzner [4], known as the *breakpoint graph*, and its decomposition into edge- or vertex-disjoint cycles, which has proved most useful in obtaining extremely tight bounds on many genome rearrangement distances, as well as formulas for computing the exact distance in several cases. The link between several genomic distances and the number of cycles in breakpoint graphs will be discussed in more detail in Section 9.

Many mathematical questions arise when studying genome rearrangement distances, particularly concerning their distributions, as well as related statistical parameters. Since quite a few such distances can be computed or approximated using the cycle decomposition of the breakpoint graph, investigating the distribution of such cycles appears as a natural, general and effective starting point to answering those questions. We will restrict our attention in this paper to the permutation model, which can be used when all genomes under comparison consist of exactly the same genes (but in a different order) without duplications.

Breakpoint graphs can be associated to permutations, and the distribution of cycles in this case was first characterised by Doignon and Labarre [1], which later led Bóna and Flynn [5] to prove a very simple expression for the expected value of the *block-interchange distance* originally introduced by Christie [6].

However, it has often been argued that *signed permutations* provide a more realistic model of evolution, since signs can be used to represent on which strand a given DNA segment is located. Using this model, Székely and Yang [7] obtained bounds for the expectation and the variance of the number of cycles in the breakpoint graph of a random signed permutation. Using the finite Markov chain embedding technique, Grusea [8] obtained the distribution of the number of cycles in the breakpoint graph of a random signed permutation in the form of a product of transition probability matrices of a certain finite Markov chain. Her method allows to derive recurrence formulas and to compute this distribution numerically, but the computational complexity is quite high and limits the practical applications.

In this work, we obtain a new expression for computing the number of unsigned permutations whose breakpoint graph contains a given number of cycles, as well as what is to the best of our knowledge the first analytic expression for computing the number of *signed* permutations whose breakpoint graph contains a given number of cycles. The formula obtained in the signed case is complicated, but we obtain simpler formulas for a couple of restricted cases. We also use our results to derive elementary proofs of previously known results, including a binomial identity and the distribution of the number of cycles in the breakpoint graph of an unsigned permutation. We prove formulas for computing the expected value and the variance of the distribution of those cycles, both in the unsigned case and in the signed case. Finally, we also discuss how the results we obtain relate to a number of widely-studied genome rearrangement distances, and in particular, how the distribution of cycles in breakpoint graphs can be used to approximate (and in some cases, to recover exactly) the distribution of those distances.

2. Notations and definitions

We recall here a few notions that will be used throughout the paper. We assume the reader is familiar with graph theory (if not, see e.g. Diestel [9]), but nevertheless review a few useful definitions, if only to agree on notation. We will work with *non-simple* graphs, i.e. graphs that may contain *loops* (edges connecting a vertex to itself) as well as parallel edges. We will also work with both undirected and directed graphs, using $\{u, v\}$ to denote edges in the former case and (u, v) to denote arcs in the latter.

Definition 2.1. A *matching* M in a graph $G = (V, E)$ is a subset of pairwise vertex-disjoint edges of E . It is a *perfect matching* of $U \subseteq V$ if every vertex in U is incident to an edge in M .

Definition 2.2. A graph is *k-regular* if each of its vertices has degree k .

In particular, if G is a 2-regular graph, then it decomposes in a unique way into a collection of edge- and vertex-disjoint cycles, up to the ordering of cycles and to rotations of elements within each cycle (i.e., $(a, b, c, d) = (b, c, d, a)$), as well as directions in which cycles are traversed if G is undirected (i.e., $(a, b, c, d) = (d, c, b, a)$). This allows us to denote unambiguously $c(G)$ the number of cycles in G . The *length* of a cycle is the number of vertices it contains, and a *k-cycle* in G is a cycle of length k .

Definition 2.3. A graph is *hamiltonian* if it contains a cycle visiting every vertex exactly once.

We now recall a few basic notions about permutations (for more details, see e.g. Björner and Brenti [10] and Wielandt [11]).

Definition 2.4. A *permutation* of $\{1, 2, \dots, n\}$ is a bijective application of $\{1, 2, \dots, n\}$ onto itself.

The *symmetric group* S_n is the set of all permutations of $\{1, 2, \dots, n\}$, together with the usual function composition \circ , applied from right to left. We use lower case Greek letters to denote permutations, typically $\pi = \langle \pi_1 \pi_2 \cdots \pi_n \rangle$, with $\pi_i = \pi(i)$, and in particular write the *identity permutation* as $\iota = \langle 1 \ 2 \ \cdots \ n \rangle$.

Definition 2.5. The *graph* $\Gamma(\pi)$ of a permutation $\pi \in S_n$ has vertex set $\{1, 2, \dots, n\}$, and contains an arc (i, j) whenever $\pi_i = j$.

Definition 2.4 implies that $\Gamma(\pi)$ is 2-regular and as such decomposes in a unique way into disjoint cycles (up to the ordering of cycles and to rotations of elements within each cycle), which we refer to as the *disjoint cycle decomposition* of π . It is also common to refer to a permutation as a k -cycle, if the only cycle of length greater than 1 that its graph contains has length k . Figure 1 shows an example of such a decomposition. To lighten the presentation, we will shorten the notation $c(\Gamma(\pi))$ into $c(\pi)$, for a given permutation π .

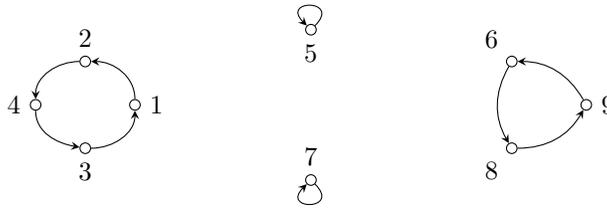


Figure 1: The graph of the permutation $\pi = \langle 2\ 4\ 1\ 3\ 5\ 8\ 7\ 9\ 6 \rangle$.

Definition 2.6. The *conjugate* of a permutation π by a permutation σ , both in S_n , is the permutation $\sigma \circ \pi \circ \sigma^{-1}$, and can be obtained by replacing every element i in the disjoint cycle decomposition of π with σ_i .

Definition 2.7. A *signed permutation* is a permutation of $\{1, 2, \dots, n\}$ where each element has an additional “+” or “-” sign.

The *hyperoctahedral group* S_n^\pm is the set of all signed permutations of n elements, together with the usual function composition \circ , applied from right to left. It is not mandatory for a signed permutation to have negative elements, so $S_n \subset S_n^\pm$ since each permutation in S_n can be viewed as a signed permutation without negative elements. To lighten the presentation, we will conform to the tradition of omitting “+” signs for positive elements.

Finally, we recall the definition of the following graph introduced by Bafna and Pevzner [4], which turned out to be an extremely useful tool for studying and solving genome rearrangement problems and which will be central to our discussions.

Definition 2.8. Given a signed permutation π in S_n^\pm , transform it into an unsigned permutation π' in S_{2n} by mapping π_i onto the sequence $(2\pi_i - 1, 2\pi_i)$ if $\pi_i > 0$, or $(2|\pi_i|, 2|\pi_i| - 1)$ if $\pi_i < 0$, for $1 \leq i \leq n$. The *breakpoint graph*

of π is the undirected bicoloured graph $BG(\pi)$ with ordered vertex set $(\pi'_0 = 0, \pi'_1, \pi'_2, \dots, \pi'_{2n}, \pi'_{2n+1} = 2n+1)$ and whose edge set is the union of the following two perfect matchings of $V(BG(\pi))$:

- black edges $\delta_B(\pi) = \{\{\pi'_{2i}, \pi'_{2i+1}\} \mid 0 \leq i \leq n\}$;
- grey edges $\delta_G = \{\{\pi'_{2i}, \pi'_{2i+1} + 1\} \mid 0 \leq i \leq n\} = \{\{2i, 2i + 1\} \mid 0 \leq i \leq n\}$.

We will often use the notation $BG(\pi) = \delta_B(\pi) \cup \delta_G$ to denote breakpoint graphs.

Genome rearrangement problems usually involve computing edit distances, i.e. the smallest number of moves needed to transform a genome into another one using only operations specified by a given set S . In the case of permutations, those distances are usually *left-invariant*, which intuitively means that genes can be relabelled so that either genome becomes ι without affecting the value of the distance to compute. Under this assumption, the pairwise genome rearrangement problem in S_n^\pm can be viewed as a constrained sorting problem, and the intuition behind the breakpoint graph construction is that black edges are meant to represent the current situation (i.e. the ordering provided by π), while grey edges are meant to represent the target situation (i.e. the ordering provided by ι). Figure 2 shows an example of a breakpoint graph. By definition, such a graph is a collection of even-length cycles that alternate black and grey edges. It can be easily seen that the example shown in Figure 2 decomposes into two such cycles.

The *length* of a cycle in a breakpoint graph differs from the traditional graph-theoretical definition that we mentioned on page 4: it is *half* the number of edges the cycle contains. Nevertheless, we will keep the terminology *k-cycle* to designate a cycle of length k , keeping in mind that its length is measured differently in the context of breakpoint graphs.

3. Cycle statistics

As is well-known (see e.g. Graham et al. [12]), the *unsigned Stirling number of the first kind* $\left[\begin{smallmatrix} n \\ k \end{smallmatrix} \right]$ counts the number of permutations in S_n which decompose into k disjoint cycles:

$$\left[\begin{smallmatrix} n \\ k \end{smallmatrix} \right] = |\{\pi \in S_n \mid c(\pi) = k\}|.$$

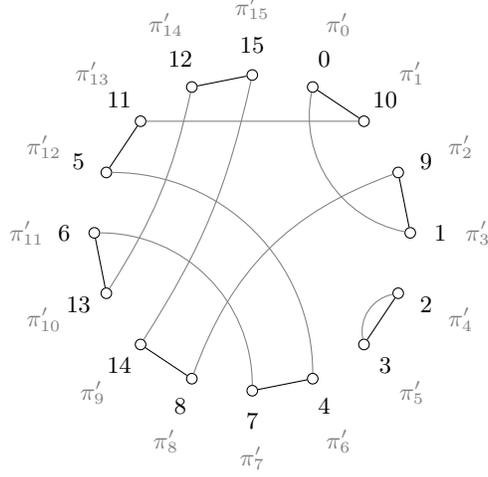


Figure 2: The breakpoint graph of $\langle -5 \ 1 \ 2 \ 4 \ -7 \ -3 \ 6 \rangle$.

Recall also that those numbers arise as coefficients in the series expansion of the *rising factorial*

$$x^{\overline{n}} = x(x+1)\cdots(x+n-1) = \sum_{k=0}^n \begin{bmatrix} n \\ k \end{bmatrix} x^k \quad (1)$$

and of the *falling factorial*

$$x^{\underline{n}} = x(x-1)\cdots(x-n+1) = \sum_{k=0}^n (-1)^{n-k} \begin{bmatrix} n \\ k \end{bmatrix} x^k. \quad (2)$$

Signing the elements of a permutation does not change its disjoint cycle decomposition, so the number of *signed* permutations that decompose into k disjoint cycles is $2^n \begin{bmatrix} n \\ k \end{bmatrix}$. We are interested in the following analogues of the Stirling number of the first kind, based on the cycle decomposition of the breakpoint graph.

Definition 3.1. The *Hultman number* $\mathcal{S}_H(n, k)$ counts the number of permutations in S_n whose breakpoint graph decomposes into k cycles:

$$\mathcal{S}_H(n, k) = |\{\pi \in S_n \mid c(BG(\pi)) = k\}|.$$

The *signed Hultman number* $\mathcal{S}_H^\pm(n, k)$ counts the number of permutations in S_n^\pm whose breakpoint graph decomposes into k cycles:

$$\mathcal{S}_H^\pm(n, k) = |\{\pi \in S_n^\pm \mid c(BG(\pi)) = k\}|.$$

It is clear from Definition 2.8 that the number of cycles in any breakpoint graph is at least one and at most $n + 1$. Hultman numbers were so named by Doignon and Labarre [1] after Axel Hultman, who first raised the question of computing those numbers [13]. The authors obtained an explicit but complicated formula for computing $\mathcal{S}_H(n, k)$, as well as formulas for enumerating permutations with a given “Hultman class” (the analogue of conjugacy classes of S_n based on the breakpoint graph). Bóna and Flynn [5] later observed that they can be computed using the following much simpler expression:

$$\mathcal{S}_H(n, k) = \begin{cases} \binom{n+2}{k} / \binom{n+2}{2} & \text{if } n - k \text{ is odd,} \\ 0 & \text{otherwise,} \end{cases} \quad (3)$$

based on a formula first obtained by Kwak and Lee [14].

In the next section, we present another way of obtaining an explicit formula for the unsigned Hultman numbers, which we will use in Section 7 to derive a new and simple proof of Equation (3). In Section 5, we will prove the first explicit formula for computing the *signed* Hultman numbers.

4. A new formula for $\mathcal{S}_H(n, k)$

We will need the following results obtained by Hanlon et al. [15], whose notation we follow. For any fixed n in \mathbb{N}_0 , let

$$Q_n^{\mathbb{C}}(h, \ell) = \mathbb{E}(\text{Re}(\text{tr}((VV^t)^n))),$$

where V is a random $h \times \ell$ matrix with independent standard complex normal entries, \mathbb{E} denotes expectation, Re denotes real part, tr denotes trace and t denotes matrix transposition. For the definition and the properties of the complex normal distribution, see for example Goodman [16].

Hanlon et al. [15] give two formulas for computing $Q_n^{\mathbb{C}}(h, \ell)$, both of which we will need. The first formula¹ is:

$$Q_n^{\mathbb{C}}(h, \ell) = \sum_{\omega \in \mathcal{S}_n} h^{c(\omega)} \ell^{c(\omega \circ \omega_{(n)})}, \quad (4)$$

¹See Corollary 2.4 p. 158 of Hanlon et al. [15].

where $\omega_{(n)}$ is a fixed n -cycle in S_n . The second formula² is:

$$Q_n^{\mathbb{C}}(h, \ell) = \frac{1}{n} \sum_{i=1}^n (-1)^{i-1} \frac{(h+n-i)^n (\ell+n-i)^n}{(n-i)!(i-1)!}. \quad (5)$$

The link between the Hultman numbers and the previous results of Hanlon et al. [15] is obtained using the following result of Daignon and Labarre [1].

Corollary 4.1. [1] $\mathcal{S}_H(n, k)$ counts the number of factorisations of a fixed $(n+1)$ -cycle β into the product $\rho \circ \omega$, where ρ is an $(n+1)$ -cycle and ω a permutation in S_{n+1} with $c(\omega) = k$.

For a polynomial $P(x)$, let $[x^k]P(x)$ denote the coefficient of the monomial x^k in $P(x)$. We derive the following new expression for computing $\mathcal{S}_H(n, k)$.

Theorem 4.1. For all n in \mathbb{N}_0 , for all k in $\{1, 2, \dots, n+1\}$:

$$\mathcal{S}_H(n, k) = \frac{1}{n+1} \sum_{i=1}^{n+1} [h^k](h+n-i+1)^{n+1}. \quad (6)$$

Proof. By Corollary 4.1, $\mathcal{S}_H(n, k)$ counts the number of factorisations of a fixed $(n+1)$ -cycle β into the product $\rho \circ \omega$, with $c(\rho) = 1$ and $c(\omega) = k$. This is clearly equivalent to enumerating factorisations of ρ^{-1} into the product $\omega \circ \beta^{-1}$ under the same conditions; therefore, setting $\omega_{(n+1)}$ to β^{-1} in Equation (4), we observe that $\mathcal{S}_H(n, k)$ is the coefficient of the monomial $h^k \ell$ in the polynomial $Q_{n+1}^{\mathbb{C}}(h, \ell)$, hence by Equation (5) equals:

$$\mathcal{S}_H(n, k) = \frac{1}{n+1} \sum_{i=1}^{n+1} (-1)^{i-1} \frac{[h^k](h+n-i+1)^{n+1} \times [\ell](\ell+n-i+1)^{n+1}}{(n-i+1)!(i-1)!}.$$

Since for every i in $\{1, 2, \dots, k+1\}$ we have

$$\begin{aligned} & [\ell](\ell+n-i+1)^{n+1} \\ &= [\ell](\ell+n-i+1)(\ell+n-i) \cdots (\ell+1)\ell(\ell-1)(\ell-2) \cdots (\ell-(i-1)) \\ &= (-1)^{i-1} (n-i+1)!(i-1)!, \end{aligned}$$

the above summation simplifies to the wanted expression, which completes the proof. \square

Besides providing a new relation involving Hultman numbers, our new formula will prove useful in obtaining simple proofs of known results, as we will see in Sections 7 and 8. Moreover, we think that the interest of our formula also lies in the fact that the method used to prove it extends to the signed case.

²See Theorem 2.5 p. 158 of Hanlon et al. [15].

5. An explicit formula for $\mathcal{S}_H^\pm(n, k)$

We now turn our attention to the problem of computing *signed* Hultman numbers, which we solve using ideas similar to those presented in the previous section. The result is obtained by characterising the 2-regular graphs that correspond to actual breakpoint graphs (Lemma 5.1 page 12), and then relating that characterisation to an enumeration result by Hanlon et al. [15].

5.1. Preliminaries

Following Hanlon et al. [15], for some fixed n in \mathbb{N}_0 , let

$$Q_n^{\mathbb{R}}(h, \ell) = \mathbb{E}(\text{tr}((VV^t)^n)),$$

where V is again a random $h \times \ell$ matrix, but this time with independent standard *real* normal entries. Hanlon et al. [15] obtain two formulas for $Q_n^{\mathbb{R}}(h, \ell)$.

Let \mathcal{F}_n denote the set of perfect matchings of $\{0, 1, 2, \dots, 2n - 1\}$. In particular, let $\varepsilon \in \mathcal{F}_n$ be the *identity perfect matching* $\{\{i, n + i\} \mid 0 \leq i \leq n - 1\}$.

The first formula³ for $Q_n^{\mathbb{R}}(h, \ell)$ is:

$$Q_n^{\mathbb{R}}(h, \ell) = \sum_{\delta \in \mathcal{F}_n} h^{c(\varepsilon \cup \delta)} \ell^{c(\delta \cup \delta_{(n)})}, \quad (7)$$

where $\delta_{(n)}$ is a fixed perfect matching such that $\varepsilon \cup \delta_{(n)}$ is hamiltonian.

The second formula is based on partitions rather than on perfect matchings.

Definition 5.1. [17] A (integer) *partition* $\lambda = (\lambda_1, \lambda_2, \dots, \lambda_l)$ is a finite sequence of integers called *parts* such that $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_l \geq 0$. Its *length* is the number of non-zero parts it contains, and if $\sum_{i=1}^l \lambda_i = n$, we call λ a *partition of n* , which we write as $\lambda \vdash n$.

We consider any two partitions to be equivalent if we obtain the same sequence when removing all parts that equal 0. The notation $\lambda = (1^{m_1} 2^{m_2} \dots r^{m_r})$ is also frequently used, and expresses the fact that exactly m_i parts of λ equal i . The reader must therefore bear in mind that when working with partitions, the notation a^b is more often to be understood in the previous meaning, and not as “ a to the power b ”.

³See Corollary 3.6 of Hanlon et al. [15].

The second formula⁴ for $Q_n^{\mathbb{R}}(h, \ell)$ is:

$$Q_n^{\mathbb{R}}(h, \ell) = \sum_{\lambda} c_{\lambda}(2) F_{\lambda}(h) F_{\lambda}(\ell), \quad (8)$$

where:

- λ ranges over all partitions of n of the form $(a, b, 1^{n-a-b})$, with either $a \geq b \geq 1$ or $a = n$ and $b = 0$,
- the function $F_{\lambda} : \mathbb{R} \rightarrow \mathbb{R}$ is defined as:

$$F_{\lambda}(x) = 2^{a-b} (x/2 + a - 1)^{a-b} (x + 2b - 2)^{n-a+b}, \quad (9)$$

- and the coefficients $c_{\lambda}(2)$ are given as follows:

$$c_{\lambda}(2) = \frac{(-1)^{n+a-b+1} 2^{a-b+1} n (2a - 2b + 1) (a - 1)!}{(n + a - b + 1)^2 (n - a + b)^2 (n - a - b)! (2a - 1)! (b - 1)!}, \quad (10)$$

if $\lambda = (a, b, 1^{n-a-b})$, with $a \geq b \geq 1$, and

$$c_{\lambda}(2) = \frac{2^n n!}{(2n)!}, \text{ if } \lambda = (n). \quad (11)$$

The numbers $c_{\lambda}(2)$ appear as coefficients in the expansion of the n^{th} power-sum function in terms of zonal polynomials. For definitions and details, see for example Macdonald [17].

5.2. Characterising valid breakpoint graphs

Recall that a breakpoint graph is a 2-regular graph that is the union of two perfect matchings of $\{0, 1, \dots, 2n + 1\}$. We now make the connection between signed Hultman numbers and the previously mentioned results explicit.

Definition 5.2. A *configuration* is the union of two perfect matchings δ_B and δ_G of $\{0, 1, \dots, 2n + 1\}$, where $\delta_G = \{\{2i, 2i + 1\} \mid 0 \leq i \leq n\}$.

Note that the above definition only slightly generalises Definition 2.8, by allowing any choice of a perfect matching for δ_B , whereas there are implicit

⁴See Theorem 5.4 of Hanlon et al. [15].

constraints on the choice of δ_B in the definition of the breakpoint graph. By definition, every breakpoint graph is a configuration, but not every configuration is a breakpoint graph, as we will see below shortly. The following notion will help us characterise configurations that are breakpoint graphs.

Definition 5.3. The *complement* of a configuration $C = \delta_B \cup \delta_G$, denoted by $\overline{C} = \delta_B \cup \overline{\delta_G}$, is obtained by replacing δ_G with $\overline{\delta_G} = \{\{2i - 1, 2i\} \mid 1 \leq i \leq n\} \cup \{0, 2n + 1\}$.

Before stating our characterisation of breakpoint graphs, we wish to stress that Elias and Hartman [18] previously used a similar but different notion of complementation (they replace δ_B with $\overline{\delta_B}$ – whose definition we will omit here – whereas we replace δ_G with $\overline{\delta_G}$) to characterise valid breakpoint graphs of *unsigned* permutations. This is not enough for our purpose, which is why we generalise their result below to encompass *signed* permutations as well.

Lemma 5.1. *A configuration $\delta_B \cup \delta_G$ is the breakpoint graph of some signed permutation π if and only if the complement configuration $\delta_B \cup \overline{\delta_G}$ is hamiltonian.*

Proof. We can easily see that the complement $\overline{BG(\pi)}$ of a breakpoint graph is hamiltonian, since its edges are $\{\{\pi'_i, \pi'_{i+1}\} \mid 0 \leq i \leq 2n\} \cup \{0, 2n + 1\}$.

Reciprocally, if the complement $\delta_B \cup \overline{\delta_G}$ of a configuration is hamiltonian, then we can recover the elements of an unsigned permutation $\pi' = \langle 0 \ \pi'_1 \ \pi'_2 \ \cdots \ \pi'_{2n} \ 2n + 1 \rangle$ by visiting the vertices along the hamiltonian cycle as follows: take $0 = \pi'_0$ as starting point, and follow the edge in δ_B that is incident to 0, setting the value of π'_1 to the other endpoint of that edge. We then keep following the cycle, assigning the label of the i^{th} encountered vertex to π'_i as we go, ending with $2n + 1 = \pi'_{2n+1}$. Note that for every $0 \leq i \leq n$, the edge $\{\pi'_{2i+1}, \pi'_{2i+2}\}$ belongs to $\overline{\delta_G}$, and therefore we have $|\pi'_{2i+1} - \pi'_{2i+2}| = 1$. From the unsigned permutation π' , we can therefore easily recover the corresponding signed permutation π in S_n^\pm , whose breakpoint graph is $\delta_B \cup \delta_G$. \square

Figure 3(a) shows the complement of the breakpoint graph of Figure 2 (page 7), which is hamiltonian. On the other hand, the complement of the configuration shown in Figure 3(b) is not hamiltonian. We now show that Equation (7) remains valid when replacing the identity perfect matching ε with the perfect matching δ_G and choosing $\overline{\delta_G}$ as the fixed perfect matching $\delta_{(n+1)}$, which clearly satisfies the condition that $\delta_G \cup \overline{\delta_G}$ is hamiltonian as required. The proof can be easily generalised to any choice of a perfect matching $\tau_{(n+1)}$ such that

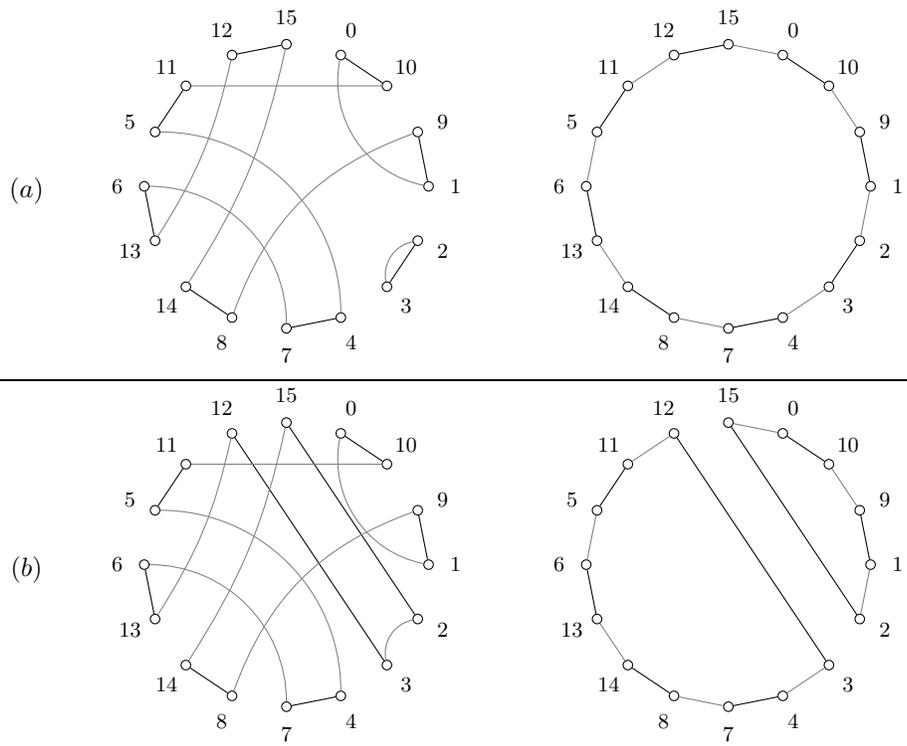


Figure 3: (a) The complement of the breakpoint graph from Figure 2 is hamiltonian; (b) a configuration whose complement is not hamiltonian.

$\delta_G \cup \tau_{(n+1)}$ is hamiltonian, but the following statement will be sufficient for our purposes.

Lemma 5.2. *For any n in \mathbb{N}_0 :*

$$Q_{n+1}^{\mathbb{R}}(h, \ell) = \sum_{\tau \in \mathcal{F}_{n+1}} h^{c(\delta_G \cup \tau)} \ell^{c(\tau \cup \overline{\delta_G})}. \quad (12)$$

Proof. First, let us note that every perfect matching ϕ in \mathcal{F}_{n+1} can be seen as a fixed-point-free involution, i.e. a permutation of $\{0, 1, 2, \dots, 2n+1\}$ that decomposes into a collection of 2-cycles only, by viewing each edge of ϕ as a 2-cycle. Therefore, conjugating ϕ by any permutation of the same number of elements is a well-defined operation that simply renames the endpoints of the given edges. Let μ be the permutation defined by

$$\mu : \{0, 1, \dots, 2n+1\} \rightarrow \{0, 1, \dots, 2n+1\} : i \mapsto \mu(i) = \begin{cases} i/2 & \text{if } i \text{ is even,} \\ \frac{i+2n+1}{2} & \text{otherwise.} \end{cases}$$

As the example in Figure 4 shows, δ_G can be mapped onto $\varepsilon = \mu \circ \delta_G \circ \mu^{-1}$, and we fix $\delta_{(n+1)} = \mu \circ \overline{\delta_G} \circ \mu^{-1}$. Finally, observe that given any two perfect matchings ϕ_1 and ϕ_2 in \mathcal{F}_{n+1} , the graphs $\mu \circ \phi_1 \circ \mu^{-1} \cup \mu \circ \phi_2 \circ \mu^{-1}$ and $\phi_1 \cup \phi_2$ are isomorphic, and hence $c(\mu \circ \phi_1 \circ \mu^{-1} \cup \mu \circ \phi_2 \circ \mu^{-1}) = c(\phi_1 \cup \phi_2)$. Taking $\delta = \mu \circ \tau \circ \mu^{-1}$, the following relations hold:

- $c(\varepsilon \cup \delta) = c(\mu \circ \delta_G \circ \mu^{-1} \cup \mu \circ \tau \circ \mu^{-1}) = c(\delta_G \cup \tau)$,
- $c(\delta \cup \delta_{(n+1)}) = c(\mu \circ \tau \circ \mu^{-1} \cup \mu \circ \overline{\delta_G} \circ \mu^{-1}) = c(\tau \cup \overline{\delta_G})$,
- $c(\varepsilon \cup \delta_{(n+1)}) = c(\mu \circ \delta_G \circ \mu^{-1} \cup \mu \circ \overline{\delta_G} \circ \mu^{-1}) = c(\delta_G \cup \overline{\delta_G}) = 1$,

and the formula in the statement follows from the above relations, the bijectivity of conjugation, and Equation (7). \square

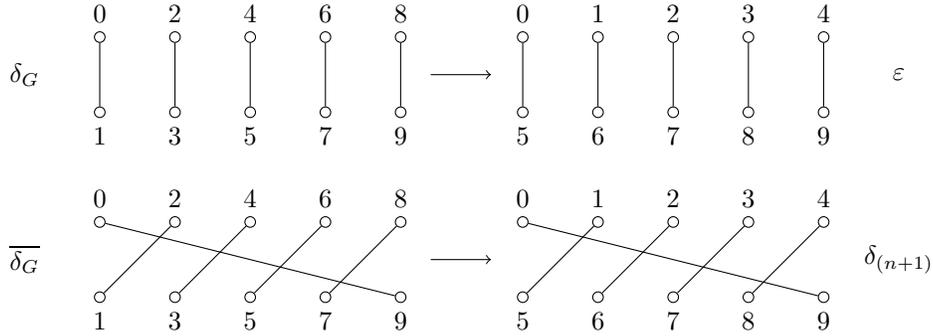


Figure 4: Mapping δ_G (resp. $\overline{\delta_G}$) onto ε (resp. $\delta_{(n+1)}$) by conjugating them by $\mu = \langle 0\ 5\ 1\ 6\ 2\ 7\ 3\ 8\ 4\ 9 \rangle$.

5.3. Enumerating breakpoint graphs with k cycles

Lemma 5.1 implies that enumerating signed permutations of n elements whose breakpoint graph decomposes into k alternating cycles is equivalent to enumerating perfect matchings τ in \mathcal{F}_{n+1} verifying $c(\delta_G \cup \tau) = k$ and $c(\tau \cup \overline{\delta_G}) = 1$, where δ_G is defined in Definition 2.8 page 5 and $\overline{\delta_G}$ is defined in Definition 5.3 page 12. Using Lemma 5.2, we thus obtain the following.

Remark 5.1. For every k in $\{1, 2, \dots, n+1\}$, $\mathcal{S}_H^\pm(n, k)$ is the coefficient of the monomial $h^k \ell$ in $Q_{n+1}^{\mathbb{R}}(h, \ell)$.

The second expression for $Q_{n+1}^{\mathbb{R}}(h, \ell)$ given in Equation (8) allows us to obtain the following explicit formula for $\mathcal{S}_H^\pm(n, k)$.

Theorem 5.1. For all n in \mathbb{N}_0 , for all k in $\{1, 2, \dots, n+1\}$:

$$\begin{aligned} \mathcal{S}_H^\pm(n, k) &= \sum_{\lambda} c_{\lambda}(2) \times [h^k] F_{\lambda}(h) \\ &\quad \times \frac{(-1)^{n-a-b} 2^{a-b-1} (2b)! (a-1)! (n-a-b+2)!}{(2b-1)b!}, \end{aligned} \quad (13)$$

where λ ranges over all partitions of $n+1$ of the form $(a, b, 1^{n-a-b+1})$, with $a \geq b \geq 1$ or $a = n+1, b = 0$, and where the function $F_{\lambda}(\cdot)$ as well as the coefficients $c_{\lambda}(2)$ follow the definitions previously given in Section 5.1⁵.

Proof. Remark 5.1 and Equation (8) yield

$$\mathcal{S}_H^\pm(n, k) = \sum_{\lambda} c_{\lambda}(2) \times [h^k] F_{\lambda}(h) \times [\ell] F_{\lambda}(\ell), \quad (14)$$

where the sum over λ , the coefficients $c_{\lambda}(2)$ and the function $F_{\lambda}(\cdot)$ are as in the statement of the present result. For a partition λ of the form $(a, b, 1^{n-a-b+1})$, with $a \geq b \geq 1$ or $a = n+1, b = 0$, it is easy to see that

$$[\ell] F_{\lambda}(\ell) = \frac{(-1)^{n-a-b} 2^{a-b-1} (2b)! (a-1)! (n-a-b+2)!}{(2b-1)b!}. \quad (15)$$

Indeed:

1. if $\lambda = (a, b, 1^{n-a-b+1})$, with $a \geq b \geq 1$, we have

$$\begin{aligned} F_{\lambda}(\ell) &= 2^{a-b} (\ell/2 + a - 1) (\ell/2 + a - 2) \cdots (\ell/2 + b) \\ &\quad \times (\ell + 2b - 2) (\ell + 2b - 3) \cdots (\ell + 1) \\ &\quad \times \ell (\ell - 1) \cdots (\ell - (n - a - b + 2)). \end{aligned}$$

⁵With the slight modification that n needs to be replaced with $n+1$.

The coefficient of ℓ in the above expression equals

$$\begin{aligned} [\ell]F_\lambda(\ell) &= 2^{a-b} \frac{(a-1)!}{(b-1)!} \times (2b-2)! (-1)^{n-a-b+2} (n-a-b+2)! \\ &= \frac{(-1)^{n-a-b} 2^{a-b-1} (2b)! (a-1)! (n-a-b+2)!}{(2b-1)b!}. \end{aligned}$$

2. if $\lambda = (n+1)$, i.e. $a = n+1$ and $b = 0$, we have

$$\begin{aligned} F_{(n+1)}(\ell) &= 2^{n+1} (\ell/2 + n)^{n+1} (\ell - 2)^0 \\ &= 2^{n+1} (\ell/2 + n) (\ell/2 + n - 1) \cdots (\ell/2 + 1) \ell/2, \end{aligned}$$

so $[\ell]F_{(n+1)}(\ell) = 2^n n!$, which verifies Equation (15).

The proof then follows from Equations (14) and (15). □

We conclude this section with Table 1, which shows a few experimental values of the signed Hultman numbers. These values were previously obtained by the first author using the method described in a previous paper of hers [8].

Note that for $k = 1$, the sequence defined by $\mathcal{S}_H^\pm(n, 1)$ for $n = 1, 2, \dots$ corresponds to sequence A001171 in the On-Line Encyclopedia of Integer Sequences [19]. As we will see in the next section, other known sequences also appear in that table.

6. Special cases

The expression obtained in Theorem 5.1 allows us to compute $\mathcal{S}_H^\pm(n, k)$ for all valid values of n and k , but we must acknowledge that even though the formula is suited for practical use, it is unfortunately quite complicated and difficult to manipulate. Simpler expressions do however exist for some particular cases, as we will show below. We will rely a lot on Lemma 5.1 in this section, and decide to use a slightly different layout for the breakpoint graph: labels are omitted for clarity, and grey edges rather than black edges are now laid out on a circle, so that computing the complement of a given configuration simply amounts to shifting grey edges sideways by one position. In order to make verifications easier for the reader, we also draw edges in the complement as dotted edges. The following particular cases are easy to verify:

Table 1: A few values of $S_H^\pm(n, k)$

$n \backslash k$	1	2	3	4	5	6	7	8	9	10	11	12
1	1	1										
2	4	3										
3	20	21	6	1								
4	148	160	65	10	1							
5	1348	1620	701	155	15	1						
6	15104	19068	9324	2247	315	21	1					
7	198144	264420	138016	38029	5908	574	28	1				
8	2998656	4166880	2325740	692088	124029	13524	966	36	1			
9	51290496	74011488	43448940	13945700	2723469	344961	27930	1530	45	1		
10	979732224	1459381440	897020784	305142068	64711856	8996295	850905	53262	2310	55	1	
11	20661458688	31674232128	20241273264	7255047116	1640552028	249029717	26004330	1910403	95304	3355	66	1

1. $\mathcal{S}_H^\pm(n, k) = 0$ for all $k < 1$ and all $k > n + 1$ (trivial);
2. $\mathcal{S}_H^\pm(n, n + 1) = 1$, since the only permutation whose breakpoint graph decomposes into $n + 1$ cycles is ι ;
3. $\mathcal{S}_H^\pm(n, n) = \binom{n+1}{2}$, since enumerating such permutations comes down to counting breakpoint graphs whose cycles all have length 1, except for one that has length 2. This in turn is equivalent to enumerating the ways in which one can connect any two of the $n + 1$ grey edges by black edges so as to obtain a valid configuration (with respect to Lemma 5.1); as can be verified on Figure 5, only one of the two possible choices of black edges (namely, configuration (b)) is valid, and the equality follows from the fact that there are $\binom{n+1}{2}$ possible ways to select two grey edges out of $n + 1$.

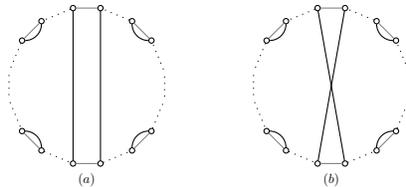


Figure 5: The two forms of 2-cycles that may arise in a breakpoint graph. Only four 1-cycles are shown in each graph, but there can be any number of them.

We now show how one can obtain a simple and explicit formula for $\mathcal{S}_H^\pm(n, n - 1)$. Although the formula is quite simple, we hope that the proof will convince the reader of the shortcomings of a case analysis in this setting.

Proposition 6.1. *For all $n \geq 1$, we have $\mathcal{S}_H^\pm(n, n - 1) = 5\binom{n+1}{4} + 4\binom{n+1}{3}$.*

Proof. Note that $\mathcal{S}_H^\pm(n, n - 1)$ is the number of permutations whose breakpoint graph contains either one 3-cycle or two 2-cycles, all other cycles having length 1 in both cases:

1. the number of permutations satisfying the first condition is the number of ways to connect three grey edges in the breakpoint graph in such a way that the complement configuration is hamiltonian (see Lemma 5.1). As Figure 6 shows, there are eight possible ways to create such a configuration, only four of which are valid (namely, configurations (a), (b), (c) and (d)). The reader can easily verify that the other configurations are invalid by replacing grey edges with dotted edges. We obtain the rightmost term in the wanted expression by noting that only four of the eight possible 3-cycles are valid, and there are $\binom{n+1}{3}$ ways to select three grey edges out of $n + 1$.

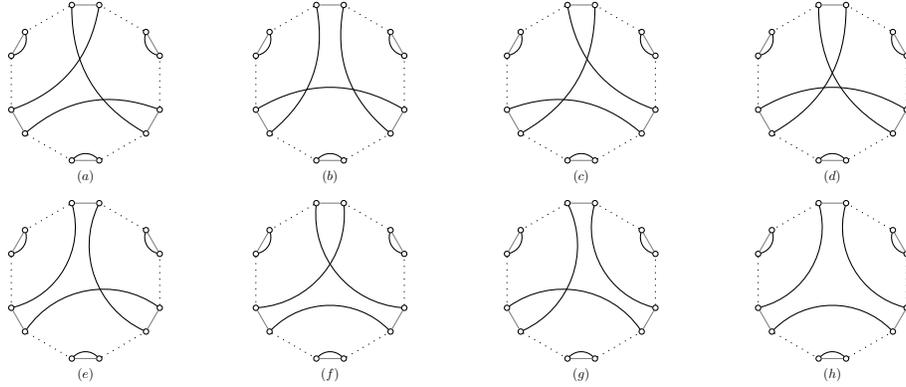


Figure 6: All possible forms of 3-cycles that may arise in a breakpoint graph. Only three 1-cycles are shown in each graph, but there can be any number of them.

2. the number of permutations satisfying the second condition can be constructed by choosing four grey edges, then connecting them by pairs while ensuring that the resulting configuration is valid. Figure 7 shows all possible configurations with two cycles of length two.

The reader can again easily verify the validity of all configurations by replacing grey edges with dotted edges. Only five possible configurations with two 2-cycles are valid (namely, configurations (b), (f), (i), (k) and (l)) out of the twelve shown in Figure 7, and there are $\binom{n+1}{4}$ ways to select two pairs of grey edges out of $n + 1$, which yields the leftmost term in the wanted expression and completes the proof.

□

7. Simpler proofs of previous results

Theorem 4.1 allows us to obtain a new proof of Bóna and Flynn's formula (Equation (3) page 8).

Corollary 7.1. [5] For all n in \mathbb{N}_0 :

$$\mathcal{S}_H(n, k) = \begin{cases} \left[\begin{smallmatrix} n+2 \\ k \end{smallmatrix} \right] / \binom{n+2}{2} & \text{if } n - k \text{ is odd,} \\ 0 & \text{otherwise.} \end{cases}$$

Proof. The key idea of the proof is the fact that, for every $i = 1, 2, \dots, n + 1$, we have

$$(h + n - i + 1)^{n+1} = \frac{1}{n + 2} \left((h - i + 1)^{\overline{n+2}} - (h - i)^{\overline{n+2}} \right), \quad (16)$$

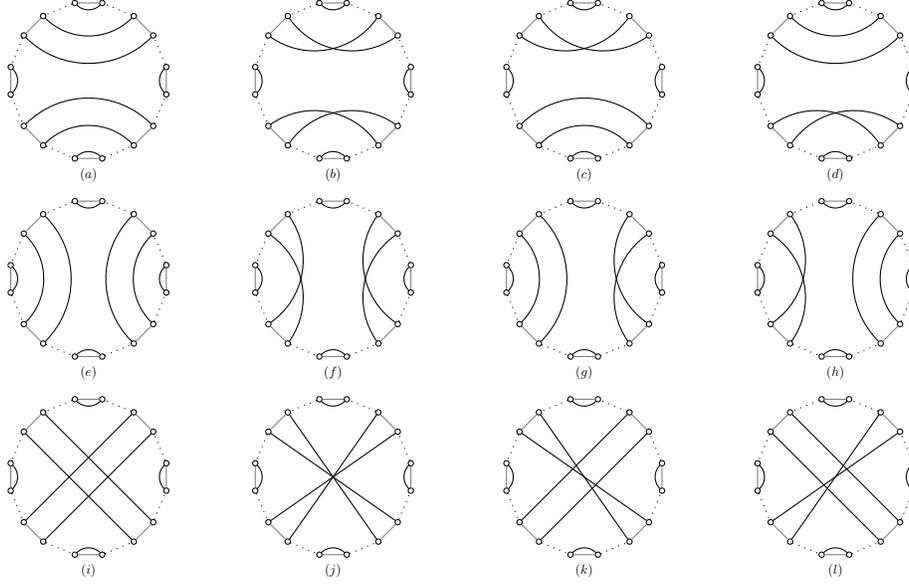


Figure 7: All possible pairs of 2-cycles that may arise in a breakpoint graph. Only four 1-cycles are shown in each graph, but there can be any number of them.

since

$$\begin{aligned}
& \frac{1}{n+2} \left((h-i+1)^{\overline{n+2}} - (h-i)^{\overline{n+2}} \right) \\
&= \frac{1}{n+2} \left((h-i+1) \cdots (h+n-i+2) - (h-i) \cdots (h+n-i+1) \right) \\
&= \frac{1}{n+2} (h-i+1) \cdots (h+n-i+1) \left((h+n-i+2) - (h-i) \right) \\
&= (h+n-i+1)^{\overline{n+1}}.
\end{aligned}$$

Summing over i in Equation (16), we obtain:

$$\begin{aligned}
& \frac{1}{n+1} \sum_{i=1}^{n+1} (h+n-i+1)^{\overline{n+1}} \\
&= \frac{1}{(n+1)(n+2)} \sum_{i=1}^{n+1} \left((h-i+1)^{\overline{n+2}} - (h-i)^{\overline{n+2}} \right) \\
&= \frac{1}{(n+1)(n+2)} \left(h^{\overline{n+2}} - (h-n-1)^{\overline{n+2}} \right) \\
&= \frac{1}{(n+1)(n+2)} \left(h^{\overline{n+2}} - h^{\overline{n+2}} \right).
\end{aligned}$$

By Equations (1) and (2), the coefficient of h^k in $h^{\overline{n+2}}$ is $\binom{n+2}{k}$ and the

coefficient of h^k in h^{n+2} is $(-1)^{n-k} \begin{bmatrix} n+2 \\ k \end{bmatrix}$. Using Equation (6), we conclude that

$$\mathcal{S}_H(n, k) = \begin{cases} \frac{2}{(n+1)(n+2)} \begin{bmatrix} n+2 \\ k \end{bmatrix} & \text{if } n-k \text{ is odd,} \\ 0 & \text{otherwise,} \end{cases}$$

which completes the proof. \square

Theorem 4.1 also allows us to obtain a simple proof of a binomial identity previously obtained by Sury et al. [20].

Corollary 7.2. [20] For all n in \mathbb{N}_0 :

$$\sum_{i=0}^n \frac{(-1)^i}{\binom{n}{i}} = (1 + (-1)^n) \frac{n+1}{n+2}.$$

Proof. Setting k to 1 in Equation (6) (page 9) yields

$$\mathcal{S}_H(n, 1) = \frac{1}{n+1} \sum_{i=1}^{n+1} (-1)^{i-1} (n-i+1)! (i-1)! = \frac{n!}{n+1} \sum_{i=0}^n \frac{(-1)^i}{\binom{n}{i}}.$$

On the other hand, as previously observed⁶ by Doignon and Labarre [1], we have:

$$\mathcal{S}_H(n, 1) = \begin{cases} \frac{2n!}{n+2} & \text{if } n \text{ is even,} \\ 0 & \text{otherwise,} \end{cases}$$

which completes the proof. \square

8. Expected value and variance of the Hultman numbers

In order to gain more insight into the distribution of the Hultman numbers, we will now investigate the question of computing the expected value and variance of the number of cycles in breakpoint graphs, both for unsigned and for signed permutations.

It will also be interesting to see how these values compare to the expected value and variance of the number of cycles in the usual disjoint cycle decomposition of a uniform random unsigned permutation π in S_n . We recall here (see e.g. Wilf [21]) the exact values of these quantities:

$$\begin{aligned} \mathbb{E}(c(\pi)) &= H_n, \\ \text{Var}(c(\pi)) &= H_n - \sum_{k=1}^n \frac{1}{k^2}, \end{aligned}$$

⁶The result can also be easily derived from Equation (3).

as well as their asymptotic behaviour when $n \rightarrow \infty$:

$$\mathbb{E}(c(\pi)) = \log(n) + \gamma + o(1), \quad (17)$$

$$\text{Var}(c(\pi)) = \log(n) + \gamma - \frac{\pi^2}{6} + o(1), \quad (18)$$

where H_n denotes the n^{th} harmonic number $H_n = \sum_{i=1}^n \frac{1}{i}$ and γ denotes the Euler-Mascheroni constant. As usual, $o(1)$ denotes a quantity that converges to 0 as $n \rightarrow \infty$.

8.1. The unsigned case

Bóna and Flynn [5] already proved a formula for computing the expected number of cycles in the breakpoint graph of a uniform random unsigned permutation. In this section we provide a new proof of their result and also give an explicit formula for the variance of this distribution. We start by computing the generating function of the Hultman numbers.

Lemma 8.1. *For all $n \in \mathbb{N}_0$, we have:*

$$F(x) = \sum_{k=0}^{n+1} \mathcal{S}_H(n, k) x^k = \frac{x^{\overline{n+2}} - x^{n+2}}{2 \binom{n+2}{2}}.$$

Proof. The derivation is straightforward:

$$\begin{aligned} \sum_{k=0}^{n+1} \mathcal{S}_H(n, k) x^k &= \frac{1}{\binom{n+2}{2}} \sum_{k=0}^{n+1} \frac{\left[\begin{smallmatrix} n+2 \\ k \end{smallmatrix} \right] - (-1)^{n+2-k} \left[\begin{smallmatrix} n+2 \\ k \end{smallmatrix} \right]}{2} x^k \quad (\text{by Equation (3)}) \\ &= \frac{1}{2 \binom{n+2}{2}} \left(\sum_{k=0}^{n+2} \left[\begin{smallmatrix} n+2 \\ k \end{smallmatrix} \right] x^k - \sum_{k=0}^{n+2} (-1)^{n+2-k} \left[\begin{smallmatrix} n+2 \\ k \end{smallmatrix} \right] x^k \right) \\ &= \frac{x^{\overline{n+2}} - x^{n+2}}{2 \binom{n+2}{2}}. \quad (\text{by Equations (1) and (2)}) \end{aligned}$$

□

Knowing the generating function allows us to easily derive the expected value and the variance of the number of cycles in the breakpoint graph of a uniform random unsigned permutation. For this purpose, we first need to compute some derivatives of the generating function.

Lemma 8.2. For all $n \in \mathbb{N}_0$, we have:

$$\begin{aligned} F(1) &= n!, \\ F'(1) &= \frac{1}{2^{\binom{n+2}{2}}} \{(n+2)!H_{n+2} + (-1)^{n-1}n!\}, \\ F''(1) &= \frac{1}{2^{\binom{n+2}{2}}} \left\{ (n+2)! \left(H_{n+2}^2 - \sum_{k=1}^{n+2} \frac{1}{k^2} \right) + 2(-1)^n n!(H_n - 1) \right\}. \end{aligned}$$

Proof. We obtain the three expressions separately.

1. For the first expression, note that, by definition, $F(1) = \sum_{k=1}^{n+1} \mathcal{S}_H(n, k)$, which is simply the total number of permutations of n elements and therefore equals $n!$.
2. We simplify the computation of $F'(x)$ by writing $x^{\overline{n+2}} = (x-1)g(x)$, with

$$g(x) = x \prod_{i=2}^{n+1} (x-i).$$

With this notation we have

$$F(x) = \frac{x^{\overline{n+2}} - (x-1)g(x)}{2^{\binom{n+2}{2}}}.$$

We thus obtain

$$F'(x) = \frac{1}{2^{\binom{n+2}{2}}} \left(x^{\overline{n+2}} \sum_{i=0}^{n+1} \frac{1}{x+i} - g(x) - (x-1)g'(x) \right).$$

At $x = 1$ we have $1^{\overline{n+2}} = (n+2)!$ and $g(1) = (-1)^n n!$, and hence the stated formula for $F'(1)$ follows.

3. Finally, the second derivative of F is given by

$$F''(x) = \frac{1}{2^{\binom{n+2}{2}}} \left(x^{\overline{n+2}} \sum_{0 \leq i \neq j \leq n+1} \frac{1}{(x+i)(x+j)} - 2g'(x) - (x-1)g''(x) \right).$$

The above sum evaluated at $x = 1$ equals

$$\begin{aligned} \sum_{0 \leq i \neq j \leq n+1} \frac{1}{(1+i)(1+j)} &= \sum_{i,j=0}^{n+1} \frac{1}{(1+i)(1+j)} - \sum_{i=0}^{n+1} \frac{1}{(1+i)^2} \\ &= \left(\sum_{i=0}^{n+1} \frac{1}{1+i} \right)^2 - \sum_{i=0}^{n+1} \frac{1}{(1+i)^2} \\ &= H_{n+2}^2 - \sum_{k=1}^{n+2} \frac{1}{k^2}. \end{aligned}$$

We also have

$$g'(x) = g(x) \left(\frac{1}{x} + \sum_{i=2}^{n+1} \frac{1}{x-i} \right),$$

and thus

$$g'(1) = g(1) \left(1 - \sum_{i=2}^{n+1} \frac{1}{i-1} \right) = (-1)^n n! (1 - H_n).$$

Using these expressions in the formula for $F''(x)$ above, evaluated at $x = 1$, gives the formula in the statement. □

The recovery of the expected value of the Hultman numbers, previously obtained by Bóna and Flynn [5], is now an easy task.

Theorem 8.1. [5] *For all $n \in \mathbb{N}_0$, the expected number of cycles in the breakpoint graph of a uniform random unsigned permutation π of n elements is*

$$\mathbb{E}(c(BG(\pi))) = H_n + \frac{1}{\lfloor (n+2)/2 \rfloor}.$$

Proof. As is well-known (see e.g. Wilf [21]), the expected value can be obtained from the generating function $F(x)$ by the formula $F'(1)/F(1)$. Using the formulas for $F(1)$ and $F'(1)$ obtained in Lemma 8.2, we obtain that the expected value of the Hultman numbers equals

$$\frac{F'(1)}{F(1)} = H_{n+2} + \frac{(-1)^{n-1}}{(n+1)(n+2)},$$

which is easily seen to be equivalent to the expression in the statement. □

Furthermore, knowing the generating function also allows us to compute the variance of the Hultman numbers. We prove the following result.

Theorem 8.2. *For all $n \in \mathbb{N}_0$, the variance of the number of cycles in the breakpoint graph of a uniform random unsigned permutation π of n elements is*

$$\text{Var}(c(BG(\pi))) = H_{n+2} - \sum_{k=1}^{n+2} \frac{1}{k^2} + \frac{(-1)^n (2H_{n+2} + 2H_n - 3)}{(n+1)(n+2)} - \frac{1}{((n+1)(n+2))^2}.$$

Proof. The variance can be obtained from the generating function $F(x)$ by the following formula (see e.g. Wilf [21]):

$$(\log F)'(1) + (\log F)''(1) = \frac{F'(1)}{F(1)} + \frac{F''(1)}{F(1)} - \left(\frac{F'(1)}{F(1)} \right)^2.$$

Using the formulas for $F(1)$, $F'(1)$ and $F''(1)$ obtained in Lemma 8.2, we obtain that the variance of the Hultman numbers equals

$$\begin{aligned}
& \frac{F'(1)}{F(1)} + \frac{F''(1)}{F(1)} - \left(\frac{F'(1)}{F(1)} \right)^2 \\
= & H_{n+2} + \frac{(-1)^{n-1}}{(n+1)(n+2)} + H_{n+2}^2 - \sum_{k=1}^{n+2} \frac{1}{k^2} + \frac{2(-1)^n(H_n - 1)}{(n+1)(n+2)} \\
& - \left(H_{n+2} + \frac{(-1)^{n-1}}{(n+1)(n+2)} \right)^2 \\
= & H_{n+2} - \sum_{k=1}^{n+2} \frac{1}{k^2} + \frac{(-1)^n(2H_{n+2} + 2H_n - 3)}{(n+1)(n+2)} - \frac{1}{((n+1)(n+2))^2}.
\end{aligned}$$

□

It is interesting to see how the mean and variance behave for large n .

Remark 8.1. *The expected value and variance of the number of cycles in the breakpoint graph of a uniform random unsigned permutation π in S_n have the following asymptotical behaviour when $n \rightarrow \infty$:*

$$\begin{aligned}
\mathbb{E}(c(BG(\pi))) &= \log(n) + \gamma + o(1), \\
\text{Var}(c(BG(\pi))) &= \log(n) + \gamma - \frac{\pi^2}{6} + o(1).
\end{aligned}$$

Proof. For the expected value, the result simply follows from the fact that $\mathbb{E}(c(BG(\pi))) = H_n + o(1)$ and $H_n = \log(n) + \gamma + o(1)$.

For the variance, first note that $\text{Var}(c(BG(\pi))) = H_{n+2} - \sum_{k=1}^{n+2} \frac{1}{k^2} + o(1)$. By further using the fact that $\log(n+2) = \log(n) + o(1)$ and the well-known result $\sum_{k=1}^{\infty} \frac{1}{k^2} = \frac{\pi^2}{6}$, the stated asymptotic formula follows. □

Interestingly, we recover exactly the same asymptotical behaviour as for the number of cycles in the usual disjoint cycle decomposition (recall Equations (17) and (18)).

8.2. The signed case

We now turn to the problem of computing the expected value and the variance of the signed Hultman numbers. As in the unsigned case, we start with the computation of the generating function for the signed Hultman numbers.

Lemma 8.3. *We have*

$$G(x) = \sum_{k=1}^{n+1} S_H^\pm(n, k)x^k = \sum_{\lambda} c_{\lambda}(2)F_{\lambda}(x)F'_{\lambda}(0),$$

where λ is subject to the same restrictions as in Theorem 5.1 page 15 and F_λ is defined as in Equation (9) page 11.

Proof. Recall (Remark 5.1 page 15) that $S_H^\pm(n, k)$ is the coefficient of the monomial $h^k \ell$ in the polynomial $Q_{n+1}^{\mathbb{R}}(h, \ell)$. If we take now $h = x$ and consider $Q_{n+1}^{\mathbb{R}}(x, \ell)$ as a polynomial only in the variable ℓ , we note that the coefficient of the monomial ℓ is obtained by summing up all the terms $S_H^\pm(n, k)x^k$, for $k = 1, \dots, n+1$. Therefore, $G(x)$ equals the coefficient of ℓ in $Q_{n+1}^{\mathbb{R}}(x, \ell)$, and hence

$$G(x) = \left. \frac{\partial}{\partial \ell} Q_{n+1}^{\mathbb{R}}(x, \ell) \right|_{\ell=0}.$$

The formula in the statement easily follows from Equation (8) page 11. \square

In order to compute the expected value and the variance of the signed Hultman numbers, we will need the following preliminary lemma.

Lemma 8.4. *Let $n \geq 1$ and λ a partition of $n+1$ of the form $(a, b, 1^{n-a-b+1})$.*

1. *In the case where $a \geq b \geq 1$, we have:*

$$\begin{aligned} F'_\lambda(0) &= \frac{(-1)^{n-a-b} 2^{a-b} (a-1)! (2b-2)! (n-a-b+2)!}{(b-1)!}, \\ F'_\lambda(1) &= \frac{(-1)^{n-a-b+1} (2a-1)! (b-1)! (n-a-b+1)!}{2^{a-b} (a-1)!}, \\ F''_\lambda(1) &= F'_\lambda(1) \{2H_{2a-1} - 2H_{n-a-b+1} - H_{a-1} + H_{b-1}\}. \end{aligned}$$

2. *In the case where $\lambda = (n+1)$, we have:*

$$\begin{aligned} F'_{(n+1)}(0) &= 2^n n!, \\ F'_{(n+1)}(1) &= \frac{(2n+1)!}{2^n n!} (H_{2n+1} - H_n/2), \\ F''_{(n+1)}(1) &= \frac{(2n+1)!}{2^n n!} \left\{ \left(H_{2n+1} - \frac{H_n}{2} \right)^2 - \sum_{k=0}^n \frac{1}{(2k+1)^2} \right\}. \end{aligned}$$

Proof. We handle both cases separately.

1. Let us first examine the case where $\lambda = (a, b, 1^{n+1-a-b})$ and $a \geq b \geq 1$. In order to simplify the proof, we write $F_\lambda(x) = x(x-1)h_\lambda(x)$, where $h_\lambda(x)$ is obtained and defined as follows:

$$\begin{aligned} F_\lambda(x) &= 2^{a-b} (x/2 + a - 1)^{\overline{a-b}} (x + 2b - 2)^{\overline{n+1-a+b}} \quad (\text{see definition}^7 \text{ page 11}) \\ &= 2^{a-b} (x/2 + a - 1)^{\overline{a-b}} (x + 2b - 2)(x + 2b - 1) \cdots (x + 1)x(x - 1) \\ &\quad \times (x - 2)(x - 3) \cdots (x - 2 + b - n + a) \\ &= x(x - 1) \underbrace{2^{a-b} (x/2 + a - 1)^{\overline{a-b}} (x + 2b - 2)^{\overline{2b-2}} (x - 2)^{\overline{n-a-b+1}}}_{=h_\lambda(x)}. \end{aligned}$$

(a) Using the above notation, we have

$$F'_\lambda(0) = -h_\lambda(0) = (-1)2^{a-b}(a-1)^{\overline{a-b}}(2b-2)!(-2)^{\overline{n-a-b+1}},$$

from which we easily obtain the wanted expression.

(b) We also have

$$\begin{aligned} F'_\lambda(1) = h_\lambda(1) &= 2^{a-b}(a-1/2)^{\overline{a-b}}(2b-1)^{\overline{2b-2}}(-1)^{\overline{n-a-b+1}} \\ &= 2^{a-b}(a-1/2)^{\overline{a-b}}(2b)!(-1)^{\overline{n-a-b+1}}, \end{aligned}$$

and obtaining the formula for $F'_\lambda(1)$ given in the statement is a simple matter, using the fact that

$$\begin{aligned} (a-1/2)^{\overline{a-b}} &= \frac{(2a-1)(2a-3)\cdots(2b+1)}{2^{a-b}} \\ &= \frac{1}{2^{a-b}} \frac{(2a-1)!}{(a-1)!2^{a-1}} \frac{(b-1)!2^{b-1}}{(2b-1)!} \\ &= \frac{(2a-1)!b!}{2^{a-b-1}(a-1)!2^{a-b}(2b)!}. \end{aligned}$$

(c) In order to simplify the computation of the second derivative, we will write $F_\lambda(x) = (x-1)g_\lambda(x)$, where

$$g_\lambda(x) = \underbrace{2^{a-b}(x/2+a-1)^{\overline{a-b}}}_{=\alpha_\lambda(x)} \underbrace{(x+2b-2)^{\overline{2b-1}}}_{=\beta_\lambda(x)} \underbrace{(x-2)^{\overline{n-a-b+1}}}_{=\gamma_\lambda(x)}.$$

With this notations, it is easy to see that $F''_\lambda(1) = 2g'_\lambda(1)$, with

$$g'_\lambda(1) = \alpha'_\lambda(1)\beta_\lambda(1)\gamma_\lambda(1) + \alpha_\lambda(1)\beta'_\lambda(1)\gamma_\lambda(1) + \alpha_\lambda(1)\beta_\lambda(1)\gamma'_\lambda(1).$$

Note that

$$\begin{aligned} \alpha'_\lambda(1) &= \alpha_\lambda(1) \left(\frac{1}{2a-1} + \frac{1}{2a-3} + \cdots + \frac{1}{2b+1} \right) \\ &= \alpha_\lambda(1) \{H_{2a-1} - H_{2b} - (H_{a-1} - H_b)/2\}, \\ \beta'_\lambda(1) &= \beta_\lambda(1) \sum_{k=1}^{2b-1} \frac{1}{k} = \beta_\lambda(1)H_{2b-1}, \\ \gamma'_\lambda(1) &= -\gamma_\lambda(1) \sum_{k=1}^{n-a-b+1} \frac{1}{k} = -\gamma_\lambda(1)H_{n-a-b+1}, \end{aligned}$$

⁷Recall, as explained in the statement of Theorem 5.1 page 15, that we must replace n with $n+1$.

and

$$\begin{aligned}\alpha_\lambda(1) &= \frac{(2a-1)!b!}{(2b)!2^{a-b-1}(a-1)!}, \\ \beta_\lambda(1) &= (2b-1)!, \\ \gamma_\lambda(1) &= (-1)^{n-a-b+1}(n-a-b+1)!\end{aligned}$$

Combining all of the above, we obtain:

$$\begin{aligned}g'_\lambda(1) &= \alpha_\lambda(1)\beta_\lambda(1)\gamma_\lambda(1) \\ &\quad \times \{H_{2a-1} - H_{2b} - (H_{a-1} - H_b)/2 + H_{2b-1} - H_{n-a-b+1}\} \\ &= \frac{(-1)^{n-a-b+1}(2a-1)!(b-1)!(n-a-b+1)!}{2^{a-b}(a-1)!} \\ &\quad \times \{H_{2a-1} - H_{n-a-b+1} - (H_{a-1} - H_{b-1})/2\}\end{aligned}$$

and we finally deduce the formula in the statement.

2. We now turn to the case where $\lambda = (n+1)$, i.e. $a = n+1$ and $b = 0$.

(a) Following the definition⁸ of $F_\lambda(x)$ given on page 11, we have

$$F_{(n+1)}(x) = 2^{n+1} (x/2 + n)^{n+1} = x \prod_{k=1}^n (x + 2k).$$

We thus obtain

$$F'_{(n+1)}(x) = \prod_{k=1}^n (x + 2k) + F_{(n+1)}(x) \sum_{k=1}^n \frac{1}{x + 2k},$$

which easily gives the wanted expressions when evaluated at $x = 0$ and $x = 1$.

(b) For the second derivative, we obtain

$$F''_{(n+1)}(x) = F_{(n+1)}(x) \sum_{0 \leq i \neq j \leq n} \frac{1}{(x + 2i)(x + 2j)},$$

hence

$$F''_{(n+1)}(1) = \frac{(2n+1)!}{2^n n!} \left\{ \left(\sum_{k=0}^n \frac{1}{2k+1} \right)^2 - \sum_{k=0}^n \frac{1}{(2k+1)^2} \right\},$$

and the formula in the statement follows. □

⁸Again, we replace n with $n+1$ in the definition.

Knowing the generating function G , we can easily obtain the expected value of the number of cycles in the breakpoint graph of a random signed permutation of n elements.

Theorem 8.3. *The expected value of the number of cycles in the breakpoint graph of a uniform random signed permutation π^\pm of n elements is*

$$\mathbb{E}(c(BG(\pi^\pm))) = H_{2n+1} - \frac{H_n}{2} - \sum_{(a,b) \in \mathcal{A}_n} r_n(a,b),$$

where $\mathcal{A}_n = \{(a,b) \in \mathbb{N}^2 : a \geq b \geq 1, a+b \leq n+1\}$ and

$$r_n(a,b) = \frac{(-1)^{n+a-b}(n+1)(2a-2b+1)(a-1)!(2b-2)!(n-a-b+2)!}{2^{n-a+b-1}n!(b-1)!(n+a-b+2)2(n-a+b+1)!}.$$

Proof. As recalled in the proof of Theorem 8.1, we have $\mathbb{E}(c(BG(\pi^\pm))) = G'(1)/G(1)$. Note that, by definition, $G(1) = \sum_{k=1}^{n+1} S_H^\pm(n,k)$, which equals the number of signed permutations of n elements, i.e. $2^n n!$. By Lemma 8.3, the expected number of cycles in the breakpoint graph of a random signed permutation is

$$\mathbb{E}(c(BG(\pi^\pm))) = \frac{1}{2^n n!} \sum_{\lambda} c_{\lambda}(2) F'_{\lambda}(1) F'_{\lambda}(0).$$

Using the formulas for $F'_{\lambda}(1)$ and $F'_{\lambda}(0)$ derived in Lemma 8.4 and the expression for the coefficients⁹ $c_{\lambda}(2)$ given in Equations (10) and (11) page 11, the formula in the statement follows. \square

The generating function G allows us also to compute the variance of the signed Hultman numbers.

Theorem 8.4. *The variance of the number of cycles in the breakpoint graph of a uniform random signed permutation π^\pm of n elements is*

$$\begin{aligned} \text{Var}(c(BG(\pi^\pm))) &= H_{2n+1} - \frac{H_n}{2} - \sum_{k=0}^n \frac{1}{(2k+1)^2} - \left(\sum_{(a,b) \in \mathcal{A}_n} r_n(a,b) \right)^2 \\ &+ \sum_{(a,b) \in \mathcal{A}_n} r_n(a,b) \{2H_{2n+1} - H_n - 2H_{2a-1} + 2H_{n-a-b+1} + H_{a-1} - H_{b-1} - 1\}, \end{aligned}$$

where \mathcal{A}_n and the coefficients $r_n(a,b)$ are as defined in Theorem 8.3.

Proof. As recalled in the proof of Theorem 8.2, the variance can be obtained from the generating function G by evaluating the function $(\log G)'(x) + (\log G)''(x)$

⁹Again, we replace n with $n+1$ in the definitions.

at $x = 1$. Therefore, the variance of the number of cycles in the breakpoint graph of a random signed permutation equals

$$\begin{aligned} & \frac{G'(1)}{G(1)} + \frac{G''(1)}{G(1)} - \left(\frac{G'(1)}{G(1)} \right)^2 \\ = & \frac{G'(1) + G''(1)}{G(1)} - (\mathbb{E}(c(BG(\pi^\pm))))^2 \\ = & \frac{1}{2^n n!} \sum_{\lambda} c_{\lambda}(2)(F'_{\lambda}(1) + F''_{\lambda}(1))F'_{\lambda}(0) - (\mathbb{E}(c(BG(\pi^\pm))))^2. \quad (\text{using Lemma 8.3}) \end{aligned}$$

Using the formulas for $F'_{\lambda}(1)$, $F''_{\lambda}(1)$ and $F'_{\lambda}(0)$ given in Lemma 8.4, we obtain that the variance equals

$$\begin{aligned} & H_{2n+1} - \frac{H_n}{2} - \sum_{k=0}^n \frac{1}{(2k+1)^2} + \left(H_{2n+1} - \frac{H_n}{2} \right)^2 - (\mathbb{E}(c(BG(\pi^\pm))))^2 \\ & - \sum_{(a,b) \in \mathcal{A}_n} r_n(a,b) \{2H_{2a-1} - 2H_{n-a-b+1} - H_{a-1} + H_{b-1} + 1\}, \end{aligned}$$

which equals the wanted expression once $\mathbb{E}(c(BG(\pi^\pm)))$ is replaced with the value derived in Theorem 8.3. \square

As in the unsigned case, we will study the behaviour of the mean and variance for large values of n . To that end, we will first prove the following lemma.

Lemma 8.5. *As $n \rightarrow \infty$, we have*

$$\sum_{(a,b) \in \mathcal{A}_n} |r_n(a,b)| = \frac{1}{\log(n)} \times o(1).$$

Proof. If we denote $k = a - b$, the above sum becomes

$$\begin{aligned} & \sum_{k=0}^{n-1} \frac{2^{k-n+1}(n+1)(2k+1)}{n!(n+k+2)2^{n-k+1}} \sum_{b=1}^{\lfloor (n-k+1)/2 \rfloor} \frac{(k+b-1)!(2b-2)!(n-k-2b+2)!}{(b-1)!} \\ = & \sum_{k=0}^{n-1} \frac{2^{k-n+1}(n+1)(2k+1)}{(n+k+2)2^{n-k+1}(k+1)\binom{n}{k+1}} \sum_{b=1}^{\lfloor (n-k+1)/2 \rfloor} \frac{\binom{k+b-1}{k}}{\binom{n-k}{2b-2}} \\ \leq & \sum_{k=0}^{n-1} \frac{2^{k-n+1}}{(n+k+2)\binom{n}{k+1}} \sum_{b=1}^{\lfloor (n-k+1)/2 \rfloor} \binom{k+b-1}{k} \\ = & \sum_{k=0}^{n-1} \frac{2^{k-n+1} \binom{k+\lfloor (n-k+1)/2 \rfloor}{k+1}}{(n+k+2)\binom{n}{k+1}}. \quad (\text{using } \sum_{j=k}^n \binom{j}{k} = \binom{n+1}{k+1}) \end{aligned}$$

We further observe that

$$\sum_{(a,b) \in \mathcal{A}_n} |r_n(a,b)| \leq \sum_{k=0}^{n-1} \frac{2^{k-n+1}}{n+k+2} \leq 2 \left(1 - \frac{1}{2^n} \right) \frac{1}{n+2},$$

and the result in the statement easily follows. \square

Based on this lemma, we can now obtain the following.

Remark 8.2. *When $n \rightarrow \infty$, the expected value and variance of the number of cycles in the breakpoint graph of a uniform random signed permutation π^\pm of n elements have the following asymptotical behaviour:*

$$\begin{aligned}\mathbb{E}(c(BG(\pi^\pm))) &= \frac{\log(n)}{2} + \frac{\gamma}{2} + \log(2) + o(1), \\ \text{Var}(c(BG(\pi^\pm))) &= \frac{\log(n)}{2} + \frac{\gamma}{2} + \log(2) - \frac{\pi^2}{8} + o(1).\end{aligned}$$

Note that, in the limit when $n \rightarrow \infty$, the mean and variance in the signed case are of the same order ($\log(n)$) as in the unsigned case, but they differ by a factor of $1/2$.

9. Applications: Distributions of rearrangement distances

As stated in the introduction of this paper, the breakpoint graph and its cycles are used in a lot of variants of genome rearrangement problems to compute evolutionary distances – either exactly or approximately. In this section, we are interested in exploring to what extent we can rely on those cycles in order to approximate the distribution of several distances that have been studied in the field of genome rearrangements, so as to obtain a better idea of how tight a particular bound on a distance is, or whether it is worth computing a distance exactly in cases where this requires solving an NP-hard problem. By “distribution of a distance”, we mean the number of (possibly signed) permutations of n elements whose distance equals k , for all possible values of k .

We will not say much about rearrangement distances or how to compute them, except for the fact that, as already stated earlier in this paper, they are based on a set S of operations that generate S_n (resp. S_n^\pm). In the following, what we mean by expressions like “the S distance of π ” is the minimum number of operations from S needed to transform a given permutation π into the identity permutation ι ; a few examples of such operations that we will consider here are summarised informally in Table 2. The reader should bear in mind that

	Distance	Operation	Description of the operation
unsigned	<i>bid</i>	block-interchange	exchanges two non-necessarily adjacent segments
	<i>td</i>	transposition	exchanges two adjacent segments
	<i>ptd</i>	prefix transposition	transposition involving $\pi_1, \pi_2, \dots, \pi_k$ for some k
	<i>rd</i>	reversal	reverses a segment
signed	<i>prd</i>	prefix reversal	reversal involving $\pi_1, \pi_2, \dots, \pi_k$ for some k
	<i>srd</i>	signed reversal	reverses a segment and flips the signs in that segment
	<i>psrd</i>	prefix signed reversal	signed reversal involving $\pi_1, \pi_2, \dots, \pi_k$ for some k

Table 2: Some abbreviations and informal definitions used throughout this section.

the discussion presented in this section focuses on experiments with relatively small amounts of data (mainly because many interesting distances are hard to compute, and because the number of (signed) permutations grows much too fast to generate the full distributions for large values of n), which is why we refrain from making any bold conjecture or actually proving any result. We will also restrict ourselves to comparing distributions for one fixed value of n , namely, the largest value for which we could obtain the distribution of the particular distance we are interested in; similar-looking plots can however be obtained for any value. We generated the distributions based on cycles of the breakpoint graph ourselves, but the distributions of the distances we consider here were computed by Galvão and Dias [22].

9.1. Unsigned distances

A few distances between unsigned permutations have been considered in the field of genome rearrangements [3]. Doignon and Labarre [1] already observed that $S_H(n, n + 1 - 2k)$ is exactly the number of permutations π in S_n whose *block-interchange distance* $bid(\pi)$ equals k , an immediate consequence of the following result.

Theorem 9.1. [6] *For all π in S_n , we have $bid(\pi) = (n + 1 - c(BG(\pi)))/2$.*

Whereas sorting by block-interchanges and computing $bid(\pi)$ can be achieved in polynomial time [6], this is not the case for any of the other unsigned operations listed in Table 2: sorting by transpositions and sorting by reversals, as well as computing the related distances, are NP-hard problems (see Bulteau et al. [23] and Caprara [24], respectively); the same problems in the context

of prefix reversals are also NP-hard [25], while their complexity in the case of prefix transpositions is open.

However, since transpositions are but a particular case of block-interchanges, the expression given in Theorem 9.1 for computing $bid(\pi)$ is also a lower bound on the *transposition distance* $td(\pi)$. Additionally, a tighter lower bound on the transposition distance was proved by Bafna and Pevzner [26].

Theorem 9.2. [26] *For all π in S_n , we have $td(\pi) \geq (n + 1 - c_{odd}(BG(\pi)))/2$, where $c_{odd}(BG(\pi))$ is the number of cycles of odd length in $BG(\pi)$.*

Consequently, it makes sense to try to approximate the distribution of the transposition distance using $\mathcal{S}_H(n, n+1-2k)$ (because of Theorem 9.1) and what could be called the *odd Hultman numbers* $S_H^{odd}(n, n+1-2k)$, i.e. the number of permutations of n elements whose breakpoint graph contains $n+1-2k$ cycles of odd length (because of Theorem 9.2). Figure 8(a) compares all three distributions for $n = 13$. To the best of our knowledge, there is no known formula for computing odd Hultman numbers.

Dias and Meidanis [27] initiated the study of *prefix transpositions*, which are transpositions that can only be applied to an initial segment of the permutation to sort. To the best of our knowledge, the complexity of sorting by prefix transpositions or computing the corresponding distance is still open. However, a lower bound on the prefix transposition distance based on the breakpoint graph is known.

Theorem 9.3. [28] *For any π in S_n , we have*

$$ptd(\pi) \geq \frac{n + 1 + c(BG(\pi))}{2} - c_1(BG(\pi)) - \begin{cases} 0 & \text{if } \pi_1 = 1, \\ 1 & \text{otherwise,} \end{cases} \quad (19)$$

where $c_1(BG(\pi))$ is the number of cycles of length 1 in $BG(\pi)$.

Figure 8(b) shows the distribution of the prefix transposition distance, together with some function of the Hultman numbers and the distribution of the number of permutations in S_n for which lower bound (19) equals k for $n = 13$. On this particular plot and the forthcoming ones, we find the offset m in $\mathcal{S}_H(n, n+1-k+m)$ experimentally by shifting the distribution of

$\mathcal{S}_H(n, n + 1 - k)$ so that it best fits the distribution of the distance we are interested in.

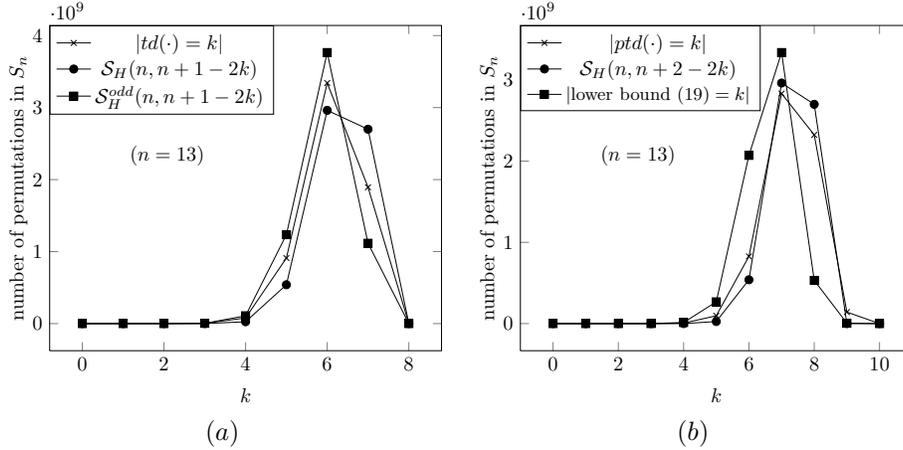


Figure 8: (a) How the distributions of the unsigned and odd Hultman numbers relate to the distribution of the transposition distance, for $n = 13$; (b) how the distributions of the unsigned Hultman numbers and the number of permutations for which lower bound (19) equals k relate to the distribution of the prefix transposition distance, for $n = 13$.

Two other distances that have received a considerable amount of attention are the *reversal distance*, where a reversal reverses the order of the elements contained in the segment of the permutation on which it acts, and the *prefix reversal distance*, where prefix reversals have the same effect as reversals but may only be applied to an initial segment of the permutation. Caprara [24] showed that computing the former is NP-hard, while Bultheau et al. [25] proved that computing the latter is NP-hard. Again, we find it interesting to examine how the distribution of the number of cycles in the breakpoint graph relates to those distances, which we do in Figure 9. We warn the reader familiar with breakpoint graphs, however, that the breakpoint graph used in our paper differs from the structure traditionally used for the study of these two distances, which admits more than one cycle decomposition; the graph we use can be seen as the result of selecting one particular decomposition among all possible decompositions. In this setting, there is a much larger difference between the distributions of both distances and of the unsigned Hultman numbers than what we have observed

for transpositions in Figure 8, which confirms that using only (our version of) the breakpoint graph in this case is not enough.

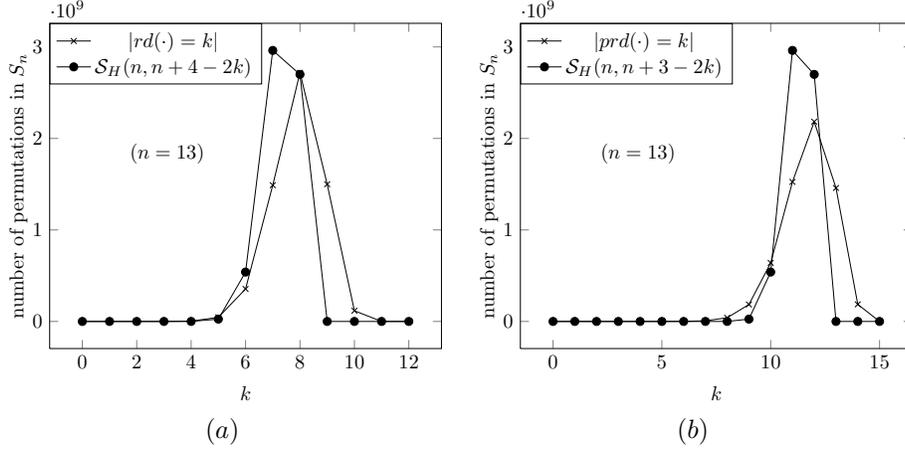


Figure 9: How the distribution of the unsigned Hultman numbers relates to the distribution of (a) the reversal distance and (b) the prefix reversal distance, for $n = 13$.

9.2. Signed distances

A number of well-studied and biologically relevant distances between signed permutations are also based on the breakpoint graph. These include the *double cut-and-join (DCJ) distance*, introduced by Yancopoulos et al. [29], who showed that its value could be computed using the formula $dcj(\pi) = n + 1 - c(BG(\pi))$. As a consequence, the number of signed permutations of n elements with DCJ distance k is exactly $\mathcal{S}_H^\pm(n, n + 1 - k)$.

Another distance whose distribution can be well approximated using the signed Hultman numbers is the *signed reversal distance* (see Table 2 for an informal definition of signed reversals). Hannenhalli and Pevzner [30] proved the following formula for computing the signed reversal distance of any permutation π , denoted by $srd(\pi)$.

Theorem 9.4. [30] *For any π in S_n^\pm , the signed reversal distance of π is*

$$srd(\pi) = n + 1 - c(BG(\pi)) + h(\pi) + f(\pi),$$

where $h(\pi)$ is the number of “hurdles” of π and $f(\pi) = 1$ if π is a “fortress”, and 0 otherwise.

We will not give more details on the terms “hurdles” and “fortress” (see Hannenhalli and Pevzner [30] for definitions), except for the fact that hurdles are particular collections of cycles in $BG(\pi)$, and that a permutation cannot be a fortress unless $h(\pi) > 0$. Our point here is that the following lower bound, first proved by Bafna and Pevzner [4], is extremely tight:

$$\forall \pi \in S_n^\pm : srd(\pi) \geq n + 1 - c(BG(\pi)). \quad (20)$$

This claim is supported by Caprara’s proof [31] of the fact that the probability that a permutation $\pi \in S_n^\pm$ is *not* tight with respect to Equation (20) is $\Theta(n^{-2})$, and by Swenson et al.’s proof [32] that the probability that π is a fortress is $\Theta(n^{-15})$. Therefore, Equation (20) provides a very good approximation of the signed reversal distance, and the distribution of $\mathcal{S}_H^\pm(n, n + 1 - k)$ closely matches that of the signed reversal distance. Figure 10 illustrates the situation for the case $n = 10$.

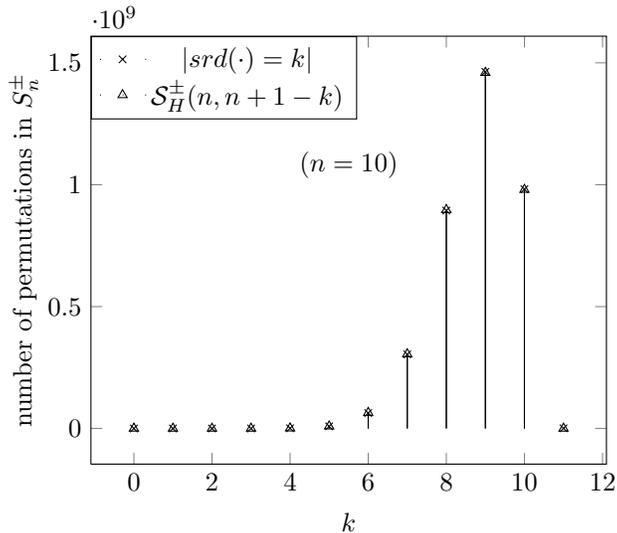


Figure 10: The distributions of the signed reversal distance and of the signed Hultman numbers, for $n = 10$.

Other distances have not been studied with that level of detail, which is why we find it interesting to try to relate their distribution to that of the Hultman numbers. A particular restriction of the signed reversal distance is the *prefix*

signed reversal distance, denoted by $psrd(\cdot)$, whose definition follows that of the signed reversal distance except that reversals can only act on an initial segment of the permutation. No formula is known for computing that distance, and the computational complexity of the problem has remained open since the first works on the subject [33]. However, a lower bound based on the breakpoint graph was recently obtained by Labarre and Cibulka [34], which naturally prompts us to wonder how exactly we can rely on the breakpoint graph to approximate that distance.

Theorem 9.5. [34] *For any π in S_n^\pm , we have*

$$psrd(\pi) \geq n + 1 + c(BG(\pi)) - 2c_1(BG(\pi)) - \begin{cases} 0 & \text{if } \pi_1 = 1, \\ 2 & \text{otherwise.} \end{cases} \quad (21)$$

Figure 11 shows a plot with the distribution of the prefix signed reversal distance and that of the signed Hultman numbers, as well as of the distribution of lower bound (21) for $n = 10$. It can be seen on that graph that the latter is quite far off from the distribution of the prefix signed reversal distance, hinting that additional work seems needed to reduce the gap between the lower bound and the actual distance.

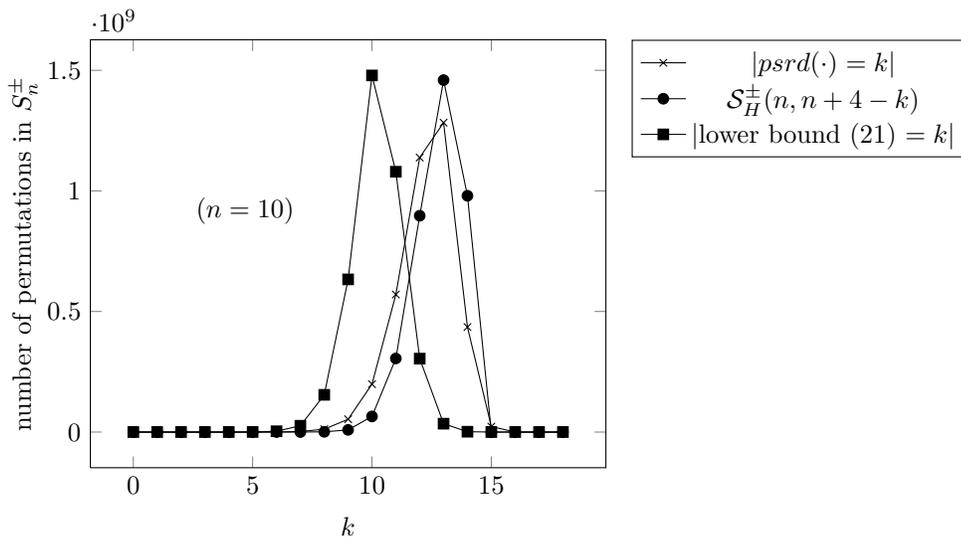


Figure 11: The distributions of the prefix signed reversal distance, of the signed Hultman numbers, and of the number of permutations for which lower bound (21) equals k , for $n = 10$.

10. Conclusions

In this paper, we proved the first explicit formula for enumerating signed permutations whose breakpoint graph contains a given number of cycles, and proved simpler expressions for particular cases. We also obtained a new expression for enumerating unsigned permutations whose breakpoint graph contains a given number of cycles, and used both formulas to derive simpler proofs of some other previously known results. Getting more insight into breakpoint graphs and their cycle decomposition is particularly relevant to edit distances used in the field of genome rearrangements, and we hope that our results can help shed light on their distributions, expected values and variances. There are several interesting directions in which our work could be extended, which we outline and motivate below.

Just like one can define conjugacy classes in the symmetric and hyperoctahedral groups, we could investigate conjugacy classes with respect to the breakpoint graph. This was already initiated by Doignon and Labarre [1], who referred to them as “Hultman classes” and provided explicit formulas for enumerating those classes in the case of unsigned permutations. More work remains to be done in the unsigned case: indeed, the work done by Bóna and Flynn [5] provides us with a very nice formula for computing the distribution of cycles, but no simpler expression than the complicated ones obtained by Doignon and Labarre [1] is yet known for enumerating Hultman classes or their cardinalities. Moreover, no work so far has been done in order to enumerate Hultman classes in the signed setting, and obtaining an expression for enumerating the so-called “simple permutations”, which are defined in this context as permutations whose breakpoint graph contains no cycle of length greater than 2, seems especially interesting (for more information about the importance of those permutations in genome rearrangements, see Hannenhalli and Pevzner [30] and Labarre and Cibulka [34]).

The expression we obtained for the signed Hultman numbers is quite useful in practice, since it allows us to obtain the distribution of those numbers for large

values of n . Unfortunately, it does not seem easy to use in order to gain insights and have an intuitive interpretation of the shape of the distribution, which would be useful in order to know how this distribution can be approximated or how it grows as n increases. Finding simpler generating functions, recurrence relations or nicer formulas would be useful in that regard and in order to obtain more information on the properties of this distribution.

The connection between the cycle structure of breakpoint graphs and factorisations of even permutations (Corollary 4.1, page 9) proved useful not only in characterising the distribution of those cycles and of the related cycle types, but also provided the foundations of a simple and generic method for obtaining lower bounds on *any* “reversible” edit distance between unsigned permutations (see Labarre [28] for more details). Is there any way to use the results and connections obtained in Section 5 in order to obtain similar results for signed permutations?

Finally, recall that permutations are just one way of modelling genomes. One natural direction would be to investigate the distribution of cycles in the breakpoint graph of other structures, like set systems or “fragmented” permutations (see again Fertin et al. [3] for an overview of existing models).

11. Acknowledgements

The first author was partially supported by the ANR MAEV under contract ANR-06-BLAN-0113. Both authors also wish to thank the group “Evolution Biologique et Modélisation”, LATP, Université de Provence, where part of this research was performed, as well as Mathilde Bouvel for bringing reference [14] to their attention.

References

- [1] J.-P. Doignon, A. Labarre, On Hultman Numbers, Journal of Integer Sequences 10 (6), article 07.6.2, 13 pages.

- [2] H. Li, N. Homer, A Survey of Sequence Alignment Algorithms for Next-generation Sequencing, *Briefings in Bioinformatics* 11 (5) (2010) 473–483, URL <http://bib.oxfordjournals.org/content/11/5/473.abstract>.
- [3] G. Fertin, A. Labarre, I. Rusu, E. Tannier, S. Vialette, *Combinatorics of Genome Rearrangements*, Computational Molecular Biology, The MIT Press, 2009.
- [4] V. Bafna, P. A. Pevzner, Genome Rearrangements and Sorting by Reversals, *SIAM Journal on Computing* 25 (2) (1996) 272–289, ISSN 0097-5397.
- [5] M. Bóna, R. Flynn, The Average Number of Block Interchanges Needed to Sort A Permutation and a Recent Result of Stanley, *Information Processing Letters* 109 (16) (2009) 927–931.
- [6] D. A. Christie, Sorting Permutations by Block-interchanges, *Information Processing Letters* 60 (4) (1996) 165–169, ISSN 0020-0190.
- [7] L. Székely, Y. Yang, On the Expectation and Variance of the Reversal Distance, *Acta Universitatis Sapientiae, Mathematica* 1 (1) (2009) 5–20.
- [8] S. Grusea, On the Distribution of the Number of Cycles in the Breakpoint Graph of a Random Signed Permutation, *IEEE/ACM Transactions on Computational Biology and Bioinformatics* 8 (5) (2010) 1411–1416, ISSN 1545-5963.
- [9] R. Diestel, *Graph Theory*, vol. 173 of *Graduate Texts in Mathematics*, Springer-Verlag, Berlin, 3rd edn., ISBN 978-3-540-26182-7; 3-540-26182-6, 2005.
- [10] A. Björner, F. Brenti, *Combinatorics of Coxeter Groups*, vol. 231 of *Graduate Texts in Mathematics*, chap. 8: Combinatorial Descriptions, Springer-Verlag, 2005.
- [11] H. Wielandt, *Finite Permutation Groups*, Translated from German by R. Bercov, Academic Press, New York, 1964.

- [12] R. L. Graham, D. E. Knuth, O. Patashnik, Concrete Mathematics: A Foundation for Computer Science, Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA, 2nd edn., ISBN 0201558025, 1994.
- [13] A. Hultman, Toric Permutations, Master's thesis, Department of Mathematics, KTH, Stockholm, Sweden, 1999.
- [14] J. H. Kwak, J. Lee, Genus polynomials of dipoles, Kyungpook Mathematical Journal 33 (1) (1993) 115–125.
- [15] P. J. Hanlon, R. P. Stanley, J. R. Stembridge, Some Combinatorial Aspects of the Spectra of Normally Distributed Random Matrices, Contemporary Mathematics 138 (1992) 151–174.
- [16] N. R. Goodman, Statistical Analysis Based on a Certain Multivariate Complex Gaussian Distribution (an Introduction), The Annals of Mathematical Statistics 34 (1) (1963) 152–177.
- [17] I. G. Macdonald, Symmetric Functions and Hall Polynomials, Oxford Mathematical Monographs, Oxford University Press, 2nd edn., 1998.
- [18] I. Elias, T. Hartman, A 1.375-Approximation Algorithm for Sorting by Transpositions, IEEE/ACM Transactions on Computational Biology and Bioinformatics 3 (4) (2006) 369–379, ISSN 1545-5963.
- [19] N. J. A. Sloane, The On-Line Encyclopedia of Integer Sequences, Published electronically at <http://oeis.org/>, 2012.
- [20] B. Sury, T. Wang, F.-Z. Zhao, Identities Involving Reciprocals of Binomial Coefficients, Journal of Integer Sequences 7, article 04.2.8, 12 pages.
- [21] H. S. Wilf, generatingfunctionology, A. K. Peters, Ltd., Natick, MA, USA, 3rd edn., ISBN 1568812795, 2006.
- [22] G. R. Galvão, Z. Dias, Rearrangement distance database, <http://mirza.ic.unicamp.br:8080/bioinfo/index.jsf>, 2011.

- [23] L. Bulteau, G. Fertin, I. Rusu, Sorting by Transpositions Is Difficult, in: L. Aceto, M. Henzinger, J. Sgall (Eds.), Proceedings of the Thirty-Eighth International Colloquium on Automata, Languages and Programming (ICALP), Part 1, vol. 6755 of *Lecture Notes in Computer Science*, Springer, ISBN 978-3-642-22005-0, 654–665, 2011.
- [24] A. Caprara, Sorting Permutations by Reversals and Eulerian Cycle Decompositions, *SIAM Journal on Discrete Mathematics* 12 (1) (1999) 91–110 (electronic), ISSN 1095-7146.
- [25] L. Bulteau, G. Fertin, I. Rusu, Pancake Flipping is Hard, in: Proceedings of the Thirty-Seventh International Symposium on Mathematical Foundations of Computer Science (MFCS), vol. 7464 of *Lecture Notes in Computer Science*, Springer-Verlag, Bratislava, Slovakia, to appear, 2012.
- [26] V. Bafna, P. A. Pevzner, Sorting by Transpositions, *SIAM Journal on Discrete Mathematics* 11 (2) (1998) 224–240 (electronic), ISSN 1095-7146.
- [27] Z. Dias, J. Meidanis, Sorting by Prefix Transpositions, in: A. H. F. Laender, A. L. Oliveira (Eds.), Proceedings of the Ninth International Symposium on String Processing and Information Retrieval (SPIRE), vol. 2476 of *Lecture Notes in Computer Science*, Springer, ISBN 3-540-44158-1, 65–76, 2002.
- [28] A. Labarre, Edit Distances and Factorisations of Even Permutations, in: D. Halperin, K. Mehlhorn (Eds.), Proceedings of the Sixteenth Annual European Symposium on Algorithms (ESA), vol. 5193 of *Lecture Notes in Computer Science*, Springer-Verlag, ISBN 978-3-540-87743-1, 635–646, 2008.
- [29] S. Yancopoulos, O. Attie, R. Friedberg, Efficient Sorting of Genomic Permutations by Translocation, Inversion and Block Interchange, *Bioinformatics* 21 (16) (2005) 3340–3346.
- [30] S. Hannenhalli, P. A. Pevzner, Transforming Cabbage into Turnip: Poly-

nomial Algorithm for Sorting Signed Permutations by Reversals, *Journal of the ACM* 46 (1) (1999) 1–27.

- [31] A. Caprara, On the Tightness of the Alternating-cycle Lower Bound for Sorting by Reversals, *Journal of Combinatorial Optimization* 3 (2-3) (1999) 149–182.
- [32] K. M. Swenson, Y. Lin, V. Rajan, B. M. Moret, Hurdles Hardly Have to Be Heeded, in: C. Nelson, S. Vialette (Eds.), *Proceedings of the Sixth International Workshop on Comparative Genomics (RECOMB-CG)*, vol. 5267 of *Lecture Notes in Bioinformatics*, Springer-Verlag, Berlin, Heidelberg, ISBN 978-3-540-87988-6, 241–251, 2008.
- [33] D. S. Cohen, M. Blum, On the Problem of Sorting Burnt Pancakes, *Discrete Applied Mathematics* 61 (1995) 105–120.
- [34] A. Labarre, J. Cibulka, Polynomial-time Sortable Stacks of Burnt Pancakes, *Theoretical Computer Science* 412 (8-10) (2011) 695–702, ISSN 0304-3975.