

Simulating Noisy Channels in DNA Storage

Mayank Keoliya, Puru Sharma, Djordje Jevdjic
School of Computing
National University of Singapore

write process:



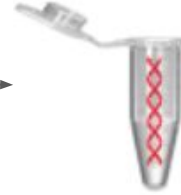
1011101...

encoding

ACGTCAG...

synthesis

Twist,
IDT



read process:



selecting



sequencing

Nanopore,
Illumina

ACGTCAG...

decoding



1011101...

DNA Storage Pipeline

Naive Method

Model described using P % error rate

$$P_{\text{insertion}} = P / 3$$

$$P_{\text{substitution}} = P / 3$$

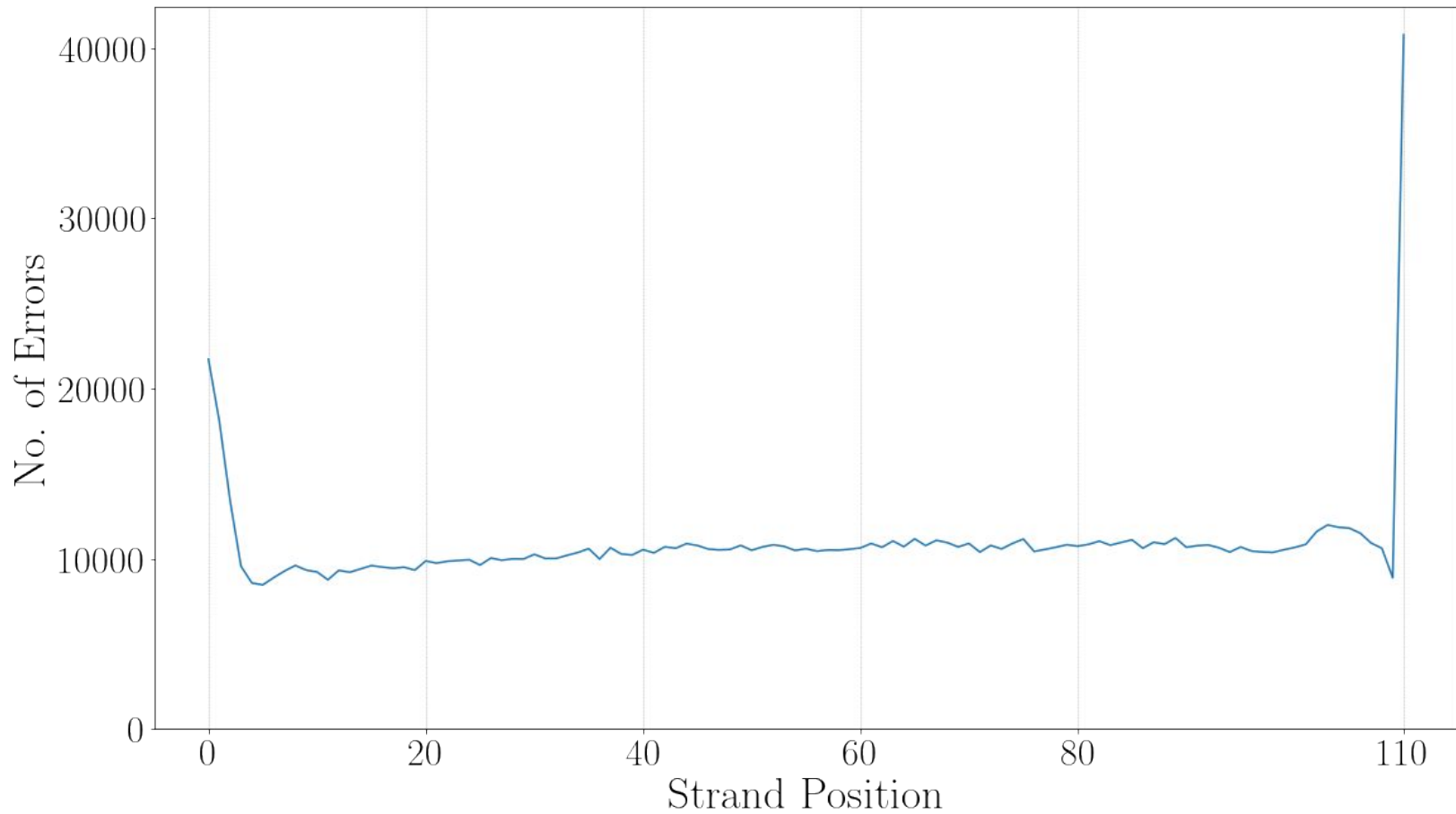
$$P_{\text{deletion}} = P / 3$$

DNASimulator

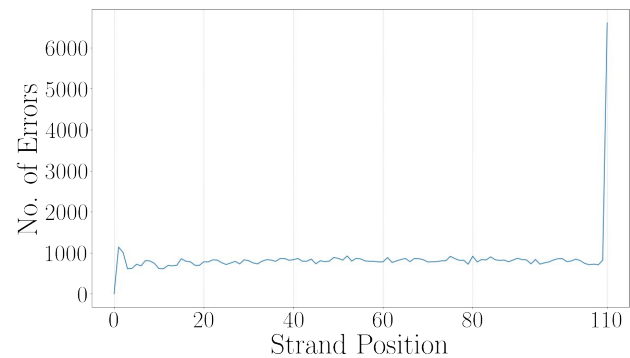
Model uses prior sequenced datasets to calculate:

for each base B (= A, G, C or T),

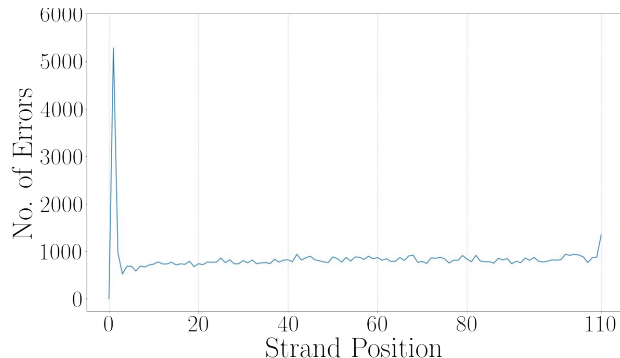
$P_{\text{ins} | B}$, $P_{\text{subs} | B}$ and $P_{\text{dels} | B}$



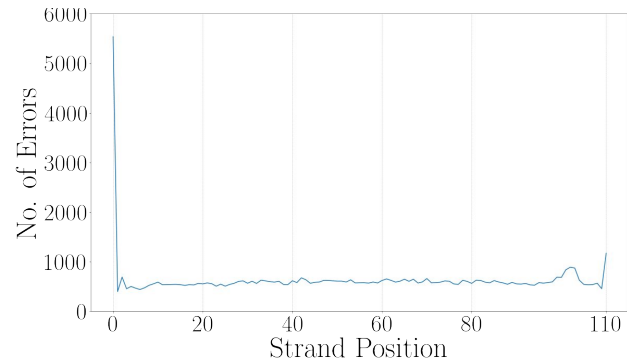
Per-Strand distribution of errors for Microsoft Nanopore dataset



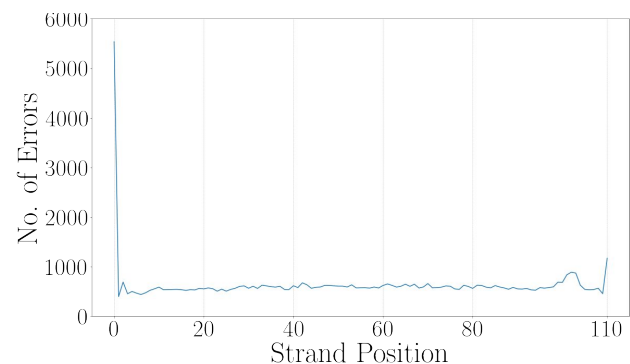
Deletion of G



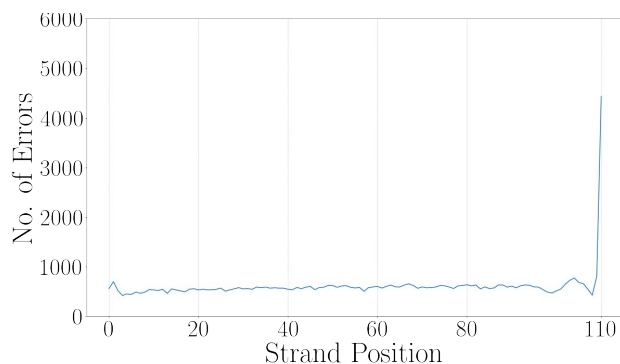
Deletion of C



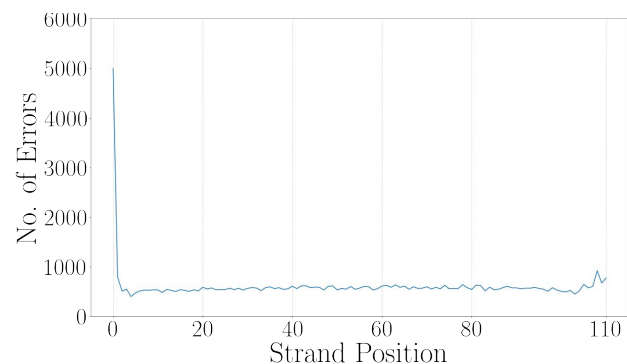
Deletion of T



Insertion of T



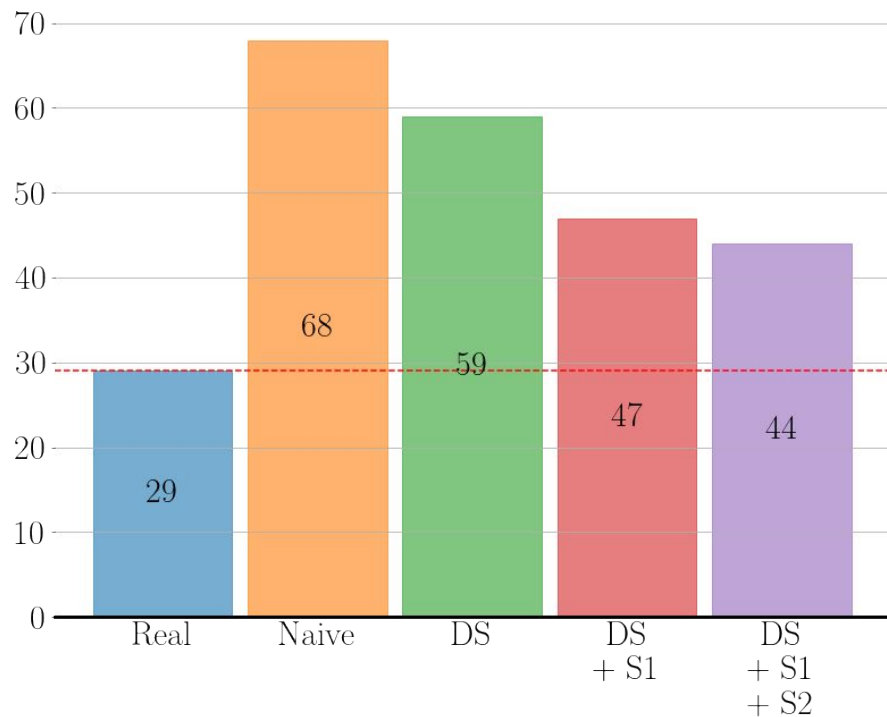
Insertion of G



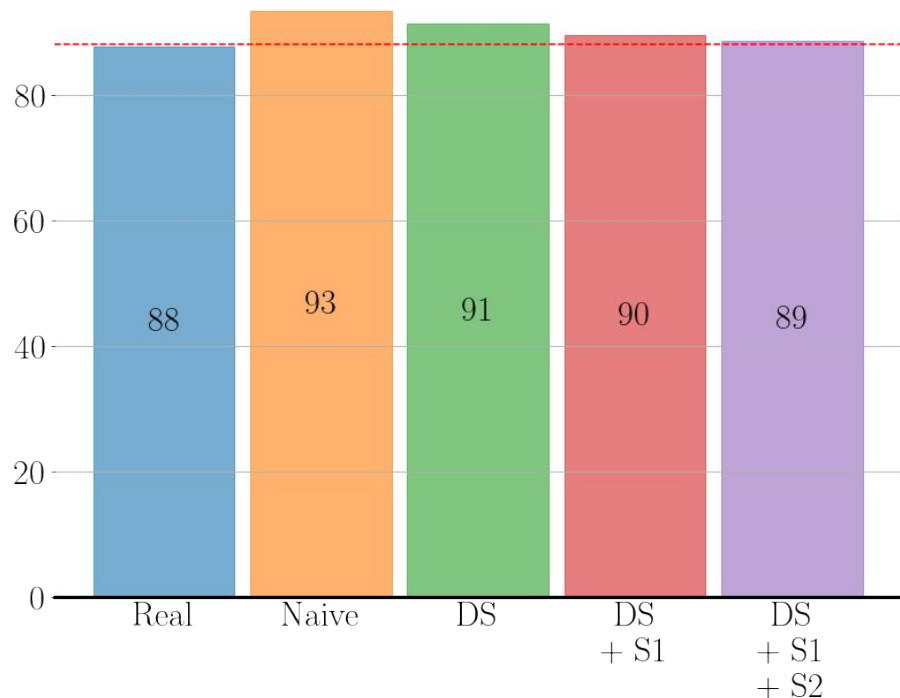
Insertion of C

Per-Strand distribution of errors **by error type** for Microsoft Nanopore dataset

Results



Per-Strand TR accuracy for datasets



Per-Char TR accuracy for datasets