

A. Tiskin

**THE CHVÁTAL–SANKOFF PROBLEM:
UNDERSTANDING RANDOM STRING COMPARISON
THROUGH STOCHASTIC PROCESSES**

ABSTRACT. Given two equally long, uniformly random binary strings, the expected length of their longest common subsequence (LCS) is asymptotically proportional to the strings' length. Finding the proportionality coefficient γ , i.e. the limit of the normalised LCS length for two random binary strings of length $n \rightarrow \infty$, is a very natural problem, first posed by Chvátal and Sankoff in 1975, and as yet unresolved. This problem has relevance to diverse fields ranging from combinatorics and algorithm analysis to coding theory and computational biology. Using methods of statistical mechanics, as well as some existing results on the combinatorial structure of LCS, we link constant γ to the parameters of a certain stochastic particle process. These parameters are determined by a specific (large) system of polynomial equations with integer coefficients, which implies that γ is an algebraic number. Short of finding an exact closed-form solution for such a polynomial system, which appears to be unlikely, our approach essentially resolves the Chvátal–Sankoff problem, albeit in a somewhat unexpected way with a rather negative flavour.

§1. INTRODUCTION

The *longest common subsequence (LCS)* for a pair of strings a, b is the longest string that is a (not necessarily consecutive) subsequence of both a and b . Given a pair of strings as input, the *LCS problem* asks for the length of their LCS (finding the actual characters of the LCS is not required). The LCS problem is a fundamental problem for both theoretical and applied computer science, and for computational molecular biology; it is also a popular programming exercise.

Key words and phrases: random strings, longest common subsequence, the Chvátal–Sankoff problem, particle processes.

This work was supported by the Russian Science Foundation under grant No. 22-21-00669, <https://www.rscf.ru/en/project/22-21-00669>.

This paper is concerned with the combinatorics of the LCS problem. Let strings a, b be of length n , uniformly random over the binary alphabet. Chvátal and Sankoff [15] (see also [48, Chapter 1]) have shown that the expected LCS length of a, b is asymptotically proportional to n . The *Chvátal–Sankoff problem* asks for the proportionality coefficient γ , i.e. the limit of the normalised expected LCS length $\frac{\mathbb{E}L_n}{n}$ as $n \rightarrow \infty$, where the random variable L_n is defined as the LCS length for strings of length n . Alexander [2] has shown that $0 \leq \gamma - \frac{\mathbb{E}L_n}{n} \leq O\left(\left(\frac{\log n}{n}\right)^{1/2}\right)$.

The Chvátal–Sankoff problem has relevance to diverse fields ranging from combinatorics and algorithm analysis to coding theory (see e.g. Bukh et al. [9]) and computational biology (see e.g. Pevzner and Waterman [37]). For such a natural and simply posed problem, it seems to be surprisingly elusive: neither an exact value nor any closed-form expression for γ are known, and the existing lower and upper numerical bounds on γ are wide apart.

Acknowledgements. I thank Gianfranco Bilardi, Chris Cox, Vassily Duzhin, Maria Fedorkina, Sergei Nechaev, Georgiy Shulga, Nikolai Vassiliev, and Anatoly Vershik for fruitful discussions. I thank my colleagues and students at the Department of Mathematics and Computer Science of St. Petersburg University for the stimulating atmosphere.

§2. RELATED WORK

LCS combinatorics. An important combinatorial feature of the LCS problem, also relevant to its computational aspect, is the problem’s close connection with transposition networks and the Hecke monoid (also called the seaweed monoid or the sticky braid monoid). This connection has been explored over decades from different angles and using greatly varying terminology. In the rest of this paper, we will describe this connection in more detail, and will use it as the first step on our path to the Chvátal–Sankoff problem.

While the computational aspect of the LCS problem is outside the scope of this paper, it should be mentioned that the problem’s computational complexity, along with that of the closely related edit distance and sequence alignment problems, has been thoroughly studied and is well-understood. Seminal work on LCS algorithms and lower bounds includes e.g. [57, 34, 1, 7].

Reference	$\gamma >$	$\gamma <$	γ
Chvátal and Sankoff [15]	0.697844	0.866595	≈ 0.8082
Deken [20]	0.7615	0.8575	
Steele [47] (conjecture attr. to Arratia)			$\stackrel{?}{=} 2(\sqrt{2}-1) \approx 0.8284$
Dančák [18]; Paterson and Dančák [36]	0.77391	0.83763	≈ 0.812
Baeza-Yates et al. [4]			≈ 0.8118
Boutet de Monvel [19]			≈ 0.812282
Bundshuh [10]			≈ 0.812653
Lueker [30]	0.788071	0.826280	
Bukh and Cox [8]			≈ 0.8122
this work	γ algebraic; exact polynomial equations		

Table 1. Bounds and estimates on γ

Random LCS on permutation strings. Apart from binary strings, a question analogous to the Chvátal–Sankoff problem can be asked about pairs of uniformly random permutations of the alphabet $\{1, \dots, n\}$. The LCS problem on such permutation strings is equivalent to finding the longest increasing subsequence (LIS) of a single permutation of length n . The LCS (respectively, LIS) length in this case turns out to be asymptotically proportional to \sqrt{n} . The proportionality constant was found to be exactly 2 in the classical works of Vershik and Kerov [56] and Logan and Shepp [29] (see also [40]), as part of a solution for the more general problem asking for the limit shape of a random Young diagram sampled from the Plancherel distribution.

Bounds and estimates for γ . Chvátal and Sankoff [15] gave the first analysis of the problem, and proved the existence of the limit γ . Properties of the convergence of the normalised LCS length to this limit were studied since then by numerous researchers. Table 1 lists some results on specific lower and upper bounds, as well as experimental numerical estimates of γ .

The best currently known analytic bounds on γ are due to Lueker [30]. Despite the ingenious methods of obtaining these bounds and numerous related results, the gap between the upper and the lower bounds remains quite wide: in particular, not a single digit of γ after decimal point is known exactly.

Stochastic evolution models. Due to the combinatorial properties of the LCS problem that will be presented in the next section, the Chvátal–Sankoff problem turns out to be closely related to the theory of stochastic

evolution models, which is a vast and actively developing field of study. Particularly relevant areas within this field include particle processes, random Young diagrams, stochastic cellular automata. Asymptotic properties of such models are studied with the help of partial differential equations (PDEs), which describe a model's evolution at the macroscopic level. In the rest of this paper, we will describe these connections in more detail, and will build upon them to obtain a solution of the Chvátal–Sankoff problem.

§3. COMBINATORICS OF THE LCS PROBLEM

LCS grid. Let strings a, b be of length m, n respectively. The *LCS grid* defined by a, b is a directed graph on an $(m + 1) \times (n + 1)$ grid of nodes; we visualise the nodes as being indexed top-to-bottom and left-to-right. Every pair of horizontally or vertically adjacent nodes are connected by an edge, directed rightwards (respectively, downwards). A pair of diagonally adjacent nodes $(i, j), (i + 1, j + 1)$, $0 \leq i < m, 0 \leq j < n$, are connected by an edge whenever $a_i = b_j$ (the two characters *match*); this edge is directed towards below-right. The LCS grid can also be viewed as an $m \times n$ grid of cells, each formed by a quadruple of adjacent nodes and their four connecting horizontal and vertical edges. The cell is called *match cell*, if the two corresponding characters match (and therefore the cell contains a diagonal edge), otherwise a *mismatch cell*. The LCS problem is equivalent to asking for the length of a path in the LCS grid from the top-left node $(0, 0)$ to the bottom-right node (m, n) , that maximises the number of diagonal edges along the path.

Example 1. Figure 1 (left) shows the LCS grid for a pair of binary strings. The horizontal and vertical edges are shown in light-blue, and the diagonal edges in solid red. The left-to-right, top-to-bottom direction of the edges is left implicit.

Sticky braids. The combinatorial structure of the LCS problem is described algebraically by the *Hecke monoid* (also known as the *sticky braid monoid*), which is defined similarly to the classical braid group, but with element inversion replaced by the idempotence relation on the monoid's generators. Given an LCS grid, strands of the corresponding sticky braid are formed by paths in the dual graph, i.e. the plane graph whose nodes are the faces of the LCS grid, and the edges go across the edges of the LCS grid. Multiplication of sticky braids in the Hecke monoid (also known

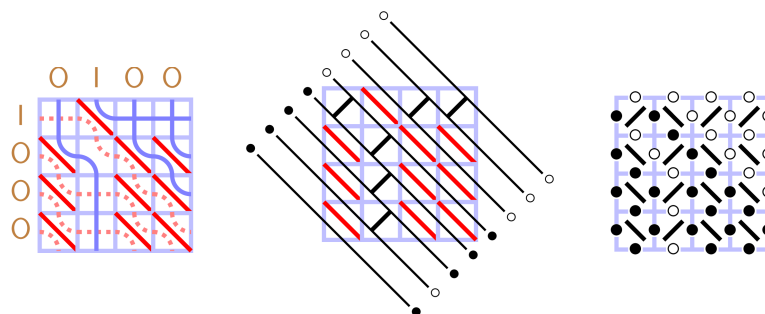


Figure 1. LCS grid with a sticky braid (left), transposition network (centre), particle evolution model (right) for strings $a = \text{"lOOO"}$, $b = \text{"OIOO"}$.

as *Demazure multiplication*) describes precisely how LCS lengths of input strings and their substrings behave under string concatenation.

Example 2. Figure 1 (left) shows a sticky braid embedded into the LCS grid of the previous example. The braid's strands are shown in darker blue and dotted red.

The connection outlined above between the LCS problem and the Hecke monoid has been rediscovered many times in different forms. In particular, it underlies implicitly the algorithms for various string comparison problems by Schmidt [46], Crochemore et al. [16, 17], Alves et al. [3], Hyvrö [24], and was made explicit by Tiskin [49, 50, 53]. More recently, new algorithmic applications of this connection were found by Sakai [42, 43], Tiskin [51, 52, 54, 55], Gawrychowski et al. [22], Hermelin et al. [23], Matarazzo et al. [35], Charalampopoulos et al. [14, 13].

Transposition networks. Another convenient tool for exposing the combinatorial structure of the LCS problem comes in the form of *transposition networks*. These are a special case of comparison networks, which are a classical type of computational circuits studied by Batchier [5], Knuth [26] and many others. In a comparison network, input values travel on an array of parallel wires; any prescribed pair of values can be sorted by a *comparator* connecting their respective wires. In a transposition network, an additional restriction is imposed that only adjacent pairs of wires can be connected by a comparator.

Given a pair of strings a, b of lengths m, n respectively, their LCS grid can be overlaid by a transposition network on $m + n$ wires, extending diagonally from above-left to below-right and passing through the midpoints of the grid's edges. These intersection points of the network's wires and the grid's edges will be called *sites*; we will distinguish horizontal and vertical edge sites. A wire passes through an alternating sequence of horizontal and vertical edge sites; the *value* of a given site is the value carried through it by the wire. A cell is crossed by two wires: one connecting its left and bottom boundary edges, the other its top and right boundary edges. The two sites at the cell's left and top boundary edges are its *entry sites*, and the two sites at its right and bottom boundary edges are its *exit sites*. The network's comparators are specified as follows: a mismatch cell always contains a comparator between the two wires that cross it, while a match cell never contains a comparator. A cell can therefore be of one of two *types*: “match” (denoted ‘ \backslash ’), containing a diagonal grid edge, and “mismatch” (denoted ‘ \nearrow ’), containing a network's comparator; the notation indicates the direction of the diagonal edge and of the comparator, respectively. Occasionally, we identify cell type ‘ \backslash ’ with value zero, and cell type ‘ \nearrow ’ with value one.

Example 3. Figure 1 (centre) shows the LCS grid of the previous example, overlaid with its respective transposition network.

Given an input of $m + n$ distinct values sorted in reverse order, the set of values' trajectories through such a transposition network forms a sticky braid corresponding to the comparison of strings a, b ; each particular value traces a strand in this braid. The network's output permutation provides detailed information about LCS lengths between various substrings of a, b . For our purposes, the above construction can be simplified as follows: instead of all distinct values, let the transposition network's input consist of m ones, followed by n zeros; note that such an input array is still sorted in reverse. In this context, value zero will be called a *hole* (denoted ‘ \circ ’), and value one a *particle* (denoted ‘ \bullet ’). This is done not only to distinguish the (binary) values in the network from (also binary) string characters and (again binary) cell types, but also to reflect in our terminology the important connection with particle interaction models, that we will develop further in the remainder of this paper.

An assignment of values/types to a subset of sites/cells of a transposition network will be called a *configuration*. In particular, the input configuration formed by m particles entering the LCS grid at its left boundary,

and n holes entering at the top boundary, will be called the *step initial condition*.

Example 4. The transposition network in Figure 1 (centre) is shown with the step initial condition input sequence at the top-left, and the corresponding output sequence of particles and holes at the bottom-right.

The LCS length of strings a, b is particularly easy to obtain from the transposition network with step initial condition: it is equal to the number of particles among the network's n outputs exiting the grid at the bottom (equivalently, the number of holes among its m outputs exiting the grid at the right). This observation underlies implicitly the bit-parallel LCS algorithms of Crochemore et al. [16] and Hyyrö [24], and was made explicitly e.g. by Majumdar and Nechaev [32] and by Krusche and Tiskin [28]. Let a, b be of equal length $m = n$; in this case, the LCS grid has the shape of a square, and the LCS length is equal to the number of particles (equivalently, the number of holes) that have never crossed the grid's main diagonal.

Example 5. In the previous example, there are three particles among the $n = 4$ outputs at the grid's bottom; the LCS length for strings a, b is also 3. In the course of the evolution of the transposition network, $4 - 3 = 1$ particle has crossed the main diagonal from left to right; accordingly, one hole has done so from top to bottom.

§4. MODEL CS

The combinatorial properties of the LCS problem allow us to reformulate the Chvátal–Sankoff problem in the language of stochastic particle interaction models. By a *network evolution model*, we will understand the evolution of site values from a given input configuration in an infinitely wide transposition network, under a certain probabilistic rule that determines the type of each of the network's cells.

Cell dependencies. Let a, b now be infinite strings, where all characters are independent uniform binary random variables. We define *model CS* (the Chvátal–Sankoff model) as a network evolution model where cell types are determined by character matches and mismatches between strings a, b , as described in the previous section.

Proposition 1. *In model CS , the types of any three distinct cells are mutually independent. The types of any three distinct cells within a \square -shape determine uniquely the type of the fourth cell.*

Proof. The first statement is straightforward by the independence and uniformity of character distribution in strings a, b . The second statement is also straightforward, since the sum of the four cells' types must be even. \square

In particular, the types of any three cells adjacent in a \square -shape are mutually independent; we shall call this property \square -independence. Note that \square -independence relies crucially the uniform distribution of string characters, and would not hold for a non-uniform character distribution, even if it were independent and identical.

Evolution. Let strings a, b be indexed by i, j respectively. The state of model CS can be thought of as evolving in several different ways — vertically, horizontally or diagonally, with the discrete time dimension indexed by i, j and $\frac{i+j}{2}$, respectively. We will focus mainly on the diagonal evolution, due to its symmetry and locality properties. The model's state under such evolution corresponds to an anti-diagonal doubly-infinite sequence of particle-hole values, alternating between horizontal and vertical edge sites. Let us index the transposition network's wires entering the grid through its top boundary with nonnegative integers $0, 1, 2, \dots$, and the wires entering the grid through its left boundary with negative integers $-1, -2, -3, \dots$; the count in both cases starts from the top-left cell. A time step under diagonal evolution then consists of two half-steps: the first involves comparators operating on pairs of adjacent sites with an odd and an even index (in that order), the second on pairs with an even and an odd index (in that order).

As discussed in the previous section, the behaviour of model CS reflects the LCS combinatorics of its underlying string pair a, b .

Proposition 2. *Let $0 \leq k \leq 2n$. Consider the prefixes of infinite strings a, b of length $k, 2n - k$ respectively, and let l be the LCS length of these prefixes. Under diagonal evolution of model CS from step initial condition after n time steps, there are $k - l$ particles at sites with indices $2n - 2k$ or greater.*

Proof. Well-known from the combinatorial properties of LCS; see e.g. [32, 28]. \square

Example 6. Figure 1 (right) shows the evolution of model CS from step initial condition on strings a, b of the previous examples. Wires with negative (respectively, nonnegative) indices are those below (respectively, above) the network’s main diagonal. Let $n = k = 4$. The LCS length of the input strings, regarded as prefixes of length $k = 2n - k = 4$ of a pair of infinite strings, is $l = 3$; as before, we note that after $n = 4$ time steps, exactly $n - l = 4 - 3 = 1$ particle has crossed over the main diagonal to wires with nonnegative indices.

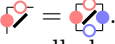

Duality. The definition of model CS is symmetric with respect to the reflection of the network about its main diagonal. A pair of configurations will be called *dual*, if one of them is obtained from the other by a reflection about an above-left to below-right axis (exchanging the directions towards below-left and above-right), with simultaneous exchange of sites’ values between particles and holes. In particular, the step initial condition is a self-dual configuration.

In the remainder of this paper, we will consider model CS with step initial condition. Our analysis will concentrate on the model’s behaviour in a small neighbourhood of the main diagonal, where the particle and hole densities should be asymptotically equal by symmetry. Duality will help to simplify the exposition, since in such a setting, a pair of dual configurations will have equal probabilities.

§5. SPECIAL NOTATION

Configuration probabilities. We consider configurations of a network evolution model as random events. The probability of an event will be denoted by its graphical representation. Thus, $\blacklozenge = 1 - \phi$ represents the probability of a given vertical edge site holding a particle, as opposed to a hole, and $\swarrow = 1 - \searrow$ represents the probability of a given cell being of type “mismatch”, as opposed to “match”.

We extend this notation to represent conditional probabilities as follows. We juxtapose the conditioning event and the conditioned event in the same picture; the elements of the conditioning event will be highlighted in red, while the elements of the conditioned event will be shown in the ordinary black. For example, the probability of a given cell being of type “mismatch”, conditioned on the cell’s left (respectively top) entry value being a particle (respectively, a hole), will be denoted by $\swarrow_{\text{red}}^{\text{red}} = \swarrow_{\text{red}}^{\text{red}} / \swarrow_{\text{black}}^{\text{black}}$.

Some events that we consider may be forced by other events: a forced event, conditioned on the forcing event, occurs with certainty. We juxtapose the forcing event and the forced event in the same picture; the elements of the forced event will be highlighted in blue, while the elements of the forcing event will be shown in either black or red, as appropriate. In the previous example, the cell's exit values are forced: . Showing forced sub-events is a notational decoration that can formally be omitted; however, it is meant to serve as an intuition aid, especially so when some non-forced sub-event becomes forced in a chain of equalities. For example, we have .

Annotated equalities. Standard annotated equality $A \stackrel{\text{def}}{=} B$ (“ A is defined as B ”) will be used to introduce new notation. Additionally, we will use some other annotations on the equality sign, as an aid to the reader. Notation $A \stackrel{r}{=} B$ (“ A and B are obtained from each other by reversal with an exchange of particles and holes”) will indicate that the equality holds by the duality property of network configurations.

Other notation. For brevity, we will denote $\bar{z} = 1 - z$ for any z , $0 \leq z \leq 1$. We will also occasionally use bracketed superscripts to denote $z^{[\bullet]} = z^{[l]} = z$ and $z^{[\circ]} = z^{[l]} = \bar{z}$. We will use subscripts and (unbracketed) superscripts to express various meanings as required; to avoid confusion, we will never use superscripts to indicate powers, not even in polynomials. Strings in the alphabet $\{\circ, \bullet\}$ will sometimes be treated as binary numbers; for brevity, we will convert such numbers to decimal where appropriate. We let $a, b, c, d, e, g \in \{\circ, \bullet\}$, $E, F, G \in \{\backslash, /\}$ for the remainder of this paper.

§6. SCALING LIMITS

Informally, the *scaling limit* of a particle evolution model is the continuous limit of the distribution of particle densities at the model's sites, as both time and space are simultaneously scaled down at appropriate rates, so that the magnitude of both time and space units tends to zero. A general introduction to the theory of scaling limits is given e.g. by Kriecherbauer and Krug [27].

Scalar conservation laws. Partial differential equations (PDEs) are an indispensable tool in studying the asymptotic behaviour of particle evolution models. Using PDEs, one can relate the global behaviour of the model, such as its non-stationary evolution from a given initial condition, with its

local behaviour, such as its stationary state in a small space-time region. A classical example of such a relationship is the asymptotic behaviour of the continuous-time totally asymmetric simple exclusion process (TASEP) with step initial condition, which was shown to be governed by the inviscid Burgers' equation by Rost [41] (see also [27, 40, 21]).

In general, the scaling limit of a conservative particle model with one spatial dimension can be associated with a *scalar conservation law* (see e.g. [27]), which is a PDE of the form

$$\frac{\partial}{\partial t}y + \frac{\partial}{\partial x}f(y) = 0$$

where $y = y(t, x)$ is the *density* function of time t and the spatial dimension x , representing the conserved quantity (typically, the mass of some fluid), and $f = f(y)$ is a strictly concave smooth function of density y called the (rightward) *flux*. We are particularly interested in the *step* initial condition:

$$y(0, x) = \begin{cases} 1 & x < 0 \\ 0 & x > 0 \end{cases}$$

In the language of PDEs, the step initial condition is a special case of the Riemann problem for a scalar conservation law. The discontinuity of y at $x = t = 0$ is known as *shock*. This initial shock dissipates over time in a *rarefaction wave*, governed by the equation's solution (see e.g. [27, 44])

$$y(t, x) = \begin{cases} (f')^{-1}(x/t) & f'(1)t \leq x \leq f'(0)t \\ y(0, x) & \text{otherwise} \end{cases}$$

where f' is the derivative of f , and superscript -1 denotes its functional inverse.

Since the solution scales linearly with t , it is sufficient for the analysis to consider a single time moment $t > 0$; a natural choice is $t = 1$. Let $y(x) = y(1, x)$. We impose further constraints $0 \leq y \leq 1$, $f(0) = f(1) = 0$, which are natural for the interpretation of y as a fluid's density. The maximum flux \tilde{f} is determined by $f'(y) = 0$, and is therefore attained at density $\tilde{y} = (f')^{-1}(0) = y(0)$; we will call these *peak flux* and *peak density*, respectively.

Recall that under the step initial condition, all the fluid's mass is concentrated in the negative half-line at time $t = 0$. The key characteristic of the system is the amount of mass transported across the origin to the positive half-line by the time $t = 1$, which turns out to be precisely the

peak flux:

$$\int_0^{+\infty} y(x) dx = \int_0^{f'(0)} (f')^{-1}(x) dx = \int_0^{\tilde{y}} f'(y) dy = f(\tilde{y}) - f(0) = f(\tilde{y}) = \tilde{f}.$$

We will call the function $1 - f = \bar{f}$ and the value $1 - \tilde{f} = \bar{\tilde{f}}$ respectively the *flux complement* function and the *peak flux complement*. A close relationship between the peak flux complement and the constant γ of the Chvátal–Sankoff problem will be exposed in the rest of this section.

Network model limit. For a network evolution model, density y in the above equations is the limiting marginal probability of a site to contain a particle (as opposed to a hole). The flux for a model X is determined as the (unconditional) probability that a particle and a hole are exchanged by a comparator within the cell. This probability, as well as its complement, have a straightforward expression in terms of marginal site probabilities:

$$f^X \stackrel{\text{def}}{=} \begin{array}{c} \bullet \\ \circ \end{array} \begin{array}{c} \bullet \\ \circ \end{array} = \begin{array}{c} \bullet \\ \bullet \end{array} \begin{array}{c} \bullet \\ \circ \end{array} = \begin{array}{c} \bullet \\ \circ \end{array} \begin{array}{c} \bullet \\ \bullet \end{array} = \begin{array}{c} \bullet \\ \circ \end{array} - \begin{array}{c} \bullet \\ \bullet \end{array} = \begin{array}{c} \bullet \\ \circ \end{array} - \begin{array}{c} \bullet \\ \bullet \end{array} \quad \bar{f}^X \stackrel{\text{def}}{=} 1 - (\begin{array}{c} \bullet \\ \circ \end{array} - \begin{array}{c} \bullet \\ \bullet \end{array}) = \begin{array}{c} \bullet \\ \bullet \end{array} + \begin{array}{c} \bullet \\ \circ \end{array} \quad (1)$$

For a model evolving vertically or horizontally, every cell is accounted for in the above expression for the flux in a given time step. For a model evolving diagonally, one half of the cells is accounted for in the first half-step of a time step, and the other half of the cells in the second half-step.

For a model that has mirror symmetry of cell type probabilities about the main diagonal (such as model *CS* and all the others considered in this paper), and that evolves diagonally from the (skew-symmetric) step initial condition, the site probabilities will be skew-symmetric about the main diagonal: particle probability at a site on one side of the main diagonal must be equal to the hole probability at the symmetrically opposite site. By symmetry, the peak density for such a model in the scaling limit is $\tilde{y} = \frac{1}{2}$, realised in a small neighbourhood of the main diagonal.

From now on, we will consider the model's state in an infinitesimally small neighbourhood of the scaling limit point $t = 1$, $x = 0$ on the main diagonal. At that point, both the model's peak flux and peak density are realised, so we will write simply y for \tilde{y} and f^X for \tilde{f}^X . The peak density y is composed from particle probabilities at horizontal and vertical sites, or, symmetrically, particle and hole probabilities at just the horizontal, or just the vertical sites: $y = u + \bar{u} = \begin{array}{c} \bullet \\ \bullet \end{array} + \begin{array}{c} \bullet \\ \circ \end{array} = \begin{array}{c} \bullet \\ \bullet \end{array} + \begin{array}{c} \bullet \\ \circ \end{array} = \begin{array}{c} \bullet \\ \circ \end{array} + \begin{array}{c} \bullet \\ \bullet \end{array} = \frac{1}{2}$. The evolution of the model in such a small neighbourhood can be considered

to be in a stationary state; we will use this stationarity to derive the joint distribution for site probabilities of our models.

A limit for model CS . In general, finding an explicit flux function for a particle evolution model may be difficult, and even the convergence to a scaling limit is not guaranteed. Fortunately, the existence of a continuous scaling limit for model CS follows directly from Proposition 2. Indeed, the model’s convergence at a point on the main diagonal is equivalent to the convergence of scaled LCS length for a pair of equally long uniformly random binary strings, i.e. to the existence of constant γ . As mentioned in the Introduction, this was established already by Chvátal and Sankoff [15] (see also [48, Chapter 1]). In much the same way, the model’s convergence at any other point is equivalent to the convergence of scaled LCS length for a pair of random binary strings with a given limiting ratio of their lengths, which can be established by a slight modification of the same proof.

The Chvátal–Sankoff problem can now be reformulated as finding the peak flux complement $\gamma = \bar{f}^{CS}$ for model CS .

§7. MODEL B

In keeping with the traditional terminology, let us define *model B* (the Bernoulli model) as the network evolution model, where a cell is assigned type “mismatch” with a fixed probability $p \stackrel{\text{def}}{=} \nearrow$, called the model’s (*jump*) *rate*, independently of any site values or types of any other cells (this initial definition will be generalised later). Intuitively, every cell tosses an independent biased coin p to determine its type.

Model B has been applied to the study of the Chvátal–Sankoff problem by Boutet de Monvel [19], Majumdar and Nechaev [32], Priezzhev and Schütz [38], Bukh and Cox [8]. It is closely related to a classical particle model known as the *totally asymmetric simple exclusion process (TASEP)*. The TASEP consists of an array of sites, each occupied by a particle or a hole. It evolves by a particle jumping at a random time into a hole on its right; symmetrically, the hole “jumps” to its left to the site previously occupied by the particle. Updates may occur in continuous time (classical TASEP, which we do not consider any further) or in discrete time (DT-TASEP). Within a time step of DT-TASEP, the update policy may be parallel (the process also known as multi-corner growth of a Young diagram, which we do not consider any further), forward-sequential,

backward-sequential, or sublattice-parallel. The latter three update policies essentially only differ by a change of coordinates, and correspond to model B evolving vertically, horizontally or diagonally, respectively. An analysis of DT-TASEP with different update policies has been given by Rajewsky et al. [39] and by Martin and Schmidt [33]. Model B and DT-TASEP can be considered as a special case of the six-vertex model analysed by Borodin et al. [6], with weights assigned according to measure $\mathcal{P}(p, 0)$ defined therein.

Model B and other network evolution models presented in this paper can also be considered as special cases of stochastic cellular automata (see e.g. [31, 11]). However, the simplifying “well-mixing” assumptions, that are usually made in that context, do not hold for our models.

Cell type probabilities. We note that a cell’s type only affects the model’s behaviour when its entry pair is $\bullet\circ$, distinguishing the events $\bullet\circ$ and $\bullet\circ$. For any other entry pairs, the cell’s exit values are forced by the entry values and are independent of the cell’s type: the corresponding events are $\circ\circ$, $\circ\bullet$, $\bullet\bullet$. In these cases, the cell’s type probability \nearrow can be set differently from p , without affecting the model’s behaviour. Therefore, we can generalise the definition of model B by introducing a formal dependency of a cell’s type on its entry pair, while making sure that the model’s new definition is still invariant with respect to duality of configurations.

Definition 1. We say that a cell’s type *depends exclusively* on a set of sites’ values in a given half-step, if, conditioned on this set, it is conditionally independent of any other site values in the same half-step.

We define $4 = 1 \cdot 2 + 2$ (one dual pair and two self-dual singletons) conditional probabilities for a cell’s type, specifying its exclusive dependence on the entry site pair:

$$p_0 \stackrel{\text{def}}{=} \circ\circ \stackrel{r}{=} p_3 \stackrel{\text{def}}{=} \bullet\bullet \quad p_1 \stackrel{\text{def}}{=} \circ\bullet \quad p_2 \stackrel{\text{def}}{=} \bullet\circ$$

The subscripts correspond to the entry pair values being read as a two-digit binary number, bottom-left to top-right: $p_0 = p_{\circ\circ}$, etc. Intuitively, a cell now has four biased coins p_0, p_1, p_2, p_3 , including a dual pair $p_0 \stackrel{r}{=} p_3$. The cell reads its entry pair (as a binary number), and then tosses the corresponding coin to determine its type; the combination of the cell’s entry pair and its chosen type then determines the cell’s exit pair.

Conditional probability p_2 corresponds to the rate p in the original definition of model B , and determines solely the model’s behaviour (in

particular, its flux). We will therefore reserve the term *rate* for p_2 , whereas the remaining conditional probabilities $p_0 \stackrel{r}{=} p_3$, p_1 will be called *pseudo-rates*. These pseudo-rates do not affect the behaviour of the model, and therefore can temporarily be left unconstrained. This leaves us the freedom to set them later, in an attempt to fit model B to the constraints of model CS .

Alternating sequences. Our models, including model B , will have time-invariant distributions satisfying the following natural property.

Definition 2. An *alternating sequence* is a doubly-infinite sequence of (generally dependent) particle-hole random variables (ξ_i) , $i \in \mathbb{Z}$, that is invariant with respect to

- a shift by 1, mapping $i \mapsto i + 1$
- a reversal about $\frac{1}{2}$, mapping $i \mapsto -i + 1$

both of these with simultaneous exchange between particles and holes.

Note that both a shift and a reversal of the given type flip the parity of indices. Definition 2 implies that an alternating sequence is also invariant with respect to arbitrary shifts and reversals, where holes and particles are exchanged if and only if the parity of indices is flipped.

We consider alternating sequences of site values in a given half-step of diagonal network evolution, identifying arbitrarily the even (respectively, odd) indices of the sequence with the horizontal (respectively, vertical) edge sites. Annotated equality $A \stackrel{r}{=} B$, when applied to such sequences, will stand for “ A and B are obtained from each other by a parity-exchanging reversal with a simultaneous exchange between particles and holes”; this is consistent with the previous usage of this notation to express duality of configurations. Furthermore, annotated equality $A \stackrel{s}{=} B$ will have similar meaning, but with a reversal replaced by a shift. Notation $A \stackrel{sr}{=} B$ will be used when $A \stackrel{s}{=} B$ and $A \stackrel{r}{=} B$ are both applicable.

We denote the marginal site probabilities by

$$u \stackrel{\text{def}}{=} \circ \stackrel{\text{sr}}{=} \bullet \quad \bar{u} \stackrel{\text{def}}{=} \bullet \stackrel{\text{sr}}{=} \circ \quad (2)$$

Substituting (2) into (1), we obtain a simple expression for the peak flux complement of our models.

Proposition 3. *Let X be a network evolution model in a stationary state, where the time-invariant distribution of site values is given by an alternating sequence. Then the peak flux complement is*

$$\bar{f}^X = \bullet + \phi = u + u = 2u$$

In the rest of this section and the next, we designate u , p_0 , p_1 , p_2 as the *main variables*; our goal is to connect them by a system of polynomial equations with integer coefficients. In principle, this could be done directly in terms of the main variables alone; however, for convenience, we will be introducing some *auxiliary variables*. Every auxiliary variable will have a separate equation expressing it in terms of previously introduced variables; thus, auxiliary variables will not add any degrees of freedom to the system, and could easily be eliminated from it, at the expense of making the equations more cumbersome.

AB sequences. We first consider the most basic special case of an alternating sequence.

Definition 3. An alternating sequence (ξ_i) , $i \in \mathbb{Z}$, is an *AB (alternating Bernoulli) sequence*, if all its elements are mutually independent.

In particular, an AB sequence of site values in a given half-step of a network evolution model is a product measure with marginal site probabilities (2).

Time invariance. Consider the evolution of model B on an AB sequence in a stationary state. The model's rate and the site densities of the sequence are connected by the *time-invariance equation*:

$$uu = \bullet\phi = \phi\phi = \bar{u}\bar{u}p_2 \quad (3)$$

We recall a well-known result on the time-invariant distribution for the diagonal evolution of model B (see e.g. Rajewsky et al. [39], Martin and Schmidt [33]).

Theorem 1. *An AB sequence with parameter u determined by (3) is a time-invariant distribution for model B with a given rate p_2 .*

Proof. It is sufficient to show that the AB property is preserved in a single half-step of the evolution of model B . The independence between site values a, b at the end of the half-step in a configuration ζ^{b-} is obvious, since these values are obtained in different cells; independence in a configuration ζ^{b-}_{-a} is established by (3). In the equation (3), the left-hand side expresses

the AB property in a configuration at the half-step's end, while the right-hand side relies on the same property at the half-step's beginning. \square

The Arratia–Steele conjecture. Since the marginal probabilities of cell types in model CS are all equal to $\frac{1}{2}$, and since these types are \boxplus -independent and, more generally, three-wise independent, at some point it was quite natural to conjecture that all cell dependence (i.e. \boxplus -dependence and, more generally, all four-wise and higher-order dependence) could also be ignored, so that model CS would be equivalent to model B with $p = p_2 = \frac{1}{2}$, which we denote $B(1/2)$. Substituting $p_2 = \frac{1}{2}$ into (3), we obtain a quadratic equation that gives us the peak site marginal probability and the peak flux complement for model $B(1/2)$, which can be considered as an approximation for γ :

$$u = \sqrt{2} - 1 = 0.414213 \dots \quad \gamma \approx \bar{f}^{B(1/2)} = 2u = 2(\sqrt{2} - 1) = 0.828427 \dots$$

The conjecture, attributed to Arratia by Steele [47], was that the above expression gives the exact value of γ . This conjecture was disproved by the upper bound $\gamma \leq 0.826280$ due to Lueker [30].

In the remainder of this paper, we will be making repairs to the Arratia–Steele conjecture by weakening the claimed type of model equivalence (local instead of global equivalence), and by replacing, in two successive steps, model B by a network evolution model from a more general class.

§8. LOCAL FITTING OF MODEL B TO MODEL CS

The Arratia–Steele conjecture makes an unsuccessful attempt to fit model $B(1/2)$ to model CS . The next natural step is to replace model $B(1/2)$ by model B with a general rate and pseudo-rates. In doing so, we can set a rate $p_2 > \frac{1}{2}$ for a better fit to the higher peak rate of model CS . We need to compensate for that by lowering (at least one of) the pseudo-rates $p_0 \stackrel{r}{=} p_3, p_1$, so that the marginal cell type probability remains at $\frac{1}{2}$, and the \boxplus -independence of cell types is maintained. Crucially, we do not need to require that the models agree across the whole network: since we are only interested in obtaining the peak flux complement, we only need to achieve the models' agreement in a small neighbourhood of the main diagonal (recall that the peak flux is precisely the flux across the main diagonal); we call this *local fitting* of the models. This attempt to obtain a local fit for model B to model CS will eventually turn out to also be unsuccessful, but somewhat less so than the Arratia–Steele conjecture: it

will give us a better approximation for γ , and will provide a stepping stone to a perfect local fit with a still more general model in subsequent sections.

For convenience, we introduce a pair of auxiliary variables for cell type probabilities conditioned on a single entry site, where the sites in a given half-step are known to form an AB sequence:

$$\begin{aligned} q_0 &\stackrel{\text{def}}{=} \text{⬢} = \text{⬢} + \text{⬢} = \bar{u}p_0 + up_1 \stackrel{r}{=} q_1 \stackrel{\text{def}}{=} \text{⬢} \\ q_1 &\stackrel{\text{def}}{=} \text{⬢} = \text{⬢} + \text{⬢} = \bar{u}p_2 + up_3 \stackrel{r}{=} q_0 \stackrel{\text{def}}{=} \text{⬢} \end{aligned}$$

Reverse cell type probabilities. Forgetting temporarily about cell types, the evolution of site values in model B can be “turned back in time” by considering a natural reverse process, where site values in half-step t are conditioned on site values (without cell types) in half-step $t + 1$. Both the forward and the reverse processes on site values can be described symmetrically as

$$1 = \text{⬢} = \text{⬢} = \text{⬢} = \text{⬢} = \text{⬢} = \text{⬢} \quad p_2 = \text{⬢} = \text{⬢} \quad \bar{p}_2 = \text{⬢} = \text{⬢}$$

Although it is not required for our results, it is remarkable, and not difficult to check via (3), that in the stationary state, it is impossible to distinguish probabilistically whether a given configuration of site values has been obtained by the forward or by the reverse process.

Reintroducing cell types breaks the symmetry between the forward and the reverse processes: a cell’s type determines its operation in the forward process only. In the reverse process, a cell’s type depends exclusively on its pair of exit sites (here, the terminology “exit sites” is still relative to the forward process). We denote the $4 = 1 \cdot 2 + 2$ resulting reverse conditional probabilities by auxiliary variables

$$r_{ab} \stackrel{\text{def}}{=} \text{⬢} \stackrel{r}{=} r_{\bar{b}\bar{a}} \stackrel{\text{def}}{=} \text{⬢}$$

which are determined by the model’s parameters via the equations

$$\begin{aligned} \bar{r}_0 &= \text{⬢} = \text{⬢} = \bar{p}_0 \stackrel{r}{=} \bar{r}_3 = \text{⬢} = \text{⬢} = \bar{p}_3 \\ \bar{r}_1 \bar{u} \bar{u} &= \text{⬢} = \text{⬢} = uu\bar{p}_1 \quad \bar{r}_2 = \text{⬢} = 1 \end{aligned} \tag{4}$$

Total probability. An individual cell in model CS takes its types \searrow and \swarrow equiprobably. Therefore, in a local fit of the models, the site probabilities, rate and pseudo-rates of model B must satisfy the total probability equation, where the pseudo-rates p_0, p_1 fulfill their purpose of balancing out the bias in the rate p_2 :

$$\bar{u}\bar{u}p_2 + 2u\bar{u}p_0 + uu\bar{p}_1 = \text{⬢} + 2\text{⬢} + \text{⬢} = \swarrow = \frac{1}{2} \tag{5}$$

Here, we have collected an equiprobable dual pair of terms $\phi \circledast \bullet \stackrel{r}{=} \bullet \circledast \phi$ into a single term $2\phi \circledast \bullet$. The remaining terms are self-dual singletons. We shall use the same shortcut subsequently without special notice; the collected mutually dual terms can always be distinguished by the leading coefficient 2.

Linking the models. We now attempt to link models B and CS . Our goal is to assign the rate and pseudo-rates for model B so that the site values and cell types at any given half-step would be probabilistically indistinguishable in both models.

In model B , a cell's type in half-step t depends exclusively on its entry site pair. Each site of this pair is an exit site for one of a pair of diagonally adjacent cells in half-step $t-1$, and depends exclusively on those cells' entry site pairs; we thus have an exclusive dependence of a cell's type in half-step t on a quadruple (two disjoint pairs) of adjacent sites in half-step $t-1$. Each site of this quadruple, in its turn, is an exit site for one of a triple of diagonally adjacent cells in half-step $t-2$, and depends exclusively on those cells' entry site pairs; we thus have an exclusive dependence of a cell's type in half-step t on a sextuple (three disjoint pairs) of adjacent sites in half-step $t-2$. On the other hand, in model CS , a cell's type in half-step t is determined uniquely by just the types of three preceding cells, two in half-step $t-1$ and one in half-step $t-2$, forming a \boxminus -shape between themselves, and a \boxplus -shape together with the current cell. Therefore, in order to relate models B and CS , a system of polynomial equations can be obtained by listing exhaustively all possible configurations of the relevant site values and cell types over three half-steps of both models' evolution. The use of duality and of the reverse cell type probabilities will provide substantial shortcuts for such an exhaustive enumeration, helping us to avoid listing dozens of configurations explicitly.

There are three linking equations: one for the rate p_2 , and two (considering duality) for the pseudo-rates $p_0 \stackrel{r}{=} p_3, p_1$. The left-hand side of every equation represents a single half-step configuration for a given rate or pseudo-rate. The right-hand side enumerates exhaustively each of the three half-step configurations of model CS that result in the configuration in the left-hand side with nonzero probability. Due to \boxminus -independence of cell types, the first two of these steps are probabilistically indistinguishable from ones of model B , and their probabilities are assigned accordingly in the equations; by using the reverse probabilities, we avoid an explicit enumeration of the site values in the first half-step, while we do enumerate

cell types. In the third half-step, the cell type is determined uniquely from the \boxplus -configuration of model CS ; the requirement that this half-step must also be probabilistically indistinguishable from one of model B provides the desired equation.

In particular, the linking equation for the rate p_2 is as follows:

$$\begin{aligned} \bar{u}\bar{u}p_2 = \text{diagram} &= \left(\text{diagram} + 2 \text{diagram} + \text{diagram} \right) + 2 \left(\text{diagram} + \text{diagram} \right) + \text{diagram} \\ &= 2u\bar{u}\bar{r}_2q_0\bar{q}_0 + 2\bar{u}\bar{u}u(r_0p_2q_0 + \bar{r}_0p_2\bar{q}_0) + \bar{u}\bar{u}\bar{u}r_1p_2p_2 \end{aligned} \quad (6a)$$

The terms representing impossible events have been crossed out and dropped from the equation; from now on, such terms will be omitted without special notice.

The remaining linking equations are as follows:

$$\begin{aligned} u\bar{u}p_0 = \text{diagram} &= \left(\text{diagram} + \text{diagram} \right) + \left(\text{diagram} + \text{diagram} \right) + \text{diagram} + \text{diagram} \\ &= \bar{u}u(\bar{r}_0\bar{q}_1q_0 + r_0\bar{q}_1\bar{q}_0) + u\bar{u}u(r_0p_0q_0 + \bar{r}_0p_0\bar{q}_0) + u\bar{u}\bar{u}r_1p_0p_2 + \bar{u}\bar{u}\bar{u}r_1\bar{q}_1p_2 \end{aligned} \quad (6b)$$

$$uup_1 = \text{diagram} = \text{diagram} + 2 \text{diagram} + \text{diagram} = u\bar{u}\bar{u}ur_1p_0p_3 + 2u\bar{u}\bar{u}r_1p_0\bar{q}_1 + \bar{u}\bar{u}r_1\bar{q}_1\bar{q}_1 \quad (6c)$$

It is important to note that the connection between the two models expressed by equations (6) is incomplete: while the equations relate the rate and the pseudo-rates of model B to model CS , they do not guarantee the preservation of the AB property on the site values. In particular, the equations' left-hand sides express site independence within a configuration of the form $\overset{a}{\underset{b}{\cdot}}$, but none of the equations implies site independence within a configuration of the form $\underset{a}{\overset{b}{\cdot}}$. Thus, there is no guarantee that our goal of probabilistic indistinguishability between the two models has been achieved: in fact, it has not, and in general model B turns out to be insufficient for a perfect fit. This limitation of model B will be overcome in subsequent sections.

Solving the equations. The resulting system has four main variables u , p_0 , p_1 , p_2 , involved in five main equations: one time-invariance equation (3), one total probability equation (5), and three linking equations (6). There are also some auxiliary variables, each of which is introduced via its own separate equation. Thus, the system is overdetermined by one equation. However, it is still consistent, since the total probability equation (5) is a consequence of the \boxplus -independence of cell types, which is implied by the time-invariance and the linking equations. While the total probability equation is formally redundant, we keep it in the system for its symmetry

and, more importantly, as an aid to computer algebra software in solving the system.

Since all the variables in the system represent probabilities, we are only interested in real solutions between 0 and 1; we call such solutions *admissible*. The system has a unique admissible solution; we denote by $B(opt)$ model B with the specific set of parameters provided by this solution.

The system's admissible solution can be obtained analytically using computer algebra software. In particular, Mathematica returns it instantly, expressed in exact radicals (for this, function `Solve` needs to be used with option `Quartics -> True`). As a result, we obtain the site marginal probability and an estimate for γ via the peak flux complement as

$$u = \sqrt{\frac{7}{3}} - \sqrt{\frac{23 - 5\sqrt{21}}{6}} - 1 = 0.407025 \dots \quad \gamma \approx \bar{f}^{B(opt)} = 2u = 0.814050 \dots$$

and the model's rate and pseudo-rates as

$$\begin{aligned} p_0 = p_3 &= -\frac{8}{3} + \frac{49}{6}u - uu - \frac{1}{2}uuu = 0.457987 \dots \\ p_1 &= \frac{29}{2} - 51u + \frac{75}{2}uu + 9uuu = 0.561206 \dots \\ p_2 &= -\frac{2}{3} + \frac{34}{3}u - 19uu - 4uuu = 0.528838 \dots \end{aligned}$$

We have thus attempted to obtain a local fit of model B to model CS , expressing the various constraints of the latter by polynomial equations with integer coefficients, and obtaining the unique admissible solution of the resulting equation system as model $B(opt)$. However, this fit is not perfect, since the AB sequence property is not preserved; enforcing its preservation by introducing additional equations would make the system truly overdetermined and inconsistent. We conclude that not only model $B(1/2)$ of the Arratia–Steele conjecture, but even the more general model B is still too rigid to provide a perfect local fit to model CS . In the next section, we will further generalise the model's definition, increasing its flexibility in order to achieve this goal.

§9. MODEL M

We now generalise model B in order to give it more flexibility to fit model CS . Following the same terminological pattern, we call this generalisation *model M* (the Markov model). While in model B , a cell's type depends just on the cell's entry pair, in model M it also depends on two

further sites, lying anti-diagonally on either side of the entry pair. In total a cell's type depends on an anti-diagonal quadruple of sites.

Consider model CS evolving diagonally, and let us take a single cell of this model in half-step t . Each of this cell's exit values becomes, in half-step $t + 1$, an entry value for one of a pair of cells, adjacent to the original cell in a \boxplus -shape. In their turn, this pair of cells have four exit values, forming an adjacent anti-diagonal site quadruple. The middle two sites of this quadruple become, in half-step $t + 2$, the entry pair for the fourth cell, completing the \boxplus -shape; the type of this fourth cell is determined uniquely by the types of the first three cells. Intuitively, as the model's evolution progresses, the distribution of site values in a given half-step stores some information about cell types in preceding half-steps. In particular, the adjacent site quadruple in half-step $t + 2$ stores the information about the types of the three cells in a \boxplus -shape in half-steps t and $t + 1$, which is the full information necessary to complete the \boxplus -shape. Thus, intuition suggests that the mutual dependence between a site quadruple and a cell in a given half-step should be precisely the right one to capture the behaviour of model CS .

Cell type probabilities. We generalise the definition of model B , replacing a cell's type exclusive dependence on its entry pair by that on its *extended entry quadruple*, which is made up by the cell's entry pair and the two sites adjacent to it antidiagonally on either side. We also make sure that the new model's definition is still invariant with respect to duality of configurations. Thus, we define $16 = 6 \cdot 2 + 4$ (six dual pairs and four self-dual singletons) conditional probabilities for a cell's type, specifying its exclusive dependence on the extended entry quadruple:

$$p_{abcd} \stackrel{\text{def}}{=} \begin{array}{c} \textcolor{red}{r} \textcolor{red}{c} \textcolor{red}{d} \\ \textcolor{red}{b} \textcolor{red}{a} \\ \textcolor{red}{-} \textcolor{red}{a} \end{array} \begin{array}{c} \textcolor{red}{r} \\ \textcolor{red}{c} \\ \textcolor{red}{d} \end{array} = p_{\bar{d}\bar{c}\bar{b}\bar{a}} \stackrel{\text{def}}{=} \begin{array}{c} \textcolor{red}{r} \textcolor{red}{c} \textcolor{red}{d} \\ \textcolor{red}{b} \textcolor{red}{a} \\ \textcolor{red}{-} \textcolor{red}{d} \end{array} \begin{array}{c} \textcolor{red}{r} \\ \textcolor{red}{c} \\ \textcolor{red}{d} \end{array}$$

The subscripts correspond to the extended entry quadruple values being read as a four-digit binary number, bottom-left to top-right: $p_0 = p_{0000}$, etc. Intuitively, a cell now has 16 biased coins p_0, \dots, p_{15} , including six dual pairs

$$p_i \stackrel{r}{=} p_j \quad (i, j) \in \{(0, 15), (1, 7), (2, 11), (4, 13), (6, 9), (8, 14)\}$$

The cell reads its extended entry quadruple (as a binary number), and then tosses the corresponding coin to determine its type; the combination of the cell's extended entry quadruple and its chosen type then determines the cell's exit pair.

$$p_5 = \text{diagram} \quad p_4 = \text{diagram} \stackrel{r}{=} p_{13} = \text{diagram} \quad p_{12} = \text{diagram}$$

AM2 sequences and space invariance. We return to considering alternating sequences of (generally dependent) binary random variables representing site values in a given half-step. Apart from fixed marginal site probability u defined by (2), an alternating sequence also has two fixed first-order and four fixed second-order conditional site probabilities

$$v_a \stackrel{\text{def}}{=} \text{diagram 1} \stackrel{s}{=} \text{diagram 2} \stackrel{r}{=} \text{diagram 3} \stackrel{s}{=} \text{diagram 4} \quad w_{ab} \stackrel{\text{def}}{=} \text{diagram 5} \stackrel{s}{=} \text{diagram 6} \stackrel{r}{=} \text{diagram 7} \stackrel{s}{=} \text{diagram 8} \quad (7)$$

$$u\bar{v}_0 = \text{diagram} = \bar{u}v_1 \quad \bar{v}_0w_0 = \text{diagram} = v_0\bar{w}_2 \quad \bar{v}_1w_1 = \text{diagram} = v_1\bar{w}_3 \quad (8)$$

Definition 4. An alternating sequence (ξ_i) , $i \in \mathbb{Z}$, is an *AM2 (alternating second-order Markov) sequence* if, given an adjacent pair (ξ_i, ξ_{i+1}) , the infinite prefix (ξ_j) , $j < i$, is conditionally independent of the infinite suffix (ξ_k) , $k > i + 1$.

We introduce auxiliary variables for unconditional probabilities of finite AM2 sequences of length 2, 4, 6. We define

$$u_{ab}^2 \stackrel{\text{def}}{=} u^{[\bar{a}]} v_a^{[b]} = \overset{r}{a} \overset{b}{b} \overset{s}{=} \overset{r}{- \bar{a}} \overset{b}{=} \overset{r}{- b} \overset{a}{=} \overset{s}{\bar{b}} \overset{a}{-}$$

$$u_{abcd}^4 \stackrel{\text{def}}{=} u_{ab}^2 w_{\bar{a}\bar{b}}^{[\bar{c}]} w_{bc}^{[d]} = \overset{r}{a} \overset{b}{b} \overset{s}{=} \overset{r}{- \bar{a}} \overset{b}{=} \overset{r}{- b} \overset{a}{=} \overset{s}{\bar{b}} \overset{a}{-}$$

$$u_{abcdeg}^6 \stackrel{\text{def}}{=} u_{abcd}^4 w_{cd}^{[e]} w_{de}^{[g]} = \underset{a}{\overset{r}{c}} \underset{b}{\overset{r}{d}} \underset{c}{\overset{r}{e}} \underset{d}{\overset{r}{g}} \underset{e}{\overset{s}{-}} \underset{a}{\overset{r}{c}} \underset{b}{\overset{r}{d}} \underset{c}{\overset{r}{e}} \underset{d}{\overset{r}{g}} \underset{e}{\overset{s}{-}} \underset{a}{\overset{r}{c}} \underset{b}{\overset{r}{d}} \underset{c}{\overset{r}{e}} \underset{d}{\overset{r}{g}} \underset{e}{\overset{s}{-}} \underset{a}{\overset{r}{c}} \underset{b}{\overset{r}{d}} \underset{c}{\overset{r}{e}} \underset{d}{\overset{r}{g}} \underset{e}{\overset{s}{-}}$$

Similar notation, which we won't require, could be introduced for AM2 sequences of any finite length.

Time invariance. Generalising Theorem 1, we now consider the evolution of model M on an AM2 sequence in a stationary state. The model's partial rates and the parameters of the sequence are linked by the *time-invariance equations*:

$$w_0 \bar{w}_2 = \text{[diagram]} = \bar{w}_3 w_1 \quad (9a)$$

$$\cancel{w_2 w_2 w_2} = \text{[diagram]} = \cancel{\bar{w}_1 \bar{w}_1 \bar{w}_1} \rightsquigarrow w_2 w_2 = \bar{w}_1 \bar{w}_1 \bar{p}_5 \quad (9b)$$

$$\cancel{w_2 w_2 w_3} = \text{[diagram]} = \cancel{\bar{w}_1 \bar{w}_1 \bar{w}_0} \rightsquigarrow w_2 w_3 = \bar{w}_1 \bar{w}_0 \bar{p}_4 \quad (9c)$$

$$\cancel{w_3 w_3 w_3} = \text{[diagram]} = \cancel{\bar{w}_0 \bar{w}_0 \bar{w}_0} \rightsquigarrow w_3 w_3 = \bar{w}_0 \bar{w}_0 \bar{p}_{12} \quad (9d)$$

In the above equations, conditioning on a pair $\circ\circ = \circ\circ$ or $\bullet\bullet = \bullet\bullet$ corresponds to cancelling out its probability from both sides of the equation. We also use (9a) to cancel some of the probabilities on either side of each of (9b)–(9d).

Theorem 2. *An AM2 sequence with parameters u, v_a, w_{ab} determined by equations (8), (9) is a time-invariant distribution for model M with given partial rates p_{abcd} .*

Proof. It is sufficient to show that the AM2 property is preserved in a single half-step of the evolution of model M . The conditional independence between site values a, d , given site values c, d , at the end of the half-step in a configuration $\underset{a}{\overset{r}{c}} \underset{b}{\overset{d}{-}}$ is established by (8), (9a), and in a configuration $\underset{a}{\overset{r}{c}} \underset{b}{\overset{d}{-}}$ by (8), (9b)–(9d). In each of the equations (9), the left-hand side expresses the AM2 property in a configuration at the half-step's end, while the right-hand side relies on the same property at the half-step's beginning. \square

Note that the time-invariance equations (9) link the four-dimensional variety of AM2 sequences with the three-dimensional variety of realisations of model M ; therefore, only a three-dimensional subvariety of AM2 sequences is realisable by model M as its time- and duality-invariant distributions.

§10. LOCAL FITTING OF MODEL M TO MODEL CS

Following the approach developed in previous sections, we now use the freedom to set the previously unconstrained pseudo-rates, in order to obtain a local fit of model M to model CS . In contrast with model B , model M turns out to have sufficient flexibility for a seamless fit.

For convenience, we introduce eight auxiliary variables for cell type probabilities conditioned on a subset of three sites in the extended entry quadruple, where the sites in a given half-step are known to form an AM2 sequence:

$$q_{abc} \stackrel{\text{def}}{=} \begin{array}{c} \text{red } \bar{a} \text{ } \text{red } \bar{b} \text{ } \text{red } \bar{c} \\ \text{red } \bar{a} \text{ } \text{red } \bar{b} \text{ } \text{red } \bar{c} \end{array} = \begin{array}{c} \text{red } \bar{a} \text{ } \text{red } \bar{b} \text{ } \text{red } \bar{c} \\ \text{red } \bar{a} \text{ } \text{red } \bar{b} \text{ } \text{red } \bar{c} \end{array} + \begin{array}{c} \text{red } \bar{a} \text{ } \text{red } \bar{b} \text{ } \text{red } \bar{c} \\ \text{red } \bar{a} \text{ } \text{red } \bar{b} \text{ } \text{red } \bar{c} \end{array} = w_{\bar{b}\bar{c}} p_{abc\circ} + \bar{w}_{\bar{b}\bar{c}} p_{abc\bullet} = \begin{array}{c} \text{red } \bar{a} \text{ } \text{red } \bar{b} \text{ } \text{red } \bar{c} \\ \text{red } \bar{a} \text{ } \text{red } \bar{b} \text{ } \text{red } \bar{c} \end{array}$$

Reverse cell type probabilities. We recall from previous discussion that model B has a remarkable property of having a reverse process defined on site values. We might expect that model M has an analogous property, but in this case reversibility is harder to establish, so we leave it as a conjecture. Just as with model B , we will define reverse probabilities on cell types of model M from first principles, without relying on any special reversibility properties of the process on site values. In order to derive the linking equations, the probability of a cell's type in half-step t will need to be conditioned on a sextuple of site values in half-step $t+1$ (without any claim that such a dependence is exclusive). We denote the $64 = 28 \cdot 2 + 8$ resulting reverse conditional probabilities by auxiliary variables

$$r_{abcdeg} \stackrel{\text{def}}{=} \begin{array}{c} \text{red } \bar{a} \text{ } \text{red } \bar{b} \text{ } \text{red } \bar{c} \\ \text{red } \bar{a} \text{ } \text{red } \bar{b} \text{ } \text{red } \bar{c} \end{array} = r_{\bar{g}\bar{e}\bar{d}\bar{c}\bar{b}\bar{a}} \stackrel{\text{def}}{=} \begin{array}{c} \text{red } \bar{a} \text{ } \text{red } \bar{b} \text{ } \text{red } \bar{c} \\ \text{red } \bar{a} \text{ } \text{red } \bar{b} \text{ } \text{red } \bar{c} \end{array}$$

We consider these probabilities in subsets of four, each subset identified by fixed site values a, b, e, g . Within each subset, we obtain four equations parameterised by the middle site pair c, d to determine each of the probabilities r_{abcdeg} :

$$u_{\circ\circ d c \circ\circ}^6 \bar{r}_{\circ\circ c d \circ\circ} = \begin{array}{c} \text{red } \bar{a} \text{ } \text{red } \bar{b} \text{ } \text{red } \bar{c} \\ \text{red } \bar{a} \text{ } \text{red } \bar{b} \text{ } \text{red } \bar{c} \end{array} = \begin{array}{c} \text{red } \bar{a} \text{ } \text{red } \bar{b} \text{ } \text{red } \bar{c} \\ \text{red } \bar{a} \text{ } \text{red } \bar{b} \text{ } \text{red } \bar{c} \end{array} = u_{\circ\circ c d \circ\circ}^6 \bar{p}_{\circ\circ c d \circ\circ}$$

$$\begin{aligned}
u_{\bullet\circ d c \circ \circ}^6 \bar{r}_{\circ \circ c d \bullet \bullet} &= \text{diagram} = \text{diagram} + \text{diagram} = u_{\circ \circ c d \bullet \bullet}^6 \bar{p}_{\circ c d \bullet} + u_{\circ \circ c d \bullet \bullet}^6 \bar{p}_{\circ c d \bullet} q_{d \bullet \circ} \\
u_{\circ \bullet d c \circ \circ}^6 \bar{r}_{\circ \circ c d \bullet \circ} &= \text{diagram} = \text{diagram} = u_{\circ \circ c d \bullet \circ}^6 \bar{p}_{\circ c d \bullet} \bar{q}_{d \bullet \circ} \\
u_{\bullet \bullet d c \circ \circ}^6 \bar{r}_{\circ \circ c d \bullet \bullet} &= \text{diagram} = \text{diagram} = u_{\circ \circ c d \bullet \bullet}^6 \bar{p}_{\circ c d \bullet} \\
u_{\circ \circ d c \bullet \bullet}^6 \bar{r}_{\bullet \bullet c d \circ \circ} &= \text{diagram} = \text{diagram} = u_{\bullet \bullet c d \circ \circ}^6 \bar{p}_{\bullet c d \circ} \\
u_{\circ \circ d c \bullet \bullet}^6 \bar{r}_{\bullet \circ c d \circ \circ} &= \text{diagram} = \text{diagram} + \text{diagram} = u_{\bullet \circ c d \circ \circ}^6 \bar{p}_{\bullet c d \circ} + u_{\bullet \circ c d \circ \circ}^6 \bar{p}_{\bullet c d \circ} q_{\bar{c} \bullet \circ} \\
u_{\circ \circ d c \bullet \bullet}^6 \bar{r}_{\bullet \bullet c d \circ \circ} &= \text{diagram} = \text{diagram} = u_{\bullet \bullet c d \circ \circ}^6 \bar{p}_{\bullet c d \circ} \bar{q}_{\bar{c} \bullet \circ} \\
u_{\bullet \bullet d c \circ \circ}^6 \bar{r}_{\circ \bullet c d \bullet \circ} &= \text{diagram} = \text{diagram} + \text{diagram} + \text{diagram} + \text{diagram} \\
&= u_{\bullet \circ c d \bullet \circ}^6 \bar{p}_{\bullet c d \bullet} + u_{\bullet \circ c d \bullet \circ}^6 \bar{p}_{\bullet c d \bullet} q_{d \bullet \circ} \\
&\quad + (u_{\bullet \circ c d \bullet \circ}^6 \bar{p}_{\bullet c d \bullet} + u_{\bullet \circ c d \bullet \circ}^6 \bar{p}_{\bullet c d \bullet} q_{d \bullet \circ}) q_{\bar{c} \bullet \circ} \\
u_{\bullet \circ d c \bullet \bullet}^6 \bar{r}_{\bullet \bullet c d \circ \circ} &= \text{diagram} = \text{diagram} + \text{diagram} = (u_{\bullet \circ c d \bullet \bullet}^6 \bar{p}_{\bullet c d \bullet} + u_{\bullet \circ c d \bullet \bullet}^6 \bar{p}_{\bullet c d \bullet} q_{d \bullet \circ}) \bar{q}_{\bar{c} \bullet \circ} \\
u_{\bullet \circ d c \bullet \bullet}^6 \bar{r}_{\bullet \circ c d \bullet \circ} &= \text{diagram} = \text{diagram} = u_{\bullet \circ c d \bullet \circ}^6 \bar{p}_{\bullet c d \bullet} \bar{q}_{\bar{c} \bullet \circ} \bar{q}_{d \bullet \circ}
\end{aligned}$$

For the sake of brevity, we do not consider various cancellations and simplifications that could be made in the above equations. For instance, the case $cd = \bullet\circ$ in each of the four-equation subsets could be written as

$$\bar{r}_{ab\bullet\circ eh} = \text{diagram} = 1 \text{ for all } a, b, e, h.$$

Total probability. Similarly to model B , in a local fit of model M to model CS the site probabilities, rate and pseudo-rates of model B must satisfy the total probability equation, where the pseudo-rates p_i , $i \in \{0, \dots, 15\} \setminus \{4, 5, 12, 13\}$ fulfill their purpose of balancing out the bias

in the partial rates p_i , $i \in \{4, 5, 12, 13\}$:

$$\sum_{a,b,c,d \in \{\circ, \bullet\}} u_{dcba}^4 p_{abcd} = \sum_{a,b,c,d \in \{\circ, \bullet\}} \frac{r_{a,b}^c d}{r_{a,b}^c} = \nearrow = \frac{1}{2} \quad (10)$$

Linking the models. Continuing the previously established pattern, we now link models M and CS . Our goal this time is to assign the partial rates and pseudo-rates for model M so that the site values and cell types at any given half-step would be probabilistically indistinguishable in both models.

In model M , a cell's type in half-step t depends exclusively on its extended entry quadruple. Each site of this quadruple is an exit site for one of a pair of diagonally adjacent cells in half-step $t-1$, and depends exclusively on those cells' extended entry site quadruples; we thus have an exclusive dependence of a cell's type in half-step t on a sextuple (two overlapping quadruples) of adjacent sites in half-step $t-1$. Each site of this sextuple, in its turn, is an exit site for one of a triple of diagonally adjacent cells in half-step $t-2$, and depends exclusively on those cells' extended entry site quadruples; we thus have an exclusive dependence of a cell's type in half-step t on an octuple (three overlapping quadruples) of adjacent sites in half-step $t-2$. On the other hand, in model CS , as discussed before, a cell's type in half-step t is determined uniquely by just the types of three preceding cells, two in half-step $t-1$ and one in half-step $t-2$, forming a \boxminus -shape between themselves, and a \boxplus -shape together with the current cell. Therefore, in order to relate models M and CS , a system of polynomial equations can be obtained by listing exhaustively all possible configurations of the relevant site values and cell types over three half-steps of both models' evolution. As before, the use of duality and of the reverse process will provide substantial shortcuts for such an exhaustive enumeration, this time helping us to avoid an explicit listing of not just dozens, as was the case with model B , but hundreds of configurations.

For convenience, we introduce 64 auxiliary variables for joint bidirectional (two forward and one reverse) cell type conditional probabilities

$$s_{abcd}^{EFG} \stackrel{\text{def}}{=} \frac{r_{a,b}^c d}{r_{a,b}^c} = s_{\bar{d}\bar{c}\bar{b}\bar{a}}^{GFE} \stackrel{\text{def}}{=} \frac{r_{\bar{c},\bar{d}}^{\bar{a}}}{r_{\bar{c},\bar{d}}^{\bar{a}}} \quad E+F+G \text{ is odd}$$

These 64 variables form four subsets of 16 variables: s_{abcd}''' , s_{abcd}^{\wedge} , s_{abcd}^{\vee} , s_{abcd}^{\cap} . We have

$$s_{abcd}^{EFG} = \frac{r_{a,b}^c d}{r_{a,b}^c} = \frac{r_{a,b}^c d}{r_{a,b}^c} + \frac{r_{a,b}^c d}{r_{a,b}^c} + \frac{r_{a,b}^c d}{r_{a,b}^c} + \frac{r_{a,b}^c d}{r_{a,b}^c}$$

$$\begin{aligned}
&= u_{\circ d c b a \circ}^6 p_{\circ a b c}^{[E]} r_{\circ a b c d \circ}^{[F]} p_{b c d \circ}^{[G]} + u_{\bullet d c b a \circ}^6 p_{\circ a b c}^{[E]} r_{\circ a b c d \bullet}^{[F]} p_{b c d \bullet}^{[G]} + \\
&\quad u_{\circ d c b a \bullet}^6 p_{\bullet a b c}^{[E]} r_{\bullet a b c d \circ}^{[F]} p_{b c d \circ}^{[G]} + u_{\bullet d c b a \bullet}^6 p_{\bullet a b c}^{[E]} r_{\bullet a b c d \bullet}^{[F]} p_{b c d \bullet}^{[G]}
\end{aligned}$$

Note that these equations rely on the exclusive dependence of cell types E , G on the respective site quadruples, and do not require the reverse dependence of F on the site sextuple also to be exclusive (which it is not).

There are 10 linking equations, one for each partial rate and pseudo-rate (considering duality). As before, the left-hand side of every equation represents a single half-step configuration for a given partial rate or pseudo-rate, while the right-hand side enumerates exhaustively each of the corresponding three half-step configurations.

$$p_0 = \text{diagram} = \text{diagram} + \text{diagram} + \text{diagram} + \text{diagram} = s_0''' + s_0^\wedge + s_0^\vee + s_0^\vee \quad (11a)$$

$$\begin{aligned}
p_1 u_8^4 &= \text{diagram} = \left(\text{diagram} + \text{diagram} + \text{diagram} + \text{diagram} \right) + \left(\text{diagram} + \text{diagram} \right) \\
&= u_1^4 (s_1''' + s_1^\wedge + s_1^\vee + s_1^\vee) + u_2^4 (s_2''' + s_2^\vee) \quad (11b)
\end{aligned}$$

$$p_2 u_4^4 = \text{diagram} = \text{diagram} = u_2^4 (s_2^\wedge + s_2^\vee) \quad (11c)$$

$$p_3 u_{12}^4 = \text{diagram} = \text{diagram} + 2 \text{diagram} = u_3^4 (s_3''' + 2s_3^\wedge + s_3^\vee) \quad (11d)$$

$$\begin{aligned}
p_4 u_2^4 &= \text{diagram} = \left(\text{diagram} + \text{diagram} \right) + \left(\text{diagram} + \text{diagram} \right) \\
&= u_8^4 (s_8''' + s_8^\wedge) + u_4^4 (s_4^\wedge + s_4^\vee) \quad (11e)
\end{aligned}$$

$$\begin{aligned}
p_5 u_{10}^4 &= \text{diagram} = 2 \text{diagram} + 2 \left(\text{diagram} + \text{diagram} \right) + \text{diagram} \\
&= 2u_5^4 s_5^\wedge + 2u_9^4 (s_9''' + s_9^\wedge) + 2u_{10}^4 s_{10}''' \quad (11f)
\end{aligned}$$

$$p_6 u_6^4 = \text{diagram} = \left(\text{diagram} + \text{diagram} \right) + \text{diagram} = u_6^4 (s_6^\wedge + s_6^\vee) + u_{10}^4 s_{10}^\wedge \quad (11g)$$

$$p_8 u_1^4 = \text{diagram} = \text{diagram} + \text{diagram} = u_8^4 (s_8^\vee + s_8^\vee) \quad (11h)$$

$$p_{10} u_5^4 = \text{diagram} = \text{diagram} = u_{10}^4 s_{10}^\vee \quad (11i)$$

$$p_{12} u_3^4 = \text{diagram} = 2 \text{diagram} = 2u_{12}^4 s_{12}^\wedge \quad (11j)$$

Theorem 3. *An AM2 sequence with parameters u , v_a , w_{ab} determined by equations (8), (9), (10), (11) is a time-invariant distribution for model CS.*

Proof. Consider the set of sequence parameters, partial rates and pseudo-rates determined by the equations. By Theorem 2, the AM2 sequence with these parameters is time-invariant for model M with the rates and pseudo-rates given by equations (8), (9). By design of equations (10), (11), the time-invariant distribution for such a model is identical to a time-invariant distribution of model CS over three successive half-steps. Furthermore, in model CS , cell types in a given half-step are completely determined by the cell types in two previous half-steps. Therefore, the invariance of a distribution over three successive half-steps of model CS is sufficient for its overall invariance. \square

Solving the equations. The resulting system has 17 main variables u , v_a , w_{ab} , p_i , $i \in \{0, \dots, 6, 8, 10, 12\}$, involved in 18 main equations: three space-invariance equations (8), four time-invariance equations (9), one total probability equation (10), and 10 linking equations (11). There also are some auxiliary variables, each of which is introduced via its own separate equation. Similarly to the analysis of model B in previous sections, the system is overdetermined by one equation, but still consistent, since the total probability equation is implied by the time-invariance and the linking equations. By the existence and uniqueness of constant γ , the system must have a unique admissible solution, providing values for the parameters of the stationary state of model CS in a small neighbourhood of the main diagonal. We are now ready to state our main result.

Theorem 4. *Constant γ is an algebraic number.*

Proof. By (1), we have $\gamma = \bar{f}^{CS} = 2u$, which is a (very simple) polynomial in the marginal site probability u . This probability, in its turn, is a variable in our (quite complicated) polynomial system expressing the local fit of models M and CS by Theorem 3. All the coefficients in these polynomials are integers 1 and 2. In an isolated real solution of a polynomial system with rational coefficients, all variables must take algebraic values¹, therefore we conclude that γ is also algebraic. \square

¹Although this statement may seem easy, it is not entirely trivial, since we are not assuming that the whole set of solutions is zero-dimensional. One way of justifying this statement is by observing that the property of a polynomial system with rational coefficients to have a unique real solution in a given neighbourhood is expressible in

On the negative side, it seems unlikely that our system has a closed-form analytic solution; obtaining even a numerical solution appears to be far beyond the capabilities of modern computer algebra software. It is conceivable that a reasonably accurate numerical approximation for the solution can be obtained either by an exhaustive enumeration of all possible configurations for the evolution of model *CS* over a sufficient number of time steps, or by Monte Carlo simulation of such an evolution over a substantially larger number of time steps, or a combination of both approaches. We leave it as an open direction for future work.

§11. CONCLUSION

In this paper, we have linked the Chvátal–Sankoff problem to the parameters of a certain stochastic particle process (model *M*), using existing results on the combinatorial structure of the LCS problem and the theory of continuous scaling limits for discrete particle processes. We have obtained a specific system of polynomial equations with integer coefficients that determines the parameters for this process, which implies that γ is an algebraic number. Short of finding a closed-form solution for such a polynomial system, which appears to be unlikely, our approach essentially resolves the Chvátal–Sankoff problem, at least in theory. Some immediate further questions arise, listed in the increasing order of their apparent difficulty.

Computational experiments. Obtaining an accurate numerical solution for our system, improving and reconciling various existing numerical estimates for γ , appears to be non-trivial, but may well be possible with some reasonable software design and programming effort, and possibly some advanced hardware.

Strings of unequal lengths. For uniformly random binary input strings of unequal lengths, the solution should be possible by a direct extension of our approach; however, it may become even more cumbersome, since we have used the symmetry between the input strings and the resulting duality properties as a substantial shortcut.

the first-order logical theory of real numbers. By Tarski’s theorem on the elementary equivalence of real closed fields (see e.g. Jensen and Lenzing [25, Theorem 2.28], Chang and Keisler [12, Theorem 5.4.4]), such a system must have a unique real algebraic solution in the same neighbourhood; the two solutions must obviously coincide.

Levenshtein distance. The Levenshtein distance problem for random binary strings has been considered by Schind and Bilardi [45]. It is well-known that the Levenshtein distance can be obtained as the LCS length between strings over the alphabet extended with an extra character \$, blowing up each character c to a two-character substring sc . Thus, the Levenshtein distance problem for binary strings becomes a special case of the LCS problem for ternary strings, where the cell types in the corresponding transposition network possess the three-wise independence and four-wise dependence properties similar to those of binary strings. Therefore, a solution should be possible to obtain by a direct extension of our approach.

Non-uniform and non-independent character distributions. The LCS problem on random input strings with more general (in particular, non-uniform and/or non-independent) character distributions seems more challenging, since in this case the three-wise independence of cell types does not hold. Such independence has been essential in obtaining our linking equations. Therefore, a generalisation to these types of character distributions seems far from straightforward.

Larger alphabets. The LCS problem on uniformly random input strings over a larger alphabet presents an opposite challenge, since in this case the cells' four-wise (and higher-order) dependencies are much looser than those with the binary alphabet. Again, such dependencies have been essential in obtaining our linking equations, so a generalisation to a larger alphabet seems far from straightforward.

More than two strings. While the LCS problem can be defined just as easily on three or more strings, none of any existing approaches, including ours, seem to be applicable for solving an analogue of the Chvátal–Sankoff problem on more than two strings. In particular, the LCS problem on three input strings does not appear to possess any of the combinatorial structure that is critical for our transposition network-based approach.

REFERENCES

1. A. Abboud, A. Backurs, V. Vassilevska Williams, *Tight hardness results for LCS and other sequence similarity measures*. — In: Proceedings of FOCS (2015), pp. 59–78.
2. K. S. Alexander, *The rate of convergence of the mean length of the longest common subsequence*. — Annal. Appl. Probab. **4**, No. 4 (1994), 1074–1082.

3. C. E. R. Alves, E. N. Cáceres, S. W. Song, *An all-substrings common subsequence algorithm*. — Discrete Appl. Math. **156**, No. 7 (2008), 1025–1035.
4. R. A. Baeza-Yates, R. Gavaldà, G. Navarro, R. Scheihing, *Bounding the expected length of longest common subsequences and forests*. — Theory of Computing Systems **32**, No. 4 (1999), 435–452.
5. K. E. Batchier, *Sorting networks and their applications*. — In: Proceedings of AFIPS, Vol. 32 (1968), pp. 307–314.
6. A. Borodin, I. Corwin, V. Gorin, *Stochastic six-vertex model*. — Duke Math. J. **165**, No. 3 (2016), 563–624.
7. K. Bringmann, M. Künnemann, *Multivariate fine-grained complexity of longest common subsequence*. — In: Proceedings of ACM-SIAM SODA (2018), pp. 1216–1235.
8. B. Bukh, C. Cox, *Periodic words, common subsequences and frogs*. — Annals Appl. Probab. **32**, No. 2 (2022), 1295–1332.
9. B. Bukh, V. Guruswami, J. Hastad, *An improved bound on the fraction of correctable deletions*. — IEEE Transactions on Information Theory **63**, No. 1 (2017), 93–103.
10. R. Bundschuh, *High precision simulations of the longest common subsequence problem*. — European Phys. J. B **22**, No. 4 (2001), 533–541.
11. J. Casse, *Probabilistic cellular automata with general alphabets possessing a Markov chain as an invariant distribution*. — Adv. Appl. Probab. **48**, No. 2 (2016), 369–391.
12. C. C. Chang, H. J. Keisler, *Model Theory*, Vol. 73 of *Studies in Logic and the Foundations of Mathematics*. North-Holland, third edition, 1990.
13. P. Charalampopoulos, P. Gawrychowski, S. Mozes, O. Weimann, *An Almost Optimal Edit Distance Oracle*. — In: Proceedings of ICALP, Vol. 198 of *Leibniz International Proceedings in Informatics*, pages 48:1–48:20, 2021.
14. P. Charalampopoulos, T. Kociumaka, Sh. Mozes, *Dynamic string alignment*. — In: Proceedings of CPM, Vol. 161 of *LIPICs*, pages 9:1–9:13, 2020.
15. V. Chvátal, D. Sankoff, *Longest common subsequences of two random sequences*. — J. Appl. Probab. **12**, No. 2 (1975), 306–315.
16. M. Crochemore, C. S. Iliopoulos, Y. J. Pinzon, J. F. Reid, *A fast and practical bit-vector algorithm for the Longest Common Subsequence problem*. — Inform. Proc. Letters **80** (2001), 279–285.
17. M. Crochemore, G. M. Landau, M. Ziv-Ukelson, *A subquadratic sequence alignment algorithm for unrestricted score matrices*. — SIAM J. Comput. **32** (2003), 1654–1673.
18. V. Dančík, *Expected Length of Longest Common Subsequences*. PhD thesis, University of Warwick, 1994.
19. J. Boutet De Monvel, *Extensive simulations for longest common subsequences: Finite size scaling, a cavity solution, and configuration space properties*. — Europ. Phys. J. B **7**, No. 2 (1999), 293–308.
20. J. G. Deken, *Some limit results for longest common subsequences*. — Discrete Math. **26**, No. 1 (1979), 17–31.
21. Pablo A. Ferrari, *TASEP hydrodynamics using microscopic characteristics*. — Probab. Surveys **15** (2018), 1–27.

22. P. Gawrychowski, *Faster algorithm for computing the edit distance between SLP-compressed strings*. — In: Proceedings of SPIRE, Vol. 7608 of *Lecture Notes in Computer Science*, pages 229–236, 2012.
23. D. Hermelin, G. M. Landau, S. Landau, O. Weimann, *Unified Compression-Based Acceleration of Edit-Distance Computation*. — *Algorithmica* **65**, No. 2 (2013), 339–353.
24. H. Hyvrö, *Mining bit-parallel LCS-length algorithms*. — In: Proceedings of SPIRE, Vol. 10508 of *Lecture Notes in Computer Science*, pages 214–220, 2017.
25. C U Jensen and H Lenzing, *Model Theoretic Algebra*, Vol. 2 of *Algebra, Logic and Applications*. Gordon and Breach Science Publishers, 1989.
26. D. E. Knuth, *The Art of Computer Programming: Sorting and Searching*, Vol. 3. Addison Wesley, 1998.
27. Th. Kriecherbauer, J. Krug, *A pedestrian's view on interacting particle systems, KPZ universality and random matrices*. — *J. Phys. A: Math. Theor.* **43**, No. 40 (2010), 403001.
28. P. Krusche, A. Tiskin, *String comparison by transposition networks*. — In: London Algorithmics 2008 Theory and Practice, Vol. 11 of *Texts in Algorithmics*. College Publications, 2009.
29. B. F. Logan, L. A. Shepp, *A variational problem for random Young tableaux*. — *Adv. Math.* **26**, No. 2 (1977), 206–222.
30. G. S. Lueker, *Improved bounds on the average length of longest common subsequences*. — *J. ACM* **56**, No. 3 (2009), 17:1–17:38.
31. J. Mairesse, I. Marcovici, *Around probabilistic cellular automata*. — *Theor. Comput. Sci.* **559** (2014), 42–72.
32. S. N. Majumdar, S. Nechaev, *Exact asymptotic results for the Bernoulli matching model of sequence alignment*. — *Phys. Review E: Statistical, Nonlinear, and Soft Matter Physics* **72**, No. 2 (2005), 020901.
33. J. Martin, P. Schmidt, *Multi-type TASEP in discrete time*. — *Latin Amer. J. Probab. Math. Statist.* **8** (2011), 303–333.
34. W. J. Masek, M. S. Paterson, *A faster algorithm computing string edit distances*. — *J. Comput. System Sci.* **20**, No. 1 (1980), 18–31.
35. U. Matarazzo, D. Tsur, M. Ziv-Ukelson, *Efficient all path score computations on grid graphs*. — *Theor Comput. Sci.* **525** (2014), 138–149.
36. M. Paterson, V. Dančík, *Longest common subsequences*. — In: Proceedings of MFCS, Vol. 841 of *Lecture Notes in Computer Science*, pages 127–142, 1994.
37. P. A. Pevzner, M. S. Waterman, *Open combinatorial problems in computational molecular biology*. — In: Proceedings of ISTCS, pp. 158–173, 1995.
38. V. B. Priezzhev, G. M. Schütz, *Exact solution of the Bernoulli matching model of sequence alignment*. — *J. Statist. Mech.: Theory and Experiment* **09** (2008), P09007.
39. N. Rajewsky, L. Santen, A. Schadschneider, M. Schreckenberg, *The asymmetric exclusion process: comparison of update procedures*. — *J. Statist. Phys.* **92** (1998), 151–194.
40. D. Romik, *The Surprising Mathematics of Longest Increasing Subsequences*. Cambridge University Press, Cambridge, 2014.

41. H. Rost, *Non-equilibrium behaviour of a many particle process: Density profile and local equilibria*. — Zeitschrift für Wahrscheinlichkeitstheorie und Verwandte Gebiete **58**, No. 1 (1981), 41–53.
42. Y. Sakai, *An Almost Quadratic Time Algorithm for Sparse Spliced Alignment*. — Theory Comput. Systems **48**, No. 1 (2011), 189–210.
43. Y. Sakai, *A substring-substring LCS data structure*. — Theor. Comput. Sci. **753**, No. 2 (2019), 16–34.
44. S. Salsa, *Partial Differential Equations in Action*, Vol. 99 of *UNITEXT*. Springer International Publishing, 2016.
45. M. Schindl, G. Bilardi, *Bounds and Estimates on the Average Edit Distance*. — In: Proceedings of SPIRE (2019), pp. 91–106.
46. J. P. Schmidt, *All highest scoring paths in weighted grid graphs and their application to finding all approximate repeats in strings*. — SIAM J. Comput. **27**, No. 4 (1998), 972–992.
47. J. M. Steele, *An Efron-Stein inequality for nonsymmetric statistics*. — Annals Statist. **14**, No. 2 (1986), 753–758.
48. J. M. Steele, *Probability Theory and Combinatorial Optimization*, Vol. 69 of *CBMS-NSF regional conference series in applied mathematics*. SIAM, 1997.
49. A. Tiskin, *Semi-local longest common subsequences in subquadratic time*. — J. Discrete Algorithms **6**, No. 4 (2008), 570–581.
50. A. Tiskin, *Semi-local string comparison: Algorithmic techniques and applications*. — Math. Comput. Sci. **1**, No. 4 (2008), 571–603.
51. A. Tiskin, *Periodic String Comparison*. — In: Proceedings of CPM, Vol. 5577 of *Lecture Notes in Computer Science*, pages 193–206, 2009.
52. A. Tiskin, *Towards Approximate Matching in Compressed Strings: Local Subsequence Recognition*. — In: Proceedings of CSR, Vol. 6651 of *Lecture Notes in Computer Science*, pages 401–414. 2011.
53. A. Tiskin, *Fast distance multiplication of unit-monge matrices*. — Algorithmica **71** (2015), 859–888.
54. A. Tiskin, *Bounded-length Smith-Waterman alignment*. — In: Proceedings of WABI, Vol. 143 of *Leibniz International Proceedings in Informatics*, pages 16:1–16:12, 2019.
55. A. Tiskin, *Communication vs synchronisation in parallel string comparison*. — In: Proceedings of SPAA, pages 479–489, 2020.
56. A. M. Vershik, S. V. Kerov, *Asymptotics of the Plancherel measure of the symmetric group and the limiting form of Young tableaux*. — Dokl. Akad. Nauk **233**, No. 6 (1977), 1024–1027.
57. R. A. Wagner, M. J. Fischer, *The string-to-string correction problem*. — J. ACM **21**, No. 1 (1974), 168–173.

Department of Mathematics
and Computer Science
St. Petersburg State University;
St. Petersburg Electrotechnical University “LETI”
E-mail: alextiskin@gmail.com

Поступило October 25, 2022