# Sequences that stall replication and shape the genome
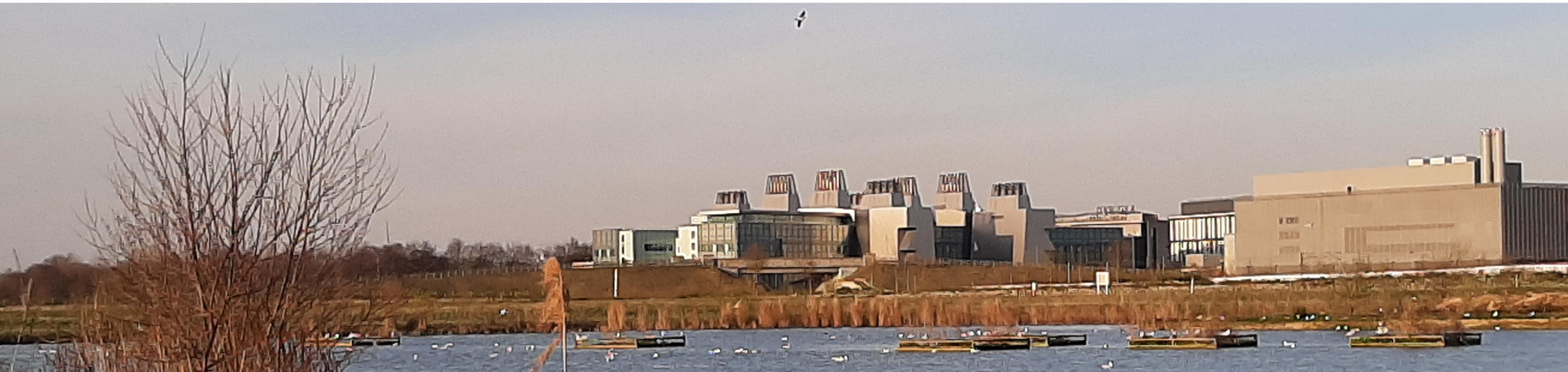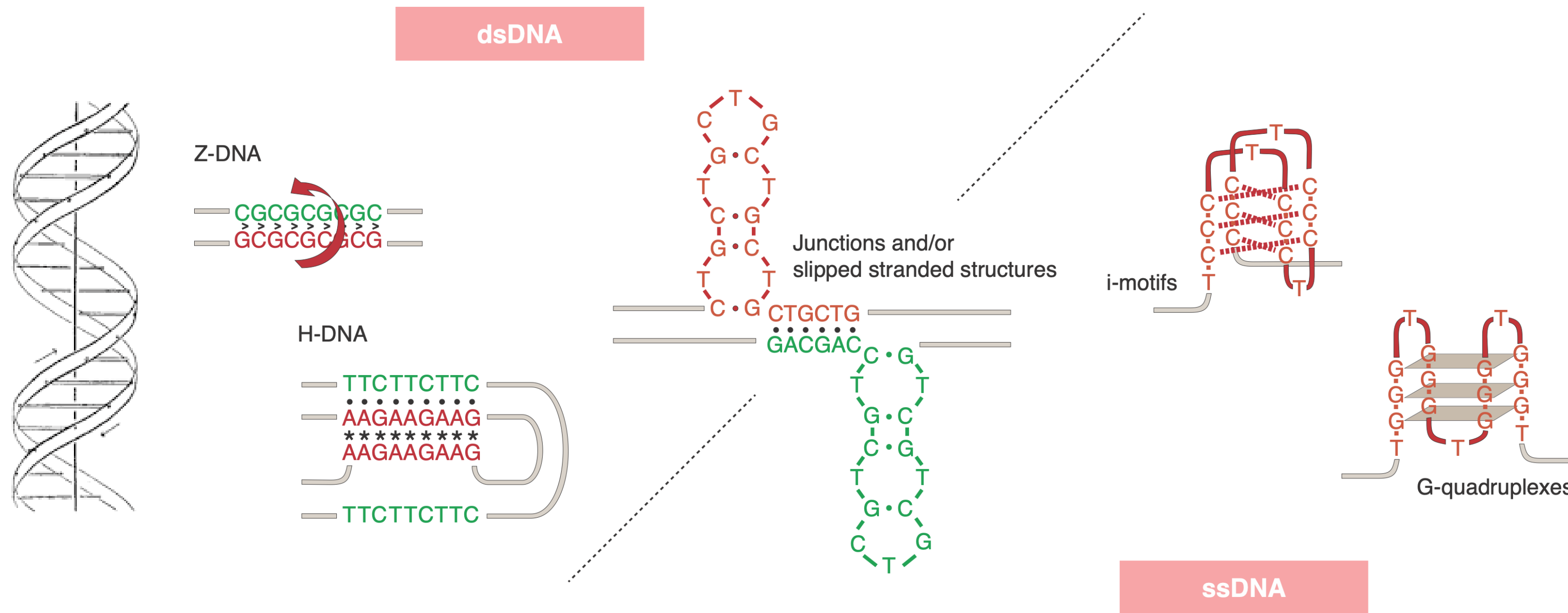
Julian E. Sale
MRC Laboratory of Molecular Biology
Cambridge

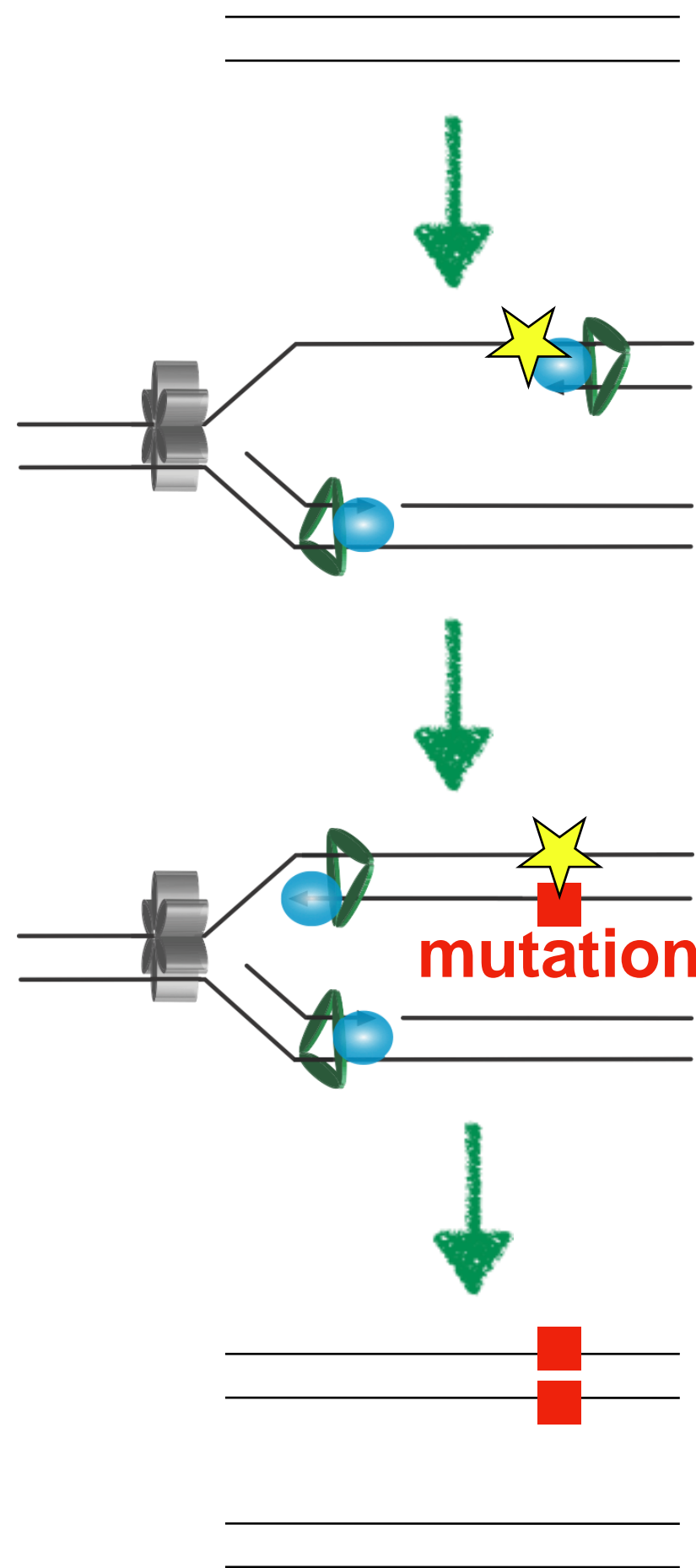# The challenge to replication posed by alternative DNA structures



- millions of loci in the human genome - often found in repetitive / low complexity sequence

- Hotspots for chromosomal rearrangements, copy number variations, mutagenesis and epigenetic instability

- Trinucleotide repeat expansion disorders

| | |
|---|---|
| Huntington's disease | (CAG)n |
| Friedreich's ataxia | (GAA)n |
| Fragile X syndrome | (CGG)n |
| etc. etc. | |

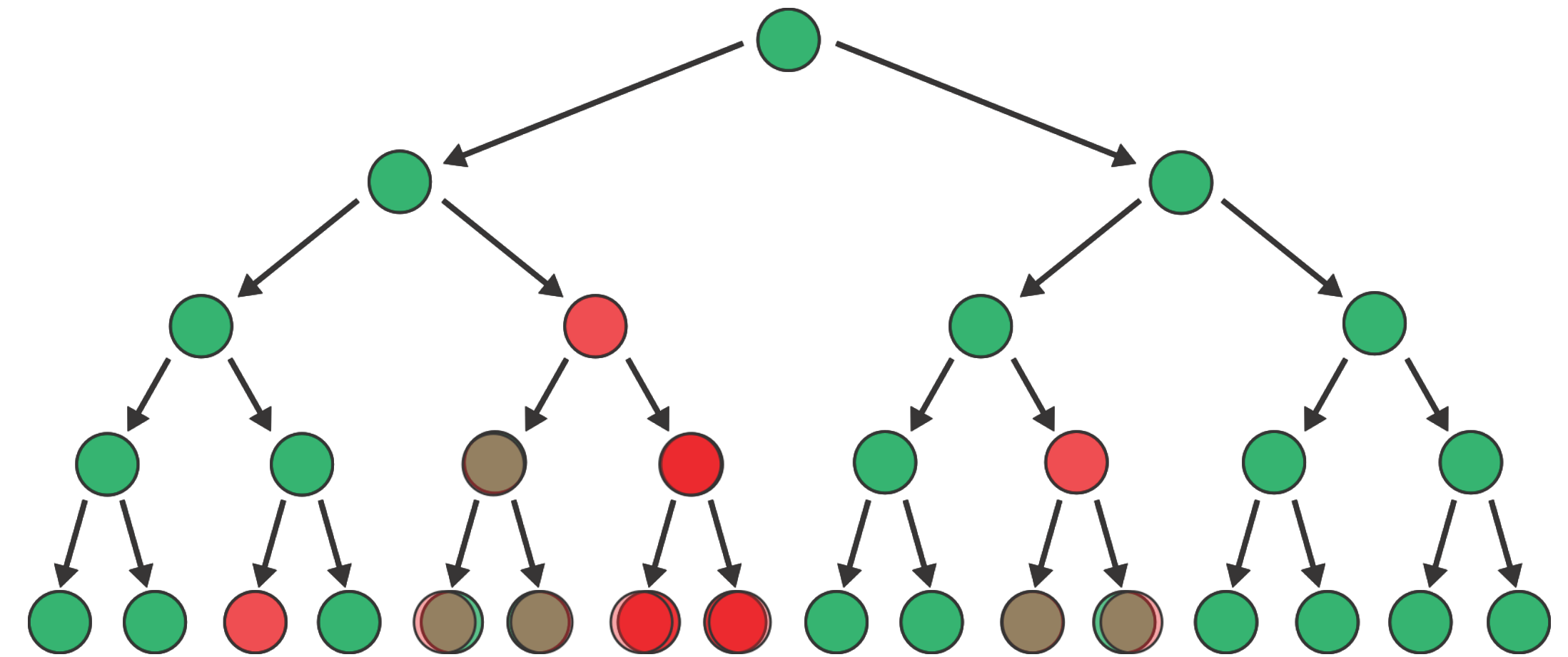**How does the replisome deal with secondary structures in the template?**

**What drives the evolution of genomic sequences with structure-forming potential?**

**polymerase stalling is transient and therefore hard to detect …**

**mutation**

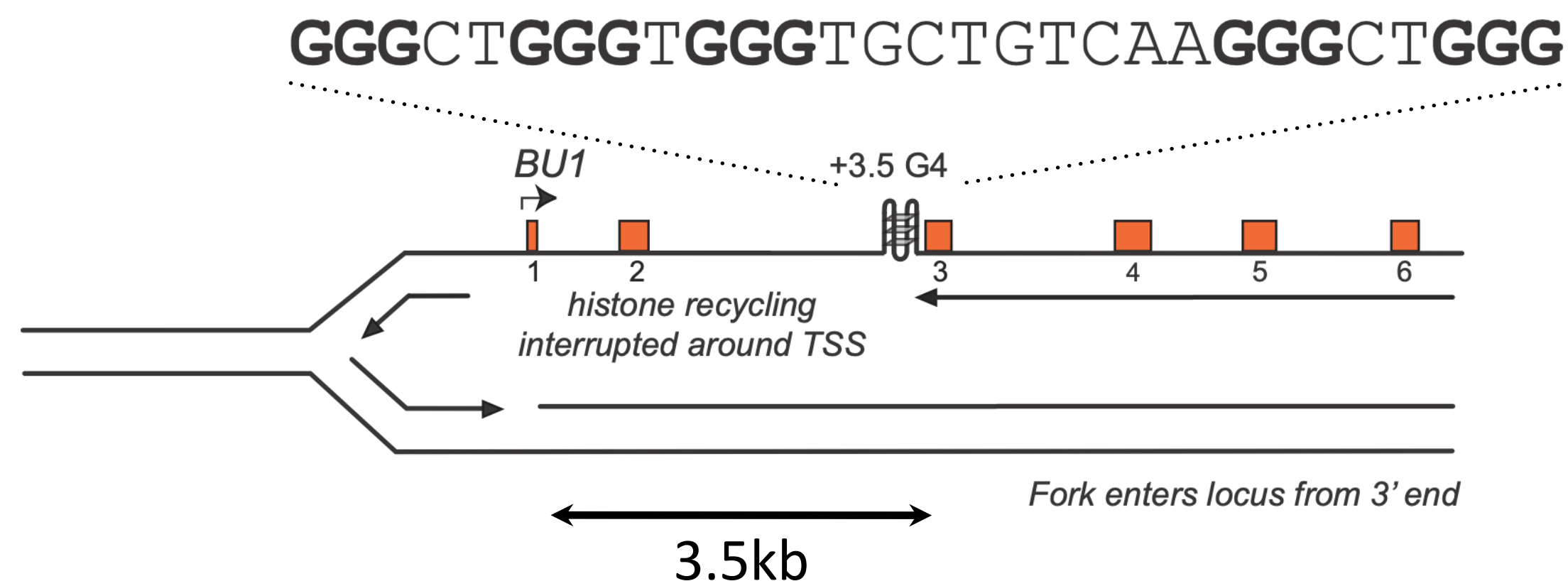**… mutagenic outcomes are traceable but relatively rare**

**How can non-mutagenic episodes of fork stalling be reported in an expanding cell population?**

# Using local loss of epigenetic memory to monitor delayed replication of DNA secondary structures
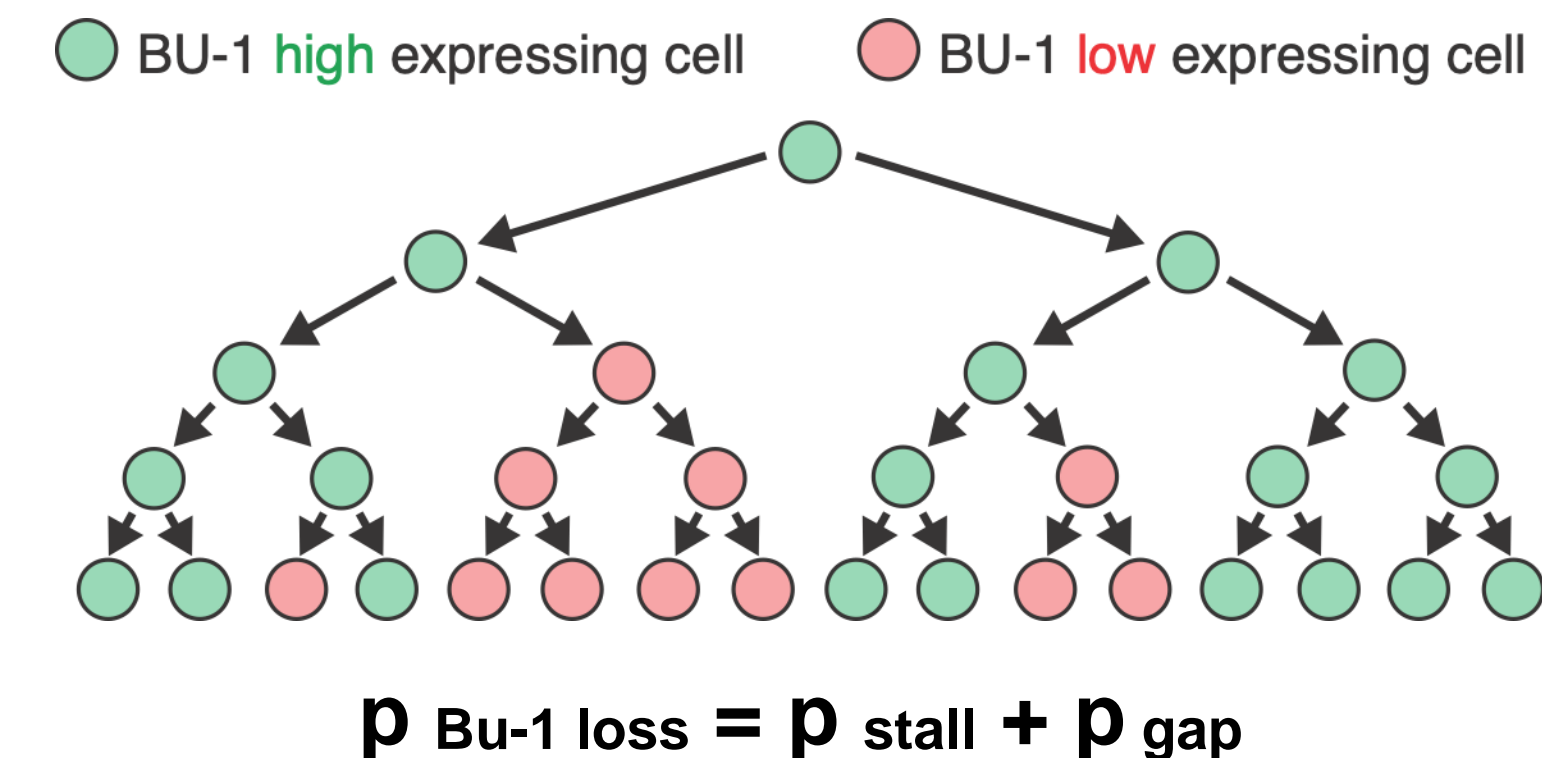


2° structure processing at the fork

normal expression

Recycled H3/4

New H3/4

aberrant expression

Delayed 2° structure processing

chromatin 'scar'

Sarkies et al. (2010) Mol Cell 40, 703-713

GGGCTGGGTGGGTGCTGTCAAGGGCTGGG

BU1     +3.5 G4

1   2       3       4   5   6

histone recycling
interrupted around TSS

Fork enters locus from 3' end

3.5kb

Sarkies et al. (2012) NAR 40, 1485-1498
Schiavone, Guilbaud et al. (2014) EMBO J 33, 2507-20

BU-1 high expressing cell     BU-1 low expressing cell

$$p_{\text{Bu-1 loss}} = p_{\text{stall}} + p_{\text{gap}}$$

# Using local loss of epigenetic memory to monitor delayed replication of DNA secondary structures

2° structure processing at the fork

normal expression

Recycled H3/4

New H3/4

Sarkies et al. (2010) Mol Cell 40, 703-713

aberrant expression

Delayed 2° structure processing

chromatin 'scar'

**GGG**CT**GGG**T**GGG**TGCTGTCAA**GGG**CT**GGG**

BU1    +3.5 G4

1    2         3      4    5    6

histone recycling
interrupted around TSS

Fork enters locus from 3' end

3.5kb

WT expn

Loss variants

WT
rev1

% of maximum

Bu-1
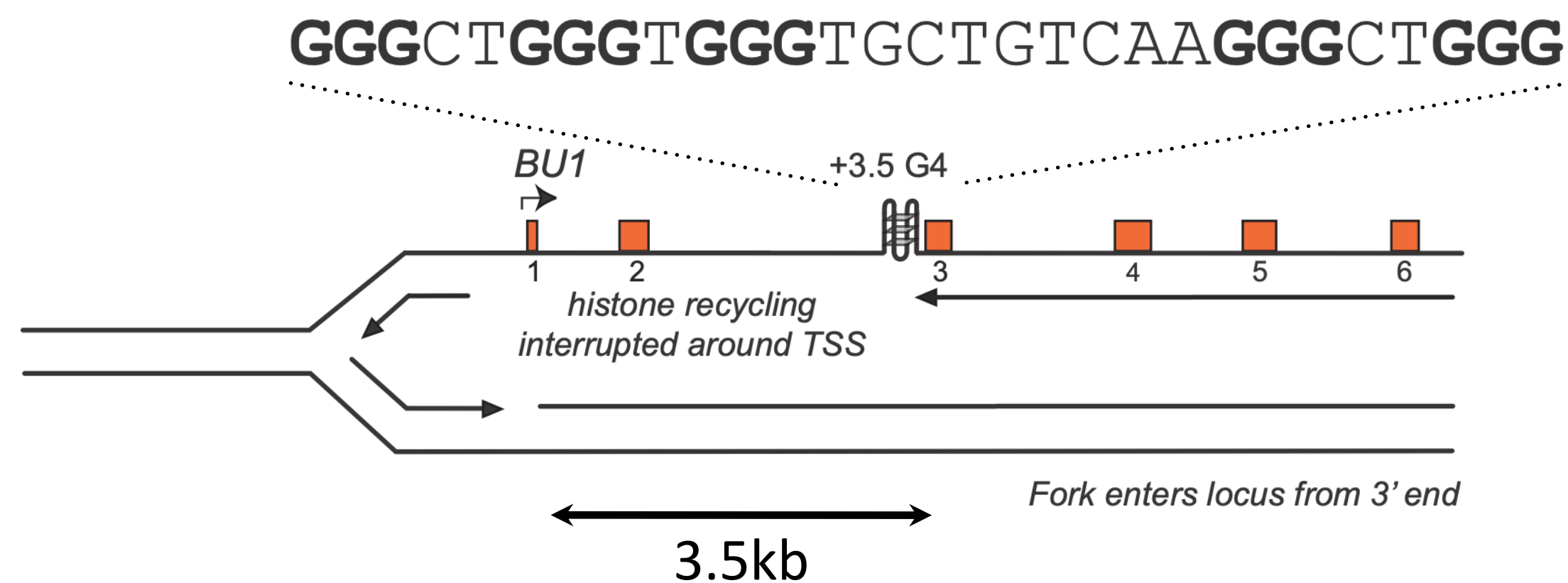
Sarkies et al. (2012) NAR 40, 1485-1498
Schiavone, Guilbaud et al. (2014) EMBO J 33, 2507-20

# Using local loss of epigenetic memory to monitor delayed replication of DNA secondary structures
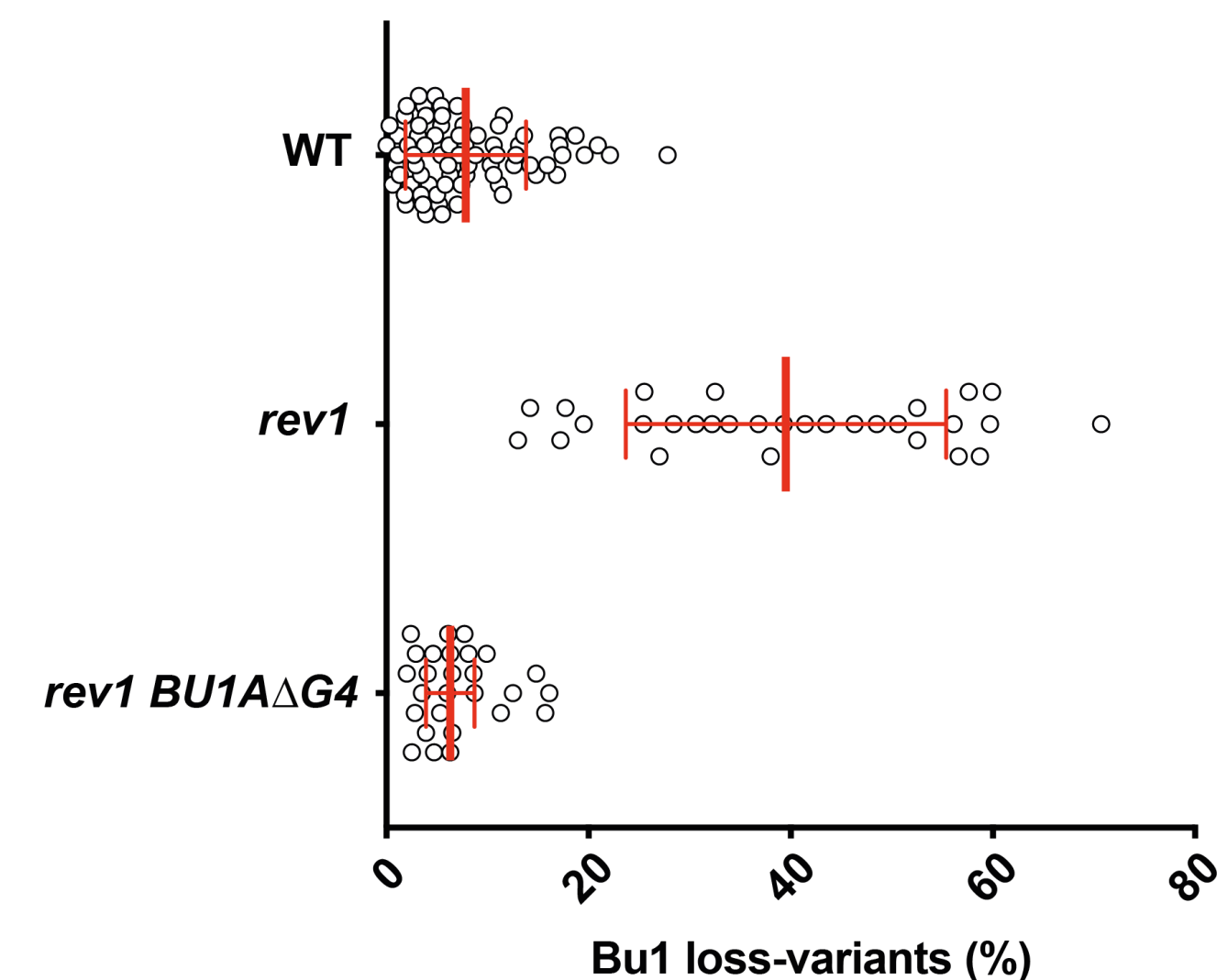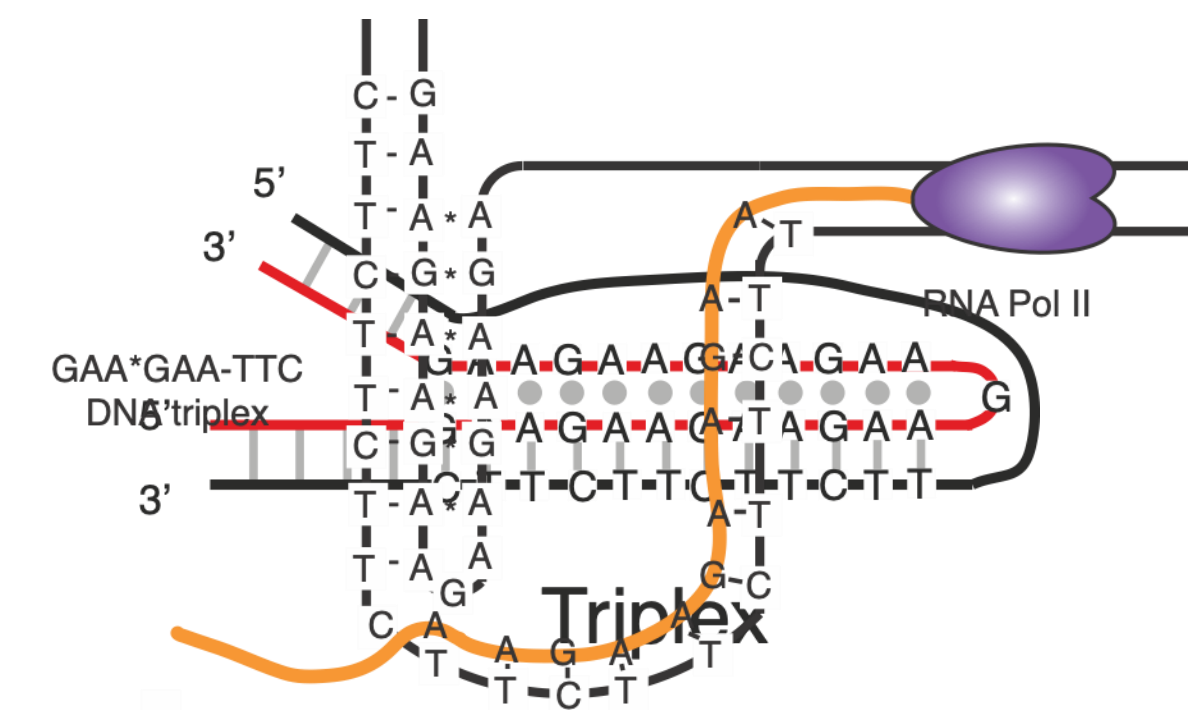


2° structure processing at the fork

Recycled H3/4

New H3/4

normal expression

aberrant expression

chromatin 'scar'

Delayed 2° structure processing

Sarkies et al. (2010) Mol Cell 40, 703-713

**GGG**CT**GGG**T**GGG**TGCTGTCAA**GGG**CT**GGG**

BU1

+3.5 G4

1   2        3      4    5      6

*histone recycling interrupted around TSS*

*Fork enters locus from 3' end*

3.5kb

Sarkies et al. (2012) NAR 40, 1485-1498
Schiavone, Guilbaud et al. (2014) EMBO J 33, 2507-20

WT

*rev1*

*rev1 BU1A△G4*

0        20       40       60       80
**Bu1 loss-variants (%)**

also:

*fancj*
*wrn*
*blm*
*pif1*

# Structure formation is likely a frequent event during replication



**PrimPol**

Repriming



AEP polymerase    ZnF  RBMs

N

1

560

- **PrimPol** is the second identified primase in vertebrates
- RNA / DNA primase
- DNA polymerase with some capacity for lesion bypass

- Unable to replicate G4s
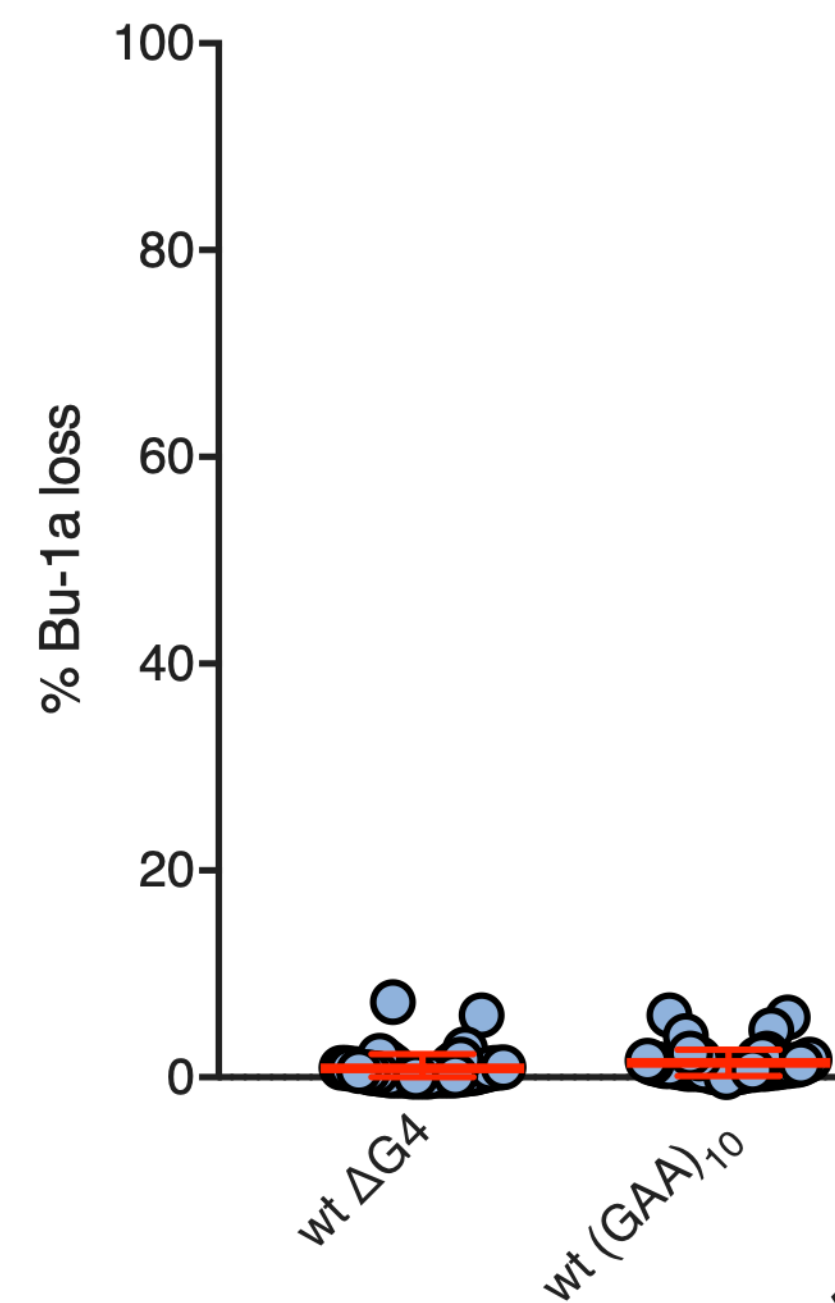- Binds to G4s and efficiently reprimes close by the structure



WT

*primpol #3*

*primpol #4*

*primpol* :
hPRIMPOL

*primpol*
*BU1$^{G4\Delta}$*

Bu-1a loss (%)

# PrimPol loss reveals that even short repeats can be replication impediments



Loss of histone modifications at BU-1 promoter

GAAGAAGAAGAAGAAGAAGAAGAAGAAGAA

ssDNA gap creating zone of interrupted histone recycling

Fork enters locus from 3' end and pauses at +3.5 G4

The 'transcriptional diode' Grabczyk & Fishman JBC 1995

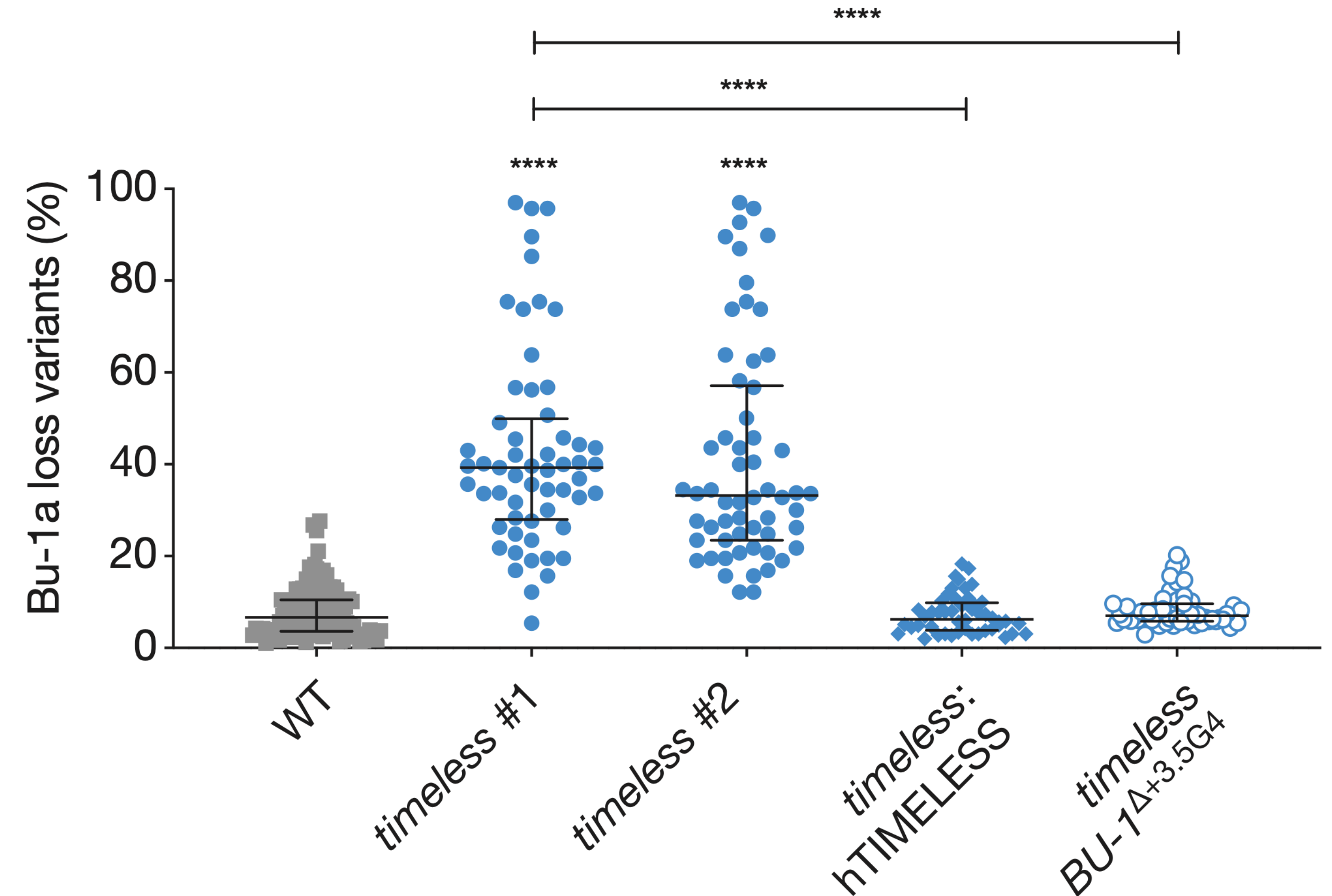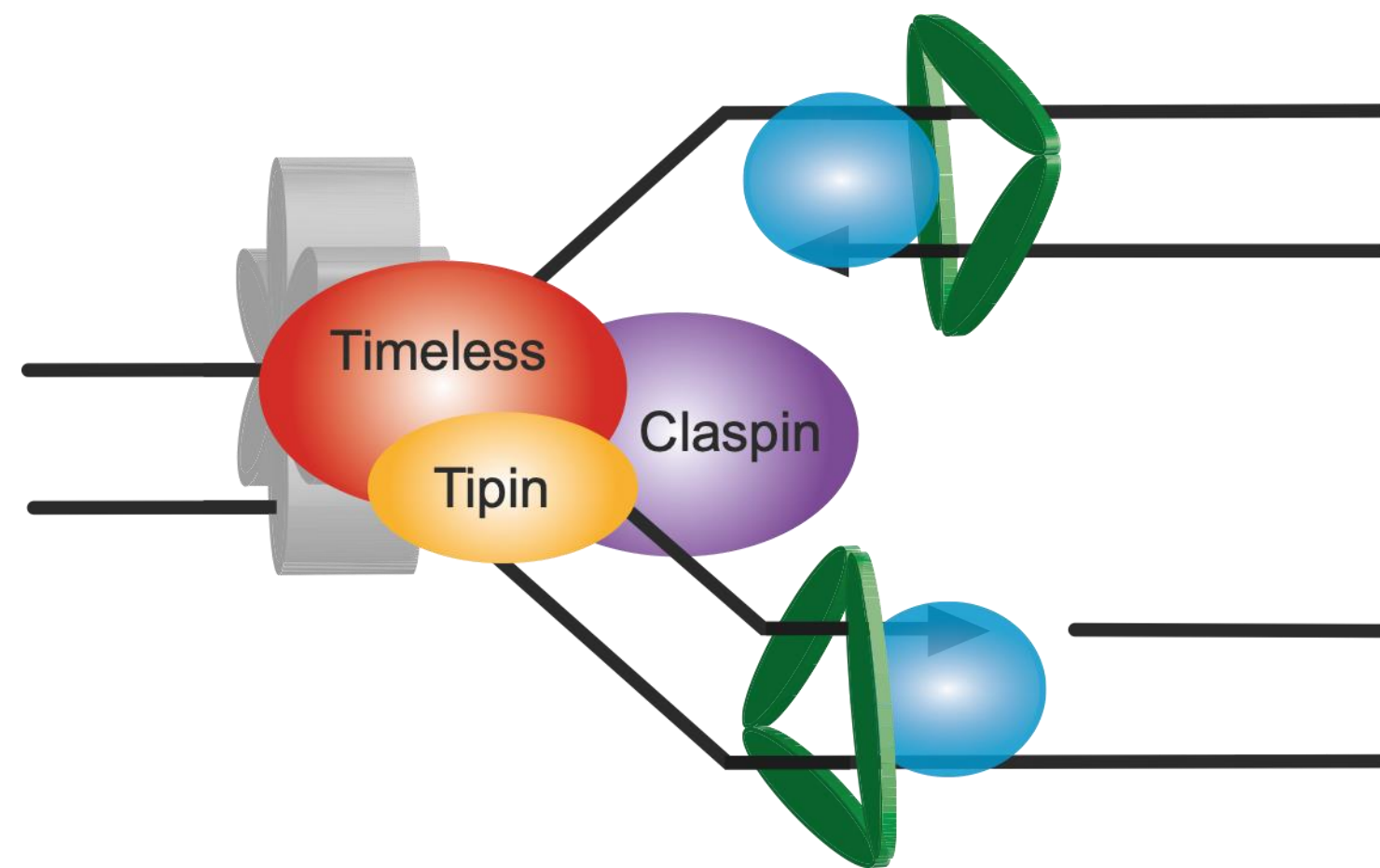(GAA)n repeats will form triplexes when transcribed as the coding strand

- **(GAA)$_{10}$ on the leading strand template is able to block DNA synthesis**

Šviković et al. (2019) *EMBO J* 38(3):e99793

- **RNaseH1 removes the need for PrimPol-mediated repriming suggesting that (GAA)10 requires RNA:DNA hybrid formation for it to become a replication impediment**

Šviković et al. (2019) *EMBO J* 38, pii: 399793

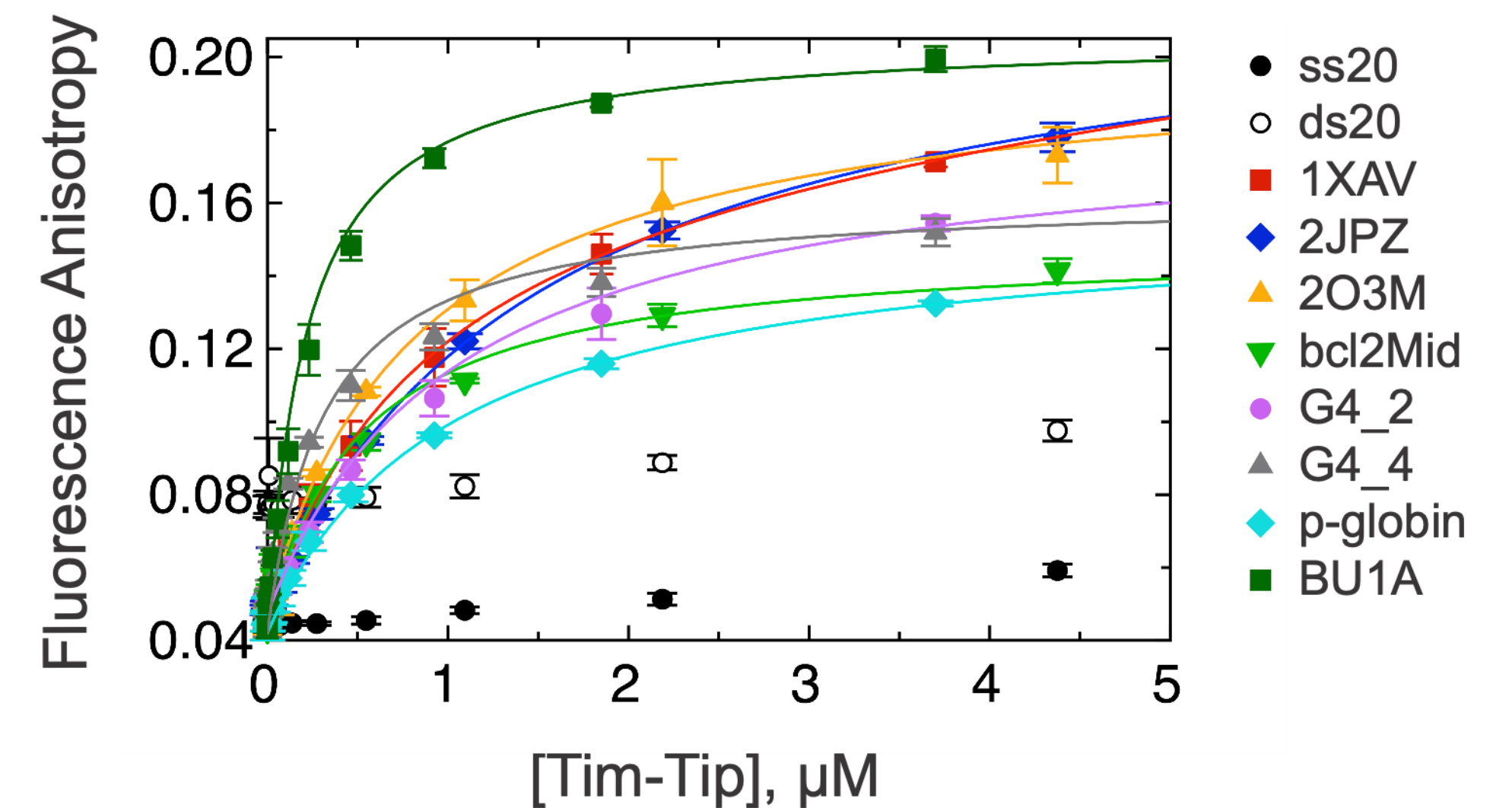# Does the replisome have specific mechanisms for surveillance of structure formation: the fork protection complex

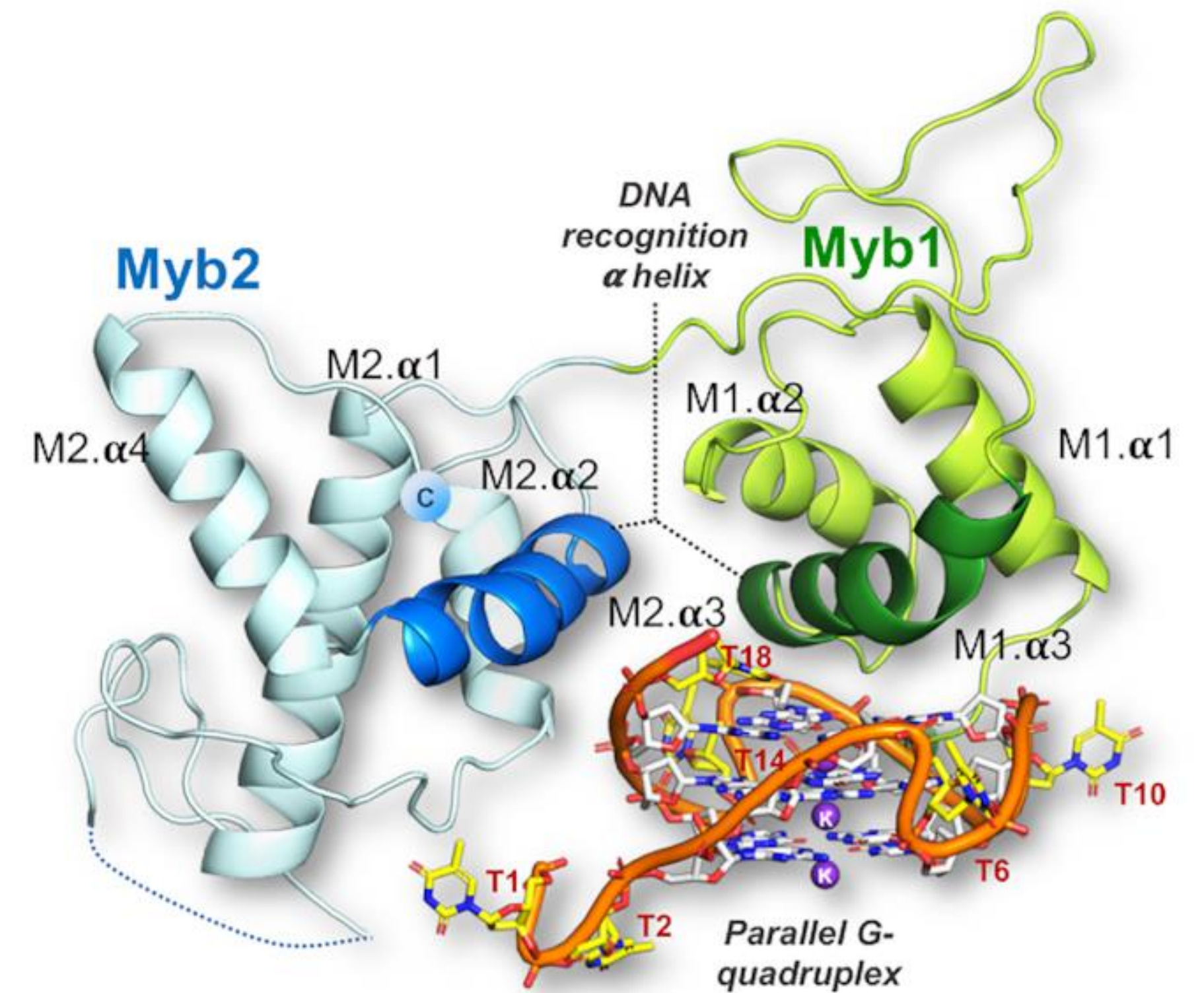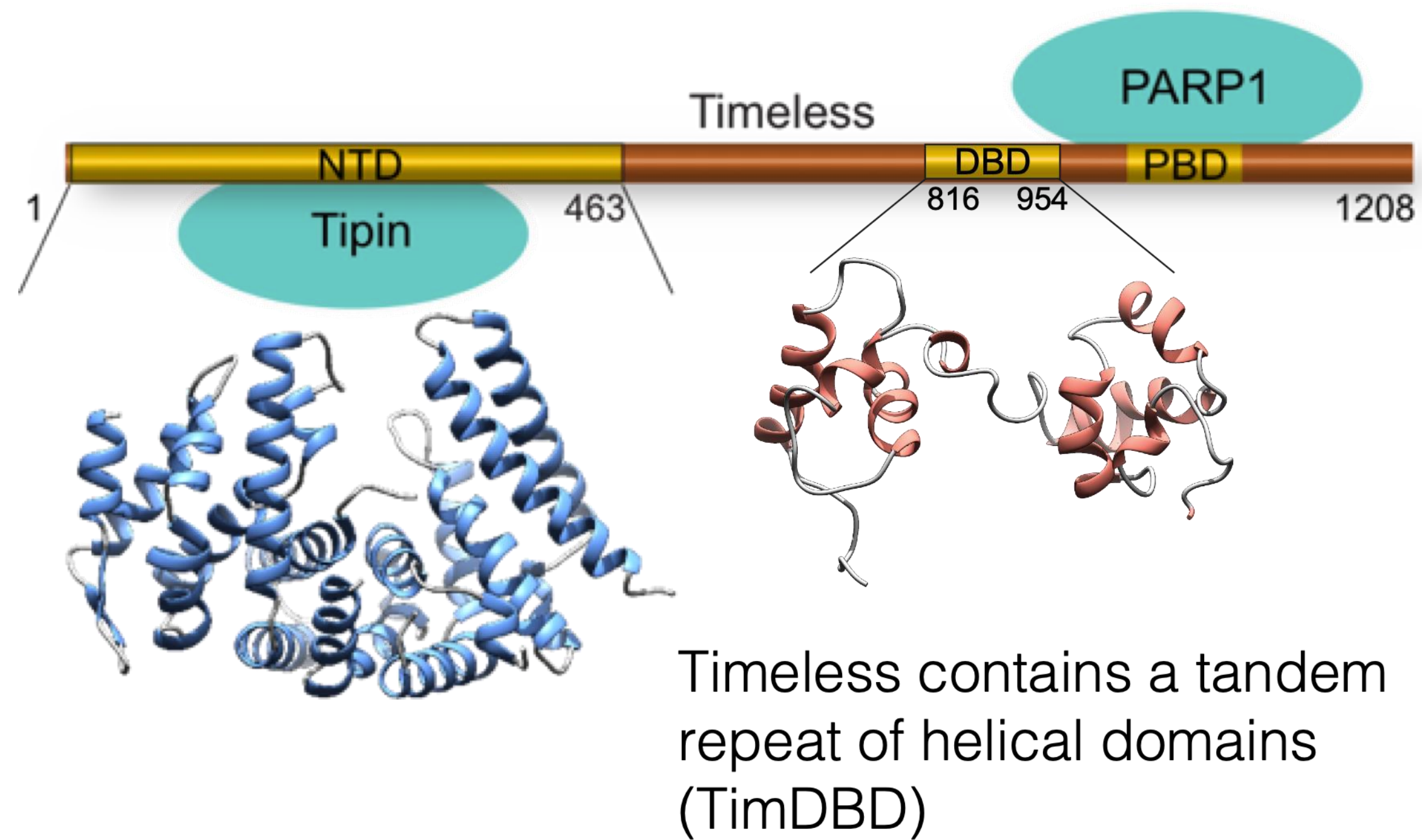# A newly identified DNA binding domain in Timeless



Timeless contains a tandem repeat of helical domains (TimDBD)
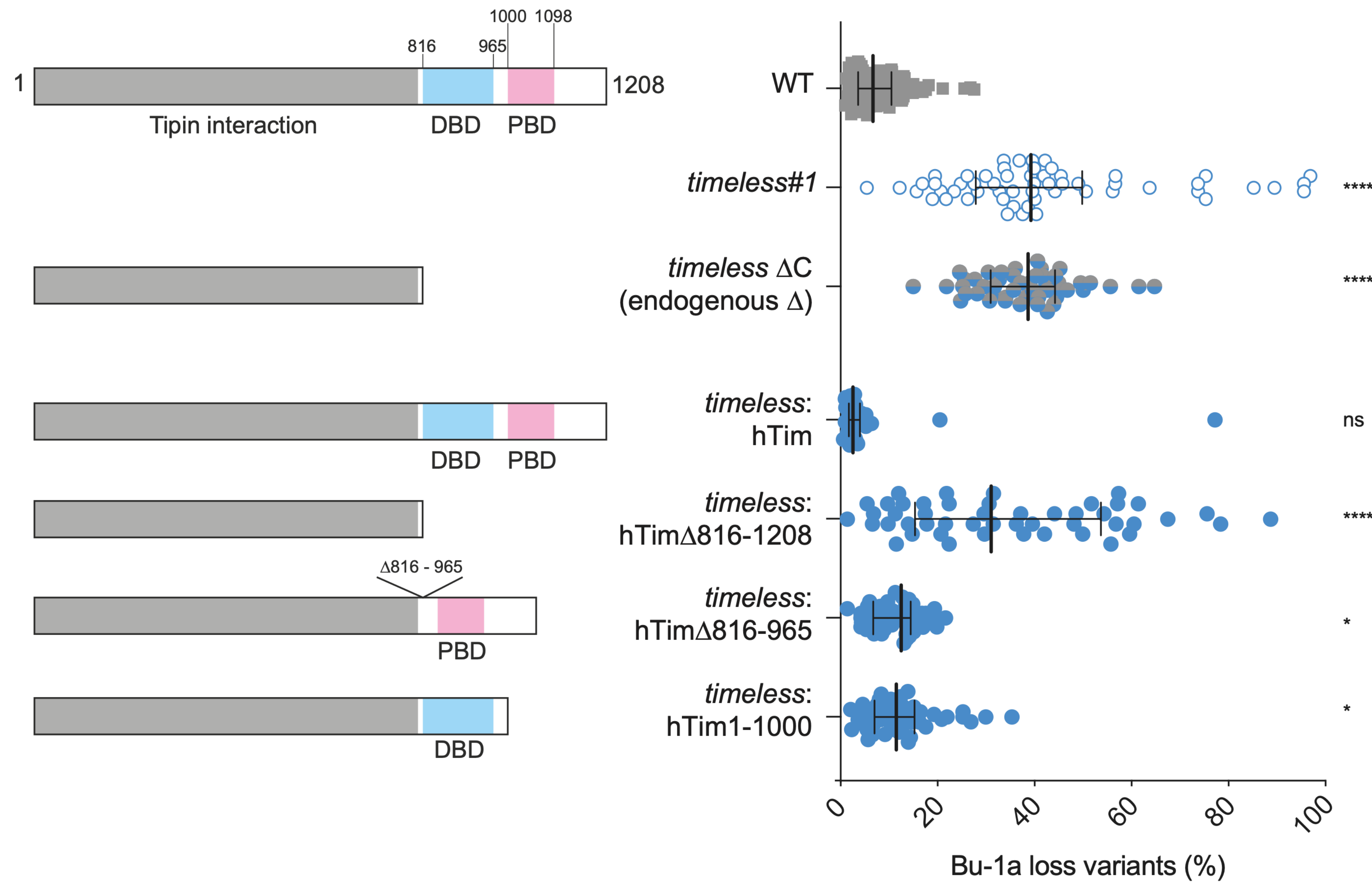
G4 DNA : **6FAM-TGAGGGTGGGTAGGGTGGGTAA**

Luca Pellegrini lab

# A newly identified DNA binding domain in Timeless



Timeless contains a tandem repeat of helical domains (TimDBD)

G4 recognition by the tandem Myb domains of RAP1

Traczyk … Rhodes *NAR 2020*

Luca Pellegrini lab

# The C-terminus of Timeless is required to prevent G4-induced instability of *BU-1*

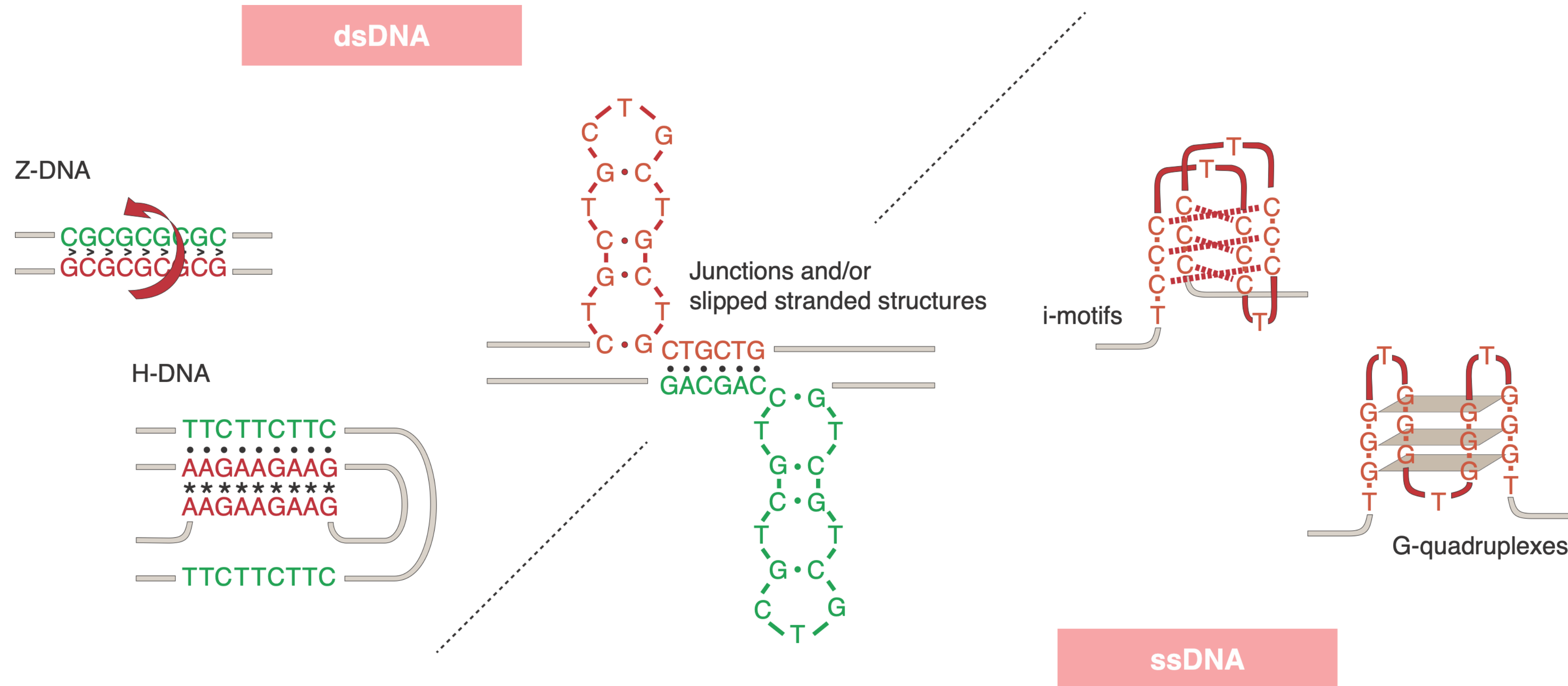# The interaction of Timeless with the helicase DDX11 is required for processing fork-stalling G4s



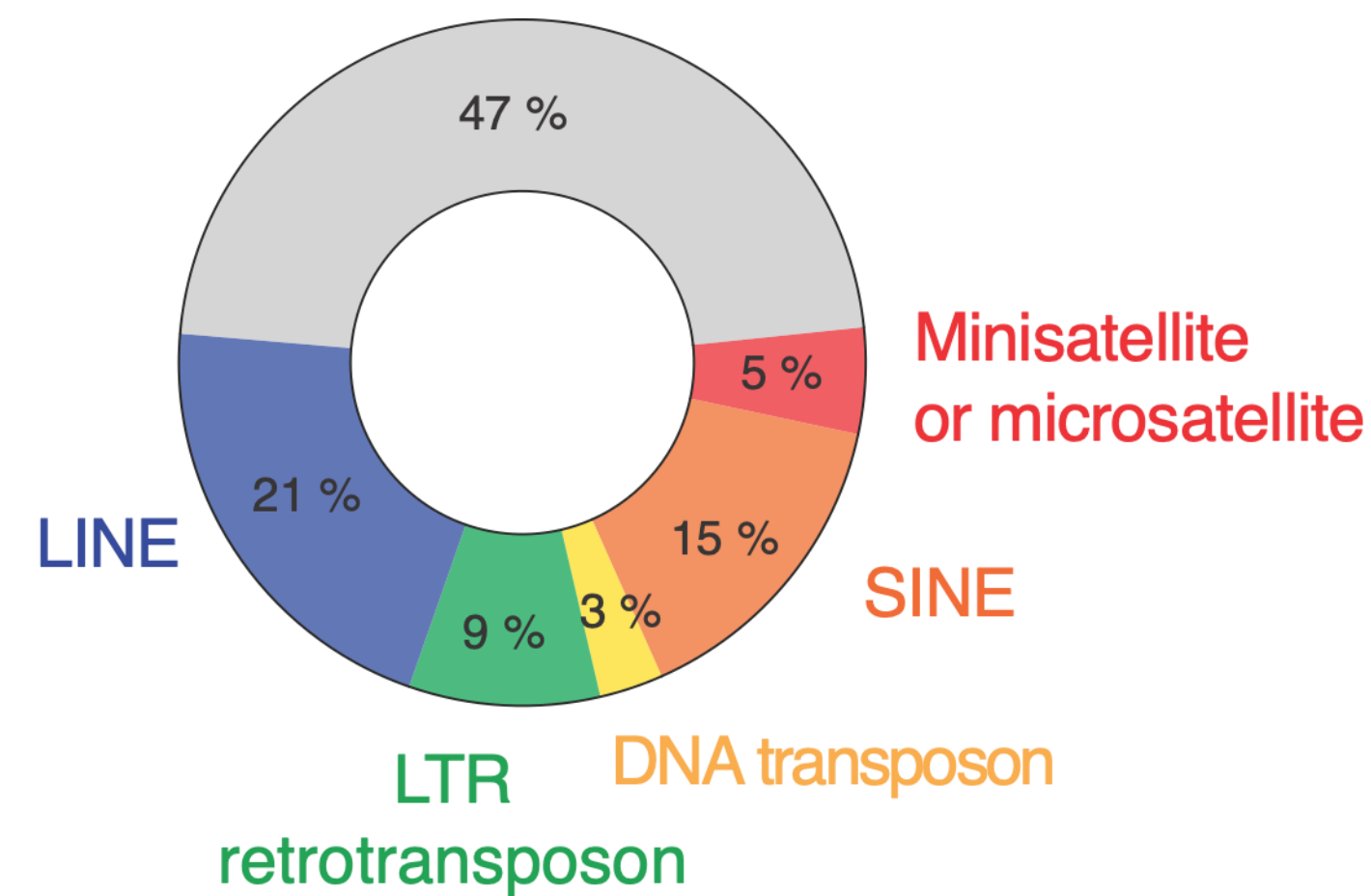Cortone et al (2018) *PLOS Genetics* 14(10):e1007622

# Timeless detects G4s at the fork and coordinates their resolution by DDX11



Unperturbed

G4-detected ahead of fork

Koch Lerner, Holzer et al (2020) *EMBO J* 39(18):e104185

# Which sequences are intrincially capable of forming secondary structures that impede DNA synthesis?

# DNA repeats contribute to gene function, genome structure and evolution, but repeat distribution is highly dependent on sequence
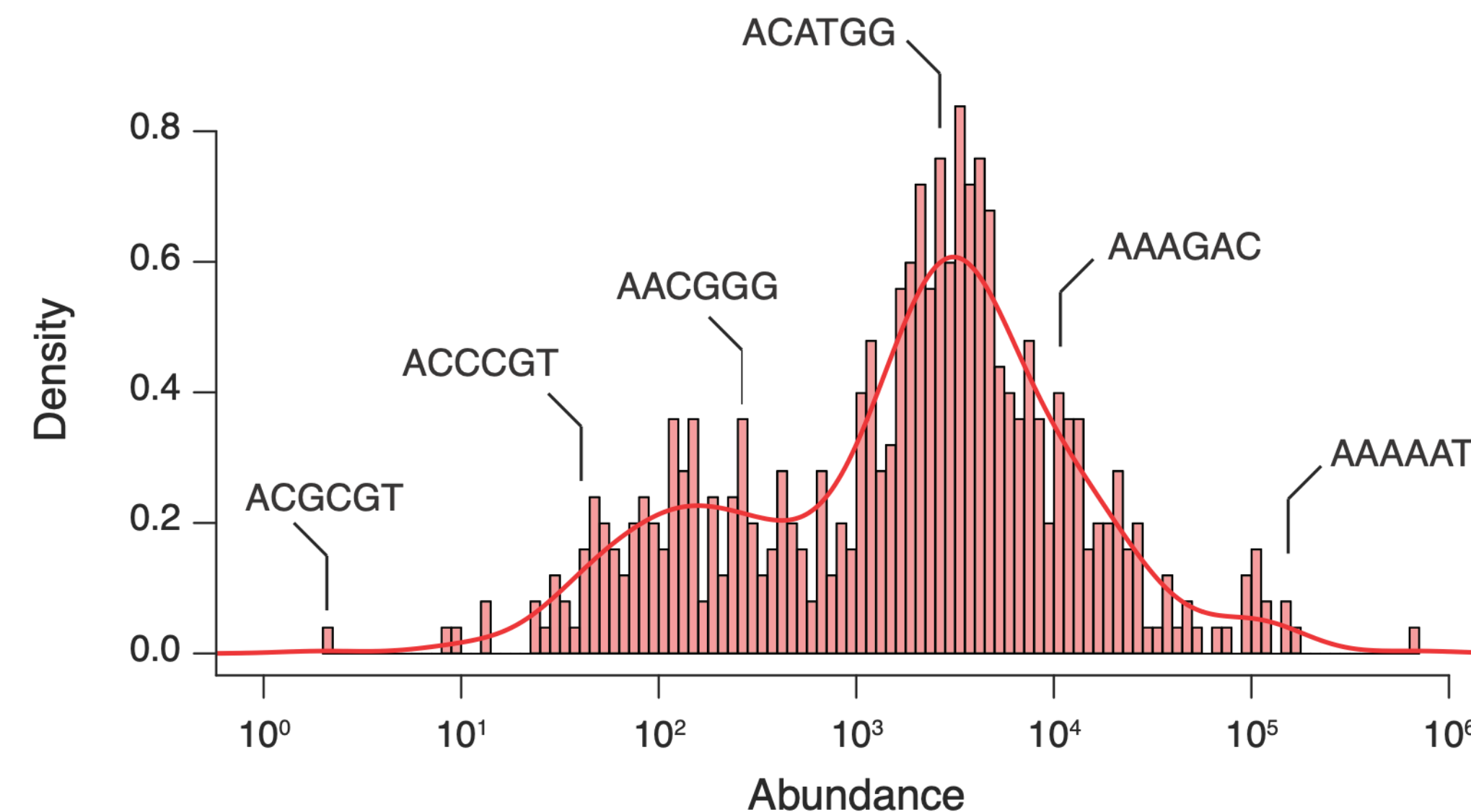
DNA polymerase slippage :



Initiation → Dissociation → Out-of-register realignement → Expansion (or contraction)
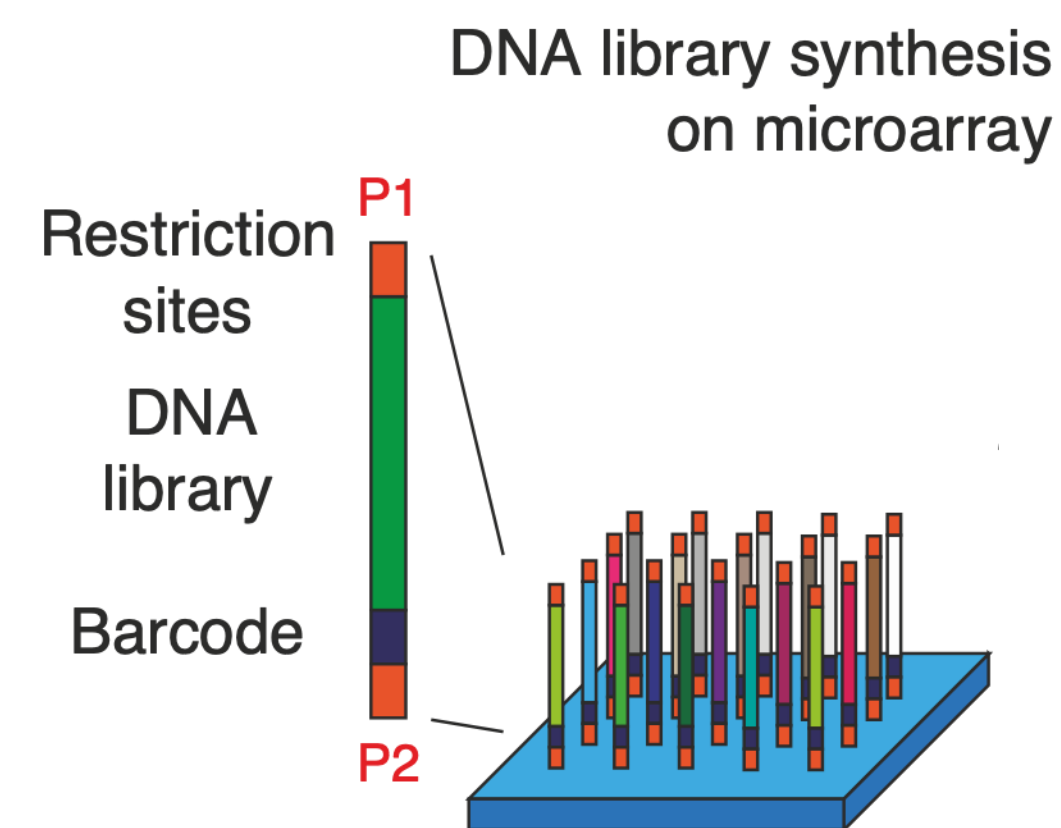
But wide distribution of STR length at equilibrium :



What determines the representation of STRs in the human genome?

Is STR abundance and length determined by DNA polymerase behaviour?

Murat, P, Guilbaud, G, & Sale, JE (2020) *Genome Biology* 21(1):209

DNA library synthesis
on microarray

Restriction sites

P1

DNA library

Barcode

P2



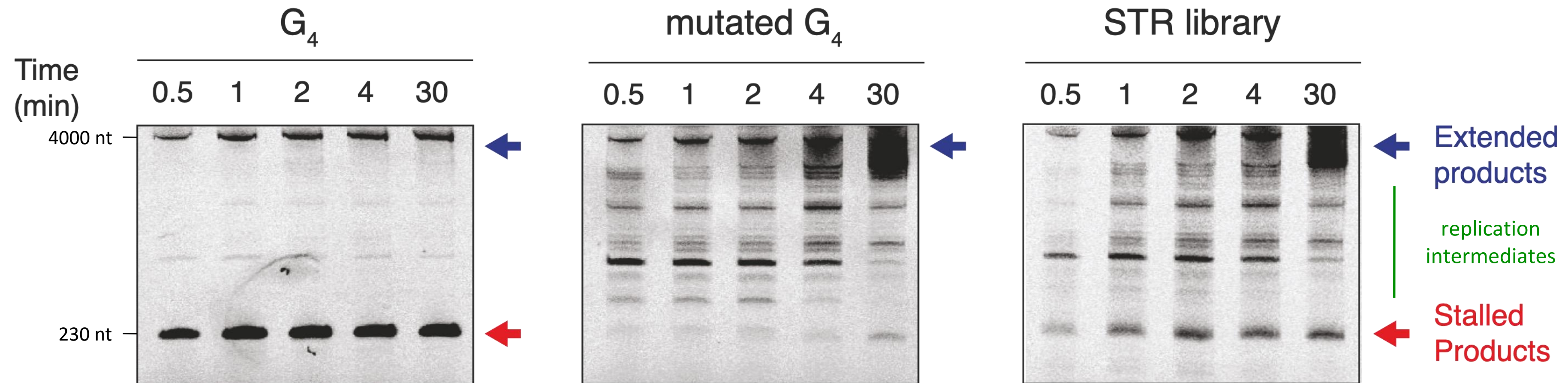*DNA library composition :*

| | |
|---|---|
| All STR permutations (2-6 bp sequences) : | 5,356 sequences |
| in three different lengths (24, 48, 72 nt) | 16,068 sequences |
| | |
| Positive controls | |
| Hairpin, G-quadruplex and I-motif forming sequences | 2,932 sequences |
| | |
| Negative controls | |
| Random sequences of varying GC content | 1,000 sequences |
| | |
| Total : | 20,000 sequences |

Murat, P,  Guilbaud, G, & Sale, JE (2020) *Genome Biology* 21(1):209

# Which sequences are intrinsically able to stall DNA synthesis?



Use of high-throughput sequencing to quantify:
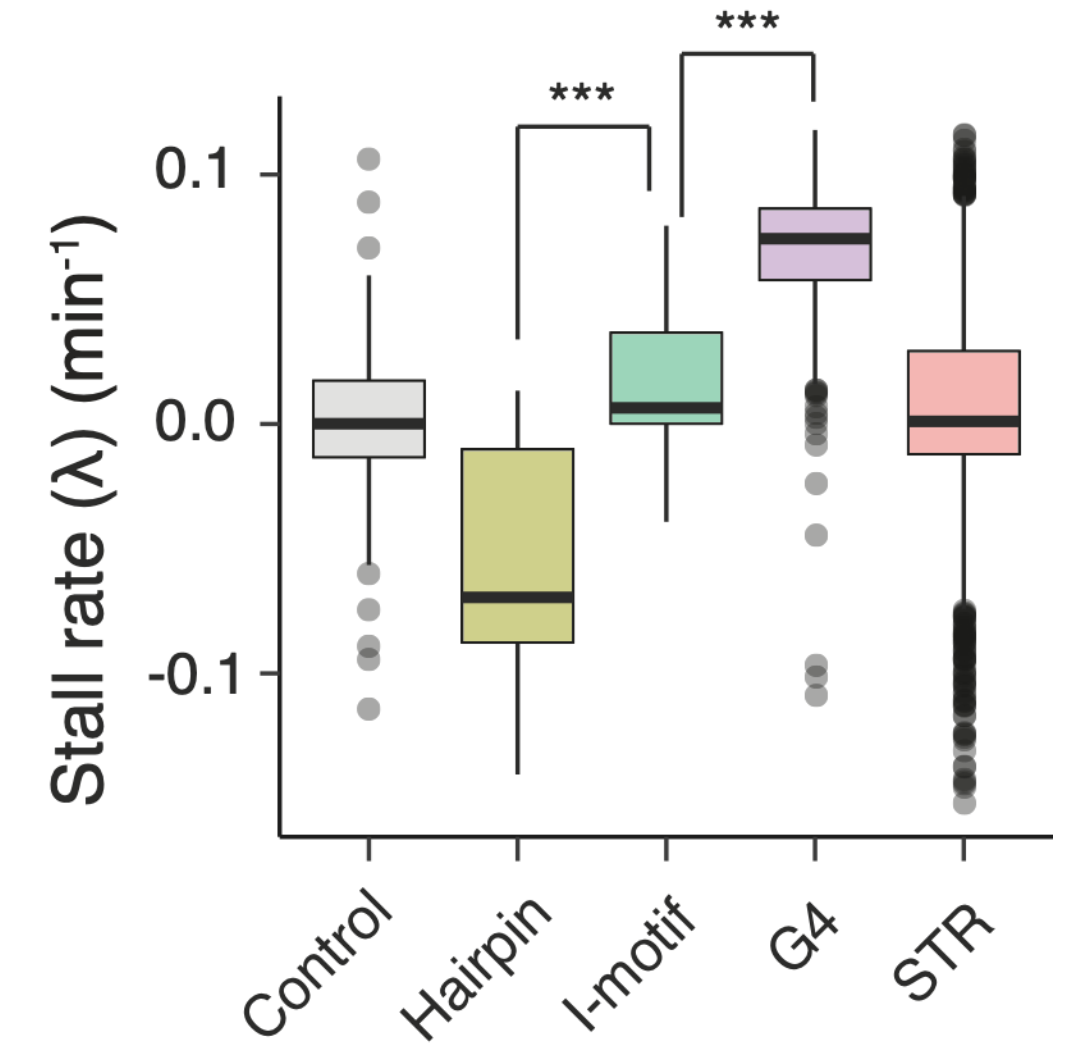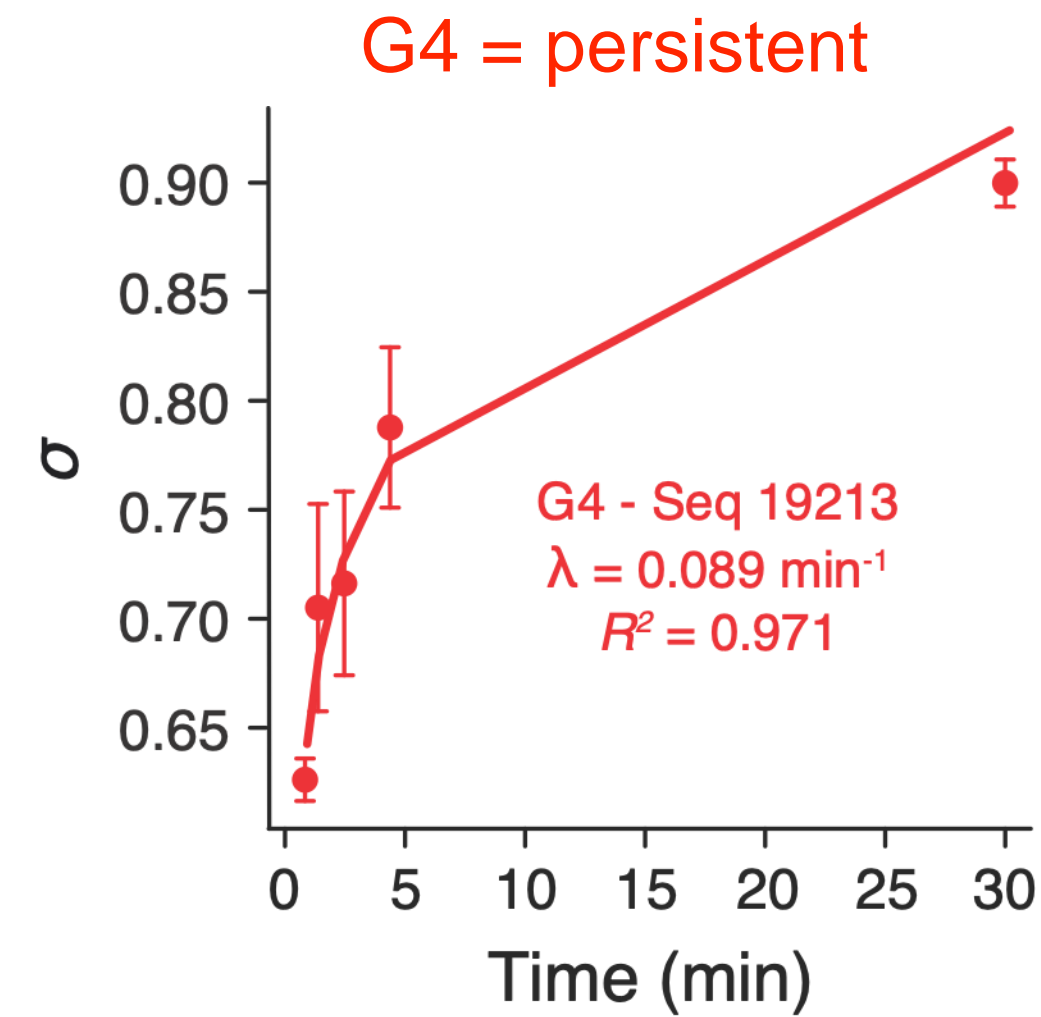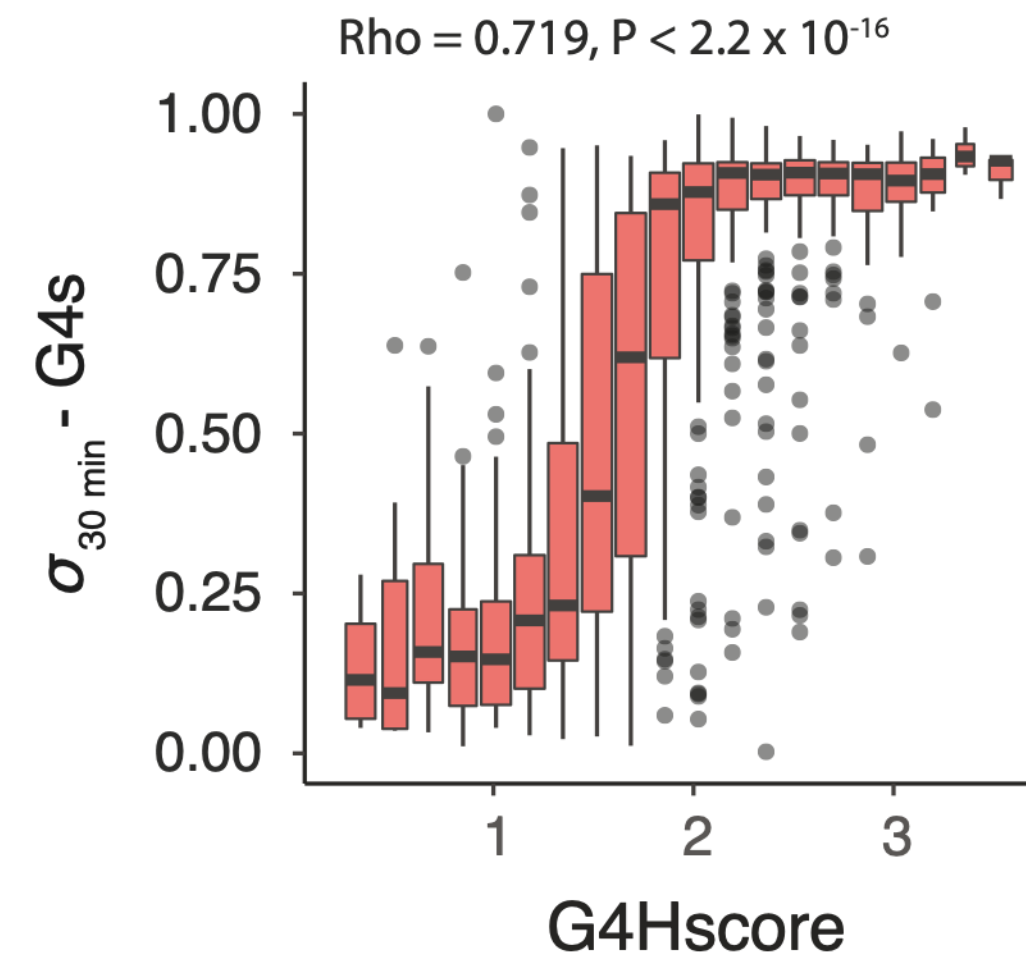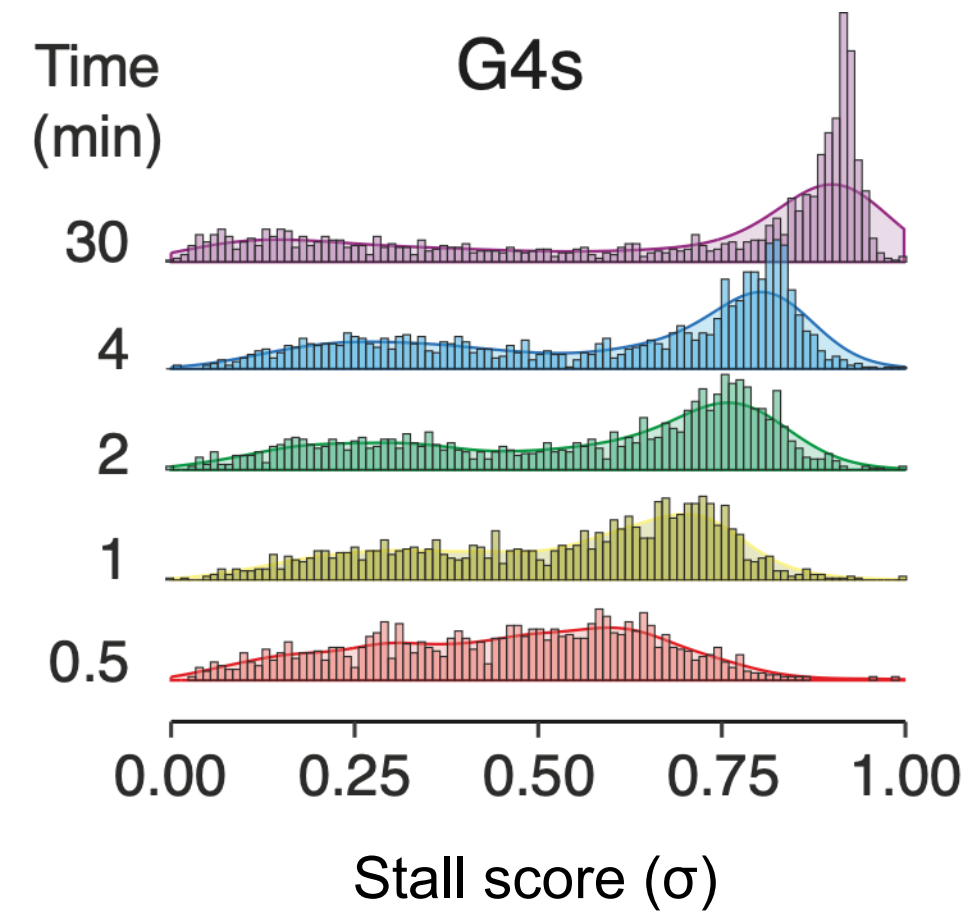
- DNA synthesis **efficiency / stalling**

$$Stall\ score\ _{(t)} = \sigma\ _{(t)} = \frac{\#\ Reads\ in\ the\ stalled\ fraction}{\#\ Reads\ in\ the\ stalled\ and\ extended\ fraction}$$

- DNA synthesis **fidelity**

$$Errors\ _{(t)} = \#\ Mutations\ in\ extended\ fraction$$

Murat, P, Guilbaud, G, & Sale, JE (2020) *Genome Biology* 21(1):209

# Structure-dependent transient and persistent stalling events



can polymerase response
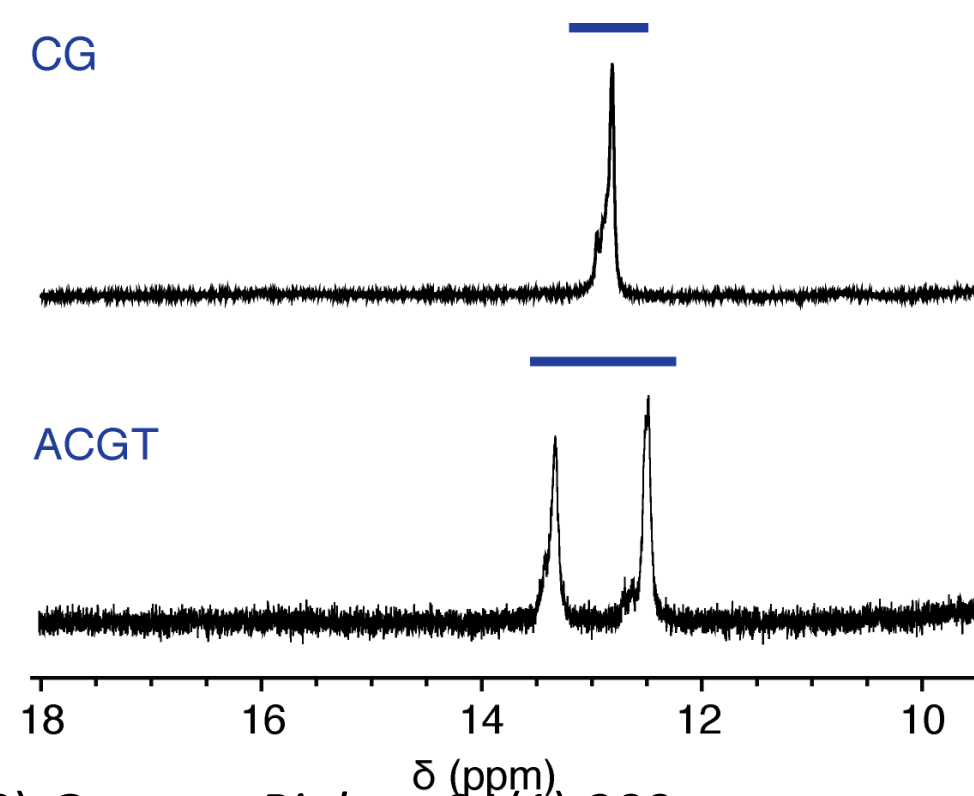be used to to categorise
STRs by structure

Murat, P, Guilbaud, G, & Sale, JE (2020) *Genome Biology* 21(1):209

# A machine learning approach allows structural categorisation of STRs based on polymerase response



accuracy 0.96 +/- 0.03 on test dataset

**Hairpin** - WC base pairing

**I-motif** - Hemiprotonated Cs

**G4** - Hoogsteen base pairing

Murat, P, Guilbaud, G, & Sale, JE (2020) *Genome Biology* 21(1):209

# Polymerase stalling promotes sequence instability



Murat, P, Guilbaud, G, & Sale, JE (2020) *Genome Biology* 21(1):209

# How does structure formation by STRs affect their evolutionary behaviour?

*A high throughput replication assay that:*

- Quantifies the efficiency and fidelity of DNA synthesis at all STR permutations
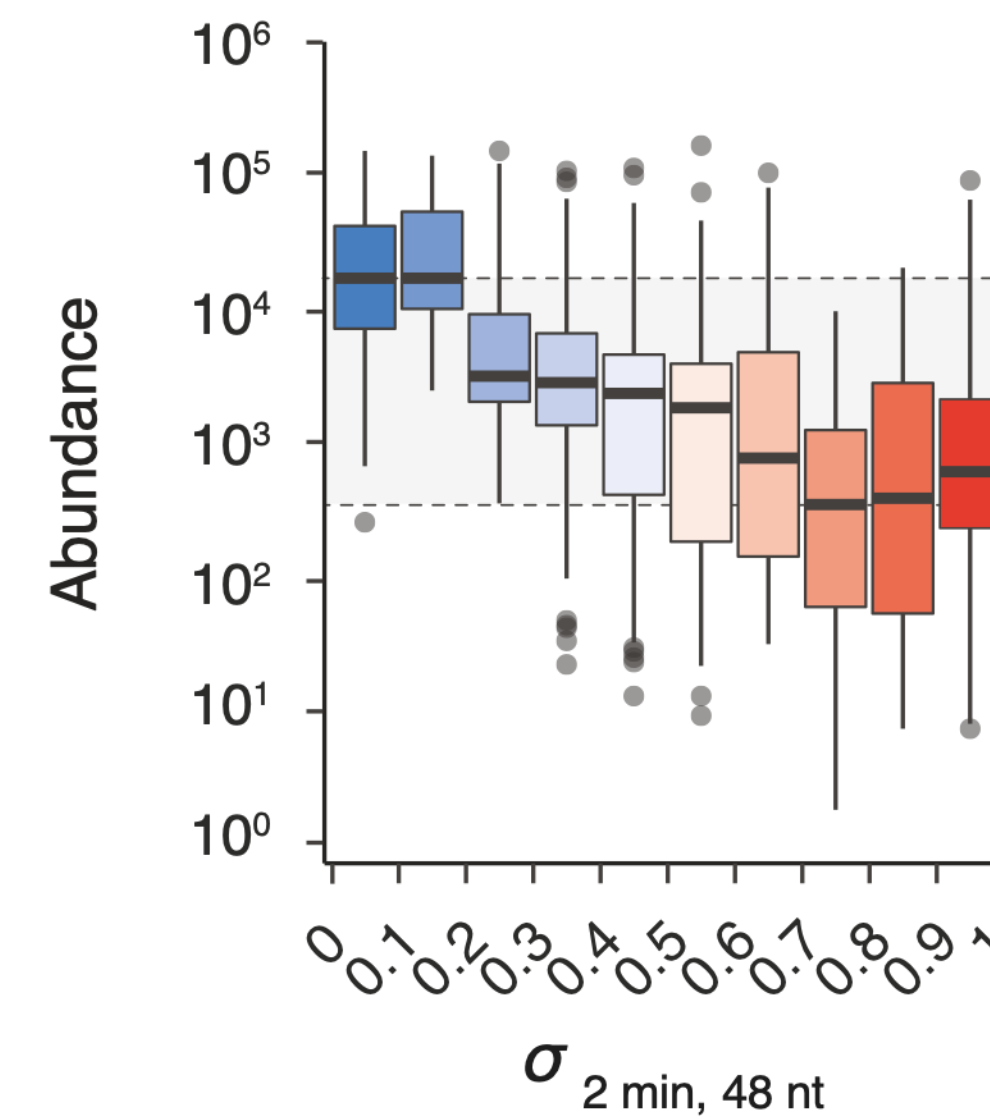
- Infers STR structure from polymerase stalling events

- Establishes general principles for synthesis-dependent STR instability:

<div style="color:red">

Unfolded
repeats
↓
Polymerase
slippage
↓
Length
instability

</div>

<div style="color:blue">

Structured
repeats
↓
Error-prone
DNA synthesis
↓
Sequence
instability

</div>

## Do these observations have any *in vivo* correlate?

Murat, P, Guilbaud, G, & Sale, JE (2020) *Genome Biology* 21(1):209

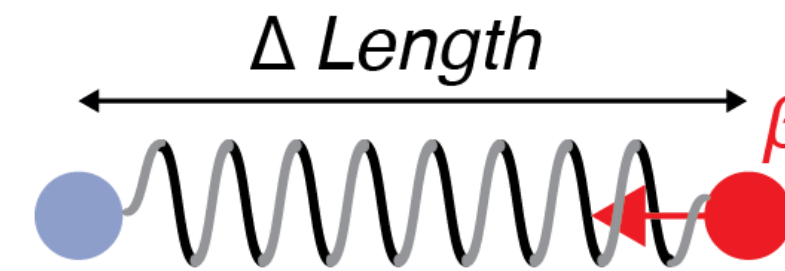# DNA polymerase stalling at DNA structures predicts STR abundance and length in eukaryotic genomes

Human genome (4,500,000 STR loci)



Same trends observed for the Mouse, Chicken, Zebrafish, Fly and Yeast genomes

Murat, P, Guilbaud, G, & Sale, JE (2020) Genome Biology 21(1):209

# Expansion is favoured and less constrained in weakly stalling repeats



Multi-step
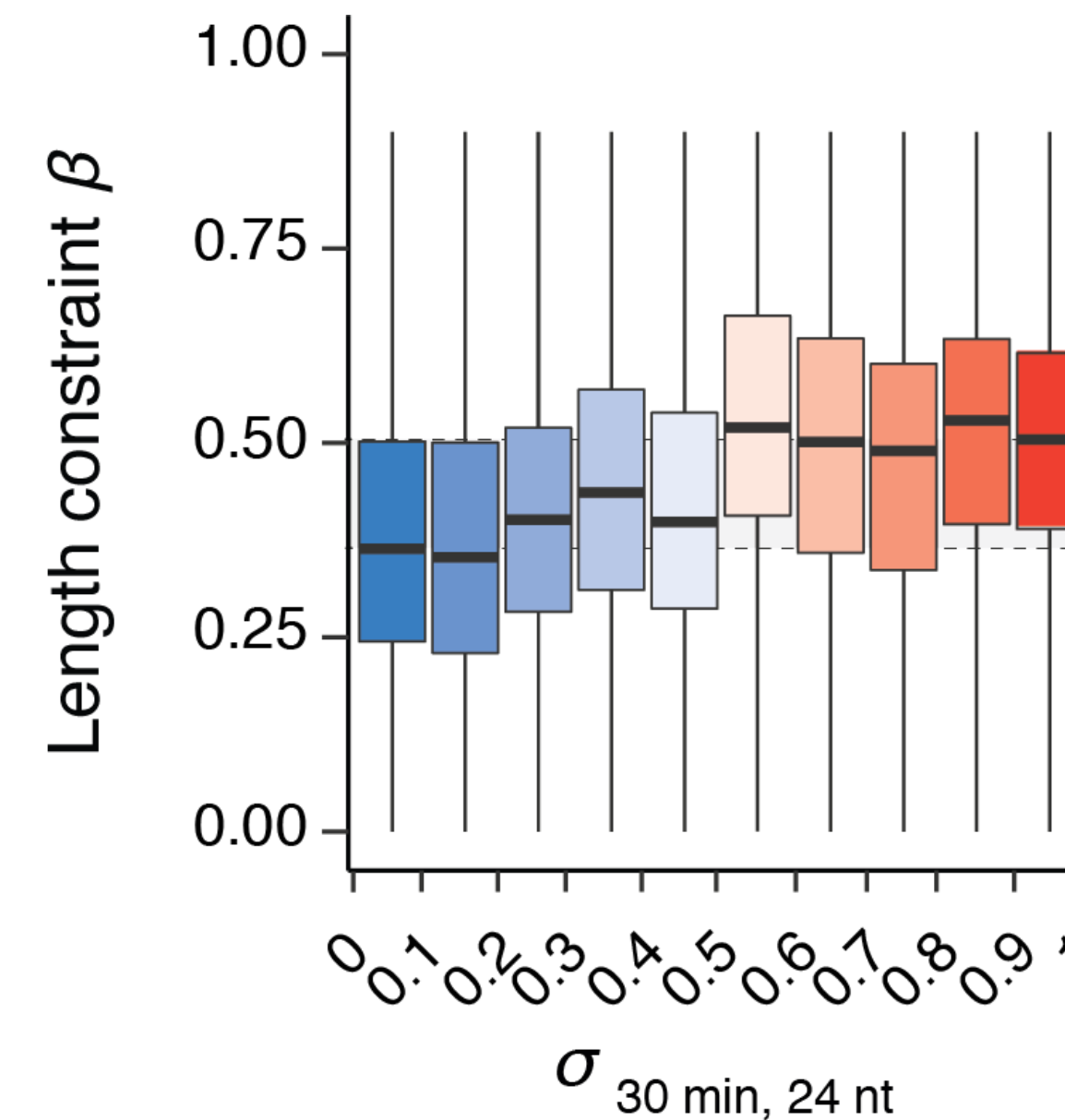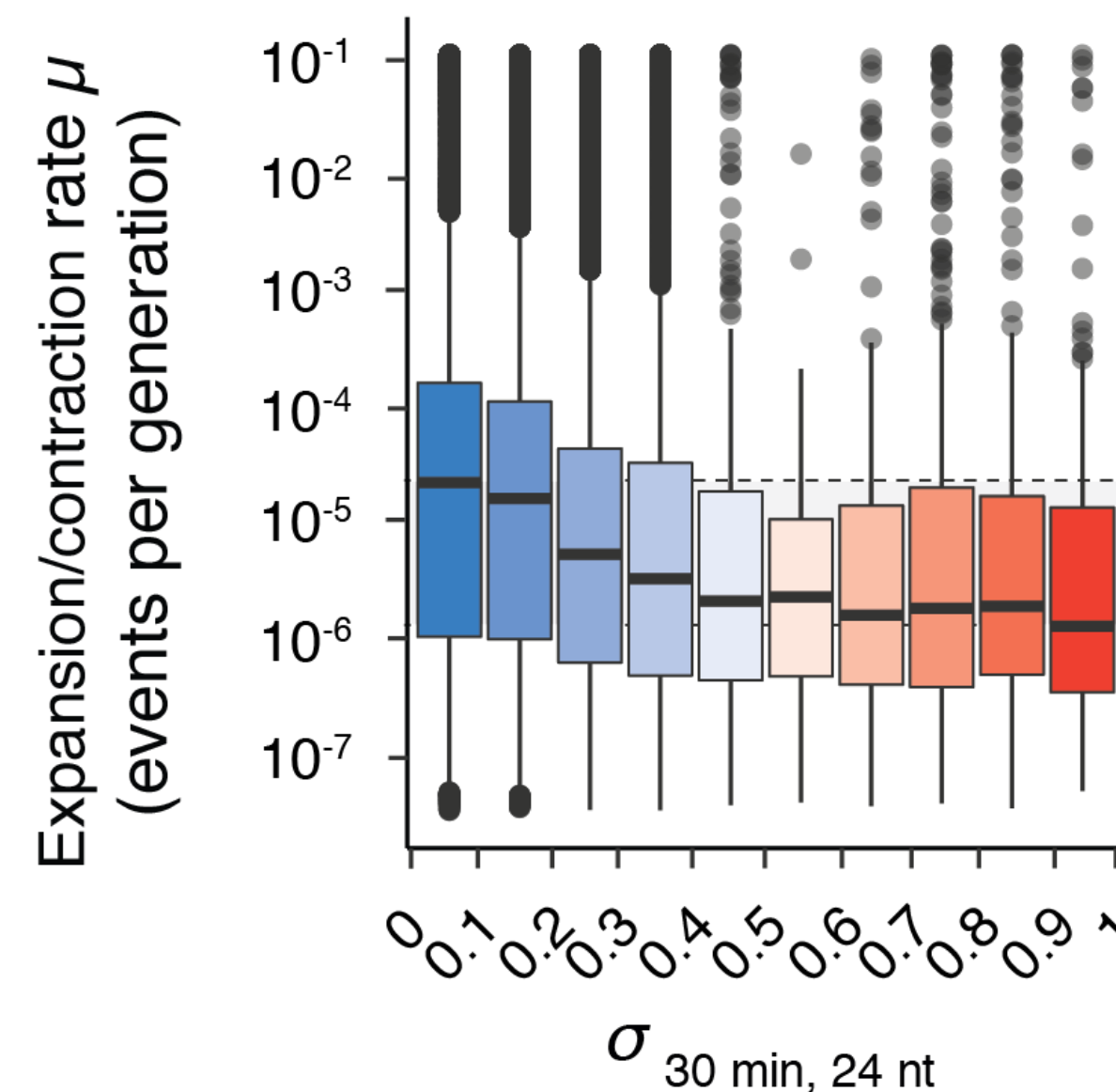Ornstein-Uhlenbeck
process

Δ *Length*

β

Length at equilibrium

Length after expansion

$\mu = \dfrac{\mathrm{d}Length}{\mathrm{d}t}$ : mutation rate

β : Length constraint

Mutation rate

Evolutionary constraint

What is the basis for the increased length constraint on structured STRs?

Murat, P, Guilbaud, G, & Sale, JE (2020) *Genome Biology* 21(1):209

Gymrek et al. Nature Gen. 2017, 49, 1495-1501

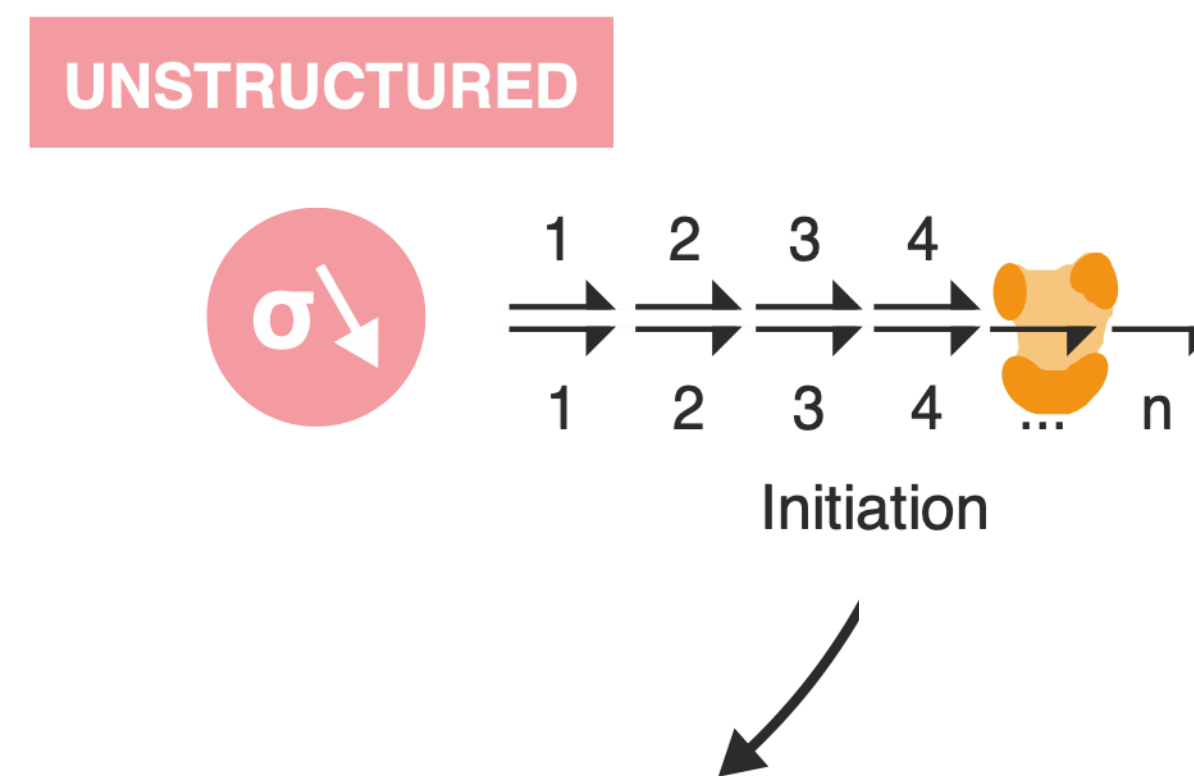# Structured STRs are prone to point mutation in the human genome: a mechanism for length constraint?



Highly stalling (high σ) sequences are more prone to point mutation
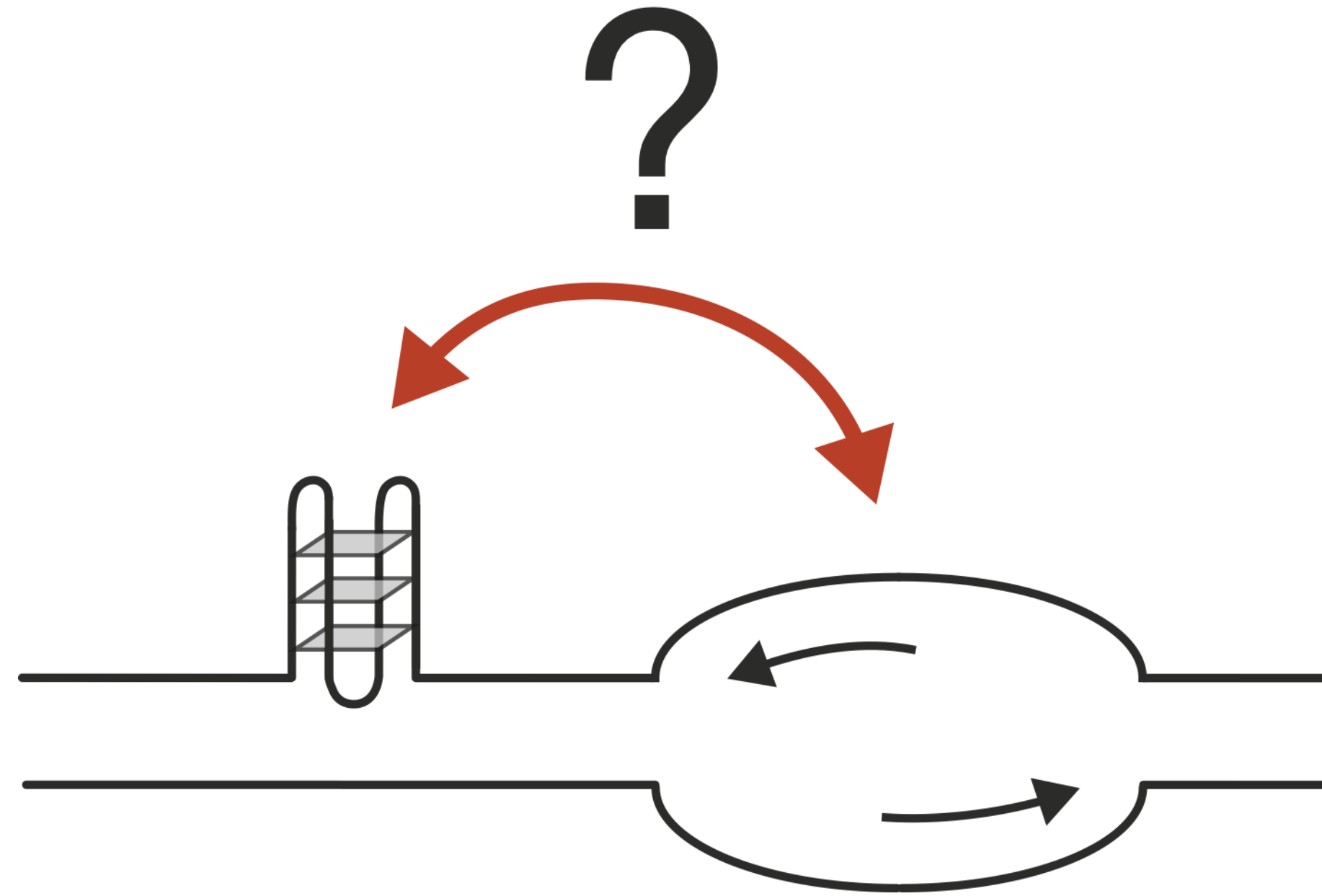
Density of germline mutations in the vicinity of STRs :

Murat, P, Guilbaud, G, & Sale, JE (2020) *Genome Biology* 21(1):209

# DNA polymerase stalling at structured DNA constrains the expansion of Short Tandem Repeats
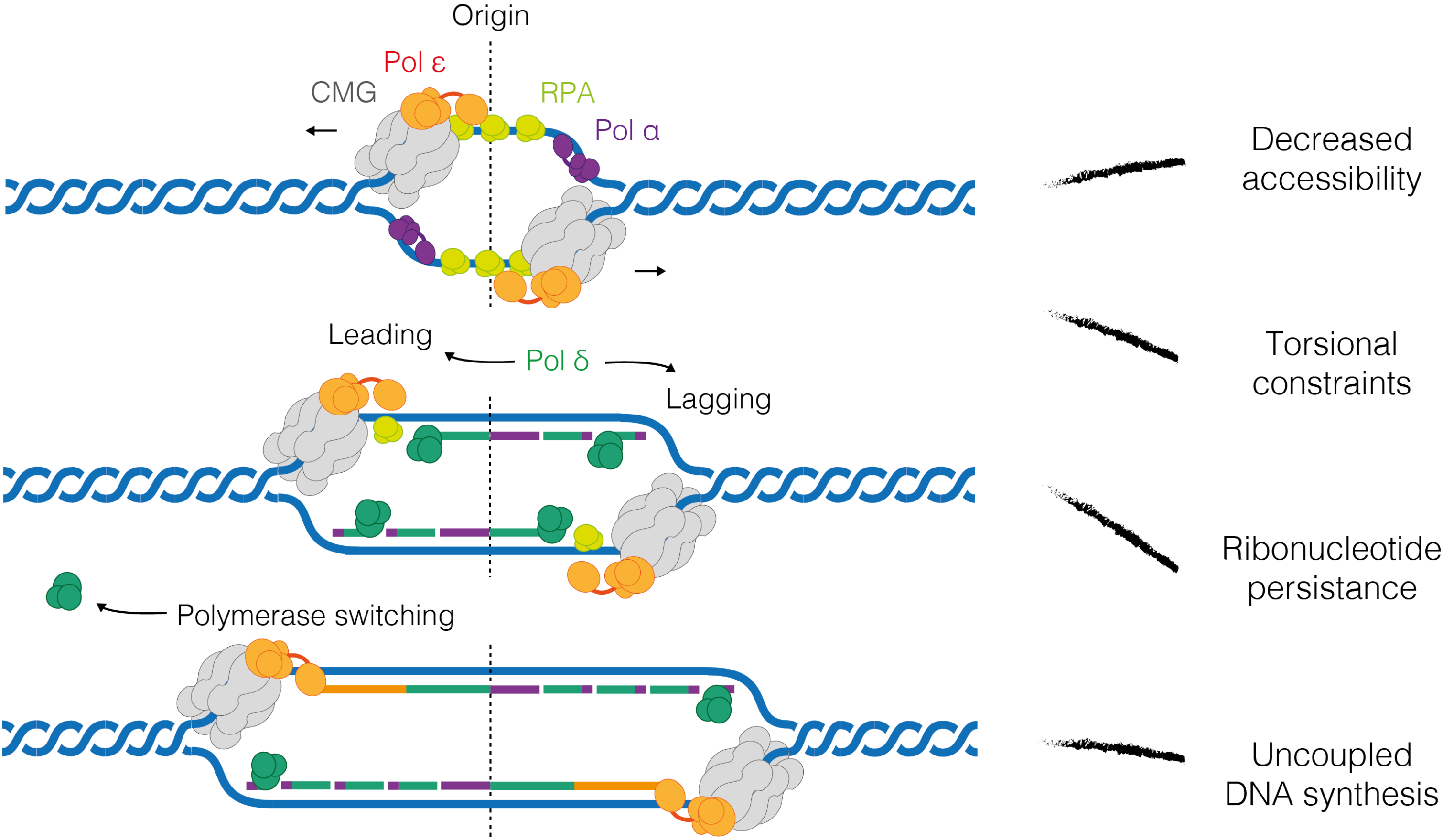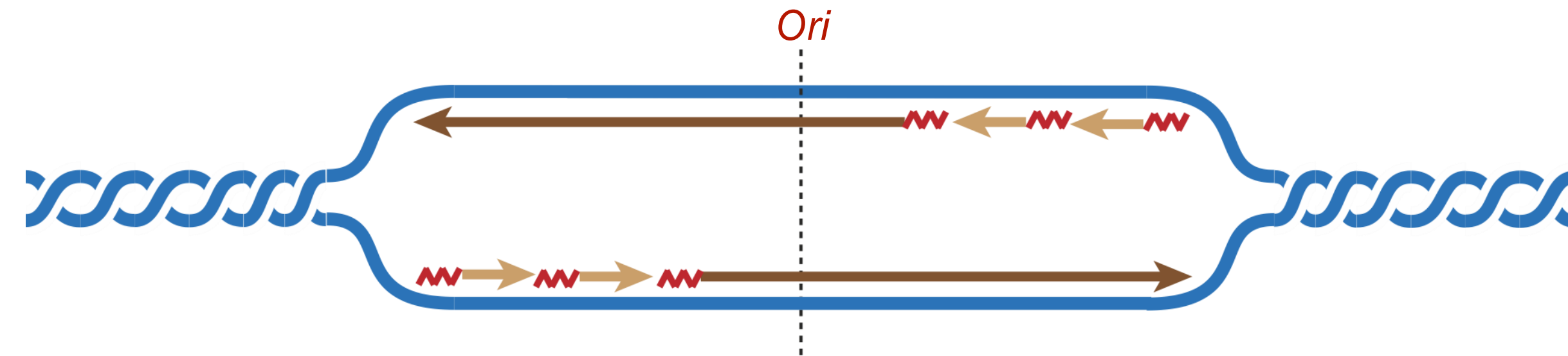
Why are G4s associated with replication origins?

# Hypothesis: replication origins are hotspots for mutagenesis

Origin

Pol ε

CMG

RPA

Pol α

Decreased accessibility

Leading

Pol δ

Lagging

Torsional constraints

Polymerase switching

Ribonucleotide persistance

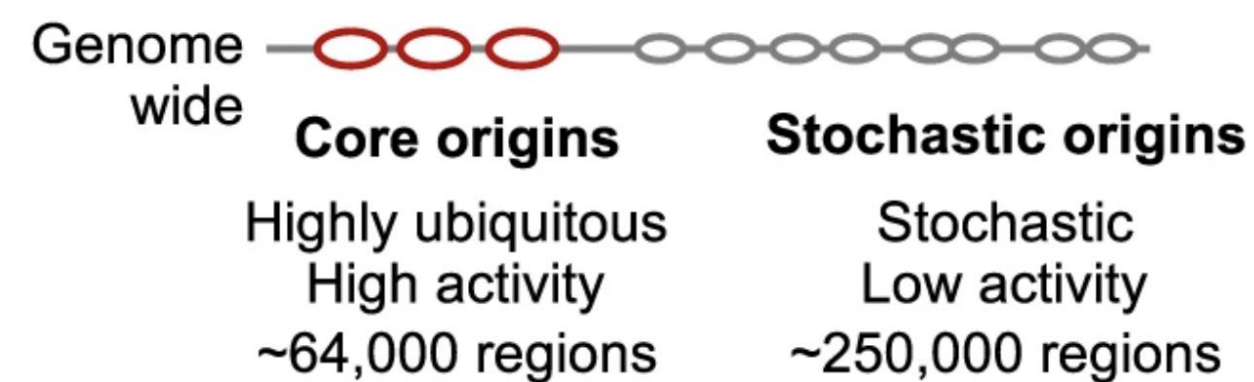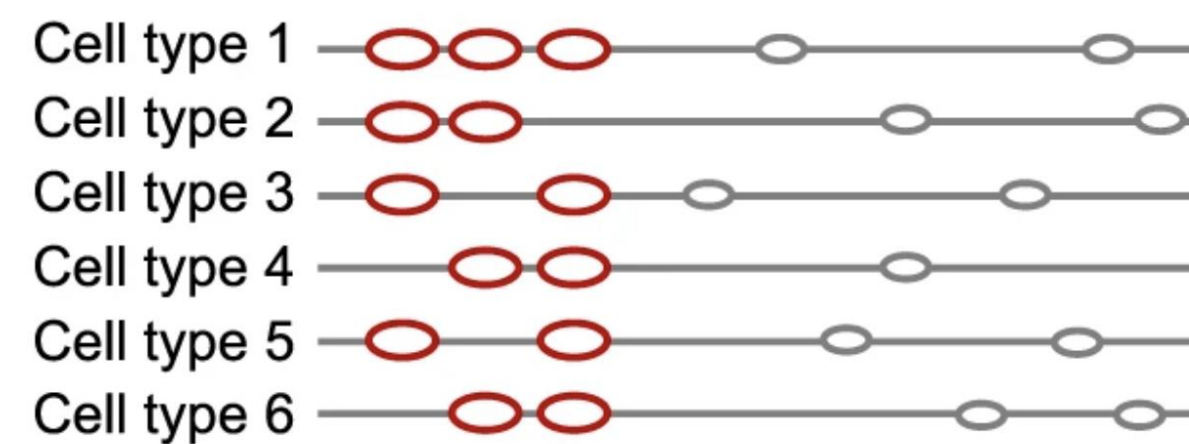Uncoupled DNA synthesis

# Identifying origins of replication in human cells



*Ori*

**Short Nascent Strand sequencing (SNS-seq)**

> 300,000 origins
1-4kb resolution

Origins identified from all cell-types

Cell type 1
Cell type 2
Cell type 3
Cell type 4
Cell type 5
Cell type 6

Genome wide

**Core origins**
Highly ubiquitous
High activity
~64,000 regions

**Stochastic origins**
Stochastic
Low activity
~250,000 regions

Akerman et al. 2020

**Okazaki fragment sequencing (Ok-seq)**

C

W

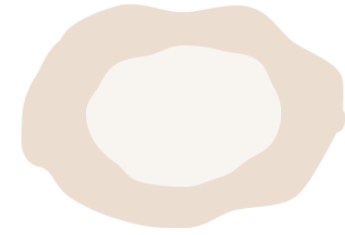1

RFD

−1

Chr2    120                    125 Mb

~ 22,000 initiation zones
~50kb resolution

Petryk et al. 2015

# Ini-seq 2: a method to map both replication origin location and efficiency



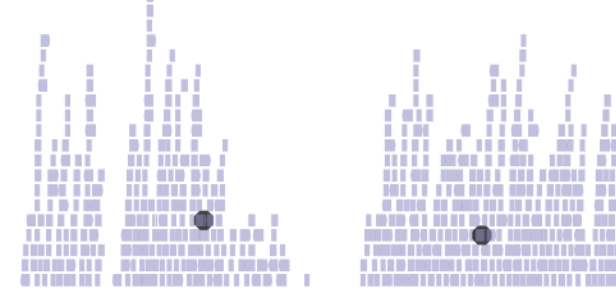Cystosolic extract from proliferating cells
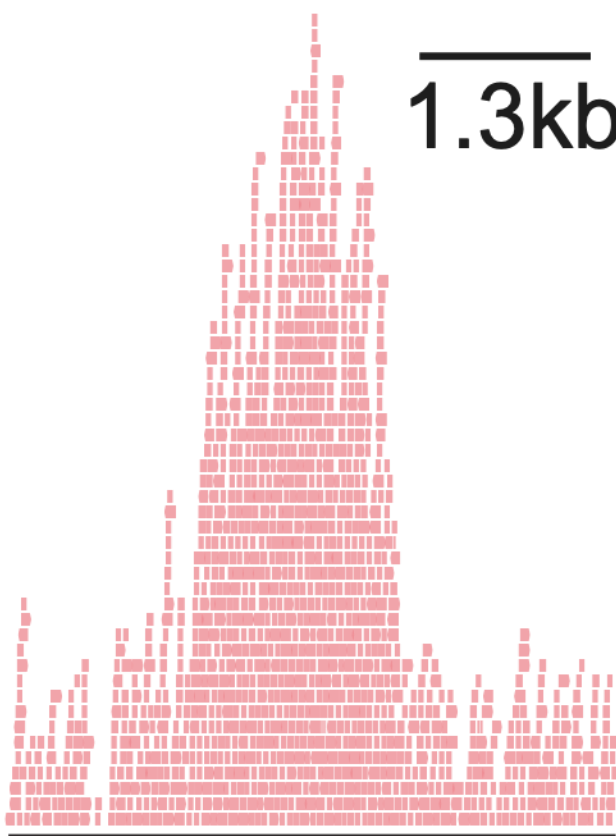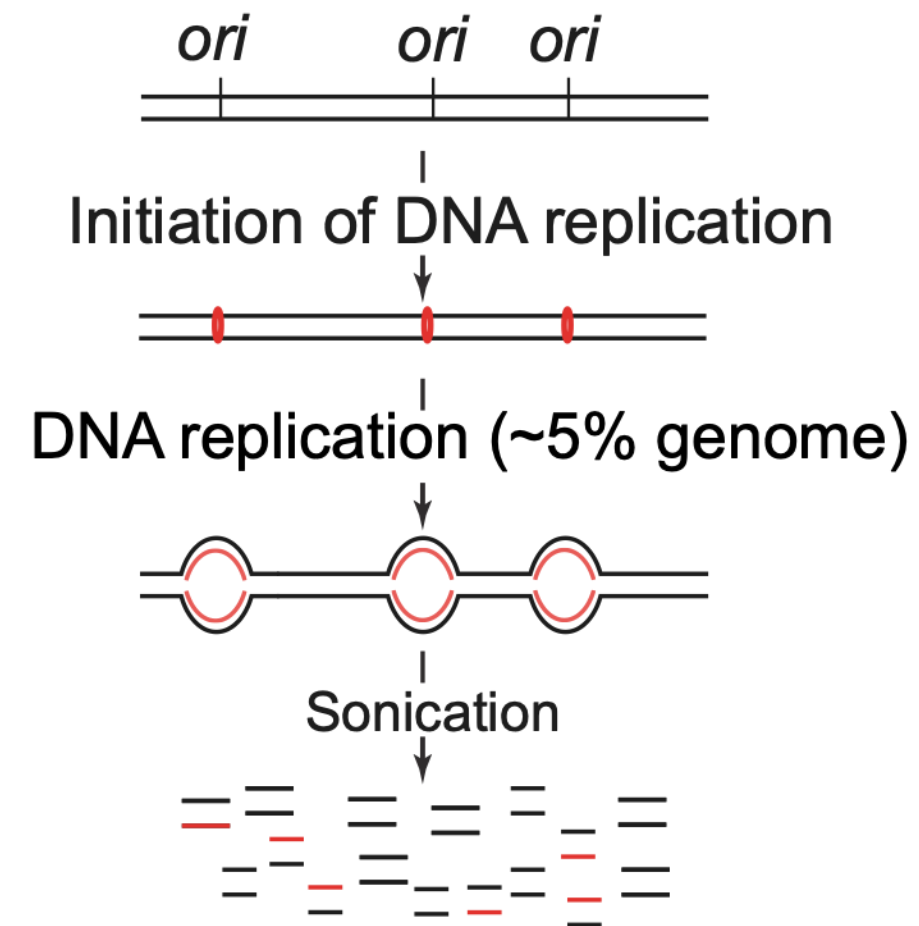
+ BrdUTP

Nuclei arrested in G1/S tranistion

*ori*   *ori*   *ori*

Initiation of DNA replication

DNA replication (~5% genome)

Sonication

$Cs_2SO_4$ refractive index

1.3660
1.3690 } LL fraction

1.3685
1.3715 } HL fraction

## *TOP1*

LL fraction

1.3kb

HL fraction

chr20   41.02

**23,905 called origins**

High: 1/3 = 7,968
Medium: 1/3 = 7,968
Low: 1/3 = 7,969

Count (x100)

15

10

5

0

-1   0   1   2

Z-score origin efficiency

**322,603 human origins**

251,956
**stochastic**
(cell type dependent)

23,905
**constitutive**
(in early domains)

46,742
**core**
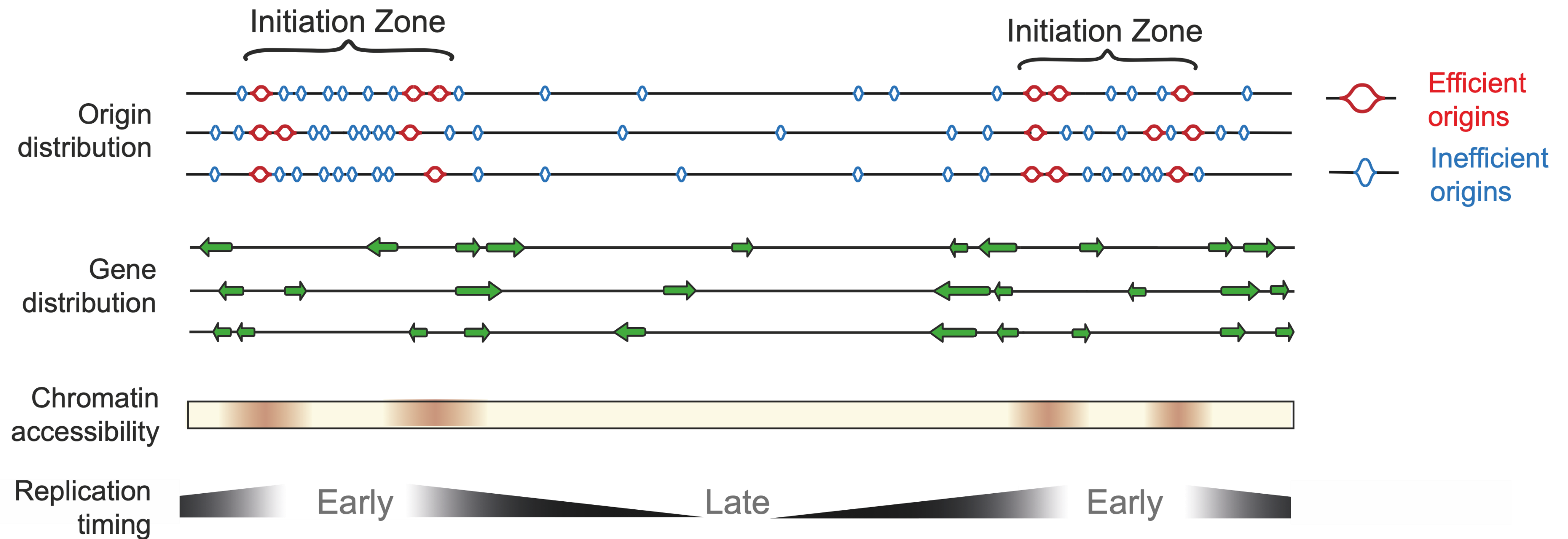(conserved in all cell types)

Datasets:
Akerman et al. Nature Commun. 2020, 11, 1-15
Guilbaud et al. (2022) Nucleic Acids Res. 50, 7436 - 7450

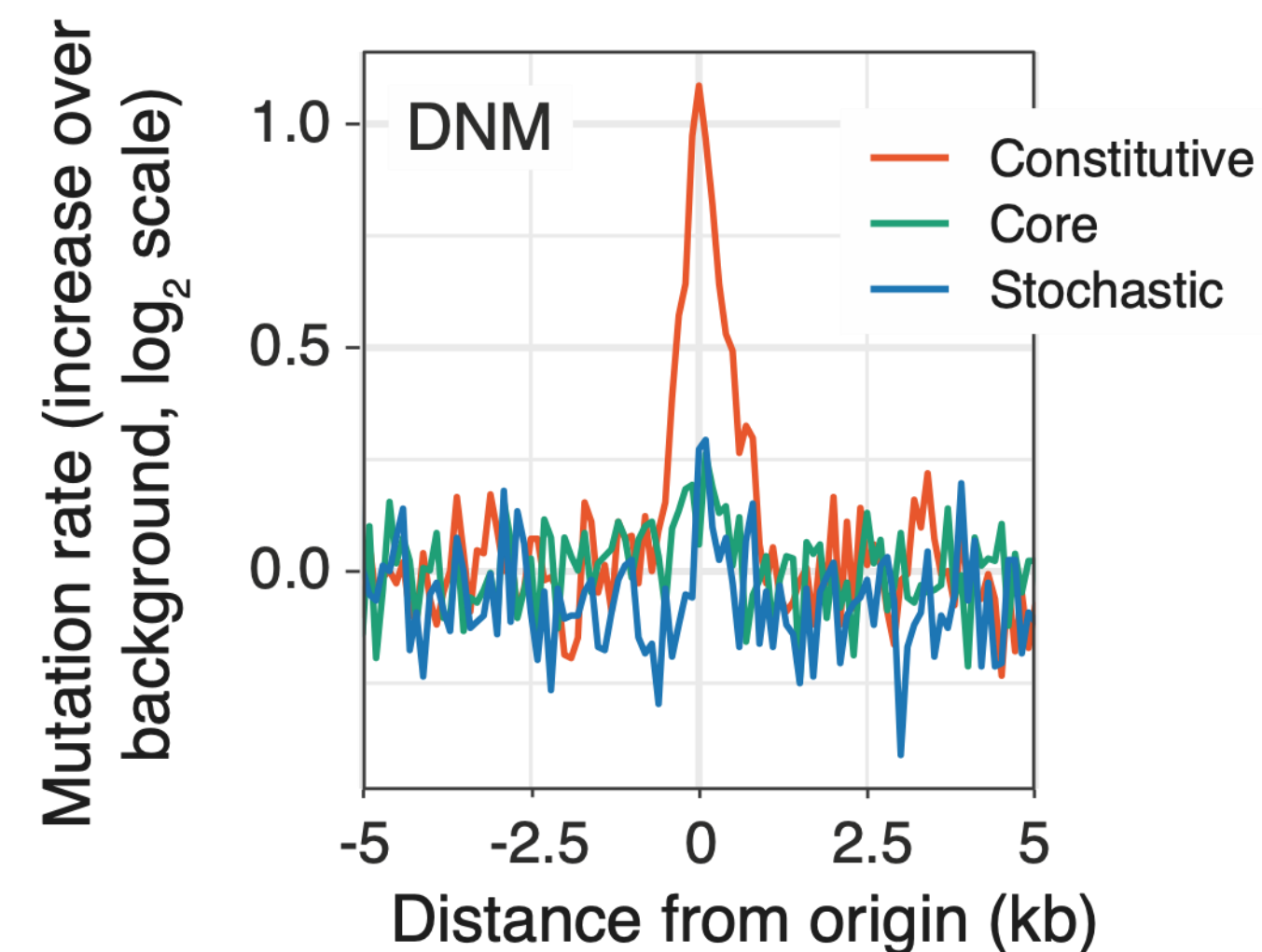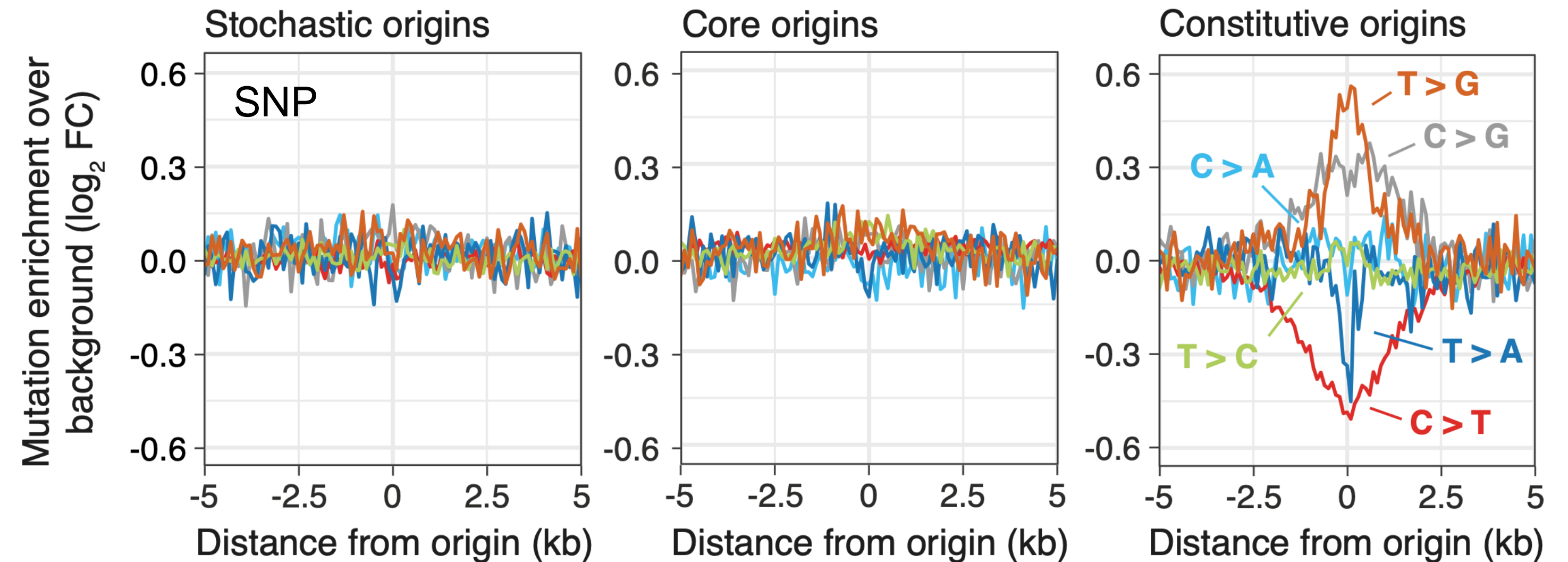Guilbaud et al. (2022) Nucleic Acids Res. 50, 7436 - 7450

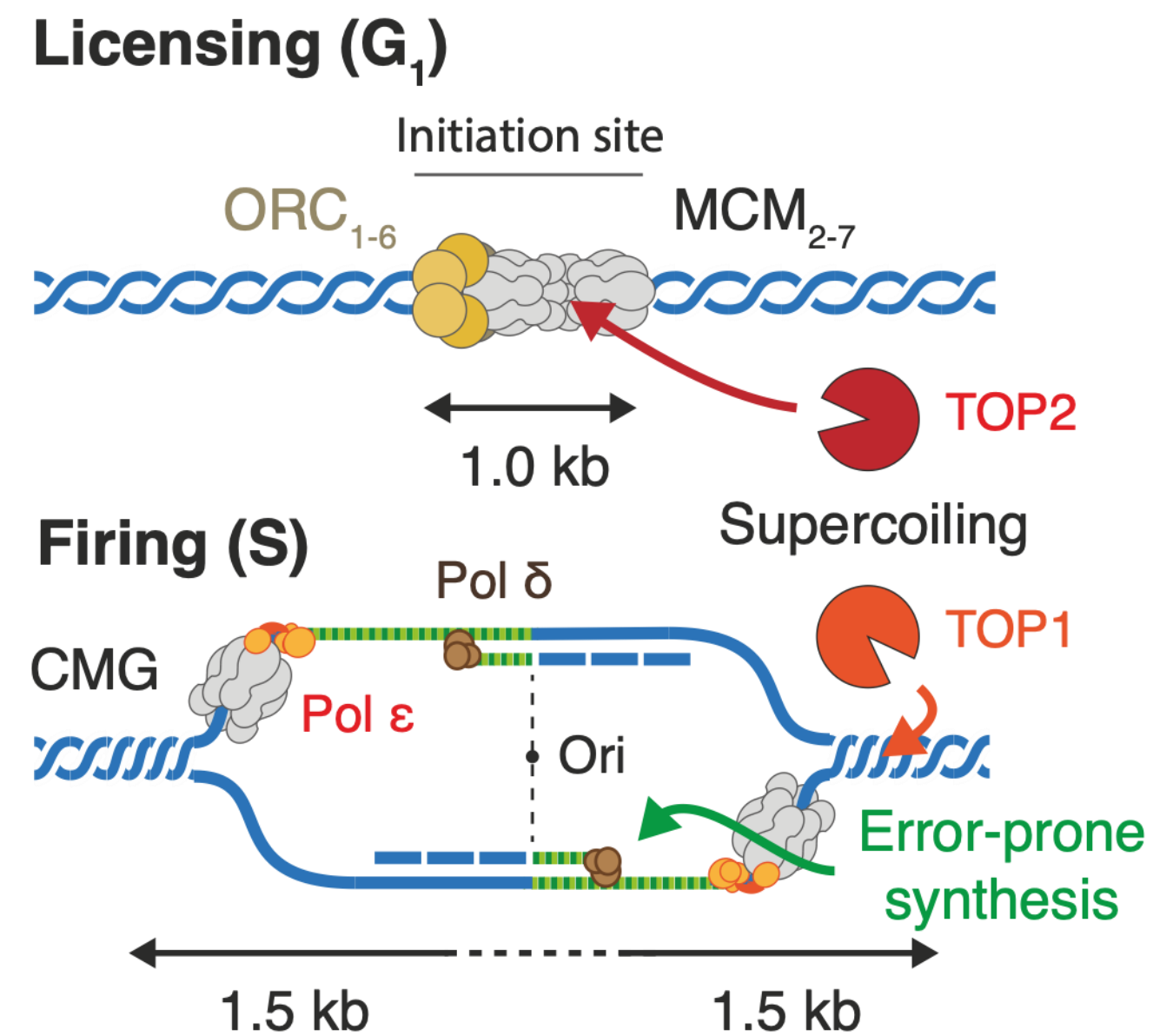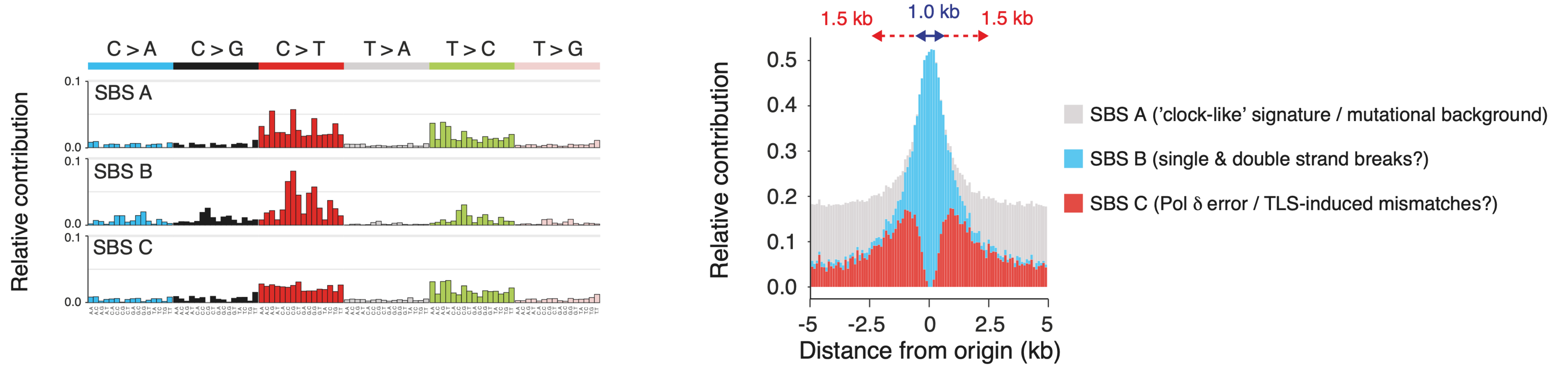# Organisation of replication origins by efficiency in the human genome



Guilbaud et al. (2022) Nucleic Acids Res. 50, 7436 - 7450

# The mutational footprint of replication initiation is revealed at 'constitutive' ini-seq origins



322,603 human origins

251,956
**stochastic**
(cell type
dependent)

23,905
**constitutive**
(in early
domains)

46,742
**core**
(conserved in all cell types)

Akerman et al. Nature Commun. 2020, 11, 1-15
Guilbaud et al. (2022) Nucleic Acids Res. 50, 7436 - 7450

Murat et al., (2022) *Sci. Adv.* 8(45):eadd3686

# The mutational signature of replication origins



Murat et al., (2022) *Sci. Adv.* 8(45):eadd3686

# Ori SBS B reflects double strand breaks at replication origins



1.5 kb 1.0 kb 1.5 kb

SBS A ('clock-like' signature / mutational background)

SBS B (single & double strand breaks?)

SBS C (Pol δ error / TLS-induced mismatches?)

**Licensing (G₁)**

Initiation site

ORC₁₋₆ MCM₂₋₇

TOP2

1.0 kb

Supercoiling

**Firing (S)**

Pol δ

CMG

Pol ε

TOP1

Ori

Error-prone synthesis

1.5 kb 1.5 kb

chr6 : 4,012,000- 4,030,000 (20 kb)

DSB Induce-seq

ATAC-seq

TOP2B CUT&RUN

Gene

Origin

PRPF4B

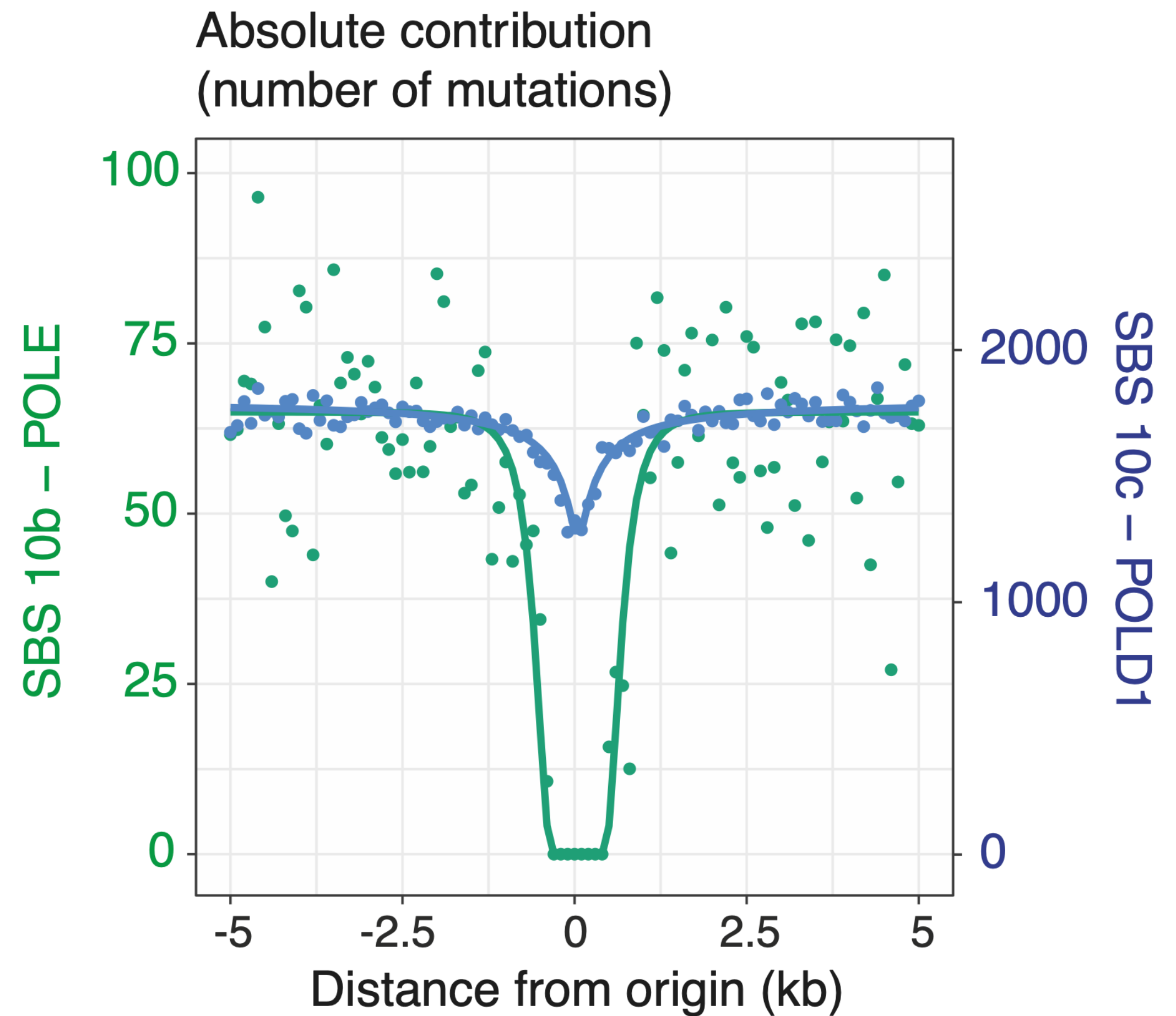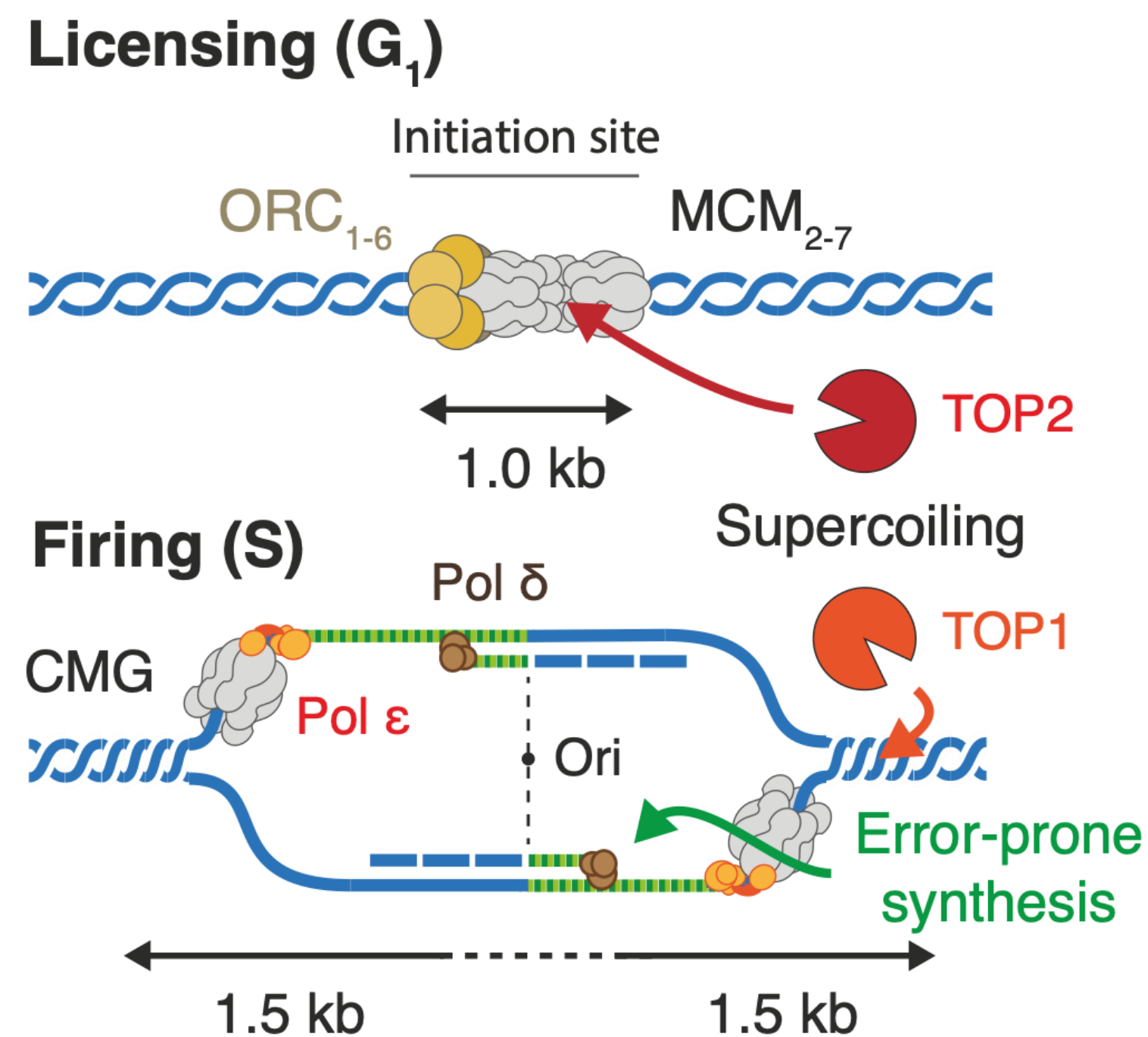## H9 hES cells

with Simon Reed & Pat van Eijk
INDUCE-seq: Dobbs et al. (2022), Nat Comms13:3989

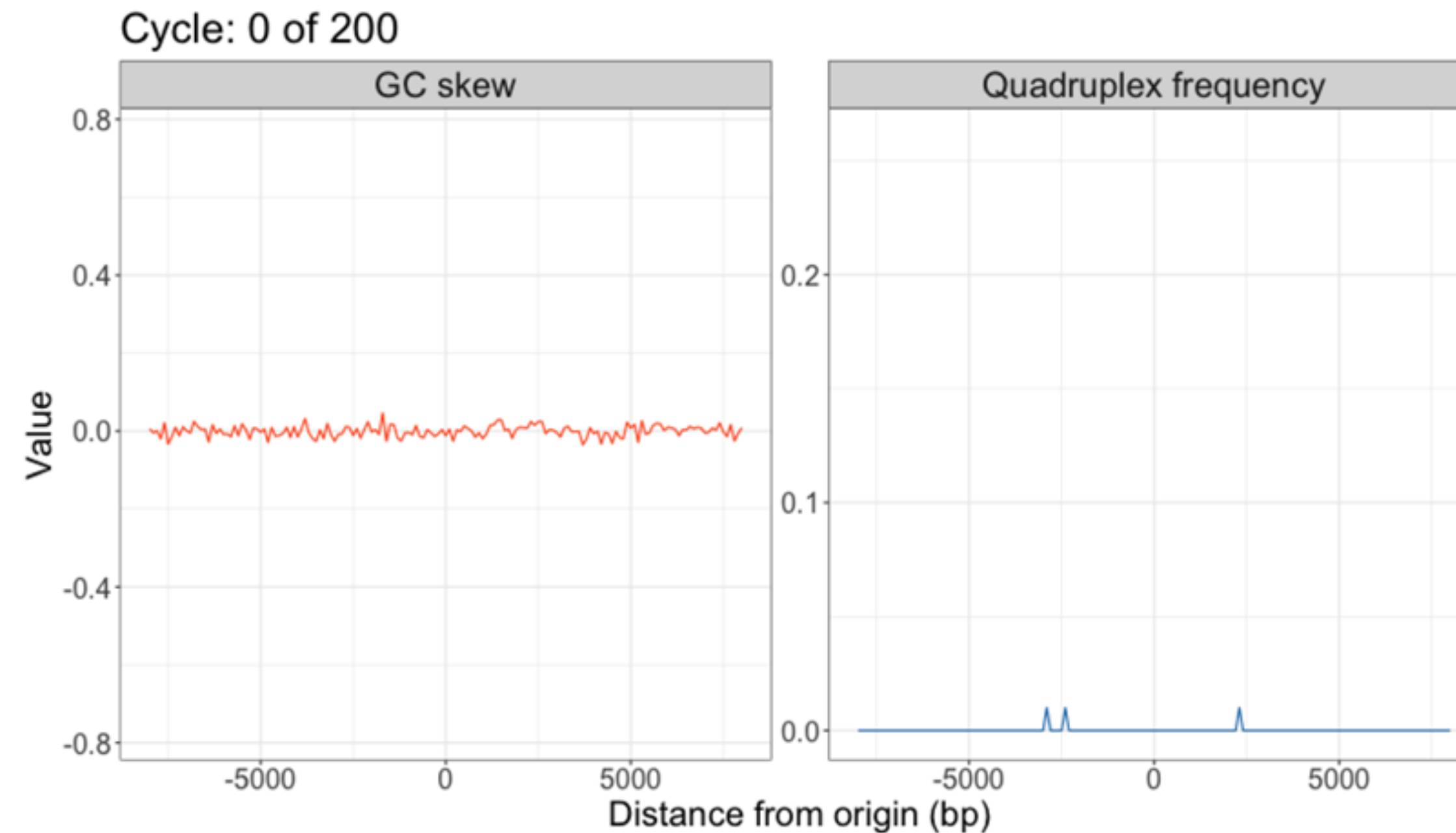DSB count (per origin)

Origin efficiency: low, medium, high

# The mutagenic signatures of replication origins are sufficient to create the sequence environment observed at efficient origins
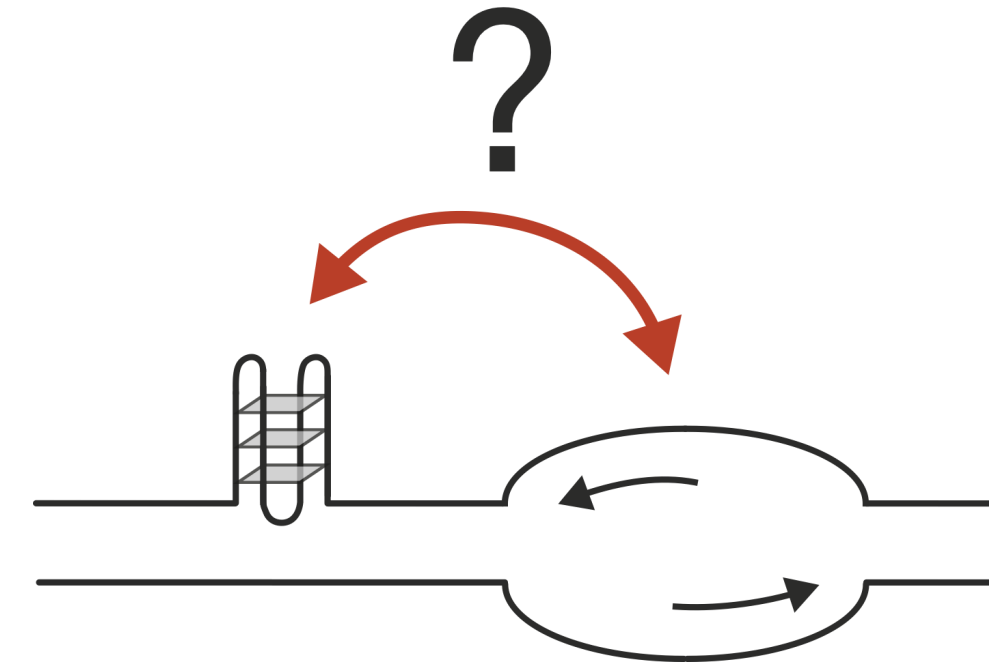


Murat et al., (2022) *Sci. Adv.* 8(45):eadd3686

# The mutagenic signatures of replication origins are sufficient to create the sequence environment observed at efficient origins
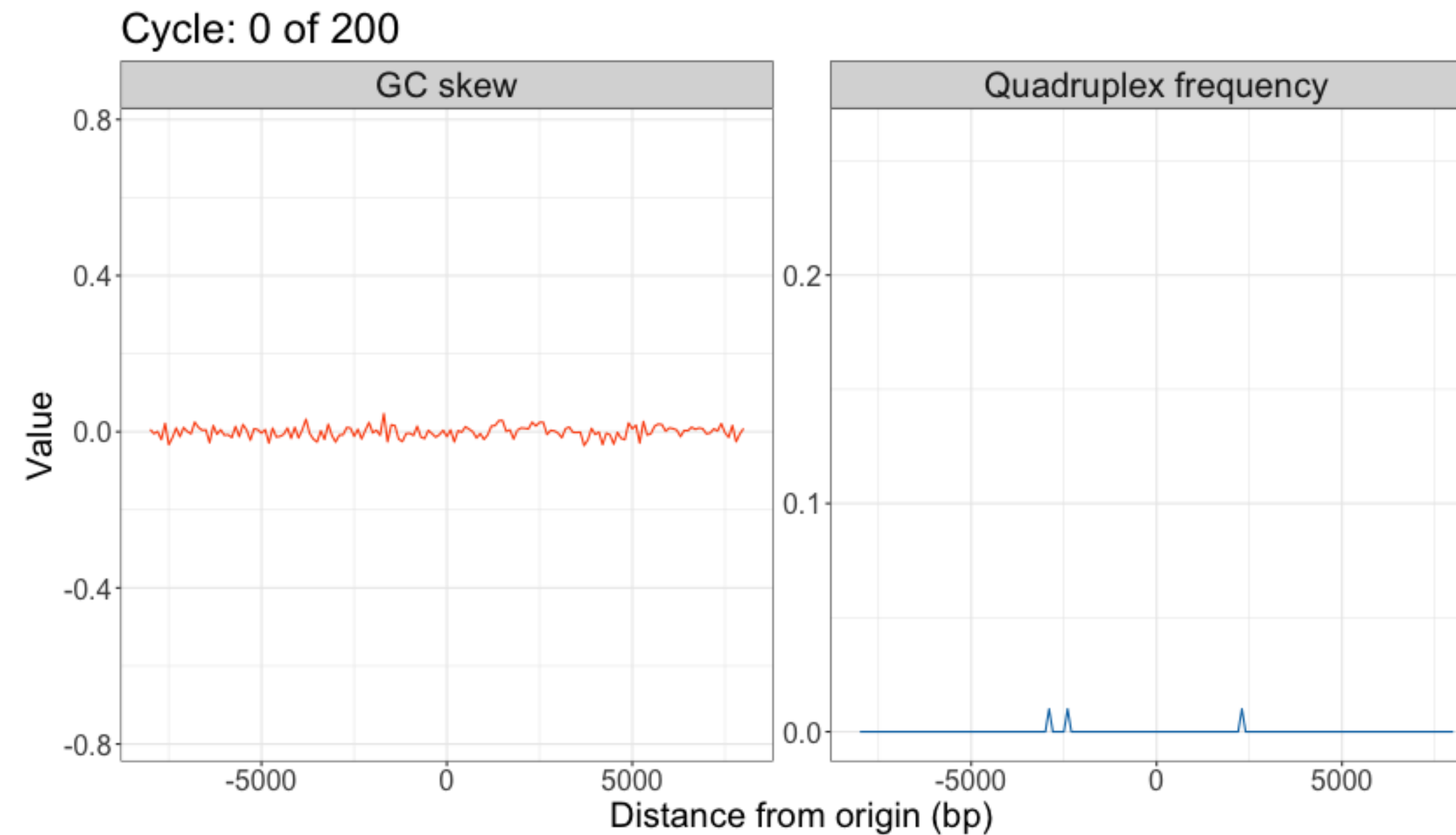


The origin of origins…
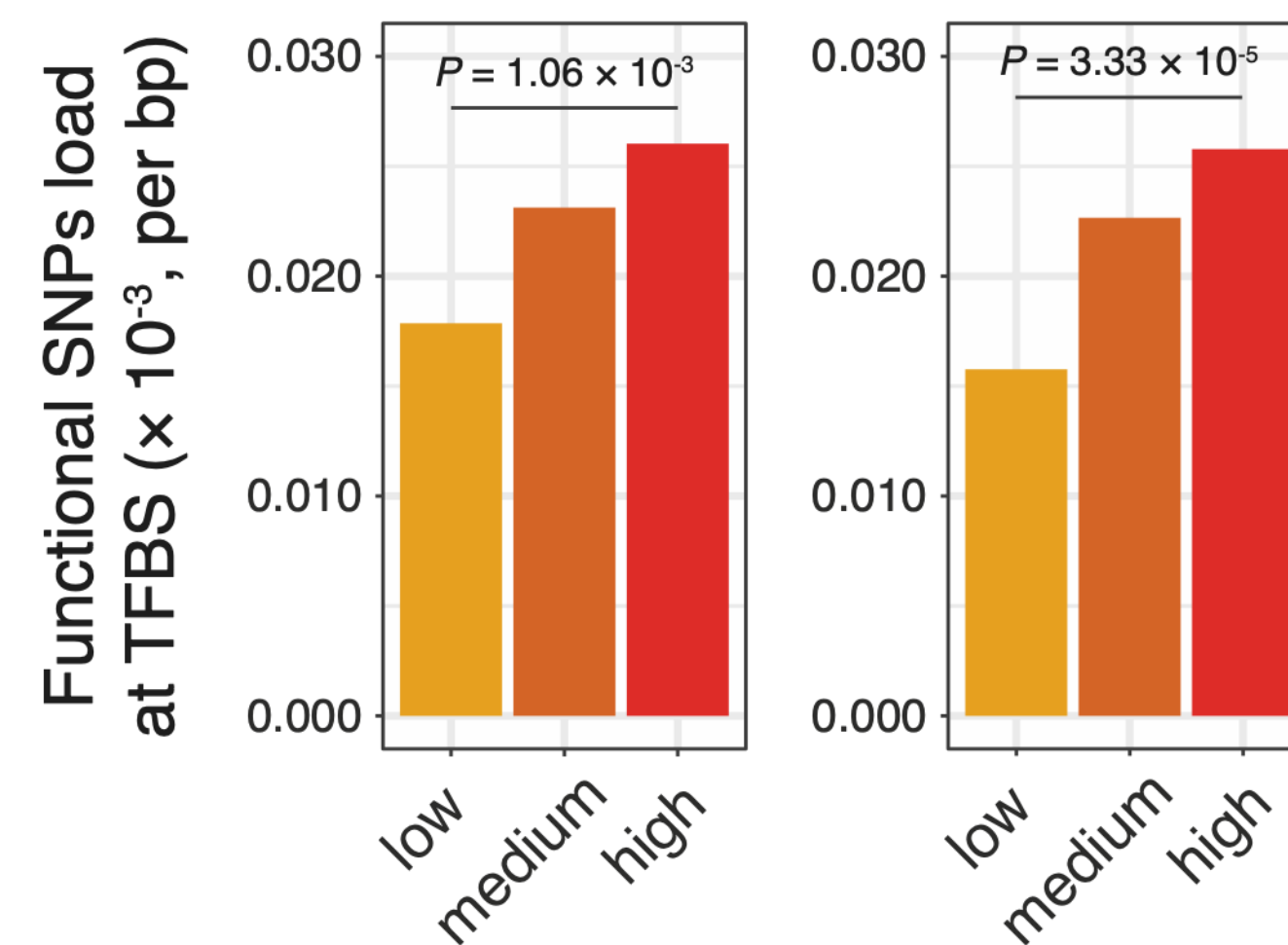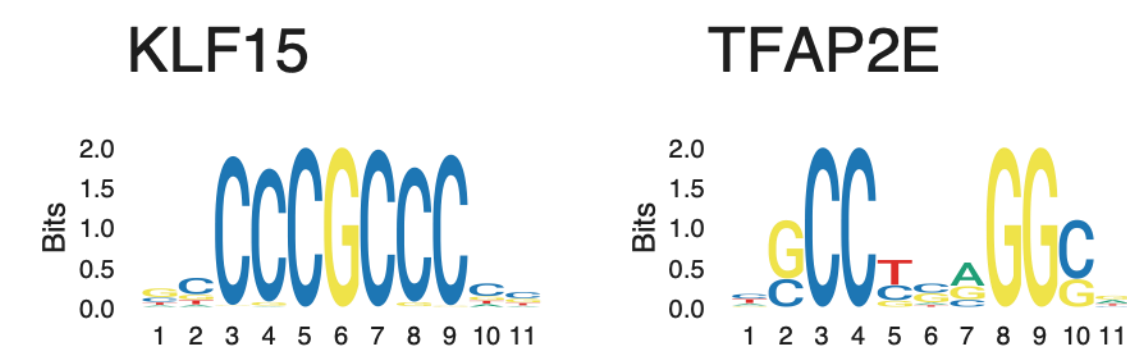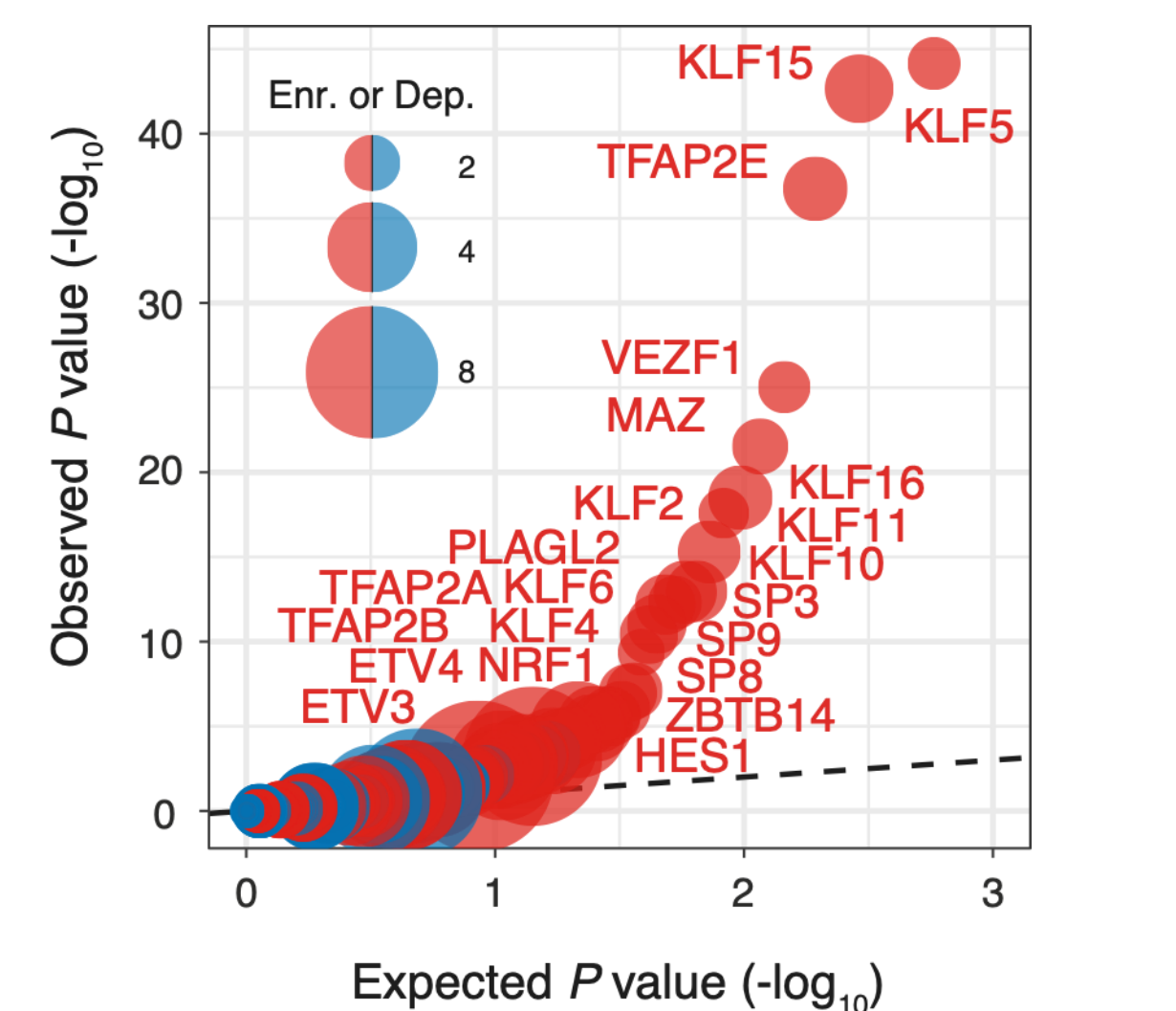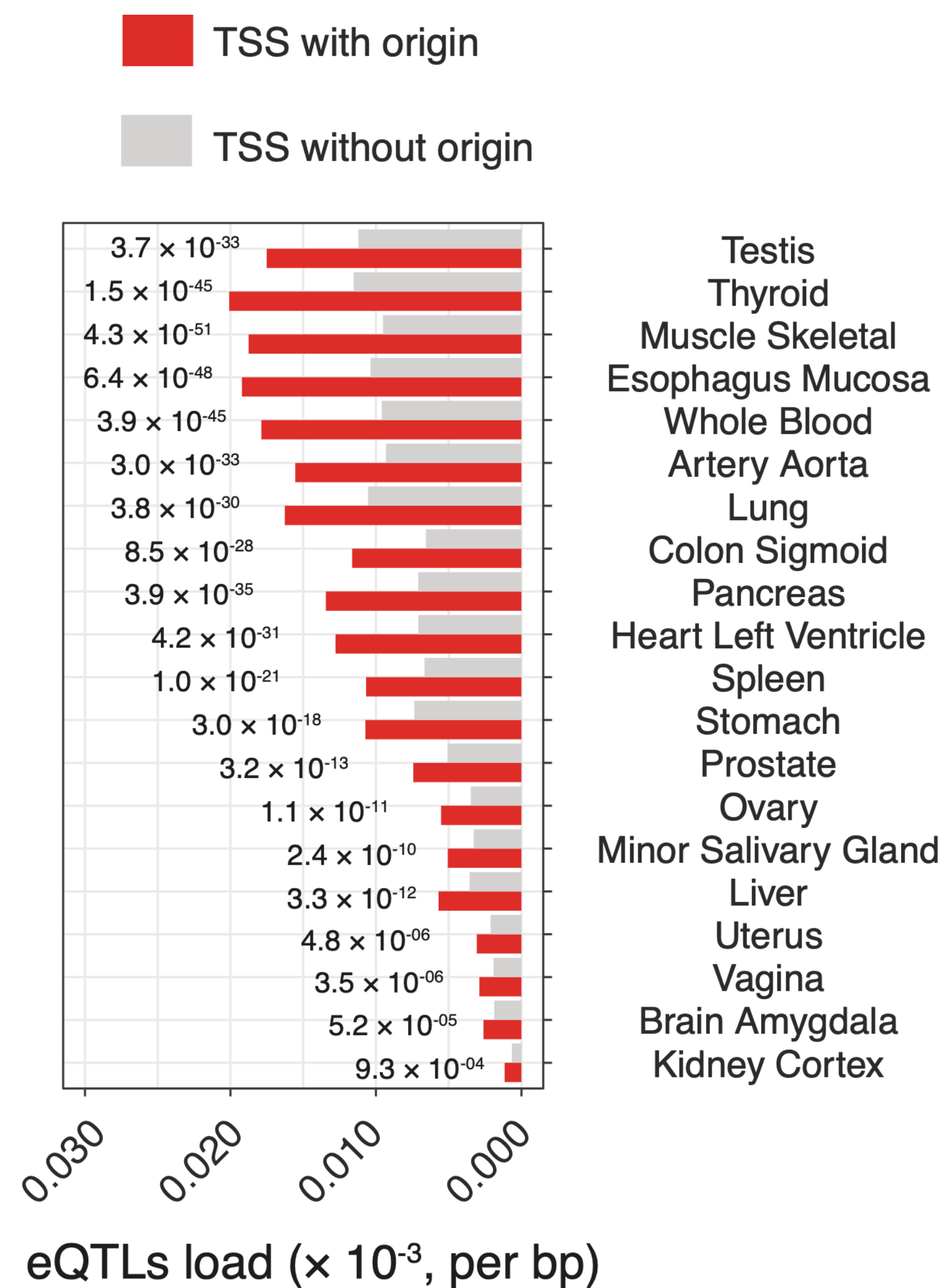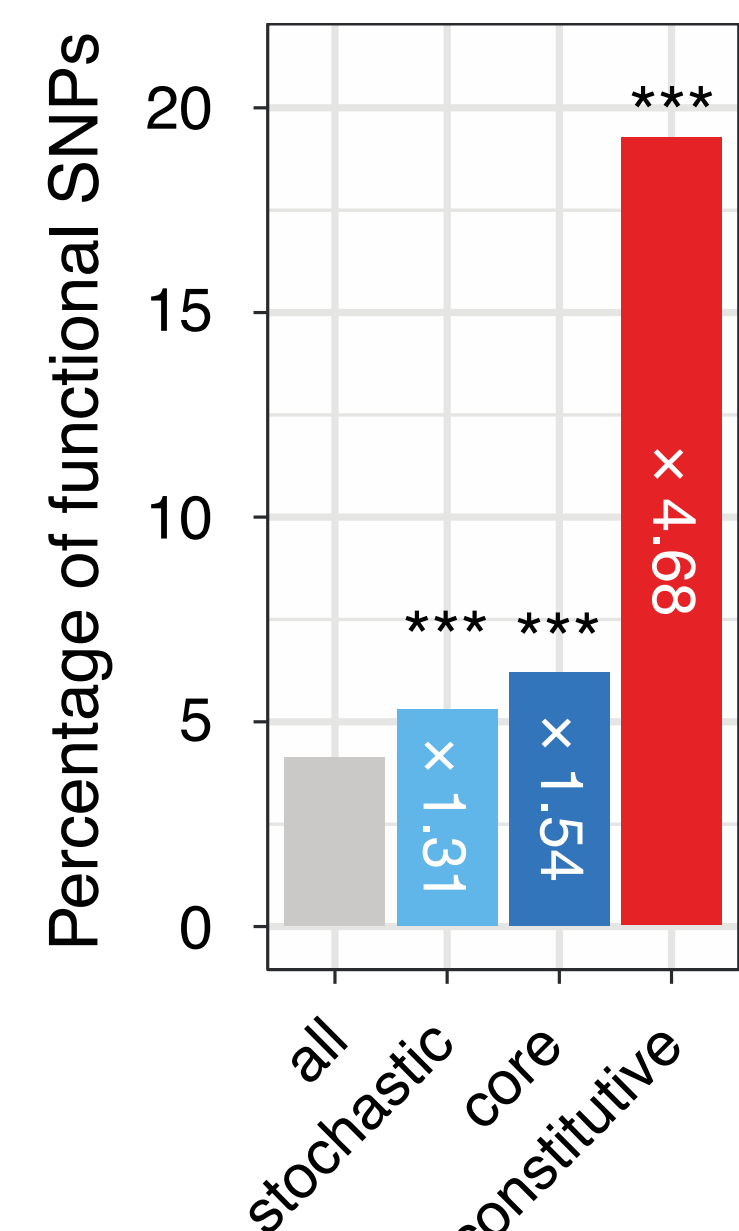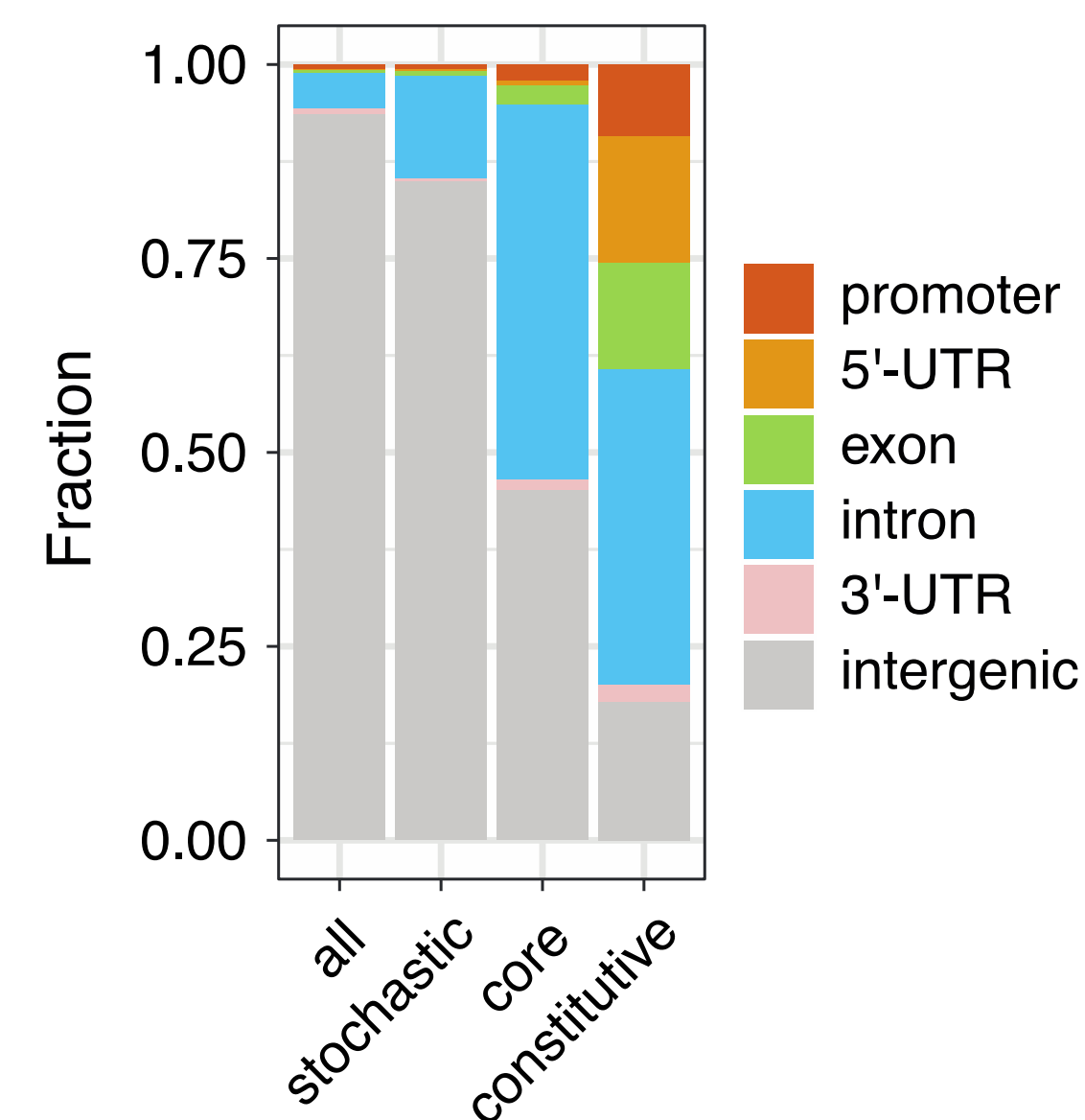Why are G4s associated with replication origins?

Murat et al., (2022) *Sci. Adv.* 8(45):eadd3686

# Constitutive origins are more likely to generate functionally important mutations



TSS with origin

TSS without origin

eQTLs load (× 10⁻³, per bp)

Testis
Thyroid
Muscle Skeletal
Esophagus Mucosa
Whole Blood
Artery Aorta
Lung
Colon Sigmoid
Pancreas
Heart Left Ventricle
Spleen
Stomach
Prostate
Ovary
Minor Salivary Gland
Liver
Uterus
Vagina
Brain Amygdala
Kidney Cortex

Murat et al., (2022) *Sci. Adv.* 8(45):eadd3686

# Summary

- Short repeat sequences with structure-forming potential are common in vertebrate genomes and have the capacity to trip up DNA replication

- Repriming is deployed frequently during replication of these sequences suggesting that they readily form replication impediments (and this can be promoted by RNA:DNA hybrid formation)

- The fork protection component Timeless links recognition of G4s via a novel DNA binding domain with recruitment of the DDX11 helicase

- The response of a model replicative polymerase to structure forming sequences *in vitro* makes powerful predictions about STR behaviour in genomes

- The identification of highly efficient sites of replication initiation has allowed the detection of replication origin dependent mutagenesis

- Highly efficient replication origins create their own sequence environment, including G4s, and are positioned to exert a significant impact on genome evolution

# Acknowledgements

MRC | Laboratory of Molecular Biology

# Features of highly efficient human replication origins

# Double strand breaks at origins occur independently of transcription



Murat et al., (2022) *Sci. Adv.* 8(45):eadd3686