

3 - Overlaps

Monday, 11 October 2021 11:05

WEEK 3

CLOSEST MATCH LOCATED IN FINAL ROW

HORIZONTAL & VERTICAL

+1 TO ACCOUNT FOR FINAL GAP

TRACEBACK - FOLLOWING MATRIX WHERE YOU STARTED

$O(P \times L)$ APPROXIMATE MATCHING

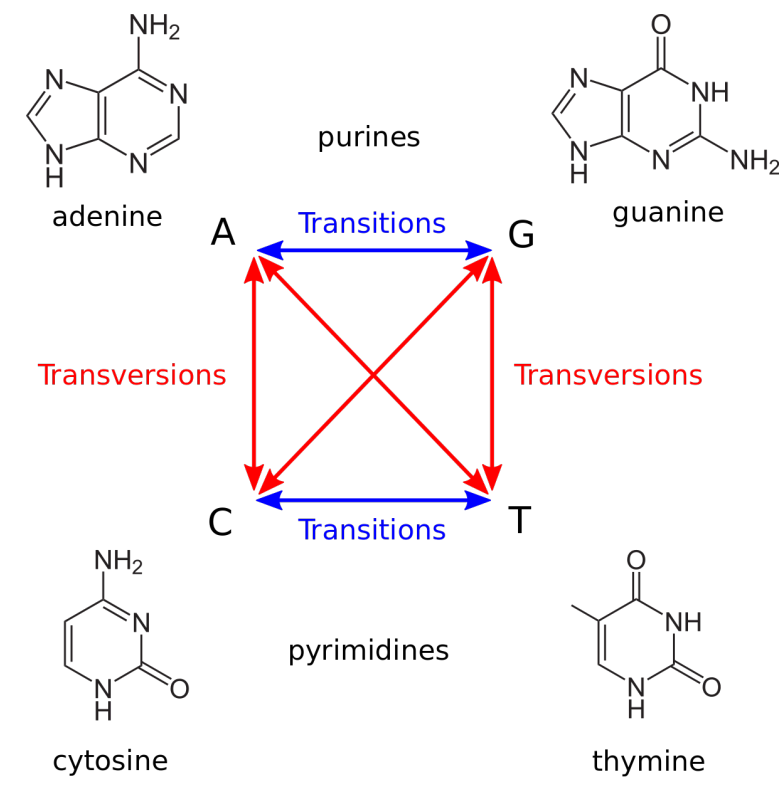
$O(L)$ EXACT MATCHING

DYNAMIC PROGRAMMING FOR EDIT DISTANCE

GLOBAL & LOCAL ALIGNMENT

IN OUR EDIT DISTANCE EVERY EDIT HAS THE SAME WEIGHT

WE CAN PENALIZE EDIT BASED ON HOW RARE THEY ARE



PROBABILISTICALLY TRANSV SHOULD BE $\times 2$ TRANSI

BUT IN REAL LIFE IT'S THE OTHER WAY AROUND

HUMAN SUBSTITUTION RATE 1/1000

SMALL GAP RATE 1/3000

INSTEAD OF +1 WE LOOK AT A $\mathbb{S}(x)$ FUNCTION THAT GIVES THE CORRESPONDING PENALTY

NEARLY LONG COMPUTATIONALLY

100 x $3.2 \cdot 10^9$ MATRIX!

WE NEED AN EFFICIENT ALGORITHM
WE CANNOT SPEND A YEAR ANALYZING
THE RESULT OF JUST 1 WEEK OF WORK

INDEXES

DON'T DEAL VERY WELL W/ MISMATCHES & GAPS

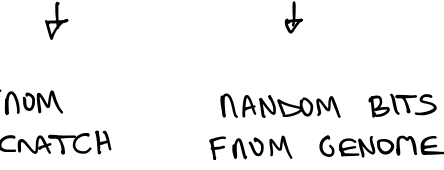
INDEX ARE BETTER SUITED FOR EXACT MATCH PROBLEM

FROM READS TO GENOME

ALIGNMENT \neq ASSEMBLY

IN ASSEMBLY YOU DON'T HAVE A REFERENCE GENOME

DE-NOVO SHOTGUN ASSEMBLY PROBLEM



THIS PROBLEM IS MORE DIFFICULT AND COMPUTATIONALLY INTENSIVE

COVERAGE - AMOUNT OF REDUNDANT INFO

NOT ALL SEQUENCES AGREE ON THE BASES

OVERAL COVERAGE WE HAVE TO KNOW $LEN(GENOME)$

IF THE 2 SEQUENCES ARE SIMILAR, IT COULD BE A HINT
THAT THE 2 READS MIGHT HAVE ORIGINATE FROM THE
SAME LOCATION IN THE GENOME

OVERLAP

IF SUFFIX & PNEFIX ARE SIMILAR :
2 SEQUENCES CAN OVERLAP

DIFFERENCES IN BASES

- MEASUREMENT ERROR
- BASECALLER SOFTWARE ERROR

POLY PLOIDY

2 CHROMOSOMES ARE DIFFERENT

MORE COVERAGE LEADS TO MORE & LONGER OVERLAPS

OVERLAPS ARE THE GWE THAT HELP US IN OUR ASSEMBLY

OVERLAPPING (DIRECT) GRAPH

NODE = CONCEPT
EDGE = RELATIONSHIP

THRESHOLD TO OVERLAP

