# Motif Finding

Transcription factor

1. ttgccacaaaataatccgccttcgcaaattgacc**TACCTCAATAGCGGTA**gaaaaacgcaccactgcctgacag
2. gtaagtacctgaaagttacggtctgcgaacgctattccac**TGCTCCTTTATAGGTA**caacagtatagtctgatgga
3. ccacacggcaaataaggag**TAACTCTTTCCGGGTA**tgggtatacttcagccaatagccgagaatactgccattccag
4. ccatacccggaaagagttactccttatttgccgtgtggttagtcgctt**TACATCGGTAAGGGTA**gggattttacagca
5. aaactattaagatttttatgcagatgggtattaagga**GTATTCCCCATGGGTA**acatattaatggctctta
6. ttacagtctgttatgtggtggctgttaa**TTATCCTAAAGGGGTA**tcttaggaatttactt

Given **p** sequences, find the most mutually similar length-**k** subsequences, one from each sequence:

$$\underset{s_1,\ldots,s_p}{\mathrm{argmin}} \sum_{i<j} \mathrm{dist}(s_i, s_j)$$

dist($s_i$,$s_j$) = Hamming distance between $s_i$ and $s_j$.

Hundreds of papers, many formulations (Tompa05)