

Relatively semi-conservative replication and a folded slippage model for simple sequence repeats

Hongxi Zhang^{a,1}, Douyue Li^{a,1}, Xiangyan Zhao^{a,1}, Saichao Pan^{a,1}, Xiaolong Wu^a, Shan Peng^a, Hanrou Huang^a, Ruixue Shi^a and Zhongyang Tan^{a,*}

^aBioinformatics Center, College of Biology, Hunan University, Changsha 410082, China

¹Co-first author

*Corresponding author: Zhongyang Tan Email: zhongyangtan@yeah.net

Abstract

Simple sequence repeats (SSRs) are found ubiquitously in almost all genome, and their formation mechanism is ambiguous yet. Here, the SSRs were analyzed in 55 randomly selected segments of genomes from a fairly wide range of species, with introducing more open standard for extensively mining repeats. A high percentage of repeats were discovered in these segments, which is inconsistent with the current theory suggested that repeats tend to disappear over long-term evolution. Therefore, a mechanism is most probably responsible for continually producing repeats during replication to balance continuous repeat disappearance, which may makes the replicating process relatively semi-conservative. To improve the current straight-line slippage model, we proposed a folded slippage model involving the geometric space of nucleotides and hydrogen bond stability to explain the high-percent SSR occurrence, which can describe SSR expansion and contraction more reasonably. And analysis of external forces in the folding template strands suggested that the microsatellites tend to expand than contract. Our research may provide implements for contributions of microsatellites to genome evolution and complement semi-conservative replication.

Introduction

Simple sequence repeats (SSRs), also referred as microsatellites, have attracted increasingly

great interests in recent decades (Chen et al, 2010; Ellegren, 2004; Mandal et al, 2019; Morgante et al, 2002; Vincens et al, 2009a; Zhao et al, 2012), and have been widely analyzed in the genome sequences of eukaryotic prokaryotic and also viral genomes (Ellegren, 2004; Lin & Kussell, 2012; Morgante et al, 2002; Zhao et al, 2012). SSRs are the most variable genomic sequences, which tend to appear frequent variations in repeat-unit number instead of nucleotide substitution. And it may be a critical power accelerate the genomic evolution (Ellegren, 2004; Li et al, 2004), have roles associate with the host-adaptation and pathogenicity (Hood et al, 1996; Li et al, 2004), be relevant with the expression of genes and activity of promoters (Hannan, 2018; Vincens et al, 2009a), have relationship with many genetic diseases (Jain & Vale, 2017; Macdonald et al, 1993; Mirkin, 2007), and be observed with microsatellite instability (MSI) in many type of cancers (Bailey et al, 2018; Chan et al, 2019; Helleday et al, 2014; Kim et al, 2013).

Though SSRs have been comprehensively researched, there is actually no precise definition or wide-convincing standard for the extraction of SSRs all the time, which is usually based on setting the minimum numbers of the iterations for the mononucleotide to hexanucleotide SSRs based on empirical criterion (Chen et al, 2010; Ellegren, 2004; Li et al, 2004; Zhao et al, 2012). Majority of previous studies showed more interesting into the relatively longer repetitive sequences (Benson, 1999; Kelkar et al, 2010; Tian et al, 2011), and most studies usually used the threshold of 6, 3, 3, 3, 3, 3 for extracting mono- to hexanucleotide SSRs (Chen et al, 2011; George et al, 2012; Rajendrakumar et al, 2007; Zhao et al, 2011), while the very short repeat-motifs with smaller iterations were almost excluded, causing the neglect of their important significance (Fungtammasan et al, 2015; Hunt et al, 2016; Schmutz et al, 2014; Teh et al, 2017). In this work, the selected SSRs were extensively extracted with a wider extracting standard for extensive repeat-motif grabbing to investigate the essential occurrences of SSRs.

It is widely accepted that DNA slippage is thought to be the primary mechanism for driving microsatellites expansion or contraction, however, slippage involves DNA polymerase pausing, dissociation and re-association (Ellegren, 2004; Gadgil et al, 2016; Viguera et al, 2001), which may help to understand the expansion and contraction of long repeat sequences; it seems difficult to

explain the remain of high percentage of short repeat sequences, and therefore, it is necessary to improve the slippage model more explicit to explain the generation of large amounts of short repeat sequences (Garcia-Diaz et al, 2006; Huang et al, 2017; Lai & Sun, 2003; Schlötterer & Tautz, 1992). It was suggested that the SSRs are most possibly born in the process of replication (Ellegren, 2004); replication is considered to be exactly semi-conservative with that the number of nucleotides in replicating chain is be precisely equal to that in template chain, and the replicating DNA molecule was shown as a straight molecule in vitro (Watson & Crick, 1953a; Watson & Crick, 1953b). Though it is well known that the DNA molecule is highly bent and packed in a super helix state within the nucleus, the replicating DNA molecule was also believed to be dragged to a straight molecule by the polymerase complex in vivo (Bell, 2011; Costa et al, 2011; Doublé et al, 1998; Kiefer et al, 1998). But there are a lot of environmental elements inside the nucleus which may disturb the polymerase complex, and these disturbances sometimes may affect the dragged straight DNA molecule returning to some extent of bent. The bent replicating DNA molecule is possibly related to the polymerase slippage for the occurrence of short SSRs. Here, we calculated the bent replicating DNA molecule with strictly considering the geometric space, the relationship between the phosphodiester bond and hydrogen bond, and also the stability of paired nucleotides; and proposed a folded replication slippage model for explaining repeats occurrence, which seems more reasonable to explain the remaining of high percentage short repeats in genomes, and also to explain the frequent microsatellite expansion and contraction. This work may also put forward some constructive suggestions for complementing the theory of semi-conservative replication.

Results and discussion

Genomes tend to produce short repeats

We analyzed 55 randomly-selected reported segment sequences covering from animal, plant, fungus, protist, bacteria, archaea and viruses (Table S1). The SSRs were extracted from all these segment sequences by using a threshold with minimum length of 3 base pairs or nucleotides. Though 2 iteration of di-, tri-, tetra-, penta- and hexa- nucleotide repeat sequence are usually ignored in most

previous studies (Ellegren, 2004; Hunt et al, 2016; Schmutz et al, 2014; Teh et al, 2017; Zhao et al, 2012), we found they occurred in a very large number. It is difficult to consider them just as random sequences but not repetitive sequences, and it is also inappropriate to consider the iteration of 3 to 5 of mononucleotide repeats just as random sequences. Therefore, the threshold was set at 3, 2, 2, 2, 2, 2 in this study for exploring more comprehensive occurrence of SSRs, which could grab shorter simple repeats that never analyzed before, and another two thresholds were used to analyze these sequences for comparison. To test whether the SSRs under this threshold are random, we generated 55 mimic sequences with same size and nucleotide composition to the corresponding 55 reported sequences.

The analyzed data showed that the reported segment sequences are averagely 44.4% constituted, with SSRs, ranging from 36.4% to 60.0% under the new threshold (Fig 1A, Table S1). And comparing analysis also show the SSR content of these segments with average of 18.8% and 5.0%. These results indicate that all these segments remained high content of SSRs, because all these segments are randomly selected from their genomes, suggesting that the remaining high content of short SSRs is a general feature of all organism genomes after long time evolution, and also suggesting that few formerly well-studied repeats may only stand for the proverbial tip of the iceberg (Chen et al, 2010; Ellegren, 2004; Lin & Kussell, 2012; Morgante et al, 2002; Zhao et al, 2012). The null hypothesis test demonstrated that the percentages of SSRs in the generated segments are all lower than those in the reported segments, indicating that the high percentages of short SSRs are not randomly remained in all reported segments.

Though the evolutionary mechanism of nucleotide sequences is still hotly debated by evolutionist, it is widely accepted that the genomic sequences are continually mutating forever; and the neutral molecular evolution and molecular clock theory suggested that the nucleotide substitution is constant over the evolution time; the thermodynamics in biology states that an isolated system will always tend to disorder (Bharadwaj et al, 2006; Kimura, 1977; Kimura, 1979; Margoliash, 1963; Zuckerkandl & Pauling, 1962; Zuckerkandl & Pauling, 1965). As the microsatellites are indeed ordered sequences, according to the former stated theories, the ordered repeats possibly tend to

mutate into disordered sequences in the long evolutionary history without any selective pressure. Therefore, the repeat sequences should tend to disappear in genomes in the long evolution history. However, the remaining high percentage of SSRs in genomes is contradicted with the ideas of repetitive sequences tend to become no repetitive sequences. Thus, it can be inferred that there is most probably a mechanism for continually produce repeats to balance continuous repeat disappearance, and be responsible for the remaining of high percentage of short repeat sequences in genomes (Fig 1B).

Furthermore, the SSRs of small iteration numbers were observed to occur largely more than those of large iteration numbers in all analyzed segments (Table 1, Table S2), and this observation indicated that the SSRs of small iteration numbers maybe the basis for forming the SSRs of large iteration numbers, otherwise, it should be that the SSRs of large iteration numbers possibly are remained in higher percent level than or at least almost same level to the SSRs of small iteration numbers. Some of the longer SSRs also possibly mutate into short SSRs by contraction and point mutation as debated by many evolutionists (Ellegren, 2004; Kelkar et al, 2011; Mirkin, 2007), and these debates are possible because of that most of short repeats were not considered in their statistics; our observations generally suggested that most of the longer SSRs possibly evolved from short SSRs by expansion. So, the genomes possibly tend to produce short repeats by a continual repeat producing mechanism with the possibility of expansion a little more than that of contraction.

Relatively semi-conservative replication

It is well known that each base pair of DNA is one-to-one correspondence without other extra residue during replication in the double-helix model (Watson & Crick, 1953a; Watson & Crick, 1953b). And Meselson and Stahl have verified the replication of DNA chains is semi-conservative by the sedimentation techniques based on the diversity differential of DNA with different isotopes, also implicating that the number of nucleotides in replicating strand is consistent with that in template strand while processing complete replication (Meselson & Stahl, 1958). However, if the remained high percentage of short repeats is produced during replication process as described above, it certainly makes the base numbers of replication strand to be unequal to those of template strand, with

one or several nucleotides/motifs being repeated and more than that in template strand. In vitro experiments also revealed the presence of repeats during DNA replication, and the nascent replication chain has a base increase (Doublié et al, 1998; Fungtammasan et al, 2015; Fungtammasan et al, 2016; Kiefer et al, 1998). And in this case, the replication process is possibly relatively semi-conservative and could be described as the following formula:

$$N_i = \text{int}[N_0(1+f_1\lambda_1)(1+f_2\lambda_2)\dots(1+f_i\lambda_i)] \quad (1)$$

$$\Delta N_i = N_i - N_{i-1} = \text{int}[N_0 f_i \lambda_i (1+f_1\lambda_1)(1+f_2\lambda_2)\dots(1+f_{i-1}\lambda_{i-1})] \geq 0 \quad (2)$$

N_0 : The number of nucleotides in the initial template strand;

N_i : The number of nucleotides in the replicating strand during No. i round replication;

$\text{int}[]$: Round the value to the lower integer;

ΔN_i : The difference for the number of nucleotides between N_i and N_{i-1} ;

λ_i ($\lambda_i \rightarrow 0$): The coefficient of occurring repeats during No. i round replication, and is most probably an infinitesimal with relating to the possibility of repeat sequence occurrence;

f_i ($0 \leq f_i \leq 1$): The fixation coefficient of repeat sequences during No. i round replication.

In general, the number of nucleotides in replicating strand is usually detected to be exactly equal to that in template strand, which is possibly because of the observed template strand being too short, for example, the total number of nucleotides in the initial template strand for stable PCR is up to two to three thousand nucleotides, in this case, we suppose $N_0 = 3000$, $\lambda_1 = 10^{-5}$, $f_1 = 1$, then the value of ΔN_1 will be 0 according to the formula (2), and therefore, $N_1 = N_0$, causing the replicating strand to be no longer (or no shorter) than template strand, and the discovery of new-born repeat is unavailable; however, when the observed strand is long enough, then ΔN_i is able to be larger than 1 at least, and it can be found that the number of nucleotides in replicating strand is different from that in template strand, for instance, we suppose $N_0 = 10^6$, $\lambda_1 = 10^{-5}$, $f_1 = 1$, then the value of ΔN_1 will be 10, in this case,

the replicating strands probably have 10 nucleotides (or repeat-motifs) more than template strand do after this replication. Thus, the increased number of nucleotides may represent newly occurred repeat sequences.

The occurrence of SSRs will possibly encounter selective pressure, though it may be different in coding or non-coding regions, then, we use f_i representing the fixation possibility of the newly born repeats facing with the selective pressure. The $f_i = 0$ when the occurrences of new repeats are the lethal mutations and unable fixation in the organism, or may be excluded by DNA repair system (Jeggo et al, 2015; Mandal et al, 2019). The fixation coefficient is $0 < f_i < 1$ when the new SSRs are the deleterious but fixed in the genome within alive individuals, like Huntington's disease (Macdonald et al, 1993). While the occurrences of new SSRs are the neutral mutations, the fixation coefficient should be $0 \leq f_i \leq 1$, and they are fixed or excluded depending on genetic drift. And the f_i of beneficial mutations is 1, representing that the new SSRs may help the organism surviving. Therefore, the remaining high percentage of short repeats suggests that the replicating process possibly produce short repeat sequences frequently which may be fixed neutrally, beneficially, or deleteriously with diseases, and also suggests that the replication may be relatively semi-conservative.

Folded slippage model

The nucleotide chains of various species tend to produce simple repeats during replication, and thus cause the number of the nucleotides in replication strand possible to be different from template strand after replication as discussed above. Moreover, how did simple repeats actually originate from is still a key argument topic (Ellegren, 2004; Kelkar et al, 2011; Torresen et al, 2019). The widely accepted mechanism of occurring SSRs is the replication slippage model, which is possibly easy to explain the expansion and contraction of longer SSRs, but possibly difficult to explain the much amounts of short repeats expansion and contraction. And the current slippage model is indeed a straight template strand model, without considering that the space is required for nucleotide base and also phosphodiester bonds are much stronger than hydrogen bond (Fig 2A) (Gao et al, 2004; Heyrovska, 2006), and also without considering what is the force to drive the replicate strand slippage. The straight replication slippage model has not given any clear suggestion, and it suggests

that the SSRs possibly occurred by slippage occasionally (Gemayel et al, 2010; Leclercq et al, 2010; Mirkin, 2007; Ohshima & Wells, 1997). Actually, there are about 33 atoms in a nucleotide (A: 33, T: 33, G: 34, C: 31) (Alberts et al, 2002), and of course the nucleotide base need a certain space in nature. According to previous reports, we simplified a nucleotide space into an intuitive plane model, whose length is about 0.489 nm (length = (distance between the double helix 1.08 - Hydrogen bond length 0.102) / 2), and with a width of 0.34 nm which is the distance between each pair of bases (Fig 2A) (Gao et al, 2004; Heyrovská, 2006; Wang, 1993). We reconstructed the linear replication slippage model with a CAD geometric calculation by considering the space of bases (Fig 2B, Fig S1); if the slippage bubble has enough geometric space to accommodate the repeat bases, the phosphodiester bond should be elongated far more than 0.34 nm, while the phosphodiester bonds in DNA is actually much stronger than hydrogen bond (Fig 2A)(Wang, 1993). So it is impossible to form a slippage bubble by a larger elongation of the phosphodiester bonds for accommodating the repeat bases. Therefore, the straight slippage model is very difficult to the occurrence of short repeats, and it is most possibly necessary to improve the slippage model.

Actually there is a fact which is widely ignored in replication slippage studies. The template strands are thought to be straight in all replication models, though it is the truth in general condition. It is also well known that the genomic DNA chains are very long and the space is too narrow in the nucleus (Fig 3A); for example, the total length of human genome is about 2 m (2×10^9 nm), while the diameter of nucleus is beneath 10^5 nm in human cell (Alberts et al, 2002); therefore, the genomic DNA chains are generally highly curved and folded in the nucleus as widely accepted. Indeed, the replicating molecule is believed to be a straight molecule (Bell, 2011; Costa et al, 2011; Doublié et al, 1998; Kiefer et al, 1998), and the replicating enzyme complexes usually straighten the template strand to be straight making the replicating strand well paired to finish the semi-conservative replication process (Costantino et al, 2014; Fragkos et al, 2015; Kiefer et al, 1998). However, there are a lot of environmental factors like temperature, viral proteins or diseases etc., which may disturb the normal works of the enzyme complexes. So, when the replicating enzyme complexes are disturbed by environmental factors, the replicating part DNA molecule may recover to some extent

of curved or folded state, and then the template strand may also be some extent of curved or folded state.

Firstly, we proposed a curved template slippage model. When the curved DNA strand is used as the template strand on inner side, the replication strand is longer than the template strand and can form more nucleotides than the template strand on the outside for during replication process. The replication strand should be longer than template strand, then, is able to provide extra spaces for accommodating the extra repeat bases (Fig 3B). However, it is well known that the links of base pairs mainly depend on 2 types of hydrogen bonds, N—H ... :N and N—H ... :O (Heyrovska, 2006), and the strengths of these hydrogen bonds are negatively correlated to the distance between every base pair; the strength of the hydrogen bond is about 3% of the 3', 5'-phosphodiester bonds (Gao et al, 2004; Griffiths et al, 2000; Luo, 2007; Wang, 1993) (Fig 2A), so the distance between the bases is fixed; even if there is space to form a slippage bubble, the hydrogen bond should be elongated to exceed the threshold of 0.167 nm (Heyrovska, 2006) and should be easy to be broken off in such condition. So, the curved slippage model is able to provide spaces for forming slippage bubble with forming unstable hydrogen bonds double-chain structures (Arm1 and Arm2) at both sides of the slippage bubble (Fig. 3B, Fig S2), indicating that the curved slippage model should be unreasonable.

Then we proposed a folded slippage model. In this case, the folded template strand forms a slippage bubble above the folding site to have sufficient space for accommodating the repeat nucleotides in replication process, the phosphodiester bonds are not elongated, but the bases are well paired with the stable hydrogen bonds at both sides of the slippage bubble (Fig 4). If folding angle is proper, thereby it is most possibly to form a very stable double-stranded folded slippage structure to provide chances for producing repeats, with considering nucleotide geometric spaces and stability of phosphodiester and hydrogen bonds. Actually, there are two conditions of the folded slippage models: When template strand is on the inner side, the repeat unit duplicated to produce new repetitive unit or repeat expansion (Fig 4); and when the template strand is on the outside, the replication strand may make the repetitive sequences to contract (Fig 5); the features of this folded slippage model can easily explain the widely observed microsatellite mutations with expansion and contraction of repeat

units (Ellegren, 2004; Gemayel et al, 2010; Gymrek et al, 2016; Kelkar et al, 2011; Mirkin, 2007). In addition, replication slippage of template strands with different folding angles may result in the expansion or contraction of repeat units with different sizes. When template chains are folded on the inner side at a rotation angle of 18° , 36° , 54° , 72° , 90° and 108° , the replication strands will produce mononucleotide to hexanucleotide repeat expanding respectively (Fig 4). So, it is necessary to break off the number of hydrogen bonds from 2 to 18 without elongating the phosphodiester bond to produce repeats; it suggested that the difficulty of formation repeats from mono- to hexanucleotide is gradually increasing, and also means the occurrence of mono-, di-, tri-, tetra-, penta- and hexanucleotide repeat is gradually decreasing; that is well consistent with our statistic data (Table 1, Table S2). Vice versa, when template chains are folded on the outside at a rotation angle of 18° , 36° , 54° , 72° , 90° and 108° , the replication strands will produce responding repeats contracting respectively (Fig 5). These features are well corresponding to the microsatellites which usually refers to the tandem repeats with repeat units from mono- to hexanucleotides (Ellegren, 2004; Kelkar et al, 2010; Zhao et al, 2011). According to this rule, we also describe the possible folded template slippage models of hepta-, octa-, nona- and decanucleotide repeats (Figs S3 and S4). In fact, the replicating strand must break off at least from 14 to 30 hydrogen bonds to make a folded slippage bubble, the energy to break off so much hydrogen bonds are almost close to energy of phosphodiester bond, then, they are very difficult to occur, and therefore, this is consistent with the observations that such long tandem repetitive sequences are often not very abundant in the genomes (Gemayel et al, 2010; Legendre et al, 2007). The (A_mT_n) repeats growing faster than (G_mC_n) repeats also suggested that the broken number of hydrogen bonds involves in the speed of repeat expansion (Katti et al, 2001; Schlötterer & Tautz, 1992; Sinai et al, 2019; Tian et al, 2011). Although this folded slippage model is just simply described in a plane form, it can still clearly simulate and explain the repeat sequences producing process. We also use the same space size to make the double-helical three-dimensional forms show the folded slippage model more intuitively (Figs 4 and 5), and the precise folding angle in the three-dimensionally double-helical forms and other issues desire further study.

There is enough geometric space in the slippage bubble of the folded template model to

accommodate repeat nucleotides without stretching the phosphodiester bonds, compared with the straight template slippage model. In contrast to the curved template model, the difference in the folded model is that the two sides of the slippage bubble are stably paired, and the Arm1 and Arm2 similar to the straight template replication model are formed at both sides (Figs 4 and 5). The folded model takes full account of the space required by nucleotides, the stability of phosphodiester bonds and the strength comparison between phosphodiester bonds and hydrogen bond, and is easy used to explain microsatellite mutations with repeat unit expansion and contraction. Therefore, we propose that the folded template chain slippage model may be considered as the most reasonable model for explaining repeats production in replicating process, and the folded template strand slippage model may be responsible for the continual producing of repeat sequences and the remaining of high percentage of repeat sequences in genomes.

Microsatellites tend to expand

As stated above, according to the folded slippage model, template chain folding on the inner side may make the replicating chain slippage for repeats expansion, vice versa, the template chain folding on the outside may make the replicating chain slippage for repeats contraction; and it seems that the possibility of repeats expansion and contraction is same. However, there are two manners for the repeat sequences contraction, one is above mentioned the template chain folds on outside, another is also above stated general mutations; the high content of the repeat sequence is still in a stable state in the genome of each species, suggesting that the possibility of repeat expansion should be higher than repeat contraction. And many reports also suggest that there is a higher possibility of repeat expansion than repeat contraction (Fungtammasan et al, 2015; Fungtammasan et al, 2016; Neil et al, 2018).

When the folded template chain slippage was deeply investigated, the replicating straight template DNA chain should return to folded under external forces from the narrow and crowded cell nucleus when the replicating enzyme complexes are disturbed, and usually the replicating enzyme complexes may provide power for balancing the external forces to drag the template DNA molecule straight. Then, we proposed an external force function for template strand returning to folded, and this function may be helpful to explore the probability of expansion and contraction. When the template

strand is on the inner side, the nucleotide bases are outward, and the space of bases at the folded site become wide and loose at outward part; while it is on the outside, the base in the folding position is squeezed inward. Comprehensively considering the small difference of the space of nucleotides at the folded site, it can be easy accepted that the external forces to make template strand folded with bases loose should be smaller than that to be squeezed; therefore, the external force required for the template strand folded on the outside (F^o) is inevitable greater than that (F^i) on the inner side, it can be described as $F^o > F^i$, suggesting that the probability for the template strand folded on the inner side is higher than that on the outside; as our folded slippage model suggested that the repeats tend to expand when the template strand on inner side and contract when the template strand on outside, therefore, the possibility of repeat expansion (P^e) is most possibly higher than that for repeat contraction (P^c), it can be described as $P^e > P^c$ (Fig 6). The SSR studies, like in Huntington disease related locus and myotonic dystrophy type 1 locus, all showed SSR expansion biased (Higham et al, 2012; Larson et al, 2015; Macdonald et al, 1993; Mirkin, 2007; Sznajder & Swanson, 2019), which proving that the expansion and of short SSRs are more frequent than that of contraction.

Thus, according to formula (2):

When the template strand on the outside, repeats tend to contract, so $\lambda^c < 0$,

$$\text{thus, } \Delta N^c = N^c_i - N^c_{i-1} = \text{int}[N_0 f^c_i \lambda^c_i (1 + f^c_1 \lambda^c_1) (1 + f^c_2 \lambda^c_2) \dots (1 + f^c_{i-1} \lambda^c_{i-1})] \leq 0.$$

When the template strand on the inner side, repeats tend to expand, so $\lambda^e > 0$,

$$\text{thus, } \Delta N^e = N^e_j - N^e_{j-1} = \text{int}[N_0 f^e_j \lambda^e_j (1 + f^e_1 \lambda^e_1) (1 + f^e_2 \lambda^e_2) \dots (1 + f^e_{j-1} \lambda^e_{j-1})] \geq 0.$$

The general repeat expansion and contraction can be described as:

$$|\sum \Delta N^e| = |\text{int}[\sum N_0 f^e_j \lambda^e_j (1 + f^e_1 \lambda^e_1) (1 + f^e_2 \lambda^e_2) \dots (1 + f^e_{j-1} \lambda^e_{j-1})]|;$$

$$|\sum \Delta N^c| = |\text{int}[\sum N_0 f^c_i \lambda^c_i (1 + f^c_1 \lambda^c_1) (1 + f^c_2 \lambda^c_2) \dots (1 + f^c_{i-1} \lambda^c_{i-1})]|;$$

$$\sum \Delta N = |\sum \Delta N^e| - |\sum \Delta N^c| = \text{int}[N_0 \sum [f^e_j \lambda^e_j (1 + f^e_1 \lambda^e_1) (1 + f^e_2 \lambda^e_2) \dots (1 + f^e_{j-1} \lambda^e_{j-1})] - [f^c_i \lambda^c_i (1 + f^c_1 \lambda^c_1) (1 + f^c_2 \lambda^c_2) \dots (1 + f^c_{i-1} \lambda^c_{i-1})]]].$$

Because λ was defined as coefficient of occurring repeats, the possibility of repeat expansion (P^e) is positively proportional to λ^e and the possibility of contraction (P^c) is positively proportional to the absolute value of λ^c ($|\lambda^c|$), if we suppose that $f^e = f^c = f$, $i = j$, and as generally $P^e > P^c$, then $\lambda^e > |\lambda^c|$,

$$\text{and also } \sum [\lambda_j^e (1 + f\lambda_{j-1}^e)(1 + f\lambda_{j-2}^e) \dots (1 + f\lambda_{j-i}^e)] \geq \sum [|\lambda_i^c| (1 + f|\lambda_{i-1}^c|)(1 + f|\lambda_{i-2}^c|) \dots (1 + f|\lambda_{i-i}^c|)],$$

$$\text{therefore, } \sum \Delta N = |\sum \Delta N^e| - |\sum \Delta N^c| \geq 0.$$

So, when the external forces for returning the folded template strand were considered, the possibility of repeat expansion should be higher than that of repeat contraction, then the revised formula (2) is also able to explain the remaining of high percentage of short repeats in genomes under a mechanism of continually producing repeats; and this mechanism might result from the folded template chain slippage model, which is possibly responsible for the widely occurring short tandem repeats, also called microsatellites or SSRs in eukaryotic, prokaryotic and also viral genomes. We improved the straight slippage model to folded slippage model by fully considering the geometric spaces of nucleotides base, the relationship between phosphodiester and hydrogen bond and the stability of these bonds. The slippage model showed that the straight replicating template DNA may return to be some extent of folded resulting from disturbed replicating enzyme complexes, and may provide chances for continually producing much amount of short repeats; though the long unit repeats may be related with the former slippage model (Gemayel et al, 2010; Viguera et al, 2001). The easily forming of folded slippage may be also responsible for the widely observed fact that repetitive part of genome is usually evolved hundred or more times than other part with only repeat units expansion and contraction (Giesselmann et al, 2019; Kelkar et al, 2011; Kim et al, 2013; Mandal et al, 2019), though the repeats occurred more in non-coding regions than in coding regions possibly because of different selective pressures (Ellegren, 2004; Gemayel et al, 2010; Mirkin, 2007). Most of new occurring repeats should be lethal mutation and may have been negatively selected to lost; some of new occurring repeats should be deleterious in genomes and responsible for a series of diseases (Arturo et al, 2010; Larson et al, 2015; Sun et al, 2018; Sznajder & Swanson, 2019); many neutral repeat expansions may be lost or fixed with no functions in genomes by genetic drift (Muller

et al, 2014); and some beneficial repeat expansions may promote the emergence of different new properties or functions, that is why the repeat sequences are reported with so many different roles (Gymrek et al, 2016; Hannan, 2018; Hood et al, 1996; Li et al, 2004; Mrazek, 2006; Sinai et al, 2019; Vences et al, 2009b). And the longer repeats might originate from short repeat expansion by the folded template slippage, and the longer genomes possibly evolved from the short genome with related to the continuous repeats producing folded slippage model in the long evolutionary replicating process.

Materials and Methods

Sequences resource

We downloaded 55 genomic sequences from Genbank of a fairly wide range of species that covering animals, plants, fungi, protozoa, bacteria, archaea and viruses. The segments for SSR analysis were randomly selected from different regions of these 55 genomic sequences, which range from 3000 to 96600 bp in length and do not contain any gaps, to verify the widespread distributions of SSRs, as the full genomic sequences are too long.

Repeat extraction

The perfect simple sequence repeats were extracted by Imperfect Microsatellite Extraction Webserver (IMEx-web, <http://imex.cdfd.org.in/IMEX/index.html>) from those 55 randomly-selected reported segments. The minimum iterations for all perfect mono- to hexanucleotide repeats were set at 3, 2, 2, 2, 2, 2 to mine the data more completely in this study, comparing with most researchers setting iterations at relatively higher self-defined values, and 3 iterations for mononucleotide repeats were confined to ensure to be commonly recognized as the SSRs.

Null hypothesis test

We also extracted perfect mono- to hexanucleotide repeats under the above threshold in the sequences that were generated by a program written in C language (Program S1) according to the nucleotide compositions and sizes of those 55 reported segments. Then, the validating test, which can verify that the short SSRs extracted in those 55 reported segments are nonrandom sequences, was

based on the comparison of the SSR percentages in the reported segments and our generated segments.

Model drawing of DNA replication

Different models were drawn to simulate the DNA replication. Normally in straight model, the hydrogen bond length between 2 paired nucleotides is reported to be 0.102 nm and the distance between 2 neighboring nucleotides is 0.34 nm, importantly, owing to the nucleotides occupying almost same space in DNA strands, the space of a nucleotide was simplified into a geometric plane form in this analysis, which was 0.489 nm in length and 0.34 nm in width. Then we applied AutoCAD to draw the straight, curved and folded slippage models according to the strict geometric calculation of the spaces of nucleotides and different strengths between hydrogen bonds and phosphodiester bonds. And the slippage models in helix structure were achieved by Rhino, which is an industrial drawing software.

Availability of data and materials

Supplementary Tables are online at

https://github.com/DooYal/Supplementary-Table-for-submitting-relatively-...-/tree/DooYal-patch-manuscript_folded/supplementary%20tables

References

Alberts B, Johnson A, Lewis J, Raff M, Roberts K, Walter P (2002) *Molecular Biology of the Cell* (4th ed.), New York: Garland Science.

Arturo LC, Cleary JD, Pearson CE (2010) Repeat instability as the basis for human diseases and as a

potential target for therapy. *Nat Rev Mol Cell Biol* **11**: 165-170

Bailey MH, Tokheim C, Porta-Pardo E, Sengupta S, Bertrand D, Weerasinghe A, Colaprico A, Wendl MC, Kim J, Reardon B, Ng PK, Jeong KJ, Cao S, Wang Z, Gao J, Gao Q, Wang F, Liu EM, Mularoni L, Rubio-Perez C, Nagarajan N, Cortes-Ciriano I, Zhou DC, Liang WW, Hess JM, Yellapantula VD, Tamborero D, Gonzalez-Perez A, Suphavilai C, Ko JY, Khurana E, Park PJ, Van Allen EM, Liang H, Group MCW, Cancer Genome Atlas Research N, Lawrence MS, Godzik A, Lopez-Bigas N, Stuart J, Wheeler D, Getz G, Chen K, Lazar AJ, Mills GB, Karchin R, Ding L (2018) Comprehensive Characterization of Cancer Driver Genes and Mutations. *Cell* **173**: 371-385 e318

Bell SD (2011) DNA replication: archaeal oriGINS. *BMC Biol* **9**: 36

Benson G (1999) Tandem repeats finder: a program to analyze DNA sequences. *Nucleic Acids Res* **27**: 573-580

Bharadwaj S, Montazeri R, Haynie DT (2006) Direct determination of the thermodynamics of polyelectrolyte complexation and implications thereof for electrostatic layer-by-layer assembly of multilayer films. *Langmuir* **22**: 6093-6101

Chan EM, Shibue T, McFarland JM, Gaeta B, Ghandi M, Dumont N, Gonzalez A, McPartlan JS, Li TX, Zhang YX, Liu JB, Lazaro JB, Gu PL, Pieltz CG, Apffel A, Ali SO, Deasy R, Keskula P, Ng RWS, Roberts EA, Reznichenko E, Leung L, Alimova M, Schenone M, Islam M, Maruvka YE, Liu Y, Roper J, Raghavan S, Giannakis M, Tseng YY, Nagel ZD, D'Andrea A, Root DE, Boehm JS, Getz G, Chang S, Golub TR, Tsherniak A, Vazquez F, Bass AJ (2019) WRN helicase is a synthetic lethal target in microsatellite unstable cancers. *Nature* **568**: 551-556

Chen M, Tan Z, Zeng G (2011) Microsatellite is an important component of complete hepatitis C virus genomes. *Infect, Genet Evol* **11**: 1646-1654

Chen M, Tan Z, Zeng G, Peng J (2010) Comprehensive analysis of simple sequence repeats in pre-miRNAs. *Mol Biol Evol* **27**: 2227-2232

Costa A, Ilves I, Tamberg N, Petojevic T, Nogales E, Botchan MR, Berger JM (2011) The structural basis for MCM2-7 helicase activation by GINS and Cdc45. *Nat Struct Mol Biol* **18**: 471-477

Costantino L, Sotiriou SK, Rantala JK, Magin S, Mladenov E, Helleday T, Haber JE, Iliakis G, Kallioniemi OP, Halazonetis TD (2014) Break-induced replication repair of damaged forks induces genomic duplications in human cells. *Science* **343**: 88-91

Doubl   S, Tabor S, Long AM, Richardson CC, Ellenberger T (1998) Crystal structure of a bacteriophage T7 DNA replication complex at 2.2   resolution. *Nature* **391**: 251-258

Ellegren H (2004) Microsatellites: Simple sequences with complex evolution. *Nature reviews Genetics* **5**: 435-445

Fragkos M, Ganier O, Coulombe P, Mechali M (2015) DNA replication origin activation in space and time. *Nature Reviews: Molecular Cell Biology* **16**: 360-374

Fungtammasan A, Ananda G, Hile SE, Su MSW, Sun C, Harris R, Medvedev P, Eckert K, Makova KD (2015) Accurate typing of short tandem repeats from genome-wide sequencing data and its applications. *Genome Res* **25**: 736-749

Fungtammasan A, Tomaszewicz M, Campos-Sanchez R, Eckert KA, DeGiorgio M, Makova KD (2016) Reverse Transcription Errors and RNA-DNA Differences at Short Tandem Repeats. *Mol Biol Evol* **33**: 2744-2758

Gadgil R, Barthelemy J, Lewis T, Leffak M (2016) Replication stalling and DNA microsatellite instability. *Biophys Chem* **225**: 38-48

Gao F, Yin C, Yang P (2004) Coordination chemistry mimics of nuclease-activity in the hydrolytic cleavage of phosphodiester bond. *Chinese Sci Bull* **49**: 1667-1680

Garcia-Diaz M, Bebenek K, Krahm JM, Pedersen LC, Kunkel TA (2006) Structural analysis of strand misalignment during DNA synthesis by a human DNA polymerase. *Cell* **124**: 331-342

Gemayel R, Vences MD, Legendre M, Verstrepen KJ (2010) Variable Tandem Repeats Accelerate Evolution of Coding and Regulatory Sequences. *Annu Rev Genet* **44**: 445-477

George B, Alam CM, Jain SK, Sharfuddin C, Chakraborty S (2012) Differential distribution and occurrence of simple sequence repeats in diverse geminivirus genomes. *Virus Genes* **45**: 556-566

Giesselmann P, Brandl B, Raimondeau E, Bowen R, Rohrandt C, Tandon R, Kretzmer H, Assum G, Galonska C, Siebert R, Ammerpohl O, Heron A, Schneider SA, Ladewig J, Koch P, Schuldt BM, Graham JE, Meissner A, Muller FJ (2019) Analysis of short tandem repeat expansions and their methylation state with nanopore sequencing. *Nat Biotechnol* **37**: 1478-1481

Griffiths AJ, Miller JH, Suzuki DT, Lewontin RC (2000) *An Introduction to Genetic Analysis*, 7th edition.

Gymrek M, Willems T, Guilmatre A, Zeng H, Markus B, Georgiev S, Daly MJ, Price AL, Pritchard JK, Sharp AJ, Erlich Y (2016) Abundant contribution of short tandem repeats to gene expression variation in humans. *Nat Genet* **48**: 22-29

Hannan AJ (2018) Tandem repeats mediating genetic plasticity in health and disease. *Nature reviews Genetics* **19**: 286-298

Helleday T, Eshtad S, Nik-Zainal S (2014) Mechanisms underlying mutational signatures in human cancers. *Nature reviews Genetics* **15**: 585-598

Heyrovská R (2006) Dependence of the length of the hydrogen bond on the covalent and cationic radii of hydrogen, and additivity of bonding distances. *Chem Phys Lett* **432**: 348-351

Higham CF, Morales F, Cobbold CA, Haydon DT, Monckton DG (2012) High levels of somatic DNA diversity at the myotonic dystrophy type 1 locus are driven by ultra-frequent expansion and contraction mutations. *Hum Mol Genet* **21**: 2450-2463

Hood DW, Deadman ME, Jennings MP, Bisercic M, Fleischmann RD, Venter JC, Moxon ER (1996) DNA repeats identify novel virulence genes in *Haemophilus influenzae*. *P Natl Acad Sci Usa* **93**: 11121-11125

Huang TY, Chang CK, Kao YF, Chin CH, Ni CW, Hsu HY, Hu NJ, Hsieh LC, Chou SH, Lee IR (2017) Parity-dependent hairpin configurations of repetitive DNA sequence promote slippage associated with DNA expansion. *P Natl Acad Sci Usa* **114**: 9535-9540

Hunt VL, Tsai IJ, Coghlan A, Reid AJ, Holroyd N, Foth BJ, Tracey A, Cotton JA, Stanley EJ, Beasley H, Bennett HM, Brooks K, Harsha B, Kajitani R, Kulkarni A, Harbecke D, Nagayasu E, Nichol S, Ogura Y, Quail MA, Randle N, Xia D, Brattig NW, Soblik H, Ribeiro DM, Sanchez-Flores A, Hayashi T, Itoh T, Denver DR, Grant W, Stoltzfus JD, Lok JB, Murayama H, Wastling J, Streit A, Kikuchi T, Viney M, Berriman M (2016) The genomic basis of parasitism in the *Strongyloides* clade of nematodes. *Nat Genet* **48**: 299-307

Jain A, Vale RD (2017) RNA phase transitions in repeat expansion disorders. *Nature* **546**: 243-247

Jeggo PA, Pearl LH, Carr AM (2015) DNA repair, genome stability and cancer: a historical perspective. *Nature Reviews: Cancer* **16**: 35

Katti MV, Ranjekar PK, Gupta VS (2001) Differential distribution of simple sequence repeats in eukaryotic genome sequences. *Mol Biol Evol* **18**: 1161-1167

Kelkar YD, Eckert KA, Chiaromonte F, Makova KD (2011) A matter of life or death: how microsatellites emerge in and vanish from the human genome. *Genome Res* **21**: 2038-2048

Kelkar YD, Strubczewski N, Hile SE, Chiaromonte F, Eckert KA, Makova KD (2010) What Is a Microsatellite: A Computational and Experimental Definition Based upon Repeat Mutational

Behavior at A/T and GT/AC Repeats. *Genome Biol Evol* **2**: 620-635

Kiefer JR, Mao C, Braman JC, Beese LS (1998) Visualizing DNA replication in a catalytically active *Bacillus* DNA polymerase crystal. *Nature* **391**: 304

Kim TM, Laird PW, Park PJ (2013) The Landscape of Microsatellite Instability in Colorectal and Endometrial Cancer Genomes. *Cell* **155**: 858-868

Kimura M (1977) Preponderance of synonymous changes as evidence for the neutral theory of molecular evolution. *Nature* **267**: 275-276

Kimura M (1979) The neutral theory of molecular evolution. *Sci Am* **241**: 98-100, 102, 108 passim

Lai YL, Sun FZ (2003) The relationship between microsatellite slippage mutation rate and the number of repeat units. *Mol Biol Evol* **20**: 2123-2131

Lamond AI (2002) Molecular Biology of the Cell, 4th edition. *Nature* **417**: 383-383

Larson E, Fyfe I, Morton AJ, Monckton DG (2015) Age-, tissue- and length-dependent bidirectional somatic CAG*CTG repeat instability in an allelic series of R6/2 Huntington disease mice. *Neurobiol Dis* **76**: 98-111

Leclercq S, Rivals E, Jarne P (2010) DNA Slippage Occurs at Microsatellite Loci without Minimal Threshold Length in Humans: A Comparative Genomic Approach. *Genome Biol Evol* **2**: 325-335

Legendre M, Pochet N, Pak T, Verstrepen KJ (2007) Sequence-based estimation of minisatellite and microsatellite repeat variability. *Genome Res* **17**: 1787-1796

Li YC, Korol AB, Fahima T, Nevo E (2004) Microsatellites Within Genes: Structure, Function, and

Evolution. *Mol Biol Evol* **21**: 991-1007

Lin WH, Kussell E (2012) Evolutionary pressures on simple sequence repeats in prokaryotic coding regions. *Nucleic Acids Res* **40**: 2399-2413

Luo YR (2007) *Comprehensive Handbook of Chemical Bond Energies*, Boca Raton, FL: CRC Press.

Macdonald ME, Ambrose CM, Duyao MP, Myers RH, Lin C, Srinidhi L, Barnes G, Taylor SA, James M, Groot N (1993) A novel gene containing a trinucleotide repeat that is expanded and unstable on Huntington's disease chromosomes. *Cell* **72**: 971-983

Mandal R, Samstein RM, Lee KW, Havel JJ, Wang H, Krishna C, Sabio EY, Makarov V, Kuo FS, Bleuca P, Ramaswamy AT, Durham JN, Bartlett B, Ma XX, Srivastava R, Middha S, Zehir A, Hechtman JF, Morris LGT, Weinhold N, Riaz N, Le DT, Diaz LA, Chan TA (2019) CANCER Genetic diversity of tumors with mismatch repair deficiency influences anti-PD-1 immunotherapy response. *Science* **364**: 485-491

Margoliash E (1963) Primary Structure And Evolution Of Cytochrome C. *Proc Natl Acad Sci U S A* **50**: 672-679

Meselson M, Stahl FW (1958) The Replication of DNA in Escherichia Coli. *Proc Natl Acad Sci U S A* **44**: 671-682

Mirkin SM (2007) Expandable DNA repeats and human disease. *Nature* **447**: 932-940

Morgante M, Hanafey M, Powell W (2002) Microsatellites are preferentially associated with nonrepetitive DNA in plant genomes. *Nat Genet* **30**: 194-200

Mrazek J (2006) Analysis of distribution indicates diverse functions of simple sequence repeats in

Mycoplasma genomes. *Mol Biol Evol* **23**: 1370-1385

Muller MJ, Neugeboren BI, Nelson DR, Murray AW (2014) Genetic drift opposes mutualism during spatial population expansion. *Proc Natl Acad Sci U S A* **111**: 1037-1042

Neil AJ, Liang MU, Khristich AN, Shah KA, Mirkin SM (2018) RNA-DNA hybrids promote the expansion of Friedreich's ataxia (GAA)(n) repeats via break-induced replication. *Nucleic Acids Res* **46**: 3487-3497

Ohshima K, Wells RD (1997) Hairpin formation during DNA synthesis primer realignment in vitro in triplet repeat sequences from human hereditary disease genes. *The Journal of biological chemistry* **272**: 16798-11806

Rajendrakumar P, Biswal AK, Balachandran SM, Srinivasarao K, Sundaram RM (2007) Simple sequence repeats in organellar genomes of rice: frequency and distribution in genic and intergenic regions. *Bioinformatics* **23**: 1-4

Schlötterer C, Tautz D (1992) Slippage synthesis of simple sequence DNA. *Nucleic Acids Res* **20**: 211-215

Schmutz J, McClean PE, Mamidi S, Wu GA, Cannon SB, Grimwood J, Jenkins J, Shu S, Song Q, Chavarro C, Torres-Torres M, Geffroy V, Moghaddam SM, Gao D, Abernathy B, Barry K, Blair M, Brick MA, Chovatia M, Gepts P, Goodstein DM, Gonzales M, Hellsten U, Hyten DL, Jia G, Kelly JD, Kudrna D, Lee R, Richard MM, Miklas PN, Osorno JM, Rodrigues J, Thareau V, Urrea CA, Wang M, Yu Y, Zhang M, Wing RA, Cregan PB, Rokhsar DS, Jackson SA (2014) A reference genome for common bean and genome-wide analysis of dual domestications. *Nat Genet* **46**: 707-713

Sinai MIT, Salamon A, Stanleigh N, Goldberg T, Weiss A, Wang YH, Kerem B (2019) AT-dinucleotide rich sequences drive fragile site formation. *Nucleic Acids Res* **47**: 9685-9695

Sun JH, Zhou LD, Emerson DJ, Phyto SA, Titus KR, Gong WF, Gilgenast TG, Beagan JA, Davidson BL, Tassone F, Phillips-Cremens JE (2018) Disease-Associated Short Tandem Repeats Co-localize

with Chromatin Domain Boundaries. *Cell* **175**: 224-238

Sznajder LJ, Swanson MS (2019) Short Tandem Repeat Expansions and RNA-Mediated Pathogenesis in Myotonic Dystrophy. *Int J Mol Sci* **20**

Teh BT, Lim K, Yong CH, Ng CCY, Rao SR, Rajasegaran V, Lim WK, Ong CK, Chan K, Cheng VKY, Soh PS, Swarup S, Rozen SG, Nagarajan N, Tan P (2017) The draft genome of tropical fruit durian (*Durio zibethinus*). *Nat Genet* **49**: 1633-1641

Tian XJ, Strassmann JE, Queller DC (2011) Genome Nucleotide Composition Shapes Variation in Simple Sequence Repeats. *Mol Biol Evol* **28**: 899-909

Torresen OK, Star B, Mier P, Andrade-Navarro MA, Bateman A, Jarnot P, Gruca A, Grynberg M, Kajava AV, Promponas VJ, Anisimova M, Jakobsen KS, Linke D (2019) Tandem repeats lead to sequence assembly errors and impose multi-level challenges for genome and protein databases. *Nucleic Acids Res* **47**: 10994-11006

Viguera E, Canceill D, Ehrlich SD (2001) Replication slippage involves DNA polymerase pausing and dissociation. *The EMBO journal* **20**: 2587-2595

Vinces MD, Legendre M, Caldara M, Hagihara M, Verstrepen KJ (2009a) Unstable tandem repeats in promoters confer transcriptional evolvability. *Science* **324**: 1213-1236

Vinces MD, Legendre M, Caldara M, Hagihara M, Verstrepen KJ (2009b) Unstable tandem repeats in promoters confer transcriptional evolvability. *Science* **324**: 1213-1216

Wang Q (1993) *Hydrogen Bond in Organic Chemistry*, Tianjin, China: Tianjin University Press.

Watson JD, Crick FHC (1953a) Genetical implications of the structure of dexoyribonucleic acid. *Nature* **171**: 964-967

Watson JD, Crick FHC (1953b) Molecular structure of deoxypentose nucleic acids. *Nature* **171**: 738-740

Zhao X, Tan Z, Feng H, Yang R, Li M, Jiang J, Shen G, Yu R (2011) Microsatellites in different Potyvirus genomes: survey and analysis. *Gene* **488**: 52-56

Zhao X, Tian Y, Yang R, Feng H, Ouyang Q, Tian Y, Tan Z, Li M, Niu Y, Jiang J (2012) Coevolution between simple sequence repeats (SSRs) and virus genome size. *Bmc Genomics* **13**: 435-435

Zuckerkandl E, Pauling LB (1962) Molecular disease, evolution, and genic heterogeneity. In *Horizons in Biochemistry*, Pullman B, Kasha M, SzentGyörgyi A (eds), pp 189-225. New York: Academic Press, New York

Zuckerkandl E, Pauling LB (1965) Evolutionary divergence and convergence in proteins. In *Evolving Genes and Proteins*, Bryson V, Vogel HJ (eds), pp 97-166. New York: Academic Press, New York

Acknowledgments

The authors thank Qijun Tian, who is from School of design in Hunan University, for his technological help at designing the folded slippage model in helix structure.

Funding

This work was jointly supported by funding from the “National Key Plan for Scientific Research and Development of China” (grant No. 2016YFC1200200 and 2016YFD0500300).

Author contributions

Z. Tan designed and directed this study. D. Li, S. Pan and H. Zhang performed the data analysis, Y. Fu, Z. Peng, H. Zhang and L. Zhang helped for performing data analysis. F. Xu helped for mathematic calculation. S. Peng, Hanrou Huang and Ruixue Shi helped for drawing maps. Z. Tan and D. Li prepared the manuscript. H. Zhang and S. Peng helped for the manuscript preparation.

Conflict of interest

The authors declare that they have no competing interests.

Figure legends

Figure 1 A high percentage of SSRs in genomes and genomes probably tend to produce repeats.

- A. SSR percentages of 55 randomly-selected reported segments and the control group, which consisted of the generated segments according to the sizes and nucleotide compositions of corresponding reported segments.
- B. Contradiction analysis of disappearance and high percentage of SSRs in the genomes.

Figure 2 Straight strand models of semi-conservative replication and slippage.

- A. The space of a nucleotides was drawn. * indicates that those number is the theoretical values (top); The stable straight model of semi-conservative replication (middle); The comparison of hydrogen bond and 3'-5' phosphodiester bonds (bottom)(Gao et al, 2004; Heyrovska, 2006; Wang, 1993). # indicates the strength ratio was calculated by the strength of hydrogen bond dividing that of phosphodiester bond.
- B. The impossible straight slippage models of mononucleotide, dinucleotide and trinucleotide repeats according to the strict geometric calculation of the space of a nucleotide and the stability of hydrogen and phosphodiester bonds.

Figure 3 The DNA molecule is highly curved or folded in the nucleus and the impossible curved slippage model.

- A. Schematic diagram of the size of the nuclear space (top) (Lamond, 2002); The normal replicating enzymes complex straighten the DNA molecule, while the disturbed replicating enzymes complex may cause the DNA molecule return to curved state (bottom).
- B. Impossible curved template slippage model according to the strict geometric calculation of the space of a nucleotide and the stability of hydrogen and phosphodiester bonds (top); Mono- and dinucleotide repeats may be impossibly produced in curved replicating strands (middle and bottom).

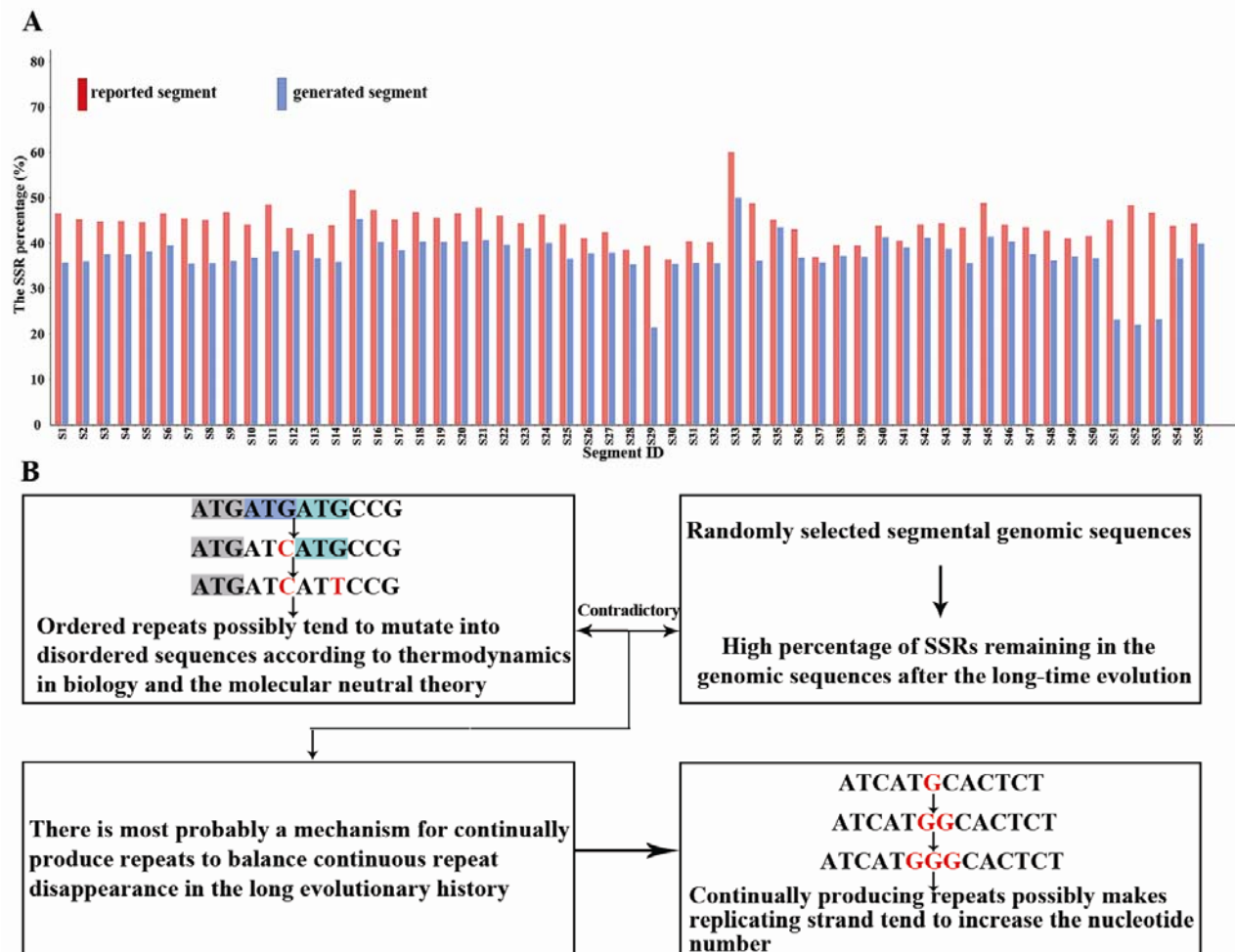
Figure 4 Stable folded slippage models of mononucleotide to hexanucleotide repeats amplification according to the strict geometric calculation of the space of a nucleotide and the stability of hydrogen and phosphodiester bonds. Repeat units tend to be expanded in the replicating strands when the template strands are on the inner side of the folded slippage models respectively. The bottom 3 figures were the folded slippage models in three-dimensional helix form.

Figure 5 Stable folded slippage models of mononucleotide to hexanucleotide repeats contraction according to the strict geometric calculation of the space of a nucleotide and the stability of hydrogen and phosphodiester bonds. Repeat units tend to be subtracted in the replicating strands when the template strands are on the outside of the folded slippage models respectively. The bottom 3 figures were the folded slippage models in three-dimensional helix form.

Figure 6 Repeats production incline to expansion. F^o , F^i refer to the force required for the two template strands to bend, respectively. $F^o > F^i$ means that the force of the template strand bending downward is greater than the bending upward, and $P^e > P^c$ means that the possibility of the template strand bending upward is greater than the downward bending.

Figures

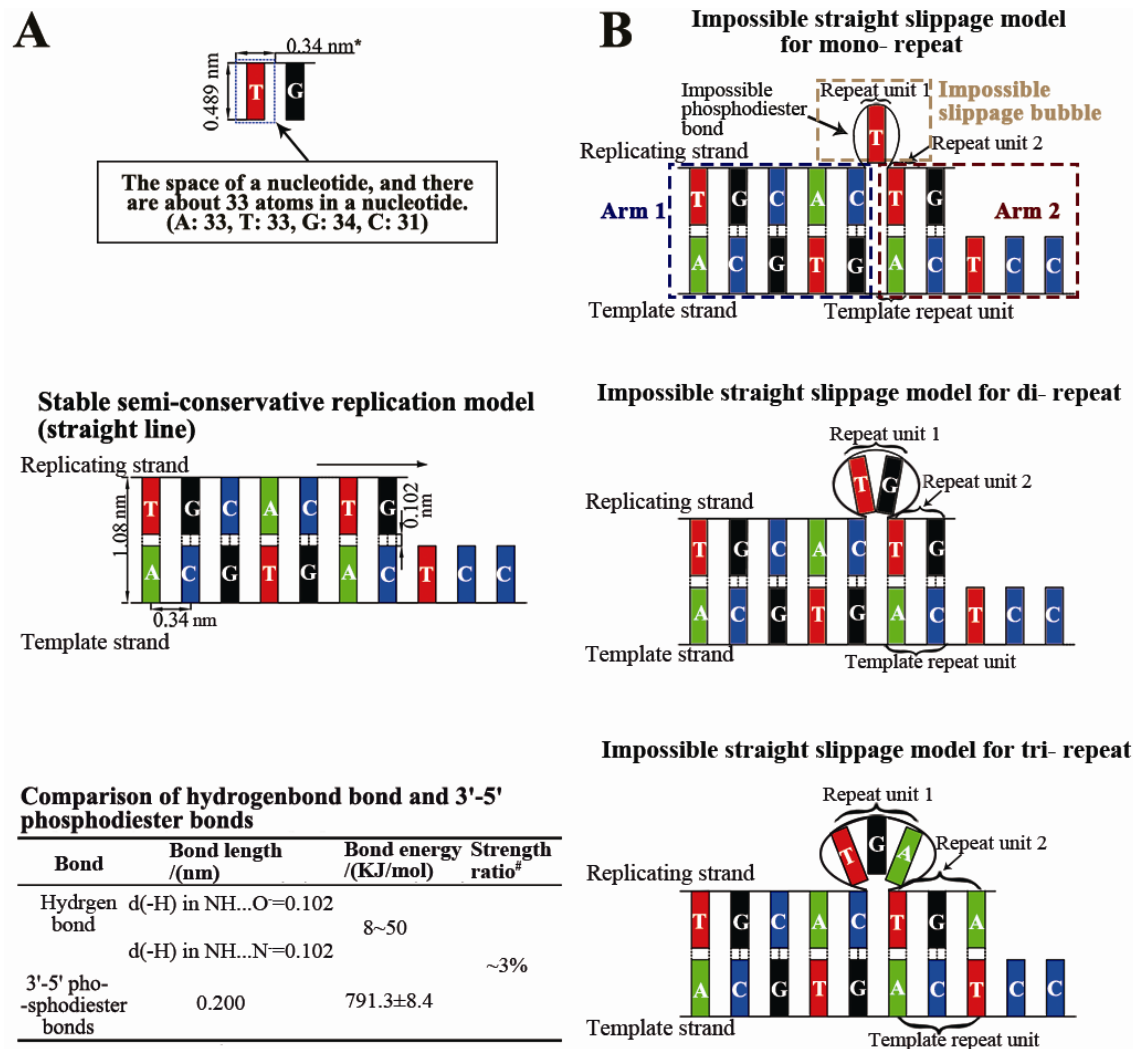
Figure 1



A high percentage of SSRs in genomes and genomes probably tend to produce repeats.

- A. SSR percentages of 55 randomly-selected reported segments and the control group, which consisted of the generated segments according to the sizes and nucleotide compositions of corresponding reported segments.
- B. Contradiction analysis of disappearance and high percentage of SSRs in the genomes.

Figure 2

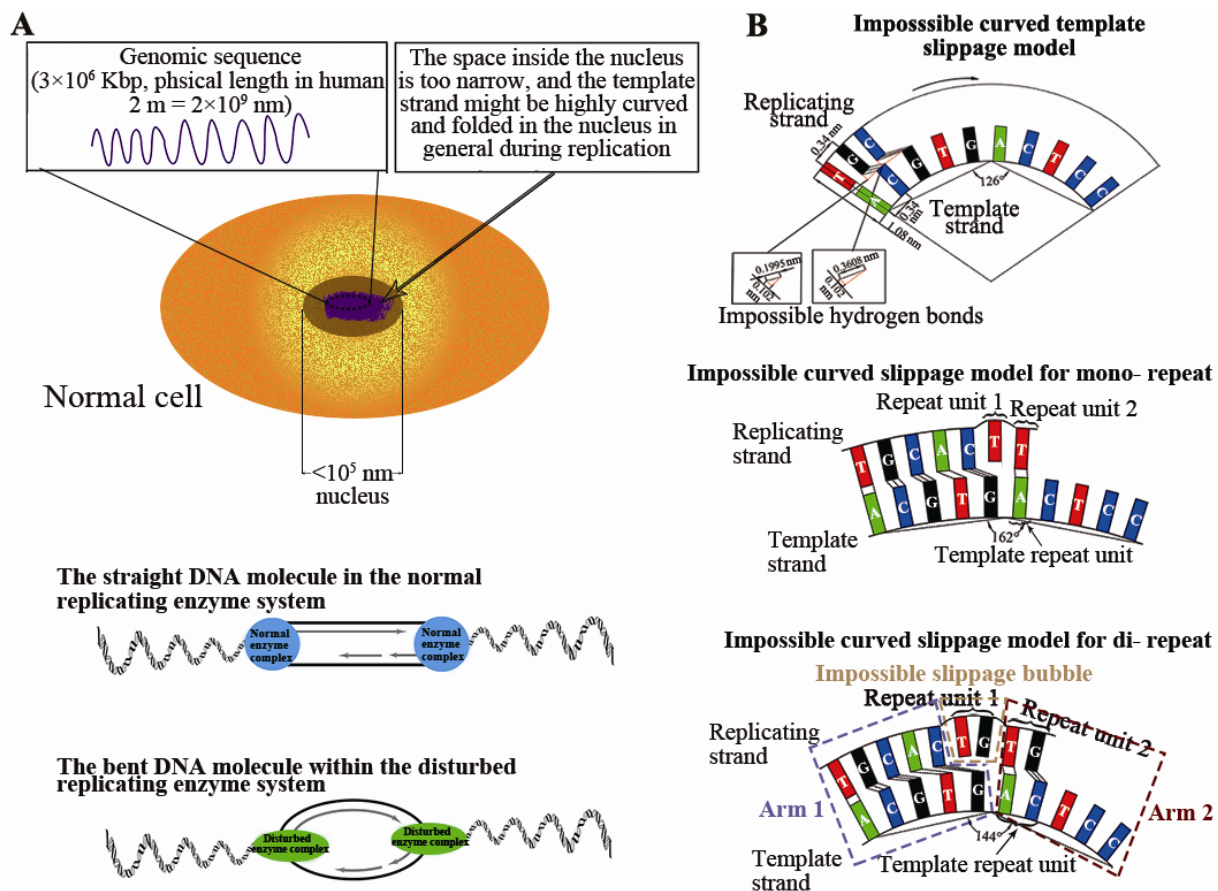


Straight strand models of semi-conservative replication and slippage.

- The space of a nucleotides was drawn. * indicates that those number is the theoretical values (top); The stable straight model of semi-conservative replication (middle); The comparison of hydrogen bond and 3'-5' phosphodiester bonds (bottom)(Gao et al, 2004; Heyrovska, 2006; Wang, 1993). [#] indicates the strength ratio was calculated by the strength of hydrogen bond dividing that of phosphodiester bond.
- The impossible straight slippage models of mononucleotide, dinucleotide and trinucleotide repeats according to the strict geometric calculation of the space of a nucleotide and the stability

of hydrogen and phosphodiester bonds.

Figure 3

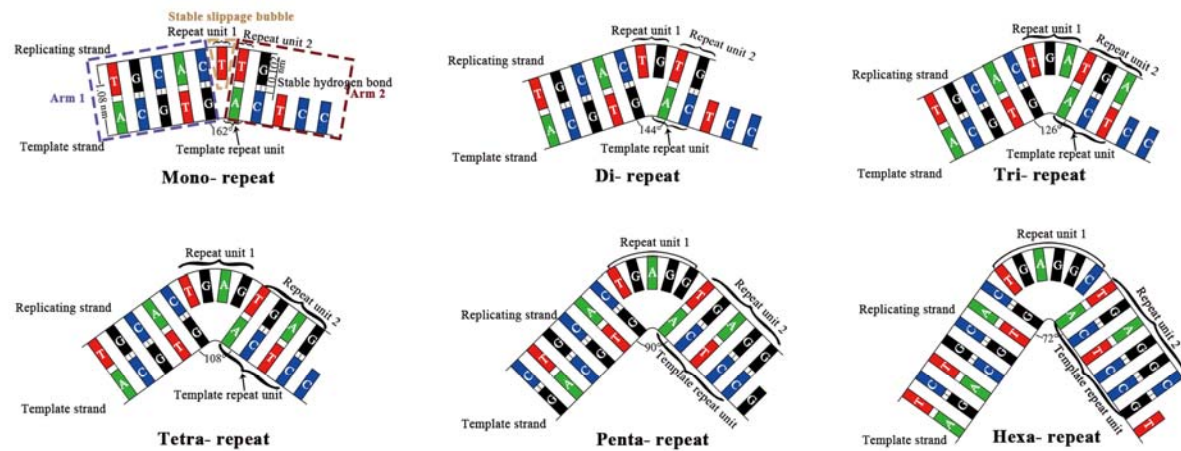


The DNA molecule is highly curved or folded in the nucleus and the impossible curved slippage model.

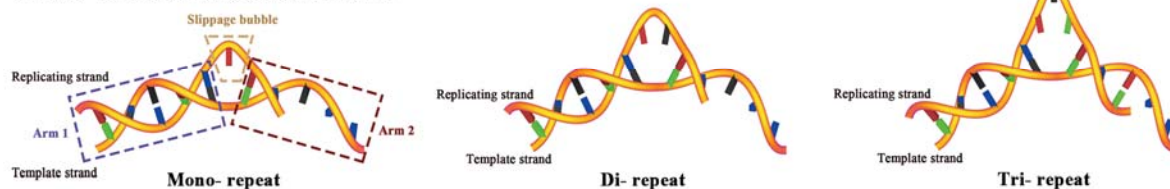
- A. Schematic diagram of the size of the nuclear space (top) (Lamond, 2002); The normal replicating enzymes complex straighten the DNA molecule, while the disturbed replicating enzymes complex may cause the DNA molecule return to curved state (bottom).
- B. Impossible curved template slippage model according to the strict geometric calculation of the space of a nucleotide and the stability of hydrogen and phosphodiester bonds (top); Mono- and dinucleotide repeats may be impossibly produced in curved replicating strands (middle and bottom).

Figure 4

Plane form



Three-dimensional helix form

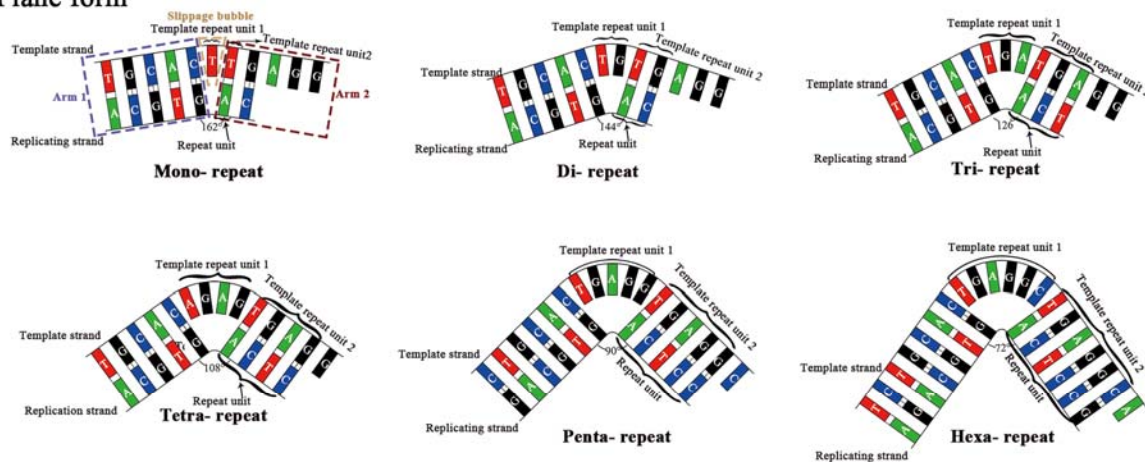


Repeat units tend to be expanded in the replicating process when the template strand on the inner side.

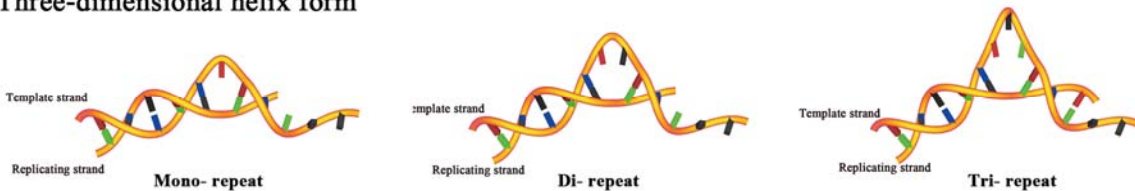
Stable folded slippage models of mononucleotide to hexanucleotide repeats amplification according to the strict geometric calculation of the space of a nucleotide and the stability of hydrogen and phosphodiester bonds. Repeat units tend to be expanded in the replicating strands when the template strands are on the inner side of the folded slippage models respectively. The bottom 3 figures were the folded slippage models in three-dimensional helix form.

Figure 5

Plane form



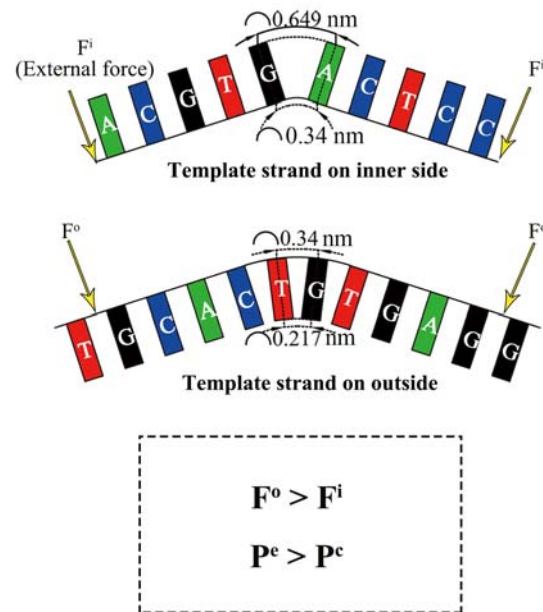
Three-dimensional helix form



Repeat units tend to be contracted in the replicating process when the template strand is on the outside.

Stable folded slippage models of mononucleotide to hexanucleotide repeats contraction according to the strict geometric calculation of the space of a nucleotide and the stability of hydrogen and phosphodiester bonds. Repeat units tend to be subtracted in the replicating strands when the template strands are on the outside of the folded slippage models respectively. The bottom 3 figures were the folded slippage models in three-dimensional helix form.

Figure 6



Repeats production incline to expansion. F^o , F^i refer to the force required for the two template strands to bend, respectively. $F^o > F^i$ means that the force of the template strand bending downward is greater than the bending upward, and $P^e > P^c$ means that the possibility of the template strand bending upward is greater than the downward bending.

Table 1. The lengths (bp) of SSRs with different repeat unit types and different iterations in the segment of the reported human reference X chromosomal sequence at the location of 144822-231384 bp.

Iteration	Mono ^a	Di	Tri	Tetra	Penta	Hexa	Total
I ₂	(18128) ^b	10040	3540	2056	1250	480	17366
I ₃	9702	1782	288	156	45	18	11991
I ₄	3844	368	12	112	-	-	4336
I ₅	2095	120	15	20	-	-	2250
I ₆	600	24	18	0	-	-	642
I ₇	182	14	- ^c	28	-	-	224
I ₈	128	16	-	0	-	-	144
I ₉	54	18	-	36	-	-	108
I ₁₀	50	0	-	-	-	-	50
I ₁₁	55	22	-	-	-	-	77
I ₁₂	24	-	-	-	-	-	24
I ₁₃	65	-	-	-	-	-	65
I ₁₄	56	-	-	-	-	-	56
I ₁₅	45	-	-	-	-	-	45
I ₁₆	64	-	-	-	-	-	64
I ₁₇	0	-	-	-	-	-	0
I ₁₈	36	-	-	-	-	-	36
I ₁₉	19	-	-	-	-	-	19
I ₂₀	0	-	-	-	-	-	0
I ₂₁	42	-	-	-	-	-	42
I ₂₂	0	-	-	-	-	-	0
I ₂₃	23	-	-	-	-	-	23
I ₂₄	0	-	-	-	-	-	0

I ₂₅	25	-	-	-	-	-	25
I ₂₆	-	-	-	-	-	-	-
I ₂₇	-	-	-	-	-	-	-
I ₂₈	-	-	-	-	-	-	-
Sum	17109	12404	3873	2408	1295	498	37587

^a Mononucleotide repeat (Mono), Dinucleotide repeat (Di), Trinucleotide repeat (Tri), Tetranucleotide repeat (Tetra), Pentanucleotide repeat (Penta), Hexanucleotide repeat (Hexa).

^b The length of mononucleotide repeats with iterations of 2 was not included in this statistics and just used as the reference here.

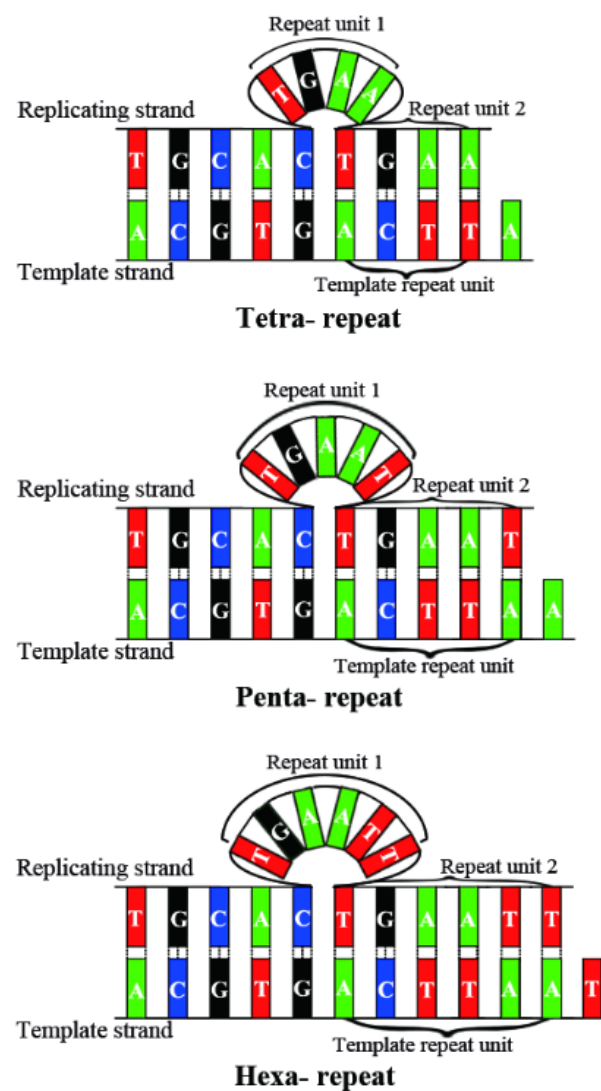
^c Beyond the largest iteration of this repeat unit type in corresponding analyzed segments were expressed as “-”.

Supplementary Tables are online at

https://github.com/DooYal/Supplementary-Table-for-submitting-relatively-...-/tree/DooYal-patch-manuscript_folded/supplementary%20tables

Supplementary Figures

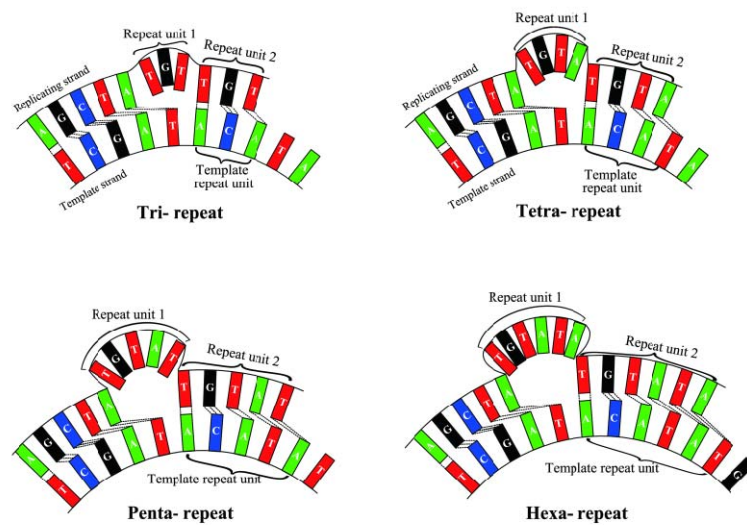
Figure S1



Impossible straight slippage models for tetra- to hexanucleotide repeats when the slippage bubble occurs at the replicating strand. The model drawing was based on the strict geometric calculation of

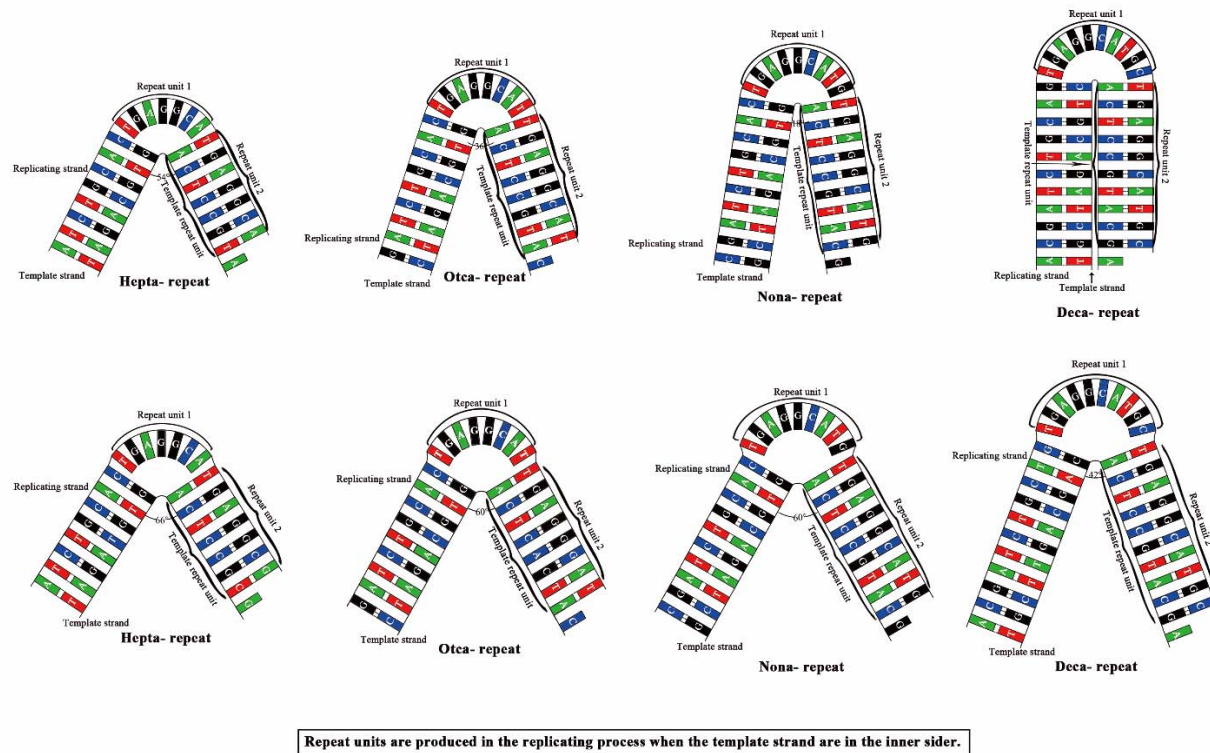
the space of a nucleotide and the stability of hydrogen and phosphodiester bonds.

Figure S2



Impossible curved slippage models for tri- to hexanucleotide repeats when the template strand in the inner side of the models.

Figure S3



The possible folded slippage models for hepta- to decanucleotide repeat amplification. Repeat units tend to be expanded in the replicating strands when the template strands are on the inner side of the folded slippage models respectively.

