

A TOPOLOGICAL MODEL OF SITE-SPECIFIC RECOMBINATION THAT PREDICTS THE KNOT AND LINK TYPE OF DNA PRODUCTS

KARIN VALENCIA

*Imperial College London, Department of Mathematics
South Kensington Campus, London SW7 2AZ
E-mail: karin.valencia06@imperial.ac.uk*

Abstract. This is a short summary of a topological model of site-specific recombination, a cellular reaction that creates knots and links out of circular double stranded DNA molecules. The model is used to predict and characterise the topology of the products of a reaction on double stranded DNA twist knots. It is shown that all such products fall into a small family of Montesinos knots and links, meaning that the knot and link type of possible products is significantly reduced, thus aiding their experimental identification. We also mention direct applications of the model.

1. Introduction. The central axis of the double helix of DNA molecules can sometimes become circular, knotted or linked (Figure 1) and this can happen naturally or artificially from topological enzymology experiments [2]. While knots and links in DNA are often detrimental to the cell, theoretically and experimentally, they can be used as probes to understand the cellular reactions that create them [1, 2, 8, 16]. These reactions are usually mediated by proteins, specialised machines that work hard to keep the bodies of living organisms functioning.

One of the most notable examples of cellular reactions that involve DNA knots and links is *site-specific recombination*, carried out by proteins called *site-specific recombinases*. This reaction mediates the process of genetic exchange (or altering of the DNA sequence) via a reciprocal exchange between defined DNA sites. Such DNA exchanges

2010 *Mathematics Subject Classification*: 57M25; 92B05.

Key words and phrases: site-specific recombination, site-specific recombinases, tertiary structure of DNA, circular DNA, DNA supercoiling, DNA knots and links, substrate molecule, twist knots, product molecule, Montesinos knots and links, specific sites contain the recognition sequences and the cross-over sites, recombinase complex, synaptic complex, productive synapse, processive recombination, distributive recombination, serine and tyrosine recombinases.

The paper is in final form and no version of it will be published elsewhere.

have important purposes including the development of drug resistance by some bacteria. Furthermore, these proteins are powerful tools for genomic engineering and have wide and important applications in the agricultural and pharmaceutical industries [13]. However, as a by-product, site-specific recombinases can often cause changes to the topology of (the *imaginary central axis* of) a circular double stranded DNA molecule, the *substrate* molecule, by changing its knot type or link type and creating the *product* molecule(s).

In this article we present a mathematical model of site-specific recombination that is used to predict and characterise the topology of all products that can arise from reactions of site-specific recombination on DNA twist knot substrates.

1.1. The biology of DNA. The primary structure of DNA refers to the structure as a polymer molecule, which consists of building blocks called nucleotides. This can be seen in the left-most image of Figure 1. Each nucleotide consists of three parts: a 5-carbon sugar (where each carbon is numbered from 1 to 5), a phosphate group and a nucleobase (one of Guanine (G), Cytosine (C), Adenine (A), Thymine (T)) attached to the 1'-carbon of the sugar. A DNA strand is made by forming a phosphodiester bond between the 3'- and the 5'-carbon atoms of adjacent sugar rings in each nucleotide.

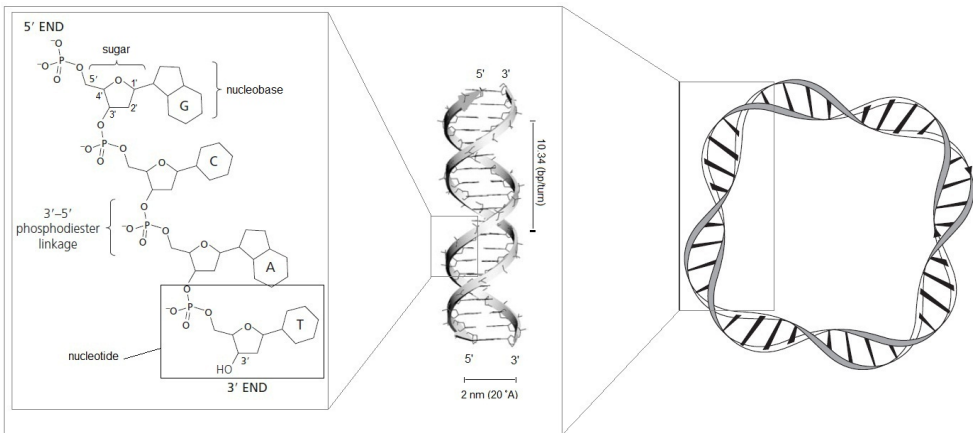


Figure 1. *Left:* The primary structure of DNA.

Middle: A secondary structure of DNA, the *double helix*.

Right: A circular conformation of the tertiary structure of DNA, a relaxed, closed circular DNA. (Modified from [2].)

The secondary structure of DNA addresses the way in which two or more DNA strands are arranged in space. In 1953, James Watson and Francis Crick published a paper determining the famous *double helix* structure (illustrated in the middle image of Figure 1) where each strand follows a right-handed helical path around the central axes of the helix and they are held together by the pairing of bases, via the formation of hydrogen bonds between complementary bases (C with G, A with T). Note that DNA can adopt a variety of conformations. The double helix is also referred to as the B-form DNA and other popular structures include the A-form DNA and the Z-form DNA. However, the B-form

is the most prevalent naturally (*in vivo*) and experimentally (*in vitro*) occurring form, so only this one is considered here.

The *tertiary structure* of DNA describes the conformation of the imaginary central axis of the double helix in space. *Supercoiled* DNA is a conformation adopted by covalently closed circular molecules that are under torsional stress and that, as a result, coil around themselves (leftmost image of Figure 2). Supercoiling is measured mathematically by the difference in linking number of the two backbone strands, between a relaxed molecule and the same molecule under torsional stress. Most naturally occurring DNA is found negatively supercoiled, and therefore supercoiling has many important implications in cell. For a more in-depth discussion see [2, 19] and references therein.

The central axis of the double helix can also be linear, closed circular, supercoiled or knotted or linked in the mathematical sense. Examples of knotted and linked DNA also appear in nature. The mitochondrial DNA of trypanosomes is naturally found as a massive non-trivial link of thousands of components, resembling a medieval mesh [11]. DNA knots and links arise more prevalently as products of topological enzymology experiments (see [9] and references therein). In these experiments small DNA plamids of length between 3 and 5 kilobasepairs (kb) are artificially constructed and acted on by a particular site-specific recombinase that leaves a footprint in the form of changing the topology of the initial molecule. The products of the reaction are then used to probe into the mechanisms of enzymes acting on DNA. Torus knots and links $T(2, m)$ and twist knots $C(2, r)$ are the most commonly occurring knots and links in DNA (Table 1 in [5]).

1.2. The biology of site-specific recombination. The main function of site-specific recombinases is to mediate DNA sequence rearrangements at specific sites. This is carried out by inserting, deleting, and inverting DNA segments.

Minimally, site-specific recombination requires one or two substrate DNA molecules, containing two short specific segments called the *specific sites* and a pair of dimer proteins that mediate the reaction. Each specific site is 30–40 bp in length, contains an inverted pair of *recognition sequences*, that bind one dimer of recombinases, and a point of breakage and rejoining of the DNA, called the *cross-over sites*. Most of the time the cross-over sites are non-palindromic, so they can be assigned an orientation. If the specific sites are on a single DNA molecule, they can either be in direct orientation or in inverted orientation. Depending on the initial arrangement of the specific sites, and recombinase used, site-specific recombination has one of three possible outcomes: integration, excision or inversion of the DNA sequence flanked by the inverted recognition sequences in the specific sites. Larger site-specific recombination systems may also require additional proteins and sites (*accessory proteins* and *enhancer sequences*).

The reaction starts when a pair of recombinases first bind at each of the two recognition sites and, possibly after trapping some number of supercoils, the cross-over sites are brought together and juxtaposed forming the *recombinase complex* (juxtaposed specific sites with proteins attached) and a *synaptic complex* (the whole substrate molecule with specific sites juxtaposed). The DNA is cleaved, exchanged and resealed. Finally, the proteins dissociate, releasing the product molecule and completing the reaction. This is illustrated in Figure 2.

During the intermediate step, once the cross-over sites have been cleaved, multiple rounds of strand exchange can occur before resealing the DNA. This is called *processive recombination*. The entire process of recombination (including releasing and rebinding) can also occur multiple times, either at the same specific sites or at different specific sites. This process is called *distributive recombination*. In this work the term substrate is used to refer specifically to the DNA prior to the first cleavage. Processive recombination is treated as one extended process, given an initial substrate with several intermediate exiting points for the reaction.

Site-specific recombinases can be broadly divided into two subfamilies based, amongst other important factors, on based on their mechanisms of exchanging the DNA sequences: *serine recombinases*, which can perform processive recombination and *tyrosine recombinases* which can not. For a more detailed exposition of the biology of site-specific recombination please refer to [9].

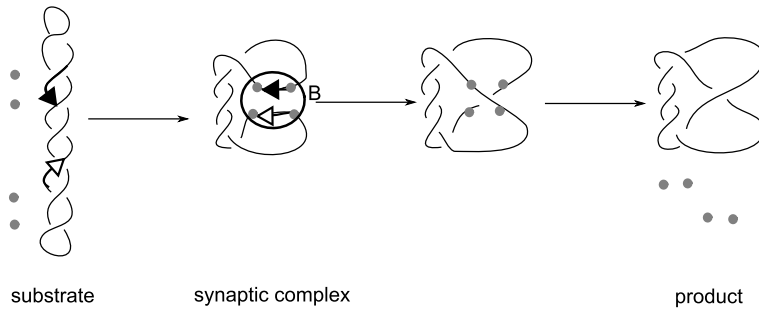


Figure 2. A simple system of site-specific recombination. Note that in every image, the black line represents the imaginary central axis of the double helix of the duplex DNA molecule. The grey dots represent the four recombinase subunits, the bold arrows represent the specific sites and the black circle represents B , the recombinase complex.

1.3. Examples of other topological tools that are useful in the investigation of site-specific recombination. A variety of topological tools have been used for analysing enzyme mechanisms and product knots and links of site-specific recombination (and other important reactions yielding DNA knots and links).

- *Linking number* [3] — used to study the structure of negatively supercoiled circular DNA in solution.
- *Schubert's classification of 4-plats* [20] — used to study interwinding in linked and knotted DNA.
- *Jones polynomial* [21] — used to work out a relationship between the polynomials associated with the substrate molecule and the product molecules obtained by site-specific recombination.
- *Tangle model* [8] — uses results in Dehn surgery to reveal possible mechanisms of recombination reactions that change the topology of an unknotted substrate molecule to a 2-bridge knot or link.

- *Band surgery* [16] — used to characterise the mechanism of unlinking torus links, mediated by Cre recombinase.
- *Rational sub tangle replacement* [1] — similar idea to the tangle model, but for a reaction that does not change the topology of the molecule.

1.4. Definitions and notation. A *twist knot*, denoted by $C(\pm 2, r)$, is a double of the unknot with v twists. It has two non-parallel rows of crossings, one with ± 2 crossings and the other with $|r| \geq 1$. A *torus knot or link*, denoted by $T(2, m)$, has one row of m crossings. It is a knot if m is odd, and a link if m is even.

Let J denote the substrate $C(2, v)$, and B denote the smallest region containing the four bound recombinase molecules and the two juxtaposed cross-over sites. B is a topological ball. It is assumed that for site-specific recombinases that utilise enhancer sequences and/or accessory proteins, these are sequestered from the recombinase complex, and call it a *productive synapse*. Figure 3 illustrates recombinase complexes that are, and that are not productive synapses.

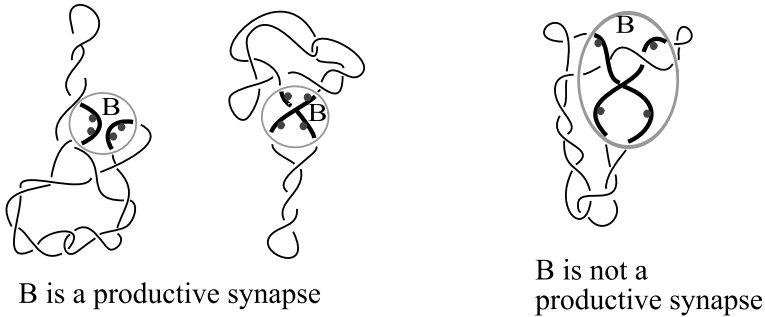


Figure 3. *Left:* Two examples of a productive synapse.
Right: A recombinase complex that is not a productive synapse.

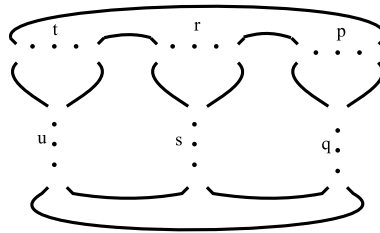


Figure 4. $F(p, q, r, s, t, u)$ is a small family of Montesinos knots and links. All predicted product knots and links fall within this family. Note that all not knots and links in this family arise as products of recombination on twist knots (for example, this family has links with up to three components, but three component links cannot arise as products of site-specific recombination on knots).

The main result of the model exposed below shows that all possible product knots and links of the reaction fall within family $F(p, q, r, s, t, u)$ of knots and links illustrated in

Figure 4, which can be obtained from the numerator closure of a Montesinos tangle that is a sum of three rational tangles ($\frac{p}{p+q} + \frac{r}{r+s} + \frac{t}{t+u}$). The variables p, q, r, s, t, u describe the number of crossings between the central axis of two DNA duplexes in that particular row or column of crossings. These variables can be positive, negative or zero. t, r, p can take horizontal and vertical zero crossings but only horizontal non-zero crossings. u, s, q can only take (both zero and non-zero) vertical crossings.

1.5. Motivation for our model. As mentioned before, site-specific recombination can change the topology of closed circular DNA. Knowing the precise nature of DNA knots and links that arise can help understand details of the mechanisms of the proteins that lead to topological changes in the DNA substrates, either by carrying out topological enzymology experiments, or using topological tools such as those mentioned in Section 1.3. In this work, given three biologically reasonable assumptions about how the enzymes mediate the reaction, for which experimental evidence is given in [5], knots and links that can be yielded as products of a non-distributive reaction starting with a twist knot substrate are predicted and characterised. Since this model predicts the exact topology of the products, e.g. chirality, it can help experimental biologists restrict the knot and link types of the products that can arise. This is especially important since it is a well known fact that there are more than 17 million knots with minimal crossing number at most 16 [7] and the available experimental tools that identify the topology of DNA molecules sometimes are not enough.

Twist knots are ubiquitous in DNA, both *in vivo* and *in vitro*. Most DNA inside prokaryotic cells are supercoiled, and in the lab most experiments done with site-specific recombinases use small supercoiled circular DNA molecules. This supercoiling promotes strand collision and DNA entanglement. A simple crossing change in such a molecule can result in the knotting of the DNA into twist knots (Figure 5).

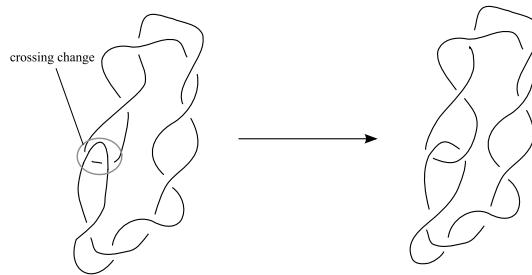


Figure 5. Twist knots are ubiquitous in DNA. It only takes one crossing change to change a supercoiled unknotted DNA molecule to a twist knot.

Together with torus knots and links, twist knots are the most common products of site-specific recombination on substrates that are unknots, unlinks and torus knots and links. This is illustrated in Table 1 in [4]. Also, in multiple rounds of processive and distributive recombination on these substrates, twist knots can become substrates of further rounds of recombination. For example, in [12] site-specific recombination on an

unknot substrate with inverted specific sites, mediated by the recombinase Gin, yields unknots, trefoils, figure-of-eight knots and the knot 5_2 (all twist knots) as products of four rounds of processive recombination on an unknotted substrate. This is illustrated in Figure 6. Also, a composite knot on six crossings, the granny knot $C(-2, 1) \sharp C(-2, 1)$ was analysed to be a product of distributive recombination on two trefoils, each a product from the first round of recombination.

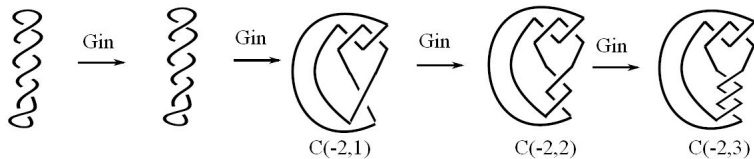


Figure 6. Twist knots arise *in vitro* (and *in vivo*) as products of site-specific recombination on substrates that are unknots, unlinks and torus knots and links. In this example, processive recombination mediated by Gin recombinase on an unknot substrate yields the knots illustrated.

2. Prediction knot and link type of products of site-specific recombination

2.1. Assumptions of our model. This model is based on three assumptions. Evidence that these assumptions are biologically reasonable is given in Section 2 of [5]. Only the mathematical statement for each assumption is presented here. For more details please refer to our papers [6, 18].

The first assumption describes possible conformations of the pre-recombinant juxtaposed cross-over sites.

MATHEMATICAL ASSUMPTION 1. $B \cap J$ consists of two arcs and there is a projection of $B \cap J$ which has at most one crossing between the two arcs, and no crossings within a single arc.

The second assumption describes the synaptic complex.

MATHEMATICAL ASSUMPTION 2. Let $C = \mathbb{R}^3 \setminus B$. J has a spanning surface D such that $D \cap \partial B$ consists of exactly two arcs, the two arcs are co-planar and $D \cap C$ is unknotted relative to ∂B .

Recall that a *spanning surface* of a knot K is a surface F with $\partial F = K$. The sentence $D \cap C$ is unknotted relative to ∂B can be understood intuitively to mean that $\partial D \cap C$ (with ∂B fixed) is not a satellite knot or link.

The third assumption describes possible conformations of the post-recombinant juxtaposed crossover sites, both with a tyrosine recombinase and a serine recombinase.

MATHEMATICAL ASSUMPTION 3 FOR TYROSINE RECOMBINASES. After non-distributive recombination mediated by a tyrosine recombinase, there is a projection of the cross-over sites which has at most one crossing (Figure 7, (1)–(8)).

MATHEMATICAL ASSUMPTION 3 FOR SERINE RECOMBINASES. After each round of processive, non-distributive recombination mediated by a serine recombinase, there is precisely one additional crossing between the cross-over sites (Figure 7, (1)–(10)).

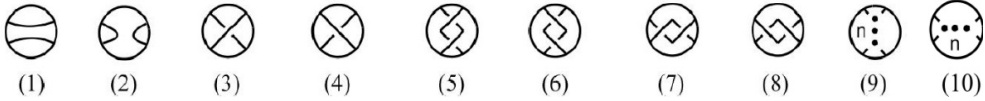


Figure 7. Projections of the post-recombinant forms of the synaptic complex: (1)–(8) for tyrosine recombinases and (1)–(10) for serine recombinases. Note that forms with hooks ((5), (6), (7), (8)) have projection with one crossing but no projections with zero crossings. The row of n vertical and horizontal crossings ((9) and (10) respectively) are products of more than one round of processive recombination.

2.2. Results. Given the three assumptions in the previous section it is shown that products of non-distributive site-specific recombination on twist knots with a tyrosine recombinase (Theorem 1) or with a serine recombinase (Theorem 2) fall within the family of knots and links $F(p, q, r, s, t, u)$, illustrated in Figure 4. Here the proofs are summarised. See [6] for more details.

2.2.1. Products of non-distributive site-specific recombination belong to one family of knots and links

THEOREM 1 (Tyrosine recombinases). *Suppose that Assumptions 1, 2 and 3 hold for a particular tyrosine recombinase-DNA complex. Then the only possible products of (non-distributive) recombination on a twist knot $C(2, v)$ are:*

- torus knots and links $T(2, m)$ for $m = v, v \pm 1, v \pm 2$,
- twist knots $C(2, s)$ for $s = v \pm 1, v \pm 2$,
- clasp knots $C(r, v)$ for $r = \pm 2, \pm 3, 4$,
- the connected sums $T(2, \pm 2) \sharp C(2, v)$,
- a member of the family $F(p, q, r, s, t, u)$ with $r = 2, |t| \leq 2, p = 0$.

THEOREM 2 (Serine recombinases). *Suppose that Assumptions 1, 2 and 3 hold for a particular serine recombinase-DNA complex. Then the only possible products of n rounds of processive (non-distributive) recombination on a twist knot $C(2, v)$ are:*

- torus knots $T(2, v \pm n)$,
- twist knots $C(2, s)$ for $s = v, v \pm n$,
- clasp knots $C(r, v)$ for $r = \pm n, \pm n + 2$,
- a connected sum $T(2, \pm n) \sharp C(2, v)$,
- a member of the family $F(p, q, r, s, t, u)$ with $r = 2, t = \pm n$ and $p = 0$.

Idea of the proof. Step 1. Use Assumption 2 and topological arguments involving spanning surfaces of twist knots to prove that the synaptic complex can only take one of five possible forms illustrated in Figure 8.

Step 2. Use Assumptions 1 and 3 to find the pre-recombinant and post-recombinant forms of the recombinase complex B for both tyrosine recombinases and serine recombinases.

Step 3. Glue the post-recombinant forms of B to each of the forms of the synaptic complex, yielding the predicted knot and link types of the products of recombination.

Step 4. Show that each of these products falls within $F(p, q, r, s, t, u)$. ■

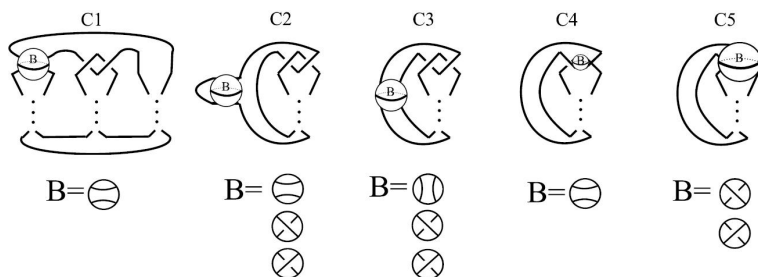


Figure 8. *Top row:* possible forms of the synaptic complex, according to the assumptions. *Bottom row:* the corresponding pre-recombinant recombinase-complex B .

NOTE. All the products predicted in these two theorems fall within family $F(p, q, r, s, t, u)$. However, not all knots and links in family $F(p, q, r, s, t, u)$ are predicted to arise as products of recombination. In particular, $F(p, q, r, s, t, u)$ contains links with up to three components. However, it is impossible to yield a three component link from recombination on a knot.

NOTE. Theorems 1 and 2 distinguish between the chirality of the product DNA molecules, since using our model we can work out the *exact conformation of all possible products* of site-specific recombination starting with a particular twist knot substrate and site-specific recombinase. For example, starting with the twist knot substrate $C(2, -1)$, which is more commonly known as a negative trefoil, then according to the model, site-specific recombination mediated by a tyrosine recombinase can yield $T(2, -5)$, which is also the negative 5_1 knot, but can never yield $T(2, +5) = (+)5_1$.

2.2.2. *Characterisation of products of distributive recombination*

COROLLARY 1. *Any products whose knot or link type is not listed in Theorems 1 and 2 must arise from distributive recombination.*

2.2.3. *The growth of product knots and links is proportional to n^5*

THEOREM 3. *The number of putative knots and links resulting from site-specific recombination on a substrate that is the twist knot $C(2, v)$ with minimal crossing number equal to n grows linearly with n^5 .*

Proof. *Step 1.* Note that while the knots and links in $F(p, q, r, s, t, u)$ have at most six non-adjacent rows containing p, q, r, s, t, u signed crossings respectively, it does not follow that the minimal crossing number of such a knot or link is $|p| + |q| + |r| + |s| + |t| + |u|$.

For example, if the knot or link does not admit, a reduced alternating diagram it is quite possible that the number of crossings can be significantly reduced. So, a priori, there is no reason to believe that the number of knots and links in this product family should grow linearly with n^5 .

Step 2. Find all distinct, non-trivial subfamilies of knots and links in $F(p, q, r, s, t, u)$.

Step 3. Find an upper bound on the number of knots and links in each subfamily as a function of its minimal crossing number. The subfamilies fall into four cases:

Reduced alternating. By Murasugi [15] and Thistlethwaite [17], these knots and links have $n = |p| + |q| + |r| + |s| + |t| + |u|$.

Reduced Montesinos. By Lickorish and Thistlethwaite [14], these knots and links have $n = |p| + |q| + |r| + |s| + |t| + |u|$.

Hara–Yamamoto. By Hara–Yamamoto [10], knots and links K that are *Hara–Yamamoto* cannot be isotoped to a diagram that is either reduced alternating or reduced Montesinos, but have $n = |p| + |q| + |r| + |s| + |t| + |u|$.

Other. Knots and links with projections that can be isotoped to be either reduced alternating or reduced Montesinos have $n \neq |p| + |q| + |r| + |s| + |t| + |u|$. ■

Since the number of prime knots and links (links with up to two components and counting chiral pairs separately) with minimal crossing number n grows exponentially as a function of n [7], Theorem 3 says that the total number of product knots and links in $F(p, q, r, s, t, u)$ with minimal crossing number $= n$ grows linearly with n^5 . Hence, the calculation n^5/e^n gives the proportion of all knots and links which are putative recombination products and as n increases, n^5/e^n decreases exponentially rapidly to zero.

2.3. Applications. A more in detail discussion of applications of this model can be found in [18]. There is also an explicit algorithm for using the model for very specific systems.

These applications fall into four broad categories: *Application 1:* theorem model can help determine the order of products of processive recombination. *Application 2:* in the common situations where the products of site-specific recombination have minimal crossing number one more than the minimal crossing number of the substrate, theorem model can help reduce the number of possibilities for these products. *Application 3:* theorem model can help characterise previously uncharacterised data. *Application 4:* theorem model can help distinguish between products of processive and distributive recombination.

References

- [1] K. L. Baker, D. Buck, *The classification of rational subtangle replacements between rational tangles*, *Algebr. Geom. Topol.* 13 (2013), 1413–1463.
- [2] A. D. Bates, A. Maxwell, *DNA Topology*, Oxford University Press, Oxford, 2005.
- [3] T. C. Boles, J. H. White, N. R. Cozzarelli, *Structure of plectonemically supercoiled DNA*, *J. Molecular Biol.* 213 (1990), 931–951.

- [4] D. Buck, E. Flapan, *Topological characterization of knots and links arising from site-specific recombination*, J. Phys. A 40 (2007), 12377–12395.
- [5] D. Buck, E. Flapan, *Predicting knot or catenane type of site-specific recombination products*, J. Molecular Biol. 374 (2007), 1186–1199.
- [6] D. Buck, K. Valencia, *Characterization of knots and links arising from site-specific recombination on twist knots*, J. Phys. A 44 (2011), 045002, 36 pp.
- [7] C. Ernst, D. W. Sumners, *The growth in the number of prime knots*, Math. Proc. Cambridge Philos. Soc. 102 (1987), 303–315.
- [8] C. Ernst, D. W. Sumners, *A calculus for rational tangles: applications to DNA recombination*, Math. Proc. Cambridge Philos. Soc. 108 (1990), 489–515.
- [9] N. D. F. Grindley, K. L. Whiteson, P. A. Rice, *Mechanisms of site-specific recombination*, Annual Rev. Biochem. 75 (2006), 567–605.
- [10] M. Hara, M. Yamamoto, *Some links with nonadequate minimal-crossing number*, Math. Proc. Cambridge Philos. Soc. 111 (1992), 283–289.
- [11] R. E. Jensen, P. T. Englund, *Network news: the replication of kinetoplast DNA*, Annual Rev. Microbiol. 66 (2012), 473–491.
- [12] R. Kanaar, A. Klippel, E. Shekhtman, J. M. Dungan, R. Kahmann, N. R. Cozzarelli, *Processive recombination by the phage Mu Gin system: implications for the mechanisms of DNA strand exchange, DNA site alignment, and enhancer action*, Cell 62 (1990), 353–366.
- [13] A. F. Kolb, *Genome engineering using site-specific recombinases*, Cloning and Stem Cells 4 (2002), 65–80.
- [14] W. B. R. Lickorish, M. B. Thistlethwaite, *Some links with nontrivial polynomials and their crossing-numbers*, Comment. Math. Helv. 63 (1988), 527–539.
- [15] K. Murasugi, *Jones polynomials and classical conjectures in knot theory*, Topology 26 (1987), 187–194.
- [16] K. Shimokawa, K. Ishihara, M. Vazquez, *Tangle analysis of DNA unlinking by the Xer/FtsK system*, Bussei Kenkyu 92 (2009), 89–92.
- [17] M. B. Thistlethwaite, *A spanning tree expansion of the Jones polynomial*, Topology 26 (1987), 297–309.
- [18] K. Valencia, D. Buck, *Predicting knot and catenanes type of site-specific recombination products of twist knot substrates*, J. Molecular Biol. 411 (2011), 350–367.
- [19] J. C. Wang, *Untangling the Double Helix. DNA Entanglement and the Action of the DNA Topoisomerases*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, 2009.
- [20] J. H. White, N. R. Cozzarelli, *A simple topological method for describing stereoisomers of DNA catenanes and knots*, Proc. Nat. Acad. Sci. U.S.A. 81 (1984), 3322–3326.
- [21] J. H. White, K. C. Millett, N. R. Cozzarelli, *Description of the topological entanglement of DNA catenanes and knots by a powerful method involving strand passage and recombination*, J. Molecular Biol. 197 (1987), 585–603.

