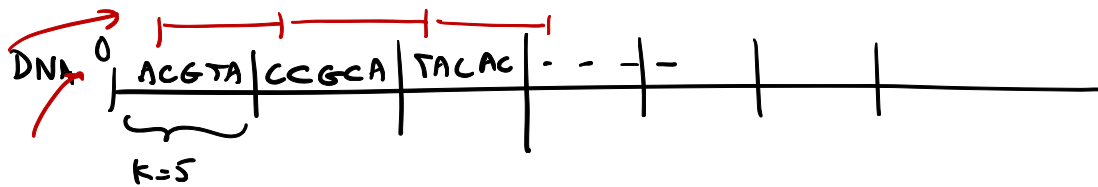


DNA/genome : 3 billion base pairs

Protein sequence : Uniprot \sim 100 million known sequence



k-mers
n-grams



k=5

I) Need 3 ORF $\begin{matrix} \nearrow 0 \\ \rightarrow 1 \\ \searrow 2 \end{matrix}$
↑
Open reading Frame \leftarrow transcription/translation

II) for k-mer
use all k offsets
 $k=0, 1, 2, \dots, k-1$

